



**HAL**  
open science

## Theoretical and data-driven models in Ecology

Àlex Giménez-Romero

► **To cite this version:**

Àlex Giménez-Romero. Theoretical and data-driven models in Ecology. Biological Physics [physics.bio-ph]. University of the Balearic Islands (UIB); Institute for Cross-Disciplinary Physics and Complex Systems, 2024. English. ⟨NNT : ⟩. ⟨tel-04963093⟩

**HAL Id: tel-04963093**

**<https://hal.science/tel-04963093v1>**

Submitted on 24 Feb 2025

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

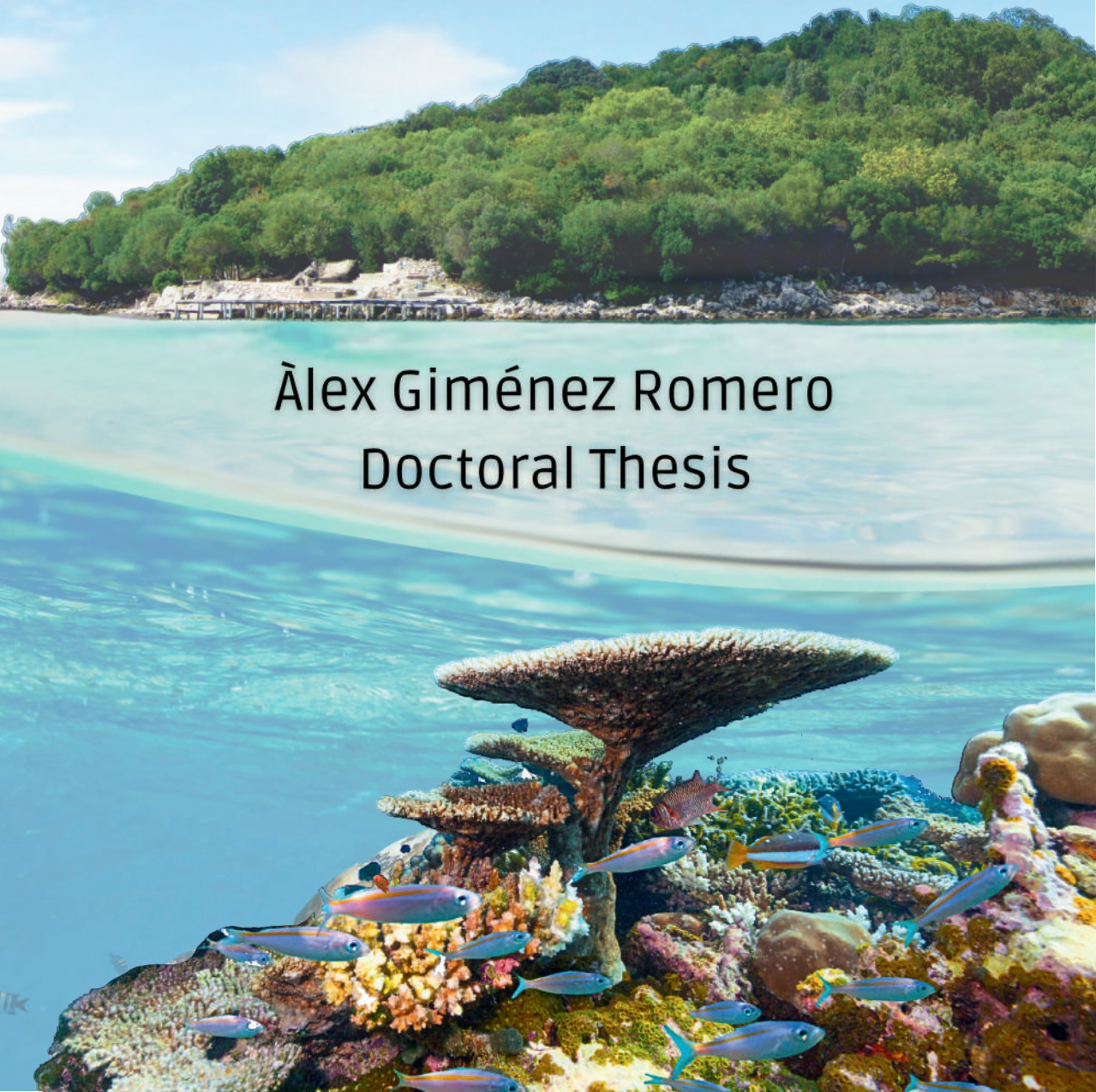
L'archive ouverte pluridisciplinaire HAL, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons CC BY-NC 4.0 - Attribution - Non-commercial use - International License

# Theoretical and data-driven models in Ecology

Àlex Giménez Romero  
Doctoral Thesis





**Universitat**  
de les Illes Balears

**DOCTORAL THESIS**  
2024

**Theoretical and data-driven  
models in Ecology**

**Àlex Giménez Romero**





**Universitat**  
de les Illes Balears

# DOCTORAL THESIS 2024

Doctoral programme in Physics

## Theoretical and data-driven models in Ecology

Àlex Giménez Romero

**Thesis Supervisor:** Manuel A. Matías  
**Thesis Tutor:** Cristóbal López Sánchez

Doctor by the Universitat de les Illes Balears

**Supervisors:**

Manuel Matías

Àlex Giménez Romero

*Theoretical and data-driven models in Ecology* ©

Palma de Mallorca, September 2024

A la meva família,  
amics, companys i a  
tothom que m'ha ajudat a  
arribar fins aquí.



Dr Manuel A. Matías of the Consejo Superior de Investigaciones Científicas (CSIC)

I DECLARE:

That the thesis title *Theoretical and data-driven models in Ecology*, presented by Àlex Giménez Romero to obtain a doctoral degree, has been completed under my supervision.

For all intents and purposes, I hereby sign this document.

Signature

Dr. Manuel A. Matías  
Thesis Supervisor

Palma de Mallorca, September 2024



## Preface

The thesis you are about to read (or probably just leaf through) is not a traditional one. I did not follow a single well-defined research line, there was not a single well-defined methodology, and I did not have a single well-defined goal. I don't think that this is inherently bad, and perhaps this situation is indeed increasingly common, but yet we are asked to write a thesis that follows a traditional structure and that tells a story that is not necessarily the one we lived. Here I would like to briefly explain the factors that led to the diversity of topics in this thesis.

The first one is of course my own personality: curious, open-minded... sorry, this is just an insider joke, look at the IFISC lemma! Now seriously, my curiosity, and perhaps my little patience, made me jump from one topic to another, from one methodology to another, and from one goal to another. Secondly, and less joyfully, the lack of funding: I never obtained a grant for my PhD. Thus, I had to constantly think about how to get funding for the next year, even if not actively, as nothing could be taken for granted. This made me go into each and almost every opportunity that appeared. But don't get me wrong, I really like working on different topics, just that this does not help to build a traditional thesis, precisely. Here I must say that I was not alone in this situation. My supervisor always had several plans for possible funding sources, and he always gave me the freedom to choose what to work on. Sadly, this is not the case for many PhD students, with or without grants. I am conscious that I was truly lucky on this matter. Finally, the nature of the scientific system that has been built over the last decades: publish or perish, impact factor, h-index, etc. I do not consider that I have made "bad" science, but perhaps I would have continued doing research on some topics that I abandoned if I had not been "worried" about the impact of my work.

At any rate, throughout the pages of this thesis and especially in the Introduction, I will try to convince you that all the work I have done is connected and that there is a common thread connecting all the topics I have worked on. Despite believing that this is true (I am of course not lying to you), if we are to be honest, this takes a secondary role.

To finish, just to comment that I have tried to write this thesis for all audiences, not only for the experts in the field. I have tried to explain the concepts in a simple way, and I have avoided the use of jargon as much as possible, especially mathematical one. I have also tried to make it enjoyable to read, with some more personal comments.

I hope you enjoy reading this thesis as much as I enjoyed when I saw it finished (I wouldn't say "as I enjoyed it writing it", precisely) and that you find it interesting and inspiring.



## Abstract

Life on Earth has evolved over billions of years, resulting in a rich diversity of species and ecosystems that provide essential services for human survival and well-being. However, this biodiversity is rapidly declining due to human activities such as habitat destruction, climate change, invasive species and emerging diseases. These interconnected drivers are causing widespread loss of species and degradation of ecosystems, threatening global ecological stability and human prosperity. Addressing this crisis requires an interdisciplinary approach to understand and mitigate its impacts, ensuring the preservation of biodiversity and the sustainability of human societies.

In this thesis we develop theoretical and data-driven methods to address pressing issues in Ecology and Conservation Biology through the lens of Complex Systems and interdisciplinary research. We address a range of contemporary challenges related to biodiversity loss driven by climate change and emerging diseases. These challenges include the spread of diseases, ocean acidification, and the decline of critical ecosystems such as coral reefs and seagrass meadows. We rely on a combination of theoretical models, computational simulations, and advanced data analysis techniques to gain a deeper understanding of these complex ecological phenomena.

In the first two parts of this thesis we develop mathematical models of disease spread to fill knowledge gaps in the transmission dynamics of marine and vector-borne plant diseases. We focus on two case studies: the Mass Mortality Event (MME) of *Pinna nobilis* and the vector-borne plant diseases caused by the bacterium *Xylella fastidiosa*. We investigate the role of key factors such as temperature or pathogen mobility in the transmission of the MME and the impact of the non-periodic seasonal abundance of insect vectors on the spread of plant diseases. These models provide insights into the mechanisms driving the dynamics of these diseases and the potential for their control and management.

In the third part, we apply this knowledge to develop a novel theoretical framework to predict the potential distribution of vector-borne plant diseases based on environmental and climatic factors. We demonstrate the utility of this model by predicting the risk of Pierce's disease of grapevines, caused by *Xylella fastidiosa*, under current and future climate scenarios. Our methodology represents a significant advancement in the field of disease biogeography, providing a way to integrate the inherent complex ecological interactions of diseases to predict their potential establishment based on environmental conditions.

Finally, in the fourth part of this thesis, we develop and apply data-driven methods to monitor and assess the health of coastal marine ecosystems. We present an innovative framework to reconstruct missing data from ocean pH

time-series using deep learning techniques, which enhance our ability to monitor ocean acidification accurately. Additionally, we employ machine learning and satellite imagery to map and evaluate the condition of seagrass meadows, offering a scalable and cost-effective approach to ecosystem monitoring. Moreover, we conduct a global analysis of the spatial properties of coral reefs using remote sensing data, uncovering universal patterns in reef size distribution and geometry. These insights are crucial for developing targeted conservation strategies to protect these vulnerable ecosystems.

This thesis underscores the importance of interdisciplinary research, integrating ecological theory, complex systems' science, and artificial intelligence to tackle ecological challenges. The findings contribute to the development of effective conservation strategies, aiming to mitigate the impacts of climate change and emergent diseases on biodiversity. Ultimately, this work supports efforts to preserve the integrity of ecosystems and ensure the sustainability of human societies in the face of ongoing environmental changes.

## Resumen

La vida en la Tierra ha evolucionado a lo largo de millones de años, produciendo una rica diversidad de especies y ecosistemas que son fundamentales para la supervivencia humana. Sin embargo, esta biodiversidad está disminuyendo rápidamente debido a actividades humanas como la destrucción de hábitats, el cambio climático, las especies invasoras y las enfermedades emergentes. Estos factores interconectados están provocando una pérdida generalizada de especies y la degradación acelerada de los ecosistemas, lo que amenaza la estabilidad ecológica global y la prosperidad humana. Para abordar esta crisis se requiere un enfoque interdisciplinario que permita comprender y mitigar sus impactos, asegurando la preservación de la biodiversidad y la sostenibilidad de las sociedades humanas.

En esta tesis desarrollamos métodos teóricos y basados en datos para estudiar problemas relacionados con la pérdida de biodiversidad causada por el cambio climático y las enfermedades emergentes. A través de la lente de los sistemas complejos, abordamos desafíos contemporáneos como la propagación de enfermedades, la acidificación de los océanos y el declive de ecosistemas críticos como los arrecifes de coral o las praderas marinas. Nos basamos en una combinación de modelos matemáticos, simulaciones computacionales y técnicas avanzadas de análisis de datos para obtener una comprensión más profunda de estos complejos fenómenos ecológicos.

En las dos primeras partes de esta tesis, desarrollamos modelos matemáticos para avanzar nuestra comprensión de la dinámica de transmisión en enfermedades marinas y transmitidas por vectores. Nos centramos en dos ejemplos: el Evento de Mortalidad Masiva (MME) de *Pinna nobilis* y las enfermedades de plantas causadas por la bacteria *Xylella fastidiosa* (Xf). Investigamos el papel de factores clave como la temperatura o la movilidad de patógenos en la transmisión del MME y el impacto de la estacionalidad en la abundancia de vectores en la propagación de enfermedades de plantas. Estos modelos brindan nuevas perspectivas sobre los mecanismos que facilitan estas enfermedades así como sobre su control y manejo.

En la tercera parte, aplicamos este conocimiento adquirido para desarrollar un nuevo marco teórico para predecir la distribución potencial de enfermedades de plantas transmitidas por vectores en función de factores climáticos. Demostramos la utilidad de este modelo al predecir el riesgo de la enfermedad de Pierce de la vid, causada por Xf, en escenarios climáticos actuales y futuros. Nuestra metodología representa un avance significativo en el campo de la biogeografía de enfermedades, permitiendo integrar las complejas interacciones ecológicas inherentes a las enfermedades para predecir su posible establecimiento en función de las condiciones ambientales.

Finalmente desarrollamos métodos basados en datos para monitorizar y evaluar la salud de los ecosistemas marinos costeros. Utilizando técnicas de aprendizaje profundo, presentamos un marco innovador para reconstruir series temporales de pH oceánico incompletas, mejorando nuestra capacidad de monitorizar la acidificación oceánica con precisión. Además, empleamos aprendizaje automático e imágenes satelitales para mapear y evaluar el estado de las praderas marinas de *Posidonia oceanica*, ofreciendo un enfoque escalable y rentable para monitorizar estos ecosistemas. Finalmente realizamos un análisis global de las propiedades espaciales de los arrecifes de coral utilizando datos de teledetección, descubriendo patrones universales en la distribución del tamaño y la geometría de los arrecifes.

Esta tesis subraya la importancia de la investigación interdisciplinar, integrando ecología, sistemas complejos e inteligencia artificial para abordar los desafíos ecológicos. Los hallazgos contribuyen al desarrollo de estrategias de conservación, con el objetivo de mitigar los impactos del cambio climático y las enfermedades emergentes. Nuestro trabajo contribuye a preservar la integridad de los ecosistemas y garantizar la sostenibilidad de las sociedades humanas frente al cambio climático.

## Resum

La vida a la Terra ha evolucionat al llarg de milions d'anys, produint una rica diversitat d'espècies i ecosistemes que proporcionen serveis essencials per a la supervivència humana. No obstant això, aquesta biodiversitat està disminuint ràpidament a causa d'activitats humanes com la destrucció d'hàbitats, el canvi climàtic, les espècies invasores i les malalties emergents. Aquests factors interconnectats estan provocant una pèrdua generalitzada d'espècies i la ràpida degradació dels ecosistemes, cosa que amenaça l'estabilitat ecològica global i la prosperitat humana. Per abordar aquesta crisi es requereix un enfocament interdisciplinari que permeti comprendre i mitigar els seus impactes, assegurant la preservació de la biodiversitat i la sostenibilitat de les societats humanes.

En aquesta tesi desenvolupem mètodes teòrics i basats en dades per estudiar problemes relacionats amb la pèrdua de biodiversitat causada pel canvi climàtic i les malalties emergents. A través de la lent dels sistemes complexos, abordem desafiaments contemporanis com la propagació de malalties, l'acidificació dels oceans i el declivi d'ecosistemes crítics com els esculls de corall o les praderies marines de *Posidonia oceanica*. Ens basem en una combinació de models matemàtics, simulacions computacionals i tècniques avançades d'anàlisi de dades per obtenir una comprensió més profunda daquests complexos fenòmens ecològics.

A les dues primeres parts d'aquesta tesi, desenvolupem models matemàtics per avançar la nostra comprensió sobre les dinàmiques de transmissió de malalties marines i transmeses per vectors. Ens centrem en dos exemples principals: l'Esdeveniment de Mortalitat Masiva (MME) de *Pinna nobilis* i les malalties de plantes causades pel bacteri *Xylella fastidiosa*. Investiguem el paper de factors clau com la temperatura o la mobilitat de patògens a la transmissió de l'MME i l'impacte de l'abundància estacional d'insectes vectors a la propagació de malalties de plantes. Aquests models proporcionen noves perspectives sobre els mecanismes que faciliten aquestes malalties així com sobre el seu control i maneig.

A la tercera part, apliquem aquest coneixement adquirit per desenvolupar un nou marc teòric per predir la distribució potencial de malalties de plantes transmeses per vectors en funció de factors ambientals i climàtics. Demostrem la utilitat d'aquest model en predir el risc de la malaltia de Pierce de la vinya, causada per *Xylella fastidiosa*, en escenaris climàtics actuals i futurs. La nostra metodologia representa un avenç significatiu en el camp de la biogeografia de malalties, ja que proporciona una forma d'integrar les interaccions ecològiques complexes inherents a les malalties per predir el seu possible establiment a funció de les condicions ambientals.

Finalment desenvolupem i apliquem mètodes basats en dades per moni-

toritzar i avaluar la salut dels ecosistemes marins costaners. Utilitzant tècniques d'aprenentatge profund, presentem un marc innovador per reconstruir sèries temporals de pH oceànic incompletes, millorant la nostra capacitat de monitoritzar la acidificació oceànica amb precisió. A més, fem servir aprenentatge automàtic i imatges satel·litàries per mapejar i avaluar l'estat de les praderies marines de *Posidonia oceanica*, oferint un enfocament escalable i rendible per a la monitorització d'aquests ecosistemes. Finalment fem una anàlisi global de les propietats espacials dels esculls de corall utilitzant dades de teledetecció, descobrint patrons universals en la distribució de la mida i la geometria dels esculls.

Aquesta tesi subratlla la importància de la investigació interdisciplinària, integrant la teoria ecològica, la ciència de sistemes complexos i la intel·ligència artificial per abordar els desafiaments ecològics. Les troballes contribueixen al desenvolupament d'estratègies de conservació efectives, amb l'objectiu de mitigar els impactes del canvi climàtic i les malalties emergents a la biodiversitat. En última instància, aquest treball recolza els esforços per preservar la integritat dels ecosistemes i garantir la sostenibilitat de les societats humanes davant dels canvis ambientals en curs.

## Agraïments

Aquesta tesi no hauria estat possible sense l'ajuda de moltes persones, així que aquesta secció també serà llarga.

En primer lloc, la persona que més ha contribuït a aquest treball és el meu director, el Dr. Manuel Matías. És gràcies a ell que he tingut l'oportunitat de realitzar aquesta tesi i començar la meua carrera investigadora. Li estic especialment agraït per confiar en mi des del primer moment, quan el meu historial acadèmic no era el més prometedor. També li estic agraït per la forma en què m'ha **guiat** durant tot el procés i tot el que m'ha ensenyat, tant en l'àmbit acadèmic com en el personal. Tinc clar que, per mi, es el millor director que hauria pogut tenir.

Els meus col·laboradors també mereixen un agraïment especial. Primer de tot vull agrair a l'Iris Hendriks, qui va ser la primera persona amb qui vaig començar a col·laborar, ja durant el meu TFM, i Amalia Grau. Juntament amb el meu director, conformen els coautors del meu primer article sobre l'esdeveniment de mortalitat massiva de les nacres. Aquest treball va ser el punt de partida de la meua tesi i va ser possible gràcies a la seva ajuda. Una menció a en Cristóbal López i en Federico Vazquez (a.k.a Fede), amb qui vam continuar treballant en aquest tema. Seguim amb Eduardo Moralejo, coautor de la majoria d'articles presents en aquesta tesi, qui em va endinsar en el món de la fitopatologia i amb qui vaig començar a treballar en el projecte de la *Xylella fastidiosa*. Aquesta col·laboració ha estat de lluny la més fructífera i puc dir que he après moltíssim treballant amb ell. Estic molt orgullós d'haver pogut mantenir una col·laboració realment interdisciplinària durant aquest anys, i realment espero que segueixi en el futur. En relació amb aquest mateix tema, també vull esmentar la Clara Lago, Arantzasu Moreno i Alberto Fereres, amb qui he tingut l'oportunitat de col·laborar. Els estic especialment agraïts per "acollir-me" durant el congrés que organitzaven a Madrid. També agrair aquí a un científic "de la casa", Jose Ramasco, no només per la part purament científica, sinó també per les converses que hem tingut. Probablement és de les persones més ocupades del IFISC, però també de les més maques. Per acabar amb aquest tema vull agrair a Jose Gutiérrez i Maialen Iturbide, qui em van acollir durant la meua estada a Santander. No tinc paraules per expressar la meua gratitud per la seva amabilitat, incloent que m'oferissin el seu cotxe per visitar la zona! Va ser una experiència molt enriquidora i espero que puguem seguir col·laborant en el futur. Un agraïment també a la Susana Flecha, amb qui vam tenir una

col·laboració molt fructífera. Va ser un plaer treballar amb ella i espero que puguem seguir treballant junts. També vull agrair a en Tomas Sintes, amb qui he tingut reunions i viatges molt divertits. Sempre recordaré el viatge a Arabia Saudita, amb en Manuel, l'Eva i en Miguel. Cal dir que es l'únic professor amb qui he compartit un llit! Ha estat un plaer treballar amb ell i espero que puguem seguir col·laborant en el futur. Finalment un agraïment especial a en Carlos Duarte, amb qui he tingut l'oportunitat de col·laborar en un projecte molt interessant sobre coralls. Primer de tot, agrair la seva invitació a Arabia Saudita, que com ja he dit va ser una experiència molt enriquidora i divertida. També agrair-li per la seva ajuda i el coneixement que m'ha transmès, encara que sigui de forma inconscient. Espero que aquesta col·laboració segueixi en el futur.

També vull agrair alguns col·laboradors que no han estat coautors dels articles d'aquesta tesi, però que han contribuït de forma significativa. Primer de tot vull agrair a en Dan Bebbler, amb qui vaig tenir l'oportunitat de treballar durant un mes en remot. També una menció especial a en Rob Salguero, amb qui vaig tenir el plaer de treballar durant 3 mesos en una estada al departament de biologia de la universitat d'Oxford. Va ser una experiència molt enriquidora poder veure com es treballa en un entorn tan diferent al de l'IFISC. Aquests tres mesos van ser també molt fructífers i vaig aprendre moltíssim. Va ser un plaer treballar amb ell i espero que puguem seguir col·laborant en el futur. També vull agrair a en Pere Colet i la Rosa López, amb qui he tingut l'oportunitat de donar classes. Han confiat en mi tant per donar classes de grau com de màster i els estic molt agraït per això. Finalment, agrair a en Luís Gordillo, a qui vaig tenir el plaer de conèixer en persona durant la seva visita a l'IFISC i amb qui segueixo en contacte telemàtic.

Un agraïment especial a tots els membres de l'IFISC i la UIB. En primer lloc, a l'equip de neteja, invisibles però imprescindibles. En especial a l'Anna Pérez, amb qui coincidí gairebé cada matí, ben d'hora. També agrair els membres de l'administració per la seva ajuda en tot moment, sobretot a la Marta Ozonas, l'Inma Carbonell, la Neus Lacomba, la Maria Quetglas i l'Alberto Sánchez. També agrair a l'Adrián García tota la seva ajuda en la part de divulgació científica. Sense ell tot el que hem fet en aquest aspecte no hauria estat possible. Finalment, un agraïment més que especial als tècnics (i ex-tècnic) de l'IFISC: l'Eduard Solivellas, l'Antònia Tugores i en Rubén Tolosa. Literalment, sense ells aquesta tesi (i totes les de l'IFISC) no haurien estat possibles. Els estic molt agraït per tota la seva ajuda, les seves explicacions i la seva paciència, a més d'altres conversacions més "informals".

També vull agrair a tots els meus amics i companys del IFISC, que han estat una part molt important de la meua vida durant aquests anys. M'agradaria començar amb els que fa més temps que conec. El primer lloc indiscutible l'ocupa l'Adrià Labay (a.k.a DJ Labay), amb qui vaig començar la carrera de Física a la UAB, amb qui casi morim al Monte Perdido (i a les festes de Martorell, les colònies o Dresden) i sense qui, probablement, no hauria acabat la carrera. Els següents a la llista són en David, Miguel, los canarios (Javi Galvan & Medi) i Jogito, qui conec des del màster. Amb ells he compartit rutillas, festes, molts exercicis i entregues i, sobretot, moltes rialles. Una menció especial a en David, qui també va ser un veí i un usuari del taxiUIB, i ha estat una ajuda imprescindible en la part d'escriptura d'aquesta tesi: gracies per anar dues setmanes adelantat. També per en Miguel, amb qui he tingut converses molt interessants i de qui admiro la seva curiositat i ganes de seguir aprenent (i la locura de entrecot que nos hizo en su casa). Seguim amb els que vaig conèixer entre el màster i el començament del doctorat: Pablo, Irene, Maria, Javi Aguilar, Medea, Rodri i Mou (o Mou i Rodri). Ells em van acollir quan vaig arribar per fer el doctorat i es van encarregar de fer la vida post-pandèmia més fàcil i divertida. En Pablo, la Irene i la Maria mereixen una menció especial per acompanyar-me en les meves rutes chill i no tan chill, a Sabotage i en general aquests anys. El següent grup el formen els increïbles companys i fundadors de ZULO: Mar Ferri, Fer (desertor, pero tt gym), Mar Cuevas, Gorka, Guillem, Pau i Manuel Miranda. Aquests 10 metres quadrats sense llum natural ni ventilació haurien estat insuportables sense vosaltres. Aquí també agrair les noves incorporacions, com en Jaume, la Sara o en Coque. Finalment tenim alguns dels integrants de la S07: Jesus, Bea, Lisa, Pepe, Dimitris i Daniele. Una menció especial a Jesus i Bea, per haver compartit mes temps amb mi i, en conseqüència, haver aguantat més chapes i gaudit de més barbacoas amb tremendo allioli. També voldria agrair als postdocs amb qui hem compartit dinars i altres moments, sobretot l'Eva. Finalment, vull tornar a agrair en Manuel Matías, a qui també considero un amic. Espero que, a part de seguir col·laborant, puguem seguir compartint sopars i copes de vi.

Per acabar, vull agrair als meus amics més antics. Als companys i amics de la carrera: Bernat, Roger, Uri, Linde, Jota i Dani. Tot i que no ens haguem vist gaire, ni jo parli massa pel grup, sempre heu estat aquí. En especial agrair a en Dani, qui sempre ha estat disponible per fer una birra i uns bons nachos al Livingstone. També agrair als amics del barri, especialment a la Lidia, el Víctor, l'Alex Jimenez i el Guti. Fora de l'àmbit acadèmic, són els amics que més han escoltat tant els meus èxits i avanços com les meves queixes, frustracions

i fracassos. També agrair a la Carmela, que ha estat una part molt important de la meua vida durant aquests anys. Ella es l'única persona que ha estat al meu costat les 24 hores del dia, 7 dies a la setmana, 365 dies a l'any. Sense ella segurament hauria acabat aquesta tesi, pero probablement també hauria acabat amb la meua salut mental. Tinc molta sort de tenir-la al meu costat, des de Barcelona, a Mallorca, a Oxford i properament a Blanes; a tot arreu on de moment m'ha seguit. Li estic especialment agraït. També he de dir que si no fos per ella moltes figures d'aquesta tesi serien apreciablement més lletges... Finalment vull agrair a la meua família per haver-me donat l'oportunitat de formar-me i per haver-me donat suport en tot moment.

# Contents

Preface .....	vii
Abstract .....	ix
Acknowledgements .....	xv
List of figures .....	xxvi
List of tables .....	xxx
List of acronyms .....	xxxvii
List of publications .....	xxxix
<b>1 Introduction .....</b>	<b>1</b>
1.1 The global biodiversity crisis .....	1
1.2 Complex systems in Ecology .....	4
1.3 Mathematical and computational models .....	7
1.4 Thesis structure .....	10
<b>2 Scientific background .....</b>	<b>13</b>
2.1 Infectious disease modeling .....	13
2.1.1 Compartmental models .....	13

2.1.2	Linear stability analysis	20
2.1.3	The next generation matrix method	22
2.1.4	Individual-based models	24
<b>2.2</b>	<b>Disease biogeography</b>	<b>27</b>
2.2.1	Ecological niches	28
2.2.2	Species distribution models	29
2.2.3	Limitations of SDMs	33
<b>2.3</b>	<b>Data-driven methods</b>	<b>34</b>
2.3.1	Machine learning	35
2.3.2	Artificial neural networks	36
2.3.3	Time series forecasting	39
2.3.4	Image segmentation	41
2.3.5	Performance metrics	45
<b>2.4</b>	<b>Ecological systems</b>	<b>47</b>
2.4.1	The Mass Mortality Event of <i>Pinna nobilis</i>	47
2.4.2	<i>Xylella fastidiosa</i> : an emerging global threat	48
2.4.3	The decline of seagrass meadows	51
2.4.4	The structure of coral reefs	53
2.4.5	Ocean acidification	54

I

**Parasite-produced marine diseases**

<b>3</b>	<b>The case of the mass mortality event of <i>Pinna nobilis</i></b>	<b>61</b>
<b>3.1</b>	<b>Introduction</b>	<b>62</b>
<b>3.2</b>	<b>The SIRP model</b>	<b>64</b>
3.2.1	Model structure and initial considerations	64
3.2.2	General SIRP model	66
3.2.3	Model reduction	70
<b>3.3</b>	<b>Numerical analysis</b>	<b>72</b>
3.3.1	The basic reproduction number $R_0$	73
3.3.2	Final state of the epidemic	74
3.3.3	Maximum of infected individuals	75
3.3.4	Numerical verification of the fast-slow approximation	76
3.3.5	Numerical verification of the exact reduction	78
<b>3.4</b>	<b>Model validation</b>	<b>79</b>
<b>3.5</b>	<b>Conclusions</b>	<b>84</b>

<b>4</b>	<b>Spatial effects in parasite-induced marine diseases</b>	<b>87</b>
4.1	Introduction	88
4.2	The SIRP spatial model	90
4.3	Results	92
4.3.1	Non-spatial limit	92
4.3.2	Approximate relation between parasites and infected hosts	94
4.3.3	Spatial threshold	96
4.3.4	Spreading speed of the infected population and time to extinction	99
4.4	Conclusions	101

## II

## Modeling vector-borne plant diseases

<b>5</b>	<b>Non-stationary vector populations</b>	<b>109</b>
5.1	Introduction	110
5.2	The model	112
5.2.1	Preliminary analysis of the model	113
5.3	Results	114
5.3.1	Epidemic threshold and disease dynamics	114
5.3.2	The basic reproduction number for non-stationary vector populations	119
5.3.3	Fast-slow approximation	121
5.3.4	Reduction to a SIR model	123
5.4	Conclusions	124
<b>6</b>	<b>A compartmental model for <i>Xylella fastidiosa</i> diseases</b>	<b>127</b>
6.1	Introduction	128
6.2	Materials and Methods	130
6.2.1	Epidemic model: the SEIR-V model	130
6.2.2	Basic reproductive number	133
6.2.3	Epidemiological data	133
6.2.4	Model fitting through Bayesian Inference	134
6.2.5	Sensitivity Analysis	135
6.3	Results	136
6.3.1	Model fit and parameter estimates	136
6.3.2	Global Sensitivity Analysis	141
6.3.3	Epidemic control through vector management	142
6.4	Discussion	143

### III Modeling the risk of vector-borne plant diseases

<b>7</b>	<b>Global predictions for Pierce's disease risk</b>	<b>151</b>
<b>7.1</b>	<b>Introduction</b>	<b>152</b>
<b>7.2</b>	<b>Methods</b>	<b>154</b>
7.2.1	Inoculation tests	154
7.2.2	Modified Growing Degree Days	155
7.2.3	Disease progress with temperature	155
7.2.4	Disease recovery through winter curing	156
7.2.5	Global climate data, MGDD/CDD computation	157
7.2.6	Disease model construction	157
7.2.7	Model calibration and validation	159
7.2.8	<i>Philaenus spumarius</i> SDM	160
7.2.9	Distribution of wine-grape production areas	160
7.2.10	Risk assessment by 2050	160
<b>7.3</b>	<b>Results</b>	<b>161</b>
7.3.1	Thermal requirements to develop PD	161
7.3.2	MGDD/CDD distribution maps	163
7.3.3	PD global risk	166
7.3.4	PD risk projections for 2050	169
7.3.5	Risk based on vector information	171
7.3.6	Combining vineyard land cover across Europe with the model output	171
<b>7.4</b>	<b>Discussion</b>	<b>172</b>
<b>8</b>	<b>Pierce's disease risk under global warming</b>	<b>177</b>
<b>8.1</b>	<b>Introduction</b>	<b>178</b>
<b>8.2</b>	<b>Methods</b>	<b>179</b>
8.2.1	Climate datasets	179
8.2.2	A climate-driven epidemiological model	180
8.2.3	Model adaptation to daily temperature data	182
8.2.4	Vector climatic suitability	182
8.2.5	Risk velocity	183
<b>8.3</b>	<b>Results</b>	<b>184</b>
8.3.1	Present and future climate suitability	184
8.3.2	Pierce's disease risk projections under climate change	186
<b>8.4</b>	<b>Discussion</b>	<b>190</b>

<b>9</b>	<b>Pierce's disease risk with high-resolution climate data ...</b>	<b>195</b>
<b>9.1</b>	<b>Introduction</b> .....	<b>196</b>
<b>9.2</b>	<b>Results</b> .....	<b>198</b>
9.2.1	Global differences in PD risk between coarse and fine-grain climate data	198
9.2.2	Pierce's disease risk surges in previously unresolved microclimates	201
<b>9.3</b>	<b>Discussion</b> .....	<b>204</b>
<b>9.4</b>	<b>Methods</b> .....	<b>206</b>
9.4.1	Climate data	206
9.4.2	Vector climatic suitability	206
9.4.3	Vineyard data	206
9.4.4	Model adaptation to daily temperature data	206

## IV Data-driven methods for ecological problems

<b>10</b>	<b>Reconstructing pH time-series with machine learning ...</b>	<b>211</b>
<b>10.1</b>	<b>Introduction</b> .....	<b>212</b>
<b>10.2</b>	<b>Results</b> .....	<b>215</b>
10.2.1	Time series data	215
10.2.2	Reconstruction pH time series with Deep Learning	217
<b>10.3</b>	<b>Discussion</b> .....	<b>220</b>
<b>10.4</b>	<b>Methods</b> .....	<b>223</b>
10.4.1	Data collection	224
10.4.2	Data processing	225
10.4.3	Computing the trend of seasonal data	225
10.4.4	Selecting the best neural network architecture	226
<b>11</b>	<b>Mapping seagrass meadows from space .....</b>	<b>229</b>
<b>11.1</b>	<b>Introduction</b> .....	<b>230</b>
<b>11.2</b>	<b>Results</b> .....	<b>232</b>
11.2.1	A deep learning framework for automated marine ecosystem labelling	232
11.2.2	A reliable AI-based solution for marine ecosystem monitoring	234
11.2.3	Towards a comprehensive model for the Mediterranean Sea	237
<b>11.3</b>	<b>Discussion</b> .....	<b>239</b>
<b>11.4</b>	<b>Methods</b> .....	<b>242</b>
11.4.1	Satellite data	242
11.4.2	Habitat data	242
11.4.3	Bathymetric data	243

11.4.4	Dataset creation	243
11.4.5	Deep learning models	244
11.4.6	Model training	244
11.4.7	Performance metrics	245
11.4.8	Consensus prediction	246
11.4.9	Model selection	246

## 12 Universal spatial properties of coral reefs ..... 249

<b>12.1</b>	<b>Introduction</b>	<b>250</b>
<b>12.2</b>	<b>Results</b>	<b>251</b>
12.2.1	Coral reefs macroecological patterns	251
12.2.2	The fractal nature of coral reefs	253
12.2.3	Fractality extends up to coral provinces	255
<b>12.3</b>	<b>Discussion</b>	<b>256</b>
<b>12.4</b>	<b>Methods</b>	<b>259</b>
12.4.1	Global coral reef data	259
12.4.2	Coral reefs as clusters of connected coral/algae class polygons	259
12.4.3	Coral reefs area, perimeter and inter-reef distance	259
12.4.4	Coral reef size distribution	260
12.4.5	Fractal dimensions from area-perimeter relation	260
12.4.6	Box-Counting fractal dimension	260
12.4.7	Compactness and elongation index	261

# V

## Discussion

## 13 Main contributions ..... 265

13.1	Marine epidemiology	265
13.2	Vector-borne plant diseases	267
13.3	Disease biogeography	268
13.4	Ecological monitoring & analysis	270

## 14 General discussion ..... 273

<b>Bibliography</b>	<b>279</b>
---------------------	------------

<b>Appendices</b>	<b>329</b>
-------------------	------------

<b>A</b>	<b>Analysis of the SIRP model</b>	<b>329</b>
A.1	Finding a conserved quantity for the SIRP model	329
A.2	Stability analysis of the fixed points of the SIRP model	331
A.3	Calculation of $R_0$ using the Next Generation Matrix method	332
A.4	Sensitivity Analysis	332
A.5	General rate change with temperature	333
A.6	Derivation of the non-spatial equation for $R_\infty$	335
<b>B</b>	<b>The SIR-V model and its application to Xf diseases</b>	<b>337</b>
B.1	Calculation of $R_0$ from standard methods	337
B.2	Calculation of $R_0$ for non-stationary vector populations	339
B.3	Determination of $R_0$ for Xf diseases	340
<b>C</b>	<b>Modeling Pierce's disease risk</b>	<b>343</b>
C.1	Inoculation tests on European grapevine varieties	343
C.2	Modeling climate suitability for PD	350
C.2.1	Modified Growing Degree Days (MGDD) from Arrhenius Equation	350
C.2.2	Relation between MGDD and within-plant bacterial population	352
C.2.3	Epidemiological and theoretical basis	353
C.2.4	Model validation	354
C.2.5	Determination of $R_0$ for Europe	356
C.2.6	Simulation details	357
C.2.7	Vector distribution influence	359
C.3	Future risk extrapolation	360
C.4	Mathematical justifications	362
C.4.1	Linear scaling of $R_0$ with vector population	363
C.4.2	Reduction to a SIR model	363
C.5	Comparison of MGDD calculations from different models	364
C.6	Analysis of MGDD and CDD time series	366
C.7	Detailed analysis of PD risk	367
C.8	Present and future climate suitability for Xf & Ps	380
C.9	Future Pierce's disease risk projections under climate change	381
C.10	Risk analysis within European wine PDOs	382
C.11	MGDD and CDD approximation	385
C.12	The effect of spatial resolution on risk projections	388
C.13	<i>Vitis vinifera</i> global distribution	391

<b>D</b>	<b>Nonlinear time-series analysis and reconstruction</b>	<b>393</b>
D.1	Seasonal adjusted fits for pH and temperature	393
D.2	Total alkalinity in the Bay of Palma	396
D.3	Bidirectional Long-Short Term Memory neural network	396
<b>E</b>	<b>Mapping marine habitats with deep learning</b>	<b>399</b>
<b>E.1</b>	<b>Satellite imagery and ground truth data</b>	<b>399</b>
E.1.1	Study region covering	399
E.1.2	Ground truth dataset composition	402
<b>E.2</b>	<b>Dataset creation</b>	<b>402</b>
<b>E.3</b>	<b>Deep learning models</b>	<b>403</b>
E.3.1	UNET	404
E.3.2	Linknet	404
E.3.3	FPN	404
E.3.4	PSPNet	405
E.3.5	Backbones	405
<b>E.4</b>	<b>Performance Metrics</b>	<b>406</b>
E.4.1	Accuracy	406
E.4.2	Precision	406
E.4.3	Recall	406
E.4.4	F1 Score	406
E.4.5	Cohen's Kappa	407
E.4.6	Intersection over Union (IoU)	407
<b>E.5</b>	<b>Spectral reflectance analysis</b>	<b>407</b>
<b>E.6</b>	<b>Architecture selection</b>	<b>409</b>
<b>E.7</b>	<b>Out-of-sample predicting power and robustness</b>	<b>409</b>
E.7.1	Understanding model performance	412
<b>E.8</b>	<b>Effect of depth on model performance</b>	<b>416</b>
<b>E.9</b>	<b>CAMELE trained with all available data</b>	<b>417</b>
<b>F</b>	<b>Coral reefs' macroecological patterns</b>	<b>419</b>
<b>F.1</b>	<b>Macroecological patterns &amp; fractal dimension</b>	<b>420</b>
<b>F.2</b>	<b>Coral reef size distribution</b>	<b>422</b>

## List of Figures

1.1	The Biodiversity Intactness Index . . . . .	2
1.2	Vegetation patterns in arid ecosystems . . . . .	5
2.1	Numerical and analytical analysis of the SIR model. . . . .	19
2.2	Historical development of the ecological niche concept . . . . .	29
2.3	Example of a species distribution model . . . . .	31
2.4	The perceptron model of a biological neuron . . . . .	37
2.5	The feedforward neural network . . . . .	38
2.6	Simple Recurrent Neural Network . . . . .	40
2.7	Convolutional Neural Network . . . . .	42
2.8	Feature maps in a convolutional neural network . . . . .	44
2.9	U-Net architecture for image segmentation . . . . .	45
2.10	The fan mussel <i>Pinna nobilis</i> . . . . .	47
2.11	Symptoms of <i>Xylella fastidiosa</i> infection . . . . .	49
2.12	Distribution of <i>Xylella fastidiosa</i> in Europe . . . . .	50
2.13	The seagrass <i>Posidonia oceanica</i> . . . . .	52
2.14	The structure of a coral reef . . . . .	53
2.15	Global trend in ocean pH . . . . .	55
3.1	SIRP model flow diagram . . . . .	67
3.2	Sensitivity analysis of $R_0$ . . . . .	74
3.3	Sensitivity analysis of the final number of dead individuals ( $R_\infty$ ) . . . . .	75
3.4	Global sensitivity analysis of the maximum of infected individuals and its time occurrence . . . . .	76

3.5	Numerical validation of the timescale approximation for the parasite population	77
3.6	Numerical verification of the exact model reduction and the subsequent approximation	78
3.7	Parameter estimation of the exact reduction of the SIRP model	81
3.8	Parameter estimation of the approximate SIR model	83
4.1	Scheme of the individual-based SIRP model	91
4.2	Comparison between the non-spatial model and the individual-based model in the high mobility limit	93
4.3	Validation of the approximate expression for the parasite population dynamics	95
4.4	Phase diagram and fit for the transition between the disease-free and propagation phases	98
4.5	Analysis of the spreading speed of the infected population	99
4.6	Analysis of the extinction time of the epidemic	101
5.1	Diagram of the model	113
5.2	Numerical verification of the predictive power of the stationary basic reproduction number	116
5.3	Numerical study of the delay induced by growing vector populations	118
5.4	Numerical verification of the expression for the non-stationary basic reproduction number	120
5.5	Numerical verification of the timescale approximation	122
5.6	Comparison between the original model and the reductions	123
6.1	Schematic representation of the model	131
6.2	Vector dynamics produced by the model compared to field-data	132
6.3	Posterior and prior distributions of the model parameters for ALSD	136
6.4	Posterior and prior distributions of the model parameters for OQDS	137
6.5	Best-fit model to the field data for ALSD and OQDS	139
6.6	Comparison of the model fit to the data for OQDS with different transmission rates	140
6.7	Global Sensitivity Analysis of the model	141
6.8	Epidemic control through vector management for ALSD in Mallorca and OQDS in Apulia	143
7.1	Climatic and transmission layers composing the epidemiological model	162
7.2	Average thermal-dependent maps for Pierce's disease (PD) development and recovery	165
7.3	Climate-driven risk maps for PD establishment in main viticulture regions worldwide	167
7.4	Temperature-driven model simulations for PD establishment from 1981 to 2019	168
7.5	Global shifts in PD risk index ( $r_j(t)$ ) from 2020 to 2050	170

7.6	PD risk in European vineyards for 2020 and 2050	172
8.1	Training presence records for modeling the distribution of <i>Philaenus spumarius</i> .	183
8.2	Changes in $X_{fPD}$ and <i>P. spumarius</i> climatic suitability under different climate projections	185
8.3	PD risk maps and associated risk velocities under different climate projections	187
8.4	Uncertainty in PD risk projections for climate change warming levels	188
8.5	Multi-scale spatial analysis of PD future risk in Europe	189
9.1	Difference in risk projections due to the climate data resolution in global viticulture areas	199
9.2	Changes in risk categories between due to the climate data resolution in global viticulture areas	200
9.3	Effect of microclimatic conditions of rivers and valleys on PD	202
9.4	Impact of high-resolution climate data on the risk of Pierce's disease for grapevines worldwide	203
10.1	Daily averaged time series data from the Bay of Palma and Cabrera stations	216
10.2	Deep learning model applied to assess the decadal $pH_T$ trend in the Bay of Palma	218
10.3	Deep learning model applied to fill the gaps in the $pH_T$ time series in Cabrera	219
10.4	Map of the stations' location in the Western Mediterranean Sea Basin	223
11.1	AI framework for seagrass monitoring from satellite imagery	233
11.2	Model performance in train and out-of-sample test datasets	236
11.3	Example of model predictions for a satellite image in the training and out-of-sample test set	237
11.4	Example of model predictions for a satellite image in the complete dataset	238
12.1	Macroecological patterns of global coral reef size, geometry and spacing	252
12.2	The fractal nature of global coral reefs	254
12.3	Box Counting Dimension of coral reef provinces surface	256
A.1	Dependence of metabolic rates on temperature	334
C.1	Factors influencing <i>Xf-Philaenus spumarius-Vitis vinifera</i> pathosystem	348
C.2	Experimental setup	349
C.3	Relationship between MGDD and temperature	351
C.4	Model validation with PD presence/absence data	355
C.5	ROC curve illustrating the model validation procedure	355
C.6	Fit of an SIR model to ALSD data	356

C.7	Average climatic suitability for <i>Philaenus spumarius</i> in Europe	359
C.8	Determination of <i>MGDD</i> and <i>CDD</i> trends and future extrapolations	360
C.9	Interannual climatic variability extrapolations of <i>MGDD</i> and <i>CDD</i>	361
C.10	Comparison of Arrhenius-based vs beta function to define <i>MGDD</i>	364
C.11	Difference on risk index using <i>MGDD</i> defined with an Arrhenius or beta function	365
C.12	Comparison of risk areas and low <i>CDD</i> areas	366
C.13	<i>MGDD</i> and <i>CDD</i> oscillations in different wine-growing regions	367
C.14	Climate suitability for $X_{f_{PD}}$ and <i>P. spumarius</i> under present and future climate conditions.	380
C.15	$\mathcal{F}(MGDD)$ , $\mathcal{G}(CDD)$ and suitability of $X_{f_{PD}}$ under different climate projections	381
C.16	Future risk of PD epidemics under different climate projections	381
C.17	Risk and risk velocity shifts as function of the climate projections	382
C.18	Validation of <i>MGDD</i> computation from daily temperature data	386
C.19	Validation of <i>CDD</i> computation from daily temperature data	387
C.20	Comparison of annual cumulative <i>MGDD</i> and <i>CDD</i> using hourly and daily temperature data	388
C.21	Difference in the increase rate of projected risk for different resolutions of the climate data used	389
C.22	Comparison of risk indices for different resolutions of the climate data used	390
C.23	Presence locations of <i>Vitis vinifera</i> obtained from GBIF	391
D.1	Seasonally adjusted fit to reconstructed and measured pH data	394
D.2	Seasonally adjusted fit to reconstructed and measured temperature data	395
D.3	Total alkalinity in the Bay of Palma	396
D.4	Scheme for the Bidirectional-LSTM Neural Network	397
E.1	Coverage of the training and out-of-sample test datasets in the Balearic Islands	400
E.2	Distribution of response values of different habitat classes	408
E.3	Confusion matrices for predicted classes and their area	413
E.4	Wasserstein distance between Train and Test datasets	414
E.5	Model performance as function of depth	416
F.1	Geographic location of the studied coral reefs	420
F.2	Computation of the fractal dimensions of each coral province using the box-counting method	421

## List of Tables

3.1	Description of the SIRP model parameters . . . . .	68
4.1	Variables and parameters of the individual-based SIRP model . . . . .	91
6.1	Estimated parameters of the model for ALSD in Mallorca . . . . .	138
6.2	Estimated parameters of the model for OQDS in Apulia . . . . .	138
7.1	Validation of model predictions . . . . .	166
7.2	Extrapolated shifts in risk areas for Pierce's disease in Europe in 2050 . . .	170
8.1	EURO-CORDEX GCM-RCM combinations considered . . . . .	180
9.1	Changes in Pierce's disease risk zones in different viticulture regions . . . .	198
9.2	Comparison of PD risk at known grapevine locations . . . . .	204
10.1	Optimal parameters used for the different RNN architectures . . . . .	226
10.2	Statistical comparison between different RNN architectures . . . . .	227
11.1	Final model performance . . . . .	239
C.1	Inoculation tests on different grapevine varieties . . . . .	345
C.2	PD risk areas in Europe with an homogeneous spatial vector distribution . .	368
C.3	PD risk areas in the US . . . . .	369
C.4	Potential distribution of PD in other world winegrowing regions . . . . .	371
C.5	Extrapolated PD risk areas in the US . . . . .	372

C.6	Extrapolated PD risk areas in Europe in 2050 with a homogeneous vector spatial distribution . . . . .	374
C.7	PD risk areas in Europe with a heterogeneous vector spatial distribution . . . . .	376
C.8	Surface of European vineyards at risk . . . . .	377
C.9	Extrapolated surface of European vineyards at risk in 2050 . . . . .	379
C.10	Risk velocity statistics for each climate projection . . . . .	382
C.11	Risk analysis in Europe under different climate projections . . . . .	383
C.12	Risk analysis for different European countries under different climate projections . . . . .	384
E.1	Metadata of the satellite images used in the study. . . . .	400
E.2	Ecological categories in the ground truth dataset . . . . .	403
E.3	Model performance comparison . . . . .	409
E.4	Number of parameters of each model . . . . .	410
E.5	Performance of all models based on Linknet architecture in the training dataset . . . . .	411
E.6	Performance of all models based on Linknet architecture in the out-of-sample test dataset . . . . .	411
E.7	Model performance per habitat class . . . . .	412
E.8	Performance of the consensus method . . . . .	415
E.9	Performance metrics for the final model in the training dataset . . . . .	417
E.10	Performance metrics for the final model in the test dataset . . . . .	418
E.11	Average performance metrics of the final model . . . . .	418
F.1	Statistics of the studied coral reefs in each province . . . . .	422
F.2	Statistical comparison of different size distributions . . . . .	423
F.3	Results of the power-law fit of the coral reef size distribution . . . . .	424

## List of acronyms

<b>ACA</b> .....	Allen Coral Atlas
<b>AEMET</b> .....	Agencia Estatal de Meteorología
<b>AI</b> .....	Artificial Intelligence
<b>AIC</b> .....	Akaike Information Criterion
<b>ALSD</b> .....	Almond Leaf Scorch Disease
<b>ARIMA</b> .....	Autoregressive Integrated Moving Average
<b>AUC</b> .....	Area Under the Curve
<b>BAM</b> .....	Biotic-Abiotic-Movement
<b>BD-LSTM</b> .....	Bidirectional Long Short-Term Memory
<b>BFGS</b> .....	Broyden-Fletcher-Goldfarb-Shanno
<b>BIC</b> .....	Bayesian Information Criterion
<b>BII</b> .....	Biodiversity Intactness Index
<b>BMR</b> .....	Basal Metabolic Rate
<b>BOATS</b> .....	Balearic Ocean Acidification Time Series

**CAMELE** ..... Consensus for Automated Marine Ecosystem Labelling and Evaluation

**CCDF** ..... Complementary Cumulative Distribution Function

**CDD** ..... Cold Degree Days

**CHELSA** ..... Climatologies at high resolution for the earth's land surface areas

**CI** ..... Confidence Interval

**CLS** ..... Coffee Leaf Scorch

**CNN** ..... Convolutional Neural Network

**CORDEX** ..... Coordinated Regional Climate Downscaling Experiment

**CORINE** ..... Coordination of Information on the Environment

**CVC** ..... Citrus Variegated Chlorosis

**DL** ..... Deep Learning

**DNA** ..... Deoxyribonucleic Acid

**DO** ..... Dissolved Oxygen

**DTM** ..... Digital Terrain Model

**EMODnet** ..... European Marine Observation and Data Network

**EPPO** ..... European and Mediterranean Plant Protection Organization

**ERA** ..... ECMWF Reanalysis

**FFNN** ..... Feedforward Neural Network

**FN** ..... False Negative

**FP** ..... False Positive

**FPN** ..... Feature Pyramid Network

**FPR** ..... False Positive Rate

**GAM** ..... Generalized Additive Model

<b>GBIF</b>	Global Biodiversity Information Facility
<b>GCM</b>	Global Climate Model
<b>GDD</b>	Growing Degree Days
<b>GOA-ON</b>	Global Ocean Acidification Observing Network
<b>GRU</b>	Gated Recurrent Unit
<b>GSA</b>	Global Sensitivity Analysis
<b>IBM</b>	Individual-Based Model
<b>IPCC</b>	Intergovernmental Panel on Climate Change
<b>IUCN</b>	International Union for Conservation of Nature
<b>JABOWA</b>	JAnak-BOtkin-WAllis
<b>LPHME</b>	Standard List of Marine Habitats of Spain
<b>LSA</b>	Local Sensitivity Analysis
<b>LSTM</b>	Long Short-Term Memory
<b>MAE</b>	Mean Absolute Error
<b>MAPE</b>	Mean Absolute Percentage Error
<b>MGDD</b>	Modified Growing Degree Days
<b>ML</b>	Machine Learning
<b>MME</b>	Mass Mortality Event
<b>MOW</b>	Mediterranean Outflow Water
<b>MSE</b>	Mean Squared Error
<b>NGM</b>	Next Generation Matrix
<b>NIR</b>	Near Infrared
<b>NUTS</b>	Nomenclature of Territorial Units for Statistics / No U-Turn Sampler
<b>OA</b>	Ocean Acidification

**ODE** ..... Ordinary Differential Equation  
**OOS** ..... Out-of-Sample  
**OQDS** ..... Olive Quick Decline Syndrome  
**PACA** ..... Provence-Alpes-Côte d'Azur  
**PCR** ..... Polymerase Chain Reaction  
**PD** ..... Pierce's Disease  
**PDE** ..... Partial Differential Equation  
**PDO** ..... Protected Designation of Origin  
**Ps** ..... *Philaenus spumarius*  
**PSPnet** ..... Pyramid Scene Parsing Network  
**RCM** ..... Regional Climate Model  
**RGB** ..... Red Green Blue  
**RMSE** ..... Root Mean Squared Error  
**RNN** ..... Recurrent Neural Network  
**ROC** ..... Receiver Operating Characteristic  
**RSS** ..... Residual Sum of Squares  
**SA** ..... Sensitivity Analysis  
**SDE** ..... Stochastic Differential Equation  
**SDM** ..... Species Distribution Model  
**SEIR** ..... Susceptible-Exposed-Infectious-Recovered  
**SI** ..... Susceptible-Infectious  
**SIP** ..... Susceptible-Infectious-Parasite  
**SIR** ..... Susceptible-Infectious-Recovered  
**SIRP** ..... Susceptible-Infectious-Recovered-Parasite  
**SIRS** ..... Susceptible-Infectious-Recovered-Susceptible

<b>SIS</b>	.....	Susceptible-Infectious-Susceptible
<b>SOP</b>	.....	Standard Operation Procedure
<b>SR</b>	.....	Surface Reflectance
<b>SRNN</b>	.....	Simple Recurrent Neural Network
<b>TA</b>	.....	Total Alkalinity
<b>TN</b>	.....	True Negative
<b>TP</b>	.....	True Positive
<b>TPR</b>	.....	True Positive Rate
<b>UNET</b>	.....	U-Net
<b>UV</b>	.....	Ultraviolet
<b>Xf</b>	.....	Xylella fastidiosa



## List of publications

### Publications in this thesis

1. À. Giménez-Romero, A. Grau, I. E. Hendriks, and M. A. Matias, “Modelling parasite-produced marine diseases: The case of the mass mortality event of *Pinna nobilis*”, [Ecological Modelling](#) **459**, 109705 (2021)
2. S. Flecha, À. Giménez-Romero, J. Tintoré, F. F. Pérez, E. Alou-Font, M. A. Matías, and I. E. Hendriks, “pH trends and seasonal cycle in the coastal Balearic Sea reconstructed through machine learning”, [Scientific Reports](#) **12**, 12956 (2022)
3. À. Giménez-Romero, F. Vazquez, C. López, and M. A. Matías, “Spatial effects in parasite-induced marine diseases of immobile hosts”, [Royal Society Open Science](#) **9**, 212023 (2022)
4. À. Giménez-Romero, R. Flaquer-Galmés, and M. A. Matías, “Vector-borne diseases with nonstationary vector populations: The case of growing and decaying populations”, [Phys. Rev. E](#) **106**, 054402 (2022)
5. À. Giménez-Romero, J. Galván, M. Montesinos, J. Bauzà, M. Godefroid, A. Fereres, J. J. Ramasco, M. A. Matías, and E. Moralejo, “Global predictions for the risk of establishment of Pierce’s disease of grapevines”, [Communications Biology](#) **5**, 1389 (2022)
6. À. Giménez-Romero, E. Moralejo, and M. A. Matías, “A Compartmental Model for *Xylella fastidiosa* Diseases with Explicit Vector Seasonal Dynamics”, [Phytopathology®](#) **113**, 1686–1696 (2023)

7. À. Giménez-Romero, M. Iturbide, E. Moralejo, J. M. Gutiérrez, and M. A. Matías, “Global warming significantly increases the risk of Pierce’s disease epidemics in European vineyards”, [Scientific Reports 14, 9648 \(2024\)](#)
8. À. Giménez-Romero, E. Moralejo, and M. A. Matías, “High-resolution climate data reveals increased risk of Pierce’s Disease for grapevines worldwide”, [bioRxiv \(2024\)](#)
9. À. Giménez-Romero, M. A. Matías and C. M. Duarte, “Universal spatial properties of coral reefs” ([Accepted in Global Ecology and Biogeography](#))
10. À. Giménez-Romero, D. Ferchichi, P. Moreno-Spiegelberg, T. Sintés, and M. A. Matías, “Mapping the distribution of seagrass meadows from space with deep convolutional neural networks”, [bioRxiv \(2024\)](#)

## Other publications

11. C. Lago, À. Giménez-Romero, M. Morente, M. A. Matías, A. Moreno, and A. Fereres, “Degree-day-based model to predict egg hatching of *Philaenus spumarius* (Hemiptera: Aphrophoridae), the main vector of *Xylella fastidiosa* in Europe”, [Environmental Entomology 52, 350–359 \(2023\)](#)
12. E. Moralejo, À. Giménez-Romero, and M. A. Matías, “Linking intercontinental biogeographic events to decipher how European vineyards escaped Pierce’s disease”, [Proceedings of the Royal Society B: Biological Sciences 291, 20241130 \(2024\)](#)
13. E. Moralejo, J. A. García-Muñoz, S. Denman, and À. Giménez-Romero, “Leaf susceptibility of Macaronesian laurel forest species to *Phytophthora ramorum*”, [bioRxiv \(2023\)](#)
14. À. Giménez-Romero, M. A. Matías and C. M. Duarte, “A comprehensive dataset on global coral reefs size and geometry” (Submitted for publication)

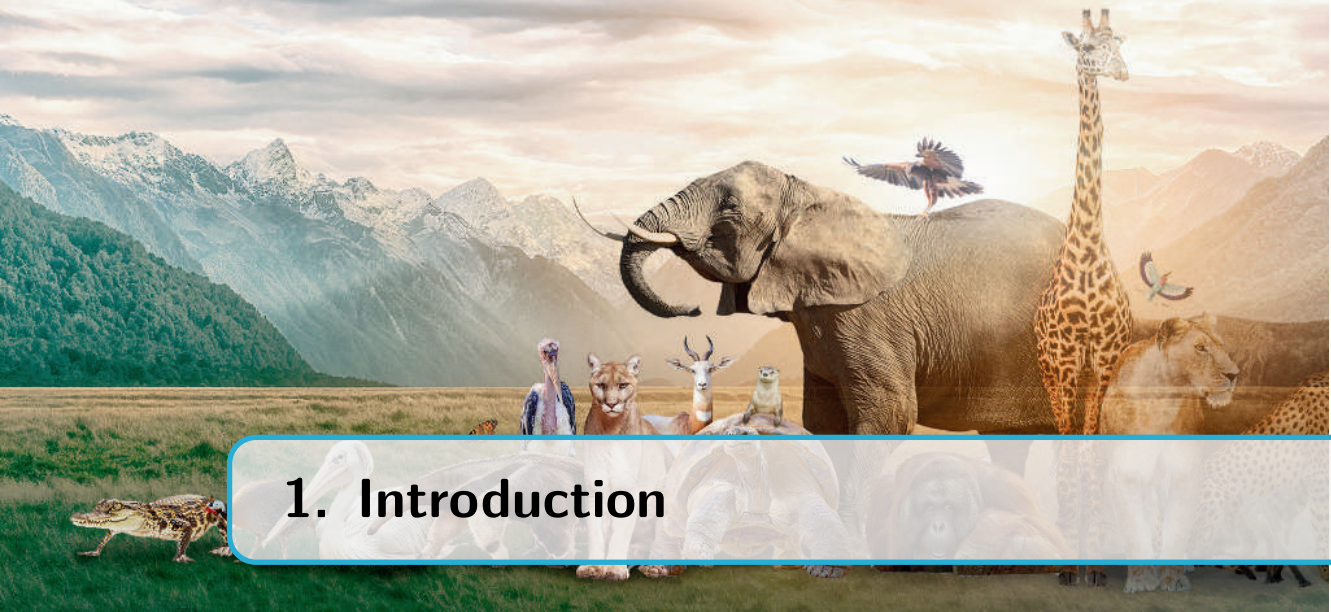
## Technical reports

1. D. P. Bebbler, S. J. Gurr, A. Karley, L. M. Lozada-Ellison, T. Beale and À. Giménez-Romero, "Interdisciplinary Analysis of Plant Health Threats to Arable and Horticultural Crops in Scotland", [Scotland's Centre of Expertise for Plant Health \(2024\)](#)

## Scientific dissemination

1. À. Giménez-Romero and M. A. Matías, "Descifrando el fondo marino desde el espacio con los ojos de la inteligencia artificial", [The Conversation \(2023\)](#)
2. À. Giménez-Romero, E. Moralejo and M. A. Matías, "Xylella fastidiosa y el cambio climático amenazan la viticultura europea: hemos calculado cuánto", [The Conversation \(2024\)](#)





# 1. Introduction

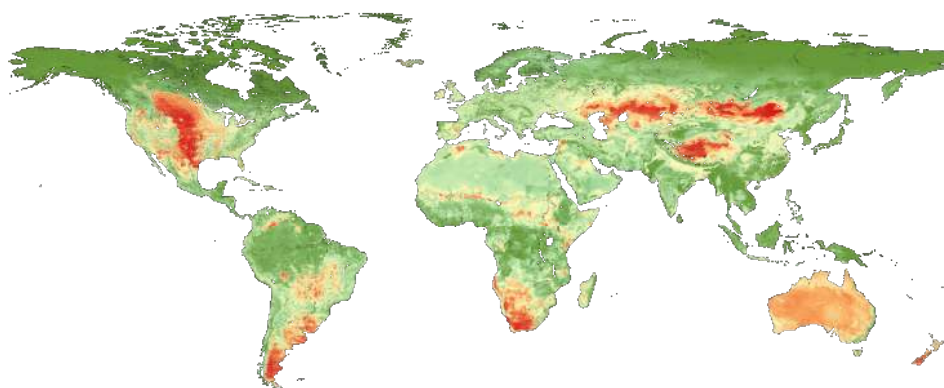
## 1.1 The global biodiversity crisis

Based on our current understanding, life on Earth appears to be unique in the Universe. Its existence is the result of a long evolutionary process that started around 3.5 billion years ago [1, 2], producing a wide variety of life forms. From the simplest unicellular organisms to the most complex multicellular forms, the diversity of life is so vast that it is estimated that there are around 9 million species of living beings on our planet [3]. This biodiversity is the result of the interaction among organisms and with their environment, which has led to the development of complex ecosystems. Ecosystems are the basic units of life on Earth, where organisms interact with each other and with the environment, forming a self-regulating system that is capable of maintaining life [4]. As Simon Levin eloquently describes [4]:

*“Ecosystems and the biosphere are complex adaptive systems, in which pattern emerges from, and feeds back to affect, the actions of adaptive individual agents, and in which cooperation and multicellularity can develop and provide the regulation of local environments, and indeed impose regularity at higher levels. The history of the biosphere is a history of coevolution between organisms and their environments, across multiple scales of space, time, and complexity.”*

Besides its intrinsic value, biodiversity is essential for the proper functioning of ecosystems [5], which in turn provide fundamental services to humankind. These ecosystem services include the production of oxygen, the regulation of climate, the provision of food and water, and many others [6]. Beyond these fundamental life-supporting benefits, biodiversity also plays critical roles in nutrient cycling and soil formation, processes integral to the sustainability of our

agricultural systems. The diversity of plant and animal life contributes to robust ecosystems that can withstand and recover from a variety of disasters, thereby ensuring ecological resilience. Moreover, biodiversity supports recreational and tourism industries, which are significant sources of income for many communities worldwide. The aesthetic and cultural values provided by diverse ecosystems foster mental and physical health, and contribute to the cultural heritage of communities, enriching our experience of the world. Without the services ecosystems provide, the stability of environments that support human life would be greatly diminished, leading to profound economic and social consequences.



**Figure 1.1: The Biodiversity Intactness Index (BII).** The BII is a measure of the relative intactness of biodiversity in different regions of the world. The index ranges from 0 to 100, with higher values indicating higher levels of biodiversity intactness. In this map, the darkest red color represents  $BII < 50$ , indicating that less than half of the original biodiversity remains in these regions. Data from the Natural History Museum data portal [7]

Unfortunately, the diversity of life on Earth is dramatically diminishing [8–10], posing a serious threat to the stability of ecosystems and the services they provide. Nowadays, wildlife extinction rates are estimated to be 100 to 1000 times higher than the natural background rate [11, 12], and up to 50% of higher taxonomic groups are already critically endangered [13]. In the last 50 years, the global population of vertebrates has declined by 69%, about 50% of the corals have disappeared due to different causes, and roughly 10 million hectares of forests are lost annually [14]. Sadly, one could continue listing more examples of biodiversity loss, but the point is clear: we are facing a global biodiversity crisis. Indeed, this has led some scientists to propose that we are entering the sixth mass extinction event in Earth’s history [15]. However, unlike the natural extinction events in Earth’s past, the current crisis is precipitated by

only one species: humans. As ecosystems falter and species vanish, the intricate web of life that sustains economies, food security, and our very existence is at risk. If left unchecked, the repercussions of this biodiversity crisis may lead to ecosystems so impaired that they no longer fulfill their roles, fundamentally altering the living conditions on our planet.

The main drivers of global biodiversity loss, as identified in the Millennium Ecosystem Assessment [16], encompass a range of direct impacts on the natural world, most of which are caused by human activities. Habitat change, exemplified by the rapid deforestation in the Amazon Rainforest, results in drastic reductions in biodiversity by stripping away the complex web of life supported by these environments [17]. Climate change brings about shifts in temperature and precipitation patterns that are markedly altering habitats. For instance, polar regions are shrinking, threatening ice-dependent species with extinction [18], while ocean acidification is disrupting marine ecosystems by impairing the ability of calcifying organisms to build their shells and skeletons [19]. Indeed, climate change may be a major threat to global biodiversity in the next 100 years [20–24], with predictions for species loss ranging from as low as 0% to as high as 54% [25]. Invasive species, such as the zebra mussel in North America, can disrupt ecosystems by outcompeting native species, leading to changes in the structure of the food webs, affecting the quality and quantity of primary production, and causing diseases, which probably have been underestimated as an ecological force [26]. Overexploitation of natural resources, as seen in the overfishing of the world's oceans, can lead to the collapse of entire ecosystems, as well as the loss of valuable food sources for human populations [27, 28]. Finally, wildlife emergent diseases threaten global biodiversity by potentially producing catastrophic declines in new and not adapted host populations. If the diseases become endemic, initial depopulation may be followed by chronic population depression, which could even lead to local extinction [29].

In addition, all these drivers are interconnected and can have cascading effects on ecosystems [30]. For example, climate change can expand the range of invasive species, which in turn can produce emerging diseases that affect native populations. Because native species have not co-evolved with the pathogens producing these diseases, they are more susceptible to them, leading to population declines and even extinctions [29]. Similarly, ocean acidification is expected to reduce the calcification rates of corals, making them more susceptible to diseases, bleaching events, and overexploitation, which can lead to the collapse of entire reef ecosystems [31]. These interactions among drivers of biodiversity loss can produce synergistic effects that amplify the impacts on ecosystems, making them more vulnerable to further disturbances and less resilient to recover from them.

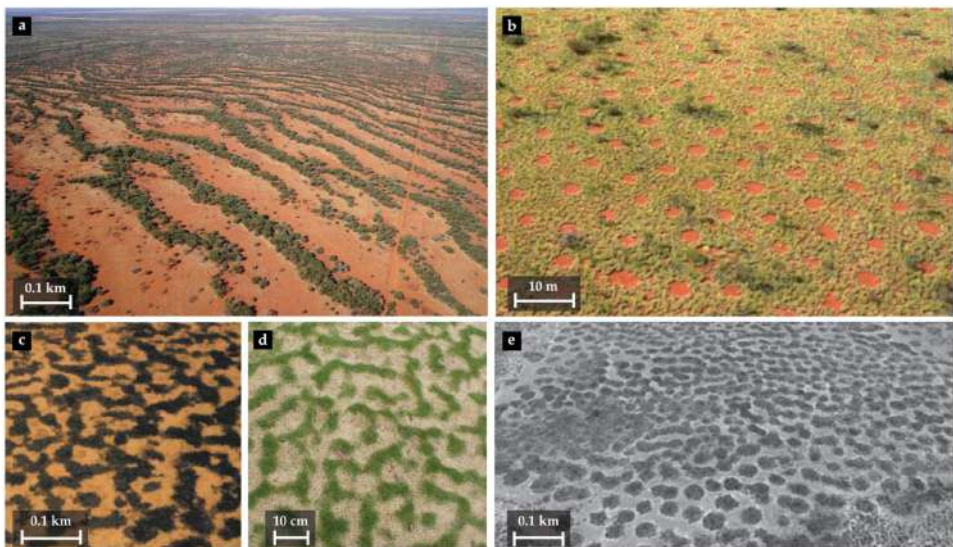
The global biodiversity crisis is a complex and multifaceted problem that requires urgent action to prevent further loss of species and ecosystems. The impacts of biodiversity loss are far-reaching, affecting not only the natural world but also human societies and economies. To address this crisis, we need to understand the underlying causes of biodiversity loss, predict the effects of environmental changes on ecosystems, and develop effective strategies for conservation. This requires a holistic approach that considers the interactions among species and with the environment, as well as the dynamics of ecosystems at different scales. In this context, complex systems science provides a powerful framework to address these challenges.

## 1.2 Complex systems in Ecology

Complex systems are composed of many interacting components whose collective features cannot be understood by simply studying the individual units in isolation [32]. The behavior of complex systems arises from the interactions among the components, which can lead to the formation of patterns, structures, and dynamics that are not present at the individual level. This emergent behavior is a hallmark of complex systems and illustrates how new properties and behaviors can arise from relatively simple interactions when viewed at a larger scale. As a common example, the flocking behavior of birds or schooling in fish is an emergent property of the interactions among the individuals, where each follows simple rules to maintain cohesion with its neighbors, leading to the formation of complex patterns at the flock level [33].

Complex systems are also characterized by other features that are ubiquitous in Ecology, such as non-linear dynamics, feedback loops, self-organization, a lack of central control, pattern formation, unpredictability linked to the presence of chaos, and the existence of multiple temporal, spatial, or organizational scales [32]. Non-linear dynamics are extremely common in ecological systems. An illustrative example are predator-prey interactions, where changes in the population of prey can lead to non-proportional changes in the population of predators, potentially causing oscillations rather than steady states [34]. Self-organization is observed through the spontaneous formation of spatiotemporal patterns that arise without any central control, often mediated by scale-dependent feedback loops, like the formation of vegetation patterns in arid ecosystems [35] (see Fig. 1.2). Contrary to this spontaneous order, the presence of chaos in ecological systems can lead to unpredictable dynamics related to the presence of chaos, as seen in the population fluctuations of some species or the change of gene frequencies in time [36, 37]. Overall, ecological processes operate across multiple spatial, temporal, and organizational scales, from local population interactions to global biogeochemical cycles, in which each level influences the

others. Indeed, it is not surprising that the principles of complex systems are that present in ecological systems, as they are composed of many interacting units whose dynamics are shaped by the interactions among them and with the environment. For example, ecosystems are composed of many species that interact with each other and with the physical environment, forming intricate food webs, nutrient cycles, and energy flows. And the same principles apply to lower organizational levels such as populations, communities, and metapopulations, where the interactions among individuals, species, and habitats give rise to patterns and dynamics that are not always intuitive.



**Figure 1.2: Vegetation patterns in arid ecosystems.** These patterns are an example of self-organized structures that emerge from the interactions among plants and with the environment. These patterns optimize resource use and enhance ecosystem stability. Adapted from [38]

Complex systems approaches have become increasingly important in Ecology, not only by providing theoretical frameworks to understand open questions in the field but also by posing new questions and challenges that were not previously considered, thus expanding the discipline in new directions [39]. A major contribution of complex systems to Ecology is the theory of self-organized systems, which provides a robust explanation for how complex ecological patterns can emerge from simple local interactions without the need for a centralized control mechanism. This principle has been instrumental in deciphering the formation of highly ordered structures such as the regular spacing in vegetation patterns, which optimize resource use in challenging environments [40].

Similarly, termite mounds are examples of natural engineering, where collective behaviors lead to the creation of sophisticated microhabitats that regulate temperature and humidity essential for colony survival [41]. In marine environments, mussel beds exhibit emergent properties such as enhanced stability and ecosystem engineering through their aggregation behavior, significantly affecting sediment dynamics [42].

Scaling relationships are another key concept in complex systems that have profound implications for understanding ecological patterns and processes. Scaling laws describe how the properties of biological or ecological systems change in a statistically regular way as a function of a parameter, temporal or spatial scale, organizational level, or all of them. These laws, when approached from a complex systems' perspective, reveal underlying principles that govern biological or ecological diversity across different scales, providing a mathematical framework to unify patterns that span across the tree of life. A well-known example is that of the allometric scaling laws in Biology, where metabolic rates, growth patterns, and lifespan are shown to scale with the size of organisms predictably (e.g., Kleiber's Law,  $BMR \propto M^{3/4}$ , where  $BMR$  is the Basal Metabolic Rate and  $M$  is the mass of the organism) [43]. For a long time, these scaling laws were considered biological curiosities until a mathematical framework was developed to explain these patterns, showing that this impressive regularity arises from the fractal-like structure of the organisms' vascular system and the invariant size of capillaries under the assumption of minimal energy expenditure [44]. Similar approaches can be followed to explain scaling relationships in ecological systems, including population density, resource distribution, and ecosystem productivity [45]. Now these scaling laws are recognized as fundamental principles that underlie the structure and function of organisms. Indeed, there is an entire field of Ecology, known as Macroecology, devoted to the study of these empirical patterns and the mechanistic processes that generate them [46].

Is it clear that one of the key tools used in the study of complex systems in Ecology is mathematical and computational modeling. Models are simplified representations of reality that capture the essential features of a system, allowing us to understand its dynamics, predict its behavior, and test hypotheses about its functioning. Models can be used to explore the consequences of different scenarios, design experiments, and inform management decisions. In Ecology, models have been used to study a wide range of topics, from population dynamics and community interactions to ecosystem processes and global biogeochemical cycles.

## 1.3 Mathematical and computational models

Mathematical and computational models have become pivotal in modern ecology, offering profound insights into the dynamics of complex ecological systems. These models serve as essential tools for simulating and understanding how ecosystems function, interact, and respond to various environmental pressures [47]. At their core, these models translate biological and ecological processes into mathematical descriptions, which can be manipulated and studied under controlled, repeatable conditions. This approach allows ecologists to test hypotheses about ecological interactions and processes in a virtual environment, where real-world experimentation would be either impractical or impossible. In addition, models provide a means to explore the consequences of different management strategies, predict the outcomes of environmental changes, and assess the impacts of human activities on ecosystems. Overall, mathematical and computational models are indispensable tools for advancing ecological research and have been recognized as some of the most powerful approaches to guide empirical work and provide a framework for synthesis, analysis, development of conservation plans, and policymaking [48–50].

One of the simplest and most widely used types of models in ecology are deterministic continuous-time models, such as ordinary differential equations (ODEs), which describe how the state of a system changes over time based on known biological or ecological processes. ODEs usually represent *mean field* descriptions of the system under study, assuming that the population or community is well-mixed and that the interactions among individuals are homogeneous. Under these assumptions, the system can be described by *rate equations* that capture the dynamics of the system in terms of rates of change of the state variables. These models have been used to study a wide range of ecological processes, from population growth and competition to predator-prey dynamics and disease transmission, providing valuable insights into the mechanisms that drive the dynamics of ecological systems [49].

Obviously, continuous-time models are not the only way to model the dynamics of ecological systems. Discrete-time models, such as difference equations, are also widely used in Ecology. Perhaps the logistic map, which describes the growth of a population over discrete generations, is by far the most famous example [36]. Indeed, this model led to the discovery of *chaotic dynamics* in simple ecological systems, in which small changes in the initial conditions can lead to large differences in the long-term behavior of the system. This discovery revolutionized the field of Ecology, highlighting the importance of non-linear dynamics in population and community dynamics. Of course, this chaotic behavior is not exclusive to discrete-time models, as it can also be observed in continuous-time models, such as the Lotka-Volterra predator-prey equations,

which exhibit complex dynamics, such as limit cycles, bifurcations, and chaos, under certain parameter values [49].

However, deterministic models may not fully capture the inherent variability and uncertainty present in ecological systems, which is not always due to chaotic dynamics. Stochastic models can provide a more realistic representation of the randomness and variability that are present in some ecological systems. For example, stochastic differential equations (SDEs) extend ODEs by including random fluctuations in the system that represent the effects of environmental variability, demographic stochasticity, or genetic drift. Similarly, discrete stochastic models simulate the dynamics of ecological systems by modeling the interactions among individuals as discrete events that occur with a certain probability. These models have been used to study the effects of environmental noise on population dynamics, the persistence of species in fluctuating environments, and the spread of diseases [49].

The dynamics of ecological systems are not only determined by their internal processes, but also by the spatial structure of the environment in which they are embedded. The natural extension of continuous-time models to incorporate spatial dynamics is the use of partial differential equations (PDEs), which describe how the state of a system changes over time and space [49]. These models have been instrumental in exploring dispersal mechanisms and predicting changes in population density across space. PDEs are also used to model the spread of invasive species and disease transmission through susceptible populations, providing insights into the impacts of spatial structure on ecological invasions [51]. In addition, nonlinear PDEs lead to pattern formation, the spontaneous emergence of spatial structure in systems with suitable local interaction and spatial coupling terms, like in the case of the Turing mechanism [52].

Environmental or demographic noise can also be incorporated into these models by means of stochastic partial differential equations (SPDEs), which provide a more realistic representation of the variability and uncertainty present in spatially structured ecological systems. Of course, spatial models can also be discrete, such as cellular automata, which divide the environment into discrete cells that can change state over time based on a set of rules. These models have been used to study the spread of forest fires, the dynamics of vegetation patterns, and the formation of spatially structured populations, providing insights into the effects of local interactions on the dynamics of ecological systems. Overall, spatial models are particularly valuable in conservation planning, helping to determine the most effective configurations of habitats to preserve species diversity and ecological function.

Both these continuous and discrete models can be used to study the dynamics of ecological systems, but they are often limited by their assumptions of ho-

mogeneity in interactions and spatial structure. Network models and individual-based models (IBMs) provide frameworks for studying more complex and heterogeneous interactions, both in a spatially extended and non-spatial context. Network models describe the interactions between the entities under study as a network of nodes and edges, where nodes represent individuals or populations and edges represent the interactions between them. These models have been used to study the structure and dynamics of various ecological networks, such as food webs and pollination networks, and provide a rich framework to tackle problems in which interactions among many species are central, such as coevolution in species-rich communities [53]. Of course, network models can be both deterministic and stochastic.

IBMs simulate the behaviors of individual organisms or entities based on a set of simple rules. These models are exceptionally useful to study heterogeneous populations, in which individuals differ in their characteristics or behaviors, and spatially explicit interactions. IBMs help in understanding how individual heterogeneity contributes to group dynamics and ecosystem-level patterns, providing insights that are often obscured in more traditional approaches. Initially, IBMs were used to model forest succession, exemplified by the JABOWA model, which described tree community dynamics in response to shading and growth patterns. Their applications then expanded to animal populations, including fish recruitment, bird nesting colonies, and predator-prey interactions, allowing for the examination of complex interspecific relationships and behaviors [54].

Data-driven techniques such as Machine Learning (ML) are also increasingly being used in ecological modeling. These methods, are particularly useful for handling complex, high-dimensional data and capturing non-linear relationships. In short, ML algorithms learn how to predict an output variable from input data by iteratively adjusting their parameters to minimize the error between the predicted and observed values. Machine learning models have been widely used to predict species distributions and habitat suitability from environmental variables, providing valuable insights into the impacts of climate change on biodiversity [55]. More recently, deep learning models, a subfield of ML based on artificial neural networks, have shown promise in ecological applications, such as identifying species in camera trap images [56] or classifying the behavior of animals captured in images and videos [55]. Specifically, the conjunction of deep learning models with remote sensing data is revolutionizing the field of landscape ecology, allowing to monitor land cover changes [57], map canopy height [58], track animal migrations [59], or detect plant diseases [60]. As computational capabilities continue to grow and as datasets become more comprehensive and accessible, the potential for mathematical and computational models in ecology expands. This progression promises not only to deepen our understanding of

ecological systems but also to improve our ability to predict and mitigate the impacts of human activities and environmental changes on the natural world.

This is only a brief overview of the wide range of mathematical and computational models that have been developed and applied in Ecology. The diversity of models reflects the complexity of ecological systems and the diverse questions that ecologists seek to answer. By combining theoretical insights with empirical data, mathematical and computational models provide a powerful framework for understanding the dynamics of ecological systems, predicting their responses to environmental changes, and informing management decisions. In the following chapters, we will explore how these models can be used to address some of the most pressing challenges in Ecology and Conservation Biology, from the spread of infectious diseases to the decline of important ecosystems under global change.

## 1.4 Thesis structure

This thesis is devoted to developing theoretical and data-driven methods to address timely problems in Ecology and Conservation Biology from the perspective of Complex Systems. We tackle a variety of current challenges related to biodiversity loss caused by climate change and emergent diseases, including expanding vector-borne plant diseases, Mass Mortality Events (MMEs) produced by marine diseases, ocean acidification, and the decline of important coastal ecosystems like coral reefs or seagrass meadows. To address these challenges, we developed a series of theoretical and data-driven models that integrate ecological theory, data analysis, mathematical models, computational simulations, and artificial intelligence. The thesis is structured as follows. In [Chapter 2](#), we present the theoretical framework and scientific background that underpin the research presented in this thesis. We review basic concepts of infectious disease modeling and disease biogeography, explain the main methods used in the data-driven approaches, and present the context of the ecological systems that will be studied in the following chapters. The main results of the thesis are divided into four parts, each formed by several chapters. In [Part I](#), we focus on the development of theoretical models to understand the dynamics of marine diseases and their impacts on marine ecosystems, exemplified by the Mass Mortality Event of *Pinna nobilis*. Similarly, in [Part II](#), we develop theoretical models to study the vector-borne plant diseases in which the vector population is non-stationary and non-periodic. Our framework is then applied to diseases caused by the bacterium *Xylella fastidiosa*. In [Part III](#), we develop and apply a climate-driven epidemic model to predict the risk of establishment of *Xylella fastidiosa* diseases globally, both in current and future climates. In [Part IV](#), we develop data-driven models to address complex ecological problems, such as

---

the decline of coral reefs and seagrass meadows or ocean acidification. Finally, in [Part V](#), we summarize our contributions, present the main conclusions of the thesis, and discuss the implications of our results for the conservation of biodiversity and the management of ecosystems under global change.





## 2. Scientific background

### 2.1 Infectious disease modeling

Infectious diseases are caused by pathogenic microorganisms, such as bacteria, viruses, parasites, or fungi, when they invade other organisms, which are called *hosts* in this context. These pathogens proliferate inside the host and can disrupt the normal functioning of the colonized organism by damaging tissues, altering physiological processes, or triggering immune responses that contribute to illness. Infectious diseases can affect a wide variety of organisms, from plants and animals to humans. The causal agents of infectious diseases can be transmitted from one host to another through different mechanisms, such as direct contact between an infected and a susceptible individual, through the air, water, or food, or through vectors. Vectors are organisms that transmit pathogens from one host to another, such as mosquitoes, that usually are not affected by the disease themselves.

Modeling the spread of infectious diseases has a long history, with the first mathematical model coming from the hand of Bernoulli in 1760 [61], who developed a model to understand the spread of smallpox. However, it was not until the early 20th century that we find the foundation of the modern theory of infectious disease dynamics, with the seminal work of Ronald Ross and Hilda P. Hudson [62–64], Kermack and McKendrick [65], and later George Macdonald [66]. Since then, a wide range of models have been developed to understand the dynamics of infectious diseases and to predict their spread. In general, these models can be classified into two main categories: compartmental models and individual-based models.

#### 2.1.1 Compartmental models

Compartmental models are based on the assumption that the population can be divided into different compartments or categories, each representing a different

state of the disease (e.g., susceptible to infection, infected). Individuals move from one compartment to another following some dynamical rules. Under the assumption of a well-mixed and sufficiently large population, one can consider that every pair of individuals has equal probability of coming into contact with one another and that fluctuations in the number of individuals in each compartment can be neglected. This is known as the mean field approximation. Under these assumptions, the dynamics of the disease can be described by a set of differential equations that govern the transitions between compartments.

Several compartmental models have been developed to describe the dynamics of infectious diseases. The most famous is the SIR model, which is a simple particular case of the general mathematical framework formulated by Kermack and McKendrick [65]. Because the original Kermack-McKendrick model is quite general and complex, I will explain the SIR model following a more modern and intuitive approach. In this model, the host population can be divided into three compartments: susceptible individuals (S), infected individuals (I), and recovered individuals (R). Other models have been developed to account for more complex scenarios, such as the SIS model (where recovered individuals become susceptible again), the SEIR model (where an exposed compartment is added to account for the incubation period of the disease), or the SIRS model (where recovered individuals can become susceptible again after some time). In addition, compartmental models have been developed to account for vector-borne diseases, including compartments for different states of the vector population.

### The SIR model

The SIR model divides the population into three compartments: susceptible individuals (S), infected individuals (I), and recovered individuals (R). Individuals come into contact with each other at a given rate  $a$  and, in the case of an  $S-I$  contact, the susceptible individual becomes infected with a probability  $b$ . We assume that the incubation period is short enough to be negligible; that is, a susceptible who contracts the disease is infective right away. Infected individuals recover at a constant rate  $\gamma$ , and once recovered, are assumed to be immune to the disease and cannot be reinfected. We finally assume that there is no entry into or departure from the population, no births or natural deaths, so that the population remains constant.

According to the mean field approximation, the probability that an infected individual contacts a susceptible one is given by  $a \cdot S/N$ . Thus, the average number of contacts between infected and susceptible individuals is given by  $a \cdot S/N \cdot I = a \cdot SI/N$ . Finally, as the probability that a susceptible individual becomes infected after a contact is  $b$ , the average number of new infections will be given by the product of this probability and the average number of contacts between infected and susceptible individuals,  $b \cdot a \cdot SI/N = \beta SI/N$ . These

considerations lead to the following description of the model given by a system of differential equations,

$$\begin{aligned}\dot{S} &= -\beta SI/N \\ \dot{I} &= \beta SI/N - \gamma I \\ \dot{R} &= \gamma I .\end{aligned}\tag{2.1}$$

Despite the apparent simplicity of the model, the non-linear nature of the differential equations makes it difficult to find analytical solutions. However, important quantitative insights can still be obtained by performing an analytical study of the model. We will now derive some important results from the SIR model, most of which can be already found in reference [49].

### Conservation of the total population

First, let's start with the simplest insight: the assumption of population conservation is inherently included in the model. We can prove this statement by just adding the three differential equations in Eq. (2.1),

$$\dot{S} + \dot{I} + \dot{R} = 0 \implies S + I + R = \text{const.} = N .\tag{2.2}$$

Physicists usually call this kind of quantity a **conserved quantity** or a **conservation law**. In this case the conservation law has a particular meaning: the total number of individuals in the population remains constant, which was already assumed in the model. In general, conservation laws can arise from other symmetries of the system and are not always so obvious. In any case, these conserved quantities allow us to reduce the number of independent variables in the system, which can be very useful to simplify the analysis of the model. In this case, the conservation law allows us to reduce the number of independent variables from three to two, so that we can study the dynamics of the system in a two-dimensional phase space given by the variables  $S$  and  $I$ , as  $R$  can be obtained from the relation  $R = N - S - I$ .

### Threshold behavior: the basic reproduction number

Now consider the starting point of an epidemic, such as our old friend the COVID-19 pandemic. At the beginning of the epidemic the number of infected individuals is not zero,  $I(0) > 0$ , while the number of recovered individuals is indeed zero,  $R(0) = 0$ . Thus, the number of susceptible individuals at the beginning of the epidemic is  $S(0) = N - I(0)$ . If an epidemic is to develop, the number of infected individuals must increase at the beginning of the epidemic, but what are the conditions for this to happen?

Just by considering the differential equation for  $I$  in Eq. (2.1) with our initial

conditions we find the following expression,

$$\frac{dI}{dt} \Big|_{t=0} = I(0) (\beta S(0)/N - \gamma) \implies \begin{cases} \frac{dI}{dt} \Big|_{t=0} < 0 \iff S(0) < \frac{\gamma N}{\beta} \\ \frac{dI}{dt} \Big|_{t=0} > 0 \iff S(0) > \frac{\gamma N}{\beta} \end{cases} . \quad (2.3)$$

The condition Eq. (2.3) defines a threshold for the development of an epidemic in the *SIR* model. This is, there is a critical number of susceptible individuals below which the epidemic will not develop,

$$S_c = \frac{\gamma N}{\beta} \equiv \rho . \quad (2.4)$$

The critical parameter  $\rho$  is sometimes called the *relative removal rate* and its reciprocal,  $\sigma = \beta/(\gamma N)$ , the *infection's contact rate*. This threshold behavior can be further simplified by considering the so-called *basic reproduction number*,  $R_0$ , defined as,

$$R_0 = \frac{S(0)}{\rho} = S(0)\sigma = \frac{\beta S(0)}{\gamma N} \simeq \frac{\beta}{\gamma} , \quad (2.5)$$

in which we have considered that  $S(0) \simeq N$  at the beginning of the epidemic.

The basic reproduction number is a very important quantity in epidemiology, and it can be proved that it measures the average number of secondary infections given by a primary infection in a *fully susceptible population*. To see it we just need to recall the definition of the rates and “read” the expression.  $\beta S(0)/N$  is the rate of new infections produced by a primary infection (as we are at time  $t = 0$ ). In other words,  $\beta S(0)/N$  is the number of secondary infections produced by a primary infection *per unit time*. Finally, we observe that, as  $\gamma$  is the rate of recovery of infected individuals,  $1/\gamma$  is the average time an individual remains infected. The product of the number of secondary infections produced by a primary infection per unit time and the average time an infected individual remains infected gives the total number of secondary infections produced by a primary infection in the whole susceptible population. Voilà! This is nothing but the basic reproduction number,  $R_0$ .

This threshold behavior agrees with our intuition given the definition of the basic reproduction number: if a primary infection produces more than 1 secondary infection an epidemic will develop, while if it does not reach to infect at least one individual it will die out. Furthermore, note that because of our mean field approach, the number of secondary infections produced by a primary one refers to the *average* number of secondary infections.

### Initial phase approximation

An approximated closed analytical result can be obtained by considering the initial phase of the epidemic, when the number of infected individuals is small compared to the number of susceptible individuals,  $I \ll S$ . In this case, the number of susceptible individuals is approximately the total population,  $S \approx N$ , as the number of recovered individuals is negligible,  $R \approx 0$ . Under these conditions, the dynamics of the infected population can be described by the following differential equation,

$$\frac{dI}{dt} = \beta S/N I - \gamma I \approx \beta I - \gamma I = (\beta - \gamma)I, \quad (2.6)$$

which has the solution,

$$I(t) = I(0)e^{(\beta - \gamma)t} = I(0)e^{\gamma(R_0 - 1)t}. \quad (2.7)$$

This result shows that the number of infected individuals grows exponentially at the beginning of the epidemic, with a growth rate related to  $R_0$ . This approximation is very useful to understand the initial phase of the epidemic, when the number of infected individuals is small, but it is only valid and useful for short times after the beginning of the epidemic.

In Fig. 2.1(a) we show the numerical solution of the SIR model for a given set of parameters. The fraction of susceptible, infected, and recovered population is shown as a function of time. The black dashed line represents the initial phase approximation of the model, which is in good agreement with the numerical solution of the model at the beginning of the epidemic. In the inset, we can see that the initial growth of the epidemic is indeed exponential (a straight line in log-linear scale), as predicted by the model.

### Maximum number of infected individuals

Another important analytical result that can be obtained from this model is the maximum number of infected individuals, which gives an idea of how severe the epidemic will be. To find this maximum we just need to find the maximum of  $I(t)$ , which is given by the condition  $\dot{I} = 0$ . From the differential equation for  $I$  in Eq. (2.1) we find the following relation,

$$\frac{dI}{dt} = 0 \implies \beta S/N - \gamma = 0 \implies S = \frac{\gamma N}{\beta} = \rho. \quad (2.8)$$

So the maximum of infected individuals, given the development of a proper epidemic, will take place when  $S(t) = \rho$ . But we still don't know the maximum number of infected individuals. To do so, we first need to go through a smart

mathematical trick. Dividing the differential equations for  $S$  and  $I$  (considering  $I \neq 0$ ) we obtain,

$$\frac{dI}{dS} = -1 + \frac{\gamma}{\beta S} = -1 + \rho/S \quad (I \neq 0).$$

Integrating this relation,

$$\int_{I(0)}^I dI = \int_{S(0)}^S dS \{-1 + \rho/S\} \implies I - I(0) = S(0) - S + \rho \ln\left(\frac{S}{S(0)}\right),$$

we obtain the phase plane trajectories for  $S$  and  $I$  given by,

$$I + S = I_0 + S_0 + \rho \ln\left(\frac{S}{S(0)}\right) = N + \rho \ln\left(\frac{S}{S(0)}\right), \quad (2.9)$$

where we have considered  $I(0) + S(0) = N$  given that  $R(0) = 0$ .

Note that Eq. (2.9) can be rewritten as  $I + S - \rho \ln(S) = N - \rho \ln(S_0) = \mathcal{C}$ , so that the right-hand side of the equation is a constant. This means that we have just found another conservation law for the system! This one is not so obvious as the previous one, but it is still a very useful result. In essence, we now can describe the dynamics of the system with only one independent variable,  $S$ , as  $I$  can be obtained from the relation in Eq. (2.9).

Finally, we find an expression for the maximum of infected individuals by substituting the condition previously found (the maximum occurs when  $S(t) = \rho$ ) in Eq. (2.9),

$$I(t) = N + \rho \ln\left(\frac{S(t)}{S(0)}\right) - S(t) \implies I_{max} = N + \rho \left[ \ln\left(\frac{\rho}{S_0}\right) - 1 \right] = N + \frac{S(0)}{R_0} \left[ \ln\left(\frac{1}{R_0}\right) - 1 \right]. \quad (2.10)$$

In Fig. 2.1(b) we can see how the analytical result compares with the numerical solution of the model for different values of  $R_0$ .

### Final number of susceptible individuals

Similarly, we can find the final number of susceptible individuals at the end of the epidemic, for sure a very important quantity. By dividing the differential equations for  $S$  and  $R$ , we obtain,

$$\frac{dS}{dR} = -\frac{\beta}{\gamma} S = -\frac{S}{\rho} \implies \int_{S(0)}^S \frac{dS}{S} = -\frac{1}{\rho} \int_{R(0)=0}^R dR \implies \ln\left(\frac{S}{S(0)}\right) = -\frac{R}{\rho}.$$

As  $I(\infty) = 0$ , necessarily  $R(\infty) = N - S(\infty)$  so that we get the following expression for the final number of susceptible individuals,

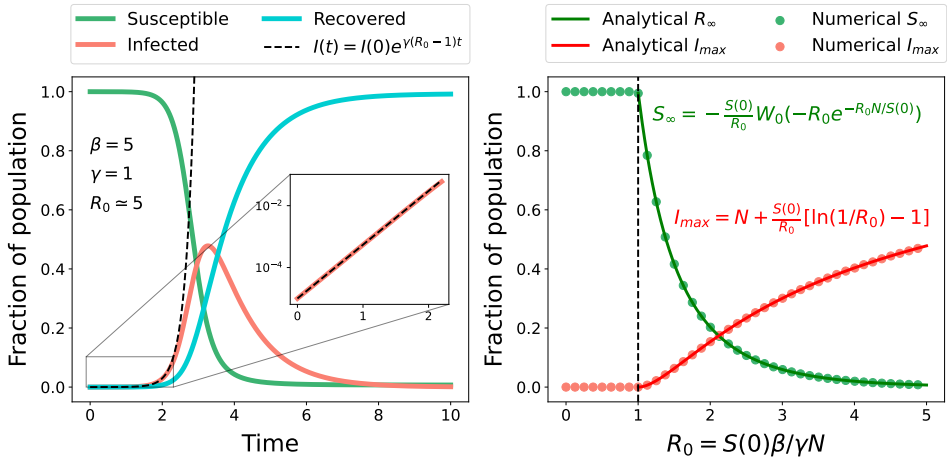
$$S(\infty) - \rho \ln(S(\infty)) = N - \rho \ln(S(0)). \quad (2.11)$$

$S(\infty)$  is nothing but the positive root of the transcendental equation Eq. (2.11). In fact, this transcendental equation can be solved by means of the Lambert's  $W$  function [67],

$$S(\infty) = -\rho \cdot W_0 \left[ -\frac{S(0)}{\rho} e^{-N/\rho} \right] = -\frac{S(0)}{R_0} \cdot W_0 \left[ -R_0 e^{-R_0 N/S(0)} \right]. \quad (2.12)$$

Of course, given that  $N = S + I + R$  and  $I(\infty) = 0$ , we can obtain the final number of recovered individuals as  $R(\infty) = N - S(\infty)$ . In Fig. 2.1(b) we show how the analytical result compares with the numerical solution of the model for different values of  $R_0$ .

**R** Note that all the analytical results that we have obtained can be expressed as a function of the initial conditions of the system (e.g.,  $S(0)$ ) and the basic reproduction number,  $R_0$ . This is a very important result, as it allows us to understand and predict the dynamics of the whole system by just knowing, in essence, the basic reproduction number.



**Figure 2.1: Numerical and analytical analysis of the SIR model.** (a) Numerical solution of the SIR model for a given set of parameters. The fraction of susceptible (green), infected (red), and recovered (blue) individuals is shown as a function of time. The black dashed line represents the initial phase approximation of the model. (b) Comparison between the analytical results obtained from the SIR model and the numerical solution of the model for different values of  $R_0$ . The maximum number of infected individuals,  $I_{max}$ , and the final number of susceptible individuals,  $S_\infty$ , are shown as a function of the basic reproduction number,  $R_0$ . The analytical results are in good agreement with the numerical solution of the model.

### Incidence functions

In the SIR model we have considered a frequency-dependent function, where the rate of new infections is proportional to the product of the number of susceptible and infected individuals divided by the total population  $\beta SI/N$ . This incidence function, known as *standard incidence*, considers that the rate of infection is bounded as the population grows, assuming that the probability of a contact between a susceptible and an infected individual is less likely as the population grows. However, other incidence functions can be considered, such as the *mass action incidence*, where the rate of new infections is proportional to the product of the number of susceptible and infected individuals,  $\beta SI$ . This incidence function can be used to describe scenarios in which all-to-all contacts are possible irrespective of the population size. The choice of the incidence function can have important consequences on the dynamics of the disease, as it can affect the threshold for the development of an epidemic, the maximum number of infected individuals, and the final number of susceptible individuals.

So far we have shown how, under minimal assumptions, one can build an epidemic model such as the SIR model. In addition, we have shown how to obtain some important analytical results from the model, such as the threshold for the development of an epidemic, i.e., the *basic reproduction number*; the maximum number of infected individuals, and the final number of susceptible individuals. However, the SIR model is one of the simplest compartmental models, and more complex models have been developed to account for more realistic scenarios. For such models, analytical results are usually not available, and numerical simulations are required to understand the dynamics of the diseases they describe. Fortunately, there are a couple of formal approaches that can be used to derive the basic reproduction number of basically any compartmental model: linear stability analysis and the next generation matrix approach.

#### 2.1.2 Linear stability analysis

Linear stability analysis is a powerful tool that can be used to study the dynamics of a **dynamical system** around its **fixed points**. The basic idea is to linearize the system of differential equations around these equilibrium points and study the stability of the system by analyzing the **eigenvalues** of the **Jacobian** of the system. A *dynamical system* is a system that evolves over time, such as the SIR model. The *fixed points* of the system are the points where it does not change over time, i.e., the points where the derivatives of the variables of the model are zero. The *Jacobian* of the system is a matrix that contains the first-order partial derivatives of the variables of the system. The *eigenvalues* of the Jacobian matrix are the roots of the characteristic polynomial of the matrix, obtained through the equation  $\det(J - \lambda I) = 0$ , and they determine the stability

of the system. If the real part of all the eigenvalues is negative, the system is stable, while if the real part of at least one eigenvalue is positive, the system is unstable. Let's see how this works in practice.

For the SIR model, the fixed points are given by the condition  $\dot{S} = \dot{I} = \dot{R} = 0$ , which implies  $I = 0$  with any value of  $S$  and  $R$ . Thus, the fixed points of the system are given by  $(S^*, I^*, R^*) = (S, 0, N - S)$ . These fixed points are also called *disease-free* state of the system, as the number of infected individuals is zero. The Jacobian matrix of the system is given by,

$$J = \begin{pmatrix} \partial\dot{S}/\partial S & \partial\dot{S}/\partial I & \partial\dot{S}/\partial R \\ \partial\dot{I}/\partial S & \partial\dot{I}/\partial I & \partial\dot{I}/\partial R \\ \partial\dot{R}/\partial S & \partial\dot{R}/\partial I & \partial\dot{R}/\partial R \end{pmatrix} = \begin{pmatrix} -\beta I/N & -\beta S/N & 0 \\ \beta I/N & \beta S/N - \gamma & 0 \\ 0 & \gamma & 0 \end{pmatrix}, \quad (2.13)$$

Now we substitute the fixed point in the Jacobian matrix and we obtain,

$$J = \begin{pmatrix} 0 & -\beta S^*/N & 0 \\ 0 & \beta S^*/N - \gamma & 0 \\ 0 & \gamma & 0 \end{pmatrix}. \quad (2.14)$$

The eigenvalues of the Jacobian matrix are the roots of the characteristic polynomial of the matrix, which is given by  $\det(J - \lambda I) = 0$ . In this case the characteristic polynomial is given by,

$$\det(J - \lambda I) = \begin{vmatrix} -\lambda & -\beta S^*/N & 0 \\ 0 & \beta S^*/N - \gamma - \lambda & 0 \\ 0 & \gamma & -\lambda \end{vmatrix} = \lambda^2(\beta S^*/N - \gamma - \lambda). \quad (2.15)$$

The roots of the characteristic polynomial are given by  $\lambda_1 = \lambda_2 = 0$  and  $\lambda_3 = \beta S^*/N - \gamma$ . As we have previously stated, the stability of the system is determined by the real part of these eigenvalues. Of course, the eigenvalues  $\lambda_1$  and  $\lambda_2$  are zero, so they do not provide any information about the stability of the system. The eigenvalue  $\lambda_3$  is negative if  $\beta S^*/N < \gamma$ , thus meaning that the fixed point  $(S^*, I^*, R^*)$  is stable if this condition is satisfied. On the other hand, the fixed point is unstable if  $\beta S^*/N > \gamma$ . Indeed, this condition defines the basic reproduction number,  $R_0$ , as we can rewrite the condition for the stability of the fixed point as  $R_0 = \beta S^*/\gamma N$  greater or smaller than 1.

Finally, note that the Jacobian matrix also provides us with another important insight about the system: the presence of conserved quantities. In this case, the sum of the first and third rows of the Jacobian matrix is zero, which indicates that the system has two conserved quantities, which we have already found: the total number of individuals in the population and [Eq. \(2.9\)](#).

**R** The stability of the fixed points of a dynamical system can be determined by analyzing the eigenvalues of the Jacobian matrix of the system. If the real part of all the eigenvalues is negative, the fixed point is stable, while if the real part of at least one eigenvalue is positive, the fixed point is unstable. In epidemic models, the basic reproduction number,  $R_0$ , determines the stability of the disease-free states of the system, which are the fixed points where the number of infected individuals is zero. If  $R_0 < 1$ , the fixed points are stable, while if  $R_0 > 1$ , the fixed points are unstable. The basic reproduction number is a threshold for the development of an epidemic: if  $R_0 < 1$ , the epidemic will die out, while if  $R_0 > 1$ , the epidemic will propagate. In addition, the Jacobian matrix of the system provides us with important insights about the system, such as the presence of conserved quantities.

### 2.1.3 The next generation matrix method

We have previously shown that the basic reproduction number,  $R_0$  of a compartmental model can be obtained by analyzing the stability of the fixed points corresponding to the disease-free state of the model. Specifically, the basic reproduction number,  $R_0$ , is related to the largest non-zero eigenvalue of a fixed point,  $\Lambda$ , such that  $R_0 > 1$  if  $\Lambda > 0$ . However, the Jacobian matrix of a compartmental model can be quite large and complicated, so the derivation of the basic reproduction number following this method can be cumbersome. The next generation matrix method is an ingenious method that can be used to derive the basic reproduction number of any compartmental model directly, without the need to analyze the stability of the fixed points. We will make use of this method in [Parts I](#) and [II](#) of this thesis.

**Method** — **The next generation matrix method.** Steps of the next generation matrix:

1. Compute the Jacobian matrix of the infected/infected compartments of the model evaluated at the disease-free state,  $J^*$ .
2. Decompose the Jacobian matrix in the form  $J^* = T + \Sigma$ , where  $T$  is the transmission part and  $\Sigma$  the transition part.
3. Compute the inverse of the transition part,  $\Sigma^{-1}$ .
4. Compute the next generation matrix,  $K = -T\Sigma^{-1}$ .
5. Solve the characteristic equation  $\det(K - \lambda I) = 0$  to obtain the eigenvalues of the next generation matrix.
6. Obtain the basic reproduction number,  $R_0$ , as the largest eigenvalue of the next generation matrix.

In the NGM method,  $R_0$  is identified as the dominant eigenvalue of a suitably defined linear operator (a linear matrix in a suitable basis). This operator is obtained by decomposing the Jacobian of the infected/infected compartments evaluated at the disease-free state,  $J^*$ , in the form  $J = T + \Sigma$ , where  $T$  is the

*transmission part*, that describes the production of new infections, and  $\Sigma$  the *transition part*, that describes changes of state. Then, it can be proved [68] that the *basic reproduction number*  $R_0$  is given by the spectral radius (i.e., the largest eigenvalue) of the next generation matrix,  $K$ , defined as,

$$R_0 = \rho(K) \quad \text{with} \quad K = -T\Sigma^{-1} \quad (2.16)$$

To appreciate the power of the NGM method, let's see how it can be applied to a slightly more complex model than the SIR model. We will now consider a compartmental model for vector-borne plant diseases.

### A compartmental model for vector-borne plant diseases

The simplest way to do this is to consider a compartmental model with five compartments: susceptible individuals ( $S$ ), infected individuals ( $I$ ), recovered individuals ( $R$ ), susceptible vectors ( $S_v$ ), and infected vectors ( $I_v$ ). For most vector-borne plant diseases, we can assume that the only mechanism of disease spread is the direct transmission from infected vectors to susceptible plants at a given rate  $\alpha$ . Infected individuals recover at a rate  $\gamma$ . The total number of plants is given by  $N$ , and the total number of vectors is given by  $N_v$ . Vectors are assumed to be born at a constant rate  $\delta$  proportional to the total number of vectors and die at a constant rate  $\mu$ . With these assumptions, the dynamics of the disease can be described by the following system of differential equations [69],

$$\begin{aligned} \dot{S} &= -\beta SI_v/N_v \\ \dot{I} &= \beta SI_v/N_v - \gamma I \\ \dot{R} &= \gamma I \\ \dot{S}_v &= \delta N_v - \alpha S_v I/N - \mu S_v \\ \dot{I}_v &= \alpha S_v I/N - \mu I_v \end{aligned} \quad (2.17)$$

We can easily check that the plant population is conserved by adding the differential equations for  $S$ ,  $I$ , and  $R$ ,

$$\dot{S} + \dot{I} + \dot{R} = 0 \implies S + I + R = N . \quad (2.18)$$

Similarly, to check if the vector population is conserved we add the differential equations for  $S_v$  and  $I_v$ ,

$$\dot{S}_v + \dot{I}_v = \delta N_v - \mu S_v - \mu I_v = 0 \implies N_v(\delta - \mu) = 0 . \quad (2.19)$$

This equation implies that the vector population is conserved only if  $\delta = \mu$ , which is a reasonable assumption given that the birth and death rates of vectors are usually similar. By now, this is the case we will consider in the following analysis, but we will see in [Part II](#) that things can get really complicated if we relax this assumption.

After this initial check, we can proceed to obtain an expression for the basic reproduction number,  $R_0$ . To compute it we can follow the next generation matrix approach. We first need to identify the disease-free state of the system, which is given by  $(S^*, I^*, R^*, S_v^*, I_v^*) = (S, 0, N - S, N_v, 0)$ . Then we need to write down the Jacobian matrix of the subsystem corresponding to the infected/infecting compartments,  $I$  and  $I_v$ , in this case and evaluate it at the disease-free state. This matrix is given by,

$$J = \begin{pmatrix} \partial I / \partial I & \partial I / \partial I_v \\ \partial I_v / \partial I & \partial I_v / \partial I_v \end{pmatrix} = \begin{pmatrix} -\gamma & \beta S / N_v \\ \alpha N_v / N & -\mu \end{pmatrix} = \begin{pmatrix} -\gamma & \beta S(0) / N_v \\ \alpha N_v / N & 0 \end{pmatrix}, \quad (2.20)$$

where we have substituted the fixed point  $(S^*, I^*, R^*, S_v^*, I_v^*)$  in the last step considering, without loss of generality, that  $(S^*, I^*, R^*, S_v^*, I_v^*) = (S(0), 0, N - S(0), N_v, 0)$

Next we decompose the matrix in the transmission and transition parts,  $J = T + \Sigma$ , and compute the inverse of the transition part,  $\Sigma^{-1}$ ,

$$T = \begin{pmatrix} 0 & \beta S(0) / N_v \\ 0 & 0 \end{pmatrix} \quad \text{and} \quad \Sigma = \begin{pmatrix} -\gamma & 0 \\ \alpha N_v / N & -\mu \end{pmatrix} \implies \Sigma^{-1} = \frac{-1}{\gamma \mu} \begin{pmatrix} \mu & 0 \\ \alpha N_v / N & \gamma \end{pmatrix}. \quad (2.21)$$

We finally obtain the next generation matrix,  $K$ , as the product of the transmission and transition parts,

$$K = -T\Sigma^{-1} = \frac{1}{\gamma} \begin{pmatrix} \frac{\beta \alpha S(0)}{\gamma \mu} \frac{S(0)}{N} & \frac{\beta S(0)}{\mu} \frac{S(0)}{N} \\ 0 & 0 \end{pmatrix}. \quad (2.22)$$

The basic reproduction number,  $R_0$ , is given by the spectral radius of the next generation matrix,  $K$ , which is the largest eigenvalue of the matrix. In this case, the matrix is a  $2 \times 2$  matrix, so the eigenvalues can be easily obtained. The characteristic polynomial of the matrix is given by,

$$\det(K - \lambda I) = \begin{vmatrix} \frac{\beta \alpha S(0)}{\gamma \mu} \frac{S(0)}{N} - \lambda & \frac{\beta S(0)}{\mu} \frac{S(0)}{N} \\ 0 & -\lambda \end{vmatrix} = \lambda \left( \lambda - \frac{\beta \alpha S(0)}{\gamma \mu} \frac{S(0)}{N} \right), \quad (2.23)$$

so the eigenvalues are  $\lambda_1 = 0$  and  $\lambda_2 = \beta \alpha S(0) / (\gamma \mu N)$ . The basic reproduction number,  $R_0$ , is given by the largest eigenvalue of the matrix, so

$$R_0 = \frac{\beta \alpha S(0)}{\gamma \mu N}. \quad (2.24)$$

### 2.1.4 Individual-based models

Individual-based models describe the dynamics of infectious diseases by considering the individuals of the population as discrete entities that can interact

with each other. These models are particularly useful when the population is not well-mixed, the interactions between individuals are heterogeneous, the individuals cannot be considered identical or when a spatial setting is considered. In individual-based models, the dynamics of the disease are described by stochastic processes that govern the interactions between individuals. These models are usually more complex than compartmental models, so analytical results are not usually available, but they can provide more realistic insights into the dynamics of infectious diseases.

In individual-based models the population is represented by a set of individuals, each of which can be in different **states**. More formally, the **system** at time  $t$  is defined by the set of states of the individuals in the population at that time. The states of the individuals can change over time according to **events** that occur in the system, which consequently change the state of the system. These events can be triggered by the interactions between individuals, the environment, internal or other external factors. The events are characterized by their **rates**, which determine the probability of occurrence per unit time. Mathematically, the dynamics of the system are described by the **master equations** [70], which are a set of differential equations that describe the probability of the system being in a given state at a given time. The master equations are usually difficult to solve analytically, so numerical methods are used to simulate the dynamics of the system.

### Gillespie's method

There are several ways to simulate the dynamics of individual-based models, but one of the most common methods is Gillespie's algorithm [71], which is indeed an exact and unbiased algorithm to simulate stochastic processes. Gillespie's algorithm is based on the idea that the system will remain in a given state for a random amount of time until an event occurs, which will then change the state of the system. Thus, the change of state of the system can be determined by two independent random processes: the time to the next event and the particular event that occurs.

Consider that we are simulating a system with  $M$  possible events, each of which has a rate  $w_i$  of occurrence. The time of the next event is given by

$$t_{\text{next}} = \frac{-\ln(r_1)}{W}, \quad (2.25)$$

where  $W = \sum_i^M w_i$  is the total rate of any event happening and  $r_1 = \hat{U}(0, 1)$  is a distributed random number in the interval  $[0, 1]$ .

Once the time of the next event has been computed, we choose the event to implement according to their conditional probabilities  $p_i = w_i/W$ . Basically, we draw another uniformly distributed random number  $r_2$  and compare it to these

conditional probabilities, finding the smallest  $i$  satisfying  $\sum_i^M p_i > v$ . The system is then updated according to the chosen event and the time is updated to  $t_{\text{next}}$ . This process is repeated until a termination condition is met. In [Algorithm 1](#) we show a pseudocode implementation of Gillespie's algorithm.

---

### Algorithm 1 Gillespie Algorithm

---

- 1: Initialize time  $t \leftarrow 0$ , and set the initial state of the system.
  - 2: **while** termination condition is not met **do**
  - 3:     Calculate the propensity functions  $w_i(t)$  for all possible events  $i$ .
  - 4:     Calculate the total propensity  $W(t) = \sum_i^M w_i(t)$ .
  - 5:     Generate two random numbers  $r_1$  and  $r_2$  uniformly distributed in  $[0, 1]$ .
  - 6:     Calculate the time until the next reaction event  $\tau = \frac{-\ln(r_1)}{W(t)}$ .
  - 7:     Select the next reaction  $j$  according to the probabilities  $P(j) = \frac{w_j(t)}{W(t)}$  compared to  $r_2$ .
  - 8:     Update the system state according to the chosen reaction.
  - 9:     Update time  $t \leftarrow t + \tau$ .
  - 10: **end while**
- 

Gillespie's method provides a statistically exact way to solve the master equations underlying the IBM and gives a physical sense to the simulated time, as it comes directly from the event rates. However, one of its major drawbacks is the computational expense of the algorithm. Nevertheless, some tricks can be implemented to overcome these limitations and still provide an exact solution to the stochastic process. Two tricks have been used in our implementation of Gillespie's algorithm of [Chapter 4](#): first, instead of comparing  $p_i$  with the random number  $r_2$  we use  $w_i$  over  $v \cdot W$ , as this way only a product is computed (instead of  $M$ ) [70]; second, we implemented the sorting direct method [72]. The method consists in dynamically sorting the events and trying to apply the most frequent ones first, which can save a considerable amount of computational time. In [Algorithm 2](#), we show a pseudocode implementation of the core of the numerical method.

We have seen how mathematical models can be used to study the dynamics of infectious diseases, which can provide important insights into the spread of diseases *at a local level*. However, another important aspect of infectious diseases is their spread and establishment *at a global level*. In the following section, we will discuss the most common methods used to study the potential distribution of diseases and their limitations, particularly in the context of climate change.

**Algorithm 2** Sorting direct method

---

```

1: U = rand() * W           ▷ Random number multiplied by total rate
2: for  $i = 1, \dots, M$  do           ▷ Iterate over events
3:   if  $U < \text{sum}(\text{rates}[\text{orders}[1:i]])$  then
4:     events[i](*args)           ▷ Execute corresponding event
5:     if  $i \neq 1$  then           ▷ Sort reactions array
6:       aux_e = copy(events)
7:       aux_o = copy(orders)
8:
9:       events[i-1] = aux_e[i]
10:      events[i] = aux_e[i-1]
11:
12:      orders[i-1] = aux_o[i]
13:      orders[i] = aux_o[i-1]
14:     end if
15:     break
16:   end if
17: end for

```

---

## 2.2 Disease biogeography

Biogeography is the study of the distribution of species and ecosystems across the Earth's surface, and it is a fundamental field in ecology. Disease biogeography is a subfield of biogeography that focuses on the geographical distribution of infectious diseases and the factors that influence their establishment, such as climate, land use, human activities, and the presence of vectors or reservoirs. Understanding how these factors influence the spread of diseases at a geographical scale is essential to predict the risk of emergent diseases developing in new areas. This, in turn, can help public health authorities design effective strategies to prevent and control the spread of diseases.

The study of disease biogeography is particularly challenging due to the complexity of the interactions between the environment, the pathogens, the vectors, and the hosts involved in the disease transmission cycle. The potential distribution of diseases has been mostly studied using correlative species distribution models, which link presence and absence data of diseases to environmental variables based on correlative methods. Despite their wide use, these models have several limitations, both conceptual and practical. In this section, we will briefly review the main concepts of biogeography and explain the basic principles of species distribution models. Finally, we will discuss the limitations of using this framework to study the potential distribution of infectious diseases.

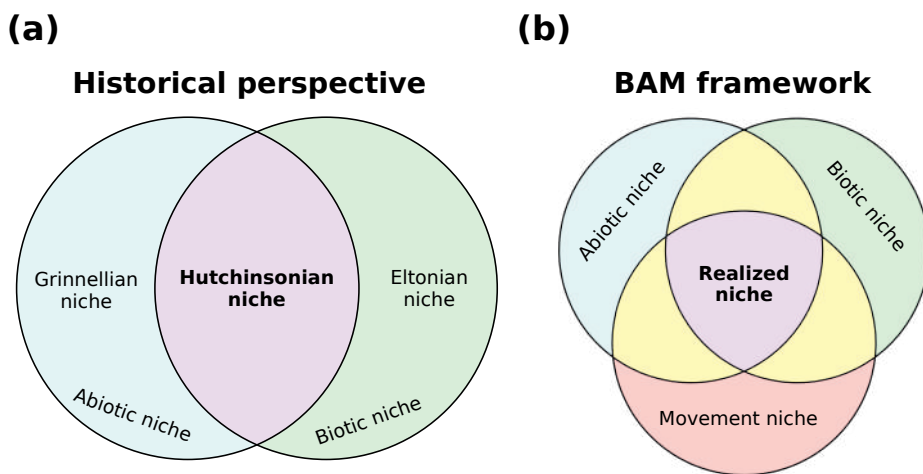
### 2.2.1 Ecological niches

The concept of “ecological niche” is central to the study of biogeography, but it has evolved significantly since its inception. The term “niche” was first introduced by Grinnell in 1917 to describe the set of environmental conditions within a species’ range that dictate its distribution [73]. This is, the niche of a species is the set of environmental conditions, such as temperature, humidity, and precipitation, that are suitable for the species to survive and reproduce. This early interpretation, known as the Grinnellian niche, emphasized the abiotic factors essential for the survival and distribution of species, but neglected the biotic interactions that also shape species’ distributions. Indeed, some years later, Charles S. Elton redefined the niche concept to encompass the roles species play within an ecosystem, including their interactions with other species. This Eltonian definition shifted the focus towards the biotic aspects of ecological niches [74].

To integrate both views, Hutchinson introduced the concepts of fundamental and realized niches [75]. The fundamental niche describes a hypervolume of environmental conditions under which a species can survive without immigration (i.e., the Grinnellian niche), while the realized niche is the subset of the fundamental niche that the species can actually exploit due to biotic interactions (e.g., competition). This distinction between the fundamental and realized niches provides a more comprehensive understanding of the ecological niche of a species, taking into account both abiotic and biotic factors. For example, a species may have broad tolerances of abiotic conditions such as wide ranges of temperature and rainfall, thus providing a large fundamental niche, which translates in a wide potential distribution of the species. However, the interactions of this species with other ones may restrict it to only a subset of the abiotically suitable areas. The presence of competitors, predators, and pathogens (e.g., negative interactions) or the absence of key mutualistic species (e.g., positive interactions) can influence the survival of a species, ultimately shaping its realized niche. In mathematical terms, we can express the realized niche  $R$  of a given species as the intersection of the abiotic niche  $A$  and the biotic niche  $B$ ,  $R = A \cap B$  (Fig. 2.2 (a)).

However, the realized niche of a species can be further constrained by dispersal limitations, which can prevent the species from fully exploiting its fundamental niche. This idea was comprehensively developed by Soberón and Peterson in 2005, who introduced the BAM (Biotic-Abiotic-Movement) framework to clarify the niche concept by emphasizing the species’ dispersal capabilities alongside biotic and abiotic factors [76]. This framework suggests that the full ecological niche of a species includes not only the environmental conditions that are abiotically suitable and biotically favorable but also the areas that the

species can physically reach and colonize. For example, a species originating from a tropical region in the Americas may have suitable environmental conditions in other tropical regions around the world, such as Southeast Asia, and even may be lucky enough to find an absence of competitors or predators in there. Yet, it will not be able to establish itself there if it is unable to disperse. So, in mathematical terms, the realized niche  $R$  of a species can be expressed as the intersection of the abiotic niche  $A$ , the biotic niche  $B$ , and the dispersal (movement) niche  $M$ ,  $R = A \cap B \cap M$  (Fig. 2.2 (b)).



**Figure 2.2: Historical development of the ecological niche concept.** The Grinnellian niche emphasizes the abiotic factors that determine the distribution of species, while the Eltonian niche focuses on the biotic interactions that shape the distribution of species. The Hutchinsonian niche integrates both abiotic and biotic factors, distinguishing between the fundamental and realized niches of a species. The BAM (Biotic-Abiotic-Movement) framework of the ecological niche concept. The full ecological niche of a species includes the abiotic conditions that are suitable for the species (abiotic niche), the biotic interactions that favor the species (biotic niche), and the areas that the species can physically reach and colonize (dispersal niche). The realized niche is the intersection of these three components.

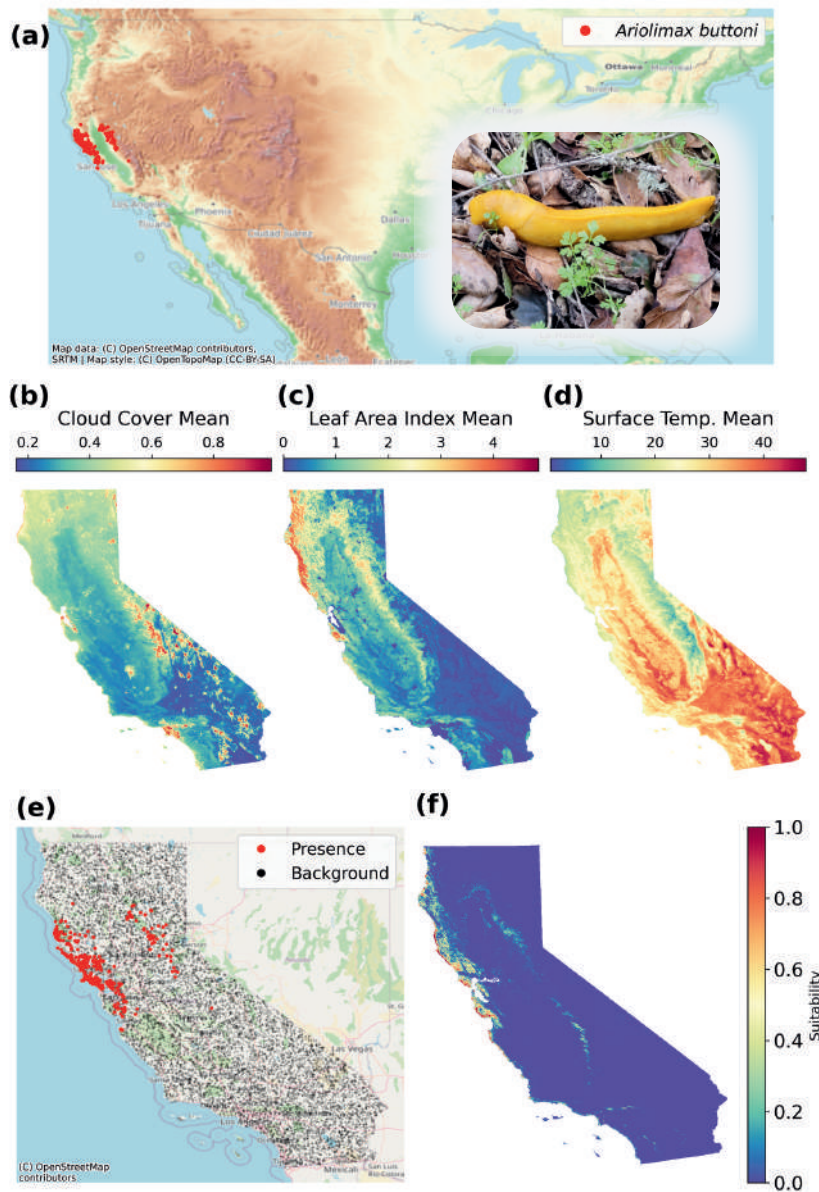
### 2.2.2 Species distribution models

Species distribution models (SDMs) are a powerful tool that can be used to predict the potential geographic distribution of species based on the environmental conditions that favor their establishment. SDMs are based on the concept of ecological niches, and they use environmental data to predict the potential dis-

tribution of species. In general, SDMs can be both mechanistic and correlative, depending on the underlying assumptions of the model. Mechanistic SDMs are based on the physiological or ecological requirements of the species, while correlative SDMs are based on statistical relationships between the presence and absence of the species in different locations and their environmental conditions [77]. Thus, correlative SDMs make use of statistical models that relate the presence and absence of a species to certain environmental variables. If this relationship renders successful, the model can be then used to predict the potential distribution of the species at different locations based on their environmental conditions. Because most of the SDMs that have been employed so far in the literature are based on correlative methods, we will refer to correlative SDMs as just SDMs in the following.

Many statistical methods can be used to build SDMs, such as logistic regression, generalized linear models, and machine learning algorithms, like random forests, support vector machines, or MaxEnt [78]. These methods have different strengths and weaknesses, and the choice of method depends on the characteristics of the data and the research question. For example, logistic regression is simple and interpretable, allowing researchers to understand the relationship between the species and the environmental variables, while machine learning algorithms are more flexible and can capture complex relationships but often lack the explanatory power of simpler models.

The overall process of building an SDM can be often divided into several steps [79]: First we need to collect data on the species of interest and the environmental variables that may influence its distribution. This data can be obtained from field surveys, remote sensing, or databases such as GBIF (Global Biodiversity Information Facility) [80] or WorldClim [81]. Once the data has been collected, it needs to be preprocessed to remove missing values, outliers, and redundant variables. In many cases, the data on the presence of the species is more abundant than the data on the absence of the species, or directly there is no absence data available. To address this issue, we can generate pseudo-absence data by randomly selecting points from areas where the species is not known to occur [82]. The data is then partitioned into a training set and a testing set. The training set is used to fit the model to the data, while the testing set is used to evaluate the performance of the model. We then train the model by fitting it to the data and estimating the parameters that best predict the presence and absence of the species based on the environmental variables. We then evaluate the model to assess its performance using metrics such as the area under the receiver operating characteristic curve (AUC) or the Kappa statistic. Finally, the model can be used to predict the potential distribution of the species at different locations based on their environmental conditions.



**Figure 2.3: Example of a species distribution model.** (a) Distribution of the banana slug *Ariolimax buttoni*, which is only present in the west of the United States. Inset: *Ariolimax (Ariolimax) buttoni* (Pilsbry & Vanatta, 1896) observed by [citysalamanders](#) (under [CC BY-NC 4.0](#)) (b-d) Environmental variables used to build the SDM. (e) Presence and modeled pseudo-absence data of the banana slug *Ariolimax buttoni* in the west of the United States. (f) Predicted potential distribution of the banana slug *Ariolimax buttoni* based on climatic variables.

To illustrate the process of building an SDM, we will consider the example of *Ariolimax buttoni*, a species of banana slug native to coastal California in the west of the United States (Fig. 2.3). This example is provided by the *elapid* Python library [83], which uses temperature, cloud cover, and leaf area index as environmental variables to predict the potential distribution of the banana slug using a MaxEnt model. The code is freely available at the documentation webpage of the library.

**Method** — **Species distribution models.** Steps of building a species distribution model:

1. Collect data on the species of interest and the environmental variables that may influence its distribution.
2. Preprocess the data to remove missing values, outliers, and redundant variables.
3. Generate pseudo-absence data to balance the presence and absence data.
4. Partition the data into a training set and a testing set.
5. Train the model and estimate the parameters that best predict the presence and absence of the species based on the environmental variables.
6. Evaluate the model using the testing set.
7. Apply the model to predict the potential distribution of the species at different locations based on their environmental conditions.

In Fig. 2.3(a) we show the distribution of the banana slug *Ariolimax buttoni* in the west of the United States. The species is only present in coastal California, where the environmental conditions are suitable for its survival. In Fig. 2.3(b-d) we show the annual mean of the environmental variables used to build the SDM<sup>1</sup>. The presence and modeled pseudo-absence data of the banana slug *Ariolimax buttoni* are shown in Fig. 2.3(e) and the predicted potential distribution of the banana slug based on the environmental variables is shown in Fig. 2.3(f).

SDMs have been widely used in the study of disease biogeography, as they can help to predict the potential distribution of diseases based on the environmental conditions that favor the establishment of the pathogens and/or vectors that transmit them [84–86]. Two main approaches have been used to apply SDMs to disease biogeography: the disease-based SDMs and the pathosystem-based SDMs. Disease-based SDMs assume that disease outbreaks are the final manifestation that results from the interaction between the environment, the pathogens, the vectors, and the hosts. Thus, the disease cases are considered as occurrences of the “species” to be modeled [87, 88]. This more pragmatic approach can be useful when the data on the pathosystem components

---

<sup>1</sup>The standard deviation of these variables are also used in the SDM but is not shown in the figure for simplicity

is scarce or disease transmission is not well understood [89]. On the other hand, pathosystem-based SDMs consider the pathogens, vectors, and hosts as separate species and model their potential distributions individually [90, 91]. This more systematic approach can provide better knowledge on the environmental factors that influence the distribution of the pathosystem components, but it requires more data and a better understanding of the interactions between the components.

However, despite the wide use of SDMs in the literature, these models have several limitations that need to be urgently addressed, especially for the case of disease biogeography. This is indeed one of the main reasons that led us to develop a new method to study the potential distribution of infectious diseases, which we will present in **Part III** of this thesis.

### 2.2.3 Limitations of SDMs

One of the key assumptions of SDMs are that species are in equilibrium with their environments and that relevant environmental gradients have been adequately sampled. However, SDMs are often used in non-equilibrium scenarios, such as to predict the future potential distribution of species after emerging invasions or climate change, which involves that current species records are unrepresentative of new conditions. In this new context, new combinations of environmental factors and biotic interactions previously unseen might play a role in the distribution of the species, thus rendering the predictions of the SDMs unreliable. In addition, the successful establishment of a species in a new area is influenced by genetic variability, phenotypic plasticity, and evolutionary changes, which are not accounted for in SDMs [92].

Another limitation of SDMs is that the results are highly dependent on the covariates used to fit the model, so that their choice can influence the predictions. In fact, the choice of environmental variables can affect the performance of the model, and the inclusion of irrelevant variables can lead to overfitting. Of course, choices can be made based on biological considerations or statistical methods based on information criteria, such as the Akaike Information Criterion (AIC) or the Bayesian Information Criterion (BIC) [93], but the final choice is often subjective.

For the particular case of the application of SDMs in disease biogeography, the limitations of the models are even more pronounced. Infectious diseases are the result of complex interactions occurring at multiple spatial and temporal scales among pathogens, vectors, hosts, and the environment, i.e., a disease is an emergent feature of the pathosystem. Indeed, basically this same idea was previously discussed by A. Townsend Peterson in 2008 [94], despite complex systems' science being a relatively new field of study. He argued that

*“Because disease processes are dynamic, taking place on extremely diverse scales of space (microscopic to continental) and time (minutes to centuries), and are the products of interactions among species (pathogens, reservoirs, vectors, etc.), their ecological and distributional dynamics may differ from those of more normal species.”*

However, current applications of SDMs to study the potential distribution of infectious diseases neglect these complex interactions, which can lead to biased predictions of the potential distribution of diseases. The disease-based SDMs assume that the distribution of disease occurrences can be seen as the joint spatial distribution of suitable ecological conditions for all the biological species involved in the transmission cycle. Yet, the current distribution of the disease may not reflect the full range of suitable conditions for its transmission. On the other hand, the pathosystem-based SDMs consider that the potential distribution of the disease is the result of the intersection of the potential distributions of the pathosystem components based on arbitrary thresholds [90, 91, 95]. This is ignoring completely the fact that the interactions between the components of the pathosystem are often non-linear and that the distribution of the disease is not the simple intersection of the distributions of the components.

In [Chapter 7](#), we will present a new method to study the potential distribution of vector-borne plant diseases that overcomes these limitations. We will consider the interactions driving disease transmission using a mathematical epidemiological model while accounting for climatic factors limiting the distribution of the pathosystem components.

## 2.3 Data-driven methods

The last part of this thesis, [Part IV](#) is focused on different global ecological problems such as ocean acidification or the loss of important coastal ecosystems like seagrass meadows and coral reefs. The causes of these problems have been extensively studied, so our goal is more practical than theoretical. In this context, we aim to develop new methods to solve some challenges that are currently faced by ecologists.

The frameworks presented in these chapters are based on deep learning, a subfield of machine learning that uses artificial neural networks to model complex patterns in large datasets. In this section, we will provide a brief overview of machine learning, specifically focusing on deep learning. After that, we will introduce the two main types of deep learning models that we will use in this thesis: long-short-term memory neural networks, which are used in time series forecasting, and convolutional neural networks, employed in image segmentation problems.

### 2.3.1 Machine learning

Machine learning is a subfield of artificial intelligence that focuses on the development of statistical algorithms that use input data to achieve a desired task without being literally programmed to produce a particular outcome [96, 97]. The concept of machine learning was pioneered by Arthur Samuel in 1959, who demonstrated that machines could learn to play checkers and improve at the game through their own experiences rather than through human instruction [98]. This capability marks a significant shift from traditional programming methods where decisions and rules are explicitly defined by human programmers. The field of machine learning integrates principles from computer science, statistics, and information theory to create algorithms capable of generalizing behaviors from input data, allowing machines to perform tasks that typically require human intelligence, such as predicting consumer behavior, filtering spam in emails, chatting with other humans, or even driving autonomous vehicles.

Machine learning is generally divided into three main types: supervised, unsupervised, and reinforcement learning. **Supervised learning** algorithms need a dataset containing several instances of the input data from which the model is expected to predict an output (features) together with the correct output (labels). Initially, the model parameters are set randomly, and the model undergoes a training process where the parameters are iteratively adjusted to minimize the difference between the predicted output and the true output [99]. On the contrary, **unsupervised learning** is used when the dataset does not contain the labels for each instance of the input data. In this case, the algorithms are designed to find patterns and structures in the data on their own [100]. A clear example is that of clustering algorithms, which group instances of the input data forming different clusters based on their similarity in a high-dimensional space. Finally, **reinforcement learning** is a type of machine learning where an agent learns to make decisions by interacting with an environment. The agent receives rewards or penalties based on its actions and learns to optimize its decisions (i.e., maximizing the reward or minimizing the penalty) to achieve a specific goal [101]. Reinforcement learning is particularly useful in applications such as robotics and gaming, being the algorithm behind the success of the AlphaGo program that defeated the world champion in the game of Go [102].

Data is the cornerstone of machine learning, and the quality and quantity of the data used to train the model are crucial to its performance. Traditional machine learning algorithms require the data to be preprocessed and transformed into a suitable format before being fed into the model, which prevented the widespread use of machine learning in many applications in its early times. However, the surge of representation learning methods, such as deep learning, has revolutionized the field of machine learning by allowing the model to learn

the features directly from the raw data without the need for manual feature engineering [103]. Deep-learning models use multiple levels of representation obtained by composing simple non-linear transformations, which allow the model to learn very complex functions if enough compositions of such transformations are considered. For example, in classification tasks, higher layers of representation amplify aspects of the input that are important for discrimination and suppress irrelevant variations, thus allowing the model to learn the features that are most relevant for the task at hand. This is particularly useful in image recognition, speech recognition, and natural language processing, where the raw data is high dimensional and complex and the features that are relevant for the task are not known a priori.

### 2.3.2 Artificial neural networks

The core of deep learning models are artificial neural networks (ANNs), which are computational models inspired by the structure and function of the human brain. ANNs consist of interconnected nodes, called neurons, which indeed model the biological neurons of a brain. The neurons are connected to each other by weighted edges, which determine the strength of the connection between neurons and mimic the synapses in the brain. The input layer receives the input data, the hidden layers process the data, and the output layer produces the final output of the network. The information is passed through the network by adjusting the weights of the connections between neurons using an optimization algorithm, such as gradient descent, which adjusts the weights to minimize the error between the predicted output of the network and the true output [104].

The basic building block of ANNs is the **perceptron**, a simple model of a biological neuron. The perceptron takes a set of input values  $x_1, x_2, \dots, x_n$  and produces an output value  $y$  based on a set of weights  $w_1, w_2, \dots, w_n$  and a bias term  $b$ . The output of the perceptron is given by the formula

$$y = \sigma\left(\sum_{i=1}^n w_i x_i + b\right), \quad (2.26)$$

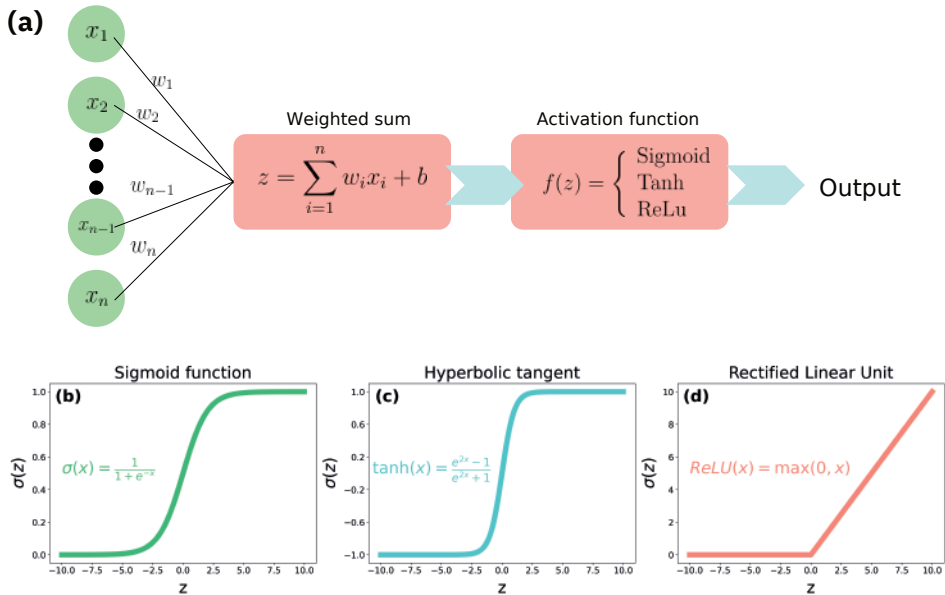
where  $\sigma$  is an **activation function**. The most common activation functions are the sigmoid function, the hyperbolic tangent function, and the rectified linear unit (ReLU) function, which are given by

$$\sigma(z) = \frac{1}{1 + e^{-z}}, \quad (2.27)$$

$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}, \quad (2.28)$$

$$\text{ReLU}(z) = \max(0, z), \quad (2.29)$$

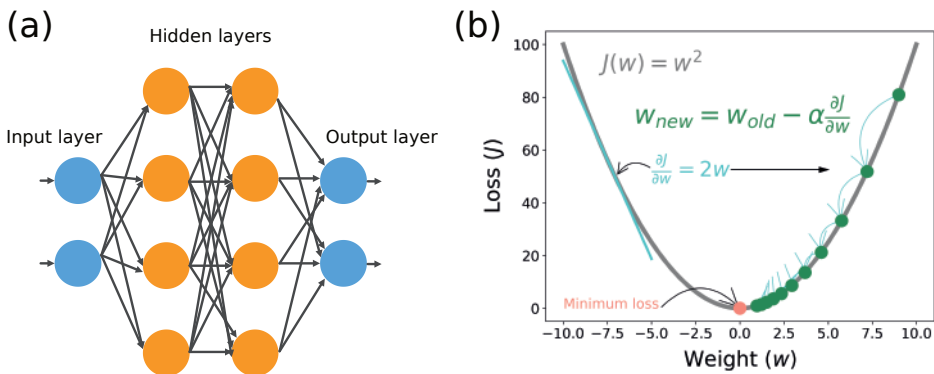
where  $z$  is the input to the activation function. The activation function is used to introduce non-linearity into the model, allowing the neural network to model complex patterns in the data. See Fig. 2.4 for a graphical representation of the perceptron model and the activation functions.



**Figure 2.4: The perceptron model of a biological neuron.** (a) The perceptron takes a set of input values  $x_1, x_2, \dots, x_n$  and produces an output value  $y$  based on a set of weights  $w_1, w_2, \dots, w_n$  and a bias term  $b$ . The output of the perceptron is given by the formula  $y = \sigma(\sum_{i=1}^n w_i x_i + b)$ , where  $\sigma$  is an activation function. (b) The sigmoid function saturates to 1 for large positive inputs and to 0 for large negative inputs. (c) The hyperbolic tangent function saturates to 1 for large positive inputs and to -1 for large negative inputs. (d) The rectified linear unit (ReLU) function is linear for positive inputs and zero for negative inputs.

The perceptron is a simple model of a biological neuron, but it can be combined to form more complex models, i.e., ANNs, which consist of multiple layers of neurons connected to each other. The first layer of the network is the input layer, which receives the input data, then the hidden layers process the data, and the output layer produces the final output of the network. The most common type of ANNs is the feedforward neural network (FFNN), in which the connections between neurons do not form cycles so that the information flows in only one direction, from the input to the output layer (Fig. 2.5(a)).

ANNs are often used in supervised learning tasks, in which the model is trained using a dataset containing input-output pairs. The weights of the connections between neurons are adjusted using an optimization algorithm to minimize the error between the predicted output of the network and the true output, which is known as the loss function. The basic idea behind these algorithms can be illustrated with the traditional gradient descent algorithm: the gradient of the loss with respect to the weights of the connections is calculated, and the weights are updated in the opposite direction of the gradient to minimize the loss, as shown in Fig. 2.5(b). The learning rate of the algorithm determines the size of the step taken in the direction of the gradient, and it is a hyperparameter that needs to be tuned to achieve good performance of the model. Nowadays, the most common optimization algorithms used to train ANNs are variations of the traditional gradient descent algorithm, such as stochastic gradient descent or the Adam optimizer [104].



**Figure 2.5: The feedforward neural network.** (a) The feedforward neural network consists of an input layer, multiple hidden layers, and an output layer. The input layer receives the input data, the hidden layers process the data, and the output layer produces the final output of the network. (b) Illustration of the gradient descent algorithm used to train the network. The weights of the connections between neurons are adjusted to minimize the loss, i.e., the error between the predicted output of the network and the true output.

Backpropagation is the algorithm used to calculate the gradient of the loss function with respect to the weights of the connections in a neural network. The algorithm works by propagating the error backwards through the network, starting from the output layer and moving towards the input layer. The gradient of the loss with respect to the weights is calculated using the chain rule of calculus, which allows the error to be decomposed into the errors of the neurons

in the previous layer. The weights are then updated using the gradient of the loss with respect to the weights, which is calculated by multiplying the error of the neuron by the input to the neuron. The process is repeated iteratively until the weights converge to a set of values that minimize the loss function. Backpropagation is a fundamental algorithm in training neural networks, and it allows the model to learn the features in the data that are most relevant for the task at hand.

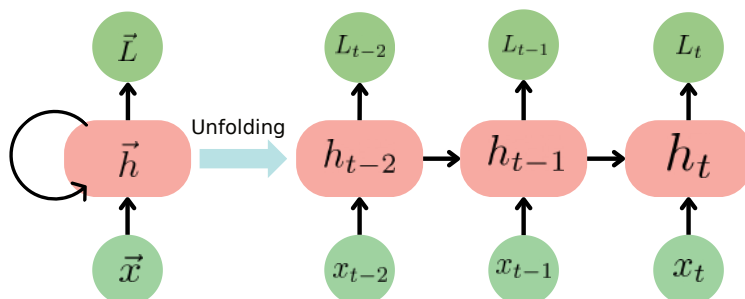
Several variations of ANNs have been developed to address specific tasks, such as image recognition and segmentation, speech recognition, natural language processing, and time series forecasting. These models have been shown to be effective at modeling complex patterns in large datasets and have achieved state-of-the-art performance in such tasks. In this thesis we have used two types of ANNs: long-short-term memory neural networks and convolutional neural networks, which are particularly well-suited for time series forecasting and image segmentation problems, respectively. In the following sections, we will provide a brief overview of these models.

### 2.3.3 Time series forecasting

Time series forecasting is a common task in machine learning that involves predicting future values of a time series based on past observations. Time series are sequences of data points collected at regular intervals over time, such as stock prices, weather data, and sales figures. Time series forecasting is used in a wide range of applications, such as predicting the stock market, forecasting the weather, and estimating the demand for products. Traditional methods for time series forecasting include autoregressive integrated moving average (ARIMA) models, exponential smoothing models, and state-space models [105]. These models are based on statistical methods and are effective at modeling linear relationships in the data, but they often struggle to capture complex patterns, such as non-linear relationships and long-term dependencies. FFNNs can be used for time series forecasting by using the  $T$  previous values of the time series as input to the network and predicting the next value of the time series as output. However, this approach does not take into account the temporal dependencies in the data, i.e., the relationship between the current value of the time series and the previous values, as the network treats each input value independently. This limitation can prevent the network from learning temporal patterns in the data, which are crucial for accurate forecasting.

To address these limitations, more advanced deep learning models have been developed for time series forecasting, recurrent neural networks (RNNs). RNNs are a class of artificial neural networks in which node connections arise along a temporal sequence, i.e., previous values in a time series are linked to current

values. In simple words, RNNs predict a point of the time series using past information. Simple Recurrent Neural Networks (SRNNs) are an extension of Feedforward Neural Networks, in which past information and learned knowledge is encoded in the network as state vectors,  $\vec{h}_t$ , that are passed from one time step to the next. In this way, the information flows not only from the input to the output layer but also from one time step to the next, allowing the network to learn temporal dependencies in the data. A schematic representation of the SRNN is shown in Fig. 2.6, where we can observe how the network receives as input a sequence of data points  $x_1, x_2, \dots, x_t$  and produces an output sequence  $L_1, L_2, \dots, L_t$ , called latent state, from a sequence of hidden states  $h_1, h_2, \dots, h_t$ , and the inputs. Notice how the current hidden state,  $h_t$ , is a function of the current input  $x_t$  and the previous hidden state  $h_{t-1}$ .



**Figure 2.6: Simple Recurrent Neural Network.** The SRNN receives as input a sequence of data points  $x_1, x_2, \dots, x_t$  and produces an output sequence  $L_1, L_2, \dots, L_t$ , called latent state, from a sequence of hidden states  $h_1, h_2, \dots, h_t$ , and the inputs. The produced latent state can be used to predict the future values of the time series.

Unfortunately, SRNNs suffer from the so-called vanishing gradient problem, i.e., distant parts of the time series do not play a role in the training process, which prevents the network from learning long-term temporal dependencies. This problem arises because the gradients of the loss function with respect to the weights of the connections tend to become very small as they are backpropagated through the network, which makes it difficult for the network to learn from distant parts of the time series. This is, the errors committed in the past are not propagated to the present. Another related issue is the exploding gradient problem, where the gradients become very large and cause the weights of the connections to diverge, leading to numerical instability.

We can visualize the problem by computing the backpropagation through time. Consider that we predict the output time series  $y_t$  at time  $t$  using the previous  $T$  time steps, i.e.,  $y_t = f(x_t, x_{t-1}, \dots, x_{t-T})$ . The change in the loss

function with respect to the weights of the connections is given by

$$\frac{\partial L}{\partial w} = \sum_{t=1}^T \frac{\partial L_t}{\partial w}, \quad (2.30)$$

where  $L$  is the loss function and  $w$  are the weights of the connections.

The contribution to the loss function of time step  $k$  can be obtained by the chain rule

$$\frac{\partial L_k}{\partial w} = \frac{\partial L_k}{\partial y_k} \cdot \frac{\partial y_k}{\partial h_k} \cdot \left( \prod_{t=2}^k \frac{\partial h_t}{\partial h_{t-1}} \right) \cdot \frac{\partial h_1}{\partial w}. \quad (2.31)$$

The term in red is the one that causes the vanishing and exploding gradient problems. Because the term is in a product, if it is smaller than 1 the gradient will vanish, while if it is larger than 1, the gradient will explode. This problem impacts the whole gradient (Eq. (2.30)), which becomes very small or very large, making the training of the network unstable and useless.

Long-Short Term Memory (LSTM) neural networks overcome this limitation by implementing three gates to update and control the cell state (forget gate, input gate, output gate), thus allowing to keep long-term dependencies [106]. Bidirectional Long-Short Term Memory (BD-LSTM) neural networks are able to encode both past and future information by implementing two LSTM layers flowing in opposite time directions. The forward layer preserves past information while the backwards layer preserves future information. Thus, using the two hidden states combined, BD-LSTM networks are able to preserve information from both the past and future at any point in time.

We will make use of BD-LSTM neural networks in Chapter 10 to predict the missing values of the ocean pH based on the past and “future” available observations of temperature, salinity, and dissolved oxygen.

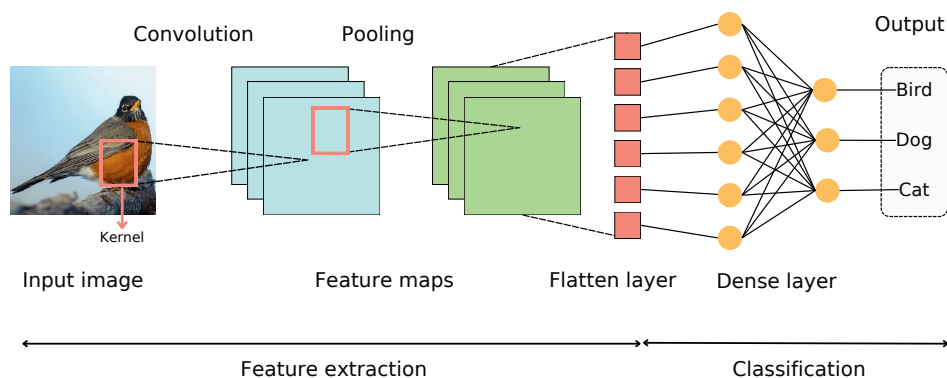
### 2.3.4 Image segmentation

Image segmentation is a computer vision task that involves partitioning an image into multiple segments, i.e., regions that share similar properties, such as color, texture, or intensity. Image segmentation is used in a wide range of applications, such as medical imaging, object detection, and autonomous driving. Traditional methods for image segmentation include thresholding, clustering, and edge detection, which are based on handcrafted features and heuristics. These methods are effective at segmenting images with simple patterns, but they often struggle to segment images with complex patterns and variations in lighting and texture.

FFNNs can be used for image segmentation. To do so the image is flattened<sup>2</sup> and fed into the network. For example, in the case of an RGB image, the input

<sup>2</sup>The image is transformed from a 2D array to a 1D array.

layer of the network would have 3 neurons, one for each color channel, and the output layer would have  $N$  neurons, where  $N$  is the number of classes to segment. However, they are not well-suited for this task because they do not take into account the spatial information in the image. Indeed, we know how the brain processes images: the visual cortex is organized in a hierarchical manner, with neurons in the lower layers detecting simple patterns, such as edges and textures, and neurons in the higher layers detecting more complex patterns, such as shapes and objects. This hierarchical organization allows the brain to model complex patterns in the image and segment the image into different regions based on the spatial information. To mimic this hierarchical organization, more advanced deep learning models have been developed for image segmentation, which are commonly known as convolutional neural networks (CNNs).



**Figure 2.7: Convolutional Neural Network.** The CNN receives an image as input and produces an output image that goes through a series of convolutional and pooling layers, applying a set of filters to the input image to extract features. The filters are learned during the training process and are used to detect patterns in the image, such as edges, textures, and shapes. The output of the convolutional and pooling layers, which are called feature maps, is passed through a series of non-linear activation functions, which introduce non-linearity into the model and allow the network to learn complex patterns in the data.

CNNs are a class of artificial neural networks that are specifically designed to process grid-like data, such as images. CNNs consist of multiple layers of neurons, called convolutional layers, which apply a set of filters (or kernels) to the input image to extract features. After each convolutional layer, a pooling layer is applied to reduce the spatial dimensions of the feature maps and make the network more computationally efficient, allowing the network to learn spatial hierarchies of features (Fig. 2.7). After the last pooling layer is applied, the

output of the network is flattened and passed through a series of fully connected layers to produce the final output of the network, e.g., the classification of the image into different classes, as shown in Fig. 2.7.

The convolutions in the convolutional layer are performed by sliding a filter over the input image and computing the dot product between the filter and the input at each position. The filter is nothing but a small matrix of weights that holds the information of the pattern that the filter is looking for in the image. Considering an image  $I$  of size  $N \times M \times L$  and a filter  $K$  of size  $n \times m \times l$ , the output of the convolutional layer,  $F$ , is given by the convolution of the image and the filter, which is computed as

$$F(i, j) = I(i, j) * K(i, j) = \sum_{x=0}^n \sum_{y=0}^m \sum_{z=0}^l I(i+x, j+y, z) \cdot K(x, y, z), \quad (2.32)$$

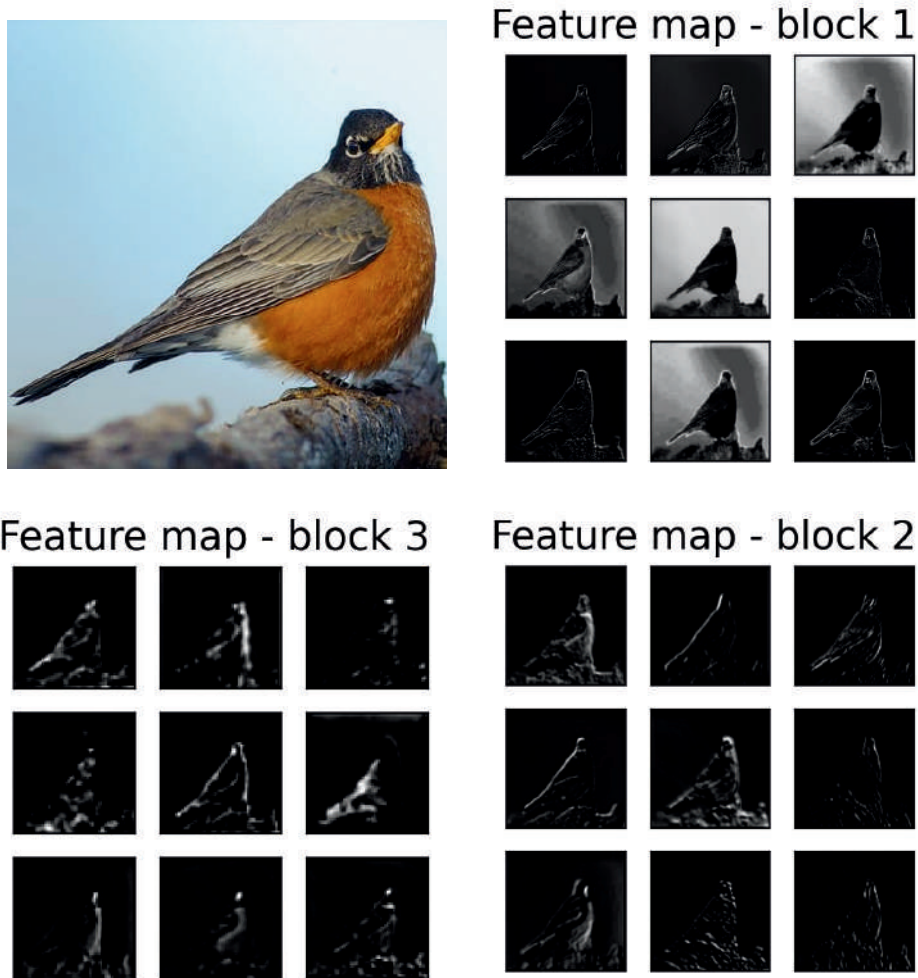
where  $i$  and  $j$  are the spatial coordinates of the output feature map  $F$ .

For a 2D image with a single channel, we can visualize the convolutional operation with the following example,

$$I = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix}, \quad K = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \implies F = I * K = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} * \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}, \quad (2.33)$$

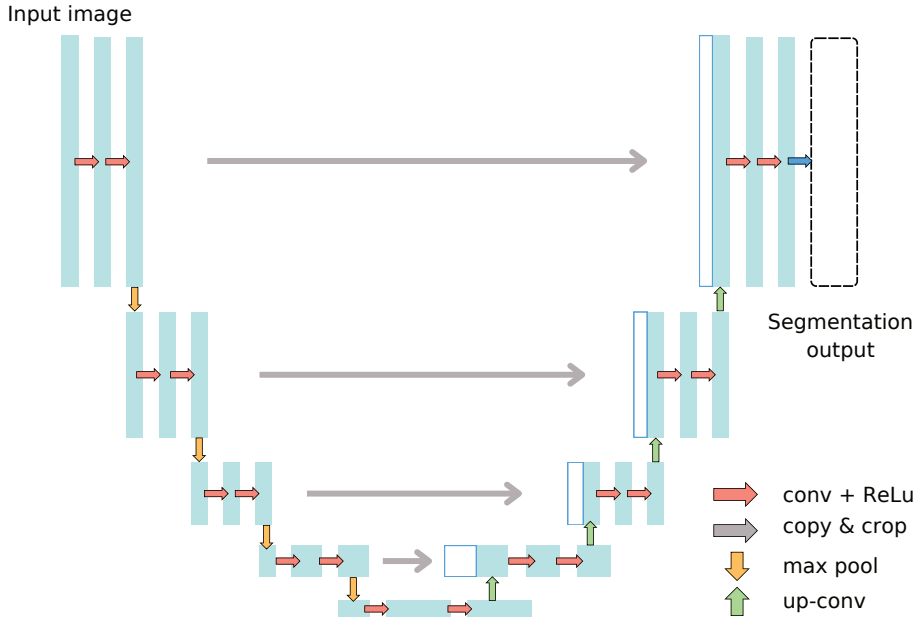
The outputs of the convolutional and pooling layers are called feature maps, which are used to detect patterns in the image, such as edges, textures, and shapes. The filters are learned during the training process and are used to detect patterns that are relevant for the task at hand. An example of feature maps produced by a CNN (VGG16) is shown in Fig. 2.8.

CNNs are the basic blocks used in more complex models specifically designed for image segmentation, such as U-Net or Linknet. These models have achieved state-of-the-art performance in image segmentation tasks and have been used in a wide range of applications, such as medical imaging, object detection, and autonomous driving. These models implement even more complex architectures called encoder-decoder networks, which consist of an encoder that extracts features from the input image and a decoder that reconstructs the image from the features. The encoder is composed of a series of convolutional and pooling layers that extract features from the input image, while the decoder is composed of a series of upsampling and convolutional layers that reconstruct the image from the features (Fig. 2.9).



**Figure 2.8: Feature maps in a convolutional neural network.** The convolutional layers of a CNN apply a set of filters to the input image to extract features. The filters are learned during the training process and are used to detect patterns in the image, such as edges, textures, and shapes. Inspired by [107]

In [Chapter 11](#), we will use different advanced models based on CNNs with encoder-decoder architectures to segment satellite images of the coastal areas of the Balearic Islands into different benthic habitats, such as seagrass meadows or sandy bottoms.



**Figure 2.9: U-Net architecture for image segmentation.** The U-Net model consists of an encoder that extracts features from the input image and a decoder that reconstructs the image from the features. The encoder is composed of a series of convolutional and pooling layers that extract features from the input image, while the decoder is composed of a series of upsampling and convolutional layers that reconstruct the image from the features.

### 2.3.5 Performance metrics

The performance of machine learning models can be evaluated using different metrics that quantify the accuracy of the predictions made by the model. These metrics are used to compare the performance of different models and to assess the quality of the predictions made by the model. The choice of performance metrics depends on the task at hand, and different metrics are used for classification, time series forecasting, or segmentation tasks.

For time series forecasting tasks, the most common performance metrics are the mean absolute error (MAE), the mean absolute percentage error (MAPE), the mean squared error (MSE), and the root mean squared error (RMSE). The MAE measures the average absolute difference between the predicted values and the true values,

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|, \quad (2.34)$$

where  $y_i$  is the true value and  $\hat{y}_i$  is the predicted value. The MAPE measures

the average percentage difference between the predicted values and the true values,

$$\text{MAPE} = \frac{1}{n} \sum_{i=1}^n \left| \frac{y_i - \hat{y}_i}{y_i} \right|. \quad (2.35)$$

The MSE measures the average squared difference between the predicted values and the true values,

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2, \quad (2.36)$$

and the RMSE is the square root of the MSE,

$$\text{RMSE} = \sqrt{\text{MSE}}. \quad (2.37)$$

For image segmentation tasks, the most common performance metric is the intersection over union (IoU), but other metrics such as the F1 score, precision and recall are also used to complement the model evaluation and provide a more comprehensive assessment of the model performance. The IoU measures the overlap between the predicted segmentation and the true segmentation, and it is defined as the intersection of the predicted and true segmentation divided by the union of the predicted and true segmentation,

$$\text{IoU} = \frac{A_{\text{pred}} \cap A_{\text{true}}}{A_{\text{pred}} \cup A_{\text{true}}}, \quad (2.38)$$

where  $A_{\text{pred}}$  is the set of pixels of the predicted segmentation and  $A_{\text{true}}$  is the set of pixels of the true segmentation. The F1 score is the harmonic mean of the precision and recall, and it is defined as

$$\text{F1} = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}, \quad (2.39)$$

where the precision and recall are defined as

$$\text{precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}, \quad (2.40)$$

$$\text{recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}, \quad (2.41)$$

this is, the precision is the ratio of true positive predictions to the total number of positive predictions, and the recall is the ratio of true positive predictions to the total number of true positive instances.

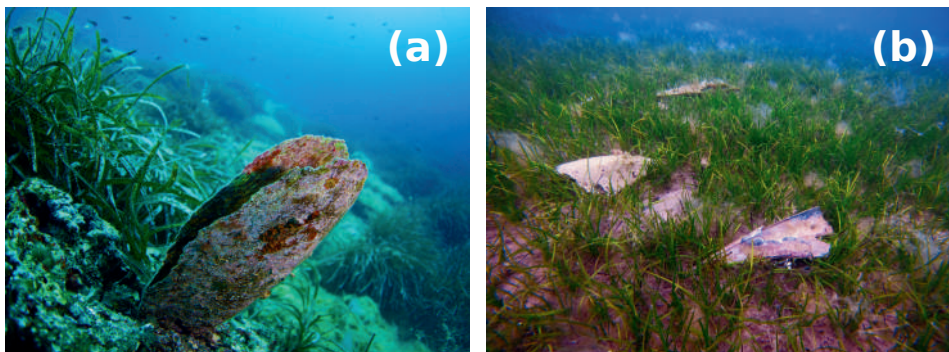
These metrics are used to evaluate the performance of the models presented in [Part IV](#) of this thesis and to compare the performance of different models.

## 2.4 Ecological systems

In this section we will introduce the basic concepts and context of the ecological systems that will be studied throughout the thesis.

### 2.4.1 The Mass Mortality Event of *Pinna nobilis*

*Pinna nobilis* is the largest endemic bivalve in the Mediterranean Sea and is under a serious extinction risk due to a Mass Mortality Event (MME) that has recently occurred throughout the whole Mediterranean basin [108–110]. Right before this MME, it was distributed across a wide type of habitats including coastal and paralic ecosystems, at depths between 0.5 to 60m [111, 112]. In open coastal waters, the distribution of the species is mainly associated with seagrass meadows, typically of *Posidonia oceanica*, which has been indicated as its optimal habitat [113]. Its lifespan is up to 50 years in favorable conditions and its size can get up to 1.2m, placing it among the largest bivalves of the world [114]. These fan mussels play a crucial ecological role in their habitat, as *P. nobilis* individuals filter water, thus retaining a large amount of organic matter from suspended detritus, contributing to water clarity [115]. Furthermore, it is a habitat-forming species, because its shell provides a hard-surface within a soft bottom ecosystem, which can be colonized by different benthic species, augmenting biodiversity [114]. In addition, at very dense populations, the species can function as an ecosystem engineer, creating biogenic reefs [116].



**Figure 2.10: The fan mussel *Pinna nobilis*.** (a) An alive individual of *Pinna nobilis* in its natural habitat. (b) Dead individuals of *Pinna nobilis* after the Mass Mortality Event, Bay of Pollença, Mallorca, Spain.

Despite *P. nobilis* populations have greatly declined due to anthropogenic activities in the 20th century [110], the ongoing MME is the most worrying and widespread threat to *P. nobilis* throughout the Mediterranean Sea. As a conse-

quence, the species has been declared as critically endangered [117]. Although different aetiological agents have been proposed, including Mycobacteria and other bacteria [118–120], there is evidence that the main cause of this mortality is the protozoan *Haplosporidium pinnae* [121–123], a new species that belongs to the genus *Haplosporidium*, one of the four genera of the protist order Haplosporida. Indeed, other Haplosporidian parasites have been previously found to be behind the extensive mortality of several oyster species [124, 125].

*Haplosporidium pinnae* exhibits a complex life cycle that includes uninucleate and binucleate cells, plasmodia, and spores. The parasite's life stages and infection mechanisms are not completely understood, but it is clear that it can cause severe pathology in infected mussels. Environmental factors, particularly temperature and salinity, significantly influence the expression and spread of the disease. Optimal conditions for the parasite's proliferation are temperatures above 13.5 °C and salinity ranges between 36.5 and 39.7 psu. Marine currents also play a vital role in dispersing the parasite, facilitating its spread across large areas. This has been evidenced by the rapid spread of the disease throughout the Mediterranean basin, affecting populations in different regions, which poses a serious threat to the survival of the species [114].

The event has spurred a concerted effort to better understand the dynamics of marine diseases, the role of environmental stressors, and the mechanisms of disease spread. In this context, the development of predictive models that can forecast the spread of the disease and assess the impact of environmental factors on the disease dynamics is crucial. In [Part I](#) of this thesis, we will develop a compartmental model to study the dynamics of the MME of *P. nobilis* based on the theory of infectious diseases explained in [Section 2.1](#).

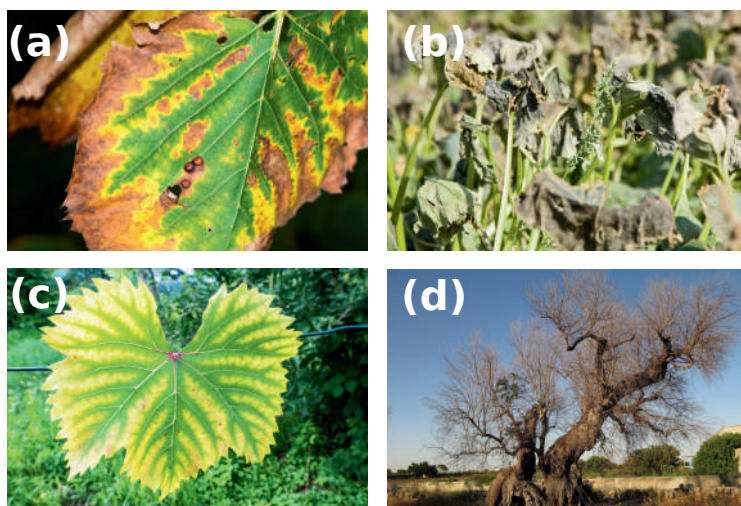
### 2.4.2 *Xylella fastidiosa*: an emerging global threat

*Xylella fastidiosa* (Xf) is a Gram-negative<sup>3</sup> xylem-limited bacterium, this is, a bacterium that colonizes the xylem vessels of plants, where it multiplies [126]. Xf is a plant pathogen with a very wide host range, being able to infect over 600 plant species [127]. The bacterium is transmitted by xylem-feeding insects, such as sharpshooters and spittlebugs, which acquire the bacterium when feeding on infected plants and transmit it to healthy plants when feeding on them [128, 129]. The bacterium causes a variety of symptoms in infected plants, including leaf scorching, chlorosis, wilting, and finally dieback, which can lead to the death of the plant. Xf is responsible for several important crop diseases, which lead to severe economic losses in agriculture and forestry and pose a serious threat to food security and biodiversity worldwide.

---

<sup>3</sup>Gram-negative bacteria have a double cell wall, which makes them more resistant, in contrast to the Gram-positive bacteria, which have a single cell wall.

As a taxonomic unit, Xf comprises three recognized subspecies, *fastidiosa*, *multiplex*, and *pauca*, and more than 90 sequence types (i.e., genetic lineages) with distinct host ranges. These lineages are associated with different diseases in different hosts, and some of them are more virulent than others. Xf subsp. *fastidiosa* is the most widespread subspecies and is responsible for several important crop diseases, such as Pierce's disease of grapevines (PD) or almond leaf scorch disease (ALSD), the latter also caused by Xf subsp. *multiplex*. Specifically, the Xf clonal lineage of the subspecies *fastidiosa* that causes Pierce's disease (hereafter Xf<sub>PD</sub>) also causes almond leaf scorch in California [130]. On the other hand, Xf subsp. *pauca* is the most virulent subspecies and is responsible for the olive quick decline syndrome (OQDS), citrus variegated chlorosis (CVC), and coffee leaf scorch (CLS).

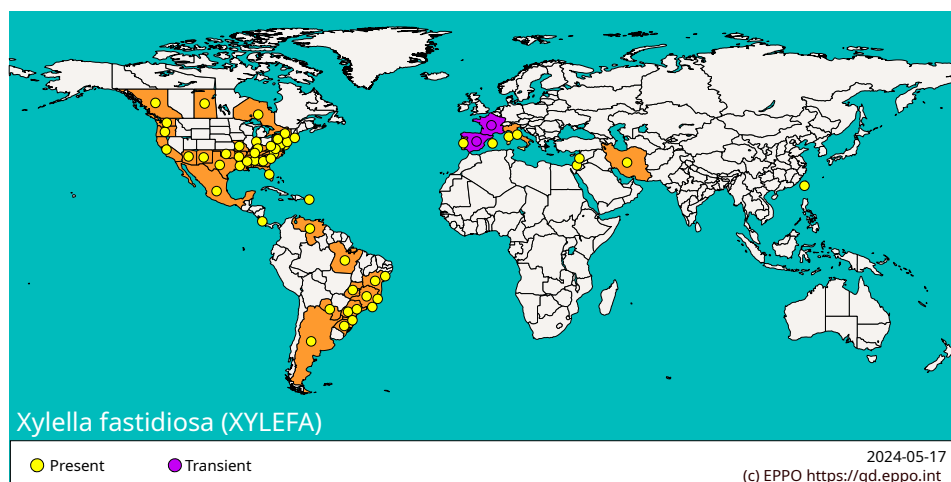


**Figure 2.11: Symptoms of *Xylella fastidiosa* infection.** (a) Leaf scorching in an American Elm leaf (b) Wilted squash plant (b) Leaf scorching (c) Chlorosis in a grapevine leaf (d) Dieback in an infected olive tree.

Until the beginning of the 21st century, Xf was a pathogen officially restricted to the American continent [131]. In 2013, the involvement of Xf subsp. *pauca* in the massive death of ancient olive trees in Apulia, Italy, and its rapid spread raised alarm in European agriculture [132]. Today, all three Xf subspecies have been detected in the Balearic Islands (Spain), including Xf<sub>PD</sub>, and several clonal lineages have been found in Corsica and the PACA region of France, Alicante (Spain), Tuscany (Italy), and Portugal [133–135]. Outside North America, Xf<sub>PD</sub> is only established on the islands of Mallorca and Taiwan and has recently been detected in Israel, Lebanon, and Portugal [136, 137]. In

February 2024, some months before the writing of this thesis, Xf<sub>PD</sub> was also found in Italy [138]. In all European outbreaks, the insect vector *Philaenus spumarius* is the main and almost unique carrier of Xf [139].

*Philaenus spumarius* (Ps) is a polyphagous xylem-feeding insect that feeds on a wide range of plant species, including many crops and ornamental plants. The insect is native to Europe and is widely distributed throughout the Mediterranean basin. Ps is a univoltine species<sup>4</sup> with a long life cycle. Eggs overwinter until hatching occurs after enough heat is accumulated, usually in early spring [140]. After egg hatching, the nymphal stage goes through five instars during about 5–6 weeks until they become adults [141]. In the fall, mature females lay masses of eggs on plant debris on the soil until they die during winter [142]. This complex life cycle, which is influenced by temperature, gives rise to a seasonal pattern of Ps abundance, with a peak in the spring and a decline in the summer and fall [143].



**Figure 2.12: Distribution of *Xylella fastidiosa* in Europe.** The map shows the distribution of *Xylella fastidiosa* in Europe, including the Balearic Islands, Corsica, the PACA region of France, Alicante, Tuscany, and Portugal. Source: EPPO Global Database.

The introduction of Xf in the Mediterranean basin has raised concerns about its possible spread to continental Europe and the potential impact it might cause on the region's agriculture and biodiversity. In this context, the development of predictive models that can forecast the spread of Xf in affected zones and assess the risk of disease establishment in new areas is crucial. Most of the efforts to

<sup>4</sup>Species that have one generation per year.

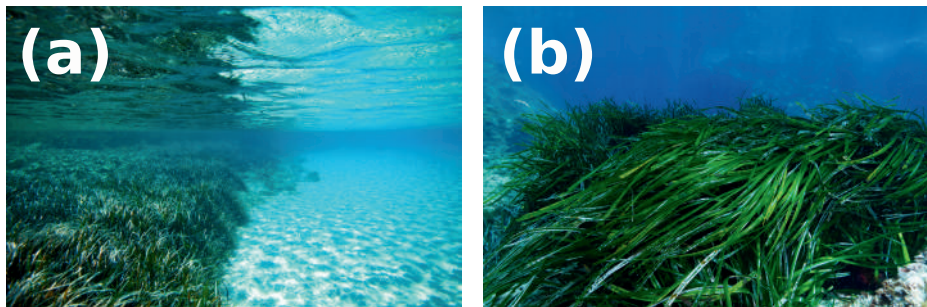
model the spread of Xf have focused on the epidemiology of the disease and the role of environmental factors in the spread of the bacterium using SDMs, which we have already argued are not appropriate for this task. This has led to biased predictions and a lack of understanding of the mechanisms driving the spread of the disease. On the other hand, most of the studies that have focused on the spread of Xf at the local level have ignored basic ecological principles, such as the role of the insect vector in the spread of the disease. Recent studies have greatly improved the models, but the characteristic seasonal patterns of the insect vector have not been taken into account, which can still lead to biased predictions and a lack of formal understanding of the mechanisms driving the spread of the disease.

In [Part II](#) of this thesis, we will develop a compartmental model for the local spread of *Xylella fastidiosa* diseases that includes all relevant ecological and epidemiological processes of the disease, including the seasonal abundance patterns of the insect vector. Later on, on [Part III](#), we will develop a mechanistic climate-driven epidemiological model to assess the risk of establishment of Xf<sub>PD</sub> in new areas. This model will be based on a simplified version of the model developed in [Part II](#) and will be used to assess the risk of establishment of Xf<sub>PD</sub> based on current and future climate conditions. Noteworthy, our methodology helps to advance the field of disease biogeography by providing a mechanistic understanding of the role of environmental factors in the spread of plant diseases.

### 2.4.3 The decline of seagrass meadows

Seagrass meadows are among the most productive and diverse ecosystems on Earth, providing habitat for a wide range of marine species [144] and playing a crucial role in the health of the oceans and the well-being of human populations worldwide. They provide vital food, shelter, and structural support, including nursery areas for commercially important species, thereby supporting both local economies and subsistence fisheries [145]. Moreover, seagrass ecosystems play a significant role in coastal erosion prevention. The dense canopies of seagrass attenuate currents and waves, facilitating particle sedimentation and mitigating sediment resuspension [146–149]. Additionally, the extensive underground network of rhizomes and roots stabilizes sediment, reducing erosion and decreasing water turbidity [150], strongly influencing coastal sedimentary dynamics [151, 152]. Specifically, the robust root systems of *Posidonia oceanica* beds act as natural barriers, protecting coastlines from the destructive force of strong waves and maintaining shoreline stability [153–155]. In addition, seagrasses keep overlying waters oxygenated and with low concentrations of nutrients and CO<sub>2</sub> [156]. Seagrasses are among the planet's most effective natural ecosystems for seques-

tering (capturing and storing) carbon, performing a rate that is 35 times faster than tropical rainforests, while their sediments never become saturated [157].



**Figure 2.13: The seagrass *Posidonia oceanica*.** (a) A seagrass meadow of *Posidonia oceanica* in the Mediterranean Sea. (b) A close-up of a *Posidonia oceanica* meadow.

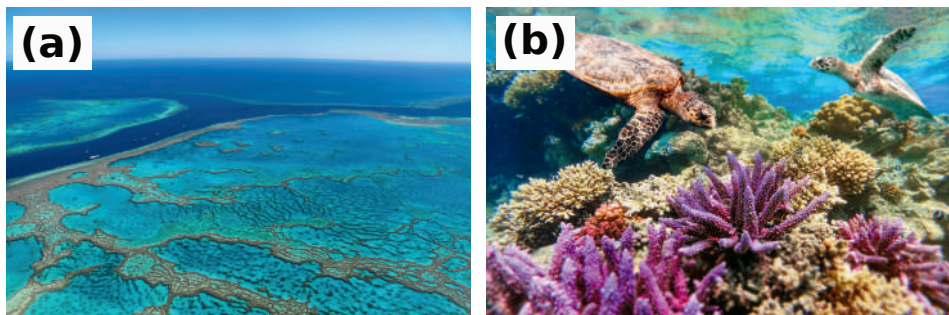
Unfortunately, seagrass meadows are facing numerous threats that are leading to their rapid decline worldwide. Human activities, such as coastal development, dredging, and land reclamation, are major contributors to habitat destruction. These activities often result in physical damage to seagrass beds and increased sedimentation, which can smother seagrass shoots and inhibit their growth. Pollution, particularly nutrient loading from agricultural runoff and sewage discharge, leads to eutrophication. This process stimulates excessive growth of algae, which can outcompete seagrasses for light and resources, ultimately causing large-scale die-offs of seagrass meadows [158]. Climate change poses additional challenges to the health and sustainability of seagrass ecosystems. Rising sea temperatures can stress seagrass plants, reducing their overall effective growth rates and increasing their susceptibility to diseases [159]. Additionally, more frequent and severe storm events can cause physical damage to seagrass meadows, further exacerbating their decline [160].

In this context, the monitoring of seagrass meadows is crucial for the conservation and management of these ecosystems. Traditional methods for monitoring seagrass meadows, such as side-scan sonar, are expensive and time-consuming, making it difficult to obtain accurate and up-to-date information on the extent and health of seagrass meadows. Remote sensing offers a cost-effective and efficient alternative for monitoring seagrass meadows, providing valuable information on the spatial distribution and health of seagrass ecosystems. In Chapter 11 of this thesis, we will develop a deep learning model for monitoring seagrass meadows based on satellite images. This model will provide accurate and cost-effective estimates of the extent of seagrass meadows,

which can be used to assess the health and sustainability of these ecosystems and inform conservation and management efforts.

#### 2.4.4 The structure of coral reefs

Coral reefs are one of the most biodiverse ecosystems on Earth, holding more than 25% of marine life with only 1% of ocean floor coverage [161]. Corals are colonies of living organisms called polyps, held together by a self-produced exoskeleton made of calcium carbonate. Polyps host photosynthetic microalgae in a mutualistic interaction: polyps obtain nutrients from microalgae while microalgae are compensated with protection and some nutrients too [162]. Consequences of global change, such as temperature increases and ocean acidification, represent a substantial threat to coral reef ecosystems [163]. High ocean temperatures promote the “bleaching” phenomenon, in which individual polyps expel the microalgae symbionts, losing their characteristic color and getting rid of their primary source of nutrients [164]. On the other hand, ocean acidification promotes the dissolution of the calcium carbonate exoskeleton of corals [31]. These two phenomena are causing widespread declines in coral reef ecosystems worldwide, leading to the loss of biodiversity and ecosystem services provided by coral reefs.



**Figure 2.14: The structure of coral reefs.** (a) Aerial view of a coral reef in the Great Barrier Reef, Australia. (b) A close-up of a coral reef in the Red Sea, Egypt.

Understanding the formation of coral reefs is crucial for the conservation and management of these ecosystems. Coral reefs are formed by the accumulation of coral exoskeletons over thousands of years, creating complex three-dimensional constructions that conform some of the largest biogenic structures in the biosphere [165]. Several forms and sizes of coral reefs can be found in the oceans, ranging from small fringing reefs to large barrier reefs, which can extend for hundreds of kilometers along the coastlines of continents and islands. Perhaps

the most famous coral barrier reef is the Great Barrier Reef, located off the coast of Queensland, Australia, which is the largest coral reef in the world and one of the most diverse ecosystems on Earth.

Some more mechanistic hypotheses on reef formation date back to the end of the 19th century, when Charles Darwin proposed that coral reefs were formed by the growth of corals on the remains of extinct volcanoes [166]. However, this would only explain the formation of atolls, which are circular coral reefs that enclose a lagoon, and not all the different observed forms. Recently, more detailed models have been proposed to explain the formation of coral reefs, which take into account the interactions between corals, algae, and other organisms, as well as the physical and chemical processes that shape the reef structure, such as hydrodynamics. To evaluate these models and ultimately understand the formation of coral reefs, it is necessary to contrast them with a comprehensive empirical dataset on the size and geometry of coral reefs. Unfortunately, such a dataset has been hitherto lacking.

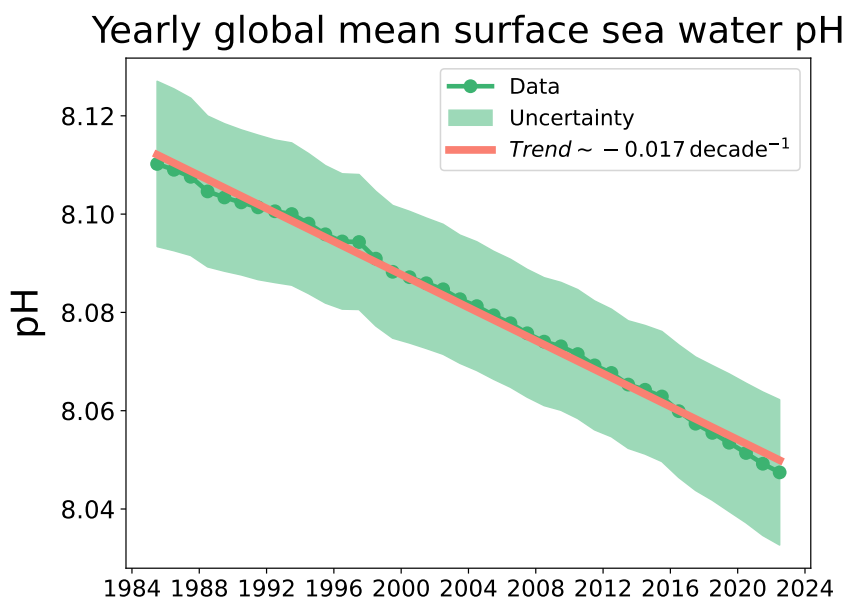
In Chapter 12 of this thesis, we will present a global analysis on the spatial properties of all tropical coral reefs worldwide, which has unraveled universal patterns in their size distribution and geometry. The results of this analysis can be used to evaluate the performance of existing models of coral reef formation and to develop new models that can explain the observed patterns in the data. In addition, our analysis provides important insights into the ecology of coral reefs and can help to design effective conservation strategies. Thus, our work contributes to the understanding of the formation of coral reefs and the conservation of these important ecosystems.

### 2.4.5 Ocean acidification

Ocean acidification is a global environmental problem that is caused by the uptake of carbon dioxide ( $\text{CO}_2$ ) from the atmosphere by the oceans. When  $\text{CO}_2$  dissolves in seawater, it reacts with water to form carbonic acid ( $\text{H}_2\text{CO}_3$ ), which dissociates into bicarbonate ( $\text{HCO}_3^-$ ) and hydrogen ions ( $\text{H}^+$ ). The increase in hydrogen ions leads to a decrease in oceanic pH, making them more acidic. The pH of the oceans has decreased by 0.1 units since the beginning of the industrial revolution, which corresponds to a 30% increase in the acidity of the oceans [167]. This trend is expected to continue in the future as atmospheric  $\text{CO}_2$  levels continue to rise, leading to further decreases in the oceanic pH.

Ocean acidification has a wide range of impacts on marine ecosystems, including changes in the growth and survival of marine organisms, the structure and function of marine communities, and the cycling of nutrients and energy in marine ecosystems. The decrease in the pH of the oceans can have negative effects on the growth and survival of marine organisms, particularly those that

rely on calcium carbonate to build their shells and skeletons, such as corals, mollusks, and some planktonic species. The increase in acidity can make it more difficult for these organisms to build and maintain their calcium carbonate structures, leading to reduced growth rates and increased mortality [19]. In addition, ocean acidification can alter the structure and function of marine communities by changing the composition of species and the interactions between them. For example, ocean acidification can favor the growth of some species over others, leading to shifts in the composition of marine communities and changes in the dynamics of marine food webs [168]. Finally, ocean acidification can affect the cycling of nutrients and energy in marine ecosystems by altering the rates of biological processes, such as photosynthesis and respiration, and the availability of nutrients, such as nitrogen and phosphorus.



**Figure 2.15: Global trend in ocean pH.** The pH of the oceans has decreased by 0.1 units since the beginning of the industrial revolution, which corresponds to a 30% increase in the acidity of the oceans. Data from Copernicus Marine Service Information [169]

Monitoring ocean acidification is crucial for understanding the impacts of this global environmental problem on marine ecosystems and for developing effective conservation and management strategies. However, monitoring ocean acidification is challenging due to the spatial and temporal variability of the pH of the oceans, which is influenced by a wide range of factors, such as tem-

perature, salinity, and biological activity. In addition, the pH of the oceans is influenced by the uptake of  $\text{CO}_2$  from the atmosphere, which varies seasonally and regionally. As a result, monitoring ocean acidification requires the collection of large amounts of data on the pH of the oceans. These challenges are exacerbated by the fact that pH sensors often fail or drift over time, leading to gaps in the time series of ocean pH data. These gaps can introduce bias in the estimates of the trends in the pH of the oceans and make it difficult to assess the impacts of ocean acidification on marine ecosystems.

In [Chapter 10](#), we will present a new framework to reconstruct the missing data in the pH time series of a coastal area of the Balearic Islands. The framework is based on deep learning, specifically recurrent neural networks, which are able to learn the temporal patterns in the data and predict the missing values in the time series. Our framework provides accurate estimates of the pH that can be used to fill the gaps in the time series of pH data, providing valuable information on the impacts of ocean acidification on marine ecosystems.



# Parasite-produced marine diseases

<b>3</b>	<b>The case of the mass mortality event of <i>Pinna nobilis</i> . . . .</b>	<b>61</b>
3.1	Introduction . . . . .	62
3.2	The SIRP model . . . . .	64
3.3	Numerical analysis . . . . .	72
3.4	Model validation . . . . .	79
3.5	Conclusions . . . . .	84
<b>4</b>	<b>Spatial effects in parasite-induced marine diseases . . . . .</b>	<b>87</b>
4.1	Introduction . . . . .	88
4.2	The SIRP spatial model . . . . .	90
4.3	Results . . . . .	92
4.4	Conclusions . . . . .	101



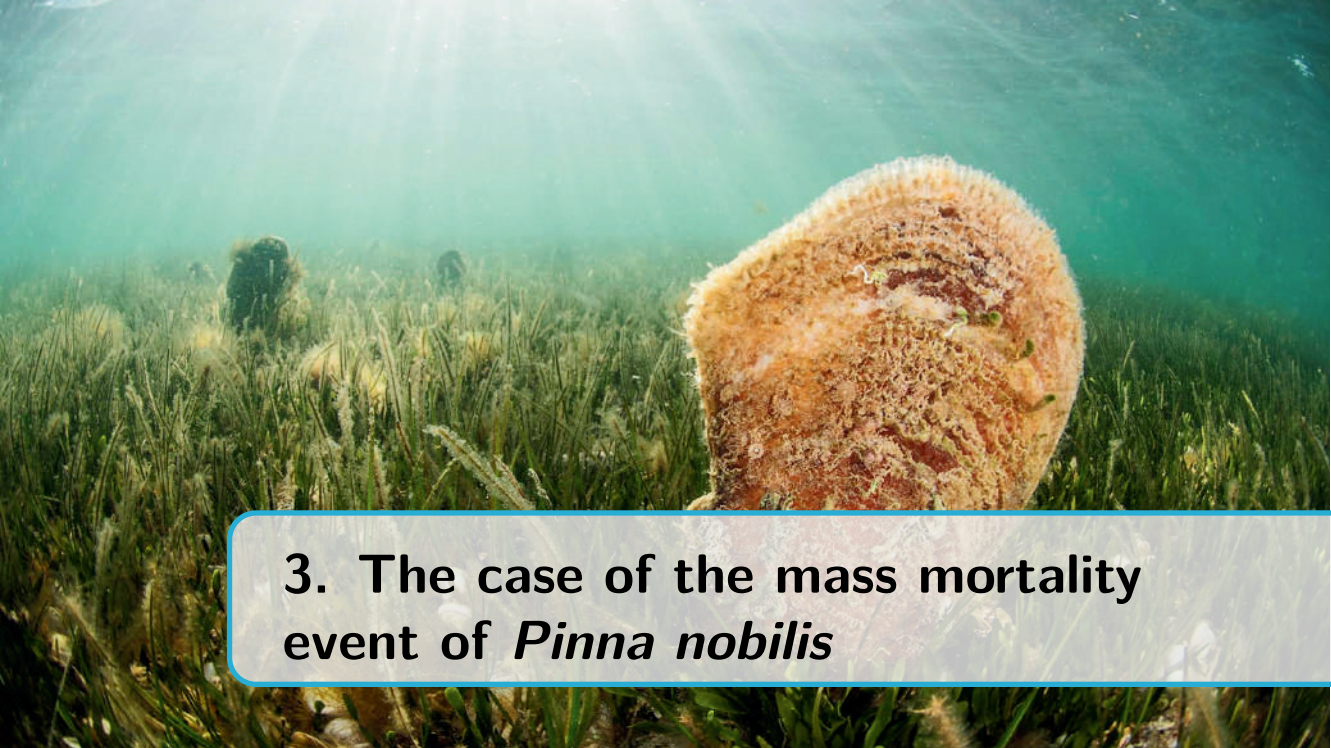
## Summary

Since 2016, the Mediterranean Sea has been experiencing a Mass Mortality Event (MME) of the fan mussel *Pinna nobilis*, caused by the parasite *Haplosporidium pinnae*. This event has raised concerns about the conservation of this species, which is listed as endangered by the International Union for Conservation of Nature (IUCN). The event has also highlighted the need for a better understanding of the dynamics of marine diseases, which are often overlooked in the literature. The state-of-the-art models for disease dynamics in marine ecosystems lag behind those for terrestrial systems. Compartmental models have been introduced only recently, lacking a detailed mathematical analysis, and spatial dynamics have been largely ignored. In this part, we study the dynamics of parasite-induced marine diseases of immobile hosts, using the MME of *Pinna nobilis* as a focal point. We develop a compartmental model to describe the transmission dynamics of diseases in sessile marine organisms. Interestingly, an in-depth mathematical analysis of the model reveals the conditions under which these diseases can be modeled as simple SIR-type systems. Validation against empirical data from the *Pinna nobilis* mortality event demonstrates the practical utility of the model for future predictions and management of marine diseases. We then extend the model to include spatial dynamics and show how the movement of parasites influences disease spread and severity. The insights gained not only contribute to our fundamental understanding of marine epidemiology but also underscore the critical role of spatial dynamics in shaping disease outcomes in marine environments.

### Objectives

- To develop a mathematical model for disease transmission among sessile marine organisms.
- To validate the developed model against real-world data from the *Pinna nobilis* mass mortality event.
- To extend the model to include spatial dynamics.
- Understand how the movement of parasites influences disease spread and severity.





### 3. The case of the mass mortality event of *Pinna nobilis*

**Published as:**

À. Giménez-Romero, A. Grau, I. E. Hendriks, and M. A. Matias, “Modelling parasite-produced marine diseases: The case of the mass mortality event of *Pinna nobilis*”, [Ecological Modelling](#) **459**, 109705 (2021)

## 3.1 Introduction

Marine organisms, like their terrestrial counterparts, can serve as hosts for a diversity of parasites and pathogens present in the ecosystem, which are directly responsible for disease outbreaks. Disease-induced mortality affects not only the host population but can cascade through the whole ecosystem, altering its structure and functioning [170]. Furthermore, climate change can increase the spread range and impact of parasites and pathogens [171]. In fact, marine infectious diseases are recently increasing due to climate change and other anthropogenic pressures, like pollution and overfishing [172]. This, in turn, threatens many ecologically valuable habitats and can also result in substantial economic losses in e.g., aquaculture [173]. Analyzing the impact of these events at appropriate scales (spatial and temporal) and biological organization levels (species, populations, and communities) is crucial to accurately anticipate future changes in marine ecosystems and propose adapted management and conservation plans [174]. Thus, there is a strong need to address the mechanism of disease propagation in marine organisms.

However, the state of the art of epidemiological studies in marine ecosystems lags behind that of terrestrial ecosystems [175]. Contact and vector-borne based infectious diseases of terrestrial vertebrates and their epidemiology are typically studied using variations of the classical formulation of Kermack and McKendrick [65, 176, 177], the SIR model. Among other things, this formalism allows understanding why epizootics spread and stop, as the propagation of a disease is a threshold phenomenon [178], regulated by the now commonplace  $R_0$  dimensionless number. Within this framework, the initiation of epidemic transmission occurs when an infected individual is in close contact with a susceptible host or through a transmission vector, as typically pathogens can only survive for a very limited time outside the host in an aerial environment. On the other hand, as air is typically a much harsher medium for pathogens than water, the sea is expected to host many pathogens (viruses, bacteria, and parasites) for a relatively long time. The longer life span of pathogens in a water medium, together with the increased buoyancy arising from the different physical properties of seawater and air, coupled to the existence of marine currents that can transmit pathogens for long distances away, allows diseases to spread faster and reach further distances in marine environments compared to epidemics in terrestrial systems [179]. As a result, the possible long-term transmission of parasites by currents in marine environments make them more prone to suffer from persistent zoonoses compared to terrestrial ecosystems, where for an epidemic outbreak to occur the presence of an initial infected host (or vector) is necessary within a susceptible population. Until quite recently, marine zoonoses were mostly studied using different models compared to terrestrial diseases, and it was not

even clear whether the same tools could be applied [180]. The abundance of pathogens in marine ecosystems is one of the reasons why proliferation models, that do not focus on transmission and assume a widespread occurrence of the pathogen and a rapid transmission problem, have been most popular in the field [181]. In fact, compartmental models are starting to be used only recently in the study of marine epizootics [182].

An important subset of marine organisms are sessile, e.g., bivalves, which means that they cannot move. In the case of sessile terrestrial organisms, disease transmission occurs mostly through vectors, insects that transmit the pathogens causing the disease. Instead, in marine ecosystems, disease transmission is most often waterborne, in particular in passive water filtering feeders, as is the case of bivalves. Recently, some compartmental models considering the pathogen population have been proposed to study particular bivalve epidemics [182–184]. Here we develop and analyze a model that describes disease transmission from an infected immobile host to a susceptible one through waterborne parasites that are explicitly described. The model is closely related to the SIP model introduced in [182]. In this first study we analyze in detail the properties of the mean-field version of the model, which aims to describe spatially homogeneous (i.e., well mixed) populations. The well-mixed approximation will be valid whenever the mean distance among hosts is smaller than the mixing length of the parasites before they get inactivated or absorbed. The model is written such that waterborne transmission is the only mechanism by which one infected immobile host can infect a healthy one and, thus, does not describe infection through direct contact. It is also assumed that the infected hosts, as invertebrates, do not have immune memory and that the probability that an infected individual recovers is small and can be neglected. Thus, the model is not adequate to study infection of highly aggregated mollusks (like some mussels) or other passive filters like corals, as for these hosts one should also include the possibility of infection through direct contact. A first very relevant question is whether the model, describing infection of immobile (sessile) hosts through waterborne parasites can be reduced to a simpler version in which the parasite compartment is not needed. One exact and two approximate reductions are presented. We believe that the model can be most useful in the rapid characterization of emergent marine epidemics if the right data from a well-mixed system are available.

A very timely case study of such emerging epidemics is the noble fan mussel (or pen shell) *Pinna nobilis*. Here we introduce and study in detail the properties of the mean-field version of a general compartmental model to study marine epidemics for bivalve populations, namely sessile, passive filter feeders invertebrate hosts that are infected through waterborne parasites. There are two main

hypotheses, the first one that a population-level description (i.e., without the consideration of spatial effects) is able to describe well the dynamics of the epidemic in a relatively dense population in small bounded regions. A second assumption is that the host becomes infected with some probability, but that there is not a critical parasite load in the infection process. After presenting the full SIRP model, then three different reductions are discussed: one exact, an approximate reduction of the former, and a third reduction based on a timescale approximation. The study is closed with a validation with the available experimental data for the infection process of *Pinna nobilis* kept in tanks. We wish to point out that, being a highly endangered and protected species, the reported data correspond to an *unintended experiment* that cannot be repeated, and maybe these data represent the only opportunity to estimate the fundamental parameters of the model. In addition, the setup in which the *Pinna nobilis* were kept in tanks, represents the ideal implementation of the conditions under which the mean-field model SIRP is valid.

## 3.2 The SIRP model

### 3.2.1 Model structure and initial considerations

In this work we analyze the SIRP model, a deterministic multi-compartmental mean-field model, continuous in time and unstructured in spatial or age terms, to study infection in bivalve populations. In particular, we stress that the model as it is written describes spatially-homogeneous populations. Compartmental models are the most frequently used class of models in terrestrial epidemiology [185], and originated in the classic study of SIR models by Kermack and McKendrick [65]. The use of compartmental models in the study of infectious processes in marine systems was quite rare until very recently [175]. As already advanced in Section 3.1, there are relevant features in the description of epidemic processes in marine ecosystems that are different with respect to the case of terrestrial ecosystems [180], and their study in marine environments is dominated so far by so-called proliferation models [181], which do not address the transmission of the pathogen. See [182] for a discussion of several compartment models for the study of marine epizootics.

Compartmental models of diseases in terrestrial ecosystems caused by microparasites (i.e., viruses, bacteria and protozoans) do not consider a compartment to describe the dynamics of the parasite [186], describing just the different stages of the host. Infection typically occurs in 2 ways: i) as a contact process, in which the microparasite is transmitted directly from an infected host,  $I$ , by contact or through air in proximity, to a susceptible host,  $S$ ; ii) through a vector that has acquired the microparasite by biting an infected host,  $I$ , and passes the

microparasite to a susceptible host,  $S$ . In the first case one can describe the infection process through some probability that the individuals come close, while in the second it is very relevant to describe the vector mobility, and at least 2 compartments, susceptible and infected vector, are typically needed. Once the microparasite enters the host, it proliferates inside it, so the infection process can be described by using compartments for susceptible individuals,  $S$ , infected individuals,  $I$ , and possible exposed individuals,  $E$ . In particular, transmission in terrestrial sessile organisms (e.g., plants) is generically vector-borne. In the case of marine ecosystems, infection typically occurs through water-borne parasites, in particular in filter-feeders sessile organisms, while vector-transmitted diseases are much less frequent. Parasites may be transported by diffusion, sea currents, or even active motion (i.e., if they have flagella). In any case, infection between sessile hosts is not through direct contact, but instead through the production and excretion of parasites by infected individuals and the assimilation by filtering of parasites by a healthy (susceptible) host. So, parasites are produced and excreted to the marine medium, in which they stay infective until they become deactivated (i.e., die) or are absorbed by hosts. In a way, in parasite-transmitted marine diseases, parasites have a dual role: they are not only agents that induce infection but also act as vectors that transmit disease from an immobile infected host to a susceptible one.

The SIRP model is a general mean-field compartmental model that describes epidemic transmission through water-borne parasites, which we think is specially adequate to describe epidemic transmission in sparsely located passive filter-feeders, like many bivalves. We exclude the case of colonies in which individuals are in proximity, e.g., mussels, corals, etc., in which direct contact could be relevant and should be included in the model. In the SIRP model hosts are described through 3 different compartments, as in the SIR model, which represent different evolution stages of the disease: a susceptible class of healthy individuals that may contract the disease,  $S$ , an infected class of individuals that may pass the disease through excretion of the parasite,  $I$ , and a class of removed (namely dead) individuals,  $R$ , which cannot be infected any more and that cannot transmit the disease, plus an extra compartment,  $P$ , for the parasite population in the medium. It is important to note that invertebrates do not develop long-term immunity in the mammalian sense [181], and so no compartment of individuals “recovered with immunity” is considered. However, bivalves have a first line of defense, with hemocytes being able to fight parasites and reduce their internal population. Nevertheless, available evidence indicates that the number of individuals that can achieve a full recovery is usually small and can be neglected, and so it is not necessary to consider a process in which individuals in the  $I$  compartment return to the  $S$  compartment at some rate, like

in the SIS model. Instead, the population's long-term response to disease, when it occurs, is through natural selection for genotypes characterized by increased resistance to or tolerance for the pathogen. As already advanced, the SIRP model includes a fourth compartment that represents the parasite population in the water medium, whose population needs to be described explicitly. An explicit compartment allows to model the situation in which the population of parasites may evolve dynamically in time, although in Section 3.2.3.3 we will consider the case in which the parasite population accommodates almost instantaneously to the infected host population, and the description of the parasite can be simplified.

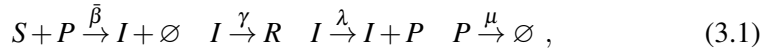
Infection occurs when the host enters in contact with the parasite in the marine medium. It involves a process of filtration of water by the bivalve, and although a detailed representation has been discussed in the literature [183], in the current SIRP model it is represented effectively. Infection is modelled through a nonlinear term, typical in compartmental models, but that now depends on the parasite population and not on the population of the infected compartment,  $I$ . In terrestrial epidemiology there are two alternative ways to represent infection [187]: i) mass action incidence, in which infection grows as the population gets larger,  $\beta SI$ ; ii) standard incidence, in which the infection is bounded as the population grows,  $\beta SI/N$ , where  $N$  is the total (host) population. One must look at these two choices as limit cases, with the possibility that in reality the system is best described by an intermediate form, closer to one of the limit cases, for example the modified SIR model in [188], in which the infection term has  $S+I$  in the denominator instead of the total population  $N$ , because the  $R$  compartment is removed. Modelling infection with an explicit representation of the parasite population encounters the same basic *dilemma* about whether the incidence grows as the (host) population increases or is bounded. We will model infection as  $\bar{\beta} SP$ , where the two possibilities are equivalent to the two different incidences just mentioned: i)  $\bar{\beta} = \beta$ ; ii)  $\bar{\beta} = \beta/N$ , where  $\beta$ , the disease transmission rate, is a *constant*, which can depend at most on external parameters, like temperature and salinity, but not on the variables defining the model (populations of host compartments or parasites). The model is valid considering both types of incidence, and in the case study we will see which incidence seems more adequate in this case.

### 3.2.2 General SIRP model

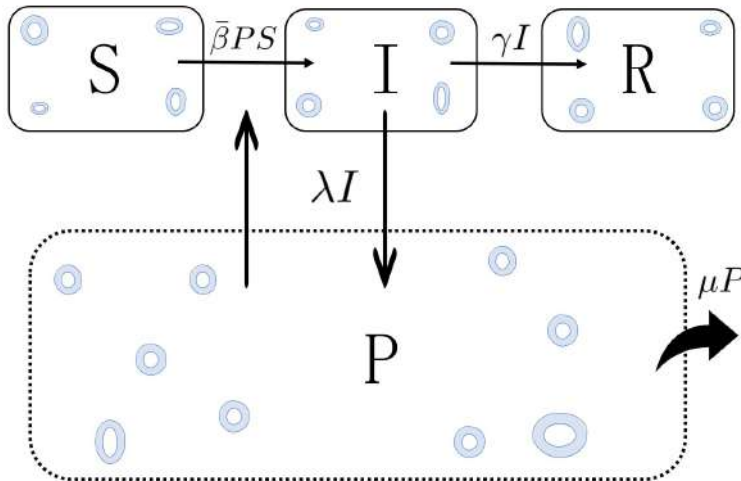
In this section we will write a mean-field compartmental model to describe epidemics of immobile (sessile) hosts in a marine medium that get infected by a water-borne parasite. Being a mean-field model implies that the model is compartmental and does not include an explicit space dependence, representing a

well-mixed system. The mean-field model describes a spatially homogeneous system, but we hope that it will be the basis for spatially inhomogeneous situations, by adding suitable terms accounting for the mobility of the parasite.

It is also assumed that the hosts become infected with some probability when exposed to a parasite, i.e., that there is not a critical parasite load needed for infection. The model is defined according to the following reaction processes,



which is graphically summarized in Fig. 3.1.



**Figure 3.1:** SIRP model flow diagram. The model variables are represented by capital letters: susceptible hosts ( $S$ ), infected hosts ( $I$ ), dead hosts ( $R$ ) and the population of parasites ( $P$ ). Arrows represent the processes in the model with their rates indicated next to them, and blue rings represent parasites. The flow follows the scheme in Eq. (3.1), that leads to the system of differential equations in Eq. (3.2).

According to the scheme in Fig. 3.1, we consider the host in 3 possible states: susceptible,  $S$ , infected by the parasite,  $I$  and removed (dead),  $R$ . Then we introduce the parasite population in the medium (water),  $P$ . In the model, the  $\bar{\beta}$ , the disease transmission rate parameter regulates the infection rate of susceptible hosts and accounts, among other mechanisms, for the parasite intake rate,  $\gamma$  the mortality of infected hosts, being the inverse of the typical mean time for an infected host to die;  $\lambda$  the production rate of parasites from infected hosts,

and  $\mu$  the inverse of the typical life time of the parasite.  $\mu$  can be related to several processes, like biological deactivation (or survival time) or other general losses, like dilution due to renewal of water in a closed experiment, natural losses in open ecosystems, or absorption by other filter feeders. We are not considering the possibility of spontaneous parasite gain, i.e. immigration, in this version of the model. A summary of the model parameters can be found in Table 3.1. We do not consider vital dynamics for the hosts, and this implies that the sum of the 3 host subpopulations is constant,  $N = S + I + R$ , as the time scale of the disease evolution is much faster than the typical life cycle of fan mussels. The model is similar to the SIP presented in [182], except for an extra term in  $\dot{P}$ ,  $-\bar{\beta}PS$ , accounting for the fact that when a parasite infects a host, it is absorbed by it. The conditions under which the SIRP model can be simplified to the SIP model are discussed in Section 3.3.5.

In order to build the deterministic model, we consider that the population is large enough to neglect fluctuations and that it is well mixed, so that spatial effects can be neglected. In this situation, we consider the infection process to be proportional to the number of parasites in the medium, so that the average number of contacts between susceptible fan mussels and the average parasite population is given by  $PS$ , and, thus, the change in the number of susceptible fan mussels takes the form  $\dot{S} = -\bar{\beta}PS$ , where the dot over a variable indicates a differentiation with respect to time:  $\dot{S} = dS/dt$ .

**Table 3.1:** Description of the SIRP model parameters.

Variable	Definition	Parameter	Definition
$S$	Susceptible host	$\beta$	Disease transmission rate
$I$	Infected host	$\gamma$	Host mortality rate
$R$	Removed host	$\lambda$	Production rate of parasites by infected hosts
$P$	Parasite in the medium	$\mu$	Parasite deactivation/dilution rate

Following this argumentation, the scheme in Eq. (3.1) and Fig. 3.1, one can write the evolution equations of the SIRP model,

$$\begin{aligned}
 \dot{S} &= -\bar{\beta}PS \\
 \dot{I} &= \bar{\beta}PS - \gamma I \\
 \dot{R} &= \gamma I \\
 \dot{P} &= \lambda I - \bar{\beta}PS - \mu P .
 \end{aligned}
 \tag{3.2}$$

Model (Eq. (3.2)) lives in the 4-dimensional (S,I,R,P) phase space, representing the variables the populations of individuals in the susceptible, infected, and removed host compartments and of parasites, respectively.

The fixed points of Eq. (3.2) are determined by the conditions<sup>1</sup>  $I = P = 0$ , to be fulfilled simultaneously. We will study the stability of the fixed point defined by  $S(0)$ ,  $I(0) = P(0) = 0$  and  $R(0) = N - S(0)$ . A linear stability analysis of this fixed point reveals that it has two null eigenvalues, that stem from the condition  $N = S + I + R$  and the conserved quantity of Appendix A.1. The first condition,  $S + I + R = N$ , implies that it is enough to consider two of the host populations, e.g.,  $S$  and  $I$ , as the third one can be obtained from the other two. The implications of the conserved quantity reported in Appendix A.1 are more subtle, as it implies that fixed points are not isolated, as it happens in ordinary dissipative dynamical systems, and there is an infinite number (a line of) fixed points for the final state of the epidemic, depending on the initial conditions. This also implies that the phase space is foliated by the conserved quantity,  $C$  of Eq. (A.5), and every initial condition,  $S_0$ , with a different value of  $C$  leads to a different asymptotic condition,  $S_\infty$ , just as shown in [49] for the SIR model (cf. Fig. 10.1 in *op. cit.*). The third eigenvalue, that is the largest of the two non-zero eigenvalues, can be positive if  $\beta S_0 \lambda > \gamma(\beta S_0 + \mu)$  and negative if the inequality is reversed, defining the conditional stability of the fixed point. The fourth eigenvalue is always negative, and all the eigenvalues are always real (Appendix A.2). The instability of the fixed point along the third eigenvalue drives the beginning of the epidemic.

An extremely important result in epidemiology is the so-called *basic reproduction number*,  $R_0$ , a dimensionless number that represents the number of secondary infections produced by a primary infection in a fully susceptible population.  $R_0 = 1$  defines the threshold for epidemic propagation: an epidemic will occur when  $R_0 > 1$ , and the number of infected individuals will grow at an exponential rate in the early phases of the epidemic [189], while if  $R_0 < 1$  the infection will wane naturally. This quantity can be formally obtained making use of the Next Generation Matrix (NGM) method [68, 190]. Applying this formal method to our system of ordinary differential equations (ODE's), one obtains the following relation for the basic reproduction number (Appendix A.3),

$$R_0 = \frac{\lambda}{\gamma \left( 1 + \frac{\mu}{\beta S(0)} \right)}. \quad (3.3)$$

The threshold condition provided by  $R_0$  (Eq. (3.3)) is equivalent to the linear stability condition for the third eigenvalue of the initial, pre-epidemic, fixed point, as  $\beta S(0)\lambda > \gamma(\beta S(0) + \mu)$  implies that this eigenvalue is positive and the disease-free equilibrium state unstable, being this equivalent to  $R_0 > 1$

<sup>1</sup>We do not consider the trivial fixed point  $S = I = R = P = 0$  that would imply  $N = 0$  and  $P = 0$  at all times.

(Appendix A.2). Thus, if  $R_0 > 1$  the fixed point is unstable, and an epidemic will ensue if infected hosts,  $I$ , or parasites,  $P$  appear in the system. An epidemic will propagate until the system reaches a stable fixed point, that signals the end of the epidemic (Appendix A.2).

### 3.2.3 Model reduction

The SIRP model lives in a 4-dimensional phase space and depends on 4 parameters, which makes difficult to confront it with experimental data. Thus, we will discuss here three alternative ways of reducing the model. The first involves an exact reduction of the model, based on the conserved quantity derived in Appendix A.1. The second reduction consists of an approximation to the previous exact reduction, that turns out to be equivalent to an exact reduction of a slightly simplified model (without the  $-\bar{\beta}SP$  term in the equation of  $\dot{P}$ ). The third one is based on an approximation valid if the system parameters fulfill certain conditions.

#### 3.2.3.1 Exact reduction of the SIRP model

From the conserved quantity derived in Appendix A.1, it is possible to write the parasite population in the SIRP model as a function of the host states as follows,

$$P(S, I) = -\frac{\lambda}{\gamma}(S + I) + \frac{\mu}{\beta} \ln(S) + S + C(0), \quad (3.4)$$

where  $C(0) = P(0) + \frac{\lambda}{\gamma}(S(0) + I(0)) - \frac{\mu}{\beta} \ln(S(0)) - S(0)$ .

Substituting Eq. (3.4) into the general SIRP model of Eq. (3.2) we obtain the following nonstandard SIR model,

$$\begin{aligned} \dot{S} &= \frac{\lambda \bar{\beta}}{\gamma} S(S + I) - \mu S \ln(S) - \bar{\beta} S^2 - S \bar{\beta} C(0) \\ \dot{I} &= -\frac{\lambda \bar{\beta}}{\gamma} S(S + I) + \mu S \ln(S) + \bar{\beta} S^2 + S \bar{\beta} C(0) - \gamma I \\ \dot{R} &= \gamma I. \end{aligned} \quad (3.5)$$

Although using the conserved quantity yields an exact reduction from a 4D dynamical system to a 3D one, the number of independent parameters and initial conditions remain unchanged, i.e., they still depend on 4 parameters and 4 initial conditions. Thus, although useful, (Eq. (3.5)) is not ideal when trying to fit experimental data, and this is the reason for trying a further approximation to Eq. (3.5) to be discussed next.

### 3.2.3.2 Further approximation to the exact reduction

A further approximation to [Section 3.2.3.1](#) that is less restrictive and expected to be valid in a broader parameter range than the one discussed in [Section 3.2.3.3](#) is possible. This approximation reduces the number of free parameters by one, which is useful in fitting available data. The approximation consists of neglecting the  $S$  term in [Eq. \(3.4\)](#), which is possible if  $\lambda/\gamma \gg 1$  and also  $\mu \ln N/(\bar{\beta}N) \gg 1$ , as  $S(t)$  decreases monotonically with time and is, at most,  $N$  at the initial time. Interestingly, this approximation is equivalent to the simplification of the equation for  $\dot{P}$  in ([Eq. \(3.2\)](#)) so that the  $-\bar{\beta}SP$  is skipped, what yields exactly the SIP model of Ref. [[182](#)]. This reduced model has an exact conserved quantity,  $\mathcal{C}$ , that differs from that of the SIRP model in the linear  $S$  term ([Appendix A.1](#)). Using this approximation one can write,

$$\begin{aligned}\dot{S} &= \frac{\lambda'}{\gamma}S(S+I) - \mu S \ln(S) - S\mathcal{C}(0) \\ \dot{I} &= -\frac{\lambda'}{\gamma}S(S+I) + \mu S \ln(S) + S\mathcal{C}(0) - \gamma I \\ \dot{R} &= \gamma I ,\end{aligned}\tag{3.6}$$

where  $\lambda' = \lambda\bar{\beta}$  and  $\mathcal{C}(0) = \bar{\beta}(P(0) + \lambda/\gamma(S(0) + I(0) - \mu/\bar{\beta} \ln S(0))) = \bar{\beta}\mathcal{C}(0)$  is a redefinition of the conserved quantity of the SIP model [Eq. \(A.9\)](#),  $\mathcal{C}(0)$ , a constant, such that it absorbs  $\bar{\beta}$  and all initial conditions of the model. The result is that [Eq. \(3.6\)](#) depends on 3 parameters and 1 constant, compared to [Eq. \(3.5\)](#) that depends on 4 parameters, facilitating, thus, the use of the model to fit experimental data.

### 3.2.3.3 Model reduction through fast-slow separation

The third reduction of the 4-D dynamical model [Eq. \(3.2\)](#) makes the assumption that the timescale in which the parasite population changes is faster than the one corresponding to the host. This means that pathogen deactivation in the medium must be faster than host mortality, which is a very reasonable assumption. In terms of the rates associated with each of these processes, this means  $\mu > \gamma$ . Taking  $\mu$  as common factor in  $\dot{P}$  one can write,

$$\varepsilon\dot{P} = \lambda I/\mu - \bar{\beta}SP/\mu - P\tag{3.7}$$

where  $\varepsilon = 1/\mu$  is small, as  $\mu$  is large. If furthermore  $\mu \gg \bar{\beta}N$  and  $\lambda \gg \beta P$  one arrives to,

$$P \approx \frac{\lambda}{\mu}I .\tag{3.8}$$

Under this approximation the slow subsystem can be written,

$$\begin{aligned}\dot{S} &= -\beta'IS \\ \dot{I} &= (\beta'S - \gamma)I \\ \dot{R} &= \gamma I ,\end{aligned}\tag{3.9}$$

that is equivalent to the classical SIR model with  $\beta' = \bar{\beta}\lambda/\mu$  instead of the infection rate  $\bar{\beta}$ . The reduced 3-D model Eq. (3.9) from the original 4-D SIRP model Eq. (3.2) depends on two parameters instead of 4 as the original model had one initial condition, e.g.,  $I(0) = N - S(0)$  if  $R(0) = 0$ , and is much more amenable to be applied to the analysis of experimental data, as shown in Section 3.4. Furthermore,  $\gamma$  could be eliminated through a time rescaling,  $t \rightarrow t' = \gamma t$  with a redefinition of  $\beta' \rightarrow \beta'' = \beta'/\gamma = \beta\lambda/(\mu\gamma)$ , leaving the model as a function of a single effective parameter. However, we will keep both  $\beta'$  and  $\gamma$  for convenience when fitting the model to experimental data in Section 3.4, as we would need to know anyhow  $\gamma$  in order to analyse the experimental data. The validity of this approximation is checked numerically in Section 3.3.

### 3.3 Numerical analysis

Due to the impossibility of solving the SIRP model analytically, in the present section we perform a numerical characterization of the model<sup>2</sup>. Moreover, we show the validity range of the performed approximations to reduce the SIRP model to an effective SIR model. We start our numerical analysis by investigating the relative influence of the model parameters on some epidemiological quantities of interest: the basic reproduction number ( $R_0$ ), related to the existence of an epidemic outbreak, continuing with the final state of the epidemic, given by the final number of dead individuals ( $R(\infty)$ ) and the maximum of infected individuals ( $I_{\max}$ ) together with the time at which it occurs ( $t_{\max}$ ).

In order to identify the most influential parameters of our model, a Sensitivity Analysis (SA) will be performed. SA can be divided into two classes: Local Sensitivity Analysis (LSA) and Global Sensitivity Analysis (GSA). LSA represents the assessment of the local impact of input factor variation on model response by concentrating on the sensitivity in the vicinity of a set of factor values. Such sensitivity is often evaluated through gradients or partial derivatives of the output functions at these factor values, such that other input factors are kept constant. Since epidemic models exhibit a threshold behavior, controlled

<sup>2</sup>All numerical simulations of the dynamical system Eq. (3.2) have been carried out using a Runge-Kutta 4th order method, with a temporal step  $\Delta t = 0.001$ . Numerically stable results are obtained with  $\Delta t \leq 0.01$ .

by the dimensionless quantity  $R_0$ , it is relevant to study its robustness with respect to small perturbations by means of the LSA explained above, as its analytical expression is known.

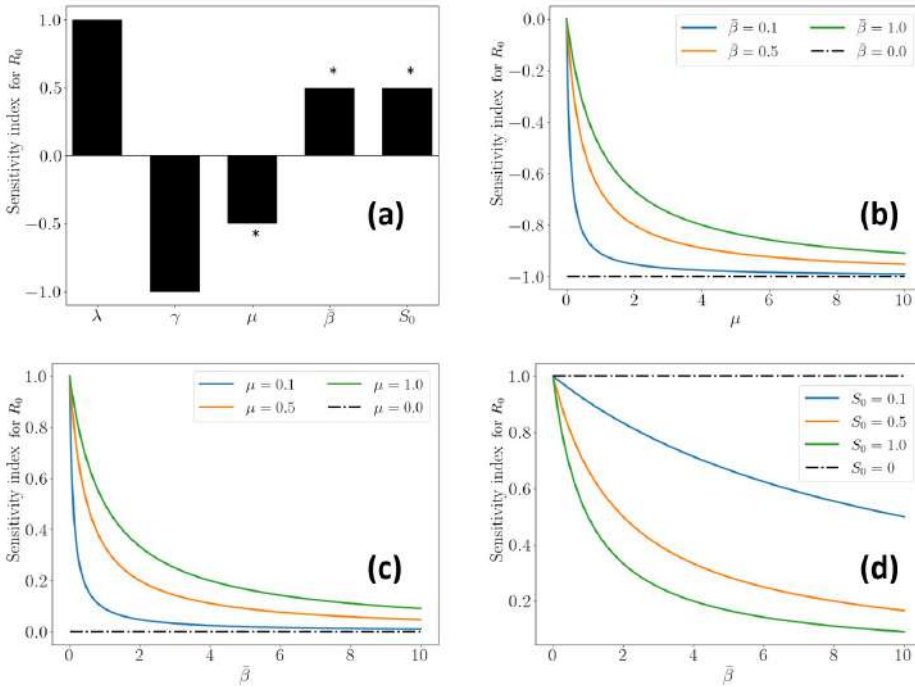
On the other hand, GSA will be applied to study the influence of the parameters in the final state of the epidemic and the epidemic peak by exploring a large domain of the parameter space. In turn, GSA is the process of apportioning the uncertainty in outputs to the uncertainty in each input factor over their entire range of interest. A sensitivity analysis is considered to be global when all the input factors are varied simultaneously and the sensitivity is evaluated over the entire range of each input factor, in clear contrast to LSA. Within GSA, first-order indices are a measure of the contribution to the output variance given by the variation of the parameter alone averaged over variations in other input parameters, while second-order indices take into account first-order interactions between parameters. While LSA is carried out analytically (if exact expressions are available), GSA is a purely numerical approach. Further mathematical details on Sensitivity Analysis can be found in [Appendix A.4](#).

For all the sensitivity analysis performed in the following sections, and in order to avoid ambiguities associated to the definition of  $\bar{\beta}$  as a function of  $N$ , we assume  $N = 1$ , so that both possible incidences yield  $\bar{\beta} = \beta$  and the numerical results are equivalent.

### 3.3.1 The basic reproduction number $R_0$

To study the relevance of parameters involved in an epidemic outbreak, a LSA was performed. We analyze the local sensitivity of  $R_0$  through the normalized sensitivity index, so that the function  $F(\vec{p})$  of [Eq. \(A.17\)](#) is substituted by the analytical expression of  $R_0$ , [Eq. \(3.3\)](#).

[Fig. 3.2 \(a\)](#) shows the sensitivity index for  $R_0$  for specific baseline parameters, where  $\lambda$ ,  $\bar{\beta}$  and  $S_0$  contribute to increase the basic reproduction number while  $\alpha$ ,  $\gamma$  and  $\mu$  contribute to decrease it, as expected. Moreover, we can see that  $\lambda$  and  $\gamma$  are the most influential parameters while  $\mu$ ,  $\bar{\beta}$  and  $S_0$  depend on each other. These dependencies cause varying influences on  $R_0$ , which are fully depicted in panels [Fig. 3.2 \(b-d\)](#). It can be seen that the influence of  $\bar{\beta}$  increases with the increase of  $\mu$  and the decrease of  $S_0$ . Similarly, the importance of  $S_0$  increases with  $\mu$  and decreases with  $\bar{\beta}$ . On the other hand, the impact of  $\mu$  increases with the decrease of both  $S_0$  or  $\bar{\beta}$ .

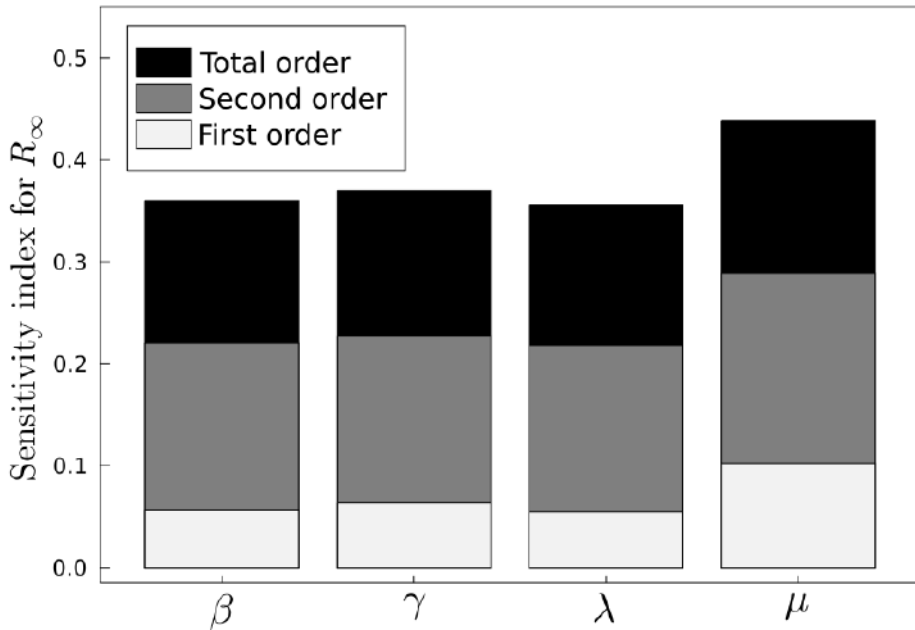


**Figure 3.2:** Sensitivity analysis of  $R_0$ . Panel (a): Local sensitivity analysis of  $R_0$  for the baseline parameters  $\lambda = 1$ ,  $\gamma = 1$  and  $\mu = \bar{\beta} = S_0 = 1$ . The asterisks mark parameters for which the sensitivity index is not constant, depending on, at least, another parameter. Panels (b-d): Local sensitivity analysis of  $R_0$  with respect to parameters with an asterisk, showing the different dependence with a second parameter and the effect on the varying sensitivity index.

### 3.3.2 Final state of the epidemic

Another important quantity in epidemiology is the final state of the epidemic, which can be characterized by the final number of dead individuals,  $R_\infty$ . Within our general SIRP model, it is not possible to find an analytical expression for  $R(t)$  so that we need to tackle the problem numerically. To this end, we perform GSA for the final number of dead individuals in order to determine the most influential parameters for this quantity. In particular, we apply the Sobol method, discussed in [Appendix A.4](#). The Confidence Interval, CI, obtained in our study is less than 1% of the index value, indicating a very high accuracy; therefore it is not shown in the figures. The results of the explained procedure are shown in [Fig. 3.3](#), where the total order (black), first-order (white) and second-order (gray) sensitivity indices for each of the model parameters are detailed. It can be observed that  $\mu$  has a slightly greater influence than the other parameters

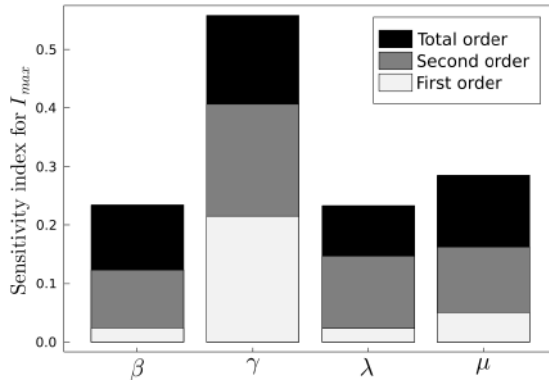
with respect to the final number of dead individuals. Note that the second-order indices are larger than the first-order ones for all the parameters, which indicates a high influence of the nonlinearities in our model, at least for the particular quantity under study.



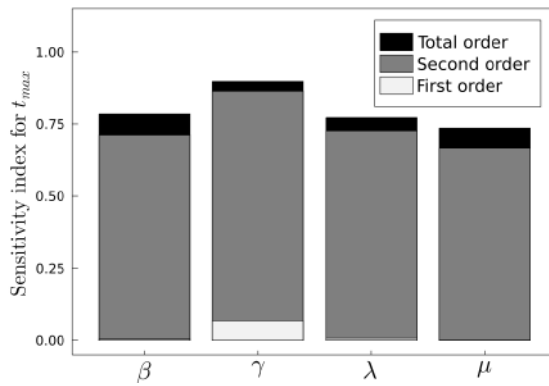
**Figure 3.3:** Sensitivity indices (LSA) for the final number of dead individuals ( $R_\infty$ ) for each one of the indicated parameters. The black bars represent the total order indices of sensitivity, while the white (gray) color represents the contribution of the first (second) order indices.

### 3.3.3 Maximum of infected individuals

A GSA of the maximum number of infected individuals,  $I_{\max}$  and the time it occurs,  $t_{\max}$  is performed to study the influence of the model parameters regarding these quantities. In this case, Fig. 3.4,  $\gamma$  has greater influence in the epidemic peak than any of the other parameters, while for the time at which the peak takes place, all the parameters have basically the same influence. Again, the second order indices (the first order interactions between parameters) account for most of the parameter sensitivity, in particular in the time of the epidemic peak, indicating the high degree of nonlinearity of this effect.



(a) Epidemic peak



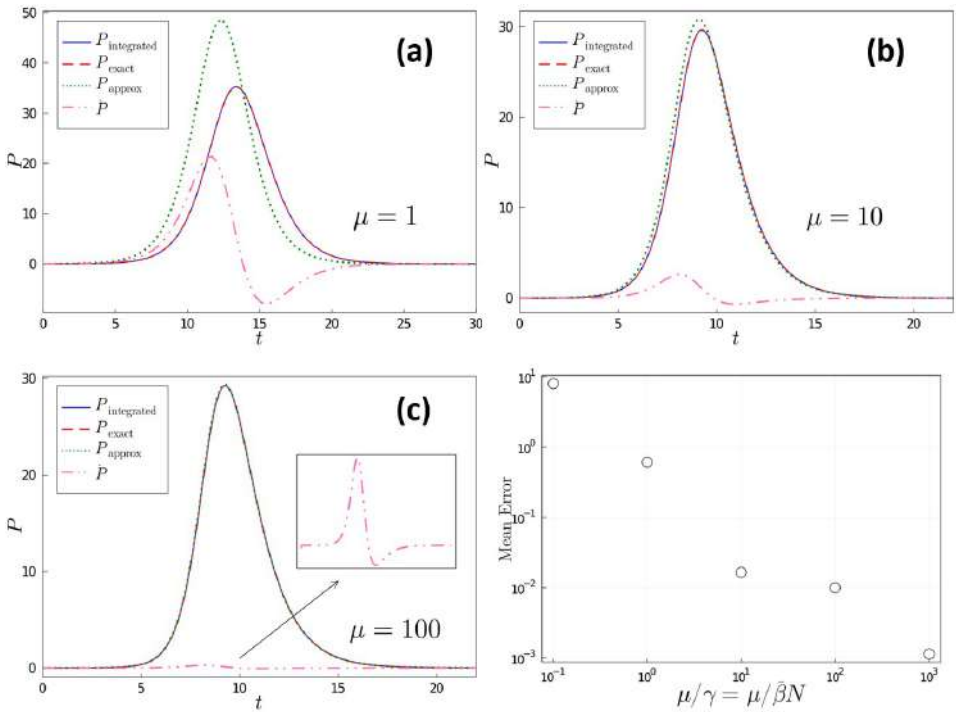
(b) Time of epidemic peak

**Figure 3.4:** Global sensitivity analysis for the maximum of infected individuals  $I_{max}$  (a) and its time occurrence  $t_{max}$  (b). The black bars represent sensitivity at all orders, while white (grey) color represents the contribution of the first (second) order indices.

### 3.3.4 Numerical verification of the fast-slow approximation

The parasite concentration approximation, based on the timescale separation discussed in Section 3.2.3.3, is now verified by computational means. The verification was performed using both mass action and standard incidence, but for the sake of simplicity we show only the results for the standard incidence case. Worth is to say that, mathematically, changing from standard incidence to mass action involves only a rescaling of the  $\bar{\beta}$  parameter, so that the numerical results are exactly the same. Fig. 3.5 contains a comparison for 3 different values of the parasite deactivation rate,  $\mu$ . It can be seen that the approximation is poor when  $\mu \sim \gamma, \bar{\beta}N$  Fig. 3.5 (a), as it could be expected. On the other

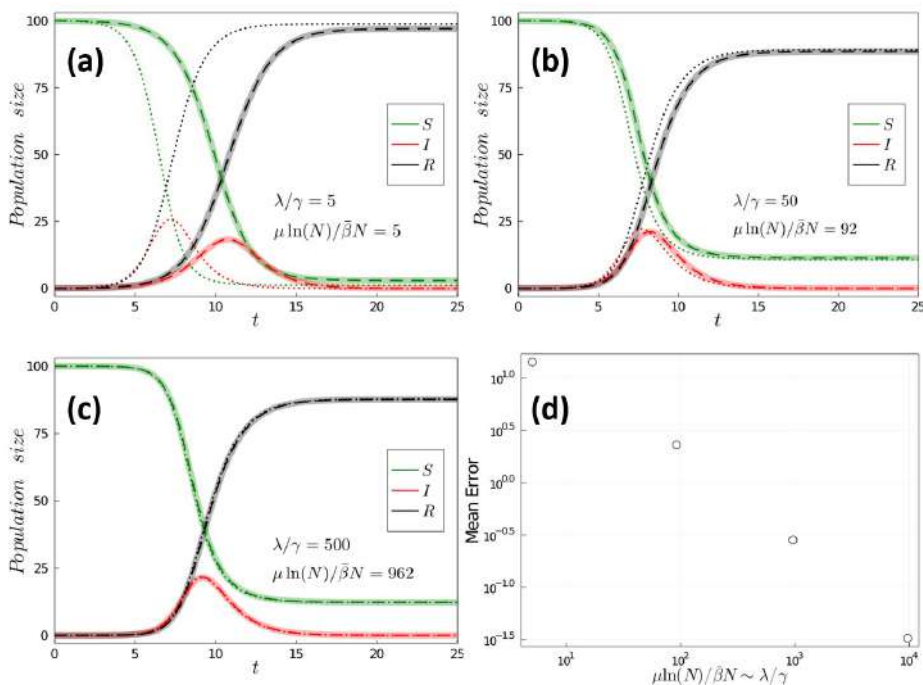
hand, the approximation is quite good when  $\mu$  is one order of magnitude larger than  $\gamma$  and  $\bar{\beta}N$  Fig. 3.5 (b), while it is extremely accurate when  $\mu$  is two orders of magnitude larger than  $\gamma, \bar{\beta}N$ , Fig. 3.5 (c). The figure also shows the numerical value of  $\dot{P}$  (pink dash dot), and it can be checked how it becomes smaller as  $\mu$  increases compared to  $\gamma, \bar{\beta}N$ , justifying the timescale separation of Section 3.2.3.3. Finally, Fig. 3.5 (a-c) also shows (dashed red line) the analytical value for  $P(S,I)$  derived in Eq. (A.7), that matches perfectly the result of the numerical integration of Eq. (3.2), as should be the case.



**Figure 3.5:** Numerical check of the approximate expression for the pathogen concentration, (Eq. (3.8)),  $\bar{\beta} = 1/50$  and  $\gamma = 1$ : (a)  $\mu = 1$ ; (b)  $\mu = 10$ ; (c)  $\mu = 100$ , while  $\lambda$  is varied to keep  $R_0 = 2.5$ , defined in (Eq. (3.3)) (with  $S_0 = N = 50$ ), i.e.,  $\lambda = 5, 27.5, 252.5$  respectively for (a)-(b)-(c), respectively. The blue solid line represents the numerically integrated quantity, the red dashed line (superimposed to the blue solid one as they are identical) is the exact solution for this quantity, (Eq. (3.4)) and the green dotted line accounts for the approximate expression from the timescale separation (Eq. (3.8)). The dash-dotted pink line represents the derivative of  $P$ ,  $\dot{P}$ , in the scaled time frame. Panel (d): Mean error between the approximate and exact solutions for increasing  $\mu/\gamma = \mu/\bar{\beta}N$ .

### 3.3.5 Numerical verification of the exact reduction

The numerical verification was performed for both mass action and standard incidence, but for the sake of simplicity in Fig. 3.6 we show only the results for the standard incidence case. First, and as it should be because it is an exact result, the exact reduction of the SIRP model discussed in Section 3.2.3.1 matches perfectly the numerical results obtained from the full model for all possible parameter values, Fig. 3.6 (a-c).



**Figure 3.6:** Numerical check of the exact model reduction along with the subsequent approximation shown in Section 3.2.3.1 with  $N = 100$ ,  $\bar{\beta} = 1/100$ ,  $\gamma = 1$ , (a)  $\lambda = 5$ ,  $\mu = 1.1$ ; (b)  $\lambda = 50$ ,  $\mu = 20$ ; (c)  $\lambda = 500$ ,  $\mu = 209$ .  $R_0 = 2.38$  for all the panels. The solid semitransparent lines represent the original 4D model, the dashed lines the exact reduction and the dotted lines the approximate model from the exact reduction. Panel (d): Mean error between the approximate and exact solutions for increasing  $\mu \ln(N)/\bar{\beta}N$  and  $\lambda/\gamma$  while  $R_0 = 2.38$  is kept constant.

Regarding the approximation to the exact reduction, one can see how the approximation converges to the exact solution as the parameters fulfil the conditions indicated in Section 3.2.3.2, namely that both  $\gamma/\lambda \gg 1$  and

$\mu \log(N)/\beta\bar{N} \gg 1$ , becoming very accurate if these ratios are larger than 1 by two orders of magnitude or more (Fig. 3.6 (c)). We recall that in this case the SIRP model converges to the SIP model of [182]. Conversely, the approximation is poor when any of these two ratios is of order 1 (Fig. 3.6 (a)), while Fig. 3.6 (b) presents the result in an intermediate case, in which the approximation is fair.

## 3.4 Model validation

In this section, the general SIRP model is validated against collected data from the *Pinna nobilis* Mass Mortality Event. As explained in Section 3.1, the disease is caused by the parasite *Haplosporidium pinnae* and the hosts, *P. nobilis*, are sessile bivalves endemic of the Mediterranean Sea. Thus, this epidemic is a perfect candidate to be described by the SIRP model. In the model, parasite production occurs only inside infected hosts, and parasites are released to the medium, either through their respiratory or digestive system. The simultaneous occurrence of the different possible stages of the parasite (uni- and bi-nucleate cells, multinucleate plasmodia, sporocysts and uninucleate spores) in the same host individual is not common among haplosporidans and makes *H. pinnae* different from previously known haplosporidan species [122]. The occurrence of uni- and binucleate stages suggest possible direct transmission from infected to healthy fan mussels, as observed in *B. ostreae* and *B. exitiosa* [191–193]. Additionally, the presence of spores (a dormant, resistant stage) could allow long persistence in the environment and the hypothetical involvement of an intermediate host as suggested for *H. nelsoni* and *H. costale* [194–196]. While uninucleate cells are always detected in infected fan mussels, sporulation has been only detected sporadically [122]. Thus, we assume that infection occurs mostly through uninucleate (or binucleate) cells by direct transmission (as the experimental observations in captivity point out, see [108]). We do not consider disease transmission through other stages. We do not consider spores, given the infrequent observation of spores and the current lack of experimental information about spore transmission (that could involve another intermediate host species). Regarding plasmodia and sporocyst stages, these stages are too large to be released through the epithelium. The distinction between uninucleate and binucleate cells seems unnecessary at this level of representation, as these phases only participate in parasite proliferation inside infected hosts, a process that we consider in an effective way. Finally, the evidence of the time course of the disease compared to the long life cycle of *P. nobilis* suggests host vital dynamics (i.e. recruitment (reproduction) and natural death) can be neglected.

After an epidemic outbreak that took place in Portlligat, in the north east of Catalonia, 215 *Pinna nobilis* individuals were extracted from their natural medium in order to be preserved as a genetic reserve in several controlled wa-

ter tanks of different institutions in Spain [108]. The institutions that participated in this preservation effort were IFAPA, IEO, IRTA, IMEDMAR-UCV and Oceanogràfic of Valencia. The original idea was to rescue the individuals before infection, however, the subsequent evolution of the rescued *Pinna nobilis* populations indicates that some individuals were already infected at the time of extraction (and/or in contact with some amount of the parasite transferred from sea water). This allowed the opportunity to use the data of the time evolution of the epidemic in the controlled water tanks, reported in [108], to evaluate the described SIRP model<sup>3</sup>. The empirical data consists of the proportion of survivors as a function of time in the controlled water tanks with a temporal resolution of one month. Despite the fact that the temperature of the water in the tanks was controlled, it was sharply lowered in most of the tanks when mortality started to appear within the population, as a last effort to keep the rest of the population safe and alive, since keeping the temperature below approximately 13.5°C is a known strategy to preserve *Pinna nobilis* individuals as disease expression is minimal [114]. Fortunately, two of the tanks kept its temperature approximately constant during the full recorded time. This is the case of the tanks in IFAPA in Huelva and the Oceanogràfic of Valencia (OCE), both Spanish institutes. These water tanks have been selected to validate our model, maintaining constant temperatures of 14°C and 17°C, respectively.

First we will fit the exact reduction of the SIRP model, assuming  $\mu \log(N) \gg \bar{\beta}N$  and  $\lambda/\gamma \gg 1$  as discussed in Section 3.2.3.2, namely Eq. (3.6). This reduced model depends on three parameters ( $\lambda'$ ,  $\mu$ ,  $\gamma$ ) and one constant,  $\mathcal{C}(0)$  (see Section 3.2.3.2), that is related to the initial conditions of the model. The order of magnitude of the mortality rate can be deduced from data, with an estimated value of  $\gamma \approx 1 \text{ month}^{-1}$ . We fix this parameter in order to give some biological information to our model prior to the computational fit. We focus on the  $R$  compartment, as it can be retrieved directly from data in [108]<sup>4</sup>. We use a box-constrained variant<sup>5</sup> of the well known BFGS optimization algorithm [197] with a common L2 loss function, also known as Residual Sum of Squares (RSS)<sup>6</sup>. By running this algorithm one observes that the optimal parameters tend to be the ones in the boundary of the box-constrained parameter space. Furthermore, if the box size is increased (or decreased) the optimal parameters

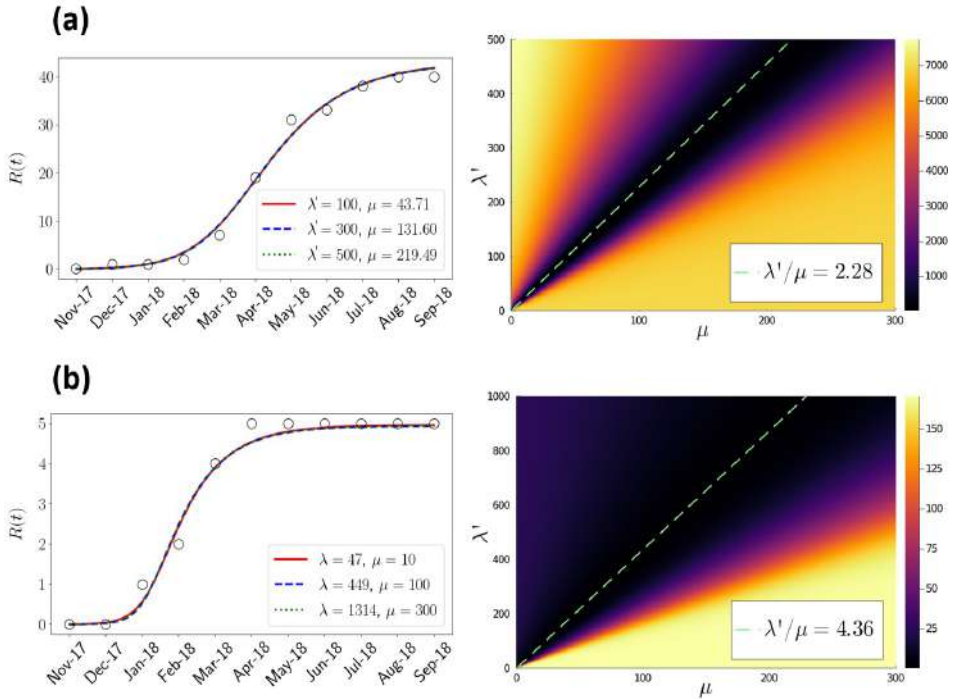
<sup>3</sup>Data use in this work with the purpose of validating and fitting parameters for the SIRP model have been taken from the Supplementary Information of [108]

<sup>4</sup>The number of dead individuals can be obtained as  $R = N - S$ , where  $S$  is the population of survivors and  $N$  is the total number of individuals in the tanks, 50 (IFAPA) and 5 (Oceanogràfic), respectively

<sup>5</sup>We constrain the optimization because the unconstrained optimization to the full range of the parameters, i.e, from 0 to  $\infty$  is not practical.

<sup>6</sup>The algorithm is implemented within the Julia high-level programming language [198] using the DifferentialEquations.jl package [199].

continue to be in the boundary of the box-constrained parameter space. This indicates that there exist several parameter combinations that optimally fit the data, and the combination parameters found by the optimization algorithm are only marginally optimal with respect to other parameter values. The locus (actually a valley) of marginal optimal parameters can be seen in the right-hand side panels of Fig. 3.7, where the cost function value of the optimization algorithm is plotted as heat map.



**Figure 3.7:** Parameter estimation for the approximation from the exact reduction of the SIRP model (Eq. (3.6)) using data from IFAPA (panel (a)) and OCE (panel (b)) water tanks, at 14°C and 17°C respectively. Left panels represent several fits of the model to empirical data of the number of dead hosts ( $R(t)$ ) using different optimal combinations of the parameters. Right panels are the RSS errors as a function of the input parameters, where the green dashed line represents the set of optimal combinations of the parameters with  $RSS = 60, 0.8$  for IFAPA and OCE, respectively.

Now we reach the point regarding the dilemma between mass action and standard incidence discussed in Section 3.2.1. If one does not correct the  $\bar{\beta}$  parameter with the size of the host population,  $N$ , that is equivalent to

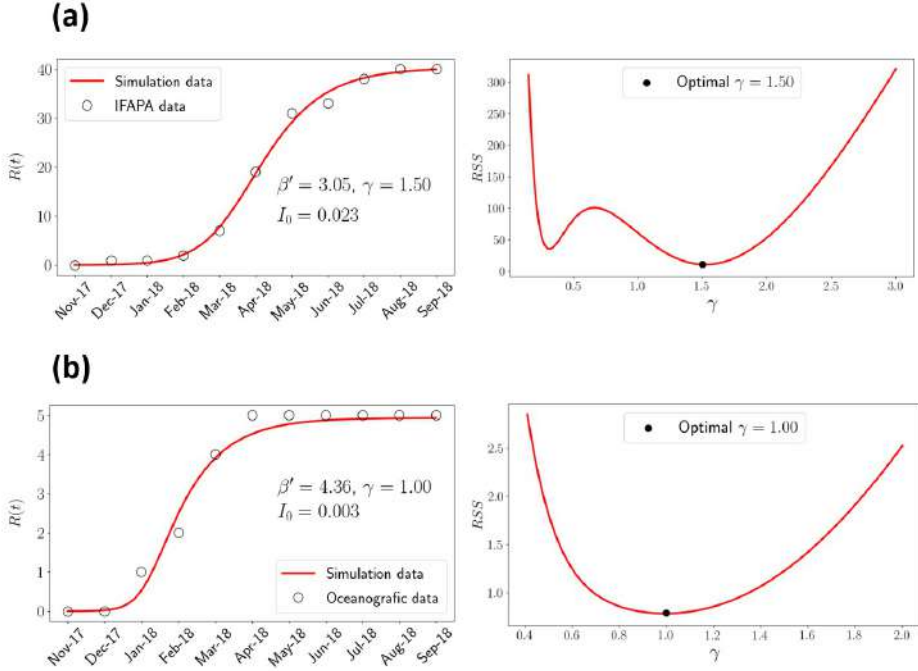
assuming the mass action incidence  $\bar{\beta} = \beta$ , the values that one would obtain for  $\beta' = \beta\lambda/\mu = \lambda'/\mu$  for both populations take disparate values in both tanks:  $\beta' = 0.046$  for the IFAPA data set and  $\beta' = 0.87$  for the Oceanogràfic (OCE) data set, a factor of 19 between them while their temperatures differ only by  $3^\circ\text{C}$ . These numbers indicate that the standard incidence is more reasonable, what amounts to choosing  $\bar{\beta} = \beta/N$ , where the final values of the reported parameter  $\beta'$  should be multiplied by  $N = 50$  for the IFAPA tank and  $N = 5$  for the OCE tank. The final result is then  $\beta' = 2.28$  and  $\beta' = 4.36$  for IFAPA and OCE tanks, that are the values reported in Fig. 3.7, implying that an almost twofold increase of the  $\beta'$  parameter corresponds to an increase of  $3^\circ\text{C}$ . This relation is in good agreement with the typical changes in rates of a wide range of organisms with a  $3^\circ\text{C}$  change in temperature, while a 19-fold change in the rate would imply at least a  $30^\circ\text{C}$  change in temperature (Appendix A.5).

The fact that there is an infinite number of combinations of the parameters that optimally fit the real data suggests that, as two parameters are slaved one to each other, that the model admits a further reduction. This reduction corresponds exactly to the approximate *SIR* model derived in Eq. (3.9), with the relationship  $\beta' = \lambda'/\mu$ , as anticipated. So, this gives further corroboration to the use of the *SIR* model Eq. (3.9) to fit  $\beta'$  as the free parameter (fixing the value of  $\gamma$  and with  $I_0$  as the initial condition determined by the fit). For consistency with the previous fitting we expect to obtain  $\beta' = 2.28$  and  $4.36$  as the optimal parameters for the IFAPA and OCE water tanks, respectively, and this is the case.

Interestingly, as reduced model Eq. (3.9) has fewer parameters to fit we can relax our initial assumption of  $\gamma = 1 \text{ month}^{-1}$  and check how the fit improves or worsens when varying  $\gamma$ . In Fig. 3.8 a fit of the reduced *SIR* model Eq. (3.9) is shown for the IFAPA (top) and Oceanogràfic (bottom) controlled water tanks<sup>7</sup>. Fig. 3.8 (c-d) shows the *RSS* error as  $\gamma$  is varied. It can be seen that for the IFAPA water tanks  $\gamma = 1.5 \text{ month}^{-1}$  yields more accurate results, while for the Oceanogràfic water tanks  $\gamma = 1 \text{ month}^{-1}$  remains optimum. This shows a decrease in the mean removal time  $1/\gamma$  for lower water temperatures, with the finite size errors inherent to the OCE tank (as  $N = 5$ ). In the left panels the simulated curve of dead individuals, *R* compartment, as a function of time for the optimal fitted parameters is confronted to the experimental data, showing a remarkable agreement. With the optimum values of  $\gamma$ , in the IFAPA tank (now with  $\gamma = 1.5 \text{ month}^{-1}$ ) a new value of  $\beta' = 3.05$  is obtained, implying a probably more reasonable ratio of 1.43 for  $\beta'$  in both tanks (it was 1.91 in the original fit). From the optimal parameters we obtain the basic reproduction number, since  $R_0 = \beta'/\gamma$  we have that  $R_0^{\text{IFAPA}} \simeq 2$  and  $R_0^{\text{OCE}} \simeq 4$ , clearly above

<sup>7</sup>The  $N$  correction corresponding to standard incidence has already been applied to these values.

the epidemic threshold.



**Figure 3.8:** Parameter fitting for the R compartment to model (Eq. (3.9)) using data from IFAPA (panel (a)) and Oceanogràfic (panel (b)). The left part of both panels of the figure shows the optimal fit of the model to empirical data with  $RSS = 10.9, 0.8$  for IFAPA and OCE, respectively. The right panels show the variation of the  $RSS$  error for some values of  $\gamma$ . The  $\beta'$  values have been obtained assuming a standard incidence, as explained in the main text.

Summarizing, the SIRP model is able to fit two sets of experimental data, agreeing with a standard incidence, according to which the infection rate depends on the amount of parasites per pen shell individual. *Pinna nobilis* individuals in the IFAPA experiment were actually distributed in 4 tanks, and the standard incidence is compatible with this experimental aspect. The temperature dependence of the fitted parameters in this range ( $14 - 17^\circ\text{C}$ , appears to be compatible (although experiments at different temperatures would be needed) with an Arrhenius dependence of the infection parameters, also known as Boltzmann-Arrhenius [45, 200], that can be extended to account for unimodal dependence on temperature, with a maximum infectivity at a characteristic temperature for the parasite [200]. Therefore, we can assume that global

change (or temperature shifts) is expected to have complex effects on infectious diseases, causing some to increase, others to decrease, and many to shift their distributions [201]. In the particular case of pen shell mortality, our model results suggest the proposed mechanism of lower disease expression at lower temperatures. This might have direct consequences for the development of the mortality event and offers a bleak perspective for the future and specifically in the eastern Mediterranean basin, where the mortality was observed later due to current patterns, but average temperatures tend to be higher than in the western part of the Mediterranean.

### 3.5 Conclusions

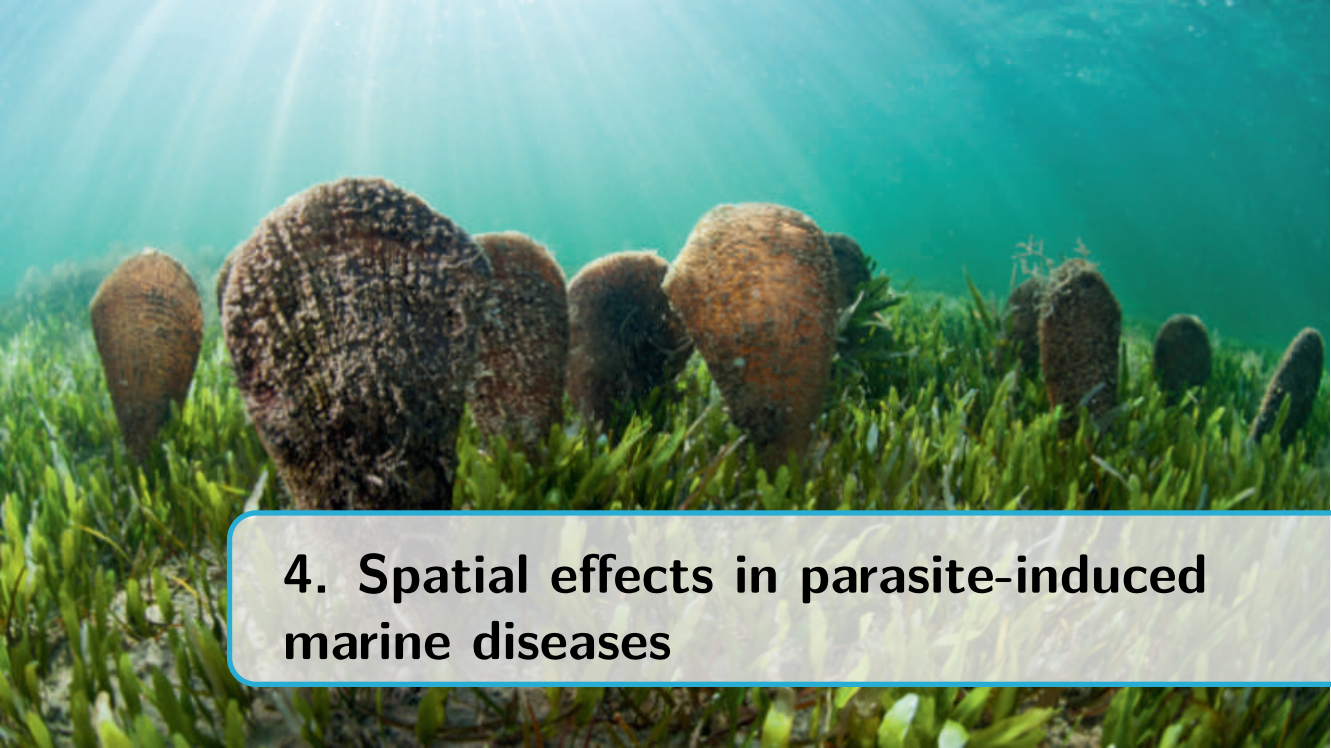
In this work we have analyzed a compartmental model to study marine epizootics for sessile hosts assuming infection by direct transmission through waterborne parasites. Moreover, we have used data from the recent mass mortality event of *Pinna Nobilis* in the Mediterranean Sea as a case study to validate our model. Compartmental models are routinely used in the study of disease infection and propagation in terrestrial ecosystems, including the study of the current Covid-19 pandemic (see, e.g., [189]). However, these models are starting to be used only recently in the study of marine epizootics [182], while proliferation models have been the most popular in the field [181]. A reason for the low popularity of compartment models in the study of marine epizootics is that there are some aspects in its modelling that differ from the now standard application to terrestrial ecosystems [180]. An important difference is that, in principle, (micro)parasites need to be modelled explicitly in marine ecosystems, while often they are not included in the description in terrestrial ecosystems [186].

The SIRP model has 4 compartments and depends on 4 parameters, so that it is not quite amenable to theoretical analysis. At the same time, due to the large number of parameters of the model, using it to analyze experimental observations can be cumbersome in practice if the parameter values are unknown. Nevertheless, we have shown three reductions of the model, one exact and two approximate ones, that can be useful to overcome these limitations that are typically present at the first stages of emergent epidemics. Indeed, the timescale approximation is able to fit the collected data of our case study for some optimal parameters, as shown in Section 3.4. This approximation is particularly useful as it only depends on 2 parameters, the death rate of infected hosts,  $\gamma$  and an effective infection rate,  $\beta'$ . Although this approximation simplifies the fitting procedure, there is a price to be paid in this analysis. The infection parameter,  $\beta$ , and the parameters regulating proliferation,  $\lambda$ , and deactivation/dilution of the parasite,  $\mu$ , become entrained into a single effective parameter,  $\beta'$ . Thus, the full understanding of the different effects at play in the system requires fur-

ther work. Furthermore, we have shown that an epidemic model for immobile hosts can be reduced to the standard SIR model, which assumes direct contact among the hosts, i.e., that the hosts are mobile. This reduction is only valid when the timescale of the parasites is much faster than that of the hosts, i.e.,  $\mu \gg \beta N, \gamma$ . Thus, our work provides a ground to apply the SIR model in marine epidemics of sessile hosts that fulfil the required conditions.

In a world with many possible new epizootics, we believe that our reduced model can be specifically useful to understand key features of those emerging diseases characterized by the spreading of waterborne parasites in a relatively fast way, provided that the temporal evolution of the disease can be determined for, at least, some set of individuals. Thus, some key parameters can be fitted to the available experimental data as shown in [Section 3.4](#). Still, the fitted relevant parameters may need to be supplemented with further information or targeted experiments. We hope that this approach can be useful in understanding emerging diseases in shellfish species of economic not only ecological value, and also, with suitable modifications, in aquaculture. It is noteworthy that our case study is a haplosporidan waterborne parasite. In fact, waterborne haplosporidans have been responsible for some of the most significant and consequential marine disease epizootics on record and are considered the major pathogens of concern for aquatic animal health shellfish industries around the world [125]. The SIRP model is the simplest model that one could think of having in mind its practical application, but could be extended to incorporate further effects that are so far described effectively.





## 4. Spatial effects in parasite-induced marine diseases

**Published as:**

À. Giménez-Romero, F. Vazquez, C. López, and M. A. Matías, “Spatial effects in parasite-induced marine diseases of immobile hosts”, [Royal Society Open Science](#) **9**, 212023 (2022)

## 4.1 Introduction

Wildlife emergent infectious diseases represent a substantial threat to ecosystems and the conservation of their biodiversity [29]. Their effects can be devastating at the ecological level, causing local extinctions [29] and in some cases pushing endemic species to the verge of extinction, as is the case of *Pinna nobilis* [114]; at the economic level, producing losses in agriculture, livestock, and aquaculture [202–204], and impacting human health, as is the case of the COVID-19 pandemic [205]. For the past decades, parasites have been continuously emerging [206, 207], while globalization and climate change have contributed to their evolution. This has allowed these parasites to enter new ecological niches and spread further the diseases they produce [208]. In particular, marine infectious diseases are recently increasing due to these and other anthropogenic pressures, like pollution and overfishing [172], inducing widespread mass mortality in several species [110, 209, 210].

An important subset of marine organisms affected by infectious diseases are sessile (i.e., they cannot move), like bivalves, sponges, or corals. An increasing number of outbreaks affecting marine mollusks have been reported, some of them causing mass mortality in commercially important bivalves [211]. Mainly due to the economic importance of some species (e.g., oysters), infectious diseases in bivalve populations have been deeply studied [196, 212–214]. Recently, deterministic compartmental models have been used to describe parasite transmitted diseases in marine sessile bivalves [183, 184, 215], showing to be able to accurately predict disease transmission in some circumstances. The main limitation of these compartmental models is the assumption of a non-spatial description of the system under study. This underlying hypothesis assumes that any pair parasite-host of the system can interact at any time, which is unrealistic in general. A non-spatial description assumes well mixed populations, which implies that the mean distance among hosts is smaller than the typical distance explored by parasites in their lifetime. This assumption can be quite realistic in some situations, as it is in [215], where the hosts were kept in tanks with water renovation. However, a non-spatial model is not expected to yield a good description of spatially extended hosts in a natural setting.

The key quantity in mathematical epidemiology is the basic reproductive number,  $R_0$ , that represents the number of infected individuals generated in one generation by the appearance of a single infected individual in a fully susceptible population. Thus,  $R_0 > 1$  ensures the onset of an epidemic, as the number of infected individuals will grow exponentially, producing a widespread disease [178]. If we first disregard spatial effects and assume a non-spatial description,  $R_0$  can be obtained from standard methods, like the Next Generation Matrix method [68], and will only depend on *intrinsic* characteristics of the pathosystem

(host-pathogen system) under study. However, this basic reproduction number is unable to characterize the threshold behavior in many situations, including spatially extended systems [216–218]. In these systems, the propagation of an epidemic to the entire system needs that a certain spatial threshold is exceeded [219]. Otherwise, the disease will only take place in suitable localized parts of the system, not being able to propagate to the total system. Thus, disease spread will be strongly affected by the host spatial distribution and pathogen mobility, which are not accounted for in non-spatial models.

In this work, we will try to unravel the transmission mechanisms of a parasite-induced disease affecting immobile hosts in a spatially extended system. We will approach the problem both theoretically and through numerical simulation. The numerical study is based on Individual-Based modeling (IBM), a method widely used to study ecological systems [220], so that individuals are treated as discrete entities, space is introduced explicitly, and the dynamics are stochastic. Representative average behaviors can be obtained by averaging over a sufficient number of realizations, and the accuracy of the approach can be calibrated by deriving the corresponding non-spatial limit, that can be confronted with the suitable compartmental model on which a particular IBM is based. The IBM approach to our problem will allow studying in depth the relation between pathogen mobility and immobile host infection. As parasites move randomly over the space, tracking the position of each parasite at different times is of fundamental importance to properly capture the stochastic dynamics of infections from parasites to hosts. Modeling parasites and hosts as individual entities allows taking into account the spatial and temporal heterogeneity of interactions between them. This heterogeneity and the level of control in microscopic interactions cannot be captured by other mathematical approaches, such as partial differential equations. On the other hand, IBMs are mathematically involved, and analytical treatments are normally cumbersome, while their numerical implementation is computationally expensive [221].

Here we introduce a spatially explicit individual-based model to study parasite-induced marine diseases of immobile hosts. The model is applied to the case of diffusing parasites and uniformly distributed hosts. The system under study is an extension of the compartmental model presented in [215]. As a main result, we find that the occurrence of an outbreak will depend on the balance between the intrinsic characteristics of the pathosystem, well represented by the above described non-spatial basic reproductive number,  $R_0$ , and features that characterize parasite mobility. We generalize the basic reproductive number, which we will refer to as  $\tilde{R}_0$ , such that it accounts for the number of hosts that get infected by the appearance of a single infected individual in a fully susceptible population in a spatially extended system.  $\tilde{R}_0$  characterizes the global epidemic

and can be written as a product between  $R_0$  and a factor describing parasite mobility. The latter factor is smaller and at most equal to 1, which implies that, as it could be expected, it is more difficult to induce a global outbreak in a spatially extended system (a two-dimensional lattice in our case) than in a well mixed (non-spatial) population.

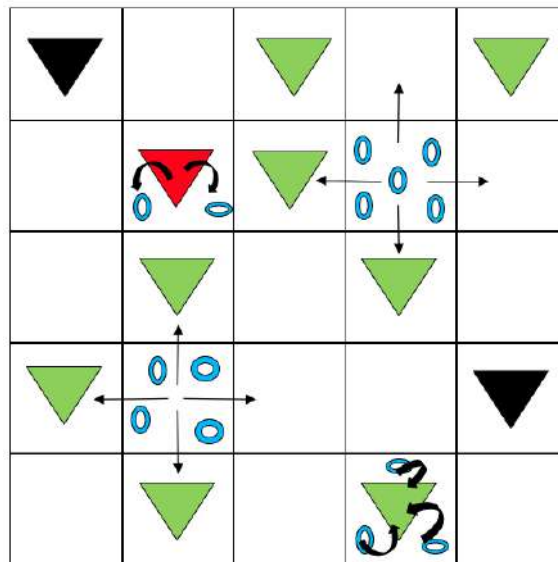
## 4.2 The SIRP spatial model

The most important biological features of the system under study are as follows: First, hosts are immobile, while the disease is transmitted by parasites produced by infected hosts. There are two mechanisms by which parasites are cleared from the medium: i) they have a finite lifetime after which they die or get inactivated; ii) they get absorbed after they infect a host and thus are no longer in the medium and cannot infect other hosts. Recruiting (birth) of hosts occurs at a very slow rate compared to other timescales in the system, and accordingly it will be considered negligible in the model. Moreover, hosts do not show long-term immunity, as is typical of invertebrates, like mollusks [181]. We also assume that recovery (healing) of infected hosts, if it occurs, can be neglected. Furthermore, we consider that dead hosts are not a source of parasites in the medium. See Chapter 3 [215] for a detailed presentation of the non-spatial SIRP model, including these biological modeling considerations.

Under these considerations, we introduce an individual-based model with explicit space characterization to study the effect of parasite mobility in disease transmission. We consider a square grid of length  $L$  with periodic boundary conditions and place a single host per site, so that there are  $N = L^2$  hosts. The hosts can be in three discrete states: susceptible,  $S$ ; infected,  $I$  and dead (or removed),  $R$ . Then, we introduce the parasite population as a new individual with a single state,  $P$ . Hosts are sessile (i.e., immobile), while parasites are allowed to move between the lattice sites. As an initial condition, we assume that the entire host population is susceptible,  $S(0) = N = L^2$ , and that a small initial number of parasites,  $P(0)$ , is introduced in the system.

Infection occurs when susceptible hosts filter parasites in their proximity. Accordingly, the infection process is implemented between parasites and susceptible hosts sharing the same lattice site. In particular, susceptible hosts in contact with a parasite become infected at rate  $\beta$ . The infection event implies the filtering of a parasite by a susceptible host so that when a new infection occurs, a parasite of that particular site is removed. Infected individuals die at rate  $\gamma$  and produce parasites at rate  $\lambda$ , while parasites die at rate  $\mu$ . Parasites move randomly between the four neighboring lattice sites at rate  $\kappa$ , which corresponds to a diffusive motion. Table 4.1 contains the definitions of the variables and parameters of the model, and Fig. 4.1 shows a schematic representation of

the dynamics.



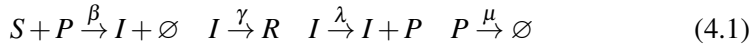
**Figure 4.1:** Scheme of the individual-based model. Green, red, and black triangles represent susceptible, infected, and dead hosts, respectively. Blue rings represent parasites, which move randomly between cells. Susceptible hosts get infected by filtering parasites, while infected hosts produce them. Dead hosts do not participate in the dynamics of the system.

**Table 4.1:** Variables and parameters of the individual-based SIRP model.

Variable/Parameter	Definition
$S$	Susceptible host
$I$	Infected host
$R$	Dead host
$P$	Parasite
$\beta$	Parasite-host transmission rate
$\gamma$	Host mortality rate
$\lambda$	Production rate of parasites by infected hosts
$\mu$	Parasite natural death rate
$\kappa$	Parasite dispersal rate (mobility)
$R_0$	Non-spatial basic reproductive number
$\tilde{R}_0$	Spatial basic reproductive number

Formally, the model is mathematically described by a system of  $N$  master equations for the probabilities of the states in each lattice site  $i$ . Eq. (4.1)

summarizes the reactive events. This is very difficult to manage analytically, so the time evolution of the model is numerically solved using Gillespie's algorithm [71] (the code can be found in [222]).



## 4.3 Results

In this section, several features of the model are studied, both numerically (from IBM simulations) and analytically. All numerical results were obtained for a square lattice of length  $L = 100$ , with  $N = S(0) = 10^4$  hosts and using a small initial condition of  $P(0) = 50$  parasites in the center site.

### 4.3.1 Non-spatial limit

An important test of the IBM implementation is to show that, under suitable circumstances, it converges to the non-spatial model on which the IBM is based. This occurs in the limit when the parasites move many times before dying or infecting a susceptible host. In this situation, each parasite typically visits all the hosts of the system and may infect any of them. This is equivalent to infecting a random host of the system, which happens with probability  $\beta S/N$ , being  $S$  the total number of susceptible hosts in the system. This corresponds to standard incidence, which represents the most accurate description of the infection process (Chapter 3, [215]). An equivalent picture is that parasites will end up uniformly distributed in the lattice, so that there will be  $P/N$  parasites in each lattice site at any time. One expects to reach these conditions when  $\kappa \gg \mu, \beta$ , and thus the system as a whole can be described by the following system of ordinary differential equations (ODE's),

$$\begin{aligned} \dot{S} &= -\beta PS/N, \\ \dot{I} &= \beta PS/N - \gamma I, \\ \dot{R} &= \gamma I, \\ \dot{P} &= \lambda I - \beta PS/N - \mu P, \end{aligned} \quad (4.2)$$

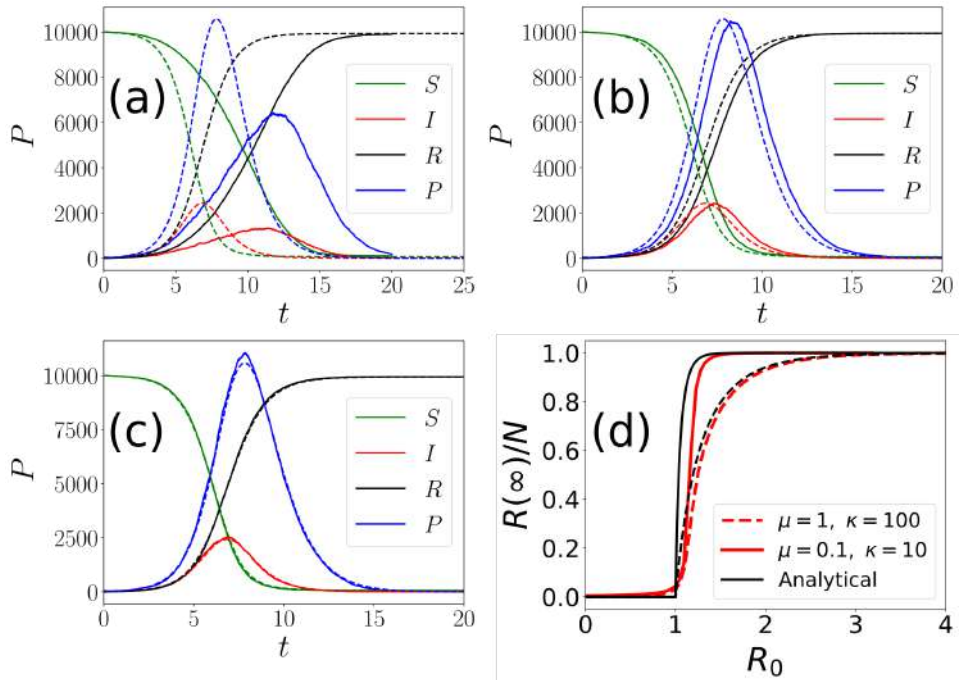
that is precisely the SIRP non-spatial model [215], where  $S$ ,  $I$ , and  $R$  are the total number of susceptible, infected, and recovered hosts in the system,  $P$  the total number of parasites, and  $N$  is the number of hosts.

The basic reproduction number,  $R_0$ , of this non-spatial model is the dimensionless quantity that yields the number of secondary infections generated by the appearance of a single infected individual in a completely susceptible population, also indicating whether the system will exhibit an epidemic outbreak,  $R_0 > 1$ , or not,  $R_0 < 1$ . In our case it can be directly computed as the mean

number of parasites produced by an infected host during its mean lifetime,  $\lambda/\gamma$ , times the mean number of susceptible hosts that get infected by parasites during their mean lifetime,  $\beta/(\mu + \beta)$ ,

$$R_0 = \frac{\lambda}{\gamma} \frac{\beta}{\mu + \beta} . \quad (4.3)$$

This result can be corroborated with standard methods such as the Next Generation Matrix method [68] (see [215]), where  $S(0) = N$  has been considered.



**Figure 4.2:** Numerical solution of the non-spatial model (Eq. (4.2), dashed lines) compared with numerical solutions of the individual-based model (solid lines) approaching the non-spatial limit with fixed  $\gamma = \mu = \beta = 1$  and  $\lambda = 6$ . (a)  $\kappa = 10^2$ , (b)  $\kappa = 10^3$ , (c)  $\kappa = 10^4$ . Panel (d) shows the final fraction of dead hosts,  $R(\infty)/N$ , as a function of  $R_0$  for  $\kappa/(\mu + \kappa) = 0.999$  with  $\mu = 1, 0.1$  compared to the analytical result.

Moreover, the model has a conserved quantity  $\mathcal{C}$  [215] that allows to find an analytical expression for the final number of dead individuals (Appendix A.6),

$$R(\infty) = N + \frac{S(0)}{\xi} W_0 \left( -\xi \exp \left( -\frac{\beta}{\mu} C \right) \right) , \quad (4.4)$$

with  $\xi = S(0) \frac{\beta(\lambda - \gamma)}{\mu\gamma}$  and  $C = P_0 + \frac{\lambda}{\gamma}(S(0) + I(0)) - S(0)$ .

The non-spatial limit of the model has been evaluated by comparing realizations of the stochastic model (in the limit  $\kappa \gg \mu, \beta$ ) with numerical solutions of the non-spatial ODE system of Eq. (4.2). Furthermore, the analytical expression for  $R(\infty)$  using the non-spatial model, Eq. (4.4), is also compared to the numerical results of the individual based model. As shown in Fig. 4.2 (a-c), as  $\kappa$  is increased compared to  $\mu$  the individual-based model approaches the non-spatial one. Fig. 4.2 (d) shows how the numerical results for  $R(\infty)$  for different  $R_0$  values approach the analytical solution in the non-spatial limit.

### 4.3.2 Approximate relation between parasites and infected hosts

In the limit  $\kappa \gg \beta, \mu$  a timescale approximation can be performed so that the parasite population dynamics directly relates to that of the infected hosts. In the non-spatial limit, it was already shown in [215] that, if  $\mu \gg \beta, \gamma$  and  $\lambda \gg \beta P/N$ , the total parasite population of the system can be well described using the approximation (see [215] for a detailed discussion),

$$P(t) \approx \frac{\lambda}{\mu} I(t) . \quad (4.5)$$

Here we extend the validity of this approximation to spatial systems far from the non-spatial limit. Consider the local dynamics of the parasite population on a lattice site  $i$ . Note that when the host in the site is susceptible, parasites in this site can either infect the host, die, or move to another site. All these processes imply that a parasite will disappear from the current site. Once the host at site  $i$  gets infected, infection can no longer occur, whereas parasite production is now possible. If  $\kappa$  is small enough compared to  $\lambda$  and  $\mu$ , the only competing processes in sites with infected individuals will be the production of parasites and their natural death, which can be fairly described by the following rate equation,

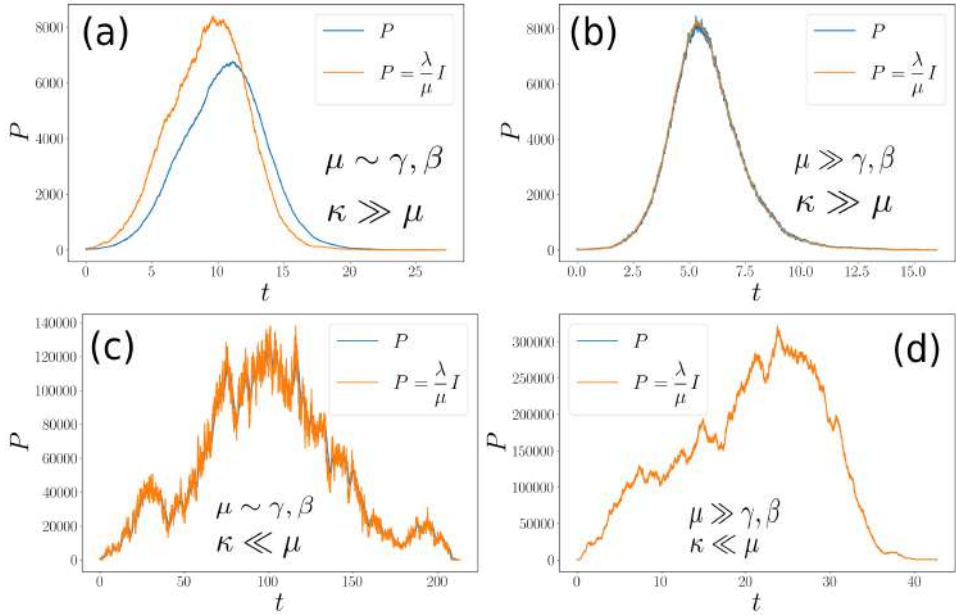
$$\frac{dP}{dt} = \lambda - \mu P , \quad (4.6)$$

with solution

$$P(t) = \frac{\lambda}{\mu} + \left[ P(0) - \frac{\lambda}{\mu} \right] e^{-\mu t} . \quad (4.7)$$

From Eq. (4.7) one may notice that the stationary value of  $P$ ,  $\lambda/\mu$ , is reached in a time proportional to  $t_{\text{eq}} \propto 1/\mu$ . This derivation allows us to find a condition for which Eq. (4.5) is valid beyond the non-spatial limit. Basically, if the mean dispersal time,  $1/\kappa$ , is greater than the equilibrium time,  $t_{\text{eq}} \propto 1/\mu$ , the parasite population in sites with infected hosts will reach its stationary level

before parasites enter or leave the sites. Thus, sites with infected hosts can be considered a closed system and the approximation holds. In other words, if the dispersal rate of parasites is small compared to the parasite deactivation rate,  $\kappa \ll \mu$ , the local parasite population of the site will reach its stationary level  $P_i = \lambda/\mu$ . It is possible to extend the result to the entire system: if there are  $I(t)$  infected sites in the system at time  $t$  and  $\kappa \ll \mu$  is fulfilled, there will be a total parasite population of  $P(t) = (\lambda/\mu)I(t)$ , which is equivalent to Eq. (4.5).



**Figure 4.3:** Numerical verification of the approximate expression for the parasite population dynamics, Eq. (4.5), for different mobility conditions. The simulations were performed fixing  $\beta = \gamma = 1$  for all panels. (a)  $\mu = 1$ ,  $\kappa = 10^2$ ,  $\lambda = 6.06$ ,  $\kappa/(\mu + \kappa) = 0.99$ ; (b)  $\mu = 100$ ,  $\kappa = 10^4$ ,  $\lambda = 306$ ,  $\kappa/(\mu + \kappa) = 0.99$ ; (c)  $\mu = 1$ ,  $\kappa = 0.01$ ,  $\lambda = 1200$ ,  $\kappa/(\mu + \kappa) = 0.01$ ; (d)  $\mu = 100$ ,  $\kappa = 1$ ,  $\lambda = 60600$ ,  $\kappa/(\mu + \kappa) = 0.01$

Thus, for the non-spatial limit ( $\kappa \gg \mu$ ) we have that if  $\mu \gg \beta, \gamma$  Eq. (4.5) is valid, while for  $\kappa \ll \mu$  the approximation is also valid regardless of the value of  $\beta, \gamma$ , as the nature of the approximation is different. Thus, in general, as  $\kappa$  decreases over  $\mu$  (the lower the parasite mobility becomes), we expect the approximation to work better.

The parasite approximation to infected host dynamics, Eq. (4.5), is numerically verified for different mobility conditions. Fig. 4.3 (a-b) shows how the approximation improves as  $\mu$  grows over  $\beta, \gamma$  (mean errors are 0.18 and 0.0081,

respectively) in the non-spatial limit, i.e.,  $\kappa \gg \mu$ , as expected. This result is in perfect agreement with that found in [215]. Then, Fig. 4.3 (c-d) show that the approximation is valid in general when  $\kappa \ll \mu$  but improves anyway when  $\mu \gg \beta, \gamma$  (mean errors are 0.04 and 0.0026, respectively). Summarizing, we see that the lower the value of  $\kappa$  is with respect to  $\mu$  the more valid Eq. (4.5) is, regardless of the value of  $\beta, \gamma$ , while in the non-spatial limit,  $\kappa \gg \mu$ , the condition  $\mu \gg \beta, \gamma$  is needed.

### 4.3.3 Spatial threshold

One of the main questions in epidemiology is to define the conditions under which an epidemic outbreak occurs, which usually is translated into the existence of a threshold. In a well-mixed (non-spatial) system, the basic reproduction number ( $R_0$ ), that characterizes this threshold  $R_0 = 1$ , can be defined exclusively from *intrinsic* parameters of the pathosystem, as the host-pathogen interaction does not depend on the host spatial structure or pathogen mobility (see Eq. (4.3)). In stochastic spatial models this formulation of  $R_0$  breaks down. First, in stochastic models, even above the threshold, there is a non-zero probability that the disease is unable to propagate initially, given by  $P_{\text{outbreak}} = 1 - (1/R_0)^{I^{(0)}}$  [223]. Furthermore, the discrete nature of the populations also modifies the estimates of  $R_0$  [224]. On the other hand, the introduction of space changes completely the nature of epidemic outbreaks, modifying the host-pathogen interactions by means of specific host spatial distributions and pathogen mobility patterns. Even if the basic reproduction number of the non-spatial model is above the threshold ( $R_0 > 1$ ), if parasite mobility is not large enough, the epidemic will stay locally confined. Thus, one expects that the threshold at which an epidemic outbreak can propagate to the rest of the system will depend on the balance between the intrinsic pathosystem parameters in  $R_0$  and parasite mobility, defining a spatial basic reproduction number,  $\tilde{R}_0$ .

Having in mind the study in Section 4.3.1, we expect that in the high mobility limit the basic reproduction number is defined by the non-spatial formula, Eq. (4.3). On the other hand, the lower the parasite mobility is, the more difficult it will be for a local outbreak to propagate through the system. Thus, it is natural to think of a spatial basic reproduction number of the form  $\tilde{R}_0 = R_0 f(\kappa)$ , where  $f(\kappa)$  is an increasing function of the parasite dispersal rate accounting for parasite mobility fulfilling  $\lim_{\kappa \rightarrow \infty} f(\kappa) = 1$ .

Indeed, some authors recently showed that the spatial basic reproduction number can be defined as the product between the non-spatial value,  $R_0$ , and a factor accounting for spatially-dependent interactions,  $f(r)$ , in the form  $\tilde{R}_0 = R_0 f(r)$  [225–227]. However, these expressions are not analytical [225, 226]

or are not directly related to pathogen mobility [227]. Here we propose a simple expression for the spatial basic reproduction number regulating the spatial propagation of the epidemic,

$$\tilde{R}_0 = \frac{\lambda}{\gamma} \frac{\beta}{\mu + \beta} \frac{\kappa}{\mu + \kappa} = R_0 \frac{\kappa}{\mu + \kappa}. \quad (4.8)$$

The derivation of Eq. (4.8) accounts for the number of secondary parasites that are able to produce new infections, or equivalently, the number of secondary infections produced by an initial infected host. If we consider an initial infected individual, on average it will produce  $\lambda/\gamma$  parasites. Then, these parasites can only move to neighboring sites or die, so that the dispersal probability is given by  $\kappa/(\mu + \kappa)$ . Finally, considering that parasites do not affect each other trying to infect the same host, the infection probability is given by  $\beta/(\mu + \beta)$ . Joining all terms, we finally obtain Eq. (4.8). This expression is valid when parasites move only to sites with susceptible individuals and do not try to infect the same host. Thus, the derived  $\tilde{R}_0$  is only an approximation to the spatial basic reproduction number for the case of an initial introduction of a small quantity of parasites in a fully susceptible population.

Note that, as expected, the spatial basic reproduction number is nothing other than the basic reproduction number of the non-spatial model multiplied by an increasing function of the parasite mobility,  $\kappa/(\mu + \kappa)$ . Taking the limit  $\kappa \gg \mu$  in Eq. (4.8) (non-spatial limit), the basic reproduction number of the non-spatial model is recovered. Conversely, in the limit of very low mobility, the  $\kappa/(\mu + \kappa)$  factor is small, and this has to be compensated with a large value of the non-spatial basic reproduction number,  $R_0$ , in order that there is an outbreak, i.e.,  $\tilde{R}_0 > 1$ .

The spatial threshold,  $\tilde{R}_0 = 1$ , given by Eq. (4.8), has been numerically checked by computing the phase diagram between the absorbing phase  $R(\infty) \approx 0$  (no infection, i.e., disease-free state) and the active phase  $R(\infty) > 0$  (in which some level of infection has occurred, i.e., propagation phase) for several values of the parasite mobility and the basic reproduction number of the non-spatial model,  $R_0$  Eq. (4.3). The transition is expected to occur at  $\tilde{R}_0 = 1$ , implying from Eq. (4.8) that the dependence of the critical value of  $R_0$ , say  $R_0^c$ , is expected to take the form,

$$R_0^c \sim \left( \frac{\kappa}{\mu + \kappa} \right)^{-1}. \quad (4.9)$$

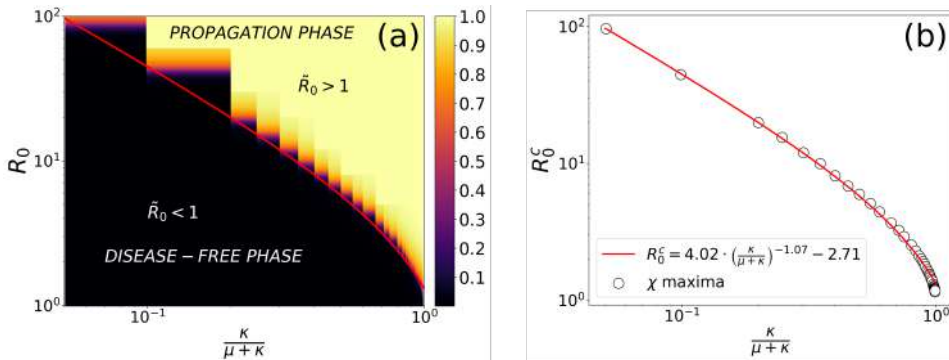
As discussed above, we expect that  $\tilde{R}_0 = 1$  with  $\tilde{R}_0$  given by Eq. (4.8) does not represent exactly the spatial threshold, and for this reason we suggest the more

general functional form,

$$R_0^c \sim A \left( \frac{\kappa}{\mu + \kappa} \right)^{-B} - C, \quad (4.10)$$

to be fitted to numerical data, where  $A = 1$ ,  $B = 1$  and  $C = 0$  would imply a perfect agreement of numerical simulations of the IBM model with Eq. (4.8).

In order to obtain the phase diagram, we compute the absorbing state of the model as an average over 1000 realizations for each value of the mobility and  $R_0$  considered. Then, the critical value  $R_0^c$  is computed for each mobility value as the  $R_0$  value for which the fluctuations of the “order parameter”  $\chi = \langle R(\infty)^2 \rangle - \langle R(\infty) \rangle^2$  are maximal, as this would be an indication of a transition between the disease-free and the propagation phases.



**Figure 4.4:** (a) Phase diagram showing the transition between the disease-free phase and the propagation phase for several values of the parasite mobility and  $R_0$ . The color code represents the fraction of dead individuals (i.e.,  $R/N$ ) in the final state of the epidemic computed by the average over 1000 realizations. (b) Fit for the transition line following Eq. (4.10), where dots are the maxima of the “order parameter” fluctuations,  $\chi = \langle R(\infty)^2 \rangle - \langle R(\infty) \rangle^2$

Fig. 4.4 (a) shows the numerical results of the computed transition between the disease-free and propagation phases. The heatmap coding represents the average value of absorbing state  $\langle R(\infty) \rangle$  for several values of the mobility factor and  $R_0$ . As expected, the lower the mobility factor is, the higher the value of  $R_0$  is needed for the disease to invade the population. Fig. 4.4 (b) shows the fit of Eq. (4.10) with less than a 1% relative error. Interestingly, we obtain  $B = 1.07 \approx 1$  which validates our expression for the spatial threshold as a first approximation. However, the values for  $A = 4.02$  and  $C = 2.7$  show a significant deviation from Eq. (4.9) and indicate that the expression Eq. (4.8) is an approximation to the

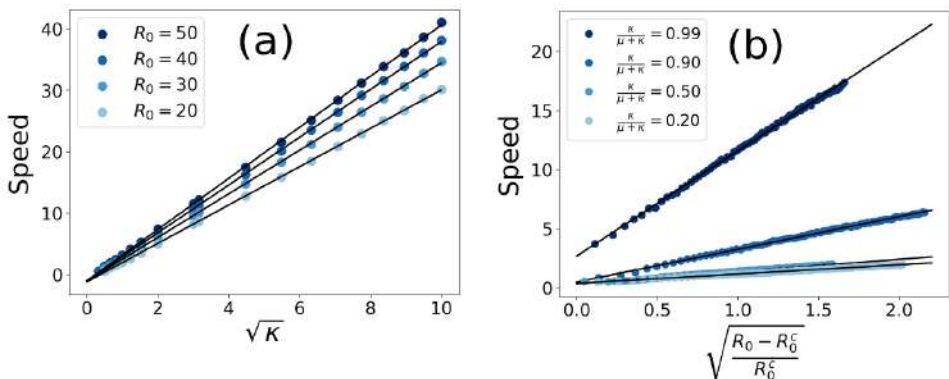
spatial basic reproduction number, which however seems to contain the right dependence on  $\kappa/(\mu + \kappa)$ , and where  $A$  could be a geometric factor for a lattice.

#### 4.3.4 Spreading speed of the infected population and time to extinction

Another relevant epidemiological question is how does an infected population spread after the onset of an epidemic. In order to obtain this spreading speed, we computed the mean time needed for an infected individual to reach the boundary of the system. More specifically, for each particular choice of the model parameters, 1000 simulations were run for several system sizes ranging from  $L = 10$  to  $L = 60$ . The computed mean time was found to depend linearly on the system size, thus allowing to compute the speed from the slope of this relation. With this procedure, the spreading speed was computed for several values of the parasite mobility and  $R_0$ , large enough to ensure an epidemic outbreak that reached the boundary of the system. In this situation, the spreading speed is expected to depend linearly on the square root of the parasite mobility,

$$v \sim \sqrt{\kappa}. \quad (4.11)$$

Fig. 4.5(a) shows this square root dependence for different values of the fixed  $R_0$ . Similarly, the speed was also computed for several values of the basic reproduction number and a fixed mobility. In this case, it varies with the square root of the distance to the critical value of  $R_0$ ,  $R_0^c$ , as shown by Fig. 4.5 (b). This is in good agreement with other mathematically similar models [228].



**Figure 4.5:** (a) Disease spreading speed as a function of the square root of the parasite mobility for several values of  $R_0$ . The plot shows a remarkable agreement with Eq. (4.11). (b) Disease spreading speed as function of the square root of the distance to the critical value of  $R_0$  for several values of the parasite mobility.

The extinction time is defined as the time elapsed from the beginning of the epidemic until the system reaches its absorbing state, that is, when no parasites or infected individuals are left. From Eq. (4.11), we expect the time to extinction to increase when the parasite mobility is decreased. Moreover, we expect the extinction time to decrease with the distance to the epidemic threshold, as we expect to reach the absorbing state faster for larger values of the spatial basic reproduction number.

In the limiting case where all (or almost all) hosts die, it is clear that the disease must have spread to the entire system. Thus, in this limit, the extinction time should be proportional to the inverse of the disease spreading speed,  $t_{\text{ext}} \sim 1/v$ . Then, in this limit, we can relate the extinction time with the parasite mobility as follows,

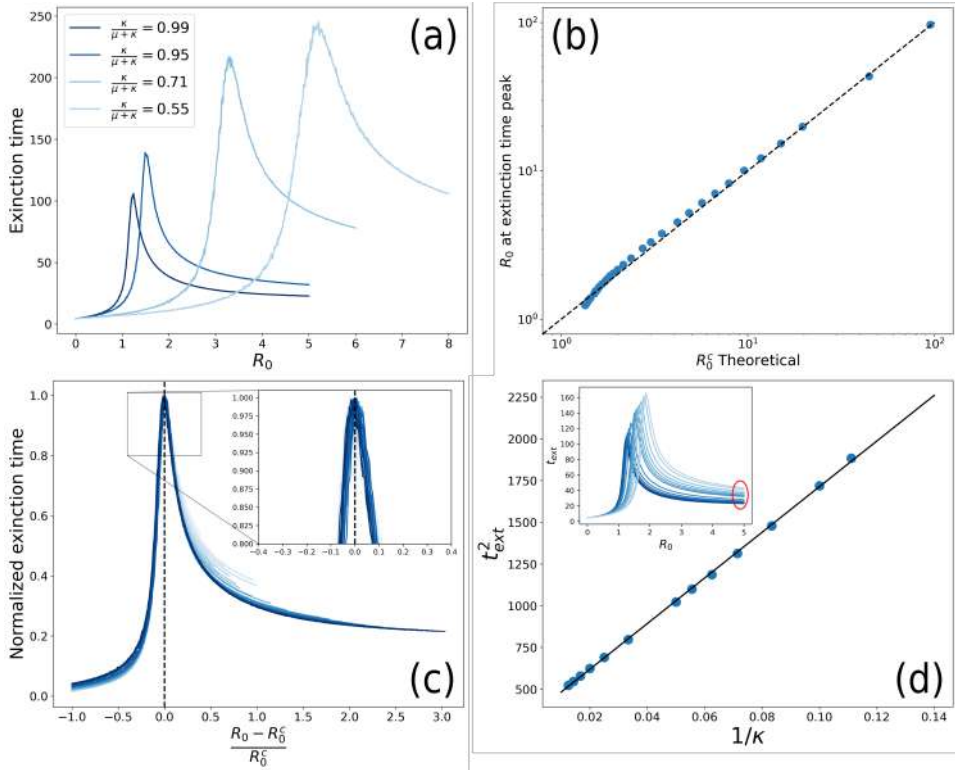
$$t_{\text{ext}} \sim \frac{1}{\sqrt{\kappa}} \quad \text{for } \tilde{R}_0 \gg 1. \quad (4.12)$$

However, the absorbing state is not always reached after all hosts become infected and Eq. (4.12) is only expected to work far from the epidemic threshold, when the disease is expected to spread to the entire system.

In Fig. 4.6 (a), the extinction time is plotted against the basic reproduction number for some values of the parasite mobility. As expected, the extinction time increases for lower values of the parasite mobility. The increasing behavior before the peak can be understood as the increasing time needed for the initial perturbations to decay to the disease-free phase. After the peak, the greater the basic reproduction number, the faster the epidemic will reach its absorbing state with a non-negligible number of dead individuals. So, with this interpretation, the peaks of the extinction time should coincide with the epidemic threshold for each value of the parasite mobility.

In Fig. 4.6 (b), we compare the numerical value of  $R_0$  at which the extinction time peaks with the theoretical value of  $R_0^c$ , computed with Eq. (4.10), showing good agreement. Thus, the dependence of the extinction time with  $R_0$  should vanish if plotted against the distance to  $R_0^c$ . Furthermore, if the extinction time is normalized (dividing each line by its maximum), all the lines should collapse near the transition point. In Fig. 4.6 (c), the normalized extinction time is plotted against the distance to the critical value of  $R_0$ . The scaling is shown to be valid only near the transition point, as expected.

In the limiting case where the epidemic dies by infecting a large part of the host population, i.e., for a large enough  $R_0$  value, the extinction time should follow Eq. (4.12), as previously discussed. In Fig. 4.6 (d), we show how the extinction time relates to the parasite mobility in this limit, following the predicted behavior.



**Figure 4.6:** (a) Extinction time for some values of the parasite mobility. (b) Comparison of the critical  $R_0$  value computed with Eq. (4.8) compared to the values obtained numerically by computing the maximum of the extinction time. (c) Scaling of the extinction time with several values of the parasite mobility. (d) Representation of the square of the extinction time as a function of the inverse of the parasite mobility. The inset shows the zone where this relation has been computed, showing a good agreement with Eq. (4.12).

## 4.4 Conclusions

In this chapter, we have developed a spatially explicit individual-based model for parasite-produced marine epidemics of immobile hosts. This study has allowed us to tackle important questions in marine epidemiology, such as how spatial constraints affect epidemic spreading in filter-feeder populations or how the infected population of hosts change in space and time. While addressing the aforementioned questions, we have shown that there exists a regime of high parasite mobility where the time progression of both host and parasite populations can be well described by the non-spatial version of the model (i.e., the system of ODE's presented in [215]). We have also shown that a fast-slow approximation

for the time progression of the parasite population, already presented in [215], can be extended for spatial systems. Interestingly, the conditions under which this approximation is valid are less restrictive than in the non-spatial case, and regimes in which this approximation is valid for low mobility and comparable time scales are reported in this contribution.

We have derived an approximate analytical expression of the *spatial* basic reproduction number able to predict the onset of a global epidemic in a spatial model. The obtained expression explicitly shows a trade-off between the intrinsic pathosystem dynamics (i.e.,  $R_0$ ) and a factor accounting for parasite mobility. Moreover, the spatial threshold defined by  $\tilde{R}_0 = 1$  separates the final state of the system in two different phases, namely a disease-free phase and a propagation phase. In the propagation phase, any initial condition of infected individuals or parasites will propagate throughout the system, causing a proper outbreak. On the other hand, in the disease-free phase, the conditions are not sufficient for a local introduction of parasites or infected individuals to spread through the system. The effect of the parasite mobility in the spatial basic reproduction number is clear: the more parasites move, the more infections they cause.

The spatio-temporal behavior of the system has been investigated in the propagation phase. First, we showed that the infected population spreads through the space with a speed directly proportional to the square root of the diffusion coefficient of parasites, showing good agreement between the derived analytical expression and numerical simulations. The time to extinction has also been studied by means of numerical simulations, showing that, if the system is far above enough of the spatial threshold, the time to extinction can be analytically computed, in good agreement with simulations. We obtained that larger values of the parasite mobility yield more severe epidemics in which there are more infections and the extinction is faster.

To summarize, in the present work we have introduced and analyzed an individual-based approach to epidemic transmission in spatially extended systems of immobile hosts. The infection mechanism is due to mobile parasites, which are in turn produced by infected hosts. The study allows to answer some biologically relevant questions, like predicting the occurrence of a global epidemic outbreak or its velocity of expansion through the system. Thus, the analytical and computational results of the model shed light on the underlying mechanisms underpinning the emergence of a global epidemic outbreak and its spatial progression. This work provides a first step into the spatial-explicit, individual-based modeling of marine epidemics of immobile hosts.

Although this work has considered the case of a spatially homogeneous distribution of hosts, we plan to extend the study to more general cases, discussing the effect of inhomogeneous spatial host distributions. Furthermore, other bio-

---

logically relevant effects could be added to the model to enhance the description of different epidemics, e.g., infected individuals could still filter parasites or the infection process could depend on the parasite-load. The model could also describe epidemics on other immobile species, such as filter-feeders like sponges or other bivalves, corals, intertidal communities or starfishes, provided that the necessary modifications in the model are properly included. Stochastic spatially explicit descriptions like the one presented here could also be extended to the study of epidemics of other immobile hosts, like vector-borne diseases of plants. However, this would imply a quite different model to describe the different epidemic compartments of the vectors and also their ecological features. We hope these studies can be useful in conservation plans or ecosystem management and could serve as a basis for more sophisticated models.





# Modeling vector-borne plant diseases

<b>5</b>	<b>Non-stationary vector populations</b> .....	<b>109</b>
5.1	Introduction .....	110
5.2	The model .....	112
5.3	Results .....	114
5.4	Conclusions .....	124
<b>6</b>	<b>A compartmental model for <i>Xylella fastidiosa</i> diseases</b> ...	<b>127</b>
6.1	Introduction .....	128
6.2	Materials and Methods .....	130
6.3	Results .....	136
6.4	Discussion .....	143




## Summary

Vector-borne plant diseases are a significant threat to agriculture, with the potential to cause widespread epidemics, food shortages, and economic losses. In particular, the bacterium *Xylella fastidiosa* has emerged as a major concern for the agricultural sector, affecting a wide range of crops worldwide. Despite extensive research, many aspects of the dynamics of vector-borne plant diseases remain poorly understood. For instance, the dynamics of the vector population, which play a crucial role in disease transmission, is often neglected in existing models. This is particularly important for *Xylella fastidiosa* diseases, where the vector population exhibits complex non-stationary dynamics that are not captured by traditional models. In this part, we focus on the modeling of vector-borne plant disease in which the vector population follows complex non-stationary dynamics. We develop a theoretical framework for modeling vector-borne diseases with growing or decaying vector populations. We show that traditional methods to predict the onset of an epidemic do not apply in this context, propose new approaches, and demonstrate that these dynamics can have a significant effect on the temporal patterns of disease spread, leading to unexpected outcomes that are not captured by traditional models. This work has enabled the construction of a model for *Xylella fastidiosa* diseases that explicitly incorporates the complex dynamics of the vector population. We validate the model using empirical data, demonstrating its predictive power and practical utility. Finally, we provide insights into the design of effective control strategies that take into account the dynamics of the vector population. We address current gaps in our understanding of how non-stationary vector dynamics influence disease spread and severity, offering new insights into the management and control of these impactful plant diseases.

### Objectives

- To develop a theoretical framework for modeling vector-borne diseases with non-stationary and non-periodic vector populations.
- To investigate the impact of this type of vector dynamics on disease spread and severity.
- To construct a model for *Xylella fastidiosa* diseases that captures the dynamics of the vector population observed in the field.
- To validate the developed model using empirical data.





## 5. Non-stationary vector populations

**Published as:**

À. Giménez-Romero, R. Flaquer-Galmés, and M. A. Matías, “Vector-borne diseases with nonstationary vector populations: The case of growing and decaying populations”, [Phys. Rev. E](#) **106**, 054402 (2022)

## 5.1 Introduction

Vector-borne diseases are caused by infectious agents transmitted by living organisms, called vectors, frequently insects. These diseases represent a significant threat to global human health [229], causing diseases such as malaria, dengue, yellow fever, Zika, trypanosomiasis, and leishmaniasis [230]. Vector-borne human diseases are responsible for more than 17% of all human infectious diseases, causing millions of cases and more than 700 000 deaths annually [231]. Moreover, crop production and farm profitability are also affected by bacterial [232] and virus [233] vector-borne diseases. Some examples are the Pierce's Disease of grapevines, which has resulted in an annual cost of approximately \$100 million in California alone [234], the olive quick decline syndrome, which could cause about 5 and 17 billion US\$ of losses in Italy and Spain, respectively, over the next 50 years in the absence of disease control measures [235], and the multiple diseases caused by viruses [236], with diseases like the tobacco mosaic or tomato spotted wilt transmitted by aphids and other vectors.

Compartmental deterministic models, e.g., the well known SIR model [65], have been widely used in the modeling of vector-borne diseases after the seminal work of Ross and Macdonald [66], which opened the way to controlling malaria outbreaks by acting on the vectors of the disease (the *Anopheles* mosquito). These models consider that both host and vector populations can be divided into different compartments describing different states of the individuals, such as susceptible, infected or removed (recovered or dead) [223], and the time-evolution of these compartments is expressed as a system of ordinary differential equations, defining a dynamical system. Compartmental models provide a mean-field description that implies well-mixed (in practice spatially homogeneous) populations. The well-mixed approximation will be valid whenever the mean distance among hosts is smaller than the mixing length of vectors before they die. In the case of vector-borne diseases, it is also equivalent to every vector effectively interacting with all the hosts and every host with all the vectors. A mean-field description is not always valid in spatially extended systems, but still, it is often the first step before writing a spatially explicit description.

The most relevant piece of information about a disease is whether an epidemic outbreak will take place or not. The *basic reproduction number*,  $R_0$ , measures the number of secondary infections caused by an initial infected individual in a fully susceptible population, defining the epidemic threshold [178, 237], that determines the emergence (or not) of an outbreak. If  $R_0 > 1$ , an epidemic outbreak will occur, while there will be no outbreak otherwise. The standard way of computing  $R_0$  in deterministic compartmental models is based on the existence of an initial disease-free (pre-pandemic) equilibrium, represented by the absence of infected hosts and vectors [238, 239]. Some standard

methods based on the linear stability condition of this equilibrium have been developed to allow the direct computation of  $R_0$ , such as the Next-Generation Matrix (NGM) method [68].

In the case of vector-borne diseases, some models assume that populations (both hosts and vectors) do not change with time (see e.g. [66, 69, 240]), thus assuming equal birth and death rates. This guarantees the existence of a disease-free equilibrium and the proper use of standard methods to determine  $R_0$ . However, this assumption could be far from reality in several pathosystems. For example, the interaction between temperature, precipitation, and other factors may lead to strong variations in the vector population [241, 242], implying that the pre-pandemic state may not be an equilibrium state and that standard methods cannot be applied.

Compartmental models of vector-borne diseases have another feature that may hinder their practical applicability. The large number of compartments and corresponding parameters in these models may lead to issues of *parameter identifiability and uncertainty* [243], which is more likely to be found in models with many compartments and parameters [244]. Usually, parameter estimation procedures are needed to connect the models with disease data, mainly using incidence or prevalence over time in the host population. Unfortunately, under many circumstances the underlying model parameters are unidentifiable, so that many sets of parameter values produce the same model fit [245]. Moreover, these parameters can be really difficult to determine from the available experimental data. Nevertheless, in some cases, mathematical manipulations can be performed to reduce the model complexity using exact or approximate relations [215]. In such cases, the number of parameters of the models can be usually reduced in terms of new parameters defined as combinations of some original parameters.

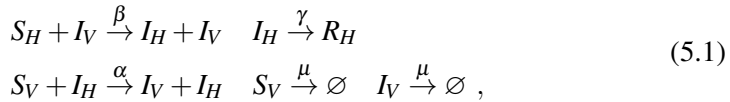
The plan of this chapter is as follows: In [Section 5.2](#), we develop a compartmental model of vector-borne transmitted diseases with constant birth and death rates for the vector population, which we will use to describe the case of growing and decaying vector populations. For simplicity, the model assumes that there is no host-to-host direct transmission and that the development of the disease is faster than host recruitment, which is also a realistic assumption in many cases, like plant diseases. [Section 5.3](#) contains the main results of the study. In particular, we show that the asymptotic approach fails to estimate the  $R_0$  of the model, overlooking outbreaks if some conditions are fulfilled. We provide an alternative method to compute  $R_0$  based on the average number of secondary host infections produced by a primary infected host in one generation. It turns out that the validity of the asymptotic approach depends, among other things, on some time-scales of the model. Furthermore, we discuss and apply

some approximations that allow to reduce the model in favor of simpler ones, with both fewer compartments and fewer parameters. In particular, we show that if some parameters fulfil certain conditions, it is possible to reduce the original model with five compartments and four parameters to a SIR model, with three compartments and two parameters. It is expected that model reductions like this one significantly help in solving possible problems of parameter unidentifiability that plague these models. It is interesting to note that a model in which hosts do not interact directly, but only through vectors, in a certain limit becomes described as if hosts would infect directly one to each other, what is assumed in some studies without suitable confirmation. Finally, the main concluding remarks of the study are presented in [Section 5.4](#).

## 5.2 The model

The compartment model for vector-borne diseases that we will use to illustrate the points to be discussed in this study consists of 5 compartments, 3 of which describe the host population (susceptible,  $S_H$ , infected,  $I_H$ , and removed,  $R_H$ ), while the other 2 describe the vector population (susceptible  $S_V$  and infected vectors,  $I_V$ ). Thus, we consider that the pathogen affects only the hosts and do not consider exposed compartments. In addition, no direct host-to-host or vertical (or mother to offspring for vectors) transmission is assumed. The model could also be generalized to include an exposed host compartment and the above-mentioned transmission modes, which would hinder the theoretical analysis without altering the qualitative conclusions of the study. We consider neither recruitment nor natural death in the host population, so that the total population is constant,  $N_H = S_H + I_H + R_H$ . Finally, we assume that infected hosts do not have a mechanism to combat the disease and become susceptible again. These assumptions are reasonable in the case of many phytopathologies.

The model is defined according to the following processes,

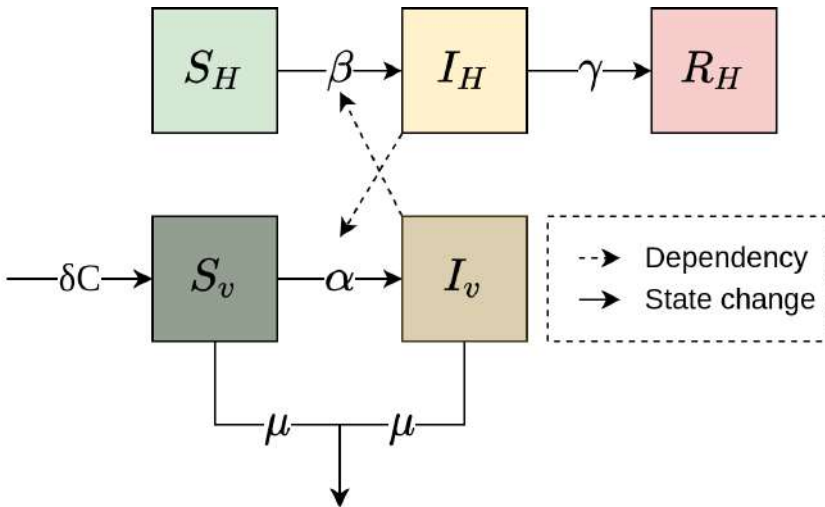


which are graphically described in [Fig. 5.1](#), being the birth of new susceptible vectors described as a source term. Thus, the host-vector compartmental model is written as,

$$\begin{aligned} \dot{S}_H &= -\beta S_H I_V / N_H \\ \dot{I}_H &= \beta S_H I_V / N_H - \gamma I_H \\ \dot{R}_H &= \gamma I_H \\ \dot{S}_V &= \delta C - \alpha S_V I_H / N_H - \mu S_V \\ \dot{I}_V &= \alpha S_V I_H / N_H - \mu I_V, \end{aligned} \quad (5.2)$$

where the crossed nonlinear terms,  $S_x I_y$ , are written divided by the total host population,  $N_H$ , which corresponds to the so-called standard incidence, which differs from the purely bilinear form known as mass action incidence [187].

The model describes infection of susceptible hosts,  $S_H$ , at a rate  $\beta$  through their interaction with infected vectors,  $I_v$ , while susceptible vectors,  $S_v$ , are infected at a rate  $\alpha$  through their interaction with infected hosts,  $I_H$ . Infected hosts exit the infected compartment at rate  $\gamma$ , while infected vectors stay infected the rest of their lifetime, as we consider that the pathogen does not affect them, as it is customary. Vectors die naturally (or disappear from the population by some mechanism) at rate  $\mu$  and are born (appear) at a constant rate  $\delta$  being susceptible. The constant term  $C$  sets the scale of the stationary value of the vector population. Fig. 5.1 shows a schematic representation of the model, and we refer to [69] for a similar model of vector-borne diseases. However, the model in [69] includes exposed compartments and direct host-to-host transmission, but assumes that the birth and death rate of vectors are identical, and thus, the population does not change with time and stays as fixed by the initial condition.



**Figure 5.1:** Schematic representation of the model in Eq. (C.26). Boxes are the compartments in which the population is divided, solid arrows represent changes in state (so transitions between compartments), and dashed arrows depict the crossed interaction between hosts and vectors.

### 5.2.1 Preliminary analysis of the model

From Eq. (C.26), it is straightforward to verify that the population of hosts remains constant over time,  $N_H = S_H + I_H + R_H$ , while the vector population

fulfills,

$$\dot{N}_v = \dot{S}_v + \dot{I}_v = -\mu(S_v + I_v) + \delta C = -\mu N_v + \delta C, \quad (5.3)$$

with solution,

$$N_v(t) = \frac{\delta}{\mu} C + \left( N_v(0) - \frac{\delta}{\mu} C \right) e^{-\mu t}. \quad (5.4)$$

From Eq. (5.4), one can write the stationary value for the vector population,  $N_v^*$ ,

$$N_v^* = \lim_{t \rightarrow \infty} N_v(t) = \frac{\delta}{\mu} C. \quad (5.5)$$

Thus, if the initial population of vectors is below (above) the stationary value, the vector population will grow (decrease) until it reaches the stationary value. On the other hand, if  $N_v(0) = N_v^* = \delta C / \mu$  the initial population of vectors is already at the stationary state. The initial condition for the vector population can be written in terms of its stationary value Eq. (5.5),  $N_v(0) = f N_v^*$ , where both  $f < 1$  and  $f > 1$  are possible, so that one gets,

$$N_v(t) = N_v^* [1 + (f - 1) e^{-\mu t}]. \quad (5.6)$$

We note that vector-borne disease models that assume constant vector populations (e.g., [69]) can be recovered by setting  $\delta = \mu$  and  $C = N_v(0)$ , so that any initial condition for the vector population is stationary, i.e.,  $\dot{N}_v = 0$  in Eq. (5.3) and  $N_v(t) = N_v(0)$ . We note that our model describes populations with an asymptotic stationary vector population and cannot describe periodic vector populations.

## 5.3 Results

### 5.3.1 Epidemic threshold and disease dynamics

Let us start with the cases in which any initial condition for the vector population is stationary and the total vector population remains unchanged. This will happen when the birth  $\delta$  and death  $\mu$  vector rates are identical, independently on the initial condition of the vector population, or in the case in which the initial condition of the vector population is already at its stationary value,  $N_v(0) = N_v^*$ , independently of the values of  $\delta$  and  $\mu$ . In such a case, the initial disease-free state of the model, given by  $I_H(0) = I_v(0) = 0$ , is a fixed point (equilibrium state) of the dynamical system Eq. (C.26), independently of the other initial conditions for the host and vector populations. This allows the definition of the basic reproduction number,  $R_0$ , using standard methods such as linear stability analysis or the Next-Generation Matrix (NGM) method [68] (see Appendix B.1).

In other cases, the total vector population will vary with time provided that the initial condition,  $N_v(0)$ , is not identical to the asymptotic value at large

times,  $N_v^*$ . In these cases, an initial disease-free state is not an equilibrium (fixed point) of the model. However, in the literature it is customary to apply the standard techniques, i.e., NGM, to compute  $R_0$  using the vector population in the asymptotic state, that is the post-pandemic disease-free equilibrium [246–250]. The use of these methods is supported by the fact that the asymptotic dynamics of the model converges to the dynamics of the subsystem where the vector population is stationary [251, 252]. In both cases, the basic reproduction number is given by,

$$R_0 = \frac{\beta \alpha S_H(0)}{\mu \gamma N_H^2} N_v^* . \quad (5.7)$$

As usual,  $R_0$  accounts for the number of secondary infections produced by an infected individual in one generation and controls the threshold behavior of the model: for  $R_0 < 1$  the epidemic dies out and for  $R_0 > 1$  an outbreak occurs. By one generation, we refer to the typical time in which new infections can be produced, being the generation time in our model,

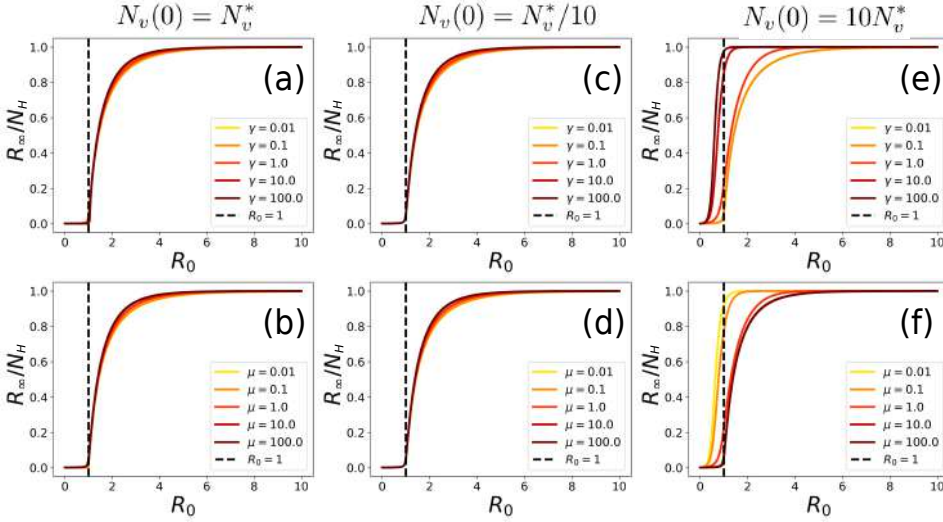
$$t_g = 1/\gamma + 1/\mu . \quad (5.8)$$

Now we will show that Eq. (5.7) is not always predictive about the onset of the epidemic. In Fig. 5.2, the final size of the epidemic,  $R_\infty/N_H$ , is plotted as a function of  $R_0$ , where  $R_\infty$  is the number of dead individuals at the end of the epidemic. Fig. 5.2 (a-d) show that Eq. (5.7) does indeed regulate the onset of an epidemic when the initial vector population is in its stationary value or below it. This result is general and does not depend on the time-scales of the system,  $1/\gamma$  and  $1/\mu$ , and so all curves in these panels behave similarly. In contrast, Fig. 5.2 (e-f) shows that Eq. (5.7) does not predict the onset of an epidemic outbreak when the initial vector population is larger than the stationary value. Thus, for  $R_0 < 1$  (computed using Eq. (5.7)), severe outbreaks appear, yielding a mortality of more than the 80% of the total population. However, one can observe that as  $\mu$  is increased, or  $\gamma$  decreased, the predictive power of Eq. (5.7) is progressively recovered.

Thus, only if the vector population reaches its stationary value before infected hosts have produced new infections can the onset of an epidemic be characterized by Eq. (5.7). Let us discuss separately the cases  $f > 1$  and  $f < 1$ , with  $N_v(0) = fN_v^*$ , namely when the initial vector population is above and below its stationary value, that is, decaying and growing vector populations towards the asymptotic state.

If  $f > 1$  Eq. (5.6), the time to approach the stationary value,  $t^*$ , is,

$$(1 + \varepsilon)N_v^* = N_v^* \left[ 1 + (f - 1)e^{-\mu t^*} \right] , \quad (5.9)$$



**Figure 5.2:** Numerical verification of the predictive power of the basic reproduction number relation Eq. (5.7), by plotting the final size of the epidemic,  $R_\infty/N_H$  as a function of  $R_0$ . In panels (a),(b) the initial vector population is in the stationary value, in panels (c),(d) is below,  $N_V^*/10$ , and in panels (e),(f) above,  $10N_V^*$ . Panels (a),(c),(e) show realizations for different  $\gamma$  values with a fixed  $\mu = 1$  baseline value. Panels (b),(d),(f) show realizations for different  $\mu$  values with a fixed  $\gamma = 1$  baseline value.

where  $\varepsilon \rightarrow 0$  is a small parameter controlling the amount by which the vector population differs from its asymptotic value at time  $t^*$ . Thus, the time to approach the stationary value, with precision  $\varepsilon$ , is given by

$$t^* = -\frac{1}{\mu} \ln\left(\frac{\varepsilon}{f-1}\right) = \frac{1}{\mu} \left| \ln \frac{\varepsilon}{f-1} \right|, \quad (5.10)$$

where the last equality assumes that the small parameter  $\varepsilon$  satisfies  $\varepsilon < (f-1) > 0$ .

If the vector population reaches its stationary value before infected hosts have had time to generate new infections, then  $R_0$ , as determined from Eq. (5.7), is a good prediction of the onset of an epidemic, which is equivalent to the condition that  $t^*$  is much smaller than the hosts infectious period,  $t^* \ll 1/\gamma$ ,

$$\frac{1}{\gamma} \gg \frac{1}{\mu} \left| \ln \frac{\varepsilon}{f-1} \right| \quad \text{or} \quad \frac{\mu}{\gamma} \gg \left| \ln \frac{\varepsilon}{f-1} \right|. \quad (5.11)$$

Otherwise, Eq. (5.7) will not be predictive of the epidemic onset, and as shown in Fig. 5.2 (e-f), one may have outbreaks with a substantial final size with

$R_0 < 1$ .

In the case of growing vector populations,  $f < 1$ , if  $R_0 < 1$  an outbreak cannot occur at all, because  $R_0$  is calculated with the asymptotic population,  $N_v^*$ , that is larger than the vector population at any finite time,  $N_v(t) < N_v^* \forall t$ , and so the threshold condition is never attained. In the  $R_0 > 1$  case, the behavior will be richer, and it will depend on the initial condition,  $N_v(0)$ . One can define an instantaneous basic reproductive number,

$$R_0^{(i)}(t) = \frac{\beta \alpha S_H(0)}{\mu \gamma N_H^2} N_v(t) = R_0 \frac{N_v(t)}{N_v^*}, \quad (5.12)$$

using  $N_v(t)$  instead of  $N_v^*$ , with  $R_0^{(i)}(t) < R_0 \forall t$  because the vector population grows. If  $R_0^{(i)}(0) > 1$  (and so  $R_0 > 1$  as well), there will be an outbreak occurring for short times, and the population of infected hosts will start growing immediately. If instead  $R_0^{(i)}(0) < 1$ , but with  $R_0 > 1$ , there must be an intermediate time, say  $t_D$ , for which  $R_0^{(i)}(t_D) = 1$ . Thus, when  $t > t_D$ , the infected host population will start growing and an outbreak will occur, inducing a delay,  $t_D$ , in the outbreak onset.

The difference between the original and the delayed dynamics stems from the waiting time to reach  $R_0^{(i)} = 1$ ,  $t_D$ , plus the non-linear effect associated with a new initial condition for the epidemic outbreak at  $t_D$ . Thus, in the case that  $R_0 > 1$  and  $R_0^{(i)}(0) < 1$ , from Eq. (B.16) and Eq. (5.6) we can analytically approximate the delay as the time needed to reach  $R_0^{(i)}(t_D) = 1$ ,

$$1 + (f - 1)e^{-\mu t_D} = \frac{1}{R_0}, \quad (5.13)$$

which yields the relation,

$$t_D = -\frac{1}{\mu} \ln \left[ \frac{1 - R_0}{(f - 1)R_0} \right], \quad (5.14)$$

where the argument of the logarithm is always positive because  $R_0 > 1$  and  $f < 1$ . Eq. (5.14) is only valid if  $f < 1/R_0$ , for  $R_0^{(i)}(0) = fR_0 < 1$ , as if otherwise  $R_0^{(i)} > 1$  the outbreak would already occur initially.

From Eq. (5.14) one can see that when the initial vector population is far enough from its stationary value,  $f \rightarrow 0$ , the delay saturates to a constant value, instead of increasing. This is,

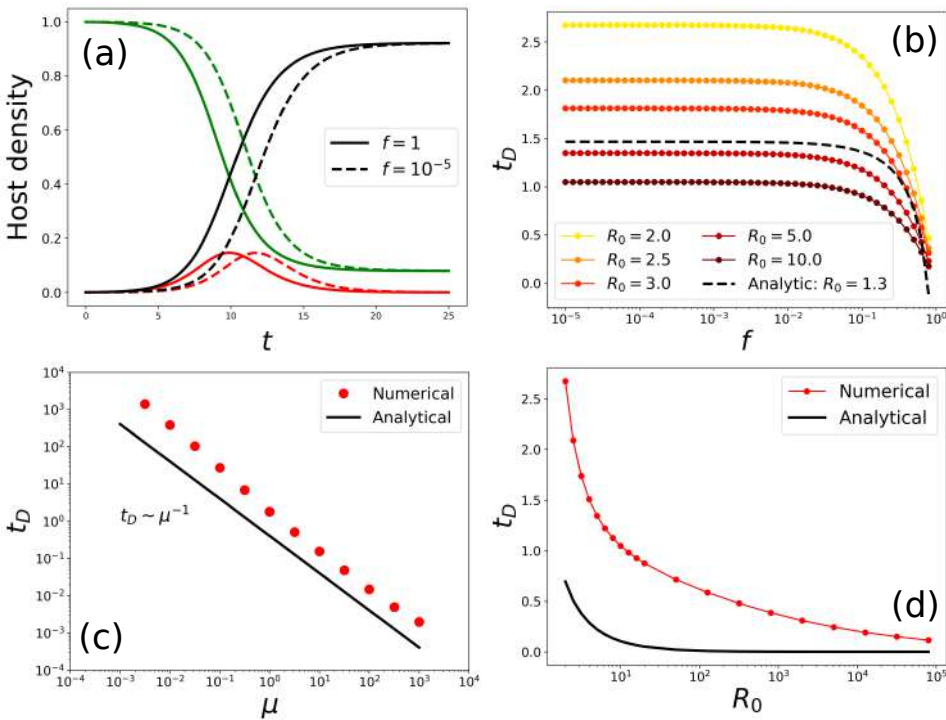
$$\lim_{f \rightarrow 0} t_D = \frac{1}{\mu} \ln \left( \frac{R_0}{R_0 - 1} \right). \quad (5.15)$$

In addition, for increasing values of the basic reproduction number,  $R_0$ , the delay tends to vanish, and from Eq. (5.15). This is,

$$\lim_{R_0 \rightarrow \infty} t_D = \frac{1}{\mu} \ln(1) = 0, \tag{5.16}$$

where the limit  $f \rightarrow 0$  is taken simultaneously to guarantee that  $R_0^{(i)}(0) = fR_0 < 1$ . On the other hand the delay,  $t_D$ , scales with the vectors lifetime,

$$t_D \sim \frac{1}{\mu} = \tau_v. \tag{5.17}$$



**Figure 5.3:** Numerical study of the delay induced by growing vector populations. (a) Comparison of hosts dynamics for a stationary vector population ( $f = 1$ ) and a growing vector population ( $f = 10^{-5}$ ). (b) Time delay as a function of  $f$  for different values of the basic reproduction number  $R_0$ . (c) Time delay as a function of the vector natural death rate. (d) Time delay as a function of the basic reproduction number,  $R_0$ , with  $f = 10^{-5}$ .

Fig. 5.3 (a) shows an example of the time delay caused in the host dynamics when the vector population grows from an initial condition far from the stationary value. In Fig. 5.3 (b) we can qualitatively observe that all the predicted

properties of the delay are fulfilled, namely, the time delay saturates for low  $f$  values and decreases with increasing  $R_0$ . Although the analytical expression (black dashed line) is clearly not exact due to nonlinear effects, Eq. (5.14) captures the basic trends of the time delay,  $t_D$ . This is clear from Fig. 5.3 (c), which shows that the delay scales with  $1/\mu$  and in Fig. 5.3 (d), which shows that the delay tends to 0 in the limit  $R_0 \rightarrow \infty$ , in agreement with the prediction of Eq. (5.16).

### 5.3.2 The basic reproduction number for non-stationary vector populations

As shown in the previous section, traditional methods to compute the basic reproduction number fail in the case of epidemic models with decaying vector populations,  $f > 1$ , unless the timescale of the vector population fulfills the strong inequality condition Eq. (5.11), as illustrated in Section 5.3.1. Here we derive an effective definition of  $R_0$  useful to predict the epidemic onset for vector-borne diseases with decaying vector populations, i.e., the case where traditional methods fail. It is defined as the *average* number of infections produced by an infected individual in *one generation* Eq. (5.8),

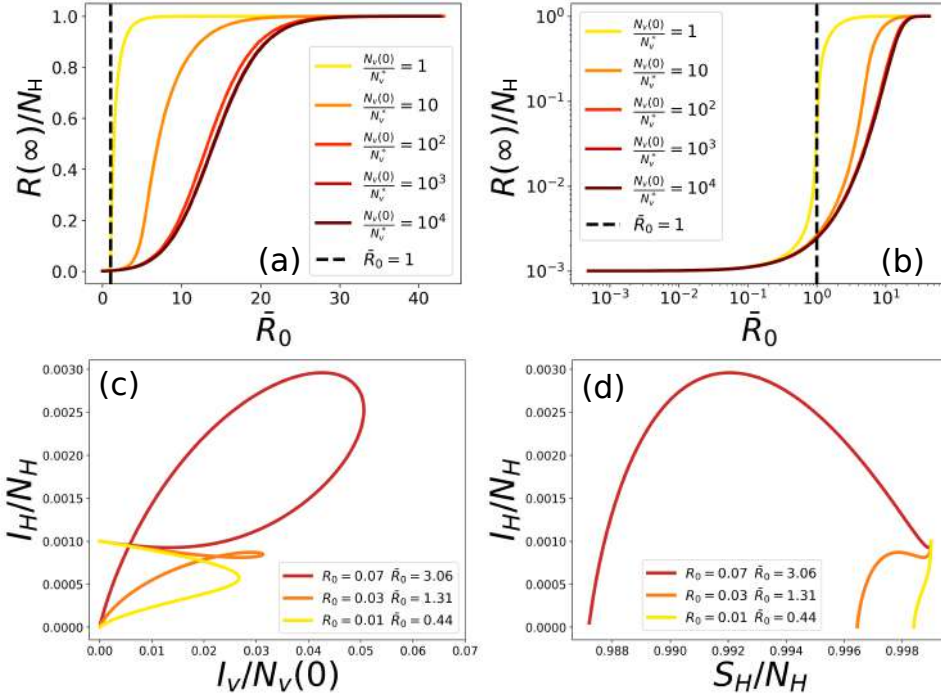
$$\bar{R}_0 = \langle R_0^i(t) \rangle \Big|_0^{t_g} = R_0 \left[ 1 - \frac{1}{\tau} (f-1) (e^{-\tau} - 1) \right] = R_0 \cdot \mathcal{F} \quad (5.18)$$

where  $\tau = 1 + \mu/\gamma$  and  $\mathcal{F}$  accounts for the effect of the decaying vector population on the stationary  $R_0$  (see Appendix B.2 for the full derivation of Eq. (5.18)).

A first observation is that  $\bar{R}_0 > R_0$  always (for  $f > 1$ ). This stems from the fact that  $\tau = 1 + \mu/\gamma > 1$ , so that  $e^{-\tau} - 1 < 0$ , and  $f - 1 > 0$ , which yields  $\mathcal{F} > 1$ . This discussion unravels why standard methods fail to predict the onset of an epidemic under decaying vector populations. Another important point is that if  $\mu/\gamma \gg 1$ , which implies  $\tau \gg 1$ ,

$$\lim_{\tau \gg 1} \mathcal{F} = \lim_{\tau \gg 1} \left[ 1 - \frac{1}{\tau} (f-1) (e^{-\tau} - 1) \right] = 1 + \frac{f-1}{\tau}, \quad (5.19)$$

and if furthermore  $\tau \sim \frac{\mu}{\gamma} \gg (f-1)$  then  $\mathcal{F} \rightarrow 1$  and  $\bar{R}_0 \rightarrow R_0$ . This is in agreement with the discussion in Section 5.3.1 showing that the  $R_0$  computed from standard methods works if  $\mu \gg \gamma$ .



**Figure 5.4:** Numerical verification of the expression for the basic reproduction number for vector-borne diseases with decaying vector populations Eq. (5.18). Final size of the epidemic as a function of the basic reproduction number in panels: (a) linear scale; (b) logarithmic scale. Phase space trajectories in panels: (c)  $I_H/N_H$  vs  $I_V/N_V(0)$  and (d)  $I_H/N_H$  vs  $S_H/N_H$ , where an initial condition  $I_H(0)/N_H = 0.01$ ,  $S_H(0)/N_H = 0.99$  and  $I_V(0)/N_V(0) = 0$  has been used for the 3 cases.  $\mu = \gamma$  has been used in all the simulations.

Fig. 5.4 (a-b) contrasts numerically the validity of Eq. (5.18) to predict the final size of the epidemic as a function of the general basic reproduction number,  $\bar{R}_0$ , in linear and logarithmic scales, respectively. We observe that, independently of the initial condition of vectors, the outbreak occurs for  $\bar{R}_0 > 1$ . However, we may notice that for large values of the initial condition of vectors, the final size of the epidemic grows more slowly, so that larger values of  $\bar{R}_0$  are needed to produce a proper outbreak. This can be explained by the fact that for  $\bar{R}_0$  slightly above the threshold ( $\bar{R}_0 = 1$ ) and with large values of  $f = N_V(0)/N_V^*$ , infections are produced only in the transient period of the dynamics, as  $R_0 < 1$ . This is, while the vector population is decaying to its stationary value, the vectors are able to produce new infections, but once the vector population reaches the stationary value, the epidemic stops. This transmission mechanism

is radically different from that of vector-borne diseases with stationary vector populations, in which the pre-pandemic disease-free state is an equilibrium of the system. The phase-space plots in Fig. 5.4 (c-d) show that the time-averaged basic reproduction number,  $\overline{R}_0$ , is able to accurately predict the conditions under which the infected host population will grow, in contrast with  $R_0$  computed in the post-pandemic fixed point. In essence, for  $\overline{R}_0 > 1$ , the infected host population,  $I_H$ , grows before reaching the absorbing state,  $I_H = I_v = 0$ , while for  $\overline{R}_0 < 1$  the infected host population is monotonically decreasing. We note that Eq. (5.18) is similar to the time-averaged basic reproduction number presented in [253] for the periodic case, which is a first-order approximation to the *true* basic reproductive number [254].

### 5.3.3 Fast-slow approximation

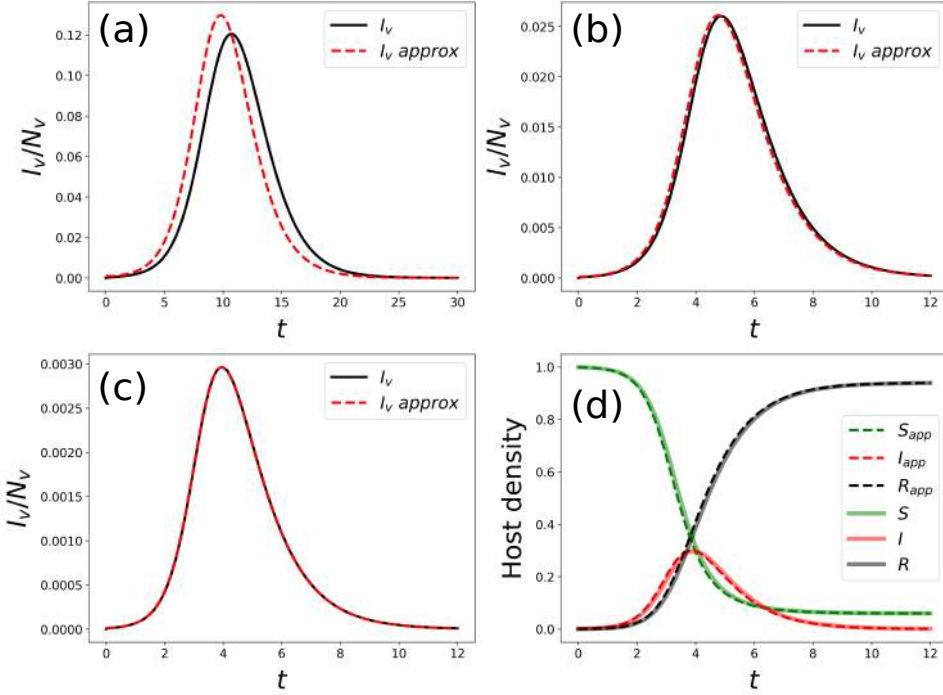
The original 5-D Eq. (C.26) model is certainly not amenable to mathematical analyses due to its high phase-space dimensionality and the fact that it depends on 4 parameters. Moreover, in a real-case application, if the parameters conforming the model are not known, the model could suffer from parameter unidentifiability. However, some approximations can be performed to reduce the mathematical complexity of the model, such as a fast-slow (or adiabatic) approximation.

If the timescale of the vector population evolution is much faster than that of the infected hosts, what is expected to be a good approximation in many practical cases, the vector population will almost instantaneously adapt to its stationary value. Thus, if  $1/\mu \ll 1/\gamma$ , or equivalently if  $\gamma/\mu \ll 1$ , we can rewrite the time derivative of the vector infected population as

$$\varepsilon \dot{I}_v = \frac{\alpha}{\mu} S_v \frac{I_H}{N_H} - I_v, \quad (5.20)$$

where time has been rescaled to  $t' \rightarrow \gamma t$  and  $\varepsilon = \gamma/\mu$  is a small parameter. Then,  $\dot{I}_v$  can be neglected, and the infected vector population can be obtained from the relationship,

$$I_v \approx \frac{\alpha}{\mu} \frac{S_v I_H}{N_H}. \quad (5.21)$$



**Figure 5.5:** Numerical verification of the timescale approximation (Eq. (C.29)) with  $N_H = 100$ ,  $\alpha = \gamma = 1$ .  $\beta$  is chosen such that  $R_0 = 3$ . (a)  $\mu = 1$ , (b)  $\mu = 10$ , (c)  $\mu = 100$ . Panel (d) shows a comparison between the approximate and original models for the parameters used in (c), where the approximated model is expected to represent well the original one.

Substituting Eq. (C.29) into the original system Eq. (C.26) and the identity  $N_v(t) = S_v(t) + I_v(t)$ , while considering that the conditions for which the timescale approximation is valid,  $\mu \gg \gamma$ , imply that the vector population will reach its stationary value almost instantaneously, so that  $N_v(t) \approx N_v^*$ , we obtain the following reduced system,

$$\begin{aligned}
 \dot{S}_H &= -\beta' \frac{S_H I_H}{\lambda N_H + I_H} \\
 \dot{I}_H &= \beta' \frac{S_H I_H}{\lambda N_H + I_H} - \gamma I_H \\
 \dot{R}_H &= \gamma I_H,
 \end{aligned} \tag{5.22}$$

where  $\beta' = \beta N_v^*/N_H$  and  $\lambda = \mu/\alpha$ .

Moreover, if  $f \neq 1$  the above-mentioned timescales' relationship must fulfill  $\frac{\mu}{\gamma} \gg \left| \ln \frac{\varepsilon}{f-1} \right|$  (cf. Eq. (5.11)) and not only  $\frac{\mu}{\gamma} \gg 1$ . It is important to notice

that the presence of direct host-to-host transmission would simply rescale the coefficient  $\beta'$ , and the SIR reduction Eq. (5.22) would keep its validity.

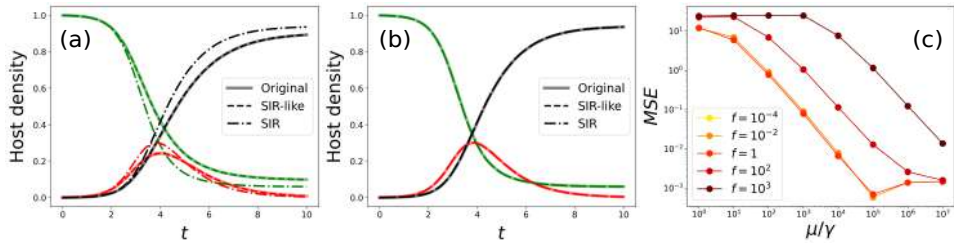
In Fig. 5.5 we numerically verify the validity of the presented fast-slow approximation. As expected, we observe that the approximation breaks down for  $\mu \sim \gamma$  (Fig. 5.5 (a)), while as  $\mu$  becomes larger than  $\gamma$  the approximation improves Fig. 5.5 (b) and it becomes quantitative when  $\mu \gg \gamma$ , Fig. 5.5 (c). Finally, we show in Fig. 5.5 (d) a comparison between the dynamics of the hosts using both the original and the approximated model using the same parameters as in Fig. 5.5 (c), where the results of both models are expected to converge.

### 5.3.4 Reduction to a SIR model

In addition to the previous condition,  $\gamma/\mu \ll 1$ , if one has that  $\lambda N_H \gg I_H$  also holds (which is indeed plausible in this limit) Eq. (5.22), then the model can be written as a standard SIR model,

$$\begin{aligned}\dot{S}_H &= -\beta_{eff} \frac{S_H I_H}{N_H} \\ \dot{I}_H &= \beta_{eff} \frac{S_H I_H}{N_H} - \gamma I_H \\ \dot{R}_H &= \gamma I_H ,\end{aligned}\tag{5.23}$$

where  $\beta_{eff} = \frac{\beta'}{\lambda} = \frac{\beta \alpha N_v^*}{\mu N_H}$ .



**Figure 5.6:** Comparison between the original model and the reductions, Eq. (5.22) (SIR-like) and Eq. (C.30) (SIR) with  $N = 100$ ,  $\mu/\gamma = 10^3$  and  $f = 1$ .  $\beta$  was chosen such that  $R_0 = 3$ . (a)  $\lambda = 1$ , (b)  $\lambda = 10^3$ , (c) Mean Squared Error between the original model and the SIR approximations as a function of the ratio  $\mu/\gamma$  and  $f$ .

In Fig. 5.6 we show the validity of the reduced models Eq. (5.22) and Eq. (C.30). Fig. 5.6 (a) shows that the SIR-like model (Eq. (5.22)) works when the timescale approximation can be performed (as  $\mu/\gamma \gg 1$ ) but the SIR model fails when the condition  $\lambda N_H \gg I_H$  is not fulfilled. Conversely, in Fig. 5.6 (b)

we show that as  $\lambda N_H \gg I_H$  is fulfilled, then the SIR model perfectly matches the original model. Finally, Fig. 5.6 (c) shows the decrease in the mean squared error of the approximation as the condition Eq. (5.11) is fulfilled for different values of  $f$ .

## 5.4 Conclusions

In the present chapter we have analyzed several features of a compartmental deterministic model for vector-borne diseases with three compartments for hosts and two for vectors that does not consider direct host-to-host nor vertical transmission. The goal is to study the behavior of the model in the case that the vector population is not stationary. In this case, the pre-pandemic disease-free state is not a fixed point (equilibrium state) of the dynamical system, and, in principle, the methods that are customarily used to determine the basic reproduction number,  $R_0$ , do not work. This is so because these methods determine the onset of an outbreak by performing a linear stability analysis of the disease-free state, assuming that it is a fixed point of the model. A common assumption made in the literature is to determine  $R_0$  from the asymptotic state for the vectors (if it is not an extinction state), a fixed point of the model.

We have analyzed several initial conditions of the vector population, characterizing different regimes. In the case that the initial condition for the number of vectors is below the asymptotic state, implying that the vector population overall grows, then  $R_0$  as determined from the asymptotic state correctly predicts the existence (or not) of an epidemic outbreak, but with a temporal delay in its appearance. This result contrasts with the situation in which the initial state is above the asymptotic state, with an overall decrease in the vector population. In this case,  $R_0$  determined from the asymptotic state may fail badly, predicting no outbreak while a large fraction of the population might get infected. We present a simple, albeit useful, generalization of  $R_0$  that is able to give a reasonable prediction of the epidemic threshold for decaying populations, including the case in which vectors become extinct, a case in which the asymptotic estimation to determine  $R_0$  cannot be applied.

Compartmental models of vector-borne diseases usually have many compartments and parameters, which can lead to a problem of parameter unidentifiability. The model analyzed here is not an exception, and when applied to real-world cases, many combinations of the parameters could be able to reproduce the available data. Thus, in order to facilitate the application of the model to experimental data, we have studied a useful fast-slow (or adiabatic) approximation that allows to reduce the model if the parameters fulfill certain conditions. In particular, our study shows that under quite realistic assumptions (the typical timescale of host infection and death is much slower than vector

timescales), it is possible to obtain a reduced SIR model. We recall that this reduction implies that, under these assumptions, the process by which hosts (that could be immobile, as they do not interact directly) get infected through the action of vectors is equivalent to a direct interaction among hosts.

The deterministic compartmental model analyzed here, with some modifications, is a clear candidate to study many vector-borne diseases, in particular phytopathologies. Furthermore, in the case of parameter unidentifiability, the model reductions performed in this work could be useful to solve this issue. In any case, this description is still idealized, as compartmental models imply a well-mixed assumption in which space is not explicitly described. This kind of representation is not always applicable to real-world scenarios, although it is useful as a first approximation. Thus, future research should focus on the integration of space and vector mobility in the model to account for more realistic situations.





## 6. A compartmental model for *Xylella fastidiosa* diseases

**Published as:**

À. Giménez-Romero, E. Moralejo, and M. A. Matías, "A Compartmental Model for *Xylella fastidiosa* Diseases with Explicit Vector Seasonal Dynamics", [Phytopathology®](#) **113**, 1686–1696 (2023)

## 6.1 Introduction

Mathematical and computational modeling in Ecology and, in particular, Epidemiology have been recently recognized as powerful approaches to guide empirical work and provide a framework for the synthesis, analysis, and development of conservation plans and policymaking [48–50, 255]. Plant epidemics, mainly plant-virus diseases, have often been described by compartmental models, which deal with the overriding importance of transmission mechanisms in determining epidemic dynamics [256–258]. These models have contributed to providing answers to some questions related to the ecology of plant diseases and have led to direct applications in disease control while guiding research directions [259].

The emergence of vector-borne plant pathogens in new areas causing huge economic impacts, such as *Xylella fastidiosa* and the *Candidatus Liberibacter* spp. (Huanglongbing or citrus greening), has sparked interest in modeling vector-transmitted plant disease epidemics [259, 260]. The vector-borne bacterium *X. fastidiosa* (Xf) is a multi-host pathogen endemic to the Americas that causes economically important diseases, mostly in woody crops [261]. Xf is a genetically diverse species with three evolutionary well-defined clades forming the *pauca*, *fastidiosa*, and *multiplex* subspecies, native from South, Central, and North America, respectively [262]. Within each subspecies, diverse genetic lineages with different host ranges are found. Xf is transmitted non-specifically by xylem-sap-feeding insects belonging to the sharpshooter leafhoppers (Hemiptera: Cicadellidae, Cicadellinae) and spittlebugs (Hemiptera: Cercopoidae) [128].

Recently, Xf has gained renewed interest due to the massive mortality of olive trees in Apulia, Italy [263]. The first focus of the olive quick decline syndrome (OQDS) was detected in 2013 around Gallipoli (Apulia, Italy) [132] and since then has spread throughout the region by the meadow spittlebug, *Philaenus spumarius*. Although this was the first official detection of Xf in Europe, it has recently been demonstrated that the pathogen arrived much earlier in Corsica [264] and in the Balearic Islands [265]. Around 1993, two strains of the subspecies *fastidiosa* (ST1) and *multiplex* (ST81) were introduced from California to Mallorca (Spain) with infected almond plants [265]. To date, over 80% of the almond trees in Mallorca show leaf scorch symptoms, and the outbreak has changed the iconic rural landscape of this Mediterranean island [133].

The meadow spittlebug, *P. spumarius* (Hemiptera: Aphrophoridae), has recently been shown to be the main vector of Xf in Europe both in transmission experiments and in field studies [129, 263, 266–268]. *P. spumarius* is a polyphagous species from the Palearctic region, presenting one generation per year (univoltine) and overwintering as eggs. Foam-forming nymphs emerge at

the end of winter, feeding on herbaceous plants. The time required for their development to the adult stage depends mainly on temperature and humidity [139, 269, 270]. In Mediterranean climates, *P. spumarius* adults generally move from the herbaceous cover to the crop canopy as evapotranspiration increases in late spring (May–June). In mid-summer, the populations of *P. spumarius* tend to decrease in the crop canopy, while the insects are captured more frequently in trees and shrubs interspersed in crops. Summer dispersal of spittlebugs to wild hosts as refugee seems a common general pattern in Mediterranean crops in Italy [269, 271] and Spain [140]. Because the bacterium has not been detected in spring on insects feeding on the herbaceous cover or in weeds in Europe [133, 139, 269], it is assumed that all spittlebug adults acquire the bacteria from the main crop (olive, almond, vine, etc.). Once infected, Xf colonizes the insect foregut in a persistent and non-circulatory manner without transovarial (parent to offspring) or transstadial (inter-stage) transmission [130, 272, 273] and without a period latency after vector acquisition [272, 274].

Several epidemic models have already been developed for Xf diseases, but they lack a realistic description of some relevant processes [259]. Some of these models assume a simple general form for infected host dynamics [275–277] or use a simplified S-I compartmental scheme for hosts, disregarding important features such as the latent period or the host mortality rate [264]. Models that do take these features into account, however, do not explicitly model the population of vectors responsible for disease transmission [278]. Other more recent models have taken a step further in explicitly modeling the vector population [279, 280], but the characterization of its dynamics is still relatively simple, as it overlooks the known seasonal patterns of vector abundance. Several recent studies have provided new insights into the ecology and temporal dynamics of the transmission of Xf by *P. spumarius* in olive plants [269, 281]. However, these experimental data of the pathosystem have not been yet integrated at the population level. Thus, there is a need to continue advancing in the modeling of Xf diseases by developing more realistic models that can elucidate the fundamental processes involved in vector-host-pathogen interactions and help to design effective control strategies.

In this work, we develop a deterministic continuous-time compartmental model to describe the general epidemiological dynamics of diseases produced by Xf in Europe. We explicitly account for key biological aspects of the disease, including the seasonal dynamics of its main vector, *P. spumarius*. Our model is able to describe field data from the two major European outbreaks: the olive quick disease syndrome (OQDS) in Apulia, Italy, caused by the subspecies *pauca*, and the almond leaf scorch disease (ALSD) in Mallorca, Spain, caused by subspecies *multiplex* and *fastidiosa*. We aimed to find the most influential

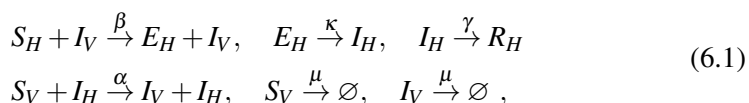
parameters in the model with respect to incidence and mortality in both diseases by performing a global sensibility analysis. With this information, the next goal was to explore control strategies acting especially on the vector population.

## 6.2 Materials and Methods

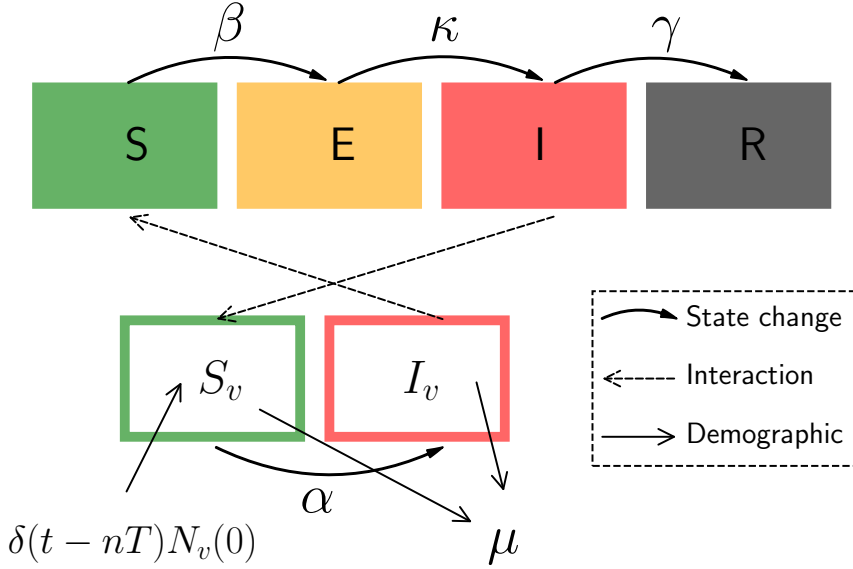
### 6.2.1 Epidemic model: the SEIR-V model

We developed a deterministic continuous-time compartmental model that incorporates the specific biological features of Xf diseases in Europe, including the dynamics of the main relevant vector, *P. spumarius* [282]. To build the model, we took the following considerations: (i) we assume there is no winter recovery of infected hosts and thus they die sometime after infection; (ii) hosts show an asymptomatic period in which they are non-infectious in practice (exposed compartment) because the bacteria are not yet systemically extended [283, 284], while vectors are infectious immediately after acquiring the bacterium [285]; (iii) vectors have an annual life cycle without mother-to-offspring disease transmission [272, 273], so we consider the annual emergence of susceptible newborn vectors and a constant death rate for both susceptible and infected vectors; (iv) infected vectors carry the bacterium during their entire lifespan without affecting their fitness; and finally, (v) we do not consider host recruitment or natural death given that the typical development time of Xf-epidemics is faster than the typical host's life cycle.

Altogether, our deterministic continuous-time compartmental model consists of six compartments, four describing the host population (susceptible,  $S_H$ , exposed,  $E_H$ , infectious,  $I_H$ , and removed,  $R_H$ ), and two describing the vector population (susceptible,  $S_V$ , and infected,  $I_V$ ). The model is defined according to the following processes,



which are illustrated in Fig. 6.1, being the birth of new susceptible vectors described as a source term.



**Figure 6.1:** Schematic representation of the model Eq. (6.2). Boxes are the compartments in which the population is divided; solid curved arrows represent changes in state, i.e., transitions between compartments; dashed arrows depict the crossed interaction between hosts and vectors; and solid straight arrows represent demographic changes in the vector population.

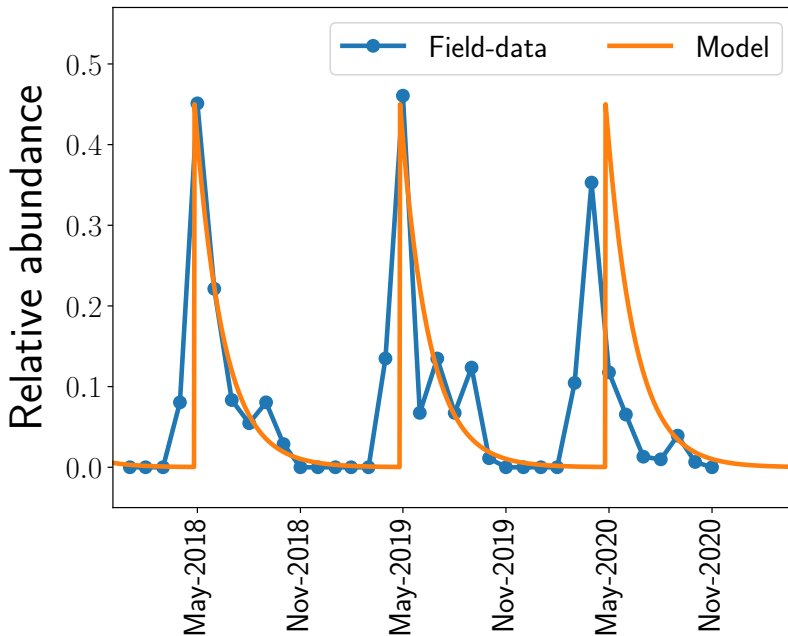
The host-vector compartmental model is written as,

$$\begin{aligned}
 \dot{S}_H &= -\beta S_H I_v / N_H \\
 \dot{E}_H &= \beta S_H I_v / N_H - \kappa E_H \\
 \dot{I}_H &= \kappa E_H - \gamma I_H \\
 \dot{R}_H &= \gamma I_H \\
 \dot{S}_v &= N_v(0) \sum_{n=1}^{\infty} \delta(t - nT) - \alpha S_v I_H / N_H - \mu S_v \\
 \dot{I}_v &= \alpha S_v I_H / N_H - \mu I_v .
 \end{aligned} \tag{6.2}$$

The model describes the exposure of susceptible hosts,  $S_H$ , at a rate  $\beta$  through their interaction with infected vectors,  $I_v$ , while susceptible vectors,  $S_v$ , get infected immediately at a rate  $\alpha$  through their interaction with infectious hosts  $I_H$ . Exposed hosts get infectious at rate  $\kappa$ , having a mean latent period  $\tau_E = 1/\kappa$ , while infectious hosts die at rate  $\gamma$ , having a mean infectious period

of  $\tau_I = 1/\gamma$ . Infected vectors stay infected and infectious for the rest of their lifetime.

Regarding the seasonal dynamics of vectors, we assume that new adults emerge synchronously each year in fields being all susceptible. This is represented by the term  $N_v(0) \sum_{n=1}^{\infty} \delta(t - nT)$  in Eq. (6.2), where  $T = 1\text{yr}$  is the period and  $\delta(t - nT)$  is the Dirac delta function, and basically implements a yearly pulse of new vectors at a certain moment in the year. Vectors are removed (die, move to herbaceous vegetation and other non-host trees, exit the field, etc.) at a given rate  $\mu$ , which we consider identical for susceptible and infected vectors. For simplicity, we consider that the quantity of annual newborn adults,  $N_v(0)$ , is constant. This outburst of new adults followed by an exponential decay resembles the temporal patterns on the abundance of *P. spumarius* observed in crop fields [143, 266, 286, 287] (Fig. 6.2).



**Figure 6.2:** Vector dynamics produced by the model compared to field-data from [143].

Fig. 6.2 shows a time series for the population of *Philaenus spumarius* in Mallorca, taken from [143] (in blue). Superimposed (in orange) is the assumption used in our model Eq. (6.2), the  $\delta(t - nT)$ , i.e., every year susceptible vectors appear in the system.

In our model (Eq. (6.2)), the crossed nonlinear terms in  $\dot{S}_H$  and  $\dot{S}_v$ ,  $S_H I_v$  and

$S_v I_H$ , are divided by the total host population,  $N_H$ . Thus, the vector-to-plant infection process is modeled using mass action incidence, which is density dependent, while the plant-to-vector infection process is modeled using standard incidence, which is frequency dependent [187]. This implies that doubling the number of vectors in the crop field would double the number of resulting exposed (or infected) hosts, as this process is population-dependent (mass action incidence), while doubling the number of hosts would not result in more vectors per unit area being infected, as this process only depends on the contact probability, being frequency dependent (standard incidence). We think this is the most reasonable assumption because, for a given plantation framework, increasing the number of hosts is expected to also increase the area of the field, while the number of vectors is an independent quantity.

### 6.2.2 Basic reproductive number

The basic reproductive number,  $R_0$ , of the model cannot be trivially computed using standard methods such as the Next Generation Matrix (NGM) [68], as there is no pre-pandemic fixed point in the system of differential equations Eq. (6.2). For periodically varying vector populations, rigorous methods have been developed [288], but not for the case of growing or decaying vector populations. Here we use the simple method developed in Chapter 5 [289] (see Appendix B.3), which effectively computes the average number of secondary infections produced by an initially infectious individual in one generation. Thus, the effective basic reproductive number is given by

$$R_0 = \frac{\beta \alpha S_H(0) N_v(0)}{\mu \gamma N_H^2 \mu \tau} (1 - e^{-\mu \tau}) , \quad (6.3)$$

where  $\tau$  corresponds to the time length of one generation, in our case one year. This  $R_0$  is calculated using the initial susceptible host population,  $S_H(0)$ . Below we will also use a time-dependent  $R_0(t)$  using  $S_H(t)$ .

### 6.2.3 Epidemiological data

Epidemiological data from an ALSD outbreak in the island of Mallorca, Balearic Islands, Spain were taken from [265]. Dated phylogenetic analysis and estimates of disease incidence showed that the introduction of both subspecies occurred around 1993, with  $\sim 79\%$  of almond trees infected by 2017 [265]. The annual proportion of infected individuals in the almond tree population between 1993 and 2017 was estimated by analyzing through qPCR the presence of Xf-DNA in the growth rings of 34 sampled trees (cf. Fig. 3 in [265]). The disease progression curve was estimated without distinguishing whether infections were caused by *multiplex* or *fastidiosa* subspecies. In addition, a two-sided bootstrap confidence interval for each data point was set using the SciPy bootstrap function in

Python [290]. On the other hand, epidemic data for OQDS were retrieved from [278]. The data consisted of 2 to 3 yearly censuses of symptom prevalence in 17 olive groves infected with Xf subsp. *pauca* in Apulia, Italy, which were aggregated to fit our model as shown in Fig. 4 in [278]. Because the compartments of our model are not in one-to-one correspondence with those shown in the work of White et al. [278], we used the sum of the symptomatic and desiccated infected trees in the dataset ( $I_S + I_D$ ) to fit the sum of the infected and dead trees ( $I + R$ ) and the sum of susceptible and asymptomatic hosts ( $S + I_A$ ) to fit the sum of susceptible and exposed hosts ( $S + E$ ). The processed data used to fit the model can be found in [291], while the raw data can be found in the supplementary data accessible online of the cited articles [265, 278].

### 6.2.4 Model fitting through Bayesian Inference

We employed an informative normal  $\mathcal{N}(\hat{\mu}, \hat{\sigma}^2)$  prior distribution, with  $\hat{\mu}$  and  $\sigma$  the mean and standard deviation, respectively, for previously measured parameters in the literature, such as the infectious and latent periods for ALS,  $\tau_I \sim \mathcal{N}(14, 4)$ ,  $\tau_E \sim \mathcal{N}(4, 1)$  [265, 283] and OQDS,  $\tau_I \sim \mathcal{N}(3.5, 1)$ ,  $\tau_E \sim \mathcal{N}(1.75, 0.5)$  [285]. The corresponding rates are given by  $\gamma = 1/\tau_I$  and  $\kappa = 1/\tau_E$ , respectively. Similarly, a prior normal distribution was used for the removal rate of vectors,  $\mu \sim \mathcal{N}(0.02, 0.0075)$ , as the mean value  $\mu = 0.02$  already captures the vector dynamics observed in field data (Fig. 6.2). Regarding the prior distribution for the transmission rates, a very wide and uninformative uniform prior distribution,  $\beta \sim \mathcal{U}(0.001, 1)$  and  $\alpha \sim \mathcal{U}(0.001, 1)$ , was used for each parameter. The number of hosts,  $N_H$ , was already provided in the datasets, while, given the lack of information about the vector population, we assumed  $N_v(0) = N_H/2$  for the initial vector population of each year. However, we tested the robustness of our results by changing  $N_v(0)$ .

The posterior distributions of the parameters were approximated using the Markov Chain Monte Carlo algorithm No U-Turn Sampler (NUTS) with the recommended target acceptance rate of 65% [292]. To ensure proper convergence, we constructed three independent Markov chains with  $10^5$  iterations each after a burn-in of  $10^4$  iterations and checked that the results were statistically equivalent. For each chain, we started at the maximum-likelihood parameters yielded by the Nelder-Mead algorithm with 1000 iterations.

The parameters of our compartmental model were determined by fitting the model to data by means of a Bayesian Inference framework using the Turing.jl package [293] in Julia [198]. The scripts used to fit the model can be found in [291].

### 6.2.5 Sensitivity Analysis

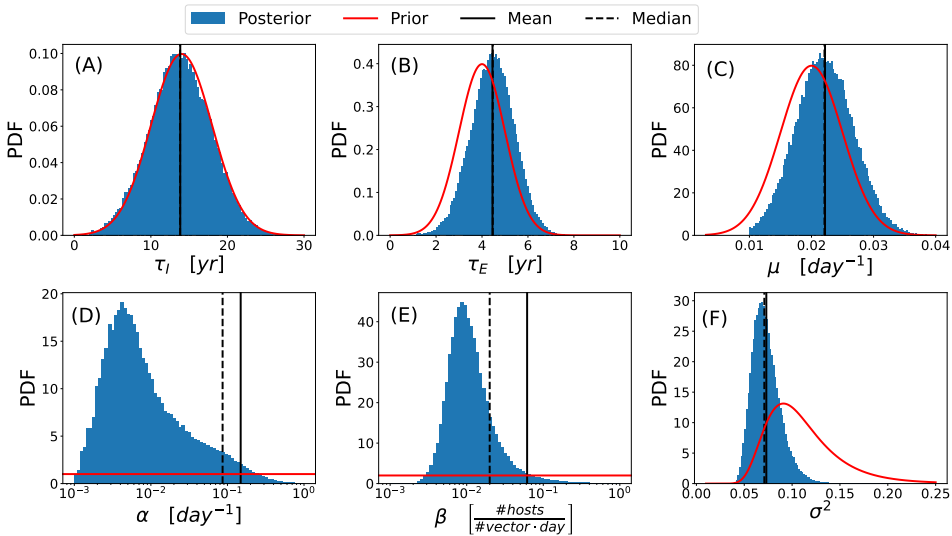
We performed a Global Sensitivity Analysis (GSA) [294] of the model to assess the relative contribution of its parameters and their interactions with different features of the epidemic. In contrast to the Local Sensitivity Analysis (LSA), the GSA assesses the influence of a large domain of the parameter space on the desired outputs of the model. We performed GSA by means of a variance-based analysis, the Sobol method [295]. This particular method provides information not only on how a particular parameter alone influences the model outputs (as happens with LSA), but also due to the nonlinear interactions among two or more parameters. Briefly, the method considers the model output,  $Y$ , as a general function of the inputs,  $f(x_1, \dots, x_n)$ , so that the variance of the output,  $Var(Y)$ , is decomposed as the sum of the variances given by the variations of the parameters alone and its interactions:  $Var(Y) = \sum_{i=1}^n Var(f(x_i)) + \sum_{i < j}^n Var(f(x_i, x_j)) + \dots$ . This information is organized in what are known as Sobol indices. The total order indices are a measure of the total variance of the output quantity caused by variations of the input parameter and its interactions,  $S_T = Var(f(x_1, \dots, x_n))/Var(Y)$ . First-order (or “main effect”) indices are a measure of the contribution to the output variance given by the variation of the parameter alone, but averaged over the variations in other input parameters,  $S_i = Var(f(x_i))/Var(Y)$ . Second-order indices take into account first-order interactions between parameters,  $S_{ij} = Var(f(x_i, x_j))/Var(Y)$ . Further indices can be obtained, describing the influence of higher-order interactions between parameters, but these are not going to be considered.

Following the Sobol method, we analyzed the variation of the time at which the infectious population peaks,  $t_{peak}$ , the magnitude of this peak,  $I_{peak}$ , and the final number of dead hosts,  $R_\infty$ , relative to variations of the model parameters. The method was implemented within the Julia high-level programming language [198] using the sub-package DiffEqSensitivity.jl in the DifferentialEquations.jl package [199].

## 6.3 Results

### 6.3.1 Model fit and parameter estimates

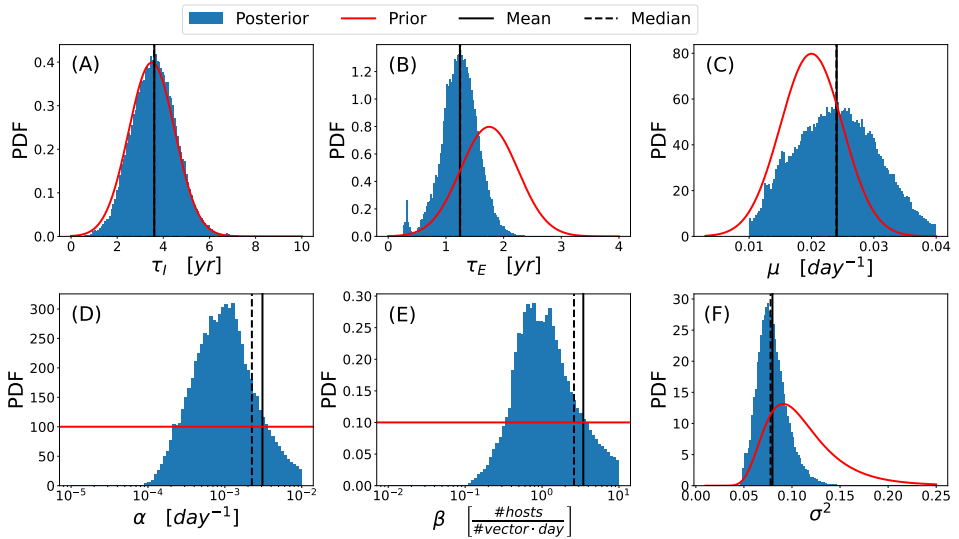
The posterior distributions of the fitted parameters, including their estimated mean and median for ALS and OQDS, are shown in Figs. 6.3 and 6.4, respectively, together with the assumed prior distributions. We observe that the literature-driven priors for the latent and infectious period,  $\tau_E$  and  $\tau_I$ , were already very good guesses and changed slightly converging to the appropriate distribution that better fitted the epidemic data for both ALS and OQDS (Fig. 6.3 (A-B) and Fig. 6.4 (A-B)). Similarly, the prior for the vector removal rate,  $\mu$ , obtained from field data, was good enough so that little changes were needed for convergence (Fig. 6.3 (C) and Fig. 6.4 (C)). On the other hand, we also observe that the completely uninformative priors for the transmission rates successfully converged to the posterior distributions (Fig. 6.3 (D-E) and Fig. 6.4 (D-E)).



**Figure 6.3:** Posterior (blue histograms) and prior (red line) distributions of the model parameters for ALS. Solid and dashed black lines correspond to the mean and median of the posterior distributions. (A) Host infectious period  $\tau_I = 1/\gamma$ . (B) Host latent period  $\tau_E = 1/\kappa$ . (C) Vector removal rate  $\mu$ . (D) Vector infection rate  $\alpha$ . (E) Host infection rate  $\beta$ . (F) The variance of the field data  $\sigma^2$ .

The latter distributions are far from a Gaussian-like shape (note that the x-axis is log-scaled), being heavy-tailed. This kind of distribution highly distorts

the statistical measures of mean, median, and standard error, indicating that the estimates for transmission rates are not as robust as the estimates for the other parameters. These rather uninformative distributions are most probably arising because of the lack of data about the vector, i.e.,  $S_v(t)$  and  $I_v(t)$ , to constrain the fits. In essence, many combinations of  $\alpha$  and  $\beta$  can similarly fit the host data while yielding quite different time series for  $S_v(t)$  and  $I_v(t)$ , which cannot be contrasted due to the lack of field data. Nevertheless, the obtained best-fit mean and median parameters, although quite different, are able to perfectly fit the data (Fig. 6.5). Finally, we also observe that the variance for the field data also converged to a bell-shaped distribution.



**Figure 6.4:** Posterior (blue histograms) and prior (red line) distributions of the model parameters for OQDS. Solid and dashed black lines correspond to the mean and median of the posterior distributions. (A) Host infectious period  $\tau_I = 1/\gamma$ . (B) Host latent period  $\tau_E = 1/\kappa$ . (C) Vector removal rate  $\mu$ . (D) Vector infection rate  $\alpha$ . (E) Host infection rate  $\beta$ . (F) Variance of the field data  $\sigma^2$ .

Mean and median parameter estimates, i.e., the best-fit parameter values for ALS and OQDS, are summarized in Tables 6.1 and 6.2, respectively. As already seen from the posterior distributions, the best-fit values for  $\tau_E$ ,  $\tau_I$ , and  $\mu$  are close to the ones given by literature and field data for both diseases. Conversely,  $\alpha$  and  $\beta$  are rather uninformative, as their 95% confidence intervals cover almost two orders of magnitude. This again indicates that without some data about the evolution of the vector states in time,  $S_v(t)$  and  $I_v(t)$ , it is nearly

impossible to derive the proper values for these parameters.

**Table 6.1:** Estimated epidemiological parameters from Bayesian model fitting to the disease progression curve of ALSD in Mallorca.

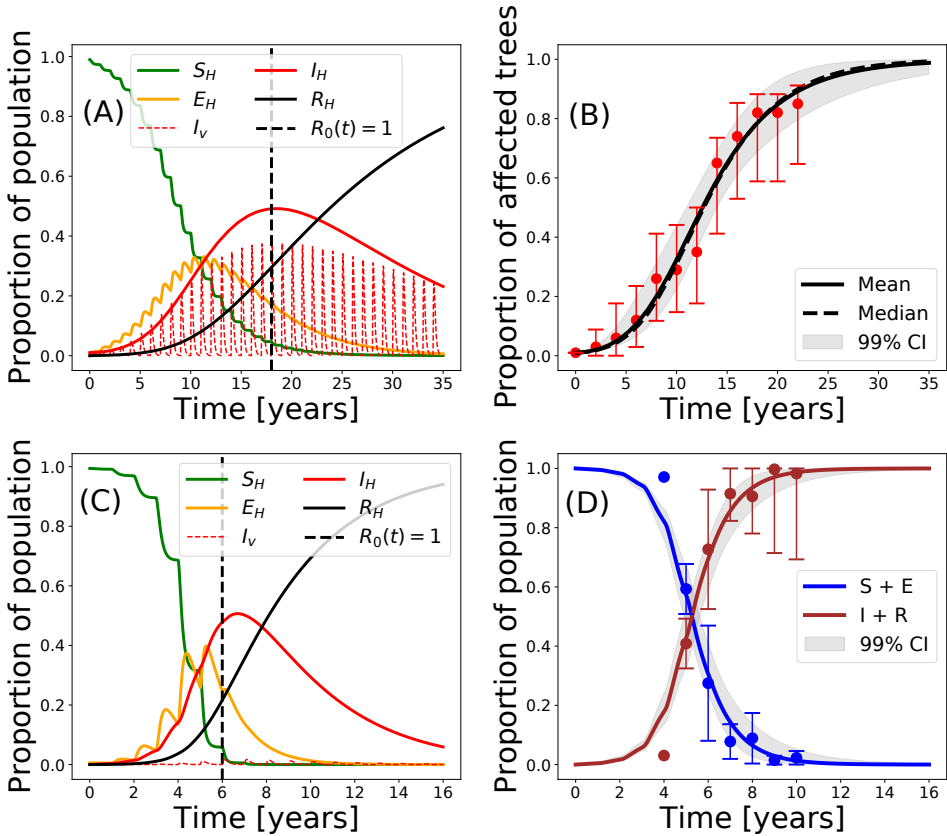
Parameter	Definition	Units	Posterior Mean	Posterior Median	95% C.I.
$\tau_I$	Host infectious period	yr	13.84	13.82	[7.12, 20.47]
$\tau_E$	Host latent period	yr	4.46	4.47	[2.88, 5.99]
$\beta$	Host infection rate	$\frac{\text{\#host}}{\text{\#vector-day}}$	0.062	0.02	[0.0061, 0.3013]
$\alpha$	Vector infection rate	$\text{day}^{-1}$	0.15	0.086	[0.0047, 0.54]
$\mu$	Vector removal rate	$\text{day}^{-1}$	0.0222	0.0221	[0.015, 0.030]
$R_0$	Basic reproductive number	-	133	25	-

Overall, the data falls within the 99% confidence limits of the fitted model for both the ALSD and OQDS outbreaks (Fig. 6.5 (B, D)). We also computed the instantaneous reproductive number,  $R_0(t)$ , by using Eq. (6.3) with  $S_H(t)$  instead of only  $S_H(0)$  along the simulation. Noteworthy,  $R_0(t) = 1$  coincides with the stopping of new infections being produced, i.e., the number of exposed hosts does not increase (Fig. 6.5 (A, C)). This supports our approximate method for computing the reproductive number for Xf diseases (Appendix B.3, Eq. (6.3)). Due to the different time scales of both epidemics ( $\tau_I^{ALSD} + \tau_E^{ALSD} > \tau_I^{OQDS} + \tau_E^{OQDS}$ ), the OQDS outbreak dies out earlier than the one for ALSD.

**Table 6.2:** Estimated epidemiological parameters from Bayesian model fitting to the disease progression curve of OQDS in Apulia.

Parameter	Definition	Units	Posterior Mean	Posterior Median	95% C.I.
$\tau_I$	Host infectious period	yr	3.61	3.60	[2.06, 5.20]
$\tau_E$	Host latent period	yr	1.24	1.25	[0.70, 1.75]
$\beta$	Host infection rate	$\frac{\text{\#host}}{\text{\#vector-day}}$	3.44	2.60	[0.55, 8.79]
$\alpha$	Vector infection rate	$\text{day}^{-1}$	0.0031	0.0022	[0.0005, 0.0084]
$\mu$	Vector removal rate	$\text{day}^{-1}$	0.0240	0.0240	[0.014, 0.035]
$R_0$	Basic reproductive number	-	33	21	-

We notice that for ALSD a large proportion of the vector population gets infected every year (Fig. 6.5 (A)), while a very small proportion is needed in OQDS to produce a lethal outbreak (Fig. 6.5 (C)). However, this last statement is rather unrealistic, as around 50% of the vectors that are captured in Apulia are indeed infected by Xf [282, 296]. Thus, the evolution of the infected vector population should be qualitatively similar to that obtained for ALSD (Fig. 6.5 (C)). As previously explained, different suitable combinations of  $\alpha$  and  $\beta$  parameters should give rise to similar progression curves for the hosts while different ones for the vectors, but the realistic values for these parameters cannot be obtained from the Bayesian fit due to the lack of data of the vector states,  $S_V(t)$ , and  $I_V(t)$ .

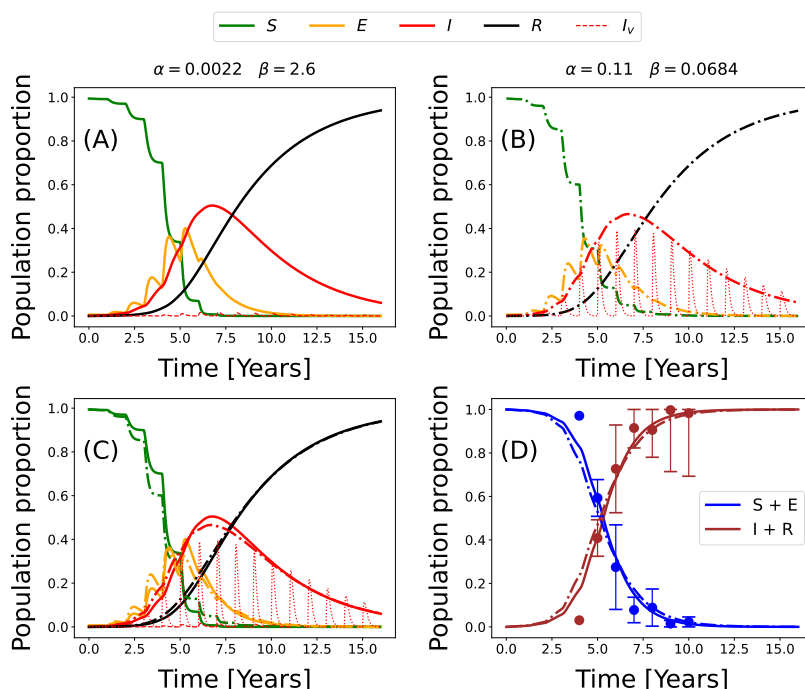


**Figure 6.5:** (A) Simulation of the model with the best-fit parameters for ALSD. (B) Model fit to field data by means of the mean and median values of the posterior distributions of the parameters for ALSD. (C) Simulation of the model with the best-fit parameters for OQDS. (D) Model fit to field data by means of the mean and median values of the posterior distributions of the parameters for OQDS. The gray-shaded area corresponds to the 99% confidence interval. The error bars for the field data correspond to their 95% confidence interval obtained with a bootstrapping technique.

Nevertheless, by manually exploring other values for  $\alpha$  and  $\beta$  parameters, we can obtain a more biologically plausible scenario for the OQDS that is still able to fit the available data for the hosts. Fig. 6.6 (A) shows a simulation of the model with previously inferred best-fit median parameters for OQDS. By changing the values of  $\alpha$  and  $\beta$ , we obtain a more realistic scenario, i.e., around 50% of the vector population getting infected during the outbreak (Fig. 6.6 (B)) [282, 296]. Noteworthy, the  $\beta$  value obtained in this way is almost identical to the transmission rate recently reported by [281] for OQDS. This change in

the transmission parameters only affects the progression curve of the infected vector population, being the progression of the host compartments practically unchanged (Fig. 6.6 (C)). Anyway, both sets of parameter values for  $\alpha$  and  $\beta$  can properly fit the field data, corresponding exclusively to the host population (Fig. 6.6 (D)).

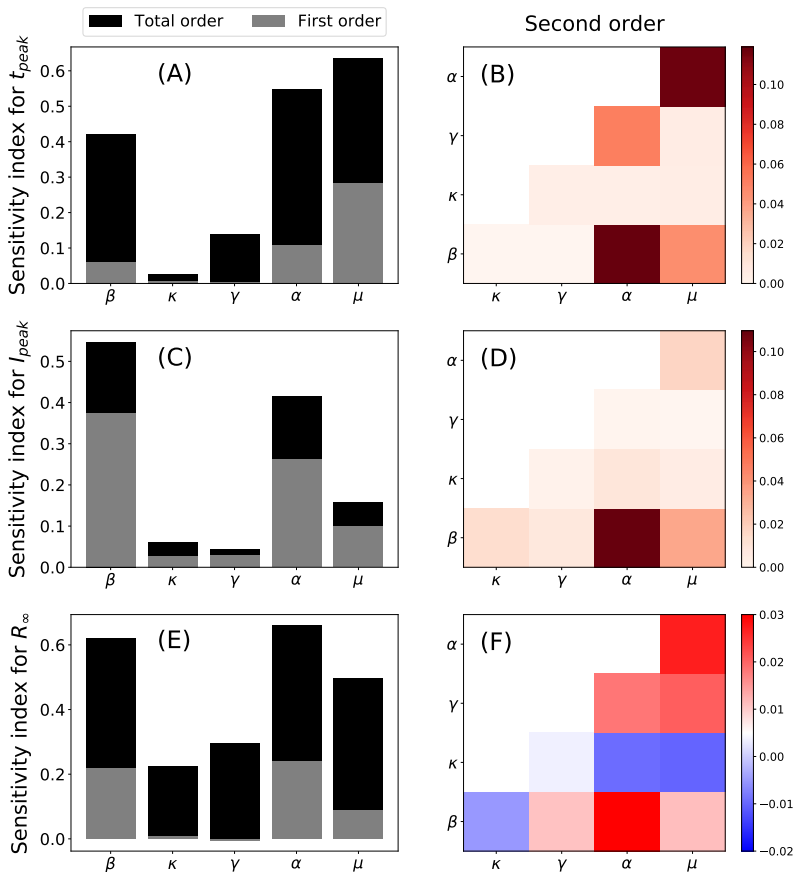
The model adjusted to the progression curves of both diseases indicates that the transmission rate  $\alpha$  must be greater than  $\beta$  when the proportion of infected vectors is relatively high ( $> 30\%$ ). We checked if the relation between  $\alpha$  and  $\beta$  held when changing the assumed  $N_v(0) = N_H/2$ , obtaining that it kept approximately the same for very different values of the initial vector population.



**Figure 6.6:** (A) Simulation of the model with the original best-fit parameters for OQDS. (B) Simulation of the model with the original best-fit parameters for OQDS but with different  $\alpha$ ,  $\beta$  values. (C) Comparison of the progression curves. Note that the curves for the hosts are very similar, while the curve for the infected vector population is very different. (D) Comparison of the model fit to the data with both simulations. Solid lines correspond to results with the original best-fit parameters, while dash-dot lines correspond to the results of the more realistic scenario with different  $\alpha$  and  $\beta$ .

### 6.3.2 Global Sensitivity Analysis

We computed the sensitivity indices for the model parameters with respect to the more relevant quantities of interest, namely, the time at which the number of infectious hosts is maximal,  $t_{\text{peak}}$ , the maximum number of infectious hosts,  $I_{\text{peak}}$ , and the final number of dead hosts,  $R_{\infty}$ . The results were obtained exploring the parameter space constrained to the intervals  $\{\beta \in (0.001, 0.1), \tau_E \in (3, 7), \tau_I \in (5, 25), \alpha \in (0.001, 1), \mu \in (0.01, 0.04)\}$  using  $10^4$  Quasi-Monte Carlo samples and are summarized in Fig. 6.7.



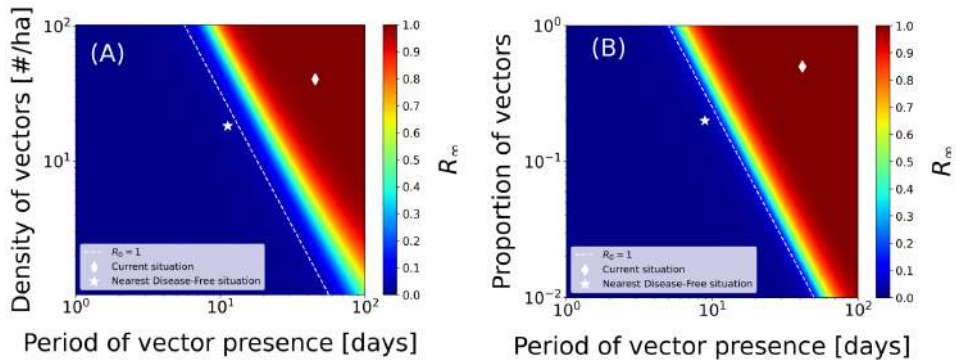
**Figure 6.7:** Global Sensitivity Analysis of the model parameters performed with the Sobolj method with respect to the time at which the infectious population peaks,  $t_{\text{peak}}$  (A-B), the magnitude of this peak,  $I_{\text{peak}}$  (C-D), and the final number of dead hosts,  $R_{\infty}$  (E-F). The left column (A,C,E) shows the total and first-order indices and the right column (B,D,F) shows the second-order indices.

Parameters  $\alpha$ ,  $\beta$ , and  $\mu$  are the most influential with regard to the time at which the infectious host population peaks,  $t_{peak}$ , the magnitude of the peak,  $I_{peak}$ , and the final number of dead hosts,  $R_\infty$ . The total output variance (total order indices) cannot be explained by the variances of the parameters alone (first order indices) (Fig. 6.7). Therefore, higher-order interactions among the parameters importantly affect the sensitivity of the quantities under study. Indeed, the contribution to the total output variance of  $\gamma$  and  $\kappa$  for  $t_{peak}$  and  $R_\infty$  come notably from higher-order interactions. This can be checked in panels (B,C,F) of Fig. 6.7, in which the contribution to the output variance from interactions between pairs of parameters (second order indices) is represented. Interactions among the parameters contribute to increasing the output variance with respect to  $t_{peak}$  and  $I_{peak}$ , while the effect is more heterogeneous in the case of  $R_\infty$ . In particular, the interactions between  $\alpha - \beta$  and  $\alpha - \mu$  produce the main contributions to the increase of output variance in all cases, while  $\kappa - \beta$ ,  $\kappa - \alpha$ , and  $\kappa - \mu$  interactions decrease the output variance.

### 6.3.3 Epidemic control through vector management

The sensitivity analysis clearly indicates that acting on  $\alpha$ ,  $\beta$ , and  $\mu$  is the best strategy to lower disease incidence and mortality. However, controlling transmission rates is cumbersome, so a different control strategy based only on vector control is considered in this section. In our model, there are two ways of implementing vector-population control: (i) decreasing the typical time,  $1/\mu$ , that vectors spend between crops each year by some mechanism (thus increasing  $\mu$ ) and (ii) reducing the initial number of vectors that invade crops each year (e.g., lowering  $N_v(0)$  via egg or nymph control [142]).

We analyzed the effect of vector management by simulating epidemic outbreaks using different values of  $\mu$  and  $N_v(0)$  and keeping the rest of the parameters as fitted for both ALS and OQDS outbreaks (Fig. 6.8). In both epidemics, decreasing the presence time and the number of vectors contribute to controlling the epidemic by lowering  $R_0$  and, consequently, the final size of the epidemic,  $R_\infty$ . Furthermore, we observe that decreasing vector presence is more efficient than decreasing its annual initial population, i.e., we further reduce  $R_\infty$ , the final size of the epidemic, by applying a similar reduction in the residence time  $1/\mu$ . This could also be anticipated as  $R_0$  depends quadratically on  $1/\mu$  while only linearly on  $N_v(0)$  (Eq. (6.3)). However, the minimal intervention strategy, starting from the current situation in the  $(1/\tau, N_v(0))$  parameter space that yields an absolute control of the epidemic,  $R_0 < 1$ , involves a mixed strategy of lowering both  $1/\mu$  and  $N_v(0)$ .



**Figure 6.8:** Epidemic control through vector management for ALSD in Mallorca (A) and OQDS in Apulia (B). The white shaded line denotes  $R_0 = 1$ , and the white diamond corresponds to the parameter values of the fitted model. The white star is the closest disease-free state to the current situation in this representation.

## 6.4 Discussion

In this work, we have developed a deterministic continuous-time compartmental model for *Xylella fastidiosa* vector-borne diseases in Europe. The model attempts to characterize the main biotic processes that lead to the development of epidemics, including the seasonal dynamics of the main vector, *P. spumarius*. We show how the model is sufficiently general to represent with some accuracy the parameters that determine the ALSD in Mallorca (Spain) and the OQDS in Apulia (Italy), both transmitted by *P. spumarius*. To our best knowledge, this is the first mathematical model describing Xf epidemics that considers the temporal pattern of vector abundance observed in field data, faithfully representing the known biological information about the pathosystem. It includes a dynamic approximation of the non-stationary populations of *P. spumarius*, mathematically represented by a sporadic source term through which vectors are born every year, and an exponential decay term. Due to the non-stationarity of the vector dynamics,  $R_0$  in the model cannot be computed with standard methods such as the Next Generation Matrix [68]. To circumvent this problem, we applied an approximate method to compute it as previously proposed by [289]. We show that this approximate  $R_0$  correctly characterizes the epidemic, further validating the method proposed by [289].

Nonlinear mathematical models of disease transmission enhance our understanding of the different mechanisms operating in an epidemic, especially compared with correlative or machine learning methods, often very useful in practice but offering very little understanding. A key aspect to render these

models useful is the determination of the parameters from available data. If this step can be properly performed, these models become very predictive and especially helpful to design disease control strategies. However, an appropriate calibration of the model relies on access to good-quality field data, which is often the bottleneck for the application of this kind of models. In the present study, the parameters have been obtained using a Bayesian inference framework, which relies on probability distributions rather than point-like measures. This way, mean or median values can be considered together with their confidence intervals able to characterize the robustness of the obtained parameters. In general, we obtained different values of the parameters for the ALSD and OQDS outbreaks in Mallorca and Apulia, respectively. The fitted values, however, are in good agreement with previous field-based measures for each disease, while the differences observed between both outbreaks may reflect differences between the Xf subspecies and crops involved (deciduous vs. evergreen).

One of the conclusions of the study is that the available data for both diseases is not enough to obtain robust estimates for all of the model parameters. The lack of data about the vector population compartments yields many possible values for the parameters that regulate transmission,  $\alpha$  and  $\beta$ , provided that the progression of the host compartments correctly fits the field data. In other words, very infectious vectors (high  $\beta$ ) that hardly ever get infected (low  $\alpha$ ) can produce a similar outbreak within the host population to that produced by very low infectious vectors (low  $\beta$ ) that get infected very often (high  $\alpha$ ). The great difference in these situations would be that, in the former, the infected vector population would be very low, while in the latter, it would be quite high. This is a manifestation of parameter unidentifiability from the fit [243, 244], which stresses the importance of transmission and calls for detailed measurements of the vector population and not just of the hosts. Furthermore, to compare transmission rates between different diseases caused by Xf (e.g.,  $\beta$ ,  $\alpha$ ), it is necessary to know the vector-host population ratio of the pathosystem ( $N_v/N_H$ ), since  $\beta$  is expressed as a number of hosts per vector per day. Although, in general, populations of *P. spumarius* in the canopy of olive trees are much larger than those found in the almond trees of the Balearic Islands during the months of July and August [143], our work is based on data from studies in which information of the vector populations is not provided. Without this information, therefore, conclusive results regarding transmission cannot be obtained.

In any case, our model shows that the vector-to-plant transmission process, mediated by  $\beta$ , is somehow different from that from the plant-to-vector one, mediated by  $\alpha$ . In essence,  $\beta$  must be smaller than  $\alpha$  in order to reproduce the observed outbreaks and have a sufficiently large vector population getting infectious, being this fact independent of the particular choice of  $N_v(0)/N_H$ . This

heterogeneity can be caused by several factors: differences in the efficiency of plant-to-vector transmission with respect to vector-to-plant transmission; differences in contact rates, i.e., susceptible vectors contact trees at a different rate than infected vectors; vector feeding preferences, i.e., differences in the probability of contacting a susceptible host compared to an infectious host, etc. Indeed, our mathematical model assumes constant contact rates with no preferences over any host state, so that under these assumptions, it indicates that the probability of effectively transmitting the pathogen from plant to vector is greater than from vector to plant. However, this interpretation is subject to this particular assumption, so that to fully disentangle this question, experimental work in the form of transmission assays should be performed. Furthermore, we found that the timing and magnitude of the infectious host peak and the final number of dead hosts are mostly controlled by the vector-to-plant transmission rate,  $\beta$ , the plant-to-vector transmission rate,  $\alpha$ , and the vector removal rate,  $\mu$ . Because these parameters are strongly related to the vector, the analysis makes clear that enhancing the knowledge about the vector, as well as obtaining precise data, is crucial to improve the modeling of Xf diseases and pose important questions to be solved in specifically designed experiments.

The fact that the most influential parameters of the model are those related to the vector can be used to design appropriate disease control strategies. Because acting on transmission rates is rather cumbersome, we argue that control strategies should focus on reducing the vector population in crop fields. In our model, this depends on two parameters,  $\mu$ , the rate at which vectors die (or move to herbaceous vegetation and other non-host trees or exit the field), and  $N_v(0)$ , the number of newborn susceptible vectors every year (assumed constant in this study). Our results show that a mixed strategy acting on both parameters is optimal to lower disease prevalence and, eventually, eradicate the disease. Interestingly, we also show that acting on the vector removal  $\mu$  is more effective than controlling the newborn vector population  $N_v(0)$ . In fact, most control strategies carried out in practice for Xf diseases focus on the latter factor, reducing  $N_v(0)$  via egg or nymph control [129, 142, 267]. However, our results indicate that alternative strategies based on increasing the removal (or dispersal) rate of vectors should be explored. Furthermore, the evolution of the population compartments of the hosts and vectors provides relevant information on the epidemiology of both diseases. In both cases, the newly defined basic reproductive number that accounts for a decaying vector population is very predictive of the moment in which new infections are not produced anymore, coinciding approximately with the peak of infectious hosts. Therefore, any intervention with control measures after this peak would have marginal effects on future disease progression.

Our mathematical model is still rather simple, implementing only a few relevant epidemic processes in contrast to the high complexity of the pathogen-vector-host interactions occurring in plant epidemics. Indeed, the model itself raises some questions about these interactions, for example, whether contact rates are homogeneous. Another simplification of the model is the fact that the spatial constraints and the intrinsic stochasticity of the transmission processes are neglected. A straightforward extension of the model would be to include a specific spatial setting and implement the explicit motion of the vector within a stochastic framework, such as individual-based models [220]. With this, the effectiveness of current and further control strategies could be tested and improved controlling for the motion of the vector. For instance, the control strategy based on the removal of symptomatic trees together with their surrounding trees at a given distance could be implemented in the model, evaluate the current effectiveness according to the present protocols, and even provide improved parameters to be implemented in the field. Of course, implementing a model in which the spatial degrees of freedom are explicitly represented would require access to further information about vector mobility and spatially resolved data to confront the model, which is not currently available.

Mathematical models tested against experimental data increase our understanding of the system under study. They also help to identify critical parameters that require better prior information to adjust functions relating to different variables and make the model predictions more accurate to suggest and test control strategies [297, 298]. Our mathematical model suggests a certain lack of knowledge of the transmission processes and reveals that the currently available data is not enough to fit complex models dealing with the explicit dynamics of the vector population.

# Modeling the risk of vector-borne plant diseases

<b>7</b>	<b>Global predictions for Pierce's disease risk</b> .....	<b>151</b>
7.1	Introduction .....	152
7.2	Methods .....	154
7.3	Results .....	161
7.4	Discussion .....	172
<b>8</b>	<b>Pierce's disease risk under global warming</b> .....	<b>177</b>
8.1	Introduction .....	178
8.2	Methods .....	179
8.3	Results .....	184
8.4	Discussion .....	190
<b>9</b>	<b>Pierce's disease risk with high-resolution climate data</b> ...	<b>195</b>
9.1	Introduction .....	196
9.2	Results .....	198
9.3	Discussion .....	204
9.4	Methods .....	206



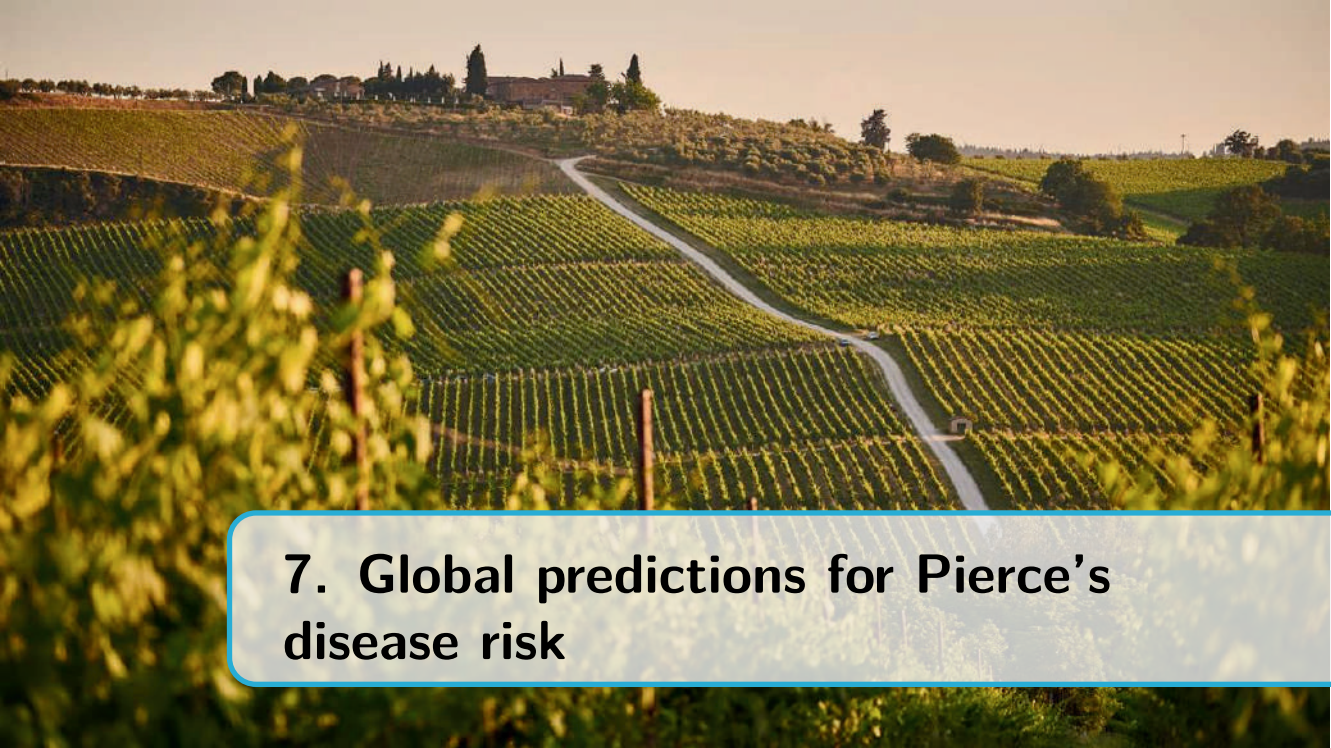
## Summary

In the face of climate change, the threat posed by vector-borne plant diseases to global agriculture and food security has become increasingly dynamic and unpredictable. Among these threats, Pierce's disease of grapevines, caused by *Xylella fastidiosa*, stands out due to its significant impact on viticulture. This part of the thesis focuses on understanding and predicting the influence of climate on the potential distribution and severity of vector-borne plant diseases, using Pierce's disease as a case study. This has been traditionally challenging due to the complex interactions between the pathogen, the vector, and the host plant, as well as the influence of environmental factors on these interactions. Species distribution models (SDMs) have been widely used to predict the potential distribution of plant diseases by focusing on individual components of the pathosystem, such as the pathogen or the vector. However, these models do not yet provide the epidemiological niche of the disease but rather the ecological niche of their constituents. Here, we develop a mechanistic climate-driven epidemiological model to address these limitations, integrating the effects of climate on the vector, the pathogen, and the host plant, as well as the interactions between these components. We validate the model using spatio-temporal data on the distribution of Pierce's disease in the United States, showing that it accurately captures the observed patterns. We then use the model to predict the potential distribution of Pierce's disease under current and future climate scenarios with available climate datasets, and study the effect of high-resolution climate data on the model predictions. Our results suggest that the potential distribution of Pierce's disease is currently constrained to climatic Mediterranean regions, but will expand globally under future climate scenarios, with significant implications for viticulture in Europe.

### Objectives

- To advance the methodologies used in modeling the potential distribution of vector-borne plant diseases.
- To predict the potential distribution of Pierce's disease under current and future climate scenarios with available climate datasets.
- To study the effect of high-resolution climate data on the predictions.
- To analyze the potential impact of Pierce's disease for viticulture worldwide, specially in Europe.





## 7. Global predictions for Pierce's disease risk

### Published as:

À. Giménez-Romero, J. Galván, M. Montesinos, J. Bauzá, M. Godefroid, A. Ferrer, J. J. Ramasco, M. A. Matías, and E. Moralejo, "Global predictions for the risk of establishment of Pierce's disease of grapevines", [Communications Biology](#) **5**, 1389 (2022)

## 7.1 Introduction

Emerging plant pathogens and pests are costly both economically and environmentally for society [299–302]. Among valuable crops recurrently affected by emerging diseases, the grapevine occupies a remarkable place in the history of plant pathology [303–306]. Nowadays, Pierce's disease (PD) is considered a potential major threat to winegrowers worldwide [261]. The annual economic burden in California alone has been estimated at over \$100 million [234], and the disease is a well-recognized limiting factor in the cultivation of *Vitis vinifera* in the southeastern US [261]. In Europe, despite strict quarantine measures to protect the wine industry (Directive 2000/29/EC), PD has recently been established for the first time in vineyards on the island of Mallorca, Spain [268, 307]. This finding, alongside the detection of PD in Taiwan [308], has raised concerns about its possible spread to continental Europe and other wine-producing regions worldwide.

The causal agent of PD [309], the bacterium *Xylella fastidiosa* (Xf) [126], is native to the Americas, where it also causes vector-borne diseases on many economically important crops, such as citrus, almond, coffee, and olive trees [274, 310]. Xf is phylogenetically subdivided into three major monophyletic clades that correspond to the three formally recognized subspecies: *fastidiosa*, *multiplex*, and *pauca*, native from Central, North, and South America, respectively [262, 311]. Although as a taxonomic unit Xf infects more than 560 plant species [127], it also shows genetic variation among subspecies and sequence types (STs) in both host specificity and host range [312]. Since 2013, diverse STs of the three subspecies have been detected in Europe mainly associated with crop and ornamental plants [132, 313, 314]; among these is the clonal lineage of the subsp. *fastidiosa* responsible for PD (hereafter termed Xf<sub>PD</sub>). The same genetic lineage also causes almond leaf scorch disease in California [130] and Mallorca (Spain) [265], where it is widespread in almond plantations and vineyards, affecting more than 23 grape varieties [268].

A key trait in the understanding of Xf's invasive potential is its capacity of being transmitted non-specifically by xylem sap-feeding insects belonging to sharpshooter leafhoppers (Hemiptera: Cicadellinae) and spittlebugs (Hemiptera: superfamily Cercopidae) [129, 315] – e.g., at least eight species transmit PD in the southeastern US [316]. Such non-specificity would have facilitated Xf<sub>PD</sub> invasion after being unwittingly brought to Mallorca around 1993 with infected almond cuttings from California and its spread thereafter to grapevines through local populations of the meadow spittlebug, *Philaenus spumarius* [265]. Recently, the role of *P. spumarius* in the transmission of PD in Mallorca has been demonstrated [268], and its involvement in epidemic outbreaks in California, previously thought marginal [128, 317], is being revisited [287, 318]. To date,

the meadow spittlebug has been confirmed as the major vector in the olive quick decline syndrome, PD, and the almond leaf scorch disease outbreaks in Europe [129, 265, 268, 319]; therefore, its geographic distribution should be taken into account when assessing the risk of Xf-related diseases [320].

The tropical origin of Xf subsp. *fastidiosa* already suggests that PD is a thermal-sensitive disease, with the temperature being a range-limiting factor [321, 322]. Thus, the accumulated heat units (i.e., growing-degree days) required to complete the process from Xf<sub>PD</sub> infection to symptom development is critical to predicting the probability of developing PD acute infections [323]. Conversely, the effect of cold-temperature exposures in the recovery of Xf-infected grapevines is a well-established phenomenon [323–325], limiting the geographic range and damage of chronic PD in vineyards in the US [261]. Such “winter curing” has been linked to the average  $T_{min}$  of the coldest month, to exposures to extreme cold temperatures for several days, or to the accumulation of chilling hours [326]. The dynamics of chronic infections –i.e., those that persist from one year to the next year– are determined by the net balance between the number of new infections during the growing season and those infected plants recovered in winter. Because new infections late in the growing season are more likely to recover during winter than early-season infections, the vector’s phenology greatly influences the dynamics of chronic infections and PD transmission [128, 277, 327, 328].

Several works have attempted to predict the potential geographic range of the subsp. *fastidiosa* [329–331] and other Xf subspecies in Europe [235, 332] and worldwide [331] using bioclimatic correlative species distribution models (SDMs). However, none of these works has explicitly included information on vectors’ distribution or disease dynamics. They hence provide little epidemiological insight into the underlying environmental causes underpinning or limiting a potential invasion. An alternative to overcome these limitations is to develop mechanistic models based on the physiology of the pathogen [77], coupled with epidemiological models that consider the disease dynamics while avoiding the difficulties of including transmission parameters for each of the PD potential vectors.

Risk maps often represent an average snapshot that overlooks interannual climate variability and the effects of climate change as limiting disease factors *per se*. This leads frequently to risk overestimation [333–336]. Increased availability of computational resources to deal with demanding climate databases now makes it possible to fit dynamic epidemiological models that include climate variability at broad spatio-temporal scales. For example, high-resolution satellite-based climate data have been employed for testing mechanistic models that relate critical physiological processes of coffee rust with climate variables

in past outbreak events [337]. Despite these important advances, no attempt at exploring mechanistic SDM has been performed yet for PD.

In this work, we present a temperature-driven dynamic epidemiological model to infer where PD would have become endemic in different wine-growing regions worldwide from 1981 onward if we forced the introduction of Xf-infected plants. We follow an invasive criterion as defined by Jeger & Bragard [259] to include, as far as we can, key plant, pathogen, and vector parameters and their interactions for estimating the risk of establishment, persistence, and subsequent epidemic development. The model assumes local  $Xf_{PD}$  spatial propagation among plants mediated by the presence of potential vectors. Due to the limited knowledge about the vectors of PD in most wine-growing regions of the world [128], we employ a fixed maximal estimate for basic reproductive numbers ( $R_0$ ) in the epidemiological models, except for Europe, where there are precise estimations of climate suitability for the main vector *P. spumarius* [320]. This heuristic approach to obtaining PD risk maps yields results that are consistent with all the relevant data available [329]. It also allows us to quantitatively approximate the current potential growth rate of PD incidence in wine-growing regions under different transmission scenarios and to extrapolate the impact of PD by 2050 [338]. By estimating a lower global risk of PD, our study casts doubts on the potential impact predicted for other Xf related diseases transmitted by *P. spumarius* [235], especially in Europe when the vector distribution is taken into account.

## 7.2 Methods

### 7.2.1 Inoculation tests

$Xf_{PD}$ -inoculation tests were conducted in 2018, 2019, and 2020. A sample of 36 local, regional, and international wine grape varieties was selected, which included nine of the 10 most cultivated wine grape varieties, representing more than 80% of the worldwide vineyard surface (<https://www.oiv.int>). Plants were randomly distributed in 12-plant rows along an insect-proof net tunnel and exposed to environmental temperature. In total, 57 rootstock-scion combinations were pin-prick mechanically inoculated [130] with two strains of Xf. subsp. *fastidiosa* (ST1) isolated from grapevines in Mallorca. Disease severity was rated by counting the number of symptomatic leaves eight weeks after inoculation in mid-May and then every two weeks until the 16th week [268]. Full details on the inoculation conditions, isolates, disease score, and statistical analysis are provided in [Appendix C.1](#).

### 7.2.2 Modified Growing Degree Days.

We generalized McMaster & Wilhelm's [339] formulation of growing-degree days to account for the growth rate of  $Xf_{PD}$  as a function of temperature under optimal culture conditions based on the well-known Arrhenius law valid in the relevant temperature range for  $Xf$  (Appendix C.2.1). Specific growth rate ( $k$ ) values at different temperatures were extracted from the publication of Feil & Purcell [323] to build the mathematical function  $f(T)$  describing the  $Xf$ 's instantaneous growth rate dependence on temperature according to

$$f(T) = \begin{cases} 0 & \text{if } T < T_{\text{base}} \\ m_1 \cdot T - b_1 & \text{if } T_{\text{base}} \leq T < T_1 \\ m_2 \cdot T + b_2 & \text{if } T_1 \leq T < T_{\text{opt}} \\ m_3 \cdot T + b_3 & \text{if } T_{\text{opt}} \leq T < T_2 \\ m_4 \cdot T + b_4 & \text{if } T_2 \leq T < T_{\text{max}} \\ 0 & \text{if } T \geq T_{\text{max}} \end{cases}$$

where  $T_{\text{base}} = 12^\circ\text{C}$ ,  $T_1 = 18$ ,  $T_{\text{opt}} = 28^\circ\text{C}$ ,  $T_2 = 32$  and  $T_{\text{max}} = 35^\circ\text{C}$ ; the slopes are  $m_1 = 0.66$ ,  $m_2 = 1$ ,  $m_3 = -1.25$  and  $m_4 = -3$  and the intercepts are  $b_1 = -8$ ,  $b_2 = -14$ ,  $b_3 = 4$  and  $b_4 = 105$ .

MGDD is then defined as:

$$MGDD(t) = \frac{1}{24} \sum_{\tau \in t} f(T(\tau)), \quad (7.1)$$

where  $\tau$  is expressed in hours,  $t$  in years, and we divide by 24 to obtain  $MGDD(t)$  in degree days. To compare whether using other functions for  $Xf$ 's growth rates in vitro could yield differences in the risk indexes, we also fitted data to a smooth beta function commonly used to represent the thermal response in biological processes [340, 341].

### 7.2.3 Disease progress with temperature

Hourly mean temperature data were recorded between April 1 and October 31 in 2018, 2019, and 2020 with an automated weather station (Quimisur, IQ2000). The temperature sensor was at a two-meter height from the bare ground and around five meters from the entrance of the insect-proof net tunnel. To characterize the progress of PD symptoms, we converted into  $MGDD$  units the cumulative hourly mean temperatures measured from the time of inoculation to the day of disease evaluation using Eq. (7.1). In total, 15  $MGDD$  levels were estimated, corresponding to weeks 8, 10, 12, 14, and 16 after inoculation in the years 2018, 2019, and 2020, respectively. Data on the number of symptomatic leaves (severity) for each plant and  $MGDD$  levels were pooled in a single database to seek a generalized average thermal response pattern among

the population of *V. vinifera* varieties (see [Appendix C.1](#)). To model the probability of chronic infections (i.e., persistent year-to-year infections), we used a survival analysis, where the event of interest depends on the cumulative MGDD rather than time. First, we defined a chronic infection cut-off point to transform the number of symptomatic leaves into binary data. Previous research had evidenced that early grapevine infections, in addition to producing more extensive and severe PD symptoms, are more likely to survive the following year than late infections [323, 324, 327]. Furthermore, susceptible cultivars generally show lower recovery percentages compared to the less susceptible ones in the field [342, 343]. Similarly, we observed in our inoculation assays that the majority of infections that reach around five or more symptomatic leaves 12 weeks after inoculation continue to develop more symptomatic leaves the following weeks, while for plants that do not exceed that threshold, symptoms tend to remain stagnant. These results indicate a low probability of survival for infections showing few symptomatic leaves during the growing season and thus support our heuristic approach of assigning five or more symptomatic leaves as a threshold for chronic infection (see [Appendix C](#) and [Fig. C.1](#) for assumptions of chronic infection). Using the “survival” package in R [344], we analyzed the cumulative probability of developing chronic infections as a function of *MGDD*.  $F(MGDD)$  was adjusted to the experimental data by the nonlinear least squares method. The 10th, 33rd, 50th, 66th, and 90th percentiles were used to scale the risk of the total *MGDD* in the logistic function,  $\mathcal{F}(MGDD)$  ([Fig. 7.1](#)).

### 7.2.4 Disease recovery through winter curing

We modeled winter curing considering the effect of temperature duration below a threshold temperature, where we assume that the bacterial killing process increases in efficiency with decreasing temperatures [324]. To adjust a probabilistic model to the accumulation of cold units, we took as reference the distribution and severity of PD in the US proposed by Purcell based on the isolines of the mean  $T_{min}$  of the coldest month (available in [326]), where PD is rare ( $T_{min}$  between  $-1.1$  °C and  $1.7$  °C), occasional ( $1.7 - 4.5$  °C), and severe ( $> 4.5$  °C). Noteworthy, the projection of these isolines in Europe has predicted with some precision the distribution of the establishment of Xf in the continent [329]. To capture the accumulation nature of the chilling process at different climatic zones, we determined the global average correlation between  $T_{min}$  and the average accumulated CDD between November 1 and March 31 in the Northern Hemisphere and between April 1 and October 31 in the Southern Hemisphere using 6,487,200 points distributed throughout the planet. The *CDD* was estimated as

$$CDD(t) = \frac{1}{24} \sum_{\tau \in t} (6 - T(\tau)) \quad \text{for } T_i \leq 6^\circ\text{C}, \quad (7.2)$$

where the threshold 6°C comes from Ref. [324].

### 7.2.5 Global climate data, MGDD/CDD computation

Global mean hourly temperature data were downloaded from the ERA5-Land dataset [345] at 0.1° spatial resolution using GRIB format. The annual average  $T_{\min}$  of the coldest month was calculated from the hourly average temperature from the ERA5-Land dataset. To calculate the annual *MGDD* and *CDD*, a simple Julia [198] library was built on top of the GRIB.jl package [346]. For the Northern Hemisphere, the accumulated *MGDDs* were computed from April 1 to October 31, whereas *CDDs* were estimated from November 1 to March 31, and the reverse for the Southern Hemisphere.

### 7.2.6 Disease model construction

We used a standard susceptible-infectious/infected-recovered (SIR) compartmental model to assess the risk of PD establishment and epidemics worldwide, represented by the following three equations in the large population limit:

$$\begin{aligned}\dot{S} &= -\beta SI/N, \\ \dot{I} &= \beta SI/N - \gamma I, \\ \dot{R} &= \gamma I,\end{aligned}\tag{7.3}$$

where  $S$  is the susceptible host population,  $I$  is the infected population,  $R$  is the dead population, and  $S + I + R = N$  is the total number of vines in the population. The reduction of a vector-borne disease model to a SIR model gives rise to a linear dependence of the basic reproductive number  $R_0$  on the vector population (see [Appendices C.2.7](#) and [C.4](#)). Vector-plant transmission of the pathogen is approximated with an effective plant-to-plant transmission rate  $\beta$  ([Appendix C.4](#)), as has been done previously for other Xf-related diseases [278], and the transition from the infected compartment to the recovered (dead) compartment is given by the recovery (mortality) rate  $\gamma$ . In a mean-field approximation of the onset of an outbreak, the basic reproductive number ( $R_0 = \beta/\gamma$ ) defines the exponential growth/decrease stage in the SIR model ([Fig. 7.1 e](#), [Appendix C.2.3](#)). Although the time from infection to vine death depends on the environmental conditions and the grape wine variety, we assigned a mortality rate of  $\gamma = 0.2 \text{ y}^{-1}$  based on the estimated median survival time of infected vines in California [130]. The maximum growth rate of the epidemic, relevant for an estimation of the risk of establishment, occurs when  $S(t = 0) \sim N$  and was approximated by the (linearized) differential equation,

$$dI/dt \approx \beta I - \gamma I = \gamma I(\beta/\gamma - 1) = \gamma I(R_0 - 1),\tag{7.4}$$

where we have assumed the initial conditions:  $S(t = 0) \approx N$ ,  $I(t = 0) = I(0) \approx 0$  and  $R(t = 0) = 0$ . This linear differential equation can be integrated exactly:

$$I(t) = I(0) \exp(\gamma(R_0 - 1)t) . \quad (7.5)$$

To account for the effect of temperature in the epidemic process, we modify the previous expression as follows

$$\begin{aligned} I(t) &= I(0) \exp(\gamma(R_0 - 1)t) \mathcal{F}(MGDD(t)) \mathcal{G}(CDD(t)) \\ &= I(0) \exp(\gamma(R_0 - 1)t) \Pi(t) , \end{aligned} \quad (7.6)$$

where  $\Pi(t) = \mathcal{F}(MGDD(t)) \mathcal{G}(CDD(t))$  is the cumulative probability of chronic infection, which depends on temperature, and  $R_0$  bears the information on the vector density.

The spatial unit of the model is given by the resolution of the climate data. In this case, it is approximately  $9 \times 9 \text{ km}^2$ , the native resolution of the ERA5-Land dataset ( $0.1^\circ \times 0.1^\circ$ ). We assume uniform conditions within each of the grid cells in terms of vector population, susceptible vines, and parameters that define the model. Risk outcome is calculated for each cell of the spatial raster individually, i.e., there is no simulated spread from one cell to another. Altogether, the equation representing the number of individuals with chronic infections in each cell  $j$  at time  $t$  is written as

$$I_j(t) = \underbrace{I_j(t-1) e^{\gamma(R_0(j)-1)}}_{\text{transmission layer}} \overbrace{\Pi_j(t)}^{\text{climatic layer}} , \quad (7.7)$$

where  $I(t-1)$  is the number of chronic infections in the previous year ( $t-1$ ) and  $\Pi_j(t) = \mathcal{F}(MGDD) \mathcal{G}(CDD)$  is the climatic layer that modulates the growth term and describes the cumulative probability of new infections becoming chronic in the time period between  $t-1$  and  $t$ . The model assumes a homogeneous distribution of the vector population among the grid cells (same  $\beta$  and then same  $R_0(j) = R_0$ ) except for Europe, where information on the spatial distribution of *P. spumarius* is available (see Methods). In this latter case, a spatially dependent  $R_0(j)$  is incorporated into the model by considering the product of the homogeneous  $R_0$  and the spatially dependent climate suitability for vectors (Appendix C.2.7).

To compute the epidemic risk maps, we carried out a simple simulation summarized in three steps: (i) at the initial condition for the first year considered,  $t_0$ , each grid cell is seeded with a single infected plant,  $I(t_0) = 1$ ; (ii) the simulation runs for a year and the incidence is calculated following Eq. (7.7); (iii) we seed again the cells for which the number of infected plants has vanished. In the last

seven years of the simulation, there is no reseeded to allow the system to relax. This process is repeated until the final year  $\tau$ . Finally, the risk index  $r_j(\tau)$  is calculated from the final number of infected plants at grid cell  $j$  as

$$r_j(\tau) = \max \left\{ \frac{\log(I_j(\tau)/I_j(t_0))}{\gamma(R_0 - 1) \tau}, -1 \right\}. \quad (7.8)$$

In this equation,  $r_j$  implicitly delimits three differential risk zones in the maps: 1) non-risk zones where  $r_j \leq -0.09$  and the number of infected plants decreases exponentially; 2) transition areas where  $-0.09 < r_j \leq 0.075$ , and 3) an epidemic risk-zone where  $r_j > 0.075$  and PD can theoretically become established and produce an outbreak, i.e., the number of infected plants increases exponentially (see [Appendix C.2.6](#) for further details.)

### 7.2.7 Model calibration and validation

Model parameters (i.e.,  $R_0$ ) were calibrated with observed records of PD presence in California and the southeast of the US, where the disease is well established. PD distribution data were collected from publications from 2001 to 2020. Publications were filtered by selecting only records where the pathogen detection on symptomatic grapevines was confirmed by PCR or Elisa. The exact coordinates of the records were taken when available in the publication or approximated to locality or county level [262, 324, 326, 347–351]. For modeling purposes and to attempt a general rough estimate of the  $R_0$  parameter valid for the entire US, we assumed a single vector with a uniform spatial distribution. We ran several model simulations with  $R_0$  ranging from 1 to 14. Model prediction performance was estimated using a ROC curve by plotting the true-positive rate (TPR), calculated as the ratio (TP/TP+FN), against the false-positive rate (FPR), calculated as the ratio (FP/TN+FP), where PD absence/presence fulfills the following conditions: true positive (TP), PD is positive and  $r > 0$ ; true negative (TN), PD is negative and  $r < 0$ ; false positive (FP), PD absent but  $r > 0$ ; and false negative (FN), PD positive and  $r < 0$  ([352]). A different approach was followed to estimate  $R_0$  for Europe given that PD is only present in Mallorca and hence spatio-temporal data on the PD distribution is limited to the island. First, we estimated the transmission rate of the main European vector *P. spumarius* from the well-studied disease progress curve of the almond leaf scorch epidemic in Mallorca. Then, using the known mortality rate of PD-infected vines  $\gamma \sim 0.2\text{y}^{-1}$  and the inferred transmission rate,  $\beta = 0.8\text{y}^{-1}$ , the basic reproduction number for PD in Mallorca yields  $R_0 = \beta/\gamma \approx 4$ . Finally, using data on the climate suitability of the vector in Mallorca,  $v = 0.8$ , and inverting the relation  $R_0(j) = R_0 v(j)$ , we estimated  $R_0 \approx 4/0.8 = 5$  as a maximal estimate baseline scenario for PD transmission in

Europe ([Appendix C.2.5](#)). This figure is not intended to be an exact estimate of  $R_0$  but rather an average reference in our model in agreement with the lesser abundance of vectors relative to the US. Furthermore, since there is no information on the distribution of the potential vectors and no PD distribution data to calibrate, we also used a conservative  $R_0 \approx 5$  scenario for the rest of the world.

### 7.2.8 *Philaenus spumarius* SDM

The potential distribution of *P. spumarius* in Europe under current and future (i.e., 2050) climatic conditions was provided by Godefroid et al. [320]. Predictions were obtained using a generalized additive model and two bioclimatic descriptors, i.e., a climatic moisture index for the coldest 8-month period of the year and the average maximum temperature in spring (March, April, and May). Both descriptors reflect physiological constraints acting on life stages of the meadow spittlebug, particularly sensitive to spring temperature and humidity (eggs and nymphs), and were identified as good predictors of *P. spumarius* distribution ([320]). We used the positive relationship between climate suitability and spittlebug adult abundance ([320]) to assume no climatic constraints on vector population sizes at optimal climatic conditions ( $v=1$ ). Climatic suitability indexes,  $v(x)$ , were used to compute a spatially-dependent basic reproduction number,  $R_0(x) = R_0 v(x)$ . The linear dependence between the basic reproduction number and climatic suitability is justified by a vector-borne epidemic compartmental model ([Chapter 5](#) and [Appendices C.2.7](#) and [C.4](#)).

### 7.2.9 Distribution of wine-grape production areas.

Risk maps were focused solely on wine-grape regions, excluding table and dried grape-producing areas. Data on the vineyard surface in Europe were obtained from the CORINE land-cover map [353] ([Fig. 7.6](#)). The Nomenclature of Territorial Units for Statistics (NUTS) was used as a geocoding for the subdivisions of European countries for statistical purposes. To visualize the locations of the main growing regions in the risk maps, we included dots representing the distribution of the main wine-growing regions collected from official statistics and maps from the countries ([Fig. 7.5](#)).

### 7.2.10 Risk assessment by 2050

Climatic variables were obtained with annual resolution by extrapolating the computed  $MGDD(t)$  and  $CDD(t)$  time series up to 2050. The observed trends of the time series were captured using a machine learning-based linear regression model, while the interannual fluctuations were modeled by Gaussian noise ([Appendix C.3](#)). Future risk extrapolations were obtained as the average of  $10^4$  simulations of this process. A correlative SDM was used to estimate vector

spatial distribution in Europe using the global circulation model MIROC5 and greenhouse gas emission scenario RCP4.5, assuming moderate climate change [320]. Afterwards, the risk was computed following the same simulation procedure previously explained.

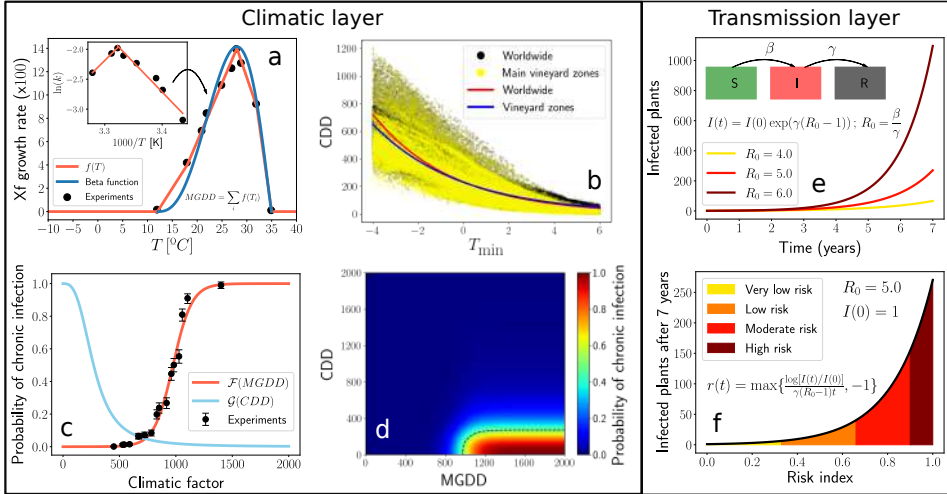
## 7.3 Results

### 7.3.1 Thermal requirements to develop PD

We examined the response of a wide spectrum of European grapevine varieties to  $Xf_{PD}$  infection in three independent experiments conducted in 2018, 2019, and 2020. Overall, 86.1% ( $n = 764$ ) of 886 inoculated plants, comprising 36 varieties and 57 unique scion/rootstock combinations, developed PD symptoms 16 weeks after inoculation. European *V. vinifera* varieties exhibited significant differences in their susceptibility to  $Xf_{PD}$  (Table C.1). All varieties, however, showed PD symptoms to some extent, confirming previous field observations of general susceptibility to  $Xf_{PD}$  [261, 268, 322]. We also found significant differences in virulence ( $\chi^2 = 68.73$ ,  $df = 1$ ,  $P = 2.2 \times 10^{-16}$ ) between two  $Xf_{PD}$  strains isolated from grapevines in Mallorca across grapevine varieties (Fig. C.1). Full details on the results of the inoculation tests are available in Methods, Appendix C.1, Table C.1.

Growing degree days (GDD) have traditionally been used to describe and predict phenological events of plants and insect pests, but rarely in plant diseases [339]. We took advantage of data collated in the inoculation trials together with temperature to relate symptom development to the accumulated heat units at weeks eight, 10, 12, 14, and 16 after inoculation (Appendix C.1). Rather than counting GDDs linearly above a threshold temperature, we consider  $Xf$ 's specific growth rate *in vitro* within its cardinal temperatures. The empirical growth rates come from the seminal work by Feil & Purcell [323] shown in the inset of Fig. 7.1 a. This Arrhenius plot was transformed, as explained in Appendix C.2.1, to obtain a piecewise function  $f(T)$  Eq. (7.1). Our model and risk maps are based on  $f(T)$  (red line in Fig. 7.1 a) because it provides the best fit to the experimental data when compared with the commonly used beta function (blue line) for representing the thermal response in biological processes [340, 341]. This Modified Growing Degree Day (MGDD) profile Eq. (7.1) enables to measure the thermal integral from hourly average temperatures, improving the prediction scale of the biological process [354]. MGDD also provides an excellent metric to link  $Xf_{PD}$  growth in culture with PD development as, once the pathogen is injected into the healthy vine, symptoms progression mainly depends upon the bacterial load (i.e., multiplication) and the movement through the xylem vessel network, which are fundamentally temperature-dependent processes [323, 355].

Moreover, MGDD can be mathematically related to the exponential or logistic growth of the pathogen within the plant (Appendix C.2.2).



**Figure 7.1: Climatic and transmission layers composing the epidemiological model.** (a) *MGDD* profile fitted to the *in vitro* data of *Xf* growth rate in Feil & Purcell 2001 [323]. The original Arrhenius plot in Kelvin degrees (inset) was converted to Celsius, as explained in (Appendix C.2.1), to obtain the fit shown in the main plot red line; the blue line represents the fit with a beta function. (b) Correlation between CDD and the average  $T_{min}$  of the coldest month between 1981 and 2019. Plotted black dots (worldwide) and yellow dots (main wine-producing zones) depict climatic data from 6,487,200 cells at  $0.1^\circ \times 0.1^\circ$  resolution, spread globally and retrieved from the ERA5-Land dataset. The red solid line depicts the fitted exponential function for worldwide data and the blue solid line for main vineyard zones. (c) Nonlinear relationship between *MGDD* (red line) and CDD (blue line) and the likelihood of developing chronic infections. Black dots depict the cumulative proportion of grapevine plants in the population of 36 inoculated varieties showing five or more symptomatic leaves at each of the 15 *MGDD* levels (see Appendix C). Vertical bars are the 95% CI. (d) Combined ranges of *MGDD* and CDD on the likelihood of developing chronic infection. (e) Transmission layer in the dynamic equation (1) of the SIR compartmental model. (f) Relationship between the exponential growth of the number of infected plants with the risk index and their ranks.

Interannual infection survival in grapevines plays a relevant role when modeling PD epidemiology. In our model, we assumed a threshold of five or more symptomatic leaves for these chronic infections based on the relationship be-

tween the timing and severity of the infection during the growing season and the likelihood of winter recovery [323, 324, 327]. This five-leaf cut-off was grounded on: (i) the bimodal distribution in the frequency of the number of symptomatic leaves among the population of inoculated grapevines (Fig. C.1), whereby vines that generally show less than five symptomatic leaves at 12 weeks after inoculation remain so in the following weeks, while those that pass that threshold continue to produce symptomatic leaves; and (ii) the observed correlation between the acropetal and basipetal movement of Xf along the cane (Fig. C.1). The likelihood of developing chronic infections as a function of accumulated MGDD among the population of grapevine varieties was modeled using survival analysis with data fitted to a logistic distribution  $\mathcal{F}(MGDD)$ . A minimum window of  $MGDD = 528$  was needed to develop chronic infections (var. Tempranillo), about 975 for a median estimate, while a cumulative  $MGDD > 1159$  indicates over 90% probability within a growing season (red curve in Fig. 7.1 c and Methods).

Next, we intended to model the probability of disease recovery by exposure to cold temperatures. Previous works had specifically modeled cold curing on Pinot Noir and Cabernet Sauvignon varieties of California as the effect of temperature and duration [324] by assuming a progressive elimination of the bacterial load with cold temperatures [327]. In the absence of appropriate empirical data to formulate a general average pattern of winter curing among grapevine varieties, we combined the approach of Lieth et al. [324] and the empirical observations of Purcell on the distribution of PD in the US related to the average minimum temperature of the coldest month,  $T_{min}$ , isolines [326]. To consider the accumulation of cold units in an analogy of the MGDD, we searched for a general correlation between  $T_{min}$  and the cold degree days (CDDs) with base temperature = 6 °C (see Methods). We found an exponential relation,  $CDD \sim 230 \exp(-0.26 \cdot T_{min})$ , where specifically,  $CDD > 306$  corresponds to  $T_{min} < -1.1^\circ\text{C}$  (Fig. 7.1 b). To transform this exponential relationship to a probabilistic function analogous to  $\mathcal{F}(MGDD)$ , hereafter denoted  $\mathcal{G}(CDD)$ , ranging between 0 and 1, we considered the sigmoid family of functions  $f(x) = \frac{A}{B + x^C}$  with  $A = 9 \cdot 10^6$ ,  $B = A$ , and  $C = 3$  (Fig. 7.1 c), fulfilling the limit  $\mathcal{G}(CDD = 0) = 1$ , i.e., no winter curing when no cold accumulated, and a conservative 75% of the infected plants recovered at  $T_{min} = -1.1^\circ\text{C}$  instead of 100% to reflect uncertainties on the effect of winter curing.

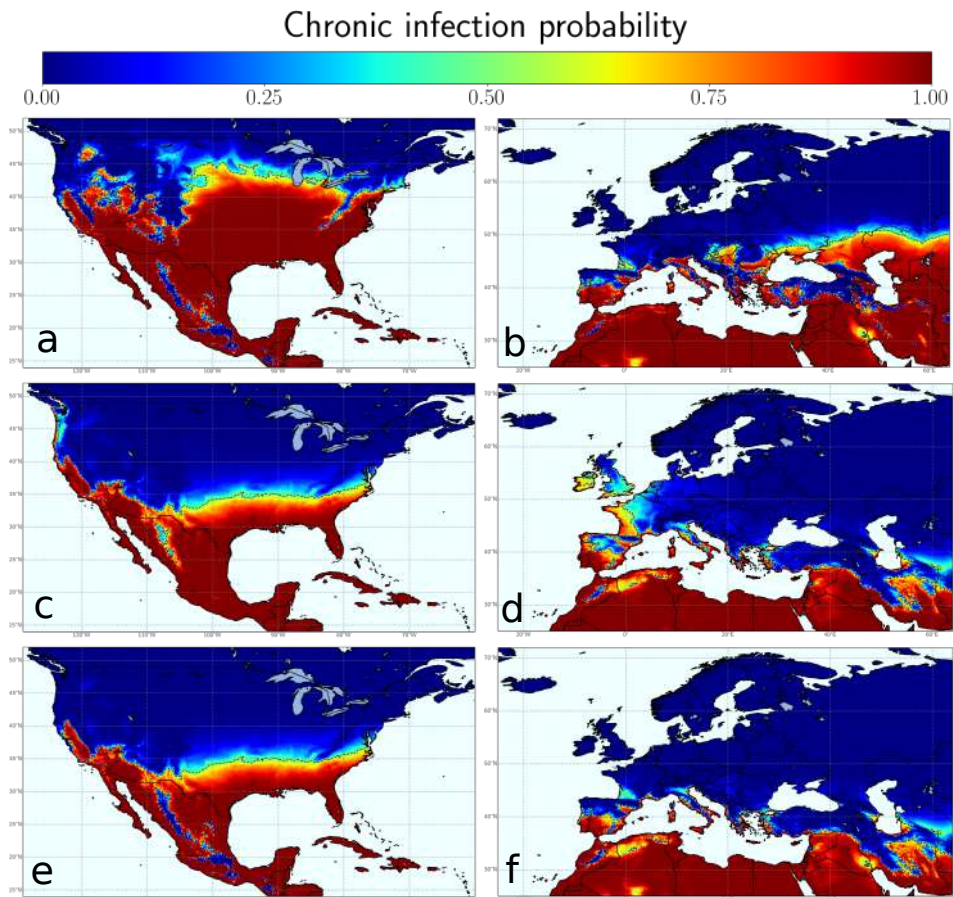
### 7.3.2 MGDD/CDD distribution maps

MGDD were used to compute annual risk maps of developing PD during the summer for the period 1981-2019 (see Methods). The resulting averaged map identifies all known areas with a recent history of severe PD in the US

corresponding to  $\mathcal{F}(MGDD) > 90\%$  (i.e., high-risk), such as the Gulf Coast states (Texas, Alabama, Mississippi, Louisiana, Florida), Georgia, and Southern California sites (e.g., Temecula Valley) (Fig. 7.2 a), while capturing areas with a steep gradation of disease endemicity in the north coast of California ( $\mathcal{F}(MGDD) > 50\%$ ). Overall, more than 95% of confirmed PD sites ( $n = 155$ ) in the US fall in grid cells with  $\mathcal{F}(MGDD) > 50\%$ . The average MGDD-projected map for Europe during 1981-2019 spots a high risk for the coast, islands, and major river valleys of the Mediterranean Basin, southern Spain, the Atlantic coast from Gibraltar to Oporto, and continental areas of central and southeast Europe (Fig. 7.2 b). Of these, however, only some Mediterranean islands, such as Cyprus and Crete, show  $\mathcal{F}(MGDD) > 99\%$  comparable to areas with high disease incidence in the Gulf Coast states of the US and California. Almost all the Atlantic coasts from Oporto (Portugal) to Denmark are below suitable *MGDD*, with an important exception in the Garonne river basin in France (Bordeaux Area) with low to moderate *MGDD* (Fig. 7.2 b).

Fig. 7.2 a shows how the area with high-risk MGDD values extends further north of the current known PD distribution in the southeastern US, suggesting that winter temperatures limit the expansion of PD northward [261]. A comparison between MGDD and CDD maps (Fig. 7.2 a vs. Fig. 7.2 c, Fig. 7.2 e) further supports the idea that winter curing is restricting PD northward migration from the southeastern US. However, consistent with growing concern among Midwest states winegrowers on PD northward migration led by climate change [356], we found a mean increase of  $0.12\% \text{ y}^{-1}$  in the areal extent with  $CDD < 306$  ( $\sim T_{\min} < -1.1 \text{ }^\circ\text{C}$ ) since 1981, comprising land areas between  $103^\circ\text{W}$  and  $70^\circ\text{W}$  of the US (Fig. C.12). Such an upward trend corresponds to  $5090\text{km}^2 \text{ y}^{-1}$  in the potential northward expansion of PD due to climate change and an accumulation of  $\sim 193420\text{km}^2$  of new areas at risk since 1981.

High CDD values would also be expected to restrict the potential PD colonization in continental Europe (Fig. 7.2 d). Unlike North America, the east-west distribution of major European mountain ranges, together with the warming effect of the Gulf Stream, decreases the likelihood of cold winter spells reaching the western Mediterranean coast.  $\mathcal{G}(CDD)$  between 100% and 95% (i.e., recovery probability  $< 5\%$  – low winter curing) are mostly prevalent below  $40^\circ\text{N}$  latitude in the southwest Iberian Peninsula and Mediterranean islands and coastlands ( $< 50\text{km}$  away). Above  $40^\circ\text{N}$  latitudes,  $CDD < 100$  are encountered mainly in the Atlantic coast and Mediterranean coast and islands (Fig. 7.2 d). In contrast, central and southeast Europe show high CDD values, likely preventing  $X_{\text{fPD}}$  winter survival on infected grapevines.



**Figure 7.2: Average thermal-dependent maps for Pierce’s disease (PD) development and recovery in North America and Europe.** PD development during the growing season based on average  $\mathcal{F}(MGDD)$  estimations between 1981 and 2019 in North America (a) and Europe (b) derived from the results of the inoculation experiments on 36 grapevine varieties. Large differences in the areal extension with favorable MGDDs can be observed between the US and Europe. The winter curing effect is reflected in the distribution of the average  $\mathcal{G}(CDD)$  for the 1981-2019 period in the United States (c) and Europe (d). A snapshot of the temperature-driven probability of chronic infection averaged for the 1981-2019 period is obtained from the joint effect of MGDD and CDD in North America (e) and Europe (f). Warmer colors indicate more favorable conditions for chronic PD, and the dashed line highlights the threshold of chronic infection probability being 0.5.

In Fig. 7.2 e-f, we show the average climatic suitability for PD establishment

only from the mechanistic relation between  $X_{f_{PD}}$  and temperature. Although all areas with current  $X_{f_{PD}}$  related outbreaks are identified, risk predictions based only on the combination of MGDD and CDD could lead to overestimation, as this approach overlooks disease transmission dynamics and climate interannual variability.

### 7.3.3 PD global risk

We ran several simulations of the model Eq. (7.7) with  $R_0$  values between 1 and 14 to validate PD spatio-temporal distribution in the US. We found  $R_0 = 8$  as the optimal parameter for maximizing the area under a ROC curve (Fig. C.5), returning an accuracy of more than 80%, except for 2006, due to data obtained from an area at the transient-risk zone (Fig. C.4 and Table 7.1). For Europe and the rest of the world, we derived an  $R_0 = 5$  as a maximal baseline estimate for modeling PD transmission (see Methods and Appendix C.2.5). These  $R_0$  values should be taken as operating estimates for the model. From the model simulations Eq. (7.7), we obtained a risk index  $r$  that measures the relative exponential growth rate in the population of infected plants at the epidemic onset with respect to the maximum growth,  $r = 1$ . This index served to rank the epidemic risk zones in high ( $> 0.9$ ), moderate (0.66–0.9), low (0.33–0.66) and very low ( $\sim 0.075 - 0.33$ ) risks (see Fig. 7.1 f, Methods, and Appendix C.2.6).

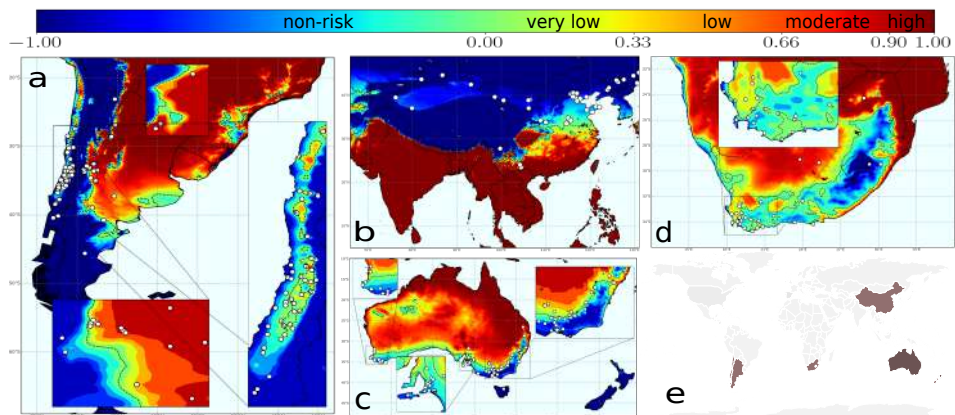
**Table 7.1: Validation of model predictions.** The items are locations where PD was present or absent. TP corresponds to true positives and TN to true negatives according to our model with  $R_0 = 8$ .

Year	Presence	Absence	TP	TN	Accuracy
2001	16	5	15	3	86%
2002	12	2	11	1	86%
2005	4	2	4	1	83%
2006	8	0	4	0	50%
2015	53	0	51	0	96%
TOTAL	93	9	85	5	88%

To date, PD is mainly restricted to the American continent, with some unrelated introductions of  $X_{f_{PD}}$  to Taiwan and Mallorca (Spain) from the United States [268, 308]. To assess the risk of PD establishment elsewhere, we projected our epidemiological model into the main wine growing regions of the Northern Hemisphere (US, Europe, and China) and Southern Hemisphere (Chile, Argentina, South Africa, Australia, and New Zealand) (Fig. 7.3 a-e). We found that emerging wine-producing areas in China are predominantly located in non-risk zones, whereas only some vineyards in the Henan and Yunnan provinces fall

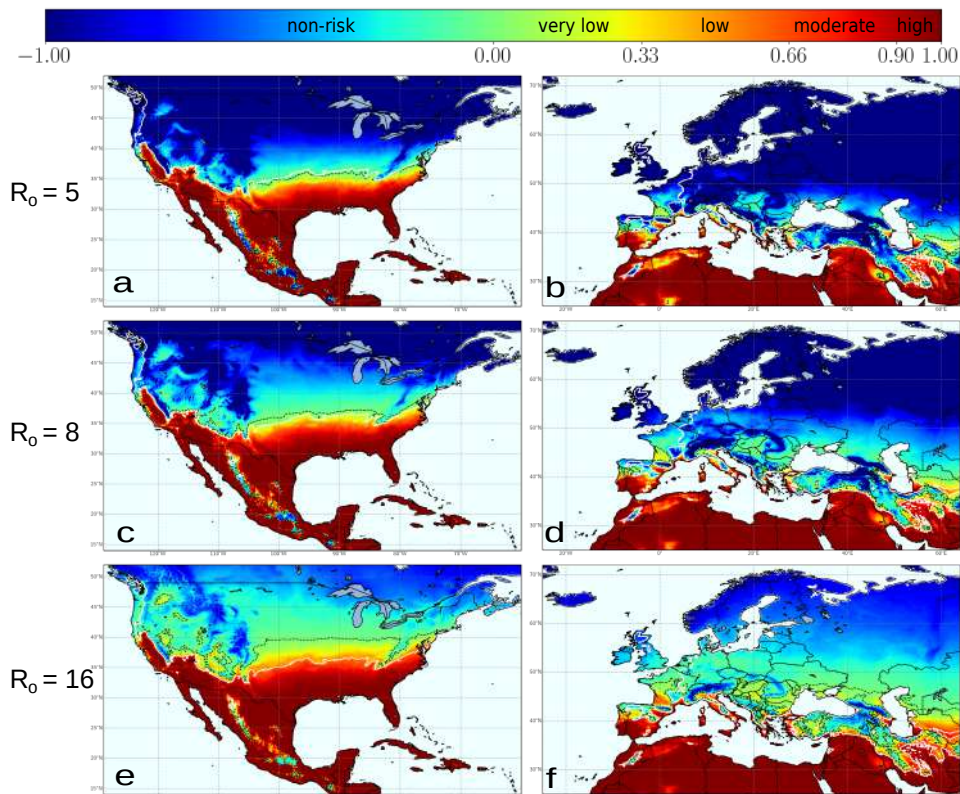
in transition and moderate-high risk zones (Fig. 7.3 b). In Europe, 92.1% of the territory is in non-risk zones and 6.1% is included in the epidemic risk zone, with only 1.9% showing a high-risk index and 1.5% a moderate risk (Table C.2). The model also reveals a progressive transition from areas without risk ( $r(t) < 0$ ) before 1990 to epidemic risk zones with low-risk indexes by 2019 ([338], see Movies), mainly affecting the basins of the rivers Po in Italy, Garonne and Rhone in France, and Douro/Duero in Portugal and Spain. This represents a mean increase of  $0.21\% \text{ y}^{-1}$  in the epidemic risk zone, a rate 3.5 times greater than that of the eastern US, which could increase the likelihood of PD establishment in Europe in the coming decades. In the US, most states around the Gulf Coast show high-risk indexes, whereas around 37.5% of California's surface is suitable for epidemics with high growth rate incidence (Table C.3).

In the Southern Hemisphere, vineyards at non-risk or transient epidemic risk zones predominate – e.g., non-risk in New Zealand and Tasmania (Fig. 7.3 c). Risk indexes in areas where PD can become established ( $r(t) > 0$ ) range from very low to low for most coastal vineyards in Australia (west, south, and east) with somehow more suitable conditions in the interior of New South Wales,



**Figure 7.3: Climate-driven risk maps for PD establishment in main viticulture regions worldwide under a baseline  $R_0 = 5$  scenario.** White dots indicate the main vineyard areas in the wine-growing regions of China and the Southern Hemisphere. (a) Chile and Argentina; (b) Asia with special attention to China; (c) Australia and New Zealand (wine areas are not marked as the whole country is without risk); and (d) South Africa. (e) Global distribution of main wine-producing areas analysed. The risk index  $r_j(t)$ , express the relative exponential growth rate of the disease incidence, and was scaled from 0.1 to 1 and ranked as very low (0.10-0.33), low (0.33-0.66), moderate (0.66-0.90) and high ( $> 0.90$ ).

Greater Perth, and Queensland (Fig. 7.3 c); a general very low or low-risk indexes are predicted in the Western Cape in South Africa (Fig. 7.3 d); overall very low but localized low to moderate risk indexes in some areas in Chile; and low to moderate growth of the number of infected vines in most of Argentina, being this the wine-growing country with the highest risk (Fig. 7.3 a). Detailed



**Figure 7.4: Temperature-driven dynamic-model simulations for PD establishment from 1981 to 2019 under different  $R_0$  scenarios with a spatially homogeneous vector distribution.** For comparison, the baseline scenario with a  $R_0 = 5$  for Europe is projected to North America (a) capturing to some extent the distribution and severity of PD in that continent. In Europe (b) high-risk areas (i.e.,  $r_j(t) > 0.90$ ) are restricted to the coastal Mediterranean and the south of the Iberian Peninsula; black dash line separate areas with  $r(t) > 0$  where theoretically PD can thrive. Under higher  $R_0$  scenarios,  $R_0 = 8$  for North America (c) and Europe (d), the dash lines tend to separate from isoline  $T_{\min} = -1.1^\circ\text{C}$  (white line); and even more in extreme transmission pressure  $R_0 = 16$  for North America (e) and Europe (f).

information on areas with non-risk, transient risk, and risk indexes (i.e., disease-incidence growth rates) in areas with the potential risk of establishment by country and regions is provided in [Table C.4](#).

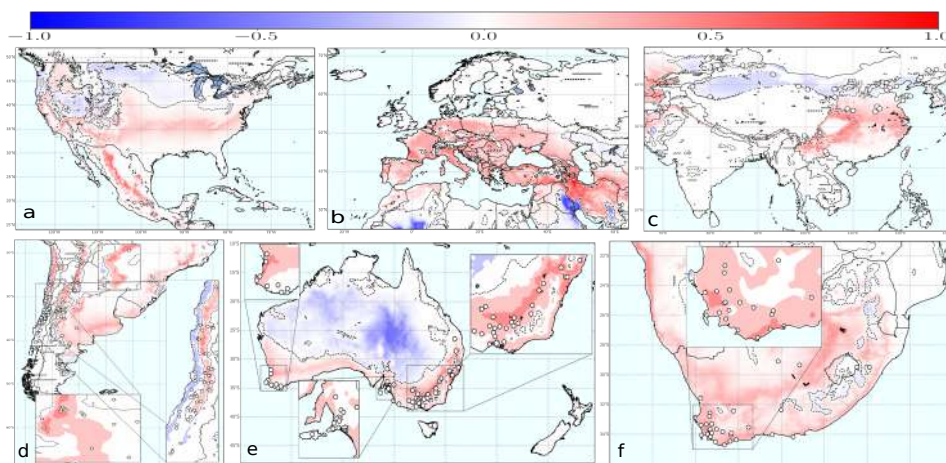
Risk indexes may vary within epidemic risk zones if any of the epidemiological parameters governing transmission change. As expected,  $I(t) < I(0)$  regions increasingly displace to northern latitudes in the US and Europe under higher transmission scenarios, increasing the risk-epidemic zones significantly ([Fig. 7.4 a-f](#)). The line representing the outbreak extinction, i.e., the non-risk zone  $r(t) < -0.09$ , in the validated  $R_0 = 8$  scenario for the US, falls at some distance above the isoline  $T_{\min} = -1.1^\circ\text{C}$  in comparison to the  $R_0 = 5$  scenario ([Fig. 7.4 c](#) vs [Fig. 7.4 a](#) and [\[338\]](#), [Movies](#)). This distribution pattern holds and moves slightly northward over time in parallel to global warming, although the trend of PD latitudinal change is moderated by high CDD values (i.e., cold accumulation). In addition, the disease extension also declines due to CDD interannual fluctuations in the simulations. Cold waves periodically occur that reach latitudes close to the Gulf, such as those that occurred in 1983, 1993, 1995, 2000, 2009, and 2013 ([Movies](#) at [\[338\]](#)), thus preventing PD expansion northward. The magnitude of this decrease is revealed after comparing the average annual increase of the areas between  $r(t) > 0$  and  $CDD < 306$  lines. From 1981 to 2019, the area with risk  $r(t) > 0$  increased at a rate of  $0.05\% \text{ y}^{-1}$ , while that of  $CDD < 306$  by  $0.12\% \text{ y}^{-1}$ , an important difference not explained alone by CDDs without considering climate fluctuations ([Fig. C.12](#)).

We checked whether using a beta function produces changes in the risk indexes with respect to the Arrhenius-based approach. Firstly, we needed to calibrate the model using the probability of developing chronic infections, as in [Fig. 7.1 c](#), with the values of MGDD obtained with the beta function. We found little differences, mainly a decrease in risk index in the transition zones between risk and non-risk zones ([Figs. C.10](#) and [C.11](#)), and non-significant differences in risk zones at the global scale.

### 7.3.4 PD risk projections for 2050

Global shifts in the risk index  $r_j(t)$  between 2019 and those projected for 2050 were calculated under the same baseline scenario ([Fig. 7.5 a-f](#), [Methods](#)). Our simulation shows a generalized increasing trend mainly due to shifts from transition zones to epidemic risk zones with very low or low-risk indexes in the main wine-growing regions, except for the US. Here the epidemic risk zone would increase by 12.8% with the greater increments in the high-risk index category (22.7%) and a decrease in the transition zones ([Table C.5](#)). Much less surface would be included in the epidemic risk zone in Europe (8.6%) compared to the US (36.5%). However, the epidemic risk zone would expand by 40.0%

with respect to 2020, a rate more than three times higher than that of the US (Table C.6).



**Figure 7.5: Global shifts in PD risk index ( $r_j(t)$ ) from 2020 to 2050.** We assume a homogeneous vector spatial distribution with  $R_0 = 5$  except for the US, where  $R_0 = 8$  has been considered. **(a)** North America; **(b)** Europe; **(c)** Asia; **(d)** South America; **(e)** Australia and New Zealand; and **(f)** South Africa. Risk-index increases are in red and decreases in blue. The dashed line represents the spatial threshold where  $r_j(t)$  difference changes from negative to positive.

**Table 7.2: Shifts in risk areas for Pierce's disease in Europe projected for 2050 under a  $R_0 = 5$  scenario.** The model was run assuming the same homogeneous spatial distribution of the vector for the whole period.

Risk	2050 km <sup>2</sup>	2020 km <sup>2</sup>	Difference km <sup>2</sup>	Difference (%)	2050 (%)	2020 (%)
No risk	8,885,300	9,334,178	-448,878	-4.8	87.6	92.1
Transition	381,081	182,872	198,208	108.3	3.8	1.8
Very low	189,025	179,225	9,799	5.5	1.9	1.8
Low	207,599	104,143	103,456	99.3	2.1	1.0
Moderate	154,780	148,621	6,159	4.1	1.5	1.5
High	322,225	190,971	131,254	68.7	3.2	1.9

Such increases are due to the emergence of previously unaffected areas in 2020 evolving into epidemic risk zones by 2050 and epidemic growth-rate increases in already epidemic risk zones in 12 of 42 countries (Table C.2). Among these 12 countries, however, there is substantial variation in the risk index in-

crements within epidemic risk zones with respect to 2019 (Table C.6). While non-risk zones still cover 87.6% of Europe's land area, epidemic risk zones with high-risk indexes are expected to be almost two-fold higher than that of 2019, comprising 3.2% of Europe (Table 7.2). Nevertheless, this is based on a simplistic linear extrapolation of *MGDD* and *CDD* values, and the potential effect of climate change on the vector's distribution is not considered. In Chapter 8, we explore the potential impact of climate change on the risk of PD establishment in Europe, taking into account different global warming scenarios.

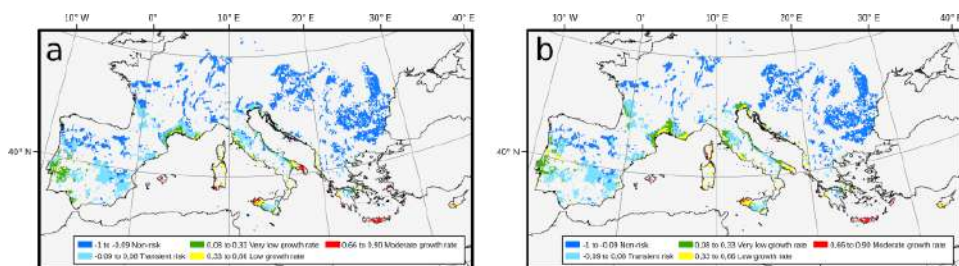
### 7.3.5 Risk based on vector information

So far, we have ignored the distribution of known and potential vector species due to their large number in the Americas and the limited quantitative information generally available. In the case of Europe, given *P. spumarius* prevalence as a potential vector and its wide distribution, we added a vector layer in a spatially dependent  $R_0(j) = R_0^{\max} v(j)$ , where  $v(j)$  is the climatic suitability for the vector (Methods);  $v = 1$  implies optimal climatic conditions with no constraints for the vector population size, while  $v = 0$  implies unsuitable climatic conditions and its absence (Fig. C.7). According to the model, no European zone shows a high-risk index, and barely 0.34% of the territory falls in areas with potential moderate exponential growth rates in disease incidence (Table C.7). Irrespective of vineyard distribution, we estimated that PD could potentially become established (i.e.,  $r(t) > 0$ ) at a maximum of 3.1% of the territory, while the area at moderate-risk index would be 5-times lesser than the model without the vector's climate suitability layer, this latter more in consonance with other proposed risk maps [329, 330]. Such differences in the projected risks are mainly concentrated in the warmest and driest Mediterranean regions and are due to uncertainties concerning temperature-humidity interactions in the ecology of the vector [320].

### 7.3.6 Combining vineyard land cover across Europe with the model output

When we integrate into the model a layer of vineyard surface from Corine-Land-Cover, we find that PD could potentially become established (i.e.,  $r(t) > 0.075$ ) in 22.3% of the vineyards in Europe. However, no vineyard is in epidemic risk zones with a high-risk index, and only 2.9% of the vineyard surface is at moderate risk (Table C.8). The areas with the highest risk index ( $r(t)$  between 0.70 and 0.88) are mainly located in the Mediterranean islands of Crete, Cyprus, and the Balearic Islands or at pronounced peninsulas like Apulia (Italy) and Peloponnese (Greece) in the continent (Fig. 7.6 a and Table C.8). Most vineyards are in non-risk zones (42.1%), whereas 35.6% are located in transition zones

that are presently non-risk but where  $Xf_{PD}$  could become established in the next decades, causing some sporadic outbreaks. In [Tables C.4](#) and [C.8](#), we provide full details of the total vineyard areas currently at risk for each country and region.



**Figure 7.6: Intersection between Corine-land-cover vineyard distribution map and PD-risk maps for 2020 and 2050.** Data were obtained from Corine-land-cover (2018) and the layer of climatic suitability for *P. spumarius* in Europe from [320]. The surface of the vineyard contour has been enlarged to improve the visualization of the risk zones and disease-incidence growth-rate ranks. (a) PD risk map for 2019 and its projection for 2050 (b). Blue colors represent non-risk zones and transient risk zones for chronic PD ( $R_0 < 1$ ). The 2050 map shows some contraction of epidemic risk zones with moderate risk indexes in Mediterranean islands and Apulia (Italy) as the climate becomes hotter and dryer.

Our model with climate and vector distribution projections for 2050 indicates a 55.8% increase in the epidemic risk zone in Europe ([Fig. 7.6 b](#)). This increment would be mainly due to the extension of epidemic risk zones with very low and low-risk indexes. However, within the epidemic risk zones, areas with moderate risk indexes would decrease from 114925ha in 2020 to 43114ha in 2050, and no vineyards would be at high risk ([Fig. 7.6 b](#); see [Table C.9](#)). Counterintuitively, our model indicates a substantial increase in the area where PD could establish and become endemic for 2050, but a moderate decline in those areas where crop damage could be expected to be significant (e.g., Balearic Islands, Crete, Cyprus, Apulia).

## 7.4 Discussion

We introduce an epidemiological approach to assess the risk of PD establishment and epidemics in vineyards worldwide. The model includes the dynamics of the infected-host population, which enables estimating the initial exponential growth/decrease rate of the disease incidence. Unlike SDM correlative studies,

Bayesian or, in general, machine learning black-box approaches, our model goes beyond by providing a mechanistic framework and thus explanatory power. In addition, it is flexible enough to simulate different climate and transmission scenarios, allowing, for instance, the incorporation of information on the spatial distribution of the vector. Comprehensive global PD risk maps result from the model simulations with historical climatic data. A web page is included, showing simulations with different parameters to estimate the risk of PD anywhere [338].

Temperature regulates key physiological processes of the ectothermic organisms involved in PD and thus limits the thermal range in which they can thrive [334].  $Xf_{PD}$  multiplication and survival within vine xylem vessels not only characterize PD but also determine the bacterial population dynamics [323, 355]. PD symptom development can therefore be characterized as a thermal-dependent continuous process within the range of  $Xf_{PD}$ 's cardinal temperatures [335]. The combination of MGDD metrics with robust experimental data provides a reliable predictor of climatic suitability and the probability of developing PD during the summer, whereas CDD accounts for the effect of cold-temperature exposure on infected-plant recovery. This opposite contribution of MGDD and CDD in the demography of infected plants shapes the impact of climate variability on the epidemic dynamics in the early stages of the invasion (Fig. 7.1 d). Given that the physiological basis of the plant- $Xf$  interaction leading to symptom development is poorly understood, we caution that other environmental factors, such as drought, nutrient status, or crop management, may modulate symptom expression and hence add an error in the MGDD parameter not measured in this work. Nonetheless, we deem the error range would be smaller than the differences in the accumulated *MGDDs* needed to reach the same disease level among varieties (i.e., regional differences) and smaller than the interannual *MGDD* oscillations found in most locations. In addition, our model is general enough to allow for other functions or adjustments of the relationship between  $Xf_{PD}$ 's growth rate and temperature in vitro as better experimental data become available. However, we deem that the differences in the risk indices would vary very little in risk zones, as we observe in PD risk maps for Europe when a beta function is applied instead of the Arrhenius-based approach to adjust the MGDD (Figs. C.10 and C.11).

Knowledge of insect distribution is crucial for predicting epidemic outbreaks of endemic diseases, as well as the risk of invasion by emerging vector-borne pathogens ([259, 357]; cf. [235]). Given the great diversity of known and potential vectors that can transmit PD [128], it has not been possible to include each region's particular vectors in the model. Therefore, when evaluating the risk of PD on a global scale, we have considered a homogeneous spatial distribution of the vector (fixed  $R_0$ ), except in Europe, where there is information on the

main vector (Fig. C.7). As expected, the European case shows how models that assume a homogeneous spatial distribution of the vector generally produce epidemic risk zones with higher risk indexes than models that include a heterogeneous spatial distribution (Table C.2 vs. Table C.7). This lack of information about vectors is one of the main reasons why the risk of vector-borne plant diseases is often overestimated.

Risk overestimation may involuntarily stem from other additional sources too. Using mean data as inputs in epidemiological models can lead to biased results when response functions are nonlinear and climate variability is not accounted for [335]. This study presents experimental evidence of a non-linear relationship between MGDDs and PD chronic infections and indirect empirical evidence of a non-linear relationship between CDDs and PD recovery (Fig. C.13). Such a non-linear response consequently greatly impacts reducing the risk of PD establishment and steeping the spatial gradients in risk maps (Figs. 7.4 and 7.6). Moreover, MGDDs and CDDs might help to explain why disease pressure is much higher in the southeastern US than in California and Europe (Figs. 7.2 and 7.4) or, for example, earlier reports of PD outbreaks in Kosovo [358]. Cooler summer nights in California and a shorter growing season compared to those found in the Gulf states in the southeastern US explain the difference in the accumulated *MGDD* for both areas. In the case of Kosovo, *CDD* values above certain thresholds could have led to the extinction of incipient outbreaks driven by several years with *MGDD* in the conducive range of PD (Fig. 7.2).

Our PD risk map for Europe confirms previous predictions for the subsp. *fastidiosa* from SDMs [329]. Both approaches make congruent predictions on PD potential distribution, providing convergent lines of independent evidence of climate suitability. However, our risk maps go further by incorporating in the epidemic risk zones information on the relative exponential growth rates in the potential disease incidence. In general terms, the epidemic risk map including vector information indicates a low risk for chronic PD. Only ~ 0.34% of European vineyard surface, mainly located in Cyprus, Crete, Sardinia, part of Sicily, and the Balearic Islands, meet climatic conditions for PD to become endemic and cause significant damage (Table C.7). Other regions, such as Bordeaux, Portugal, the Rhône Valley, and the Veneto region, would be included in epidemic risk zones but with very low to low exponential growth rates in disease incidence. By contrast, notorious wine-growing regions in Spain (e.g., Rioja, Ribera del Duero), France (e.g., Burgundy), and Italy (e.g., Piedmont) currently fall within areas considered non-risk zones, transient-epidemic zones, or epidemic risk zones with very low risk indexes (Fig. 7.6).

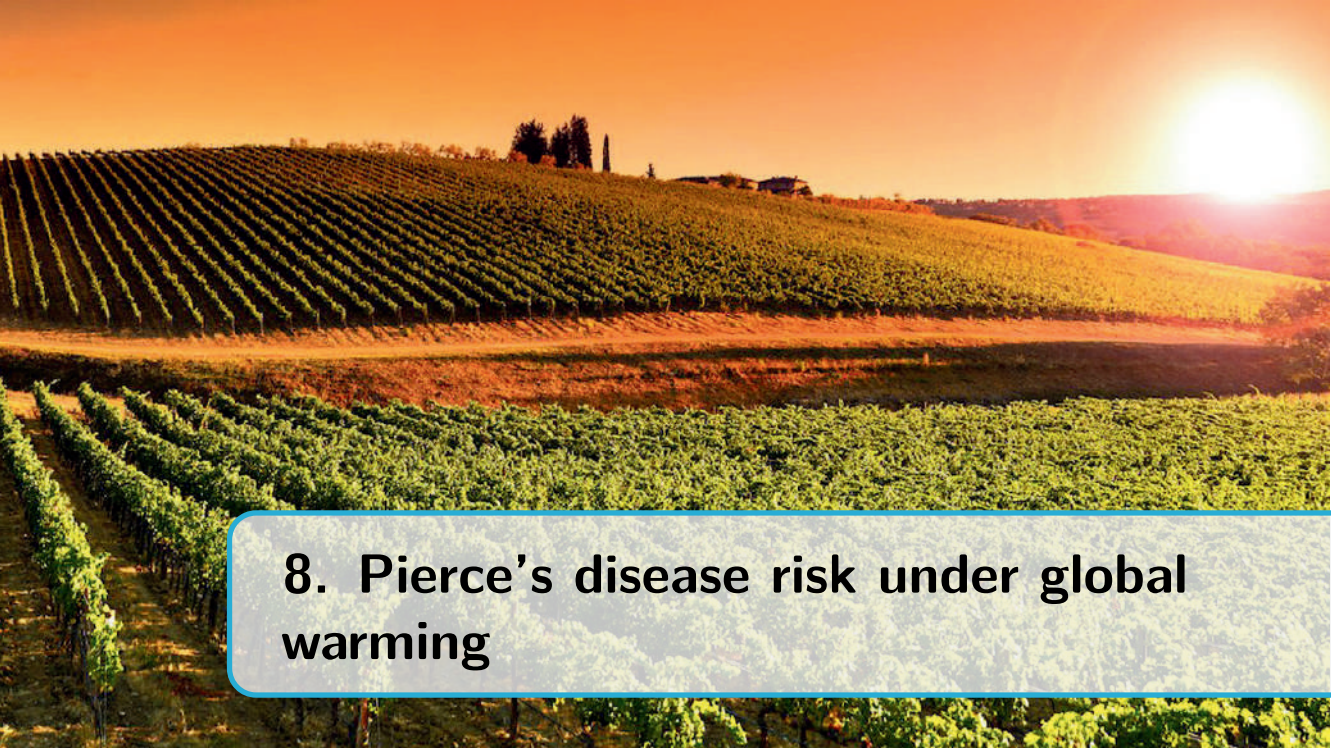
The dynamic nature of the simulation outputs already points to a progressive

global increase in the extension of PD epidemic risk zones ( $r(t) > 0$ ) in the last decade, irrespective of vineyard distribution (see movies on [338]). This is even more accentuated in the model projections for 2050, which point out a global expansion of PD epidemic risk zones at different velocities among continents due to climate change (Fig. 7.5). For example, many important viticulture areas in western Europe included in non-risk or transition zones before 1990 are progressively shifting to hotter summers and milder winters and hence would be increasingly suitable for the disease within the extrapolated current scenario. This is further illustrated by a 40% increase of the potential epidemic risk zone by 2050 concerning 2020 for Europe and more moderate increases in the United States and the Southern Hemisphere (Fig. 7.5). Nonetheless, our model projection for 2050 that includes spatial heterogeneity in the vector distribution, as in Europe, would indicate lower transmission rates because global change is predicted to have negative effects on *P. spumarius* abundance in Europe [320, 359]. At the global scale, there is certainly scientific consensus that climate change will follow a general pattern summarized in the paradigm "dry gets drier, wet gets wetter" [360]. In agreement, our model projection for PD on the vineyards of Mallorca (Spain) suggests shifts to slightly less favorable conditions for  $Xf_{PD}$  transmission and an expected progressive decrease in the impact of the disease by 2050. This example and others in Mediterranean islands advocate for certain caution when interpreting climate change projections, especially in other Mediterranean climates of the world, where the complex interactions between humidity and temperature can limit the presence and abundance of vectors (Fig. C.7).

The scope of our study excludes location-specific complexities surrounding PD ecology due to scale limitations. The spatial distribution of the vector is considered only for the *V. vinifera*- $Xf_{PD}$ -*P. spumarius* pathosystem in Europe, so  $R_0$  estimations could locally differ in other wine-producing regions elsewhere (Fig. 7.3). Disease incidence thus could locally vary where the climate is conducive to PD. Such variation is because transmission rates tend to increase exponentially rather than linearly under environmental conditions favoring vector abundance [328], as has been observed at a local scale on vineyards in Mallorca [268]. Our study also does not contemplate likely changes within the PD pathosystem. To date, PD is caused by  $Xf_{PD}$  (i.e., ST1/ST2), but other genotypes of the subsp. *fastidiosa* or other subspecies and their recombination could arise in the future with different ecological and virulence traits [262]. On the other hand, new vector species could be accidentally brought in [128], as exemplified with the introduction of the glassy-winged sharpshooter (*Homalodisca vitripennis*) in California, modifying transmission rates and disease incidence in new areas [277]. To capture these uncertainties in relation to the vector, we

have performed simulations with  $R_0 = 8$  and  $R_0 = 16$  (Fig. 7.4). Remarkably, a comparison of PD risk maps for Europe with different  $R_0$  suggests for non-Mediterranean areas the need to stress more surveillance on the introduction of alien vectors rather than in the pathogen itself. This is because, under the current scenario ( $R_0 = 5$ ) with *P. spumarius* as the main vector, most of the non-Mediterranean vineyards would not support the establishment of PD, but the introduction of new insect vectors with greater transmission efficiency ( $R_0 = 8$ ) could compensate for the climatic layer and increase the risk index above 0. In addition, differences in grapevine varietal response alongside virulence variation among Xf strains may slightly modify PD thermal tolerance limits and therefore locally modulate epidemic intensity (see details in Appendix C). Such an effect could be seen with cv. Tempranillo, a widely planted variety in northern Spain (Table C.1); the rate of symptom progress and systemic movement is higher than the average varietal response to Xf<sub>PD</sub> (i.e., lower MGDD), which in addition might imply higher survival rates. This point calls for further testing in the field.

Our model partially explains why PD has not become established in continental Europe and other main wine-growing regions worldwide during the last 150 years, in contrast to other exotic diseases and pests brought in with native vines from the US [303–306]. We suggest that the underlying causes of this low-invasiveness risk in Europe are fundamentally two: (i) low climatic suitability for chronic PD and (ii) a climatic mismatch between environment conditions suitable for both the vector and the pathogen and their interplay in disease dynamics, similar to the situation recently described for the *V. vinifera*-Xf<sub>PD</sub>-*P. spumarius* pathosystem in northern California [287]. Currently, suitable conditions for the pathogen's invasion mostly concur in Mediterranean islands and coastlands. Likewise, similar results would be expected in other Mediterranean climates of the main winegrowing regions of the Southern Hemisphere if a vector spatial distribution layer is incorporated in the model simulations (see [338]). Finally, although increasing global warming will extend epidemic risk zones in all continents, some caution is recommended to not incur risk overestimation, as we show in the PD risk projections for 2050 in Europe when taking into account the vector spatial distribution; complex interactions between temperature and humidity in the ecology of the vectors may have a great effect in their distribution, abundance, and thus transmission capacity [320]. There is an urgent need to fill the knowledge gap on the ecophysiology for each potential vector to downscale PD model predictions to local and regional situations.



## 8. Pierce's disease risk under global warming

### Published as:

À. Giménez-Romero, M. Iturbide, E. Moralejo, J. M. Gutiérrez, and M. A. Matías, "Global warming significantly increases the risk of Pierce's disease epidemics in European vineyards", [Scientific Reports 14, 9648 \(2024\)](#)

## 8.1 Introduction

Climate change is widely recognized as an important driver of shifts in the distribution and prevalence of plant diseases worldwide [242, 361–365]. Although the impact of climate change on the distribution of plant diseases has been approached from various perspectives [366, 367], few studies have considered epidemiological dynamics in climate projections [368, 369]. Modeling disease epidemics is a complex task, as they are emergent phenomena resulting from non-linear interactions between disease components. In addition, many of the processes involved in disease development also exhibit non-linear responses to changes in environmental variables [335, 370]. This complexity is further exacerbated in the case of vector-borne plant diseases [259]. While climate primarily determines the potential geographic range of each organism in the pathosystem, the development of epidemic outbreaks depends on favorable host-pathogen-vector-climate interactions that drive transmission chains. Consequently, modeling the risk of vector-borne plant diseases implies delimiting their epidemiological niche rather than the ecological niche of their parts, as is commonly done.

The emergence of *X. fastidiosa* in Europe has renewed interest in modeling vector-borne plant diseases, particularly for the risk that Pierce's Disease (PD) poses to the European wine industry. Despite recent studies agreeing that the current risk of PD establishment in Europe is primarily confined to the Mediterranean basin [280, 330, 371], its potential future progression is not yet clear. Some efforts have been made to characterize the geographical distribution of Xf-induced diseases in Europe under climate change, but these are limited to the use of species distribution models (SDMs) for the pathogen [235, 372, 373] and the vector [371], which have led to conflicting results. On the one hand, higher temperatures are expected to promote bacterial growth in susceptible crops in continental southern Europe, while on the other hand, these areas are progressively experiencing drier environmental conditions detrimental to vector populations [371]. Furthermore, the use of SDMs to predict the potential distribution of vector-borne plant diseases, while capable of providing good approximations, is generally inadequate. The observed distribution of the pathogen cannot be separated from that of the vector, especially in the case of obligate pathogens such as Xf. Furthermore, the potential distribution of the pathogen or the vector alone has no epidemiological meaning, which implies that they do not provide quantitative predictions of the severity of the disease. In addition, a larger set of available climate models is desirable to properly deal with the inherent uncertainty in the predictions.

To overcome these limitations, here we use a novel climate-driven epidemiological model of PD epidemic risk [280]. The model determines the spatio-

temporal epidemic risk based on the spatial distribution of vectors, temperature-dependent bacterial growth and survival within hosts, and subsequent epidemiological dynamics. The model forces the introduction of the pathogen and examines whether the disease can establish and spread from previous states under climatic conditions of the location. In [280], the model was used to determine the risk of PD under current climatic conditions in wine-growing regions worldwide, while a linear regression was used to obtain a crude first estimate of the risk in the future. This simple estimation is not expected to be reliable because it overlooks the role of nonlinearities in the model and also does not take into account climate change scenarios (as described in the manuscript). To assess the potential distribution and relative impact of PD on European vineyards under different levels of global warming, here we used state-of-the-art regional climate projections from the EURO-CORDEX initiative [374]. Our study takes into account uncertainties in climate projections and provides an updated and comprehensive assessment of PD risk in European wine-growing regions, addressing previous limitations. We posit that pest risk maps constructed from projections of epidemiological models driven by climate data provide more realistic, quantitative, and explanatory predictions than correlative and probabilistic models. Additionally, they offer valuable insights for anticipating and managing the potential impacts of PD and thus ensuring the resilience of viticulture despite future climate challenges.

## 8.2 Methods

### 8.2.1 Climate datasets

We used E-OBS version v21e [375] as the reference observational climatic dataset, providing daily gridded data for Europe at a resolution of 0.1 degrees ( $\sim 10$  km). Maximum and minimum temperature data was used to compute the MGDD and CDD indices involved in the growth and survival processes of the  $X_{fPD}$  pathogen (see “Climate-driven epidemiological model” section below). To calibrate the distribution models of *P. spumarius* capturing the widest possible range (North America and Europe), we used the ERA5-Land reanalysis [376] due to its global (land) coverage and high resolution (0.1 degrees, as E-OBS). Daily precipitation and daily minimum and maximum temperature data were retrieved to calculate the moisture index and maximum temperature during spring index required for the vector suitability model (see “Vector suitability” section below). Historical and future projections of both indexes were calculated using regional climate simulations from the state-of-the-art large high-resolution (0.11 degrees) ensemble provided by EURO-CORDEX [377]. This dataset includes daily simulations of precipitation and temperatures from a large ensemble of

Regional Climate Models (RCMs) driven by Global Climate Models (GCMs) from the CMIP5 project [378]. For this, we considered the RCP8.5 simulations for 40 combinations of GCMs-RCMs (Table 8.1). In order to calculate 20-year mean climatic indexes across the different global warming levels (+1.5°C, +2°C, +3°C, and +4°C), we relied on the time periods during which each CMIP5 driving model reaches the designated level within the RCP8.5 scenario (see [379]). This information is available at the IPCC WGI Atlas GitHub repository [380].

**Table 8.1: EURO-CORDEX GCM-RCM combinations used in this study. Numbers indicate the number of runs in each combination.**

	CNRM-CM5	EC-EARTH	HadGEM2-ES	IPSL-CM5A-MR	MPI-ESM-LR	NorESM1-M
CLMcom-CCLM4-8-17_v1	1	1	1		1	
DMI-HIRHAM5_v2	1	1	1			
GERICS-REMO2015_v2	1					
IPSL-WRF381P_v2	1					
KNMI-RACMO22E_v2	1		1			
SMHI-RCA4_v1	1	2	1	1		1
CLMcom-ETH-COSMO-crCLIM-v1-1_v1		2	1		1	1
DMI-HIRHAM5_v1		1		1	1	
IPSL-WRF381P_v1		1	1	1		1
KNMI-RACMO22E_v1		2		1	1	1
MOHC-HadREM3-GA7-05_v1		1	1		1	1
GERICS-REMO2015_v1				1		1
MPI-CSC-REMO2009_v1					1	
SMHI-RCA4_v1a					1	
DMI-HIRHAM5_v3						1

### 8.2.2 A climate-driven epidemiological model

For the sake of clarity, we review the model developed in Chapter 7. Our model describes the initial exponential rise (or decrease) of infected plants at the onset of an epidemic based on the spatial distribution of the vector and the bacterial growth and survival processes mediated by temperature. The density of vectors at a given cell controls the number of new plants that will be inoculated with the bacterium, while the local temperature mediates the growth and survival processes of the in-plant bacterial population, leading the initial inoculation to an infection or not. These temperature-driven growth and survival processes are described with the accumulation of two metrics denoted *Modified Growing Degree Days* (MGDD) and *Cold Degree Days* (CDD). The base function to compute the MGDD is proportional to the  $X_f$  temperature-dependent growth rate and is defined by

$$f(T) = \begin{cases} 0 & \text{if } T < T_{\text{base}} \\ m_1 \cdot T - b_1 & \text{if } T_{\text{base}} \leq T < T_1 \\ m_2 \cdot T + b_2 & \text{if } T_1 \leq T < T_{\text{opt}} \\ m_3 \cdot T + b_3 & \text{if } T_{\text{opt}} \leq T < T_2 \\ m_4 \cdot T + b_4 & \text{if } T_2 \leq T < T_{\text{max}} \\ 0 & \text{if } T \geq T_{\text{max}} \end{cases}$$

where  $T_{\text{base}} = 12^\circ\text{C}$ ,  $T_1 = 18$ ,  $T_{\text{opt}} = 28^\circ\text{C}$ ,  $T_2 = 32$  and  $T_{\text{max}} = 35^\circ\text{C}$ ; the slopes are  $m_1 = 0.66$ ,  $m_2 = 1$ ,  $m_3 = -1.25$  and  $m_4 = -3$  and the intercepts are  $b_1 = -8$ ,  $b_2 = -14$ ,  $b_3 = 4$  and  $b_4 = 105$ . MGDD are then computed as

$$MGDD(t) = \frac{1}{24} \sum_{\tau \in t} f(T(\tau)),$$

where  $\tau$  is expressed in hours,  $t$  in years, and we divide by 24 to obtain  $MGDD(t)$  in degree days. The accumulation period goes from the 1<sup>st</sup> of April to the 31<sup>st</sup> of October in the Northern Hemisphere and from the 1<sup>st</sup> of November to the 31<sup>st</sup> of March in the Southern Hemisphere.

CDD are computed between 1<sup>st</sup> November and 31<sup>st</sup> March in the Northern Hemisphere and between 1<sup>st</sup> April and 31<sup>st</sup> October in the Southern Hemisphere as

$$CDD(t) = \frac{1}{24} \sum_{\tau \in t} (6 - T(\tau)) \quad \text{for } T_i \leq 6^\circ\text{C}.$$

Altogether, the number of infected hosts is described by the following recurrence relation,

$$I(t) = I(t-1)e^{\gamma(R_0-1)} \mathcal{F}(MGDD(t)) \mathcal{G}(CDD(t)),$$

where  $\gamma$  is the death rate of infected vines,  $R_0$  is the basic reproduction number of the disease, and  $\mathcal{F}(\cdot)$  and  $\mathcal{G}(\cdot)$  are sigmoidal-like functions that relate the MGDD and CDD metrics to the probability of developing an infection from a given inoculation.

To incorporate the spatial heterogeneity of the vector population,  $R_0$  in each cell  $j$  can be related to the climatic suitability of the vector, such that

$$R_0^j = R_0^* \cdot s_j.$$

For Europe, we use the climatic suitability of the main vector, *Philaenus spumarius*, and  $R_0^* = 5.0$  for all simulations. For the United States, we use a homogeneous basic reproductive number  $R_0 = 8.0$  in all cells, as this choice of the parameter reproduces spatio-temporal data on PD epidemics in the US with an area under the curve of  $\sim 90\%$ . For the rest of the zones, we use  $R_0 = 5$  because of the lower presence of vectors compared to the United States.

We use  $\gamma = 0.2$ , as previously estimated in [130], and the specific form of  $\mathcal{F}(\cdot)$  and  $\mathcal{G}(\cdot)$  is given by

$$\mathcal{F}(x) = \frac{1}{1 + e^{-0.012(x-975)}} \quad (8.1)$$

$$\mathcal{G}(x) = \frac{2 \cdot 10^7}{2 \cdot 10^7 + x^3} \quad (8.2)$$

Finally, the risk index is derived as the effective growth rate of the infected population over the simulated time,

$$r_j = \max \left\{ -1, \frac{\ln(I_j(T)/I(0))}{\gamma(R_0 - 1) \cdot T} \right\}. \quad (8.3)$$

Because the typical timescale of the disease is 5 years ( $1/\gamma$ ), we simulate periods of 7 years. If more years are available to simulate, we perform a re-introduction of the disease as a single infected plant in each cell after each 7-year period.

The code used to run the model is freely accessible at GitHub [381].

### 8.2.3 Model adaptation to daily temperature data

We used the model developed in Chapter 7 [280], in which MGDD and CDD metrics were defined using hourly temperature data (Eqs. (7.1) and (7.2)). However, the E-OBS and CORDEX datasets only provide daily granularity. To overcome this limitation, we use a basic sinusoidal extrapolation relating maximum and minimum daily temperature to hourly temperatures,

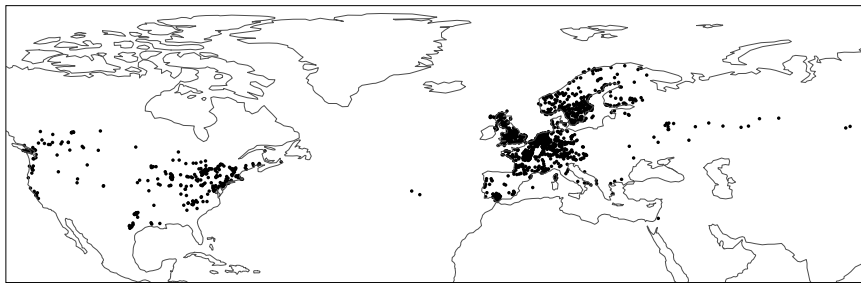
$$T_h = \frac{T_{max} + T_{min}}{2} + \frac{T_{max} - T_{min}}{2} \sin(w \cdot h), \quad (8.4)$$

with  $w = 2\pi/24$  and  $h$  ranging from 0 to 23. This approximation was validated with data from the national meteorological agency in Spain (AEMET). Basically, we used hourly temperature data obtained from 50 meteorological stations in the period 2010-2020 and computed both MGDD and CDD using the full hourly data and only the daily maximum and minimum temperatures, in the latter case using Eq. (C.31). The results showed no differences between hourly or daily temperatures' computation to estimate MGDD and CDD (Figs. C.18 and C.19). Because the temporal resolution of the E-OBS and ERA-5 land data sets is not the same and the data are acquired using different methodologies, we evaluated the possible divergence between the MGDD and CDD estimates [280]. These metrics calculated with both data agreed, showing a mean difference of 54 and 17 units for MGDD and CDD, respectively, and a standard deviation of 200 units for both metrics (Fig. C.20).

### 8.2.4 Vector climatic suitability

Following [371], we used the MaxEnt [382] algorithm to calibrate the relationship of *P. spumarius* global occurrence (predictand) with moisture index and maximum temperatures during summer index (predictors) estimated from 2003 to 2022. Data of the presence records of *P. spumarius* were obtained from The Global Biodiversity Information Facility (GBIF) [80, 383] and different Spanish plant protection agencies and research institutions ("Instituto de Ciencias

Agrarias” at CSIC, Madrid, Spain; “Servicio de Sanidad Vegetal de la Junta de Andalucía” based in Sevilla and Jaén, Andalucía, Spain; Sanidad Agrícola Econex S.L. based in Murcia, Spain), as reported in [371]. A total of 1652 presence records were used (Fig. 8.1), ensuring that there were no duplicated records within each cell of the climate layer grid.



**Figure 8.1: Training presence records for modeling the distribution of *Philaenus spumarius*.**

In addition, we randomly generated pseudo-absences, also known as background points, using “The Three-Step” method proposed in [82]. This method incorporates a model performance criterion to determine the optimal sampling background extent, thereby ensuring that the model fitting was not adversely affected by the pseudo-absence sampling. Nevertheless, we accounted for the potential variability introduced by randomly selecting points from the background by performing 10 realizations of this sampling process. A total of 4956 pseudo-absences (three times the number of presences) were used in each realization.

Model evaluation was performed using a  $k$ -fold cross-validation approach (where  $k = 10$ ) and the resulting AUCs (Area Under the ROC Curve) consistently exceeded 0.9 within the range of 0 to 1, with a value of 1 indicating perfect prediction and 0.5 indicating no discriminatory power (i.e., random guessing). Finally, the calibrated models were used to predict the suitability of *P. spumarius* in the reference historical period (2003-2022) and under increasing global warming scenarios (panels b, d, f, and h in Fig. 8.2).

### 8.2.5 Risk velocity

To assess the dynamic nature of the risk index and its spatial propagation, we introduced the concept of risk velocity, a metric analogous to the recently proposed concept of climatic velocity [21]. The risk velocity represents the rate at which the risk index changes over time and spreads across different locations. From an epidemiological perspective, risk velocity can be thought of as the

speed and direction the host would need to move to maintain its current risk conditions under climate change. Risk velocities were defined following the definition of climate velocity as the ratio of the risk temporal trend and the risk spatial gradient in each cell. Thus, the units for the risk velocity correspond to kilometers per year ( $km/year$ ). Risk velocities were computed using the VoCC R package [384, 385].

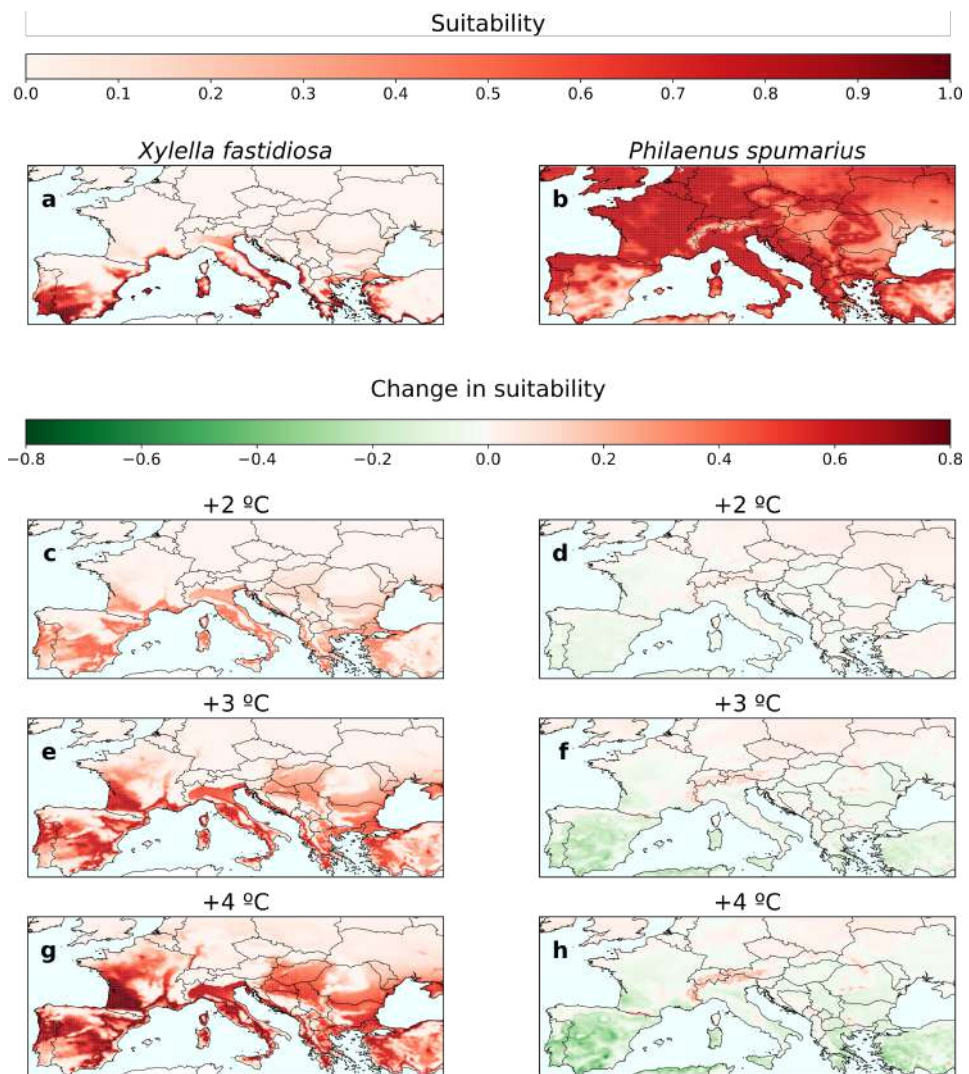
## 8.3 Results

### 8.3.1 Present and future climate suitability

To gain a deeper understanding of how climate change affects each component of the pathosystem, we performed a separate analysis of climatic suitability conditions for  $Xf_{PD}$  and *P. spumarius*. The thermal dependence of  $Xf_{PD}$  growth and survival within the infected vine was mechanistically modeled by probability functions relating the accumulation of modified degree days (*MGDD*) and cold degree days (*CDD*) to symptom development and recovery,  $\mathcal{F}(MGDD)$  and  $\mathcal{G}(CDD)$ , respectively (see Methods). Climatic suitability for pathogen establishment was then determined by  $\mathcal{F}(MGDD) \cdot \mathcal{G}(CDD)$ , i.e., the overall probability of symptom development during the growing season and subsequent survival for overwintering infection (see Methods). For *P. spumarius*, climatic suitability was modeled using an SDM based on a previous study [371], with the climatic moisture index [386] and spring maximum temperatures as key predictors (see Methods). Both analyses were evaluated under current (2003-2022) and future climate conditions, considering scenarios of increasing global warming (+1.5°C, +2°C, +3°C, and +4°C) based on the latest generation of regional climate projections covering Europe [374] (see Methods).

Progressive global warming increases the accumulation of *MGDD* during the growing season and reduces the recovery rate (i.e., *CDD*) during winter, thus favoring the geographic expansion of the pathogen (Fig. 8.2 and Fig. C.14). Conversely, increasing temperatures tend to reduce the climatic suitability of vectors in more arid areas of southern Europe, leading to a progressive migration to higher areas and latitudes in continental regions in search of climatic refuge. These general trends hold for both organisms under the +2, +3, and +4 °C temperature increase scenarios (Fig. 8.2).

The mechanistic approach to modeling pathogen establishment risk enables each of the two opposing directional processes of growth and survival (*MGDD* vs. *CDD*) to be appropriately weighted in the final result. For example, the Bordeaux region in western France has not been at risk due to low cumulative *MGDD* and low winter protection effect. In the transition from the +1.5°C scenario to the +4°C scenario, this area will experience a spectacular increase in



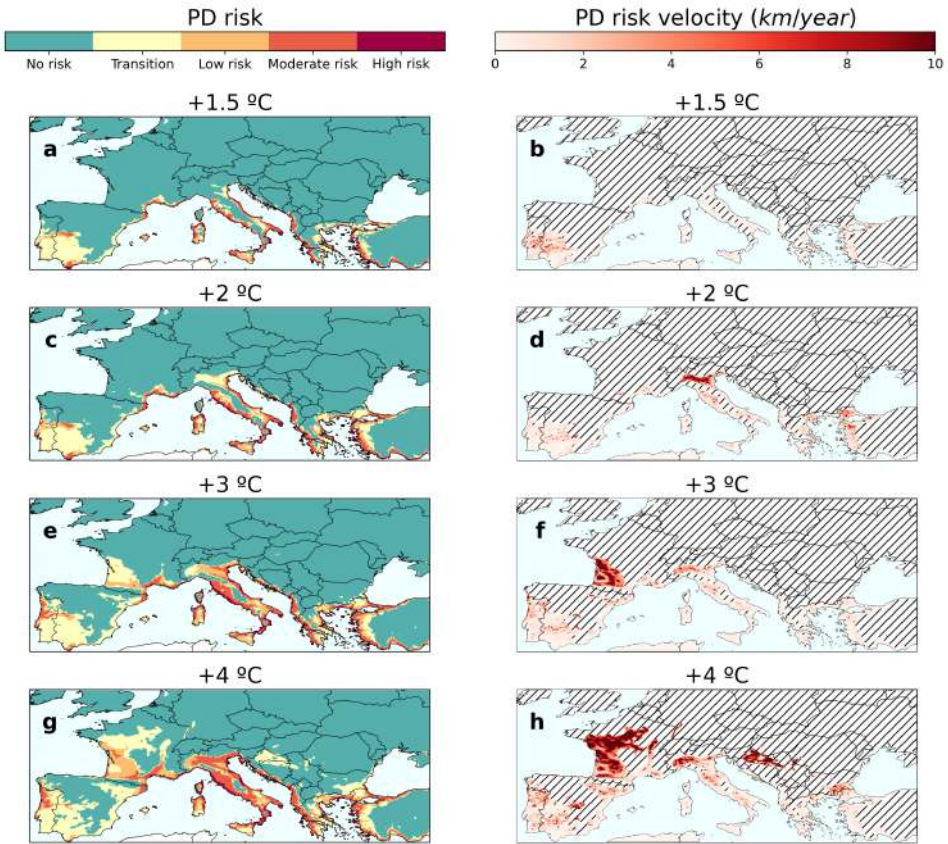
**Figure 8.2: Changes in  $Xf_{PD}$  and *P. spumarius* climatic suitability (i.e., probability of occurrence) under different climate projections compared to the current scenario (2003-2022).** Current climatic suitability for the pathogen (a) and the vector (b). In an increasing temperature scenario (+2°C, +3°C, and +4°C), the climatic suitability for the pathogen geographically expands in southern Europe and moves northwards (c, e, and g), while the climatic suitability for the vector decreases (d, f, and h). The suitability values for each scenario correspond to a 20-year average.

risk mainly due to the expected summer warming (Fig. C.15). Conversely, areas of Central Europe such as Hungary and Serbia already experience suitable conditions for pathogen growth in a  $+1.5^{\circ}\text{C}$  scenario ( $\mathcal{F}(MGDD) > 0.6$ ); however, cold winters tend to eliminate any potential summer infection [ $\mathcal{G}(CDD) < 0.3$ ] (Fig. C.15). Climate change would further increase the growth of the pathogen and reduce the winter curing effect in Central Europe, ultimately exposing the region to  $Xf_{PD}$  (Fig. C.15).

### 8.3.2 Pierce's disease risk projections under climate change

The limited intersection between the climatic suitability ranges for the pathogen and the vector (Fig. 8.2 a and b) suggests a marginal risk of PD epidemics in Europe. Since disease transmission requires a vector, the climate-suitability maps for *P. spumarius* indicate a lower risk and potential economic impact of Xf induced diseases on any host in southern Europe, particularly Spain, than previously predicted [235]. Realistic risk maps require a defined epidemiological framework to account for inter-annual climate variation and transmission in disease dynamics, in addition to accounting for changes in the distribution of climatic conditions favorable to the pathogen and vector (i.e., climatic suitability). Our epidemic risk model focuses on delimiting the disease dynamics by simulating an epidemic process in which the emergence of newly exposed hosts is influenced by the climatic suitability of the vector, while the transition to the infectious state is driven by the climatic suitability for  $Xf_{PD}$  chronic infections. The effective growth rate of the infected host population over the simulated period is used to derive a risk index  $r$ , bounded between  $-1$  and  $1$ . Within this modeling framework, different risk categories naturally emerge: no risk ( $r < -0.1$ ), transition zone ( $-0.1 \leq r < 0.1$ ), low risk ( $0.1 \leq r < 0.33$ ), moderate risk ( $0.33 \leq r < 0.66$ ) and high risk ( $r \geq 0.66$ ). For further details, see the Methods section and the original paper [280].

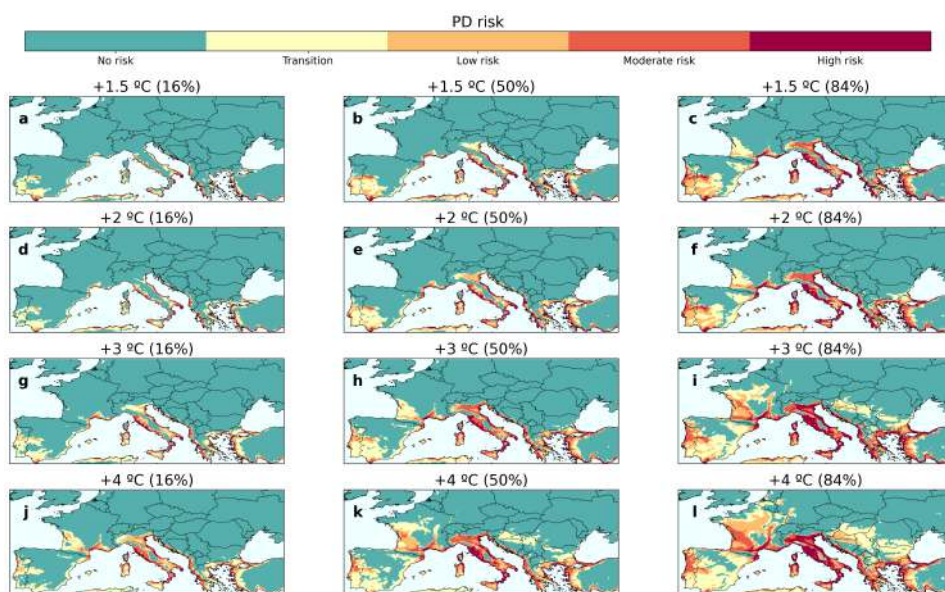
Global warming ( $+1.5^{\circ}\text{C}$ ,  $+2^{\circ}\text{C}$ ,  $+3^{\circ}\text{C}$ , and  $+4^{\circ}\text{C}$ ) is expected to increase the risk of PD epidemics in southern Europe, with France, Italy, and Portugal being particularly affected (Fig. 8.3 and Fig. C.16). This general trend affects each of the risk categories under the different climate change scenarios (Fig. C.17). Furthermore, we observed that a global temperature increase above  $+3^{\circ}\text{C}$  represents a tipping point for the possible spread of PD beyond the Mediterranean (Fig. 8.3 and Fig. C.16). To quantify the potential spread of PD, we calculated risk velocity, an index that allows us to identify areas where risk is changing or spreading rapidly (see Methods). We found a consistent and notable increase in the mean risk velocity within most of the identified risk zones (Fig. 8.3 and Fig. C.17), increasing from almost  $1 \text{ km y}^{-1}$  to  $5 \text{ km y}^{-1}$  as the temperature rises from a  $+1.5^{\circ}\text{C}$  to a  $+4^{\circ}\text{C}$  scenario (Table C.10). This



**Figure 8.3: PD risk maps and associated risk velocities under different climate projections.** (a, b) +1.5°C climate projection. (c, d) +2°C climate projection. (e, f) +3°C climate projection. (g, h) +4°C climate projection. Risk velocities have been calculated only in risk zones,  $r > 0$ , in each scenario. Hatched lines in panels (b, d, f, and h) indicate no risk zones where risk velocities have not been calculated.

acceleration is evident when we compare that in the +1.5°C scenario, approximately 6% of the grid cells have risk velocities greater than  $5 \text{ km y}^{-1}$ , while this value increases to 50% in the +4°C scenario (Table C.10). Furthermore, our estimates of PD risk velocity are broadly consistent with estimates of the velocity of temperature change [21], indicating that shifts in PD risk in our model adequately track climate change.

Fig. 8.4 shows the uncertainty in the projections of the PD risk map, comparing the 16th and 84th percentiles ( $1\sigma$ ) from the set of 40 regional climate models to the median risk map. The spatial distribution of PD risk is robust

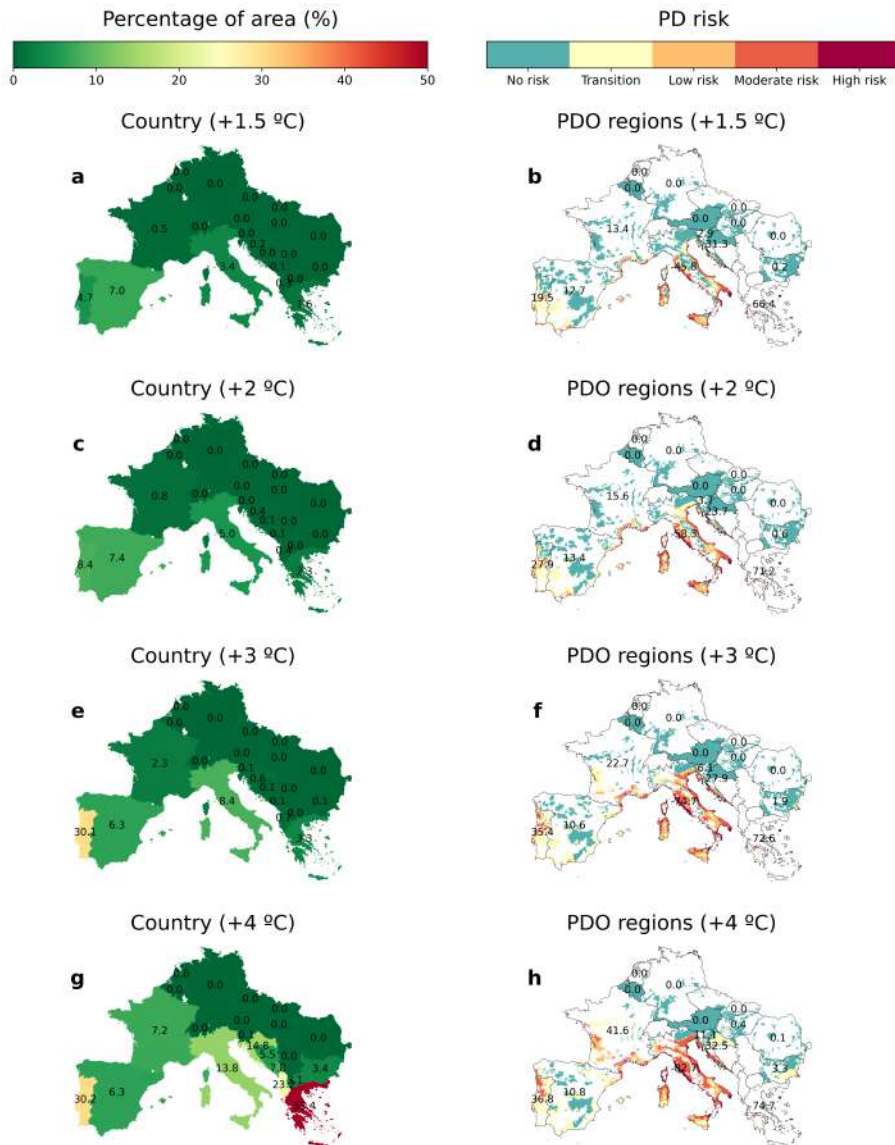


**Figure 8.4: Uncertainty in PD risk projections for climate change warming levels.** The maps show the result of the median risk values of the set of 40 regional climate model projections for each level of temperature increase (b, e, h, k) and the uncertainty of the projections considering the  $1\sigma$  deviations from the median, this is, the 16% (a, d, g, j) and 84% (c, f, i, l) percentiles. For each warming level (each row), the projected risk of PD is comprised between the results shown in the 16% percentile (first column) and those shown in the 84% percentile (last column) with 68% probability ( $1\sigma$ ).

across models, while the uncertainty in the level of warming is bounded by a  $\pm 1^\circ\text{C}$  increase (Fig. 8.4), e.g., the median spatial distribution of risk obtained for a  $+3^\circ\text{C}$  warming level is expected to occur between the  $+2^\circ\text{C}$  and the  $+4^\circ\text{C}$  scenarios under a  $1\sigma$  confidence level. This means that, depending on the specific model, a given spatial distribution of risk (e.g., Fig. 8.4 h) may manifest under a  $+4^\circ\text{C}$  warming scenario in more conservative projections or a  $+2^\circ\text{C}$  warming scenario in others, while most models project it for a  $+3^\circ\text{C}$  warming level. This indicates a high degree of confidence in the location and projected severity of future outbreaks, but greater uncertainty in their timing.

In order to improve the accuracy of the risk maps and the relative impact of PD, we carried out a comprehensive analysis at multiple scales, from the national level to the regions with Protected Designation of Origin (PDO) and finally taking into account the distribution of vineyards. This approach allows us to disaggregate the results at different administrative levels to facilitate the

design of risk management and the implementation of appropriate phytosanitary measures.



**Figure 8.5: Multi-scale spatial analysis of PD future risk in Europe.** (a, c, e, g) Percentage of country areas at risk ( $r > 0.1$ ) for each climate projection. (b, d, f, h) PD risk zones in Protected Designation of Origin (PDO) wine regions for each climate projection. PDO data was obtained from [387]. The corresponding interactive analysis at the vineyard level can be found in [338]

Overall, our model simulations show a consistent increasing trend in the risk of PD in Europe under all climate change scenarios. The percentage of land area at risk in Europe increases from 0.32% in the +1.5°C scenario to 1.87% in the +4°C, while the number of regions with PDO at risk increases from 18.17% to 47.32%. The vineyard area increases from 18.67% to 40.35% (Table C.11). At country level, Portugal and Greece face the highest overall risk, escalating from 12% and 2% of their area in the +1.5°C scenario to a striking 47% and 63%, respectively, in a +4°C scenario. In contrast, countries such as France and Italy experience a smaller but still significant increase in risk areas, never exceeding the 20% threshold, while Spain, the third-largest wine producer, shows a decreasing trend in risk areas above the +2°C scenario (Fig. 8.5 and Table C.12). Such contrasting patterns in PD risk between countries only emerge when using our modeling framework.

A different picture emerges when looking at the spatial distribution of PDO regions and vineyards. For example, PD risk within French and Italian PDO regions increases significantly from 13.4% and 45.8% in a +1.5°C scenario to 41.6% and 82.7% in a +4°C scenario (Fig. 8.5 and Table C.12), while the percentage of vineyard area at risk rises from 24.21% and 57.49% in a +1.5°C scenario to an astonishing 80% in a +4°C scenario. Important European PDOs would be at risk from a warming of +2°C, such as parts of the southern Rhône Valley (Châteaufort du Pape), Provence and Languedoc in France, Penedés in Spain, Bairrada in Portugal, and Chianti and Brunello di Montalcino, among others, in Italy (see Supplementary Information). A detailed interactive analysis of the impact of PD in European PDO regions and vineyards is available on our website [338].

## 8.4 Discussion

Previous research has attempted to assess the potential geographic distribution of  $X_{fPD}$  subspecies, the insect vector *P. spumarius*, and PD using species distribution models (SDMs) under future climates. While these climate-suitability-based predictions provide insights into the ecological niche of key disease players, bioclimatic correlative models neglect disease dynamics, a key factor in avoiding disease overestimation [280]. Unlike previous attempts, our approach integrates the compound effect of climate change in the pathosystem using a mechanistic epidemiological model to overcome these limitations. Unlike dimensionless climatic suitability indexes or disease probabilities used in SDMs, the risk index,  $r$ , in our model provides information on the expected growth rate in the event of an outbreak. Furthermore, the risk index is not fixed but varies annually depending on the weather conditions of the previous years. Inter-annual climate variability thus has an impact on disease dynamics, especially in areas where the

risk index is lower. Another feature of our risk predictions is their lack of ambiguity, which should not be confused with certainty. Risk estimates are based on  $R_0$ , which depends on the insect vector population [289], among other factors. Areas where  $r < -0.1$  permanently cannot theoretically support an outbreak, and the population of infected plants will decline over time. For example, our model clearly indicates that there is no risk of PD in the UK. This is not an arbitrary threshold; it is given by the epidemiological model. It is therefore very likely that the absence of PD in continental Europe is a consequence of low risk indices and that it has only become established in certain coastal areas since the late 1990s. On the contrary, the risk index in the Mediterranean islands has remained moderately high with little variation over the last 40 years [280].

Despite Pierce's disease having been found in Spain, Portugal, France and Germany [329], it has only affected vineyards on the island of Mallorca [268]. Perhaps for this reason, little attention has been paid to the risk of it reaching continental vineyards. Our risk model indicates why this possibility was very low until the mid-90s, and what the conditions were for it to occur on the Mediterranean islands [280]. In this work, we clearly show that with increasing temperatures, PD will become a serious threat to important wine-growing areas in southern Europe that were not previously at risk. A key finding of our study is the identification of a tipping point for the risk of PD establishment at a global mean temperature increase of  $+3^\circ\text{C}$ . Beyond this threshold, the risk of PD spreading north of the Mediterranean region becomes remarkably higher, while the risk of PD epidemics in Portugal, Italy, and France (Fig. 8.3) undergoes a significant quantitative leap. This suggests that as global temperatures continue to rise, the range of PD may expand into new territories. Indeed, the projected increase in risk velocities under higher warming scenarios further emphasizes the potential for rapid spread of PD into previously unaffected regions (Fig. 8.3).

Pest risk map projections are subject to uncertainties inherent in the variability of climate model predictions [388]. While previous studies on pathogen and vector distributions have been based on a limited number of climate models, our risk maps are based on the most modern set of regional climate projections produced by the EURO-CORDEX initiative, reflecting the state-of-the-art knowledge (Table 8.1). This allows us to adequately estimate the uncertainty of the resulting PD risk map projections for each temperature rise scenario. This confirms that although the spatial distribution of the risk of establishment is robust, there is an uncertainty of  $\pm 1^\circ\text{C}$  in the level of warming (Fig. 8.4). The models are therefore fairly good at pointing where the increased risk will occur, but it is more difficult to know when it will be reached.

Overall, our results highlight the contrasting effect of climate change on PD risk distribution in Europe, revealing it as a multifactor and multiscale pro-

cess (Table C.12). Climate change has an opposite effect on each component of the pathosystem, enhancing areas of potential chronic PD infections while diminishing the suitable geographic range for the vector. At the same time, the characteristic spatial scale at which risk is assessed strongly influences conclusions. At the country level, there are significant variations in the extent of accumulated risk between different projections. However, when analyzed at a finer scale, such as at the level of PDO regions or vineyards, the results change completely. Countries that previously had marginal areas at risk now show a higher percentage of PDO regions and vineyard at risk. These results underlie the urgency of tailored mitigation and adaptation strategies to protect vineyards and PDOs, considering their specific spatial distribution and risk index, as well as the potential impacts of climate change.

Our results are influenced by the intrinsic uncertainty associated with the correlative models used to determine the spatial distribution of the vector, the epidemiological parameters, and the uncertainties in the climatic projections. Although the spatial resolution of our climate projections is considered to be high, it may not capture the complex microclimate structure found in certain European wine-producing regions. Therefore, risk assessment results could differ locally with higher-resolution data. In addition, we have not considered the possible influence of climate change on latitudinal and altitudinal shifts in the distribution of European vineyards [389, 390], as this would only affect the calculation of the percentage of vineyard surface at risk but not the actual spatial distribution of risk. In any case, the risk estimates for the PDO regions include areas much larger than the areas of planted vines, which allows some margin in the adaptation and migration of the vineyards to different microclimatic conditions. In addition, the PDO and vineyard databases used in this study also have their own limitations. Future studies incorporating more refined modeling techniques, specific regional grapevine varieties, crop management and improved data resolution would enable a more nuanced understanding of PD risk and its potential impact at the local scale.

It is noteworthy that the mathematical framework employed in this study could be applied to other *Xylella fastidiosa* diseases, such as Almond Leaf Scorch Disease or Olive Quick Decline Syndrome and, more generally, to other vector-borne plant diseases. However, this requires the availability of some specific data and conducting some experiments. First, data for the temperature-dependent growth rate of the pathogen is needed to build the function that computes the MGDD. Then, symptom development experiments need to be carried out to build the  $\mathcal{F}(MGDD)$  and  $\mathcal{G}(CDD)$  functions that relate symptom development with temperature. Finally, the spatial distribution of the agent responsible for disease transmission is desired. Of course, using presence/absence data one can

use SDM to obtain this spatial distribution.

Climate change is currently one of the biggest challenges for EU agricultural policy [391]. Quantitative regional predictions of climate change on emerging diseases, such as this one, provide a valuable and unambiguous tool for decision-making. In our approach to the problem, risk indexes not only include information on where or where not PD may become established, but also reflect the exponential growth rate of potential epidemics, which are directly related to their potential economic impact. In addition, risk indices and velocities provide a dynamic framework for assessing the feasibility of eradication efforts when  $X_{fPD}$  is detected in a new area, providing critical information for strategic crop protection. Our study evidences the need to selectively allocate more resources to surveillance and research on PD in southern European countries, considering the associated uncertainties. This strategic allocation of resources based on risk assessment can help to prioritize proactive measures and effectively manage the potential impact of PD in different European countries.

Our research highlights the complex dynamics of PD and its relationship with climate change. By adopting an interdisciplinary approach that integrates climate projections, epidemiological modeling, and spatial analysis, we provide valuable insights into the potential establishment and spread of PD in European wine-growing regions from the country to the vineyard levels. Our study demonstrates that an accurate assessment of the risk of PD establishment requires a nuanced understanding of the vector-plant-pathogen-climate system and the explicit consideration of the vineyard spatial setting. These findings can inform decision-making processes and support the development of effective strategies to mitigate the risks posed by PD and safeguard the future of viticulture in the face of a changing climate.





## 9. Pierce's disease risk with high-resolution climate data

### Published as:

À. Giménez-Romero, E. Moralejo, and M. A. Matias, "High-resolution climate data reveals increased risk of Pierce's Disease for grapevines worldwide", [bioRxiv \(2024\)](#)

## 9.1 Introduction

Climate plays a pivotal role in shaping the distribution and dynamics of agricultural pests and pathogens [333, 361, 362, 392, 393], with implications for global food security [394, 395]. As our climate undergoes unprecedented changes due to anthropogenic activities, agriculture faces multifaceted threats ranging from alterations in temperature and precipitation patterns to increased frequency of extreme weather events [396]. Such shifts create novel environments that may favor the proliferation of certain pests or pathogens while posing challenges to the survival of others [333, 363]. The consequences of these changes extend beyond immediate agricultural landscapes, reverberating through global food systems and posing significant challenges to the sustainability and resilience of food production [397].

Understanding the intricate relationships between climatic conditions, the pathosystem components, and the subsequent epidemiological dynamics is essential for developing effective strategies to mitigate and manage emerging agricultural challenges, especially in the face of changing environmental conditions. However, modeling disease epidemics is a complex task, as they are emergent phenomena resulting from non-linear interactions between disease components that also exhibit non-linear responses to changes in environmental variables [259, 335, 370]. Thus, while climate primarily determines the potential geographic range of each organism in the pathosystem, the development of epidemic outbreaks depends on favorable host-pathogen-vector-climate interactions that drive transmission chains.

It has long been recognized that ecological phenomena typically depend on the scale of description, particularly with regard to the effects of climate [398]. Climatic databases with finer spatial resolution are continuously being developed with the goal of allowing more accurate predictions [399]. Some recent studies have shown that the local climate experienced by individuals might deviate substantially from regional averages, with implications for the population dynamics of a forest herb [400]. Likewise, the choice of climate data affects the predictions of species distribution models (SDMs) [401]. In particular, the spatial resolution of the data can influence the predictions of invasion risk for some species [402]. It is therefore clear that the resolution of climate data will have a significant impact on predicting the risk of plant diseases and pests.

Among emerging pathogens, *Xylella fastidiosa* (Xf) is considered one of the most dangerous phytopathogenic bacteria worldwide [261, 403]. It is naturally transmitted by xylem sap-feeding insects, such as sharpshooters and spittlebugs, and exhibits a broad host range that encompasses economically important crops such as grapevines, citrus, almond trees, and olive trees [128, 403]. The con-

sequences of Xf diseases are devastating: about 200 million citrus trees are infected annually in Brazil [404], there are losses over \$100 million annually in the grape industry in California, [234] and approximately 21 million olive trees have been killed by the bacterium in the Apulia region of Italy [405]. Assuming massive spread throughout Europe, Xf has been projected to potentially contribute up to €5.2 billion of annual losses in the olive sector alone [235]. Overall, Xf diseases pose a major threat to agrosystems worldwide, highlighting the need for precise and predictive models to guide effective management practices.

Previous research has provided insights into the potential geographic range of Xf subspecies through SDMs [372, 373]. These models, however, have led to overestimates of risk by failing to account for the distribution and abundance of potential vectors necessary for disease transmission [371]. A quite different approach to mapping the risk of Pierce's disease of grapevines has been developed based on climate-driven epidemiological models with the option to integrate vector's distribution information and the specificity of the Xf subsp. *fastidiosa* strain responsible for the disease, (hereafter Xf<sub>PD</sub>) [280]. This model correctly identifies areas in the United States with recurrent PD outbreaks and forecasts increasing epidemic risk in Mediterranean islands and coastlines with ongoing climate change.

Although risk maps based on hourly temperature data from the ERA5 have allowed fine adjustments in the calibration of the thermal response to Xf infection, these achievements have entailed losses in spatial resolution (0.1° spatial resolution) [280, 345]. Such limitation is particularly significant when dealing with vector-borne plant diseases like PD, where the interactions between the pathogen, vector, and host plants exhibit non-linear responses to climatic conditions. Subtle variations in temperature, humidity, or precipitation at the local scale thus can have profound effects on the reproduction and life cycles of the organisms involved and, hence, on the dynamics of disease transmission.

Topographical heterogeneity is a recognized issue in invasion biology but has received little attention in crop science. Vineyards are increasingly located in valleys, ridges, hillsides, and riverbanks, usually with altitudinal and microclimatic gradients in short transects. They are therefore a remarkable example of a crop subject to scaling problems when studying ecological or epidemiological processes at regional and global scales. In this work, we address this spatial resolution limitation by modeling the risk of PD using high-resolution climate data from the CHELSA dataset [406]. The study period was deliberately chosen to include real data on temperature increases due to ongoing climate change. Our study shows a greater global risk of PD and a higher rate of risk increase, underscoring the urgency of reevaluating global strategies to prevent the spread of the pathogen with international trade in plant diseases.

## 9.2 Results

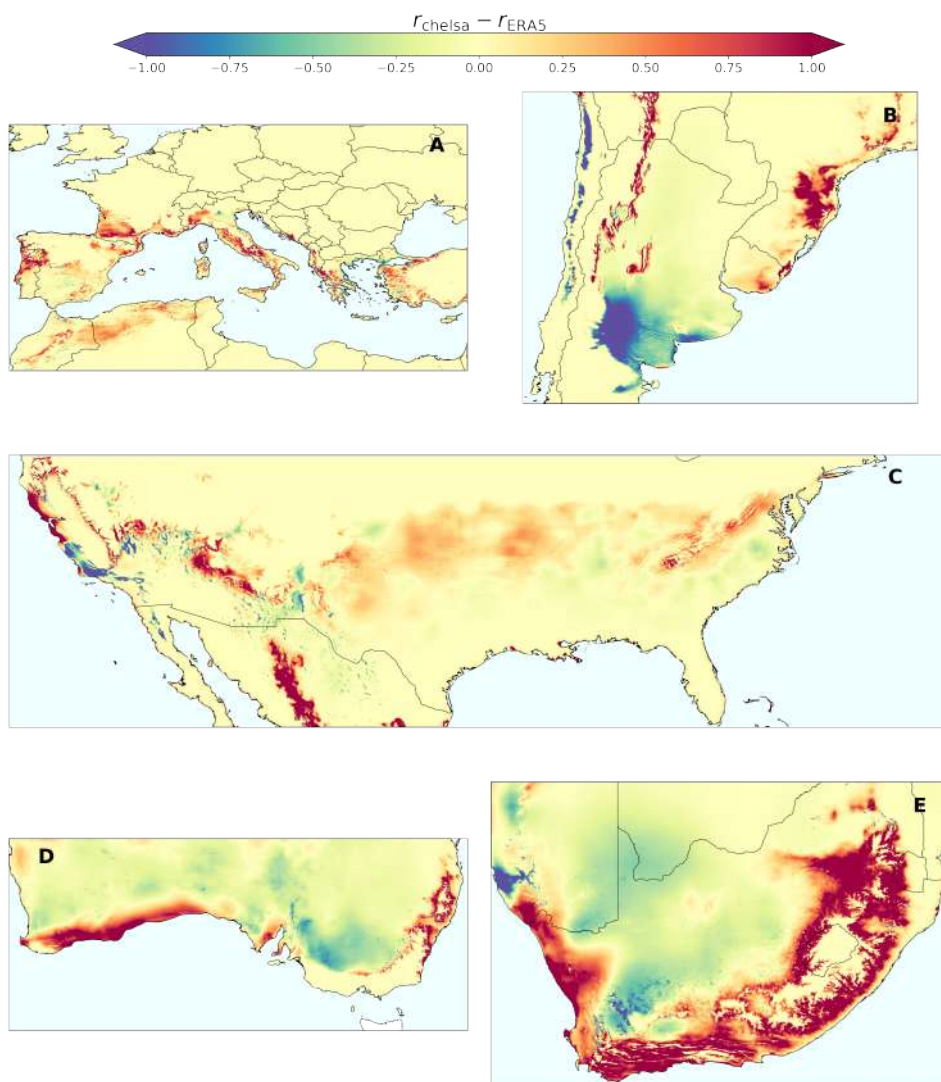
### 9.2.1 Global differences in PD risk between coarse and fine-grain climate data

We computed the risk of PD using the previously developed climate-driven epidemiological model in Chapter 7 [280] coupled with the CHELSA dataset [406], which features key climate variables (e.g., temperature and precipitation) at a high spatial resolution of 1 km and daily temporal resolution covering the period 1979-2016. The resulting spatial and temporal patterns of disease risk in the main wine-growing regions were compared with the previous risk projections derived from the ERA5 dataset [376] in Chapter 7, characterized by an intermediate spatial resolution of 10 km and hourly temporal resolution [280]. Risk projections in Europe use the climatic suitability,  $s$ , of the main European vector, *P. spumarius* (see Methods), while for the rest of the world it is assumed that there are no risk-limiting effects due to the vector ( $s = 1$ ), but only due to climatic conditions.

**Table 9.1: Changes in Pierce's disease risk zones in different viticulture regions.** The table illustrates transitions between risk and no-risk categories, as well as transitions among risk categories, highlighting the dynamic shifts in risk patterns across viticulture areas in Europe, the USA, South Africa, South America, and Australia. Risk increase refers to changes from low to moderate risk or from moderate to high risk. Likewise, risk decrease refers to changes from moderate to low risk or high to moderate risk.

	Europe	USA	South Africa	South America	Australia
<b>Risk to no-risk (km<sup>2</sup>)</b>	1.91e+04	3.93e+04	2.26e+04	2.49e+05	5.81e+04
<b>Transition to risk (km<sup>2</sup>)</b>	6.37e+04	1.37e+05	5.53e+04	3.79e+04	8.49e+04
<b>Risk decrease (km<sup>2</sup>)</b>	3.58e+04	1.46e+05	3.56e+05	1.76e+05	1.28e+06
<b>Risk increase (km<sup>2</sup>)</b>	6.28e+04	2.37e+05	1.05e+05	1.50e+05	1.55e+05
<b>No risk to risk (km<sup>2</sup>)</b>	2.04e+05	2.36e+05	2.77e+05	1.90e+05	2.15e+05
<b>Total changes (km<sup>2</sup>)</b>	3.85e+05	7.95e+05	8.15e+05	8.03e+05	1.79e+06

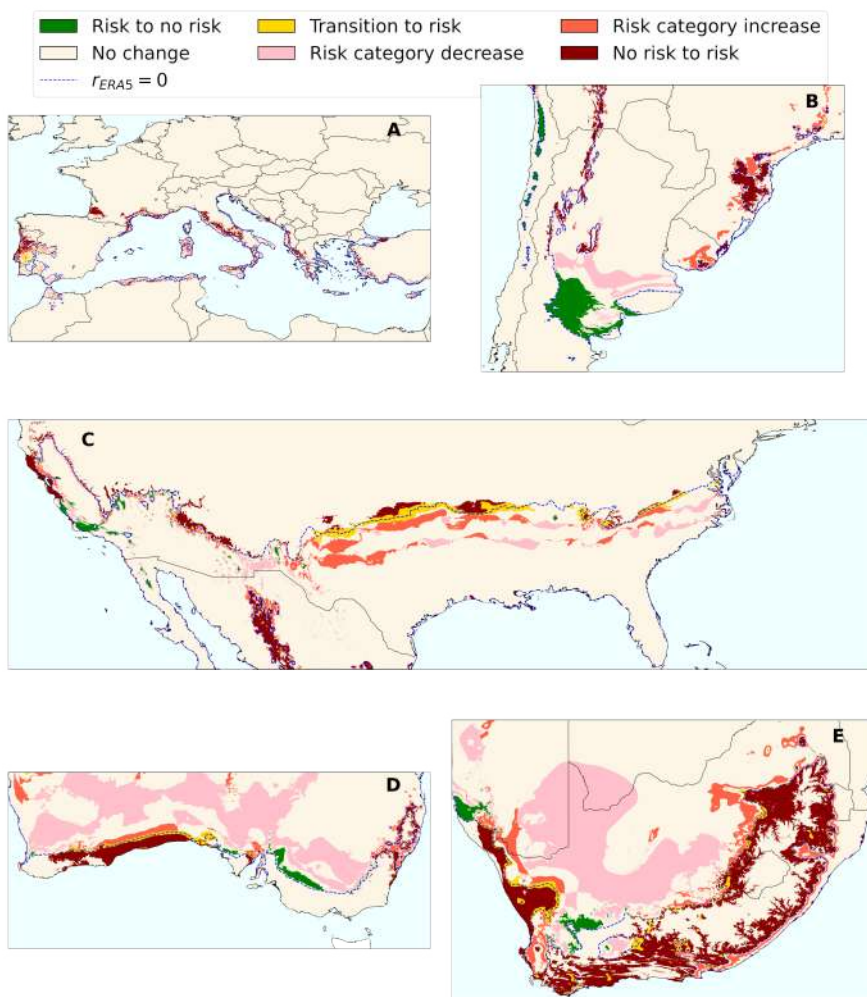
When contrasting model results derived from high- and medium-resolution data for the latest available time (2016), the disparity in risk projections extends beyond regional differences, showing a global increase in risk indices across wine-growing areas (Figs. C.22 and 9.1). Overall, these increases in the extension of PD risk areas ranged from 100,000 to 1 million km<sup>2</sup> across viticulture regions worldwide (Fig. 9.2). Transitions from no-risk to risk zones covered an area one order of magnitude larger than those in the opposite direction –from risk to no-risk (Fig. 9.2 and Table 9.1). In total, a surface of 4.6 million km<sup>2</sup> changed its risk category with the CHELSA database, representing about 16%



**Figure 9.1:** Difference in risk projections based on CHELSA (high-resolution; 1 km) and ERA5 (mid-resolution; 10 km) datasets in global viticulture areas. (A) Europe (B) South America (C) United States (D) Australia (E) South Africa.

of the land area studied. In contrast, the largest decreases in the risk indices occurred mainly in the Southern Hemisphere, although with few exceptions most of these decreases remained within the risk zones (Fig. 9.1), while similar land expansions were observed to increase their risk category (low to moderate or moderate to high) (Fig. 9.2 and Table 9.1). The largest changes in risk

indices occur in ecotones on both sides of the  $r = 0$  line, as is clearly seen in the south-eastern United States, in coastal areas (e.g., southern Australia and northern California) due to higher resolution that better distinguishes between land and coast, and finally in the river valleys and slopes of mountain systems (Fig. 9.2 and Table 9.1).



**Figure 9.2:** Changes in risk categories between CHELSA (high-resolution; 1 km) and ERA5 (mid-resolution; 10 km) projections in global viticulture areas. (A) Europe (B) South America (C) United States (D) Australia (E) South Africa. Risk category increase refers to changes from low to moderate risk or from moderate to high risk. Likewise, risk category decrease refers to changes from moderate to low risk or high to moderate risk.

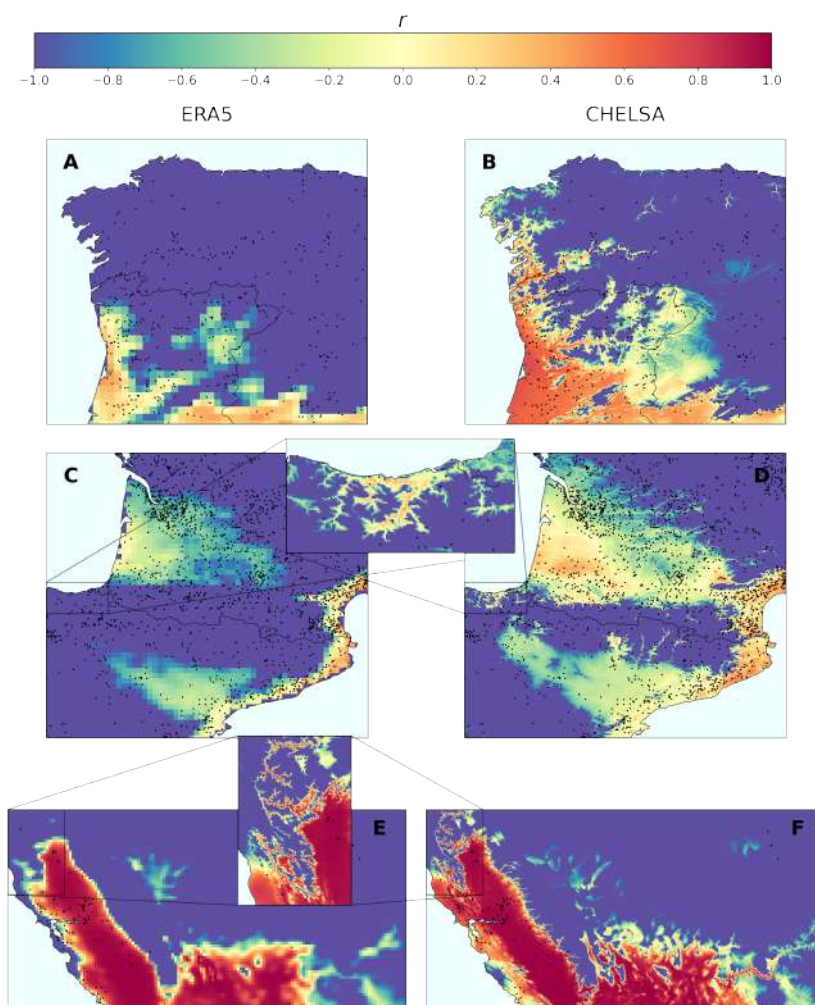
Next, we compared the temporal progression of the area at risk using both high and mid-resolution data over the entire available time span (1986-2016), considering that the risk for each year is computed based on the preceding seven years. We found a notable increase in the rates of expansion of the area at risk within viticulture zones worldwide, practically doubling previous estimates (Fig. C.21). These results point to an accelerated pace at which the risk of PD is growing, compatible with the predictions of different global warming scenarios [407].

### 9.2.2 Pierce's disease risk surges in previously unresolved microclimates

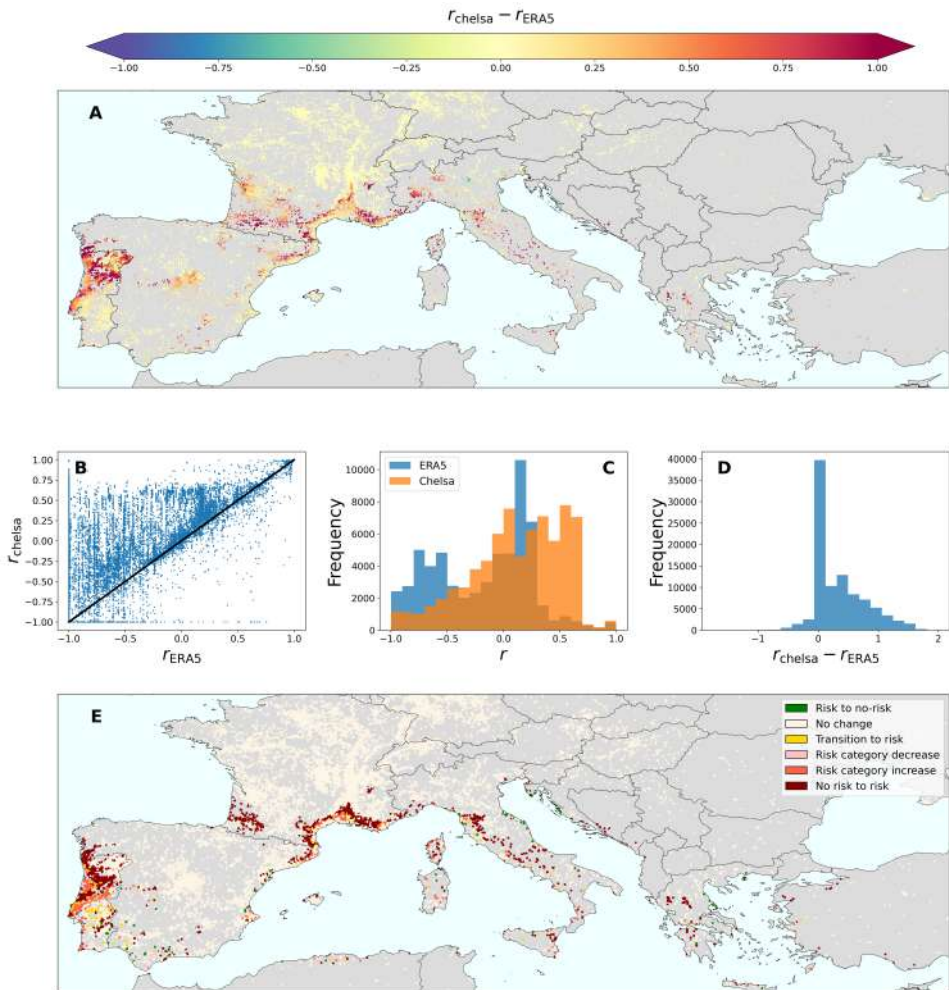
River valley vineyards are renowned for their high quality wines, such as Douro, Napa, Rhone, and many others. It is therefore important to understand the risk of PD with climate change at a more detailed level. In our analysis, we have identified rivers and valleys as specific relief areas where a greater increase in PD risk is observed when employing CHELSA's finer-scale climate data (Fig. 9.3). In some important wine-growing areas of Southern Europe, we observed an abrupt emergence of risk zones previously classified as no-risk when using lower resolution climate data (Fig. 9.3). Such pronounced differences in risk patterns are highlighted, for example, in the fairly steep valleys and hillsides along the Douro River in Portugal, where the specific microclimatic conditions were previously obscured by the coarser resolution of the ERA5 data. These findings are particularly significant for PD, as vineyards are often located in proximity to rivers or valleys and their surroundings, creating microclimates that attenuate cold winters (black dots in Fig. 9.3). A gradual increase in the climatic suitability for PD in some river basins may thus favor the spread of the pathogen from coastal to interior areas of the continents, allowing interconnection between areas that would otherwise remain isolated. Coastal areas close to cool water masses may also undergo an increase in risk when using higher resolution data, as exemplified in California (Fig. 9.3 E, F).

Finally, to obtain a comprehensive assessment of the impact of microclimatic conditions on the risk of PD establishment, we collated a dataset of over 100,000 *Vitis vinifera* presence locations worldwide from GBIF [80], with a predominant concentration of points from Europe (Fig. C.23). Each data point was assigned a risk index based on the ERA5 and CHELSA projections, respectively, using the nearest pixel from each database. This approach revealed an increase in the risk indices associated with the vine locations (Fig. 9.4 A-D), mostly showing shifts towards higher risk indices (Fig. 9.4 A, E) from no risk to risk or increases in risk category (low to moderate or moderate to high), while a negligible number of points decreased in risk category (Fig. 9.4 E). Such behavior was common to all key viticulture regions studied, although the extent of increases differed between

continents, with substantial expansion of vineyard areas at risk in Europe and South Africa (Table 9.2). Overall, our results emphasize the global relevance of microclimatic conditions in influencing the risk landscape for PD in viticulture areas (Table 9.2).



**Figure 9.3:** Effect of microclimatic conditions of rivers and valleys on Pierce's disease of grapevines. Comparison of the risk predicted using ERA5 mid-resolution dataset (A, C, and E) and CHELSA high-resolution dataset (B, D, and F). (A-B) north-western Iberian Peninsula. (C-D) Southern France and north-eastern Spain. (E-F) Western United States. Black dots represent grapevines (*Vitis vinifera*) presence data obtained from GBIF (see Methods).



**Figure 9.4:** Impact of high-resolution climate data on the risk of Pierce's disease for grapevines worldwide. (A) Difference in risk indices in Europe, which accounts for the 96% of the points in the dataset. (B) Comparison of the risk indices derived from CHelsa and ERA5 datasets. Points with perfect agreement would lie in the solid black diagonal curve. (C) Histogram of risk indices derived from ERA5 (blue) and CHelsa (orange). (D) Histogram of the differences in risk indices between CHelsa and ERA5 datasets. (E) Changes in risk categories when using high-resolution climate data (CHelsa) with respect to mid-resolution data (ERA5). Risk category increase refers to changes from low to moderate risk or from moderate to high risk. Likewise, risk category decrease refers to changes from moderate to low risk or high to moderate risk.

**Table 9.2: Comparison of PD risk at known grapevine locations.** Comparison of grapevine presence locations at risk in key viticulture regions using CHELSA and ERA5 datasets

	N° points	risk CHELSA (%)	risk ERA5 (%)
Europe	96102	41.2	21.8
USA	792	69.8	66.3
South Africa	36	47.2	5.6
South America	112	77.7	74.1
Australia	186	51.6	45.7

### 9.3 Discussion

Our study sheds light on the relevance of the spatial scale of observation in the intricate interplay between microclimatic conditions and the risk of PD for grapevines on a global scale. The use of high-resolution climate data reveals previously unrecognized local areas with microclimates conducive to the establishment of PD worldwide. Contrary to the simplistic assumption that higher resolution data might yield only marginal distinctions at regional levels, our study demonstrates that slight variations in climate data at local scales can lead to a global surge in disease risk. These increases not only affect the spatial distribution of risk but also its temporal dimension, as suggested by the rate of increase in the surface area at risk. In the case of PD, we show that this rate nearly doubles when high-resolution climate data is considered compared to previous estimates obtained with mid-resolution data. Thus, our findings indicate a critical need for the use of local or high-resolution climate data in the assessment of disease risk, especially in areas characterized by diverse topography and even when only attempting to global estimates.

Such observed differences arise from the non-linear nature of disease dynamics and the response of the pathosystem components to environmental shifts [335, 363]. Therefore, models dependent on broader climate data may not capture the complexities of microclimates, resulting in an underestimation of disease risk. While this is not inherently negative, recognizing these limitations helps to assume such risk estimates as a conservative lower bound until proven otherwise. Acknowledging these constraints is crucial for refining our understanding of disease dynamics and ensuring that our risk assessments are sufficiently cautious in the absence of more reliable data. Likewise, data coarsening procedures should be avoided, if possible, when modeling climate-driven disease dynamics, even in spite of computational efficiency. This recommendation applies not only to disease risk predictions but to all those in which non-linear functions depending on climate variables are present, such as species distribution models or phenological models [408].

Despite the valuable insights gained, our analysis heavily relies on the qual-

ity and resolution of the climate data from the CHELSA dataset [409]. While this dataset offers information at a high spatial resolution, the temporal dimension is limited to a daily frequency, which forces us to apply an approximation to infer hourly data. Furthermore, the data may still be subject to biases or uncertainties inherent to the nature of the methodology employed in their construction. On the other hand, vector presence data is only accurately obtained for Europe, while a homogeneous presence is assumed in other viticulture areas. Additionally, the study primarily focuses on the effect of temperature conditions and the presence of potential vectors to determine the risk of Xf establishment, which may not encompass all possible contributing factors. Other variables, such as soil characteristics or vineyard management practices, were not explicitly considered in this analysis, leaving room for additional complexities in the disease dynamics. Furthermore, the study predominantly examines the risk at a global scale, and the applicability of the findings to specific local contexts may vary.

Future research should aim to address the aforementioned limitations and provide a more comprehensive understanding of the multiple interactions influencing PD development in viticulture regions. Other factors influencing disease spread, such as human behavior, land use changes, and ecological shifts, should also be explored, offering a more comprehensive and holistic view of the interplay between environmental conditions and disease vulnerability. The acceleration in the rate at which the risk of PD is growing calls for more research into control strategies to mitigate its impact on grapevine crops worldwide.

Although PD is currently restricted to North America and recently introduced in Taiwan [410], Mallorca (Balearic Islands, Spain) [268, 411], and Israel [136], since the mid-1990s climatic conditions are increasingly conducive to the establishment of PD in Southern Europe [280]. For example, with the increase in the resolution of climate data, our model predicts the recent detection of PD in Portugal [412], which was not anticipated using the ERA5 data [280]. In a short time, it is foreseeable that there will be more epidemic outbreaks in vineyards in Southern Europe if the entry of infested plants is not controlled. This does not necessarily have to be vines but can also include other plants such as almond trees or ornamental plants [265].

Overall, our study contributes to the growing body of knowledge on the impact of climate on agricultural pests and pathogens, emphasizing the importance of considering microclimatic conditions for a more deep understanding of disease dynamics. Future research should focus on developing comprehensive models that integrate high-resolution climate data, considering both the global and local factors that influence disease dynamics. This holistic approach will enable a more accurate prediction of disease risk, allowing for the development of

targeted management strategies and the enhancement of global food security.

## 9.4 Methods

### 9.4.1 Climate data

Climate data was downloaded from two datasets for our analysis: the ERA5 dataset [345, 376] and the CHELSA dataset [406, 409]. ERA5 offers mid-resolution climate data with a spatial resolution of 10 km and hourly temporal resolution, while CHELSA provides high-resolution data with a spatial resolution of 1 km and daily temporal resolution. Both datasets exhibit global coverage and encompass crucial climate variables, such as temperature and precipitation. For our simulations, we used the mean hourly temperature data from the ERA5 dataset and the maximum and minimum daily temperature data from the CHELSA dataset.

### 9.4.2 Vector climatic suitability

Vector climatic suitability data was obtained from [373], in which a Generalized Additive Model (GAM) is employed to calibrate the relationship of *P. spumarius* global occurrence with moisture index and maximum temperatures during summer index estimated from 1979 to 2013 using the CHELSA dataset.

### 9.4.3 Vineyard data

To assess the risk of Pierce's disease in locations where grapevines are present, we collected a comprehensive dataset of over 100,000 *Vitis vinifera* presence data records from the Global Biodiversity Information Facility (GBIF) [80, 383]. We note that while the dataset spans the globe, 96% of the points are located in Europe (Fig. C.23).

### 9.4.4 Model adaptation to daily temperature data

We used the model developed in Chapter 7 [280], in which MGDD and CDD metrics were defined using hourly temperature data (Eqs. (7.1) and (7.2)). However, the CHELSA dataset only provides daily granularity, so we use a basic sinusoidal extrapolation relating maximum and minimum daily temperature to hourly temperatures,

$$T_h = \frac{T_{max} + T_{min}}{2} + \frac{T_{max} - T_{min}}{2} \sin(w \cdot h), \quad (9.1)$$

with  $w = 2\pi/24$  and  $h$  ranging from 0 to 23. This approximation was validated in Chapter 8 [407].

# 10 Data-driven methods for ecological problems

<b>10</b>	<b>Reconstructing pH time-series with machine learning . . . .</b>	<b>211</b>
10.1	Introduction . . . . .	212
10.2	Results . . . . .	215
10.3	Discussion . . . . .	220
10.4	Methods . . . . .	223
<b>11</b>	<b>Mapping seagrass meadows from space . . . . .</b>	<b>229</b>
11.1	Introduction . . . . .	230
11.2	Results . . . . .	232
11.3	Discussion . . . . .	239
11.4	Methods . . . . .	242
<b>12</b>	<b>Universal spatial properties of coral reefs . . . . .</b>	<b>249</b>
12.1	Introduction . . . . .	250
12.2	Results . . . . .	251
12.3	Discussion . . . . .	256
12.4	Methods . . . . .	259



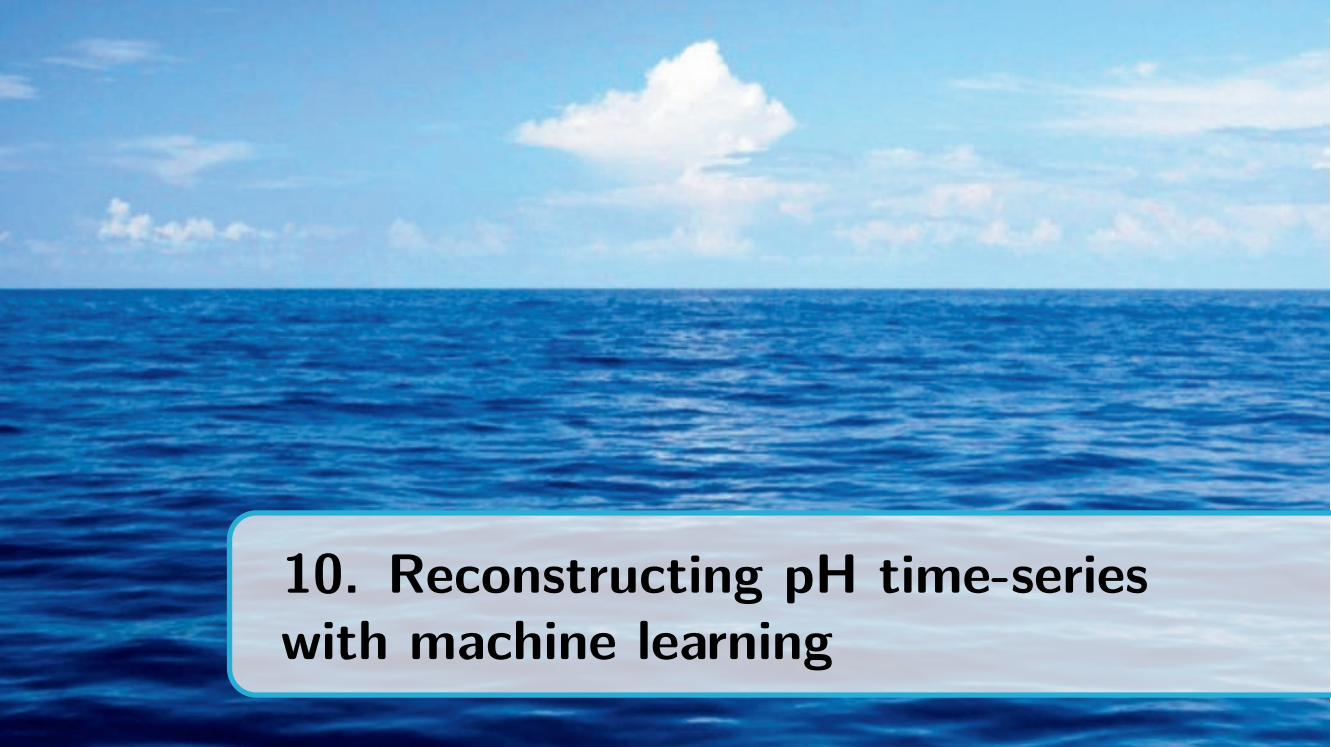
## Summary

In an era where ecological issues are becoming increasingly complex and globalized, the integration of data-driven methodologies into ecological research offers unprecedented opportunities for addressing these challenges. The last part of this thesis explores the application of such methods to study global ecological problems related to coastal marine ecosystems: the decline of coral reefs, the acidification of coastal waters, and the loss of seagrass meadows. Coral reefs are among the most biodiverse ecosystems on Earth, providing essential services to marine life and coastal communities. However, they are threatened by a combination of human activities and climate change, which have led to widespread coral bleaching events and the loss of coral cover. The acidification of coastal waters is another pressing issue posing significant risks to marine life, particularly species that rely on calcium carbonate for their skeletal structures like coral reefs. Seagrass meadows are also under threat from human activities, such as coastal development and pollution, which have led to the loss of seagrass habitats and the decline of associated biodiversity. The understanding of the spatio-temporal dynamics of these ecosystems together with the monitoring of their health and factors affecting their resilience is crucial for effective conservation and management strategies. However, this is often challenging due to the complexity of these ecosystems and the difficulty of collecting data at the necessary spatial and temporal scales. By leveraging large existing datasets, machine learning algorithms, and spatial analysis, here we provide novel insights into the spatial properties of coral reefs, the reconstruction of pH time series in coastal waters, and the mapping of seagrass meadows from satellite imagery. These case studies exemplify the potential of data-driven approaches to enhance our understanding of ecological dynamics, improve environmental monitoring, and inform conservation strategies.

### Objectives

- To develop a machine learning model able to reconstruct pH time series that present missing data due to sensor failures.
- To develop a machine learning model able to map seagrass meadows from multispectral satellite imagery.
- To investigate the spatial properties of coral reefs at a global scale.





## 10. Reconstructing pH time-series with machine learning

**Published as:**

S. Flecha, À. Giménez-Romero, J. Tintoré, F. F. Pérez, E. Alou-Font, M. A. Matías, and I. E. Hendriks, “pH trends and seasonal cycle in the coastal Balearic Sea reconstructed through machine learning”, [Scientific Reports](#) **12**, 12956 (2022)

## 10.1 Introduction

Atmospheric carbon dioxide ( $\text{CO}_2$ ) emissions are exponentially increasing since the industrial revolution, principally due to fossil fuel use, industry, and land-use change. Around 46% of this  $\text{CO}_2$  remains in the atmosphere, while the rest is captured by natural compartments: the terrestrial biosphere and the ocean [413]. At present, the oceans have absorbed around an estimated 26% of the total anthropogenic  $\text{CO}_2$  released from 2011 to 2020 [413]. Once  $\text{CO}_2$  dissolves in seawater, a sequence of chemical reactions occurs that derives in an increase of  $[\text{H}^+]$  ions, which results in a decrease in seawater pH. This process, a consequence of increasing atmospheric  $\text{CO}_2$ , is termed Ocean Acidification (OA) [414]. In addition to the pH decrease,  $[\text{H}^+]$  ions react with carbonate ions  $[\text{CO}_3^{2-}]$  to form  $[\text{HCO}_3^-]$ , leading to a reduction of the  $[\text{CO}_3^{2-}]$  ion levels [415].

Low carbonate levels affect the saturation state of calcium carbonate minerals, increasing difficulties in shell-forming for calcifying marine organisms (e.g., plankton, mollusks, echinoderms, and corals). Consequences of OA are an important threat to marine ecosystems visible in higher levels of the trophic chain, with complex and wide-ranging impacts on the physiology of different species and therefore on their metabolism [19, 416]. These metabolic effects will have numerous consequences for organisms; in particular, they can cause a decrease in growth, locomotion, reproductive capacity, and homeostasis if they are not capable of controlling the conditions for calcification [168]. Negative effects of this magnitude could cause an unexpected cascade effect impacting on the structure and functions of ecosystems and trophic networks [417] and cannot be easily generalized.

Also, ocean  $\text{CO}_2$  uptake and derived OA are not homogeneous at the global scale, with some areas more affected. For instance, the Mediterranean Basin is an area where effects are stronger compared to the global ocean [418]. The Mediterranean Sea, constituting only 0.82% of the surface and 0.32% of the volume of the global ocean, is cataloged as one of the most complex marine ecosystems, defined as a “miniature ocean” [419], inhabited by an extensive and diverse biota that represents between 4 and 18 % of the world’s total marine species [420] and serves as a model [419] to anticipate the responses of the global ocean to different types of pressures. It has also been defined as a climate change “hot spot” [418], with OA and its derived consequences characterized as one of the climatic threats with the greatest potential impact, followed by the temperature and UV radiation increase [421]. The temperature rise in this semi-enclosed sea is expected to be two to four-fold times higher than that in the global ocean [422, 423]. In addition, the sixth assessment report (AR6) of the IPCC places a high level of confidence on the increase in frequency of heatwaves and ongoing ocean acidification [424]. Recent studies have confirmed that there

is a trend of around 0.34 °C warming per decade in the Mediterranean Outflow Water (MOW) through the Strait of Gibraltar towards the Atlantic Ocean [425], associated with decreasing values of pH. Furthermore, in the Mediterranean Sea, due to its biogeochemical and hydrodynamic characteristics, such as the high alkalinity of its waters and the active thermohaline circulation [426], there is a larger absorption of atmospheric CO<sub>2</sub> and intense transport of this CO<sub>2</sub> from the oceanic surface to deep areas [427, 428], already observed in the MOW [429, 430], with estimated OA trends of -0.0044 pH units per year in the Strait of Gibraltar [429] and ranging from -0.0017 to -0.003 in the Mediterranean Basin [431, 432].

The Mediterranean Sea has an extensive coastline, which extends for 46 000 km and is shared by 21 countries [433]. Coastal zones, as transitional areas, are inherently complex systems due to the strong biogeochemical-physical coupling producing relevant biogeochemical exchanges. Interactions in coastal areas involve terrestrial inputs of nutrients and particulate matter from river runoff and groundwater discharges, oceanic forcing (waves, tides, and currents), and atmospheric exchange of aerosols and trace gases, all of which are influenced by the intense human activity on the coastline [434]. Hence, processes related to the carbon system in coastal areas are more dynamic and complex than in the open ocean [435], and the range of pH change, from -0.023 to 0.023 pH units per year [436], is larger than in the open ocean, where pH trends have been estimated from -0.0013 to -0.0026 pH units per [437]. In particular, anthropogenic CO<sub>2</sub> inputs appear to play a minor role compared to other sources of variability in coastal zones [438]. Therefore, it is difficult to foresee how the pH conditions in the coastal areas in the year 2100 will differ from the present, due to the lack of knowledge on precise current pH values in the different coastal ecosystems and their variability obtained from long time series. Carbonate chemistry and in particular pH fluctuations are characterized by a wide spatial heterogeneity and temporal variability (daily and seasonal oscillations) in coastal ecosystems [438–440]. The variability of pH is determined by a wide range of physical and biogeochemical processes, from mesoscale hydrological processes to small-scale metabolic processes [441].

The primary production in the western Mediterranean Sea is characterized by seasonal variability induced by the increase of the surface layer nutrients by the winter vertical mixing in the water column [442]. In addition, the presence of macrophytes [443] in the coastal areas of the northern Mediterranean Sea, mainly the endemic *Posidonia oceanica* with meadows that extend from the surface to 30–40 m depth, render these areas as highly productive habitats. In these ecosystems, variability tends to follow daily and seasonal cycles, since biological metabolism is responsible for variations in the concentrations of oxygen

(O<sub>2</sub>) and CO<sub>2</sub> [439, 444], so increasing pH values are expected for autotrophic ecosystems during daylight hours (production > respiration). Indeed, recent studies indicate that seagrass meadows can locally alleviate low pH conditions for extended periods of time with important implications for the conservation and management of coastal ecosystems [445].

Nevertheless, changes in pH can appear idiosyncratic and display a diversity of patterns depending on the coastal area under consideration, as many drivers of the carbon system can influence these variable ecosystems, including temperature variability, biological activity, and terrestrial and open ocean inputs [436]. Therefore, the properties of the carbon system have to be evaluated while taking into account the different interactions in every area. To the present day, there is still a lack of understanding of how coastal areas behave and how they contribute to the global carbon budget, also in part due to the intensive effort necessary to obtain representative time series of the carbon system data according to standard practices. The Global Ocean Acidification Observing Network (GOA-ON) defined that the accuracy included in the “weather goal” should be better than 0.02 and for the “climate goal” <0.003 pH units [446]. Instrumentation for autonomous pH measurements has improved in recent years, and production costs have come down, but they remain complex and relatively expensive. For climate change studies, commercial oceanographic instrumentation barely accomplishes the GOA-ON “climate goal” accuracy recommendation, with only spectrophotometric devices and Ion Sensitive Field-Effect Transistors (ISFETs) based pH probes reaching the standards. In this sense, the SAMI-pH sensor (Submersible Autonomous Moored Instrument, Sunburst Sensors, LCC), based on spectrophotometric techniques, has been denoted as an excellent pH sensor for OA studies [429]. However, the maintenance of oceanographic time series stations entails several operational and non-operational difficulties, involving financial costs, meteorological risks (i.e., bad conditions for navigation, instrumentation loss, etc.), deployment in areas with high transit, issues essentially related to the sensor itself (i.e., instrumental failure), and possible human errors. Therefore, the appearance of data gaps is common, implying the lack of pH data obtained using high-quality instrumentation for global carbon studies.

Currently, novel computational methods based on Machine Learning (ML) are allowing us to tackle these data absence difficulties. Machine Learning is a part of Artificial Intelligence that has attained a mature status in the last decade or so, particularly through the so-called Deep Learning (DL) model [104], with major advances in solving problems that have resisted the best attempts of the artificial intelligence community for many years [103]. In particular, some DL techniques are useful in time series forecasting [447] and also in the reconstruction of coupled time series [448], such as Recurrent Neural Network

(RNN) architectures like Long Short-Term Memory (LSTM) [106] or Gated Recurrent Unit (GRU).

Nowadays, there is an increasing number of studies that use DL to understand the processes involved in the carbon system variability, but mainly focused on the open ocean [449–453], while relatively few studies focused on coastal seas [454, 455] and none specifically in the Mediterranean coastal Sea, perhaps because of the complexity and heterogeneity of the basin and its continental shelves. Therefore, the main objective of this study is to obtain the trend for pH decrease in the coastal Balearic Sea by applying Machine Learning techniques. In addition, this study aims to provide a useful tool to fill gaps in pH time series and to reconstruct pH data when additional environmental variables are available.

## 10.2 Results

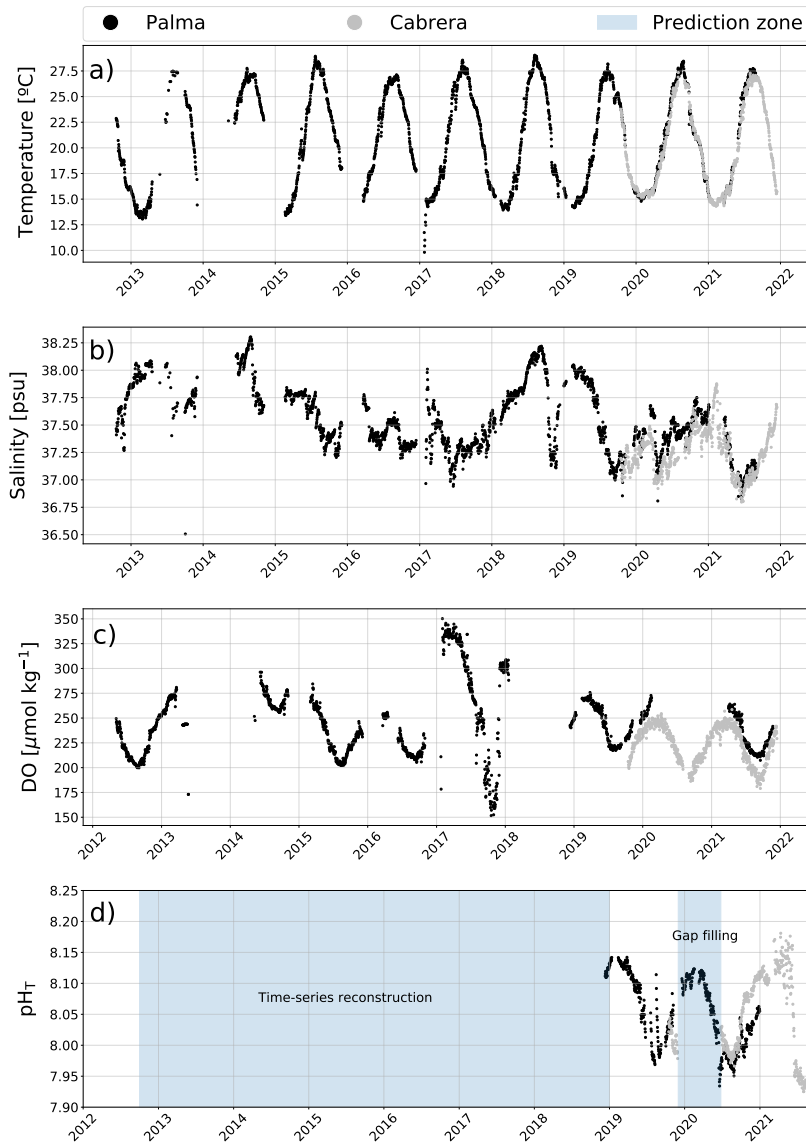
### 10.2.1 Time series data

The collection of pH values, in total scale ( $\text{pH}_T$ ), started in December 2018 in the Bay of Palma, recording data almost continuously until the end of 2021. In the Cabrera station,  $\text{pH}_T$  was obtained from November 2019 to December 2021, with a relevant data gap from December 2019 to June 2020 (Fig. 10.1 (d)) due to a sensor malfunction with a reparation prolonged for an extended period of time owing to the Covid-19 lockdown. Additional environmental parameters like temperature, salinity, and dissolved oxygen (DO) concentration are available from the Bay of Palma station since 2012, while only a limited time series of these variables (since 2019) is available for Cabrera (Fig. 10.1 (a-c)).

In both stations, temperature ranged from a minimum of 12.99 °C to maximum values of 29.07°C from 2012 to 2021, with no observed differences between the stations in Cabrera and the Bay of Palma (Fig. 10.1 (a)). The surface water temperatures are a clear representation of the typical Mediterranean climate seasonality with mild winters and warm to hot summers. Salinity did not show a repetitive seasonal pattern between years in either station. However, in Cabrera salinity is slightly lower than in the Bay of Palma. During the data acquisition period, the lowest salinity value of 36.83 was found in Cabrera and the highest of 38.30 in the Bay of Palma (Fig. 10.1 (b)).

The surface water of the coastal sites in the Balearic Sea in the Palma Bay and the Cabrera stations was highly saturated with oxygen during all the seasons, with DO concentrations up to 348.94  $\mu\text{mol kg}^{-1}$  during winter and of 169.66  $\mu\text{mol kg}^{-1}$  during the summer and early autumn (Fig. 10.1 (c)).  $\text{pH}_T$  values obtained starting in December 2018 to December 2021 increased during winter, reaching up to 8.18 pH units at *in situ* temperature and decreasing

to 7.91 pH units in summer, with the highest variability and maximum and minimum values measured in Cabrera (Fig. 10.1 (d)).



**Figure 10.1:** Daily averaged time series data from the Bay of Palma (black dots) and Cabrera stations (grey dots): a) Temperature ( $^{\circ}\text{C}$ ), b) Salinity (psu), c) Dissolved oxygen (DO) values ( $\mu\text{mol kg}^{-1}$ ) and d)  $\text{pH}_T$  in pH units. The pH time series of the Bay of Palma will be reconstructed in the period 2012-2021 while only gaps will be filled in Cabrera, as marked in blue in the figure.

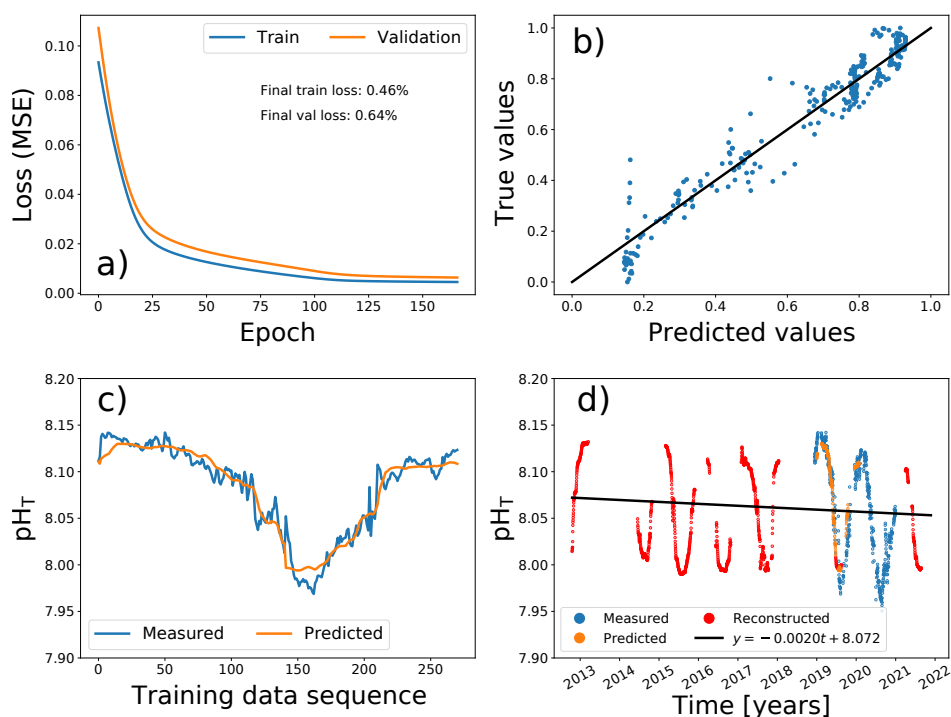
Considering that the sampling period of the additional (temperature and salinity) variables and calculated parameters (Total Alkalinity; TA) was larger in the Bay of Palma compared to Cabrera, we evaluated the linear tendencies with time for the Bay of Palma variables. The sea surface temperature in the Bay of Palma increased with a rate of  $0.05 \pm 0.03^\circ\text{C}$  per year ( $R^2=0.0008$ ,  $p\text{-value}=0.09$ ) from 2012 to 2021, whereas the salinity decreased significantly with  $-0.059 \pm 0.002$  psu per year ( $R^2=0.2456$ ,  $p\text{-value}<0.001$ ). The annual trend for TA, clearly related to the decrease in surface salinity, showed a relevant decrease of  $-4.0 \pm 0.4 \mu\text{mol kg}^{-1}$  ( $R^2=0.0379$ ,  $p\text{-value}<0.001$ ), supported by the discrete water samples for TA obtained during the period from 2019 to 2021 (Fig. D.3).

### 10.2.2 Reconstruction pH time series with Deep Learning

The amount of available  $\text{pH}_T$  data from both Palma Bay and Cabrera stations is comparable and relatively short (mostly in Cabrera), but the length of the additional ambient data (temperature, salinity, and DO) differs enormously between stations. Thus, there is a need to approach the time series prediction problem for both sites with different objectives. Common to both sites, a DL model with a RNN architecture will be developed to predict the  $\text{pH}_T$  time series from the accompanying ambient data (temperature, salinity, and dissolved oxygen), which are expected to be correlated with  $\text{pH}_T$  [449, 454]. To avoid the effect of site-specific correlations between ambient data and  $\text{pH}_T$  time series, the model will be trained independently with the dataset of each location. In this way, a proper model calibration is ensured, and the prediction power of the model is enhanced. In the Bay of Palma, the model will be used to reconstruct the  $\text{pH}_T$  time series from 2012, exclusively from the points for which the full set of ambient time series data are available (Fig. 10.1 (d)). This is not possible in Cabrera due to the fact that no temperature, salinity, or DO concentration is available before 2019. Fortunately, these time series do not have the same gaps that the  $\text{pH}_T$  time series exhibits. Thus, we will use the model to fill the gaps in the  $\text{pH}_T$  time series from 2019 to present, as shown in Fig. 10.1 (d).

A BiDireccional Long Short-Term Memory (BD-LSTM) neural network (see Fig. D.4) was selected as the best recurrent neural network architecture to reconstruct the  $\text{pH}_T$  time series in the Bay of Palma. The training process was successfully completed with no signs of overfitting, achieving less than 1% error in both training and validation sets (Fig. 10.2 (a)). The BD-LSTM neural network was able to fairly predict the majority of the individual pH data points in the time series, although there are some deviations (Fig. 10.2 (b)). Furthermore, the time series pattern is perfectly captured by the neural network (Fig. 10.2 (c)). Notice the gaps in the reconstructed pH points (in red) in

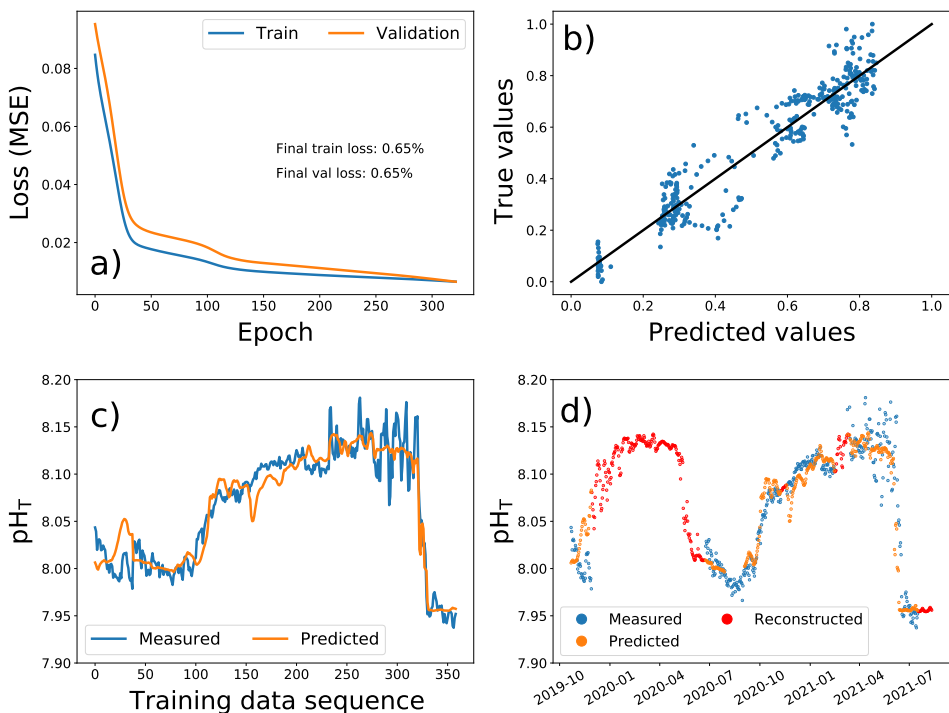
Fig. 10.2 (d), that are those for which the full ambient time series is not available. Finally, the reconstructed pH data using the BD-LSTM model was used to assess the decadal trend of acidification in Palma Bay, which yielded  $-0.0020$  units of pH per year (black line in Fig. 10.2 (d)). Indeed, to further characterize this decadal trend, 1000 independent training-prediction processes were carried out using a BD-LSTM neural network. The results showed a mean slope of  $-0.0020 \pm 0.00054$  for the decadal acidification trend (see methods).



**Figure 10.2:** Bidirectional LSTM neural network model applied to assess the decadal pH<sub>T</sub> trend in the Bay of Palma. a) Training process monitoring loss for both training and validation sets. b) Predicted pH<sub>T</sub> values against their true values, where the black line is the reference for a perfect prediction. c) Predicted pH<sub>T</sub> time series in the training process (orange) and ground truth series (blue) d) Final prediction for the decadal pH<sub>T</sub> time series using the output data of the trained model and the measured data. Measured pH data are shown in blue, predicted data in the training process are shown in orange and reconstructed data are shown in red. The black line represents the decadal pH trend.

Regarding the Cabrera data set, with the available ambient time series, it is only possible to fill the data gaps, a task for which a BD-LSTM neural

network was also used. As for the Bay of Palma the training process was successfully completed with no signs of overfitting, yielding less than 1% error in both training and validation dataset (Fig. 10.3 (a)). The model fairly predicts most of the individual  $\text{pH}_T$  data points in the training dataset, showing some deviations as usual (Fig. 10.3 (b)). The tendency of the time series is perfectly captured by the model (Fig. 10.3 (c)) and thus the gap can be filled with reliable data, red points in (Fig. 10.3 (d)).



**Figure 10.3:** Bidirectional LSTM Neural Network model applied to fill the gaps in the  $\text{pH}_T$  time series in Cabrera. a) Training process monitoring loss for both training and validation sets. b) Predicted  $\text{pH}_T$  values against their true values where the black line is the reference for perfect prediction, c) Predicted  $\text{pH}_T$  time series in the training process (orange) and ground truth series (blue) and d) Gaps in the  $\text{pH}_T$  time series filled with the trained model (red), while measured  $\text{pH}_T$  are shown in blue and predicted data in the training process shown in orange.

## 10.3 Discussion

The achievement of long-term oceanographic data series suitable to evaluate the effects of climate change constitutes a great operational effort that is unequivocally accompanied by partial data loss due to multiple factors (human and instrumental). The advances in the development of pH sensors are enabling the acquisition of precise pH data without identified drift through highly accurate indicator-based spectrophotometric methods [456]. However, in order to determine OA trends, several years of quality seawater pH data are needed, adding more difficulty to the vicissitudes inherent to field work. Recently, the application of computational methods based on Deep Learning (DL) is becoming a useful tool to fill the gaps due to data loss. Several studies have implemented the DL methodology and successfully predicted bio-optical and biogeochemical parameters [451, 453–455, 457–461].

Here, the application of a BiDireccional Long Short-Term Memory (BD-LSTM) neural network to predict  $\text{pH}_T$  from physical data, namely temperature, salinity, and dissolved oxygen, the latter as a key indicator of biological activity, permitted the reconstruction of gaps in the time series of  $\text{pH}_T$  and allowed the reconstruction of nine years of  $\text{pH}_T$  data. The BD-LSTM architecture has been proved extremely effective in predicting sequence data, such as time series, as they combine the information for both front and back directions of time [462] and is more effective (accurate and stable) compared to unidirectional Long Short-Term Memory neural networks. In this study, the BD-LSTM offered better estimation results over the other neural networks considered, both in time series reconstruction and missing data filling.

In the Cabrera station, the BD-LSTM applied permitted a reliable reconstruction of the gaps in  $\text{pH}_T$  data from December 2019 to June 2020 (Fig. 10.1 (d)), constituting an advantageous methodology to support the acquisition of long time series data without losing accuracy, as the model can reproduce  $\text{pH}_T$  data with an error lower than 1% (Fig. 10.3 (b)), closely following the annual variability of the observations (Fig. 10.3 (d)).

The ability of the BD-LSTM to reconstruct time series was observed through the reconstruction of nine years of  $\text{pH}_T$  data in the Bay of Palma station (Fig. 10.2 (d)). The modeled  $\text{pH}_T$  data combined with the observations allowed the accomplishment of a long pH time series in order to estimate a pH trend, seasonally adjusted through a sinusoidal fitting, with a rate of decrease of  $0.0020 \pm 0.00054$  pH units per year ( $R^2 = 0.1$ ,  $p$ -value < 0.001, Fig. D.1), and represents the first estimate of pH trend obtained in the coastal Balearic Sea. Additionally, we applied a linear fit on the reconstructed pH time series, obtaining a trend of  $-0.0025 \pm 0.00053$   $\text{y}^{-1}$  ( $R^2 = 0.01$ ,  $p$ -value < 0.001). This fit was discarded because it was shown to introduce a bias in the pH decrease

trend.

The observed decrease in pH in the Balearic Sea coastal area is well aligned with OA trends reported for open ocean areas, from  $-0.0013$  pH units  $\text{yr}^{-1}$  in the Munida station (New Zealand) to the high trend found in the Cariaco Basin station up to  $-0.0026$  pH units per year [437]. The processes associated with the increased pH decline in the Cariaco Basin were related to the upwelling of subtropical underwater, rich in dissolved inorganic carbon, thus lowering the pH.

In the Mediterranean Sea, previous annual estimates in open ocean areas ranged from  $-0.003$  to  $-0.0044$  [429, 432], reflecting the effect of the hydrodynamical and biogeochemical characteristics of the basin on the seawater pH variability [428, 463, 464]. However, it can be assumed that differences in physical oceanography and ecological processes between areas may modulate local changes of pH. In a coastal Mediterranean area located in the northwestern basin, close to Villefranche-sur-Mer, a rate of pH change of  $-0.0028 \pm 0.0003$  pH units  $\text{yr}^{-1}$  was observed [431] and attributed principally to atmospheric forcing and secondly to increased warming. The calculated trend of pH decrease due to the atmospheric  $\text{CO}_2$  growth during the period of this study, from 2013 to 2021, was of  $0.0025 \pm 0.0002$  pH units per year ( $R^2=0.95$ ,  $p\text{-value}<0.001$ ), consistently related to the seawater pH decline. Therefore, these analyses suggest that the atmospheric forcing is the main driver responsible for the pH-decreasing trend found in the surface coastal Balearic Sea. Subsequently, the difference between the seawater pH decreasing trend obtained and the pH trend calculated from the atmospheric levels could be related to natural biogeochemical processes, not distinctly quantifiable with the available length of the Bay of Palma pH time series.

In addition, the effect of temperature on surface ocean pH can be considered. This occurs directly through the temperature dependence of the seawater  $\text{CO}_2$  chemistry as changes in temperature and salinity influence the equilibrium constants of the oceanic  $\text{CO}_2$  system and indirectly through air-sea exchange of  $\text{CO}_2$ . The influences of these two temperature processes on surface ocean pH have been found responsible for 50% of the increase in  $[\text{H}^+]$  ions, thus a pH decrease, in the surface layers of the Iceland and Irminger Seas [465]. In the Mediterranean Sea northwestern basin, a temperature increase of  $0.072 \pm 0.022$   $^\circ\text{C yr}^{-1}$  was estimated to be responsible for 40% of the pH decrease [431]. The obtained temperature variability in the Balearic Sea coastal area during this study was of  $0.035 \pm 0.008$  ( $R^2 = 0.008$ ,  $p\text{-value} < 0.001$ , Fig. D.2), indicating that temperature-driven changes could also be assumed to affect the pH trend.

The observed seasonal variability of the data presented a  $\text{pH}_T$  increase from 7.91 during summer up to 8.18 pH units (Fig. 10.1 (d)) in winter seasons,

clearly followed by the TA values (Fig. D.3). Seasonal changes in TA levels in the study area are ranging from around 2350 to 2550  $\mu\text{mol kg}^{-1}$  (Fig. D.3), largely overtaking the seasonal differences reported previously in the Balearic Sea of up to 50  $\mu\text{mol kg}^{-1}$  in total [466]. This discrepancy in variability could be explained by the intense metabolic processes at the coastal location of the Bay of Palma station. This shallow area has a strong coverage of *Posidonia oceanica*, which due to its high ecosystem production [467] could be triggering an increase of pH and TA levels, as seen in salinity normalized TA values (NTA, not shown) during winter-spring, due to the uptake of nitrate and phosphate and the calcium carbonate dissolution [466, 468], and during summer, related to the lower community production [469] an NTA-pH decrease [466].

Another result from this study worth mentioning is the obtained decreasing TA trend in the Bay of Palma of  $-4.0 \pm 0.4 \mu\text{mol kg}^{-1}$  per year. The Western Mediterranean is characterized by lower total alkalinity values with respect to the rest of the basin, which are less salty with low-alkalinity water [470, 471] as a result of the nearby influence of Atlantic waters, which was not expected to influence decreasing decadal TA values. In the northwestern basin, TA values increased over time at a rate of  $2.08 \pm 0.19 \mu\text{mol kg}^{-1} \text{ yr}^{-1}$ . In the Balearic Sea, the decreasing TA confirms the Atlantic forcing on the alkalinity values and the negligible TA discharges due to rivers in the Balearic Islands. There is a marked south-to-north surface gradient in the western region coupled with the west-to-east gradient of alkalinity in the Mediterranean Sea related to the Atlantic influence [466, 472]. Due to a well established linear relationship of TA and salinity [473] and the calculated origin of our values [472] we cannot neglect the strong TA related to the salinity decrease in the study area of  $-0.059 \pm 0.002$  psu per year ( $R^2=0.25$ ,  $p\text{-value}<0.001$ ). This rate is in agreement with the salinity decrease found at the coastal site at Villefranche-sur-Mer ( $-0.0017 \pm 0.0044$  psu  $\text{yr}^{-1}$ ) [431]. Notwithstanding, the intense salinity decrease observed in the Bay of Palma can be linked to a decrease in the intensity of the southern spreading of the Balearic Current through the Ibiza channel (located between Ibiza and Mallorca Islands) driven by mesoscale processes and the prevalence of new Atlantic Water coming from the Strait of Gibraltar [474]. At any rate, this observation is out of the scope of this study, so that further investigation is needed.

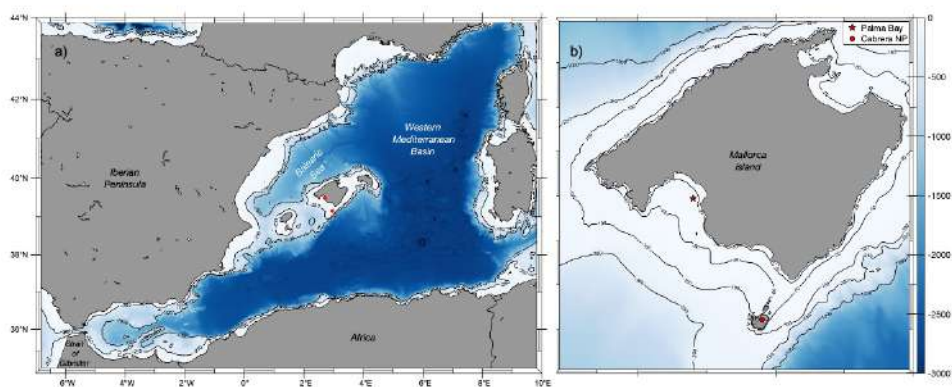
In summary, this work pointed out the useful use of DL techniques, specifically the BD-LSTM architecture, to reconstruct pH data relevant to evaluate seasonal pH variability and to elucidate the climate change consequences, as the OA effect, in a coastal area of the Balearic Sea, which can be extended to the coastal areas of the Western Mediterranean Sea Basin. Nevertheless, future research is necessary to assess and confirm these regional trends, which

highlights the importance of maintaining the time series monitoring networks whose data are the base of this study.

## 10.4 Methods

### Study area

We monitored two coastal stations located in the Balearic Islands in the Western Mediterranean Sea (Fig. 10.4). One site was positioned within the Bay of Palma ( $39^{\circ}29.57088'N$ ;  $2^{\circ}42.02430'E$  Fig. 10.4 (b)) at a fixed station consisting of an oceanographic buoy managed by the Balearic Islands Coastal Ocean Observing and Forecasting System (SOCIB), where meteorological, hydrological, and hydrodynamic data are collected with an hourly frequency since October 2012.



**Figure 10.4:** a) Map of the stations' location in the Western Mediterranean Sea Basin (red dots) and b) detailed location of the Bay of Palma (red star) and the Cabrera National Park (Cabrera NP, red dot) study sites. Maps were developed with the MATLAB<sup>®</sup> R2010b software (<https://mathworks.com>) by using the M\_Map toolbox [475].

The buoy is located at the surface over 20 m bottom depth. The Bay of Palma is a large bay with a surface area of 217 km<sup>2</sup> and approximately 30% seagrass cover [476]. The second station was located at 4 m depth on a mooring line over 8 m bottom depth deployed in the Marine and Terrestrial National Park of the Archipelago of Cabrera ( $39^{\circ} 9.08217'N$ ;  $2^{\circ} 57.04767'E$ ; Fig. 10.4 (b)). The mooring line is in a small bay of just under 1 km<sup>2</sup> and full protection with the largest meadow of the archipelago, covering 89.1% of the surface area between 0-10 m depth [477]. Neither site has important freshwater inputs.

Both stations are part of the Balearic Ocean Acidification Time Series (BOATS) network included in the Interdisciplinary Thematic Platform: Water:iOS.

### 10.4.1 Data collection

In both stations a SAMI-pH (Sunburst Sensors LCC) was attached, at 1 m in the Bay of Palma and at 4 m depth in Cabrera. The pH sensors were measuring pH on the total scale ( $\text{pH}_T$ ) at an hourly rate since December 2018 in the Bay of Palma and since November 2019 in Cabrera. The sensor precision and accuracy are  $<0.001$  pH and  $\pm 0.003$  pH units, respectively. Monthly maintenance of the sensors was performed, including data download and surface cleaning.

Temperature and salinity from the Bay of Palma oceanographic buoy were obtained from October 2012 and for the Cabrera mooring line from November 2019 with a CT SBE37 (Sea-Bird Scientific©) in both stations. The accuracy of the CT is  $0.002$  °C for temperature and  $0.003$  mS  $\text{cm}^{-1}$  for conductivity. Additionally, oxygen data from a SBE 63 (Sea-Bird Scientific ©) sensor attached to the CT in Cabrera and from a YSI 6600V2-4 Multiparameter Water Quality Sonde with a 6450 ROX DO sensor (Yellow Spring Instruments Inc. ©) [478] and a miniDot (PME, Inc. ©) in the Bay of Palma were used. The accuracy of oxygen sensors is  $\pm 2\%$ ,  $\pm 1\%$ , and  $\pm 5\%$  for the SBE 63, the YSI, and the miniDot, respectively.

Periodically, water samplings for dissolved oxygen (DO), pH in the total scale at  $25$  °C ( $\text{pH}_{T25}$ ), and total alkalinity (TA) were obtained during the sensor maintenance campaigns. DO and  $\text{pH}_{T25}$  samples were collected in order to validate the data obtained by the sensors.

DO concentrations were evaluated with the Winkler method modified by Benson and Krause (1983) [479] by potentiometric titration with a Metrohm 808 Titrando with an accuracy of the method of  $\pm 2.9 \mu\text{mol kg}^{-1}$  and with an obtained standard deviation from the sensors data and the water samples collected of  $\pm 5.9 \mu\text{mol kg}^{-1}$ .

$\text{pH}_{T25}$  data was obtained by the spectrophotometric method with a Shimadzu UV-2501 spectrophotometer containing  $25$  °C-thermostated cells with unpurified *m*-cresol purple as an indicator, following the methodology established by Clayton and Byrne (1993) [480] by using Certified Reference Material (CRM Batch #176 supplied by Prof. Andrew Dickson, Scripps Institution of Oceanography, La Jolla, CA, USA). The accuracy obtained from the CRM Batch was  $\pm 0.0051$  pH units and the precision of the method of  $\pm 0.0034$  pH units. The mean difference between the SAMI-pH and discrete samples was  $0.0017$  pH units.

TA samples were collected in 50 ml Falcom vials and poisoned with  $20 \mu\text{L}$  of  $\text{HgCl}_2$  and determined by open cell potentiometric titration with a Titrando 808

(Metrohm) following the Standard Operation Procedure (SOP) 3b [481]. TA values were also calculated from the temperature and salinity values obtained in the Bay of Palma from 2012 by using a second-order polynomial model for TA specifically described for the Mediterranean Basin [472].

pH values due to the atmospheric CO<sub>2</sub> levels were estimated by using the CO2SYSv3 program [482], with the most internally consistent and preferred carbon [483, 484] and sulfate dissociation constants [485] for current surface ocean studies [486], with the Bay of Palma *in situ* temperature and salinity, the calculated TA values, and the atmospheric CO<sub>2</sub> levels converted from dry air to wet [487] as inputs. The carbon dioxide (CO<sub>2</sub>) atmospheric molar fraction used was obtained from the monitoring station of Lampedusa (LMP), Italy, of the NOAA (National Oceanic and Atmospheric Administration, USA) monitoring network [488].

### 10.4.2 Data processing

Once the data were validated, several processing steps were performed to ensure an optimal training process for the neural network models. First, all the data of the time series were re-sampled by averaging the data points, obtaining a daily frequency. Afterwards, a standard feature-scaling procedure (min-max normalization) was applied to every feature (temperature, salinity and oxygen) and to pH<sub>T</sub>. Finally, we built our training and validation sets as tensors with dimensions ( $\text{batch}_{\text{size}}, \text{window}_{\text{size}}, N_{\text{features}}$ ), where  $\text{batch}_{\text{size}}$  is the number of examples to train per iteration,  $\text{batch}_{\text{size}}$  is the number of past and future points considered, and  $N_{\text{features}}$  is the number of features used to predict the target series. Temperature values below  $T = 12.5$  °C were discarded as they are considered outliers in sensor data outside the normal range in the study area.

### 10.4.3 Computing the trend of seasonal data

The trend of seasonal time series is often computed by means of statistical methods based on moving averages or more advanced techniques such as the Seasonal Trend Decomposition Loess [489]. Nevertheless, these procedures do not work with gappy time series, so a different approach is needed. In this work, we fitted the following oscillatory function with trend to our data:

$$y(t) = A \sin(\omega t + \phi) + Bt + C, \quad (10.1)$$

where the parameter  $B$  corresponds to the trend of the data.

Moreover, after this fit, the seasonal component ( $A \sin(\omega t + \phi)$ ) can be removed from the original time series, and a standard linear regression can be performed on the transformed data to obtain the trend (which is exactly  $B$ ) with the  $R^2$  and  $p$ -value estimates given by the linear regression (Figs. D.1 and D.2).

### 10.4.4 Selecting the best neural network architecture

Several recurrent neural network (RNN) architectures were considered as candidates to reconstruct the pH time series, including a Simple Recurrent neural network (SRNN), Long-Short Term Memory (LSTM), BiDirectional LSTM (BD-LSTM; Fig. D.4), and BiDirectional Gated Recurrent Unit (BD-GRU).

Initially, manual tests were performed on each architecture to determine the optimal set of parameters that yielded the best possible results. These tests were based on minimizing the errors in both the training and validation sets while avoiding overfitting. To avoid overfitting, we implemented automated callbacks to stop the training process whenever the validation loss increased or crossed the training loss. During this test, we determined the minimum number of nodes, which helps in avoiding overfitting, and the minimum window size, which allows to use the most possible number of data points for training and prediction. All the RNNs were trained in batches of size 32. To enhance clarity and accessibility, the optimal values obtained for the more relevant parameters are summarized in Table 10.1.

**Table 10.1:** Optimal parameters used for the different RNN architectures

	Hidden layers	Nodes/Cells	Window size	Activation function	Output function	Loss	Learning rate	Optimizer
SRNN	1	3	6	Tanh	Sigmoid	MSE	0.01	Adam
LSTM	1	3	6	Tanh	Sigmoid	MSE	0.01	Adam
BD-LSTM	1	3	6	Tanh	Sigmoid	MSE	0.01	Adam
BD-GRU	1	1	6	Tanh	Sigmoid	MSE	0.01	Adam

In order to identify the best-performing architecture, an automated procedure was developed to statistically compare the outputs of each model. Each architecture was trained in 1000 independent processes, ensuring a final training mean-squared error of less than 0.8% while avoiding overfitting implementing the previously mentioned callbacks. The code used for the analysis can be found in [291].

In Table 10.2, a summary of the statistical results obtained for each architecture is presented. All architectures provide similar training and validation errors and provide similar results for the decadal  $\text{pH}_T$  trend, predicting a slope of around  $-0.0030$  pH units per year with an intercept of 8.07 pH units. However, the BD-LSTM turns out to be the architecture providing most accurate (the smallest training and validation errors) and precise (the smallest statistical error) results (Fig. D.4). Thus, we selected the BD-LSTM neural network as the best architecture to reconstruct the  $\text{pH}_T$  time series. Data corresponding to the Bay of Palma were used in the selection of the best neural network architecture.

**Table 10.2:** Statistical comparison between different RNN architectures

	<b>Slope</b>	<b>Intercept</b>	<b>Training error</b>	<b>Validation error</b>	<b>Training epochs</b>	<b>Training time</b>
<b>RNN</b>	$-0.0021 \pm 0.00077$	$8.07 \pm 0.006$	$0.54 \pm 0.08$	$0.72 \pm 0.12$	$293 \pm 95$	$15.52 \pm 4.75$
<b>LSTM</b>	$-0.0018 \pm 0.00067$	$8.06 \pm 0.005$	$0.49 \pm 0.03$	$0.68 \pm 0.05$	$245 \pm 68$	$17.55 \pm 4.21$
<b>BD-LSTM</b>	$-0.0020 \pm 0.00054$	$8.07 \pm 0.004$	$0.46 \pm 0.03$	$0.64 \pm 0.04$	$167 \pm 45$	$15.13 \pm 3.00$
<b>BD-GRU</b>	$-0.0020 \pm 0.00066$	$8.07 \pm 0.005$	$0.51 \pm 0.07$	$0.74 \pm 0.10$	$347 \pm 95$	$27.68 \pm 6.84$

The code and data used to determine the best neural network architecture can be found in a GitHub repository [490].



An underwater photograph showing a dense field of green seagrass meadows. The water is clear and blue, with sunlight filtering through from above, creating a bright, slightly hazy atmosphere. The seagrass blades are long and narrow, swaying gently in the water.

## 11. Mapping seagrass meadows from space

**Published as:**

À. Giménez-Romero, D. Ferchichi, P. Moreno-Spiegelberg, T. Sintes, and M. A. Matías, "Mapping the distribution of seagrass meadows from space with deep convolutional neural networks", [bioRxiv \(2024\)](#)

## 11.1 Introduction

Coastal ecosystems, encompassing seagrasses, mangroves, salt marshes, and coral reefs, among others, provide invaluable services that contribute to supporting the livelihoods of coastal communities, impacting the well-being of the residents [16, 491]. Seagrass meadows, in particular, are crucial for enhancing coastal biodiversity, protecting shorelines from erosion, and mitigating climate change by sequestering large quantities of carbon [157, 492]. However, if these habitats are degraded, they could leak stored carbon into the atmosphere and further accelerate global warming [492, 493].

In fact, despite the considerable uncertainty surrounding global seagrass extent values, it is estimated that about one third of seagrass global extent has been lost since World War II [492]. Seagrass declines are primarily attributed to eutrophication, water quality degradation, habitat destruction, and climate change, particularly global warming [159]. Furthermore, the sensitivity of seagrasses to future ocean temperatures under different emission scenarios poses significant concerns. Models project a decline in the global suitable habitat for these ecosystems throughout the current century, both latitudinal and across water depth, with a notable compression of suitable habitat toward the lower distribution limit imposed by light availability [494].

In this context, the United Nations (UN) recognized the severity of global biodiversity loss and degradation of ecosystems and stressed the negative impact that this situation has on food security, nutrition, access to water, and the health of the rural poor and people worldwide. Accordingly, the UN declared the period 2021-2030 as the “Decade of Ocean Science for Sustainable Development” and the “Decade of Ecosystem Restoration” [495, 496], underscoring the urgency and importance of safeguarding marine ecosystems, including *Posidonia oceanica* meadows. Achieving these targets, particularly concerning the preservation and restoration of coastal ecosystems, requires a rigorous, evidence-based approach to conservation practice and policy. This entails conducting thorough analyses of high-quality monitoring data to inform decision-making and validate intervention strategies.

The comprehensive mapping of several marine habitats, such as coral reefs, kelp forests, deep-sea vent communities, and seagrass beds, has been successfully achieved through the use of side-scan sonar systems [497–500]. This methodology has provided valuable insights into the structure and distribution of these ecosystems and helped to design informed conservation strategies and management practices. However, the cost and time-intensive nature of these methods present challenges in deploying continuous monitoring systems for marine environments. As a result, practical monitoring of biodiversity often occurs infrequently rather than in real-time, preventing a constant spatio-temporal

evaluation of the status of these ecosystems.

A recently emerging possibility is to combine remote sensing technologies with available georeferenced habitat data to develop correlative or mechanistic models that are then capable of monitoring biodiversity at finer temporal scales. Among various methodologies, Machine Learning (ML) models trained with multi-spectral satellite imagery data appear to be the most promising [501–504]. In particular, many efforts to determine the spatial distribution of *Posidonia oceanica* from airborne imagery using ML have been recently made [505–522]. However, despite the seminal insights of many of these works, they present numerous limitations that hinder the delivery of functional models suitable for a real-case deployment, serving merely as potent proof of concept for the methodologies studied. For instance, many of these studies rely on inadequate metrics for evaluating model performance in image segmentation problems, such as accuracy, leading to an overestimation of the model's performance and neglecting more suitable and demanding metrics like Intersection over Union. Additionally, ground truth data predominantly rely on photointerpretation, often with a limited number of validation points obtained from field data, undermining the models' robustness. Several studies employ only a single or few satellite images, limiting the models' generalizability, while there's a prevalence of simplistic ML methods like Supported Vector Machines and Random Forests for image segmentation tasks, despite the suitability of Deep Convolutional Neural Networks (CNN) for such tasks being well-established [523, 524]. Thus, while these studies lay the groundwork for innovative methodologies, further research is imperative to develop robust and scalable models capable of meeting the demands of real-world applications in biodiversity monitoring and management.

To reach this goal, three key considerations need to be addressed. Firstly, an extensive georeferenced habitat dataset must be employed, acquired through a meticulous and consistent methodology. This dataset should cover a broad geographical area and encompass various spatial scales to ensure the representation of diverse ecological conditions. Secondly, deep learning models, preferably based on convolutional neural networks (CNNs), should be trained using a diverse set of satellite images. These images should incorporate variations in acquisition dates, geographic locations, and the positions of satellites relative to Earth and the sun. This approach enables learning under real-world conditions and enhances the robustness of the models. Lastly, the generalization capability of the models must be evaluated by testing their predictive performance across regions that are geographically distinct from the training dataset. These regions should be characterized by different environmental conditions, ensuring the reliability and applicability of the models in varied real-world scenarios.

Here, we present a comprehensive framework that addresses these consider-

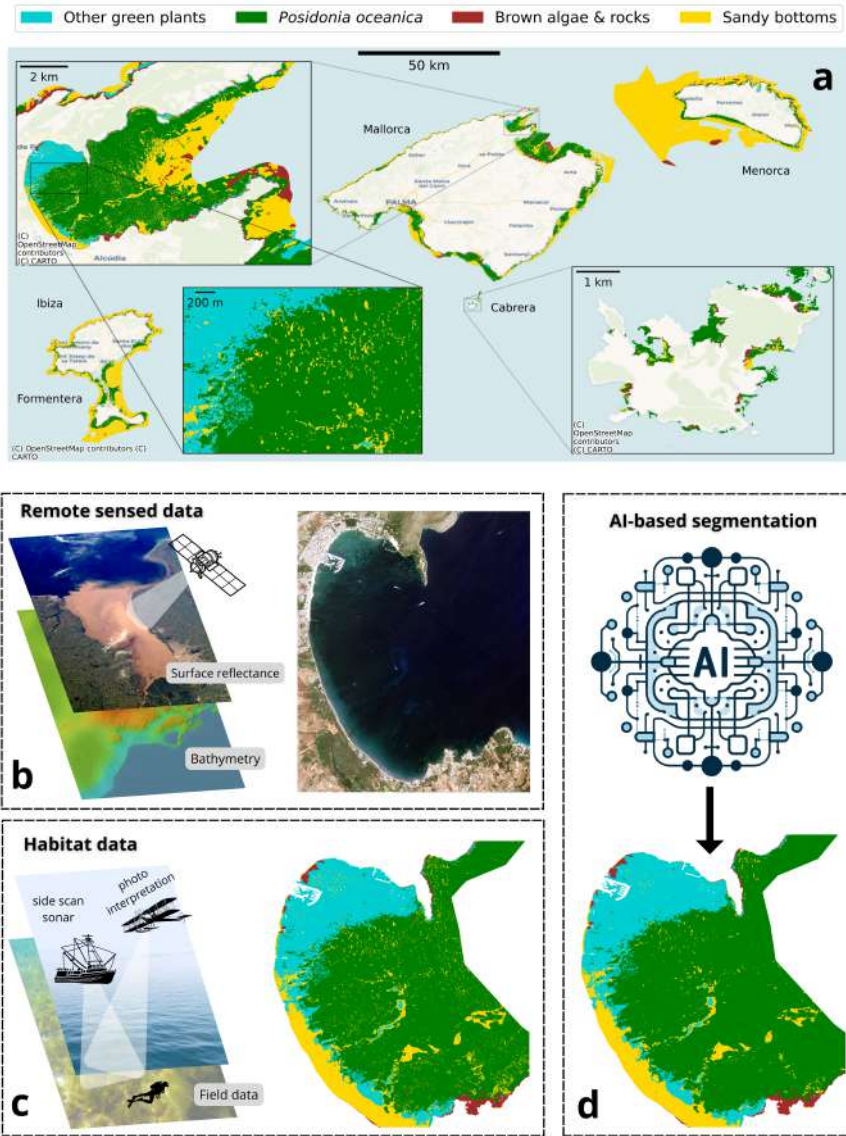
ations, providing a robust and reliable model for the classification of *Posidonia oceanica* meadows and related habitats in the Mediterranean Sea using satellite imagery. We demonstrate the model's generalization capability and robustness by training only with data from a particular region of our extensive dataset and evaluating its performance on the other regions. We show that our model is capable of providing reliable estimates for the distribution of the considered habitats and accurate measures for their extension area. In addition, we measure the model's loss of accuracy in its estimates for new regions, providing a lower bound for model performance. This is a crucial step to advance in the development of a reliable map of the distribution of *Posidonia oceanica* meadows in the Mediterranean Sea. Finally, we train the model with all available data so that the resulting model can be used to classify *posidonia* meadows in other regions of the Mediterranean Sea.

## 11.2 Results

### 11.2.1 A deep learning framework for automated marine ecosystem labelling

We developed a deep learning framework based on convolutional neural networks to accurately classify benthic habitats in the Mediterranean Sea using satellite imagery (Fig. 11.1). We used a comprehensive and extensive habitat dataset of the Balearic Sea, comprising a 20-year effort of data acquisition based on side-scan sonar supported by photointerpretation of high-resolution airborne imagery and in-situ observations (Fig. 11.1 a). The dataset covers about 2,500 km<sup>2</sup> of the coastal habitats of the Balearic Islands at high spatial resolution and contains 28 different classes, including the ecologically significant species *Posidonia oceanica*, which were aggregated in 4 major ecological groups: *Posidonia oceanica*, Other green plants, Rocks & brown algae, and Sandy bottoms (Fig. 11.1 a, c, Methods & Appendix E.1). This dataset was combined with satellite imagery of the coastal areas of the Balearic Islands acquired from PlanetScope [525], covering around 1200 km<sup>2</sup> with different dates and satellite positions (Fig. 11.1 b, c, Methods & Appendix E.1).

We trained 40 different deep learning models using 4 different state-of-the-art architectures and 10 different backbones for each architecture (Fig. 11.1 d, Methods & Appendix E.3). Furthermore, we implemented a consensus prediction approach to enhance the robustness and reliability of model predictions, which involves aggregating the results from multiple deep learning models to mitigate potential biases introduced by individual models (Methods). To evaluate the models, we trained them with the data from only one island and performed a posterior systematic study of their performance on the other islands.



**Figure 11.1:** (a) Spatial distribution of the 4 main ecological benthic habitats in the Balearic Sea, present in the whole Mediterranean. The dataset provides detailed information at multiple spatial scales up to 3 m resolution. (b-d) Scheme of the pipeline to train the CAMELE model. Satellite-based surface reflectance data is merged with bathymetric estimates to produce the inputs (features) of the model. Habitat data obtained with side-scan sonar, photointerpretation and field observations were used as ground truth data (labels). These features and labels are used to train deep convolutional neural networks to perform image segmentation.

Thus, the train-test split was roughly 50%-50% rather than the traditional 80%-20% split, with the test set representing diverse environmental conditions and benthic habitats formed by slightly different species than the training set (Methods & [Appendix E.2](#)). This approach was chosen to simulate real-world scenarios in which one cannot control for specific environmental conditions, constrained dates for image acquisition, the position of the satellite with respect to the sun and the earth, or even find new species not contained in the original training dataset. We thereafter refer to our test set as “out-of-sample” test set and to the model as “Half model” to emphasize this idea.

We performed an extensive evaluation of the models’ performance in the training and out-of-sample test datasets, using a variety of metrics such as Intersection over Union (IoU), Precision, Recall, F1-score, Kappa and Accuracy (Methods). Our results show that the best performing framework was to use the 10 models defined by the Linknet architecture together with the consensus prediction approach (Methods and [Appendices E.6](#) and [E.7](#)), which hereafter we refer to as CAMELE (Consensus for Automated Marine Ecosystem Labelling and Evaluation).

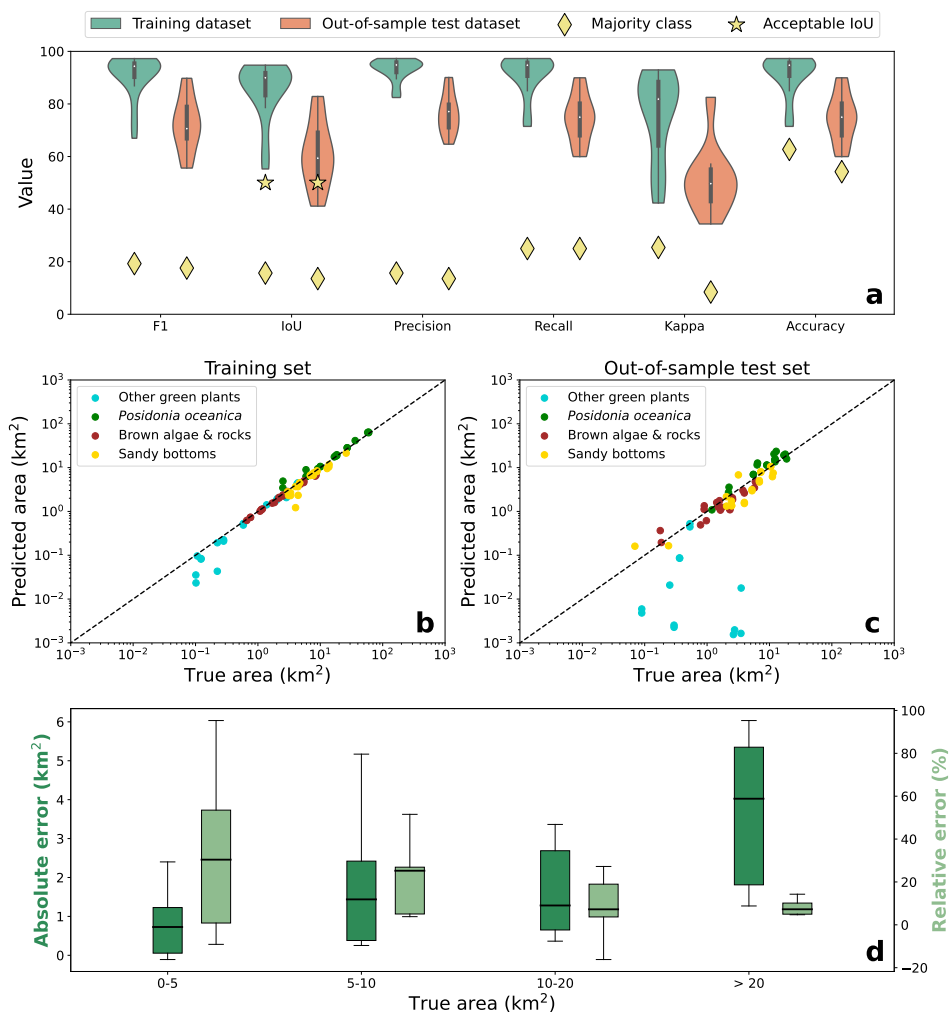
### 11.2.2 A reliable AI-based solution for marine ecosystem monitoring

CAMELE’s performance in both training and out-of-sample test datasets was highly notable, with a mean IoU score of 88.22% and a mean F1-score of 93.13% in the training dataset, compared with a mean IoU score of 61.97% and a mean F1-score of 72.77% in the out-of-sample test dataset ([Fig. 11.2 a](#) and [Tables E.5](#) and [E.6](#)). We note that an IoU score greater than 50% is considered an acceptable prediction in image segmentation tasks [526] (yellow stars in [Fig. 11.2 a](#)). Furthermore, the model outperforms by a large margin the naive baseline of predicting only the majority class (yellow diamonds in [Fig. 11.2 a](#)). We observe an overlap between the distribution of the performance metrics in the training and out-of-sample test datasets, showing that model performance is consistent in both sets ([Fig. 11.2 a](#)). The model was able to segment some images in the out-of-sample test dataset with notable performance (e.g., 15% of the images with an IoU score higher than 80% and 20% of the images with an IoU score higher than 70%), while only 10% of the images had an IoU score lower than 50% ([Fig. 11.2 a](#) and [Table E.8](#)). This demonstrates that the model is able to generalize to some extent to new regions, with different environmental conditions and the presence of some different benthic habitats.

The decrease in performance in the out-of-sample test dataset can be attributed to the different environmental conditions, including the presence of different benthic habitats not included in the test dataset. However, we observed that a significant part of the pixels categorized by the ground truth data

as Other green plants, Brown algae & rocks, or Sandy bottoms were being classified by the model as *Posidonia oceanica*, substantially affecting the overall performance of the model (Appendix E.7 and Fig. E.3). Surprisingly, we found that the distribution of response values for the Other green plants class in the test dataset is much more similar to the distribution of response values of the *Posidonia oceanica* class in the train set than to its own class (Appendix E.7 and Fig. E.4). Then it is not surprising that the model classifies all those samples as *Posidonia oceanica*. In contrast, the model achieved notable performance in segmenting the *Posidonia oceanica* class, with a mean IoU of 77.30% compared with the mean IoU of 91.97% achieved in the training dataset (Table E.7). At any rate, the model still achieved an overall notable performance in the out-of-sample test dataset, demonstrating the generalization capability and robustness of CAMELE and highlighting its potential for real-world applications in biodiversity monitoring and management.

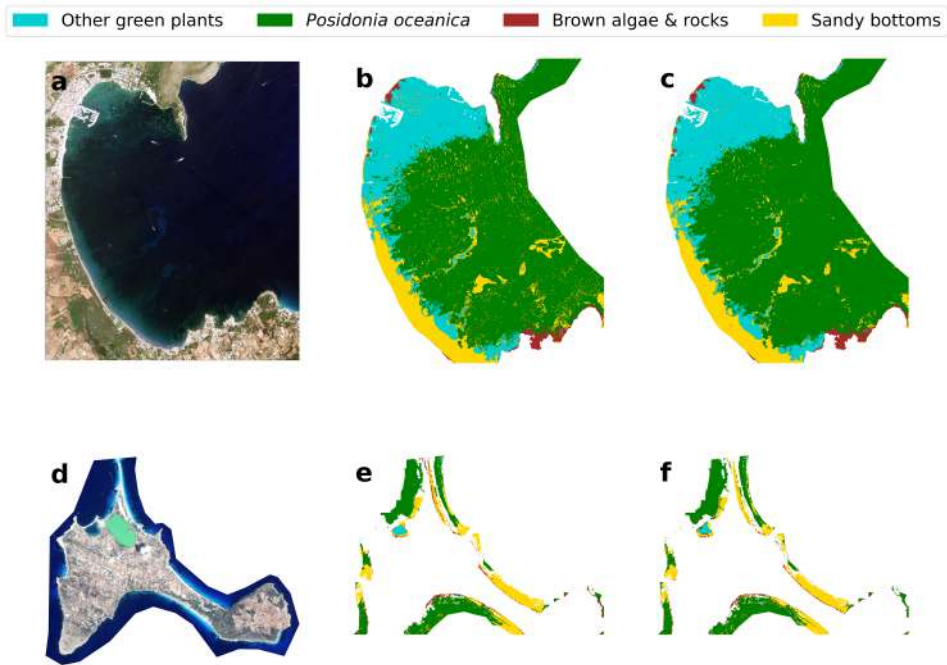
CAMELE's performance was further evaluated by comparing the true and predicted area for each habitat class in each image of the training and out-of-sample test datasets. The model achieved a notable performance in predicting the area of the different habitat classes except for the Other green plants class, as expected from the previous analysis (Fig. 11.2 b, c). Specifically, the median absolute errors committed in the prediction of the area of the different habitat classes were 0.98 km<sup>2</sup>, 0.06 km<sup>2</sup>, 0.15 km<sup>2</sup>, and 0.70 km<sup>2</sup> in the training dataset, compared with 2.24 km<sup>2</sup>, 0.26 km<sup>2</sup>, 0.47 km<sup>2</sup>, and 1.12 km<sup>2</sup> in the out-of-sample test dataset for the *Posidonia oceanica*, Other green plants, Rocks & brown algae, and Sandy bottoms classes, respectively. However, the relative errors were 5.61%, 11.77%, 6.77%, and 14.34% in the training dataset, compared with 24.92%, 99.20%, 28.16%, and 35.05% in the out-of-sample test dataset. Of course, the high relative errors for the Other green plants class in the out-of-sample test dataset are nonsensical, as the true area of this class is small, leading to a high relative error. Thus, we observe that the absolute errors doubled in the out-of-sample test dataset, while the relative errors increased by a factor of 5. At any rate, the models' performance in predicting the area of the *Posidonia oceanica* class was particularly notable, with relative errors significantly decreasing with the extent of the area to be predicted, linked to the wider spatial context available (Fig. 11.2 d). For instance, the median relative error for true extent areas between 1 and 5 km<sup>2</sup> was 30% (86% for 95% confidence interval) compared with 7% (86% for a 95% confidence interval) for areas larger than 20 km<sup>2</sup>. This finding underscores the importance of considering the spatial context when predicting the area of benthic habitats, highlighting the potential of CAMELE to provide reliable estimates for the distribution and extension of the considered habitats in the Mediterranean Sea.



**Figure 11.2:** Model performance in train and out-of-sample test datasets. (a) Violin plots for F1-score, IoU, Precision, Recall, Kappa and Accuracy in both training and out-of-sample test datasets. (b-c) True vs predicted area for each habitat class in the training (b) and out-of-sample test (c) datasets. The diagonal dashed line indicates perfect prediction. (d) Box plot for the absolute and relative errors committed in the prediction of *Posidonia oceanica* area as function of the true area. Relative errors significantly decrease with the extent of the area to be predicted, linked to the wider spatial context available.

To further illustrate the model's performance, we present an example of model predictions for a satellite image in the training and out-of-sample test sets (Fig. 11.3). The model accurately classified the different benthic habitats in

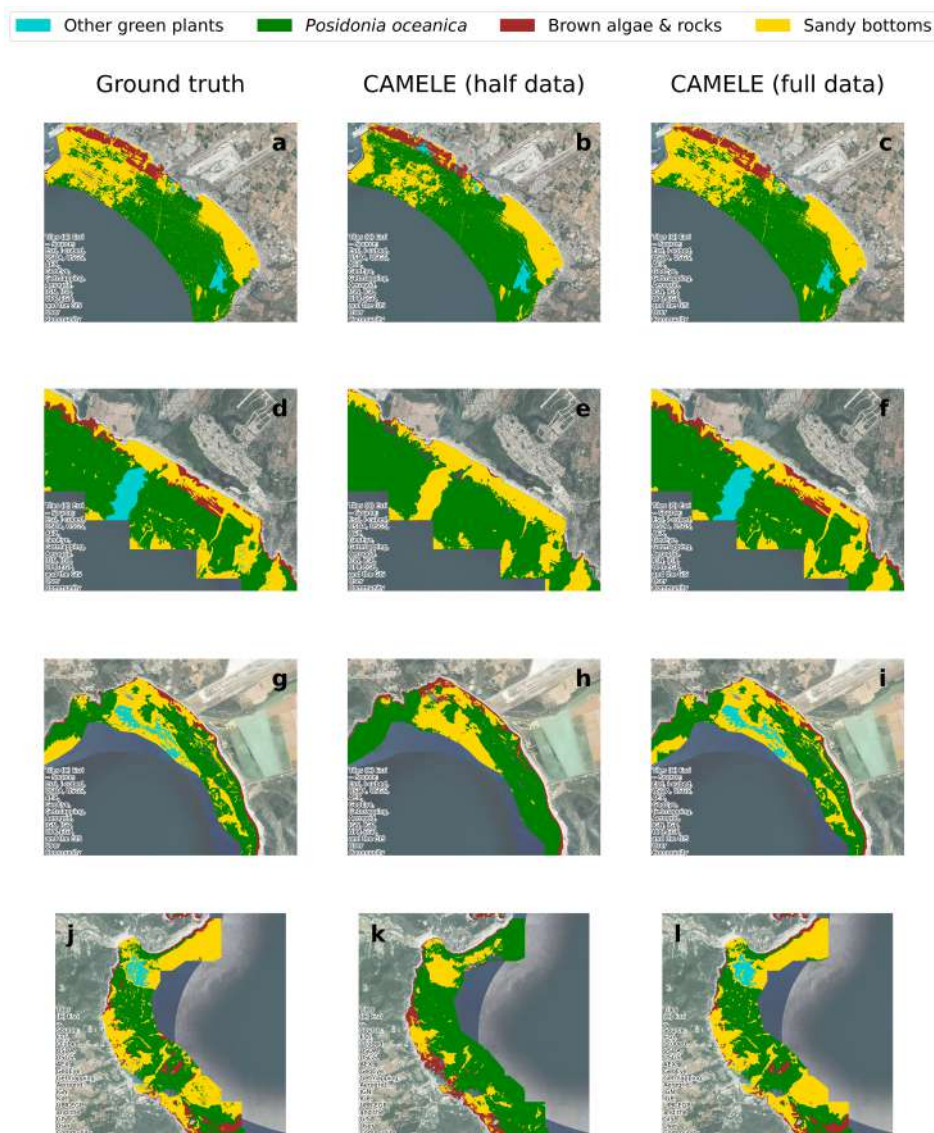
both cases with 92.78% and 82.54% IoU scores, respectively, providing reliable estimates for the distribution and extension of the considered habitats. We note that although there is a 10 point difference in the IoU score between each prediction, this is almost unobservable in the visual inspection of the predictions, highlighting the extreme sensitivity of this metric to small differences.



**Figure 11.3:** Example of model predictions for a satellite image in the training and out-of-sample test set. (a) Satellite image from Pollença Bay on the island of Mallorca, a part of the training set. Image © 2022 Planet Labs PBC (b) Ground truth data for the benthic habitats in Pollença Bay. (c) Habitat classification from the CAMELE model in Pollença Bay (92.78% IoU). (d) Satellite image from Formentera, part of the out-of-sample test set. Image © 2022 Planet Labs PBC (e) Ground truth data for the benthic habitats in Formentera. (f) Habitat classification from the CAMELE model in Formentera (82.54% IoU).

### 11.2.3 Towards a comprehensive model for the Mediterranean Sea

Finally, we trained CAMELE with all available data (using 13144 patches for the actual training and 3286 for the validation set), providing the scientific community with our final trained models that are freely accessible at [527]. We thereafter refer to this model as the “final” model.



**Figure 11.4:** Example of model predictions for a satellite image in the complete dataset. The first column shows the ground truth data for the benthic habitats, the second column shows the habitat classification from CAMELE model trained only with half of the data (from Mallorca island), and the third column shows the habitat classification from CAMELE model trained with all available data. (a-c) Palma bay, Mallorca. (d-f) Son Bou beach, south-east of Menorca. (g-i) Es Còdols, south of Ibiza. (j-l) Cala San Vicente, east of Ibiza.

We evaluated the performance of the final model in the complete dataset and, for comparison, also in the previous training and out-of-sample test datasets, achieving a median IoU score of 95.22%, 94.73%, and 96.22%, respectively (Table 11.1). Notably, the model's segmented all the images in the complete dataset with an IoU score higher than 90%, with a mean, median, and maximum IoU score of 94.64% 95.22%, and 98.5%, respectively (Table E.9). Finally, we assessed the model's robustness by predicting on a new set of images from the Balearic Islands in 2023, not used in the training or out-of-sample test datasets. The model achieved a remarkable mean IoU score of 80% in this new dataset (Tables E.10 and E.11).

**Table 11.1: Final model performance.** Performance (IoU score) of the final model in the complete dataset and comparison with the previous model performance and previous data splits conforming the training and out-of-sample (OOS) test datasets.

	<b>Train</b>	<b>OOS test</b>	<b>Complete dataset</b>
<b>Half model</b>	88.22	61.97	79.21
<b>Final model</b>	94.73	96.22	95.22

In Fig. 11.4, we present some examples of model predictions in different regions of the Balearic Islands, showing the ground-truth data for the benthic habitats, the habitat classification from the CAMELE model trained only with half of the data (from Mallorca island), and the habitat classification from the CAMELE model trained with all available data. We observe that the main differences between the predictions of the two models occur in the areas with more complex habitat distribution. In these areas, the model trained with all available data is able to capture the complexity of the habitat distribution more accurately, providing a more detailed and reliable classification of the benthic habitats. In any case, the model trained only with data from Mallorca still provides a notable performance in segmenting the *Posidonia oceanica* meadows from the other islands, which in the end is the most important habitat to be monitored. A web-based application for interactively visualizing all model predictions, together with the ground truth data, is available at [528].

## 11.3 Discussion

The present study represents a significant advancement in the field of marine habitat monitoring, leveraging the synergistic combination of remote sensing data and machine learning algorithms to address critical challenges in biodiversity conservation. Unlike previous studies, we have developed a comprehen-

sive framework for classifying *Posidonia oceanica* meadows from multi-spectral satellite imagery based on deep convolutional neural networks, which in the last years have been shown to be highly effective in image segmentation tasks. This approach takes advantage of both the rich spectral and spatial information provided by satellite imagery and the capacity of deep learning models to learn complex patterns and relationships in the data. In addition, our analyses do not rely on a single or few satellite images but rather on a comprehensive dataset of satellite images, covering a wide geographical area and various spatial scales, to ensure representation of diverse ecological conditions. This approach is crucial for training models under real-world conditions and enhancing their robustness and generalization capability. Indeed, here we make a substantial effort to evaluate the model's generalization capability by testing its predictive performance across regions geographically distinct from the training dataset, characterized by different environmental conditions, to ensure its reliability and applicability in varied scenarios. Of course, this would not have been possible without the extensive and detailed georeferenced habitat dataset used in this study.

One of the key findings of our study is the remarkable generalization ability of our deep learning model across different regions of the Balearic Islands. Despite being trained exclusively on data from a specific area, the model demonstrated the capacity to provide reliable estimates of *Posidonia oceanica* distribution in other islands, where environmental conditions and benthic habitats vary. This highlights the robustness of our approach and its potential for broader applicability in marine habitat monitoring. The ability of the model to generalize across regions is crucial for its practical utility in conservation efforts beyond the training domain. By accurately classifying marine habitats in unseen environments, our model showcases its capacity for real-world application, particularly in scenarios where comprehensive training data may be limited. Furthermore, the provided metrics, including median absolute and relative errors in *Posidonia oceanica* area prediction, offer valuable insights into the model's performance and potential limitations. By quantifying prediction errors and their relationship to true area estimates, we establish a foundational understanding of the model's accuracy and reliability, enabling more informed decision-making in habitat monitoring and conservation efforts.

The combination of remote sensing and Machine Learning holds promise for revolutionizing biodiversity monitoring, offering unprecedented opportunities for conservationists, policymakers, and researchers to make informed decisions and address pressing environmental challenges. With its ability to provide near-real-time data on ecosystem dynamics, our approach offers a cost-efficient and scalable solution for continuously assessing biodiversity distribution in the Mediterranean Sea, allowing to identify habitat degradation and monitor ecosys-

tem resilience in the face of environmental change. This is particularly crucial in the context of the ongoing climate crisis, where the ability to rapidly detect and respond to habitat loss and degradation is essential for preserving marine biodiversity and ecosystem services. Thus, our model represents a powerful tool in the conservation toolbox, enabling timely interventions supporting the sustainable management of coastal ecosystems and the preservation of biodiversity for future generations.

Despite the notable advancements and potential of our study, several challenges and limitations remain to be addressed. Firstly, the reliance on satellite imagery for habitat classification presents inherent limitations, including cloud cover, atmospheric interference, and limited spatial and spectral resolution, which may impact the accuracy and detail of habitat classification, especially in complex coastal environments. Additionally, the availability and quality of training data pose challenges, as incomplete or biased datasets can impact model performance and generalization capabilities. Moreover, our analysis focused on 4 major aggregated ecological groups, including *Posidonia oceanica* habitats, neglecting other important benthic species and ecosystems that contribute to overall marine biodiversity. Furthermore, while our models exhibit robustness in cross-regional generalization within the Balearic Islands, their applicability to other geographical regions with distinct environmental conditions remains untested. Thus, it may exhibit limitations in delineating finer-scale habitat features in regions with distinct ecological characteristics. Addressing these limitations through continued data collection, model refinement, and validation efforts will be crucial for advancing the reliability and applicability of remote sensing-based approaches in marine habitat monitoring and conservation.

Looking ahead, several avenues for future research and practical applications emerge from our study. Firstly, continued efforts to expand and refine training datasets, incorporating data from diverse geographic regions and ecosystem types, will further enhance the accuracy and generalization capabilities of machine learning models for marine habitat monitoring. Additionally, ongoing advancements in remote sensing technology, including the development of higher spatial and spectral resolution sensors, hold promise for improving habitat classification accuracy and detail. Integration with emerging techniques such as drone-based imaging and LiDAR further expands the scope and resolution of habitat monitoring efforts, enabling finer-scale analysis and management. Moreover, the development of user-friendly tools and platforms for data visualization and decision support facilitates the translation of scientific findings into actionable conservation strategies. By empowering stakeholders with accessible and interpretable information, we can foster greater engagement and collaboration in marine conservation initiatives, ultimately contributing to the sustainable

management of coastal ecosystems and the preservation of biodiversity.

In conclusion, our study represents a significant step forward in the field of marine habitat monitoring, showcasing the transformative potential of remote sensing and machine learning technologies. By advancing our understanding of ecosystem dynamics and supporting evidence-based conservation practices, our work contributes to the broader mission of safeguarding marine biodiversity for future generations.

## 11.4 Methods

### 11.4.1 Satellite data

Satellite imagery was obtained from Planet under the Education and Research Program, which provides limited, non-commercial access to PlanetScope and RapidEye imagery [525]. In particular, we acquired PlanetScope images obtained through the Super Dove (PSB.SD) instrument, which consists of Coastal Blue, Blue, Green I, Green, Yellow, Red, Red Edge and NIR and operates from 2020. Surface Reflectance (SR) products were selected to ensure consistency across localized atmospheric conditions, minimizing uncertainty in spectral response across time and location. SR is derived from the standard Analytic Product (Radiance) and is processed to top-of-atmosphere reflectance and then atmospherically corrected to bottom-of-atmosphere reflectance.

We obtained a total of 60 satellite images along the coast of the Balearic Islands, covering up to 1200 km<sup>2</sup> of surface area for the years 2020 to 2023 (Appendix E.1 and Fig. E.1). The images were acquired for days with clear sky conditions comprised between June and September, as these are the months in which the biomass of seagrass and algae is more abundant. No other filters were applied to the images, as we aimed to train the model in real-case scenarios in which one cannot control for specific environmental conditions, constrained dates for image acquisition, the position of the satellite with respect to the sun and the earth, etc. Thus, we obtained images with different conditions (Appendix E.1 and Table E.1).

### 11.4.2 Habitat data

We used georeferenced habitat data from the government of the Balearic Islands, corresponding to the outcome of different European and national projects comprising a ~ 20-year effort of data acquisition covering a total of 2,500 km<sup>2</sup> [529, 530]. The first project started around the year 2000, and there has been a major recent update around 2018. The data consist of 28 different habitat classes following the nomenclature and coding of the Standard List of Marine Habitats of Spain (LPHME), obtained from side scan sonar, photointerpretation of airborne

imagery, and in-situ observations. We aggregated the different habitat classes into 4 major ecological groups present in the whole Mediterranean Sea. The aggregation was based on feature similarity and ecological function, resulting in the following classes: *Posidonia oceanica*, green algae, brown algae & rocks, and sandy bottoms (Appendix E.1 and Table E.2).

Fig. 11.1 a shows the spatial distribution of the considered ecological groups in the Balearic Sea. This comprehensive dataset provides detailed information at multiple spatial scales, offering a valuable insight into the intricate spatial patterns of the different habitat classes across the region. The detailed information captured by the high-resolution data contributes significantly to the reliability of our model and the precision of our predictions, particularly vital in the context of habitat conservation for *Posidonia oceanica* meadows.

### 11.4.3 Bathymetric data

Bathymetric data was obtained from the European Marine Observation and Data Network (EMODNET) [531]. EMODnet-Bathymetry provides a service for viewing and downloading a harmonised Digital Terrain Model (DTM) for the European sea regions. The data consist of GeoTIFF layers with  $\sim 100$  m pixel resolution of mean depth values.

### 11.4.4 Dataset creation

Satellite images were processed together with habitat and bathymetric data to construct the final training and testing datasets. The NIR band, which is strongly attenuated by water, was used to mask out pixels corresponding to land using a simple clustering algorithm (i.e., K-means) and subsequently substituted for bathymetric data. The resulting processed satellite images comprise the primary source for the input data of our model. The ground truth (or label) dataset consists of raster files analogous to the satellite image with single-band values representing the benthic class corresponding to each pixel. To construct it, all pixels of each processed satellite image were associated to a given benthic class using the aggregated habitat data. Pixels for which a class was not available were masked out in both the processed satellite image and label data. Similarly, pixels that were already masked in the satellite image were masked in the label image. Finally, patches measuring  $256 \times 256$  pixels ( $\sim 750\text{m} \times 750\text{m}$ ) were extracted from each satellite and label image for the years 2020 to 2022 (keeping the year 2023 as a final test), resulting in a final dataset comprising up to 16,430 patches.

Despite train-test data split usually follows an 80%-20% ratio, training was conducted exclusively with data from Mallorca (8488 patches, of which 1698 were allocated for validation) and tested with the remaining data from Menorca,

Ibiza, Formentera, and Cabrera islands (7942 patches), following a roughly 50%-50% train-test split. Furthermore, our test set represented diverse environmental conditions and even benthic habitats formed by slightly different species than the training set, thereby simulating real-world scenarios that the model may encounter in operational settings. By testing the model on data from regions beyond its training domain, we aimed to assess its generalization ability and robustness. Specifically, we sought to determine whether the model could accurately classify marine habitats and benthic features in unseen environments, probably slightly different from the ones at training, thereby demonstrating its capacity for real-world application. We thereafter refer to our test set as “out-of-sample” test set, to emphasize this idea. All images from 2023 were kept for a final test to analyze model robustness once it is trained with data from all regions in the years 2020 to 2022.

#### 11.4.5 Deep learning models

We trained different state-of-the-art deep learning models for semantic image segmentation, such as UNET [532], Linknet [533], FPN [534], and PSPNet [535]. Each architecture is formed by repeating convolutional blocks, which are usually referred to as the “backbone” of the model. We tested 10 different backbone models for each architecture, leading to the training and evaluation of 40 deep learning models, which represents an unprecedented effort in the field. We used the segmentation-models Python library [536] to define and train all models.

#### 11.4.6 Model training

Before training, the input data was standardized using the mean and variance of the training data,

$$z = \frac{X - \mu}{\sigma} . \quad (11.1)$$

This standard scaling procedure ensures that all input features have a consistent scale, preventing the dominance of certain features during the training process. It is crucial to note that the same standard scaling procedure must be applied during predictions, using the mean and standard deviation computed from the training set. This consistency ensures that the model interprets new data comparably to the training data.

Additionally, for the categorical nature of the output, labels were one-hot encoded. This encoding converts categorical labels into binary vectors, where each class is represented by a unique binary value, facilitating the model's interpretation of the multi-class classification task.

In terms of loss function selection, we opted for dice loss due to its effective-

ness in handling imbalanced datasets, a common characteristic in tasks involving semantic segmentation [537]. Dice loss, also known as the Sørensen–Dice coefficient [538, 539], measures the similarity between predicted and ground truth data by computing the intersection over the union of the two. To further account for class imbalance, we applied loss weights inversely proportional to the proportion of examples of each class. The learning rate was set to 0.001 to ensure a smoother training process.

The models were trained in a computing cluster, using 10 cores and a maximum of 400 GB of RAM for each model. The training process was performed for 1000 epochs with a batch size of 32. The total training time was approximately 1 month for the 40 initial models using the data from the island of Mallorca and about 3 months for the 10 final models using all available data.

#### 11.4.7 Performance metrics

The performance of our trained models was primarily evaluated by the Intersection over Union (IoU) score (Eq. (11.7)), which evaluates the spatial overlap between the predicted image and the ground truth, as this is a suitable metric for image segmentation problems [537]. Additionally, we considered other metrics such as accuracy (Eq. (11.2)), precision (Eq. (11.3)), recall (Eq. (11.4)), F-1 score (Eq. (11.5)), and Cohen’s Kappa (Eq. (11.6)) to perform a comprehensive evaluation of the model. Accuracy gauges the overall correctness of predictions, precision measures the accuracy of positive predictions, recall assesses the model’s ability to capture all positive instances, and the F-1 score provides a balanced assessment of precision and recall. For a binary classification problem, the metrics can be defined from the confusion matrix, where TP, TN, FP, and FN represent true positives, true negatives, false positives, and false negatives, respectively, and  $N$  is the total number of pixels in the image.

$$\text{Accuracy} = \frac{TP + TN}{N} \quad (11.2)$$

$$\text{Precision} = \frac{TP}{TP + FP} \quad (11.3)$$

$$\text{Recall} = \frac{TP}{TP + FN} \quad (11.4)$$

$$\text{F1 Score} = \frac{2 \times TP}{2 \times TP + FN + FP} \quad (11.5)$$

$$\kappa = \frac{2(TP \times TN - FP \times FN)}{(TP + FN)(FN + TN) + (TP + FP)(FP + TN)} \quad (11.6)$$

$$\text{IoU} = \frac{|P \cap L|}{|P \cup L|} = \frac{TP}{TP + FP + FN} \quad (11.7)$$

These metrics collectively offer a holistic understanding of the model's effectiveness in classifying benthic habitats. See Supplementary Information for further details.

### 11.4.8 Consensus prediction

We implemented a consensus prediction approach to enhance the robustness and reliability of model predictions. The consensus prediction involves aggregating the results from multiple deep learning models, each trained with different architectures and backbones. By combining predictions from diverse models, we aim to mitigate potential biases introduced by individual models and enhance the overall accuracy and generalization capabilities.

For each input patch, predictions from all trained models were collected, and a voting mechanism was employed to determine the final consensus prediction. Specifically, the class label with the highest frequency across all model predictions was assigned to each pixel. This ensemble-based strategy leverages the diversity of information captured by different models, leading to a more robust and reliable classification outcome.

### 11.4.9 Model selection

To filter among the 4 different architectures, we evaluated the models' performance in the training dataset. Despite all models achieved high IoU scores ( $> 0.8$ ), Linknet and UNET architectures were the best performing models with a mean IoU of 90.98% and 90.90%, respectively (Table E.3). Because Linknet architecture has less trainable parameters than UNET, being more efficient in terms of computational resources and less prone to overfitting, we selected it as the final architecture to build CAMELE, finally consisting of 10 different models (Appendix E.6).

We then evaluated the performance of each of the 10 models and the consensus prediction approach in the training and out-of-sample test datasets (Appendix E.7). All models achieved high performance, with a median IoU of 88.12% and F1-score of 93.16% in the training dataset (Table E.5), while the models' performance in the out-of-sample test dataset was significantly reduced, with a median IoU of 60.73% and F1-score of 71.87% (Table E.6). When analyzing the performance of the models individually, we observed the emergence of "specialists", which performed significantly better than the rest of the models in segmenting specific classes. This finding underscores the importance of the consensus prediction approach, which leverages the diversity of information

---

captured by different models. The consensus prediction approach significantly improved the models' performance in the out-of-sample test dataset, with an IoU score of 61.97% and a F1-score of 72.77% (Table E.6), highlighting the effectiveness of the ensemble-based strategy in mitigating potential biases introduced by individual models and enhancing the overall accuracy and generalization capabilities of CAMELE. Thus, the consensus prediction approach was selected as the final model for CAMELE.



An aerial photograph of a coral reef system, showing a complex, interconnected pattern of light-colored coral patches and darker blue water channels. A semi-transparent white text box with a thin blue border is overlaid on the lower portion of the image.

## 12. Universal spatial properties of coral reefs

**Published as:**

Àlex Giménez-Romero, Manuel A. Matías, Carlos M. Duarte. "Unraveling the universal spatial properties of coral reefs". *Submitted*

## 12.1 Introduction

Coral reefs form some of the largest biogenic structures in the biosphere [165] and are prevalent structures across tropical coastal waters. Coral reefs often form complex, labyrinthine structures that render sailing in tropical waters challenging but protect the shorelines of tropical coastal nations while supporting biodiversity and providing food supply supporting local communities [540]. The formation of coral reefs has intrigued scientists for a long time, with Darwin formulating a model based on coral reefs accreting on volcanic structures [166]. Darwin's model continues to stir discussion [541], as it would represent a class of coral reefs, at best, among the broad range of configurations and possible origins [542].

Coral reefs can form isolated oval or ring-shaped structures (e.g., coral atolls), linear structures parallel to the shoreline (fringing or barrier reefs) or convoluted or highly-branched structures (e.g., [543]), and often present nested structures, such as smaller reefs within large coral reef lagoons [544]. Whereas previous studies have examined the geometry of coral reefs [543–548], finding evidence of fractality both at individual reef [544, 548] or regional scales [543, 545], a global assessment of coral reef size and geometry was hitherto precluded by the lack of data on reef form and size at the global scale. Resolving the size and geometry of coral reefs has become, however, a matter of urgency, as coral reefs are rapidly declining, with an estimated 50% of coral cover lost globally from 1957–2007 [549]. The IPCC projects that 70% to 99% of coral cover may be lost if climate change reaches 1.5 to 2.0 °C above pre-industrial levels, respectively [550]. Yet, the Kunming-Montreal Biodiversity Framework [551] calls for stopping all biodiversity losses and restoring 30% of degraded habitats, including coral reefs, by 2030, which requires action at individual reefs and, therefore, an understanding of coral reef form and size as an underpinning for the necessary conservation action.

The release of the Allen Coral Atlas (ACA) [500], a worldwide mapping initiative that provides benthic habitat data of shallow-water (above 10 m deep) tropical reefs, presents an unprecedented opportunity to characterize the form and size of coral reefs on a global scale. In contrast to previous efforts, such as NOAA's Coral Reef Information System [552], Khaled bin Sultan Living Oceans Foundation [553], or Millennium Coral Reefs [554], the ACA dataset combines extensive global coverage of reef areas with the direct availability of processed habitat data obtained from recent high-resolution (3 m) satellite imagery using deep learning models. This invaluable resource offers a unique vantage point for researchers, scientists, and conservationists to explore and understand coral reef ecosystems in ways that were previously challenging or impossible. This requires processing the classified images to extract a canonical inventory of individual

coral reefs that can be used to derive macroecological laws and principles from the patterns and trends observed across coral reefs. In turn, this enables us to advance and consolidate our understanding of these complex ecosystems and the scale of the effort required to conserve and restore them.

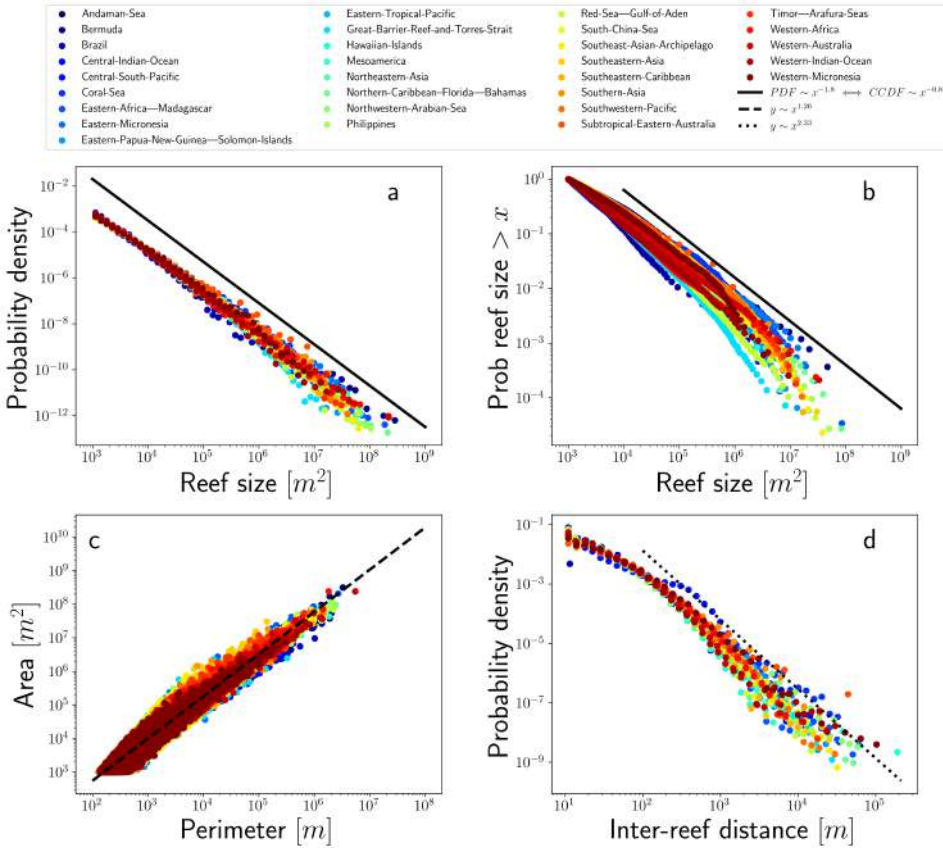
Here we report universal macroecological patterns [46] in shallow-water tropical coral reef size and geometry using data from the 2022 release of the ACA [500]. Specifically, we produced a global inventory of individual reefs [555], allowing the analysis of the size and inter-reef distance distributions of individual coral reef areas within each coral province of the Atlas. We then examined the relationship between reef size, shape, and fractal geometry of coral reefs.

## 12.2 Results

### 12.2.1 Coral reefs macroecological patterns

The area of all coral reefs within each province of the ACA was computed after processing the data to identify individual reefs (see Methods). We identified a total of 1,579,772 individual shallow-water tropical reefs, extending over a total of 52,423  $km^2$  of ocean area. The canonical nature of this openly available dataset (see Methods), which includes all shallow-water tropical coral reefs worldwide, allows the mean and median size of individual reefs to be estimated at 3.32 ha and 0.3 ha, respectively, across all coral reef provinces in the Atlas (see Table F.1 for statistics in each province).

Our analysis reveals that the size distribution of the area of coral reefs in all provinces converges to conform a power-law distribution,  $y \sim x^{-\alpha}$ , that holds over multiple scales, where  $y$  reflects the probability of occurrence of reefs of area class  $x \text{ km}^2$  (Fig. 12.1 a-b). This behavior was consistent across all provinces within a range of 3 to 5 decades in the coral reef area, yielding an average exponent of  $\langle \alpha \rangle = 1.84$  (95% CI: 1.55 to 2.12) (see Methods). The global size-frequency distribution of coral reefs, fitted to all data independently on the province, conforms to a power-law with an exponent of  $\alpha = 1.8$ , consistent with the mean exponent  $\langle \alpha \rangle$ , providing evidence for a universal scaling law governing the size distribution of coral reefs. This suggests that the observed distribution of reef areas arises mainly due to biological processes contributing to reef formation, common among the provinces. Local geophysical processes interacting with coral reef growth, which might vary among provinces, are therefore likely to have a limited role, provided the universality of the power-law scaling of the reef size distribution.



**Figure 12.1: Macroecological patterns of global coral reef size, geometry, and spacing.** (a) Size distribution. The black line corresponds to a fitted power-law of exponent 1.8. (b) Complementary cumulative distribution function (CCDF) with corresponding exponent of 1.8. (c) Area-Perimeter relation. The black dashed line corresponds to a fitted power-law with exponent 1.23. (d) Inter-reef distance distribution. The black dotted line corresponds to a fitted power-law to the distribution tail with an exponent of 2.33.

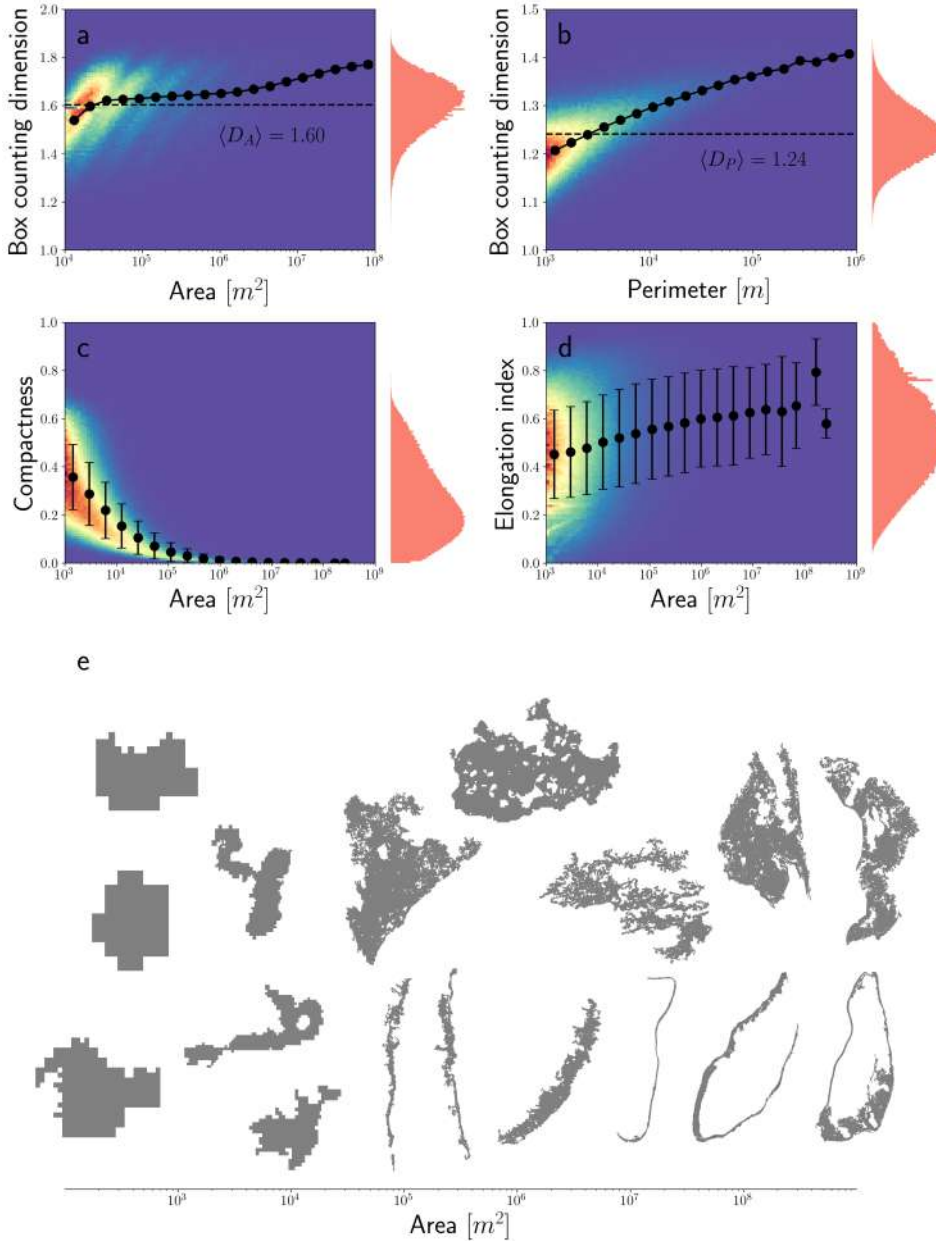
The presence of a power-law size distribution also suggests that the object studied may be fractal in nature [556–560], at least along a certain range. A simple way to test whether coral reefs are fractals is using the area-perimeter relation [561, 562], a method of fractal analysis that characterizes the complexity of irregular shapes by examining the relationship between their area and perimeter (see Methods). The scaling of coral reef area to perimeter of all individual coral reefs also converges into a single power-law with an exponent of  $\alpha = 1.2578$  (95% CI: 1.2573 to 1.2583), indicating again a universal behavior

(Fig. 12.1 c). When the relationship is fitted for each province independently, we found a mean exponent of  $\langle \alpha \rangle = 1.2574$  (95% CI: 1.1757 to 1.3391), practically identical to the general exponent. The exponent for the area-perimeter relationship is significantly lower than the value of 2 that would be found for a smooth Euclidean geometry. According to fractal analysis, the fractal dimensions of the perimeter and area of coral reefs are given by  $D_P = 1.2950$  and  $D_A = 1.6289$  (see Methods), respectively, which are clearly different from the putative euclidean dimensions of  $D_P = 1$  and  $D_A = 2$ . Taking each province independently, these fractal dimensions would vary between  $D_P^{\max} = 1.3506$ ,  $D_P^{\min} = 1.2468$  and  $D_A^{\max} = 1.6696$  and  $D_A^{\min} = 1.5879$  within a 95% confidence interval, further ensuring that the obtained fractal dimensions are significantly different from the euclidean geometry. Coral reefs develop fractal-like geometries, exhibiting complex, self-similar structures across different scales. This feature might arise from the complex physical and ecological processes that shape the distribution and growth of coral reefs. However, how this occurs remains largely unknown, as we lack mechanistic models able to generate coral reef landscapes across scales and time that can be challenged to reproduce these patterns.

The spatial distribution of coral reefs within each province was also investigated by means of the inter-reef distance, defined as the minimum distance between a reef and its nearest neighbor. We find a heavy-tailed relation where the tail conforms to a power-law with an exponent of 2.33 (Fig. 12.1 d). This reveals that most of the reefs are close to each other, while a non-negligible number of them are isolated. This finding is again mostly independent on the geographic location of the analyzed coral reefs, arising as a universal property of coral reef provinces.

### 12.2.2 The fractal nature of coral reefs

We computed the fractal dimension of the area and perimeter of each individual reef from all provinces using the well-known box-counting algorithm (see Methods). The mean values for the fractal dimension of the perimeter,  $D_P = 1.24$  (95% CI: 1.13 to 1.35) and the fractal dimension of the area,  $D_A = 1.60$  (95% CI: 1.39 to 1.81), are well-defined and consistent with those obtained from the area-perimeter relationship (Fig. 12.2 a, b). The fractal dimension for the reef areas is quite stable around the mean, although it increases slightly for large coral reefs (Fig. 12.2 a). On the other hand, the fractal dimension for the perimeter shows a more pronounced increase as a function of reef size. This indicates that as coral reefs grow, their contour gets more and more convoluted, increasing its complexity, while its surface remains geometrically more stable.



**Figure 12.2: The fractal nature of global coral reefs.** 2D histograms from all shallow-water coral reefs worldwide for: (a) the surface fractal dimension and area; (b) the perimeter fractal dimension and perimeter; (c) the compactness and area and (d) the elongation index and area. The black line corresponds to the mean values of the Y axis measure as a function of the X axis measure. The red histogram corresponds to the distribution of the Y axis measure. (e) Example of coral reefs shape as a function of their surface.

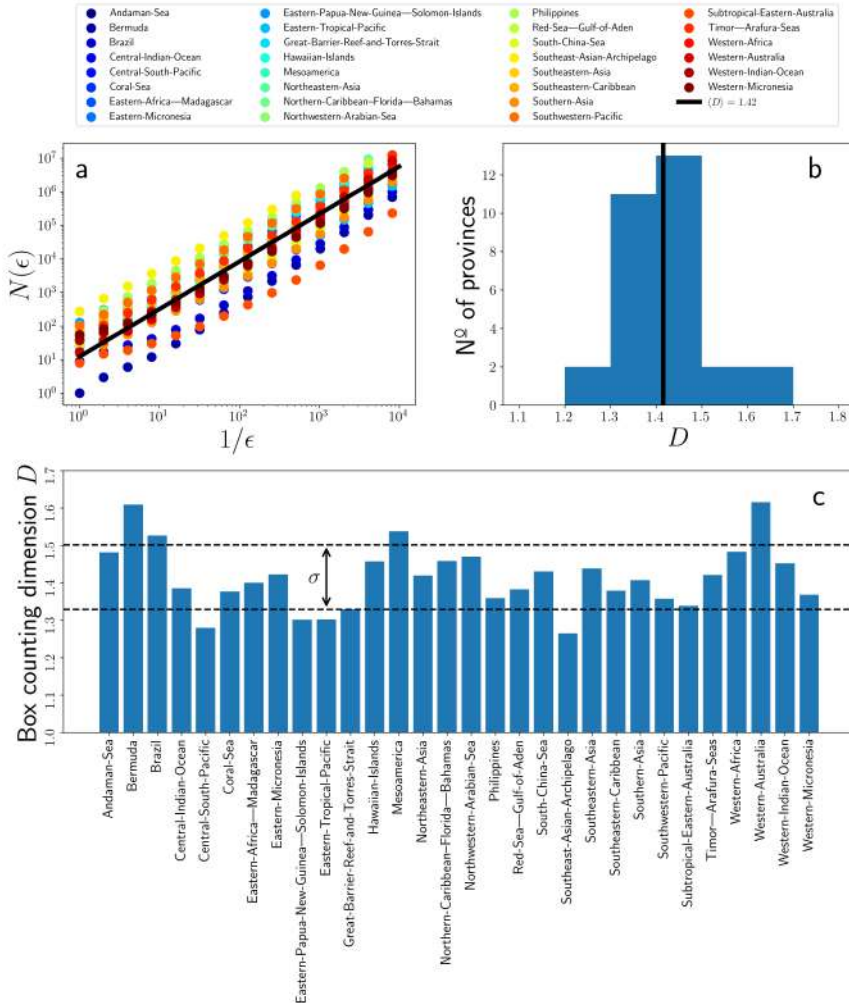
To further understand the reef formation process from a geometric perspective, we computed other shape measurements such as compactness and elongation indices (see Methods). We observe that the compactness of coral reefs decreases rapidly with increasing size (Fig. 12.2 c). Two changes of shape are consistent with this result: transitioning from round to elongated shapes or keeping the rounded shape while developing holes, reef lagoons, within their surface. These two processes are not necessarily mutually exclusive, as elongated shapes could appear from the evolution of empty, rounded shapes. The results obtained for the elongation index measurement show that both processes must occur, as reef elongation increases with size while showing very high variance. This hypothesis can be contrasted by examining the different shapes that coral reefs have formed (Fig. 12.2 e).

Altogether, our analyses indicate that coral reefs evolve from simple rounded filled shapes (high compactness and low elongation index) to more complex elongated and less compact forms (low compactness and high elongation index), giving rise to fractal objects with stable surface fractal dimension and increasing perimeter fractal dimension as they grow (Fig. 12.2 e).

### 12.2.3 Fractality extends up to coral provinces

The fractal dimension of coral reefs varies across a range of spatial scales, spanning from individual colonies to entire reef systems [563]. This variability arises from the different physical, biological, and ecological processes that come into play during the development of the different structures present at each organization level. The different coral provinces can be understood as the largest organizational level of these organisms, and the processes involved in maintaining such big structures should be different from those of individual reefs, giving rise to a different fractal dimension.

To investigate this hypothesis, we computed the surface fractal dimension for each coral province as a whole using the well-known box-counting algorithm (Methods, Fig. F.2). We found a mean fractal dimension of  $\langle D \rangle = 1.42$  (95% CI: 1.24 to 1.59), consistent across the different coral provinces assuming a normal distribution (Fig. 12.3). The fractal dimension of coral reef landscapes is similar to those expected from sizes expanding along a Fibonacci series, which yield a fractal dimension of 1.44 [564]. This suggests, again, that coral reefs are self-organized systems that exhibit similar patterns of complexity and irregularity at different scales and that the fractality is an intrinsic property of coral reefs that likely arises from the underlying biological, physical, and ecological processes involved in their growth dynamics.



**Figure 12.3: Box Counting Dimension of coral reef provinces surface.** (a) Scaling of the measure (number of boxes of length  $\epsilon$ ,  $N(\epsilon)$ ) as a function of the ruler used (box length  $\epsilon$ ). The slope of the fit for each province corresponds to its Box-Counting surface fractal dimension,  $D$ . (b) Histogram of the obtained Box-Counting fractal dimensions. The solid black line represents the mean Box-Counting fractal dimension,  $\langle D \rangle$ . (c) Box-counting surface dimension for each coral province. Black dashed lines correspond to a  $1\sigma$  deviation from the mean.

### 12.3 Discussion

Coral reefs self-organize to form macroecological patterns that are largely independent of the geographical location of the reefs, suggesting that the specific

physical conditions to which each coral province is subject have little influence on the development of these scaling laws. Coral reef geometries conform fractal structures that follow a power-law size distribution, with many reefs close to each other while others are completely isolated across coral reef provinces. Overall, our findings provide strong evidence that the power-law size distribution of coral reefs and their fractal geometries are fundamental features of these ecosystems, which reflect the underlying ecological and physical processes that shape them across scales. Furthermore, because the exponent of the power-law size distribution and the fractal dimensions are consistent across different regions, these features are not only fundamental but universal laws reflecting coral reef growth processes.

These universal scaling properties must be used to challenge mechanistic models aimed at reflecting coral reef growth, which hitherto lacked such constraints, both as individual reefs as well as coral reef provinces. For example, at the scale of coral reef province, the characteristic spacing (inter-reef distance distribution) likely influences the interaction between coral reefs, hydrodynamic flows, and the dynamics of the limiting nutrients transported to support photosynthesis and calcification, further constrained by sea level change and available vertical accommodation space [565–567]. Realistic models of coral reef growth and dynamics should be able to reproduce the universal features described here, such as the power law in coral reef size distribution, the fractal geometries, and the changes in reef shape with size.

Power-law distributions have also been found in many natural systems [568–573]. Different mechanisms can produce such emergent power-laws [574]: Self-Organized Criticality, in which the system evolves naturally towards a critical state [575]; Highly Optimized Tolerance, in which the system evolves following a trade-off between yield, cost of resources, and tolerance to risk [576, 577], or correlated noise models, in which external drives and internal dynamics compete on similar time scales, yielding a non-critical steady state characterized by heavy-tailed distributions [578]. Ecological power-laws can arise through a combination of several ecological and physical processes, including competition for resources and random disturbances leading to failure or loss. In [579, 580], it was shown that the emergence of power-laws in vegetation patterns requires the interplay between global competitive and locally facilitative interactions [579, 580]. These power-law distributions are different from those obtained in classical critical systems, where power laws occur exclusively at the transition point [581, 582]. Actually, the ecological power-law reported in [580] only occurs far enough from the (true) critical transition to extinction. It has been conjectured that living systems exhibiting this behavior could draw important functional advantages from operating close to an emergent critical point, namely an opti-

mal balance between robustness and flexibility [583]. In the case of corals, this could reflect a balance among self-organization dynamics driven by competition for space and resources with other species, environmental factors like ocean currents and water quality, predation interactions within the reef ecosystem, and local facilitative interactions, such as energy dissipation. The combination of these mechanisms should be explored in developing models of coral reef formation.

Fractal models have also been used to investigate universal principles that govern the structure and dynamics of complex ecological systems [570], such as the geometry of benthic ecosystems or power-law scaling [584]. In coral reef ecosystems, fractal geometry might emerge as an efficient structure for nutrient acquisition [548]. The fact that  $D_p > 1$  implies that the reef contour is highly convoluted, while  $D_A < 2$  is the result of the development of multiple holes in the reef surface. The combination of these mechanisms maximizes the chance that coral individuals living at the reef surface are in contact with an external flow, thus being able to obtain nutrients. According to [563], coral colonies with a larger space-filling surface and smaller perimeters increase energy gain while reducing the exposure to competitors. A similar argument could apply in the case of coral reefs, which can be included in modeling approaches.

The shape of coral reefs varies with size, and thereby during growth, from more compact, circular structures at small size, observed at the mean area of 3.32 ha, to increasingly elongated structures, which may break the closed shape to form long, linear, fringing reefs. The change in geometry from circular to fringing has been postulated to result mostly from reduced accommodation space as coral reefs grow, and the interaction with sea level [585], but has also been explained as resulting from the interaction between reefs and nutrients transported along hydrodynamic flows from a prevalent direction [566]. Reefs inside reef structures can also be observed, suggesting that as coral reefs expand in size, smaller reefs appear inside them, starting as compact coral heads to then develop empty inner spaces. Overall, this phenomenology suggests that a Turing instability [52, 586] might be present, which indeed has been previously suggested [566], arising due to the interplay between diffusion of nutrient species and nutrient uptake and recycling processes within reefs. However, the Turing mechanism would yield a normal distribution of inter-reef distances, not compatible with the heavy-tailed inter-reef distance distribution found in this study, neither with the obtained power-law scaling.

Overall, our findings have important implications for the understanding of the structure and function of coral reefs, as well as for their conservation and management. The macroecological characterization of universal laws in the geometry of coral reefs, along with the dataset of unique coral reefs globally,

which are reported here for the first time, should help design effective coral reef restoration projects as well as optimize and quantify the effort and resources required.

## 12.4 Methods

### 12.4.1 Global coral reef data

Global-scale coral reef benthic data were obtained from the Allen Coral Atlas (ACA) [500], a publicly available dataset of high-resolution satellite imagery and machine learning-based coral reef classifications. We downloaded the data from the ACA website, which is already divided into the different coral reef provinces. The downloaded dataset consists of GeoJSON files for each coral province with several Polygons and Multipolygons forming the different benthic classes, from which we selected the “coral/algae” class (see [500]). Despite the data is already provided in vector format, the reefs are not identified as individual entities, i.e., a single reef can be formed by many polygons or multipolygon objects. Thus, we processed the dataset with the methods explained below to obtain a representation of individual reefs.

### 12.4.2 Coral reefs as clusters of connected coral/algae class polygons

A label assignment algorithm was developed to identify the different independent (not connected) components forming the coral reefs. Basically, we followed an iterative process in which connected components were assigned the same label, thus being identified as forming the same component. Polygons were considered to be connected if they intersected. To efficiently compute the intersections among polygons we used the Sort-Tile-Recursive algorithm [587] implemented in Python Shapely library [588]. The implementation of the algorithm can be found in the `Preprocessing.py` module at [589].

Coral reefs of less than  $10^3 \text{ m}^2$  were considered as possible noise in the dataset, and thus disregarded. We made this choice based in the fact that the ACA is obtained from satellite imagery of about 3m resolution. Thus, coral reefs of less than this area would represent less than 100 pixels.

### 12.4.3 Coral reefs area, perimeter and inter-reef distance

We computed coral reef area and perimeter using the `geopandas.GeoSeries.area` and `geopandas.GeoSeries.length` methods in `geopandas` Python’s library [590]. For each coral reef, we defined the inter-reef distance as the distance to its nearest neighbor. Thus, the inter-reef distance distribution is obtained after obtaining the nearest neighbor to each reef and computing that distance. To make this computation efficient, we used the Sort-Tile-Recursive algorithm [587]

implemented in Python Shapely library [588]. The implementation of the algorithm can be found in [589].

We note that our estimates of the total area for each coral reef provinces (Table F.1) slightly differ from that directly provided by the ACA because we removed “reefs” smaller than  $10^3 \text{ m}^2$ . Of course, if the area is computed before this data cleaning step the results are identical.

#### 12.4.4 Coral reef size distribution

We fitted the coral reef size data using the powerlaw package in Python [591, 592]. We performed goodness-of-fit tests using a range of alternative distribution models, including log-normal, exponential, and stretched exponential distributions. We found that the power-law distribution (including its truncated form) provided a significantly better fit to the data than any of the alternative models with  $x_{\min}$  ranging from  $10^3 \text{ m}^2$  to  $10^4 \text{ m}^2$  (Tables F.2 and F.3).

#### 12.4.5 Fractal dimensions from area-perimeter relation

The area of regular objects such as squares or circles scale as the square of the perimeter  $A \sim P^2$ , while the area of irregular fractal objects scale more generally as  $A \sim P^\sigma$ , where  $\sigma = D_A/D_P$  with  $D_A$  and  $D_P$  being the fractal dimension of the area and the perimeter, respectively [561, 593]. These fractal dimensions can be easily computed from the area-perimeter scaling exponent,  $\sigma$ , as  $D_P = (2 + \sigma)/2$  and  $D_A = (2 + \sigma)/2$  [593].

#### 12.4.6 Box-Counting fractal dimension

We computed the Box-Counting fractal dimension of all mapped areas following a box-counting algorithm [561]. Briefly, the method computes the number of boxes of length  $\varepsilon$ ,  $N(\varepsilon)$ , needed to cover the underlying object. Then, the fractal dimension is simply defined as,

$$D = \lim_{\varepsilon \rightarrow 0} \frac{\ln N(\varepsilon)}{\ln 1/\varepsilon} . \quad (12.1)$$

In practice, the mathematical limit  $\varepsilon \rightarrow 0$  is unreachable, and the fractal dimension is computed from the slope obtained in the plot of  $\ln N(\varepsilon)$  versus  $\ln(1/\varepsilon)$ . To efficiently compute the number of overlapping boxes we used the Sort-Tile-Recursive algorithm [587] implemented in Python Shapely library [588]. The implementation of the algorithm can be found in [589].

### 12.4.7 Compactness and elongation index

The compactness measurement is defined as the isoperimetric quotient,

$$C = \frac{4\pi A}{P^2} , \quad (12.2)$$

where  $A$  and  $P$  are the area and perimeter of the object under study, respectively.

The elongation index is defined as the Flaherty & Crumplin (1992) length-width measure, stated as measure LW\_7 in [594] and implemented in PySAL Python's library [595].



# Discussion

<b>13</b>	<b>Main contributions</b> .....	<b>265</b>
13.1	Marine epidemiology .....	265
13.2	Vector-borne plant diseases .....	267
13.3	Disease biogeography .....	268
13.4	Ecological monitoring & analysis .....	270
<b>14</b>	<b>General discussion</b> .....	<b>273</b>



## 13. Main contributions

In this thesis, we have presented original contributions to Ecology and Conservation Biology through the lens of complex systems, specifically focusing on the development of theoretical models and computational methods to address key challenges in the field. We have contributed to the advancement of marine epidemiology, vector-borne plant diseases, disease biogeography, and ecological monitoring, providing new insights into the dynamics of emergent diseases and the conservation of critical coastal ecosystems.

Before presenting a general discussion and outlook, we summarize the main contributions of each part of this thesis. Here we briefly describe the main results and their implications, and we discuss the remaining questions and challenges that remain unexplored.

### 13.1 Marine epidemiology

In [Chapter 3](#), we have developed the Susceptible-Infected-Recovered-Parasite (SIRP) model, a compartmental model that describes the transmission dynamics of diseases affecting sessile marine organisms [215]. The model has been applied to study the Mass Mortality Event (MME) of the fan mussel *Pinna nobilis* in the Mediterranean Sea. Empirical data on the spread of the disease in controlled water tanks were used to estimate the parameters of the model and the basic reproduction number, showing that the model is able to reproduce the observed disease dynamics. Interestingly, we observed that if the parasites die fast enough compared to infected hosts, the dynamics of the disease can be described by a simple SIR model (timescale approximation). This means that, in this limit, a parasite-produced disease in immobile host populations can be described by a model in which there is a direct interaction between mobile hosts. Thus, the parasites “connect” the hosts in a way that is formally equivalent to a direct interaction. In addition, we found an indication that the transmission of the disease follows an Arrhenius-like temperature dependence, which allowed us to discern between mass action law and frequency-dependent transmission.

The natural extension of the SIRP model is the development of a spatially explicit version. In [Chapter 4](#), we have developed an Individual-Based Model (IBM) to study the spatial spread of the disease in a realistic marine environment, showing several interesting results. First, we have shown that the non-spatial SIRP model describes the spatially extended system when the dispersal rate of

the parasite is fast compared to its death and absorption rate (by susceptible hosts) [596]. Second, counterintuitively, we have shown that the timescale approximation also works in the spatial model in the low dispersal limit (i.e., when the non-spatial model is not valid). Third, we have derived an expression for the threshold for disease invasion in the spatial setting. It consists of the non-spatial threshold,  $R_0$ , times a spatial factor that depends on the lifetime and mobility of the parasite. Finally, we have found that the spreading speed of the infected population and the time to extinction scale with the dispersal rate of the parasite.

The development and detailed mathematical analysis of these models has provided new insights into the dynamics of emergent diseases in marine ecosystems and has the potential to inform conservation and management efforts. For instance, the fact that the dynamics of the disease can be described by a simple SIR model in the limit of fast parasite death or that the spatial spread of the disease can be described by the non-spatial model in the limit of fast parasite dispersal is a result that can be generalized to other diseases affecting sessile marine organisms. This can be crucial when facing emergent diseases in marine ecosystems, as it provides an easier way to predict the dynamics of the disease with limited data. Similarly, if the dispersal rate of the parasite is known, the closed mathematical form obtained for the threshold for disease invasion in the spatial model provides a way to estimate the death rate of the parasite from empirical data on the spread of the disease. The results obtained can be compared with the death rate of the parasite in laboratory conditions. If both measures are inconsistent, this would suggest that other transmission mechanisms are at play (e.g., an intermediate vector). In addition, the analysis of the spreading speed of the infected population and the time to extinction as a function of the parasite dispersal rate can be used to assess the impact of the disease on the host population as a function of the ocean currents. This can be crucial to inform conservation and management efforts.

Our research on parasite-produced diseases in marine ecosystems has also provided new questions and challenges that remain unexplored. The MME of the fan mussel *Pinna nobilis* in the Mediterranean Sea has occurred in the whole Mediterranean basin, traveling from west to east, in a relatively short period of time. The mechanisms that have allowed the disease to spread over such a large area in such a short period of time are still largely unknown. The SIRP model developed in this thesis can be used to study the spread of the disease in the whole Mediterranean basin and to identify the key factors that have allowed the disease to spread so fast. For this purpose, metapopulation models that integrate a network of marine currents can be used to model the spread of the disease in the whole Mediterranean basin. Another interesting

question that remains unexplored is the fact that a PDE model of the SIRP model does not show the same results as the IBM model with respect to the threshold for disease invasion. This is an interesting area of research that has not been explored in this thesis. The lack of stable funding and the difficulty in obtaining quantitative experimental data about the spread of the disease in the field, which is crucial to estimate the parameters of the model, have prevented us from pursuing this line of research. In any case, we believe that these are interesting and important questions that could be addressed in future work.

## 13.2 Vector-borne plant diseases

Moving forward to [Chapter 5](#), another key contribution is the development of a theoretical framework for modeling vector-borne plant diseases where vector populations follow non-periodic seasonal dynamics, such as growing or decaying populations towards a stationary state [289]. This study was motivated by the need to develop more realistic dynamical models for vector-borne plant diseases, such as those caused by *Xylella fastidiosa*. Our main finding was that the threshold for disease invasion in this type of models cannot be determined with traditional methods, such as the Next Generation Matrix (NGM) method. This is because the initial stage of the pandemic, when the pathogen is introduced (formally the disease-free state), is not a fixed point of the system, a necessary condition for these methods to work. Thus, we developed a new method based on the concept of the basic reproduction number, which provides a more accurate estimate of the threshold for disease invasion in these models. In addition, we found that, if the vectors died fast enough compared to the infected hosts, the dynamics of the disease could be described by a simple SIR model.

Once the theoretical framework was developed, we were ready to apply it to the case of *Xylella fastidiosa* diseases. In [Chapter 6](#), we adapted the model to apply it to the case of *Xylella fastidiosa* diseases, in which the vector population shows non-periodic seasonal dynamics due to its complex life cycle [597]. The model was contrasted with empirical data from the two main European outbreaks of *Xylella fastidiosa* diseases: Almond Leaf Scorch Disease (ALSD) in Mallorca and Olive Quick Decline Syndrome (OQDS) in Puglia. The model was able to reproduce the observed dynamics of the disease in both cases and provided valuable insights into new possible control strategies for the disease. However, we found that the cross-transmission rates (from infected hosts to vectors and vice versa) were not identifiable from the data, which only provided information about the dynamics of the disease in the host population and not in the vector population. This is an important conclusion of our work: if we want to study vector-borne plant diseases using more realistic (and complex) models, we need to have data on the vector population.

Again, the development and detailed mathematical analysis of these models has provided new insights into the dynamics of vector-borne plant diseases. The NGM method is widely used to estimate the threshold for disease invasion in epidemiological models, but we have shown that it is not always applicable [289], although a generalization is available in the case of periodically varying vector populations [254]. Indeed, it has already been wrongly used in the literature, which can lead to incorrect conclusions. Our new method provides a correct estimate of the threshold for disease invasion in models where the vector population is non-stationary and non-periodic, a case in which the NGM method is not applicable. We have shown again that under certain conditions the dynamics of the disease can be described by a simple SIR model, which can be crucial to predicting the dynamics of the disease with limited data. Interestingly, these conditions are exactly analogous to those found in the SIRP model developed in Chapter 3 for parasite-produced diseases in marine ecosystems of immobile hosts. This suggests that the dynamics of diseases in immobile host populations, under some conditions and to some extent, can be described by the same model regardless of the specific transmission mechanism (in our case, a parasite with no epidemiological states or a vector that can take the susceptible or infected state). Finally, we have shown that if we want to study vector-borne plant diseases using more realistic and complex models, we need to have data on the states of the vector population in addition to hosts.

Our research in this part has also provided new questions and challenges that remain unexplored. Despite having successfully obtained a method to estimate the threshold for disease invasion in models where the vector population is non-stationary and non-periodic, the values we obtain for  $R_0$  are not “realistic”, being usually much higher than 1. We observed that the way in which the outbreaks were produced was not the same as in the usual epidemic models. Basically, we did not seem to have a second-order phase transition between the disease-free and endemic states. This is an interesting area of research that has not been explored in this thesis.

### 13.3 Disease biogeography

The last contributions to epidemiology come from the development in Chapter 7 of a mechanistic climate-driven epidemiological model to assess the risk of establishment of Pierce’s disease (PD), caused by *Xylella fastidiosa* [280]. We fitted the model to empirical data on the distribution of PD in the United States and showed that the model was able to reproduce the observed distribution of the disease. We then used the model to assess the risk of PD establishment under current climate conditions in all viticulture areas worldwide. The most important result is that, in Europe, the potential distribution of the disease

is currently confined to the Mediterranean basin. In the framework of our mechanistic risk model, this can be understood by looking at the interplay between the modified degree days (MGDDs) and cold degree days (CDDs) [280]. While low temperatures in winter (high CDD) prevent the establishment of the disease in continental areas, the mild temperatures (low CDD) close to coastal areas do not decrease the risk of establishment of the disease.

The next logical step was to use the model to assess the risk of PD establishment under future climate conditions. Despite non-coastal Mediterranean areas in Europe being somehow protected from the establishment of the disease by now, climate change could change this scenario. In [Chapter 8](#), we have used the latest regional climate change projections for Europe to assess the risk of PD establishment under future climate conditions [407]. Our results have shown that there is a critical warming level above which PD could become established in continental Europe. However, the risk of PD establishment and its potential impact is not uniform across Europe, with some regions being more vulnerable than others. Similarly, conclusions are not uniform across different administrative levels. Some countries can experience low overall areas at risk, but these areas can be highly productive, thus translating into a high vineyard area at risk.

These results must depend on the spatial scale of the climate data employed. Up to that point, we have used climate data at a coarse spatial resolution ( $0.1^\circ \sim 100\text{km}^2$ ), which is the standard in the field. However, the use of high-resolution climate data can provide more accurate predictions of the risk of PD establishment, taking into account previously neglected microclimates. In [Chapter 9](#), we have used high-resolution climate data to assess the risk of PD establishment in viticulture areas worldwide [598]. We expected a small correction to our previous results, but we found that the potential distribution of the disease is much larger than previously estimated. The differences appeared mostly in valleys and rivers, which are precisely the areas where vineyards are often located. We then analyzed the change in vineyard locations, showing a significant increase of the vineyard area at risk worldwide.

Our model advances the field of disease biogeography by providing a mechanistic understanding of the role of environmental factors in the establishment of plant diseases and provides a more reliable methodology than traditional methods, such as Species Distribution Models (SDMs). Our framework provides a robust way to map the suitability of the pathosystem components to the risk of disease establishment by considering the epidemiological dynamics of the disease. In addition, our model outcomes naturally provide a measure of the potential impact of the disease as well as zones with high uncertainty in the predictions. With this framework, it is straightforward to assess the impact

of different climate change scenarios on the risk of disease establishment and to identify the most vulnerable regions. We have provided an answer to where the disease could become established in Europe under future climate conditions, which was an open question in the field. In addition, the analysis of the impact of high-resolution climate data on the risk of disease establishment has shown that previous estimates of the potential distribution of the disease were underestimated. All our results are crucial to inform management, prevention, and control efforts of plant diseases, such as *Xylella fastidiosa* diseases.

Further extensions of the model could go in the direction of incorporating other environmental factors that could affect the establishment of the disease, such as the presence of alternative hosts or the presence of other vectors. Similarly, the model could be extended to consider the effect of control measures, such as vector control or host removal, or include expected mobility of vectors and infected hosts by plant-trade networks. Finally, the model could be extended to consider other vector-borne plant diseases than Pierce's disease.

## 13.4 Ecological monitoring & analysis

In the last part of this thesis, we have developed different data-driven methods for ecological monitoring and analysis. In [Chapter 10](#), a deep learning framework based on recurrent neural networks was developed to reconstruct missing data in ocean pH time series, addressing data gaps in monitoring ocean acidification and providing reliable pH trend estimates [599]. This has allowed to derive the decadal trend of pH in a coastal area of the Mediterranean Sea, showing a significant decrease in pH over the last decades.

This framework can be used to fill gaps in monitoring ocean acidification, providing reliable estimates of the pH trend that are crucial to properly assess the impact of ocean acidification on marine ecosystems. The model can also be used to reconstruct long chunks of missing data in past ocean pH time series, but one should proceed with caution as the whole context at play could have changed. Thus, the conclusions derived from the model results should be taken with care. Future extensions of the model could go in the direction of incorporating other environmental factors that could affect the pH of the ocean and developing a generalized approach that could be applied to any region of the world.

In [Chapter 11](#), another deep learning framework based on convolutional neural networks was developed to monitor seagrass meadows using satellite imagery, providing accurate and cost-effective estimates of the extent of seagrass meadows in the Balearic Islands [600]. We have provided a big step forward in the field. Previous studies were only a proof of concept for this methodology, often based on a single image or a small dataset and inappropriate performance

measures or model architectures. Furthermore, these works did not assess the role of image variability and geographic context in the model performance. In simpler words, these studies were not thought to be used in practice. Our model, on the contrary, has been trained and validated on a large dataset of satellite images from the whole Balearic Islands, showing that the model is able to generalize its learning, at least to the whole Balearic Islands.

Our results hold significant importance for the conservation of seagrass meadows, as they provide the basis to monitor these critical ecosystems at a large scale. The model can be used to obtain the most up-to-date distribution of *Posidonia oceanica* in the Balearic Islands and, to some extent, in the whole Mediterranean Sea. This is crucial to assess the impact of global change on seagrass meadows and to inform conservation and management efforts. Several future extensions of the model could be considered: monitor seagrass meadows in other regions of the world, include other types of marine habitats or species, etc. By now, we will retrain our model with other satellite sources to obtain the distribution of *Posidonia oceanica* some decades ago. Both results will be crucial to assessing the impact of global change on seagrass meadows.

Finally, in [Chapter 12](#), we conducted a global analysis of the spatial properties of tropical coral reefs from mapped remote sensing data, revealing universal patterns in coral reef size distribution and geometry. We found that the size distribution of coral reefs follows a power-law distribution and that the geometry of coral reefs can be described by a fractal dimension. Interestingly, we found that both the exponent of the power-law distribution and the fractal dimension of coral reefs are independent of the reefs' location, suggesting that these patterns are universal properties of coral reefs.

Despite previous work having already shown the presence of power-law size distributions and fractality in some selected coral reefs, our work is the first to show that these patterns are indeed universal properties of coral reefs. It was suggested that fractal dimensions could vary with the availability of resources, but we have not found significant variations in the fractal dimension of coral reefs in different regions of the world. Indeed, our work connects both theoretically and empirically the fractal dimension of coral reefs with their power-law size distribution, which indicates that these patterns are related to the same underlying processes, such as the growth and breakage of coral colonies.

Our results are crucial to challenging existing and future models of coral reef formation. The power-law size distribution and fractality of coral reefs are key properties that must be reproduced by any model of coral reef formation. The universal nature of these patterns suggests that they are related to fundamental processes that are common to all coral reefs. Indeed, the macroecological patterns observed for coral reefs are by now unexplained. This is an interesting

area of research that has not been explored in this thesis, clearly deserving further investigation.

## 14. General discussion

Here we present a general discussion and outlook of the work presented in this thesis in a broader context: interdisciplinary science at the intersection of Ecology with complex systems and artificial intelligence.

Our research during these four years has been focused on interdisciplinary (also known as cross-disciplinary) research. And here I mean really **interdisciplinary**, not *multidisciplinary*. Nowadays, the term “interdisciplinary science” is used very often, but in many cases it refers to research that is done by a group of researchers from different disciplines, each one working on their own part of the problem (i.e., multidisciplinary). In Ecology, it is quite typical that ecologists (biologists, etc.) focus on gathering data and performing relatively simple statistical analyses. After that, they seek collaboration by passing the data to a group of physicists (mathematicians, etc.), who will perform more complex analyses or even develop models. And the same happens in the opposite direction: there are plenty of physicists who develop models and then seek collaboration with ecologists to validate the models with real data. Both cases are examples of multidisciplinary research, where the researchers from different disciplines work on their own part of the problem and then collaborate to put the pieces together without really integrating the different disciplines in the research process. In any case, this is a very good way of doing research, and it has been very successful in many cases. However, following complex systems’ philosophy, we believe that the whole is more than the sum of its parts.

Most of our research has been really **interdisciplinary**, or at least we have tried to make it so. The biologists and ecologists that we have worked with have been involved in the whole research process, from the very beginning to the end. They have been involved in the development of the models, in the analysis of the results, and in their interpretation. To me, this is a crucial step in interdisciplinary research, and it is the only way to really integrate different disciplines in the research process. The trade-off between *simple enough* mathematical models and *what is indeed enough* to capture appropriately the complexity of the system is a difficult one, and it can only be solved by working together with the experts in the system that we are studying. Physicists will usually tend to think on spherical cows in a vacuum, while biologists and ecologists will usually tend to think on the complexity of the real world, not believing that a simple model can capture the complexity of the system. So, probably quite often, physicists

will make wrong assumptions, leading to unrealistic conclusions, while biologists will overcomplicate models, making them intractable and being unable to generalize the results. The power of interdisciplinary research based on complex systems' science is that it can bridge this gap, integrating the knowledge of the experts in the system with the knowledge of the expert modelers. This can lead to rather general and realistic conclusions that can be applied to a wide range of systems. In my opinion, this interdisciplinary approach is the way to address the pressing challenges that we face in the 21st century, specifically those related to Ecology and Conservation Biology.

Complex system science has proven to be a very powerful tool to study the complexity of the real world. It has been applied to a wide range of systems, from social systems to ecosystems. The main idea behind complex systems is that the whole is more than the sum of its parts. This means that the **interactions** between the parts of the system are crucial to understanding the system as a whole. Perhaps the best example within this thesis is our contributions to disease biogeography. As explained in [Section 2.2](#), renowned researchers had already mentioned the importance of interactions to understand the biogeography of diseases [94]. Despite these insights, the assessment of the risk of Pierce's disease had been hitherto based on the probability of presence of each of the pathosystem components (i.e., the pathogen and the vector) independently, neglecting their interactions. Indeed, this led to contrasting conclusions about the future risk of the disease, as the bacterium is expected to expand its range due to climate change, while the opposite happens to the vector. A recent study somehow went in the right direction by assessing the risk of the disease by considering the overlap of the distributions of the bacterium and the vector with *enough probability of presence* [95]. Unfortunately, this was done in a very simplistic way, setting an arbitrary threshold to define risk areas and, again, neglecting the host-bacterium-plant interactions. It is clear that a formal framework to integrate the role of interactions into the modeling of disease biogeography was lacking.

Our mindset based on complex systems' science, together with the collaboration with biologists and entomologists experts on this pathosystem, has allowed us to address this issue. We have developed a formal framework to integrate the role of interactions into the modeling of disease biogeography. Of course, this has been specifically applied to the case of Pierce's disease, but the framework is indeed general and can be applied to other diseases and systems. A project left for the future is to apply this framework to other systems and formally develop a general theory of disease biogeography based on our approach.

The universal spatial properties of coral reefs uncovered in Chapter 12 are another example of the intersection of Ecology and Complex Systems. As com-

---

mented in [Chapter 1](#), power-laws and fractality are common features of complex systems, usually indicating that the system is self-organized and adaptive. Our results are a great example of how very general and regular phenomena can be found in ecological systems regardless of the specific species that compose the reef. This somehow resembles the universality of the metabolic theory of ecology, in which the metabolic rate of organisms scales with body mass with an exponent close to  $3/4$ , regardless of the specific species. This is a very general and regular pattern that expands across several orders of magnitude of the body mass. In the case of metabolism, this universality is explained by the fractal-like structure of the vascular system of organisms and the invariant size of the capillaries [44]. In the case of coral reefs, the universality of the spatial properties that we found remains to be explained.

In this thesis, we have developed several mathematical models to study the complexity of ecological systems, mostly for epidemiology. Following with the idea of universality, it is noteworthy how, under some conditions, different diseases can be described by exactly the same mathematical framework that was developed almost a century ago, the SIR model. A disease caused by a parasite that affects filter feeders in a marine environment can be mathematically described in the exact same way that a disease transmitted by an insect vector that affects plants in a terrestrial environment. From the biological point of view, these two diseases are completely different: the hosts, the pathogen, the transmission mechanism, etc. However, from the mathematical point of view, they can be described by the same equations whenever there exists a timescale separation between the lifetimes of hosts and parasites/vectors. This is a very powerful tool, as it allows us to apply the same mathematical framework to a wide range of systems. This is the power of mathematical modeling in Ecology from the perspective of complex systems' science.

Of course, these two diseases can be described by exactly the same equations only under some conditions, namely when the infectious agent (parasite or infected vector) die fast enough compared to infected hosts. In any case, the mathematical framework that we use to describe these diseases in the general case is still the same: compartmental models. This means that, contrary to what we might expect, the differences between the components of each pathosystem are not relevant to the mathematical description of the disease, at least at our scale of description (i.e., the population level). To characterize the dynamics of the disease at this level, we don't need to know the details of the transmission at a more detailed level, if the pathogen is a bacterium or a virus, etc. But what we do need to know is how the different components of the system interact with each other.

Mathematical modeling has been a very powerful tool to study the complex-

ity of ecological systems. However, it is important to keep in mind that models are just models, and they are always an approximation of the real world. The strength of a model is not in its complexity but in its simplicity. A model should be as simple as possible, but not simpler. This means that we should always try to develop the simplest model that can capture the complexity of the system. This is a very difficult task, and it requires a deep understanding of the system under study. And that is why interdisciplinary research is crucial. This is the direction that we have followed in this thesis, and we believe that this is the way to go in the future. Nevertheless, the “complexity” of some challenges is beyond the reach of traditional mathematical models. Some problems are naturally embedded in high-dimensional spaces, depending on many variables and parameters that are difficult to measure and quantify. In these cases, the use of artificial intelligence techniques can be very useful, although it is important to keep in mind that these techniques are not a panacea.

Artificial intelligence (AI) techniques have been used in a wide range of fields, from computer science to social sciences. In Ecology, the use of AI techniques is still quite limited, but it is growing. In this thesis, we have developed machine learning models to reconstruct time series of environmental variables and map benthic habitats from satellite imagery. However, the use of these models to obtain ecological insights

In essence, we have used AI to perform specific tasks, focused on solving a given technical problem rather than **understanding** the problem and its solution. This is a very different approach from the one that we have followed with mathematical modeling. Nowadays, AI techniques are usually used to solve specific problems, and they are usually seen as a “black box” that gives an answer to a given problem. This is a very powerful tool, but it is important to understand what they are meant for.

A clear limitation of AI techniques is that they are really data-hungry. A typical deep learning model requires thousands of examples to learn a given task. This is a clear limitation when we are working with ecological systems, as data are usually scarce and difficult to obtain. Of course, there are ways to overcome this limitation, such as data augmentation, transfer learning, etc. To me, these are just workarounds (which can work pretty well in most cases), as high-quality, high-quantity data is a pivotal requirement for AI techniques to work properly. This is a clear limitation of AI techniques, and it is important to keep it in mind when we are using them.

At any rate, research in AI is growing very fast, and it is likely that in the future we will see more and more applications of AI techniques in Ecology. This is a very exciting field, and it is likely that it will lead to very important discoveries. However, it is important to keep in mind that AI techniques are just

---

tools, and they should be used as such. They are not, by any means, a substitute for a deep understanding of the system that we are studying. However, this could change in the future with the development of explainable AI. This is a very active field of research, and it is likely that in the future we will see AI techniques that are able to explain the reasons behind their predictions. This would be a game changer, as it would allow us to use AI techniques to understand the systems that we are studying, not just to solve specific problems.

Overall, our research has contributed to the fields of Ecology, complex systems' science, and artificial intelligence. We have developed different theoretical and data-driven approaches to advance our understanding of current challenges related to biodiversity loss. By developing mathematical models of disease spreading, we have been able to advance our understanding of the Mass Mortality Event of *Pinna nobilis* and the diseases caused by *Xylella fastidiosa*. We have built on this knowledge to develop a formal framework that integrates the role of pathosystem interactions into the modeling of disease biogeography, which has proven to be crucial to accurately assess the risk of disease at a global scale. Nevertheless, mathematical models are not suited to answer all the questions that we face in Ecology, and we have complemented our research with data-driven approaches. To advance our understanding of ocean acidification and the impact of climate change on *Posidonia oceanica* meadows, we have developed machine learning techniques to reconstruct pH time series and map benthic habitats from satellite imagery. Finally, we have investigated the general properties of coral reefs (which indeed we have found to be universal) by using advanced data analysis techniques. The general patterns unraveled by our research will help develop models to better understand the dynamics of coral reefs and their response to global change.

The future of interdisciplinary research at the intersection of Ecology, complex systems' science, and artificial intelligence is exceedingly promising. The pressing challenges of the 21st century demand a profound understanding of the ecological systems we study, making interdisciplinary research indispensable. Integrating different disciplines throughout the research process is crucial for effectively addressing these challenges. Complex systems' science holds the unique capability to bridge gaps between disciplines, merging the expertise of system specialists with that of model experts. This integration fosters the development of comprehensive and realistic conclusions that can be broadly applied across various systems. As interdisciplinary research continues to evolve, we can anticipate a growing number of applications of complex systems' science and artificial intelligence in Ecology. This thesis aspires to contribute to this dynamic field and aid in tackling the urgent environmental challenges of our time.



# Bibliography

- [1] T. N. Taylor and E. L. Taylor, “The Biology and Evolution of Fossil Plants”, [Geological Magazine](#) **130**, 547–547 (1993) (cited on page 1).
- [2] J. W. Schopf, “Fossil evidence of Archaean life”, [Philosophical Transactions of the Royal Society B: Biological Sciences](#) **361**, 869–885 (2006) (cited on page 1).
- [3] B. J. Cardinale et al., “Biodiversity loss and its impact on humanity”, [Nature](#) **486**, 59–67 (2012) (cited on page 1).
- [4] S. A. Levin, “Self-organization and the Emergence of Complexity in Ecological Systems”, [BioScience](#) **55**, 1075–1079 (2005) (cited on page 1).
- [5] L. Gamfeldt et al., “Multiple functions increase the importance of biodiversity for overall ecosystem functioning”, [Ecology](#) **89**, 1223–1231 (2008) (cited on page 1).
- [6] G. C. Daily, “Introduction: What Are Ecosystem Services?”, in *Nature’s Services: Societal Dependence on Natural Ecosystems*, edited by G. C. Daily (Island Press, Washington, DC, 1997), pages 1–10 (cited on page 1).
- [7] T. Newbold et al., “Global map of the Biodiversity Intactness Index, from Newbold et al. (2016) Science”, [Natural History Museum](#) (2016) (cited on page 2).
- [8] J. B. Hughes et al., “Population Diversity: Its Extent and Extinction”, [Science](#) **278**, 689–692 (1997) (cited on page 2).
- [9] G. Ceballos and P. R. Ehrlich, “Mammal Population Losses and the Extinction Crisis”, [Science](#) **296**, 904–907 (2002) (cited on page 2).
- [10] H. M. Pereira et al., “Scenarios for Global Biodiversity in the 21st Century”, [Science](#) **330**, 1496–1501 (2010) (cited on page 2).
- [11] G. Ceballos et al., “Accelerated modern human-induced species losses: Entering the sixth mass extinction”, [Science Advances](#) **1**, e1400253 (2015) (cited on page 2).
- [12] S. L. Pimm et al., “The biodiversity of species and their rates of extinction, distribution, and protection”, [Science](#) **344**, 1246752 (2014) (cited on page 2).
- [13] K. F. Smith et al., “The role of infectious diseases in biological conservation”, [Animal Conservation](#) **12**, 1–12 (2009) (cited on page 2).

- [14] WWF, *Living Planet Report 2022 - Building a Nature-Positive Society*, edited by R. Almond et al. (WWF, Gland, Switzerland, 2022) (cited on page 2).
- [15] A. D. Barnosky et al., “Has the Earth’s sixth mass extinction already arrived?”, *Nature* **471**, 51–57 (2011) (cited on page 2).
- [16] Millennium Ecosystem Assessment, *Ecosystems and Human Well-being: a Framework Working Group for Assessment Report of the Millennium Ecosystem Assessment* (Island Press, Washington, 2005) (cited on pages 3, 230).
- [17] W. F. Laurance et al., “Averting biodiversity collapse in tropical forest protected areas”, *Nature* **489**, 290–294 (2012) (cited on page 3).
- [18] E. Post et al., “Ecological Consequences of Sea-Ice Decline”, *Science* **341**, 519–524 (2013) (cited on page 3).
- [19] K. J. Kroeker et al., “Impacts of ocean acidification on marine organisms: quantifying sensitivities and interaction with warming”, *Global Change Biology* **19**, 1884–1896 (2013) (cited on pages 3, 55, 212).
- [20] C. D. Thomas et al., “Extinction risk from climate change”, *Nature* **427**, 145–148 (2004) (cited on page 3).
- [21] S. R. Loarie et al., “The velocity of climate change”, *Nature* **462**, 1052–1055 (2009) (cited on pages 3, 183, 187).
- [22] S. L. Pimm, “Climate Disruption and Biodiversity”, *Current Biology* **19**, R595–R601 (2009) (cited on page 3).
- [23] R. Warren et al., “Quantifying the benefit of early climate change mitigation in avoiding biodiversity loss”, *Nature Climate Change* **3**, 678–682 (2013) (cited on page 3).
- [24] R. Warren et al., “The projected effect on insects, vertebrates, and plants of limiting global warming to 1.5°C rather than 2°C”, *Science* **360**, 791–795 (2018) (cited on page 3).
- [25] M. C. Urban, “Accelerating extinction risk from climate change”, *Science* **348**, 571–573 (2015) (cited on page 3).
- [26] D. L. Strayer, “Alien species in fresh waters: ecological effects, interactions with other stressors, and prospects for the future”, *Freshwater Biology* **55**, 152–174 (2010) (cited on page 3).
- [27] P. K. Dayton et al., “Environmental effects of marine fishing”, *Aquatic conservation: marine and freshwater ecosystems* **5**, 205–232 (1995) (cited on page 3).

- [28] F. C. Coleman and S. L. Williams, “Overexploiting marine ecosystem engineers: potential consequences for biodiversity”, *Trends in Ecology & Evolution* **17**, 40–44 (2002) (cited on page 3).
- [29] P. Daszak et al., “Emerging Infectious Diseases of Wildlife– Threats to Biodiversity and Human Health”, *Science* **287**, 443–449 (2000) (cited on pages 3, 88).
- [30] C. Mora et al., “Experimental simulations about the effects of overexploitation and habitat fragmentation on populations facing environmental warming”, *Proceedings of the Royal Society B: Biological Sciences* **274**, 1023–1028 (2007) (cited on page 3).
- [31] O. Hoegh-Guldberg et al., “Coral Reefs Under Rapid Climate Change and Ocean Acidification”, *Science* **318**, 1737–1742 (2007) (cited on pages 3, 53).
- [32] G. Bianconi et al., “Complex systems in the spotlight: next steps after the 2021 Nobel Prize in Physics”, *Journal of Physics: Complexity* **4**, 010201 (2023) (cited on page 4).
- [33] T. Vicsek et al., “Novel Type of Phase Transition in a System of Self-Driven Particles”, *Phys. Rev. Lett.* **75**, 1226–1229 (1995) (cited on page 4).
- [34] A. J. Lotka, *Elements of physical biology* (Williams & Wilkins, 1925) (cited on page 4).
- [35] M. Rietkerk and J. van de Koppel, “Regular pattern formation in real ecosystems”, *Trends in Ecology & Evolution* **23**, 169–175 (2008) (cited on page 4).
- [36] R. M. May, “Biological Populations with Nonoverlapping Generations: Stable Points, Stable Cycles, and Chaos”, *Science* **186**, 645–647 (1974) (cited on pages 4, 7).
- [37] R. M. May, “Simple mathematical models with very complicated dynamics”, *Nature* **261**, 459–467 (1976) (cited on page 4).
- [38] E. Meron, “Vegetation pattern formation: The mechanisms behind the forms”, *Physics Today* **72**, 30–36 (2019) (cited on page 5).
- [39] B. T. Milne, “Motivation and Benefits of Complex Systems Approaches in Ecology”, *Ecosystems* **1**, 449–456 (1998) (cited on page 5).
- [40] C. E. Tarnita et al., “A theoretical foundation for multi-scale regular vegetation patterns”, *Nature* **541**, 398–401 (2017) (cited on page 5).
- [41] A. Heyde et al., “Self-organized biotectonics of termite nests”, *Proceedings of the National Academy of Sciences* **118**, e2006985118 (2021) (cited on page 6).

- [42] J. van de Koppel et al., “Experimental Evidence for Spatial Self-Organization and Its Emergent Effects in Mussel Bed Ecosystems”, *Science* **322**, 739–742 (2008) (cited on page 6).
- [43] R. H. Peters, *The Ecological Implications of Body Size*, Cambridge Studies in Ecology (Cambridge University Press, 1983) (cited on page 6).
- [44] G. B. West et al., “A General Model for the Origin of Allometric Scaling Laws in Biology”, *Science* **276**, 122–126 (1997) (cited on pages 6, 275).
- [45] J. H. Brown et al., “Toward a metabolic theory of ecology”, *Ecology* **85**, 1771–1789 (2004) (cited on pages 6, 83).
- [46] J. H. Brown, *Macroecology*, Macroecology (University of Chicago Press, 1995) (cited on pages 6, 251).
- [47] R. Hoch et al., “Towards a standard for documentation of mathematical models in ecology”, *Ecological Modelling* **113**, 3–12 (1998) (cited on page 7).
- [48] S. A. Levin, editor, *Mathematics and Biology: The interface, challenges and opportunities* (Lawrence Berkeley Lab., CA (USA), 1992) (cited on pages 7, 128).
- [49] J. D. Murray, *Mathematical Biology. I. An Introduction*, 3rd edition (Springer, New York, NY, 2002) (cited on pages 7, 8, 15, 69, 128, 330).
- [50] S. Sarkar et al., “Biodiversity conservation planning tools”, *Annual Review of Environment and Resources* **31**, 123–59 (2006) (cited on pages 7, 128).
- [51] D. Tilman and P. Kareiva, *Spatial ecology: the role of space in population dynamics and interspecific interactions* (Princeton University Press, 1997) (cited on page 8).
- [52] A. M. Turing, “The Chemical Basis of Morphogenesis”, *Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences* **237**, 37–72 (1952) (cited on pages 8, 258).
- [53] J. Bascompte, “Networks in ecology”, *Basic and Applied Ecology* **8**, 485–490 (2007) (cited on page 9).
- [54] D. L. DeAngelis and V. Grimm, “Individual-based models in ecology after four decades”, *F1000prime reports* **6** (2014) (cited on page 9).
- [55] S. Christin et al., “Applications for deep learning in ecology”, *Methods in Ecology and Evolution* **10**, 1632–1644 (2019) (cited on page 9).
- [56] M. A. Tabak et al., “Machine learning to classify animal species in camera trap images: Applications in ecology”, *Methods in Ecology and Evolution* **10**, 585–590 (2019) (cited on page 9).

- [57] N. Kussul et al., “Deep Learning Classification of Land Cover and Crop Types Using Remote Sensing Data”, *IEEE Geoscience and Remote Sensing Letters* **14**, 778–782 (2017) (cited on page 9).
- [58] N. Lang et al., “A high-resolution canopy height model of the Earth”, *Nature Ecology & Evolution* **7**, 1778–1789 (2023) (cited on page 9).
- [59] Z. Wu et al., “Deep learning enables satellite-based monitoring of large populations of terrestrial mammals across heterogeneous landscape”, *Nature Communications* **14**, 3072 (2023) (cited on page 9).
- [60] P. J. Zarco-Tejada et al., “Divergent abiotic spectral pathways unravel pathogen stress signals across species”, *Nature Communications* **12**, 6088 (2021) (cited on page 9).
- [61] D. Bernoulli and D. Chapelle, “Essai d’une nouvelle analyse de la mortalité causée par la petite vérole, et des avantages de l’inoculation pour la prévenir”, (1760) (cited on page 13).
- [62] R. Ross, “An application of the theory of probabilities to the study of a priori pathometry.—Part I”, *Proceedings of the Royal Society of London. Series A, Containing Papers of a Mathematical and Physical Character* **92**, 204–230 (1916) (cited on page 13).
- [63] R. Ross and H. P. Hudson, “An application of the theory of probabilities to the study of a priori pathometry.—Part II”, *Proceedings of the Royal Society of London. Series A, Containing Papers of a Mathematical and Physical Character* **93**, 212–225 (1917) (cited on page 13).
- [64] R. Ross and H. P. Hudson, “An application of the theory of probabilities to the study of a priori pathometry.—Part III”, *Proceedings of the Royal Society of London. Series A, Containing Papers of a Mathematical and Physical Character* **93**, 225–240 (1917) (cited on page 13).
- [65] W. O. Kermack and A. G. McKendrick, “A contribution to the mathematical theory of epidemics”, *Proceedings of the Royal Society of London. Series A, Containing Papers of a Mathematical and Physical Character* **115**, 700–721 (1927) (cited on pages 13, 14, 62, 64, 110).
- [66] G. Macdonald et al., “The epidemiology and control of malaria.”, *The Epidemiology and Control of Malaria*. (1957) (cited on pages 13, 110, 111).
- [67] J. Lehtonen, “The Lambert W function in ecological and evolutionary models”, *Methods in Ecology and Evolution* **7**, 1110–1118 (2016) (cited on page 19).

- [68] O. Diekmann et al., “The construction of next-generation matrices for compartmental epidemic models”, *Journal of The Royal Society Interface* **7**, 873–885 (2010) (cited on pages 23, 69, 88, 93, 111, 114, 133, 143, 337–340, 363).
- [69] F. Brauer et al., “Some models for epidemics of vector-transmitted diseases”, *Infectious Disease Modelling* **1**, 79–87 (2016) (cited on pages 23, 111, 113, 114, 341).
- [70] R. Toral and P. Colet, “Introduction to Master Equations”, in *Stochastic Numerical Methods* (John Wiley & Sons, Ltd, 2014) Chap. 8, pages 235–260 (cited on pages 25, 26).
- [71] D. T. Gillespie, “Exact stochastic simulation of coupled chemical reactions”, *The Journal of Physical Chemistry* **81**, 2340–2361 (1977) (cited on pages 25, 92).
- [72] J. M. McCollum et al., “The sorting direct method for stochastic simulation of biochemical systems with varying reaction execution behavior”, *Computational Biology and Chemistry* **30**, 39–49 (2006) (cited on page 26).
- [73] J. Grinnell, “The Niche-Relationships of the California Thrasher”, *The Auk* **34**, 427–433 (1917) (cited on page 28).
- [74] C. Elton, “The animal community”, *Animal ecology*, 239–256 (1927) (cited on page 28).
- [75] G. E. Hutchinson, “Concluding remarks”, in *Cold Spring Harbor symposia on quantitative biology*, Vol. 22 (Cold Spring Harbor Laboratory Press, 1957), pages 415–427 (cited on page 28).
- [76] J. Soberon and A. T. Peterson, “Interpretation of Models of Fundamental Ecological Niches and Species’ Distributional Areas”, *Biodiversity Informatics* **2**, 10.17161/bi.v2i0.4 (2005) (cited on page 28).
- [77] M. Kearney and W. Porter, “Mechanistic niche modelling: combining physiological and spatial data to predict species’ ranges”, *Ecology Letters* **12**, 334–350 (2009) (cited on pages 30, 153).
- [78] J. Franklin and J. Miller, *Mapping species distributions: Spatial inference and prediction*, English (US) (Cambridge University Press, 2010) (cited on page 30).
- [79] J. Elith et al., “Novel methods improve prediction of species’ distributions from occurrence data”, *Ecography* **29**, 129–151 (2006) (cited on page 30).
- [80] GBIF.Org, *GBIF Occurrence Download*, 2023 (cited on pages 30, 182, 201, 206).

- [81] S. E. Fick and R. J. Hijmans, “WorldClim 2: new 1-km spatial resolution climate surfaces for global land areas”, *International Journal of Climatology* **37**, 4302–4315 (2017) (cited on page 30).
- [82] M. Iturbide et al., “A framework for species distribution modelling with improved pseudo-absence generation”, *Ecological Modelling* **312**, 166–174 (2015) (cited on pages 30, 183).
- [83] C. B. Anderson, “elapid: Species distribution modeling tools for Python”, *Journal of Open Source Software* **8**, 4930 (2023) (cited on page 32).
- [84] D. M. Pigott et al., “Mapping the zoonotic niche of Ebola virus disease in Africa”, *eLife* **3**, edited by P. Jha, e04395 (2014) (cited on page 32).
- [85] A. S. Barro et al., “Redefining the Australian Anthrax Belt: Modeling the Ecological Niche and Predicting the Geographic Distribution of *Bacillus anthracis*”, *PLOS Neglected Tropical Diseases* **10**, 1–16 (2016) (cited on page 32).
- [86] T. O. Alimi et al., “Predicting potential ranges of primary malaria vectors and malaria in northern South America based on projected changes in climate, land cover and human population”, *Parasites & Vectors* **8**, 431 (2015) (cited on page 32).
- [87] D. M. Pigott et al., “Mapping the zoonotic niche of Marburg virus disease in Africa”, *Transactions of The Royal Society of Tropical Medicine and Hygiene* **109**, 366–378 (2015) (cited on page 32).
- [88] C. A. Quiner and Y. Nakazawa, “Ecological niche modeling to determine potential niche of Vaccinia virus: a case only study”, *International Journal of Health Geographics* **16**, 28 (2017) (cited on page 32).
- [89] E. E. Johnson et al., “An Ecological Framework for Modeling the Geography of Disease Transmission”, *Trends in Ecology & Evolution* **34**, 655–668 (2019) (cited on page 33).
- [90] A. M. Samy et al., “Leishmaniasis transmission: distribution and coarse-resolution ecology of two vectors and two parasites in Egypt”, *Revista da Sociedade Brasileira de Medicina Tropical* **47**, 57–62 (2014) (cited on pages 33, 34).
- [91] C. M. Baak-Baak et al., “Ecological niche model for predicting distribution of disease-vector mosquitoes in Yucatán State, México”, *Journal of medical entomology* **54**, 854–861 (2017) (cited on pages 33, 34).
- [92] J. Elith and J. R. Leathwick, “Species Distribution Models: Ecological Explanation and Prediction Across Space and Time”, *Annual Review of Ecology, Evolution, and Systematics* **40**, 677–697 (2009) (cited on page 33).

- [93] K. P. Burnham and D. R. Anderson, *Model selection and multimodel inference: a practical information-theoretic approach* (Springer, 2002) (cited on page 33).
- [94] A. T. Peterson, “Biogeography of diseases: a framework for analysis”, *Naturwissenschaften* **95**, 483–491 (2008) (cited on pages 33, 274).
- [95] S. Yoon and W. Lee, “Assessing potential European areas of Pierce’s disease mediated by insect vectors by using spatial ensemble model”, *Frontiers in Plant Science* **14**, 1209694 (2023) (cited on pages 34, 274).
- [96] R. S. Michalski et al., *Machine learning: An artificial intelligence approach* (Springer Science & Business Media, 2013) (cited on page 35).
- [97] I. El Naqa and M. J. Murphy, *What is machine learning?* (Springer, 2015) (cited on page 35).
- [98] A. L. Samuel, “Some studies in machine learning using the game of checkers”, *IBM Journal of Research and Development* **44**, 206–226 (2000) (cited on page 35).
- [99] M. A. Nielsen, *Neural networks and deep learning*, Vol. 25 (Determination press San Francisco, CA, USA, 2015) (cited on page 35).
- [100] R. O. Duda, P. E. Hart, et al., *Pattern classification* (John Wiley & Sons, 2006) (cited on page 35).
- [101] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction* (MIT press, 2018) (cited on page 35).
- [102] D. Silver et al., “Mastering the game of Go without human knowledge”, *Nature* **550**, 354–359 (2017) (cited on page 35).
- [103] Y. LeCun et al., “Deep learning”, *Nature* **521**, 436–444 (2015) (cited on pages 36, 214).
- [104] I. Goodfellow et al., *Deep Learning* (MIT Press, 2016) (cited on pages 36, 38, 214).
- [105] R. J. Hyndman and G. Athanasopoulos, *Forecasting: principles and practice* (OTexts, 2018) (cited on page 39).
- [106] S. Hochreiter and J. Schmidhuber, “Long Short-Term Memory”, *Neural Computation* **9**, 1735–1780 (1997) (cited on pages 41, 215, 396).
- [107] J. Brownlee, *How to visualize filters and feature maps in convolutional neural networks*, 2019 (cited on page 44).
- [108] J. R. García-March et al., “Can we save a marine species affected by a highly infective, highly lethal, waterborne disease from extinction?”, *Biological Conservation* **243**, 108498 (2020) (cited on pages 47, 79, 80).

- [109] M. Zotou et al., “*Pinna nobilis* in the Greek seas (NE Mediterranean): on the brink of extinction?”, *Mediterranean Marine Science* **21**, 575–591 (2020) (cited on page 47).
- [110] M. Vázquez-Luis et al., “S.O.S. *Pinna nobilis*: A Mass Mortality Event in Western Mediterranean Sea”, *Frontiers in Marine Science* **4**, 220 (2017) (cited on pages 47, 88).
- [111] A. Butler et al., “Ecology of the pteroid bivalves *Pinna bicolor* Gmelin and *Pinna nobilis* L”, *Marine Life* **3**, 37–45 (1993) (cited on page 47).
- [112] P. Prado et al., “*Pinna nobilis* in suboptimal environments are more tolerant to disease but more vulnerable to severe weather phenomena”, *Marine Environmental Research*, 105220 (2020) (cited on page 47).
- [113] I. E. Hendriks et al., “Seagrass Meadows Modify Drag Forces on the Shell of the Fan Mussel *Pinna nobilis*”, *Estuaries and Coasts* **34**, 60–67 (2011) (cited on page 47).
- [114] M. Cabanellas-Reboredo et al., “Tracking a mass mortality outbreak of pen shell *Pinna nobilis* populations: A collaborative effort of scientists and citizens”, *Scientific Reports* **9**, 13355 (2019) (cited on pages 47, 48, 80, 88).
- [115] S. Trigos et al., “Utilization of muddy detritus as organic matter source by the fan mussel *Pinna nobilis*”, *Mediterranean Marine Science* **15**, 667–674 (2014) (cited on page 47).
- [116] S. Katsanevakis, “Transplantation as a conservation action to protect the Mediterranean fan mussel *Pinna nobilis*”, *Marine Ecology Progress Series* **546**, 113–122 (2016) (cited on page 47).
- [117] IUCN, “The Red List of Threatened Species. Version 2019-3”, (2019) (cited on page 48).
- [118] F. Carella et al., “A mycobacterial disease is associated with the silent mass mortality of the pen shell *Pinna nobilis* along the Tyrrhenian coastline of Italy”, *Scientific Reports* **9**, 2725 (2019) (cited on page 48).
- [119] T. Šarić et al., “Epidemiology of Noble Pen Shell (*Pinna nobilis* L. 1758) Mass Mortality Events in Adriatic Sea Is Characterised with Rapid Spreading and Acute Disease Progression”, *Pathogens* **9**, 776 (2020) (cited on page 48).
- [120] F. Scarpa et al., “Multiple Non-Species-Specific Pathogens Possibly Triggered the Mass Mortality in *Pinna nobilis*”, *Life* **10**, 238 (2020) (cited on page 48).

- [121] S. Darriba, “First haplosporidan parasite reported infecting a member of the Superfamily Pinnoidea (*Pinna nobilis*) during a mortality event in Alicante (Spain, Western Mediterranean)”, *Journal of Invertebrate Pathology* **148**, 14–19 (2017) (cited on page 48).
- [122] G. Catanese et al., “*Haplosporidium pinnae* sp. nov., a haplosporidan parasite associated with mass mortalities of the fan mussel, *Pinna nobilis*, in the Western Mediterranean Sea”, *Journal of Invertebrate Pathology* **157**, 9–24 (2018) (cited on pages 48, 79).
- [123] A. Box et al., “Reduced Antioxidant Response of the Fan Mussel *Pinna nobilis* Related to the Presence of *Haplosporidium pinnae*”, *Pathogens* **9**, 932 (2020) (cited on page 48).
- [124] E. M. Burrenson and S. E. Ford, “A review of recent information on the Haplosporidia, with special reference to *Haplosporidium nelsoni* (MSX disease)”, *Aquatic Living Resources* **17**, 499–517 (2004) (cited on page 48).
- [125] I. Arzul and R. B. Carnegie, “New perspective on the haplosporidian parasites of molluscs”, *Journal of Invertebrate Pathology* **131**, 32–42 (2015) (cited on pages 48, 85).
- [126] J. M. Wells et al., “*Xylella fastidiosa* gen. nov., sp. nov: Gram-Negative, Xylem-Limited, Fastidious Plant Bacteria Related to *Xanthomonas* spp.”, *International Journal of Systematic and Evolutionary Microbiology* **37**, 136–143 (1987) (cited on pages 48, 152).
- [127] A. Delbianco et al., “A new resource for research and risk analysis: the updated European food safety authority database of *Xylella* spp. host plant species”, *Phytopathology* **109**, 213–215 (2019) (cited on pages 48, 152).
- [128] R. A. Redak et al., “The biology of xylem fluid-feeding insect vectors of *Xylella fastidiosa* and their relation to disease epidemiology”, *Annual Review of Entomology* **49**, 243–270 (2004) (cited on pages 48, 128, 152–154, 173, 175, 196).
- [129] D. Cornara et al., “EPG combined with micro-CT and video recording reveals new insights on the feeding behavior of *Philaenus spumarius*”, *PLoS One* **13**, 1–20 (2018) (cited on pages 48, 128, 145, 152, 153).
- [130] R. P. P. Almeida and A. H. Purcell, “Biological Traits of *Xylella fastidiosa* Strains from Grapes and Almonds”, *Applied and Environmental Microbiology* **69**, 7447–7452 (2003) (cited on pages 49, 129, 152, 154, 157, 181, 343).
- [131] R. P. P. Almeida and L. Nunney, “How do plant diseases caused by *Xylella fastidiosa* emerge?”, *Plant Disease* **99**, 1457–1467 (2015) (cited on page 49).

- [132] M. Saponari et al., “Identification of DNA sequences related to *Xylella fastidiosa* in oleander, almond and olive trees exhibiting leaf scorch symptoms in Southern Italy”, *Journal of Plant Pathology* **55**, 668 (2013) (cited on pages 49, 128, 152).
- [133] D. Olmo et al., “Landscape Epidemiology of *Xylella fastidiosa* in the Balearic Islands”, *Agronomy* **11**, 473 (2021) (cited on pages 49, 128, 129).
- [134] N. Denancé et al., “Several subspecies and sequence types are associated with the emergence of *Xylella fastidiosa* in natural settings in France”, *Plant Pathology* **66**, 1054–1064 (2017) (cited on page 49).
- [135] E. Marco-Noales et al., “Evidence that *Xylella fastidiosa* is the Causal Agent of Almond Leaf Scorch Disease in Alicante, Mainland Spain (Iberian Peninsula)”, *Plant Disease* **105**, 3349–3352 (2021) (cited on page 49).
- [136] N. Zecharia et al., “*Xylella fastidiosa* Outbreak in Israel: Population Genetics, Host Range, and Temporal and Spatial Distribution Analysis”, *Phytopathology* **112**, 2296–2309 (2022) (cited on pages 49, 205).
- [137] C. Carvalho-Luis et al., “Dispersion of the bacterium *Xylella fastidiosa* in Portugal”, *Journal of Agricultural Science and Technology A* **12**, 35–41 (2022) (cited on page 49).
- [138] *Xylella, ritrovati 6 alberi di mandorlo infetti in agro di Triggiano. Pentassuglia: “Fondamentale il lavoro di monitoraggio sui vettori” - PRESS REGIONE* (cited on page 50).
- [139] D. Cornara et al., “*Philaenus spumarius*: when an old acquaintance becomes a new threat to European agriculture”, *Journal of Pest Science* **91**, 957–972 (2018) (cited on pages 50, 129).
- [140] M. Morente et al., “Distribution and relative abundance of insect vectors of *Xylella fastidiosa* in olive groves of the Iberian Peninsula”, *Insects* **9**, 175 (2018) (cited on pages 50, 129).
- [141] C. R. Weaver and D. R. King, “Meadow spittlebug, *Philaenus leucophthalmus* (L.)”, *Ohio Agricultural Experiment Station* **741**, 1–99 (1954) (cited on page 50).
- [142] C. Lago et al., “Degree-day-based model to predict egg hatching of *Philaenus spumarius* (Hemiptera: Aphrophoridae), the main vector of *Xylella fastidiosa* in Europe”, *Environmental Entomology* **52**, 350–359 (2023) (cited on pages 50, 142, 145).
- [143] J. López-Mercadal et al., “Collection of data and information in Balearic Islands on biology of vectors and potential vectors of *Xylella fastidiosa* (GP/EFSA/ALPHA/017/01)”, *EFSA Supporting Publications* **18**, 6925E (2021) (cited on pages 50, 132, 144).

- [144] M. W. Beck et al., “The identification, conservation, and management of estuarine and marine nurseries for fish and invertebrates”, *Bioscience* **51**, 633–641 (2001) (cited on page 51).
- [145] K. L. Heck Jr et al., “Critical evaluation of the nursery role hypothesis for seagrass meadows”, *Marine Ecology Progress Series* **253**, 123–136 (2003) (cited on page 51).
- [146] T. C. Granata et al., “Flow and particle distributions in a nearshore seagrass meadow before and after a storm”, *Marine Ecology Progress Series* **218**, 95–106 (2001) (cited on page 51).
- [147] E. W. Koch et al., “Fluid Dynamics in Seagrass Ecology—from Molecules to Ecosystems”, in *Seagrasses: Biology, Ecology and Conservation*, edited by A. W. D. Larkum et al. (Springer, Dordrecht, NL, 2006), pages 193–225 (cited on page 51).
- [148] A. R. Bos et al., “Ecosystem engineering by annual intertidal seagrass beds: Sediment accretion and modification”, *Estuarine, Coastal and Shelf Science* **74**, 344–348 (2007) (cited on page 51).
- [149] E. Gacia and C. Duarte, “Sediment Retention by a Mediterranean *Posidonia oceanica* Meadow: The Balance between Deposition and Resuspension”, *Estuarine, Coastal and Shelf Science* **52**, 505–514 (2001) (cited on page 51).
- [150] J. D. Madsen et al., “The interaction between water movement, sediment dynamics and submersed macrophytes”, *Hydrobiologia* **444**, 71–84 (2001) (cited on page 51).
- [151] N. Marbà et al., “Effectiveness of protection of seagrass (*Posidonia oceanica*) populations in Cabrera National Park (Spain)”, *Environmental Conservation* **29**, 509–518 (2002) (cited on page 51).
- [152] T. van der Heide et al., “Positive Feedbacks in Seagrass Ecosystems: Implications for Success in Conservation and Restoration”, *Ecosystems* **10**, 1311–1322 (2007) (cited on page 51).
- [153] M. S. Fonseca and J. A. Cahalan, “A preliminary evaluation of wave attenuation by four species of seagrass”, *Estuarine, Coastal and Shelf Science* **35**, 565–576 (1992) (cited on page 51).
- [154] J. F. Sánchez-González et al., “Wave attenuation due to *Posidonia oceanica* meadows”, *Journal of Hydraulic Research* **49**, 503–514 (2011) (cited on page 51).
- [155] R. C. van de Vijssel et al., “Optimal wave reflection as a mechanism for seagrass self-organization”, *Scientific Reports* **13**, 20278 (2023) (cited on page 51).

- [156] C. M. Duarte and C. L. Chiscano, “Seagrass biomass and production: a reassessment”, *Aquatic Botany* **65**, 159–174 (1999) (cited on page 51).
- [157] E. Mcleod et al., “A blueprint for blue carbon: toward an improved understanding of the role of vegetated coastal habitats in sequestering CO<sub>2</sub>”, *Frontiers in Ecology and the Environment* **9**, 552–560 (2011) (cited on pages 52, 230).
- [158] R. J. Orth et al., “A global crisis for seagrass ecosystems”, *Bioscience* **56**, 987–996 (2006) (cited on page 52).
- [159] M. Waycott et al., “Accelerating loss of seagrasses across the globe threatens coastal ecosystems”, *Proceedings of the National Academy of Sciences of the U.S.A.* **106**, 12377–12381 (2009) (cited on pages 52, 230).
- [160] M. Björk et al., *Managing seagrasses for resilience to climate change* (IUCN, 2008) (cited on page 52).
- [161] M. L. Reaka-Kudla, “The global biodiversity of coral reefs: a comparison with rain forests”, in *Biodiversity II: Understanding and protecting our biological resources*, Vol. 2 (Joseph Henry Press Washington, DC, 1997), page 551 (cited on page 53).
- [162] G. Muller-Parker et al., “Interactions Between Corals and Their Symbiotic Algae”, in *Coral Reefs in the Anthropocene*, edited by C. Birkeland (Springer Netherlands, Dordrecht, 2015), pages 99–116 (cited on page 53).
- [163] E. Couce et al., “Future habitat suitability for coral reef ecosystems under global warming and ocean acidification”, *Global Change Biology* **19**, 3592–3606 (2013) (cited on page 53).
- [164] B. E. Brown, “Coral bleaching: causes and consequences”, *Coral Reefs* **16**, S129–S138 (1997) (cited on page 53).
- [165] C. Wiener and A. Davis, “Exploration Down Under”, *Oceanography* **34**, 58–67 (2021) (cited on pages 53, 250).
- [166] C. Darwin, *The Structure and Distribution of Coral Reefs* (Smith, Elder & Co., London, 1874) (cited on pages 54, 250).
- [167] J. C. Orr et al., “Anthropogenic ocean acidification over the twenty-first century and its impact on calcifying organisms”, *Nature* **437**, 681–686 (2005) (cited on page 54).
- [168] I. E. Hendriks et al., “Biological mechanisms supporting adaptation to ocean acidification in coastal ecosystems”, *Estuarine, Coastal and Shelf Science* **152**, A1–A8 (2015) (cited on pages 55, 212).

- [169] E.U. Copernicus Marine Service Information (CMEMS), “Global Ocean acidification - mean sea water pH time series and trend from Multi-Observations Reprocessing”, *Marine Data Store*, 10.48670/moi-00224 (2024) (cited on page 55).
- [170] J. R. Ward and K. D. Lafferty, “The Elusive Baseline of Marine Disease: Are Diseases in Ocean Ecosystems Increasing?”, *PLoS Biology* 2, edited by Larry Crowder, e120 (2004) (cited on page 62).
- [171] C. A. Burge et al., “Climate change influences on marine infectious diseases: Implications for management and society”, *Annual Review of Marine Science* 6, 249–277 (2014) (cited on page 62).
- [172] K. D. Lafferty et al., “Are Diseases Increasing in the Ocean?”, *Annual Review of Ecology, Evolution, and Systematics* 35, 31–54 (2004) (cited on pages 62, 88).
- [173] K. D. Lafferty et al., “Infectious diseases affect marine fisheries and aquaculture economics”, *Annual Review of Marine Science* 7, 471–496 (2015) (cited on page 62).
- [174] I. L. Pairaud et al., “Impacts of climate change on coastal benthic ecosystems: assessing the current risk of mortality outbreaks associated with thermal stress in NW Mediterranean coastal areas”, *Ocean Dynamics* 64, 103–115 (2014) (cited on page 62).
- [175] D. Harvell et al., “The rising tide of ocean diseases: Unsolved problems and research priorities”, *Frontiers in Ecology and the Environment* 2, 375–382 (2004) (cited on pages 62, 64).
- [176] W. O. Kermack and A. G. McKendrick, “Contributions to the mathematical theory of epidemics-II. The problem of endemicity”, *Proceedings of the Royal Society of London. Series A, Containing Papers of a Mathematical and Physical Character* 138, 55–83 (1932) (cited on page 62).
- [177] W. O. Kermack and A. G. McKendrick, “Contributions to the mathematical theory of epidemics-III. Further studies of the problem of endemicity”, *Bulletin of Mathematical Biology* 53, 89–118 (1933) (cited on page 62).
- [178] R. M. Anderson, “Discussion: The Kermack-McKendrick epidemic threshold theorem”, *Bulletin of Mathematical Biology* 53, 3–32 (1991) (cited on pages 62, 88, 110).
- [179] D. L. Cantrell et al., “Modeling Pathogen Dispersal in Marine Fish and Shellfish”, *Trends in Parasitology* 36, 239–249 (2020) (cited on page 62).
- [180] H. I. McCallum et al., “Does terrestrial epidemiology apply to marine systems?”, *Trends in Ecology & Evolution* 19, 585–591 (2004) (cited on pages 63, 64, 84).

- [181] E. N. Powell and E. E. Hofmann, “Models of marine molluscan diseases: Trends and challenges”, *Journal of Invertebrate Pathology* **131**, 212–225 (2015) (cited on pages 63–65, 84, 90).
- [182] G. Bidegain et al., “Marine infectious disease dynamics and outbreak thresholds: contact transmission, pandemic infection, and the potential role of filter feeders”, *Ecosphere* **7**, e01286 (2016) (cited on pages 63, 64, 68, 71, 79, 84, 330).
- [183] G. Bidegain et al., “Microparasitic disease dynamics in benthic suspension feeders: Infective dose, non-focal hosts, and particle diffusion”, *Ecological Modelling* **328**, 44–61 (2016) (cited on pages 63, 66, 88).
- [184] G. Bidegain et al., “Modeling the transmission of *Perkinsus marinus* in the Eastern oyster *Crassostrea virginica*”, *Fisheries Research* **186**, 82–93 (2017) (cited on pages 63, 88).
- [185] O. Diekmann et al., *Mathematical Tools for Understanding Infectious Disease Dynamics* (Princeton U.P., Princeton, 2013) (cited on page 64).
- [186] R. M. May and R. M. Anderson, “Population biology of infectious diseases: Part II”, *Nature* **280**, 455–461 (1979) (cited on pages 64, 84).
- [187] M. Martcheva, *An Introduction to Mathematical Epidemiology* (Springer, New York, 2015) (cited on pages 66, 113, 133, 362).
- [188] F. Brauer, “Models for the spread of universally fatal diseases”, *Journal of Mathematical Biology* **28**, 451–462 (1990) (cited on page 66).
- [189] M. Castro et al., “The turning point and end of an expanding epidemic cannot be precisely forecast”, *Proceedings of the National Academy of Sciences* **117**, 26190–26196 (2020) (cited on pages 69, 84).
- [190] P. van den Driessche and J. Watmough, “Further Notes on the Basic Reproduction Number”, in *Mathematical Epidemiology*, edited by F. Brauer et al. (Springer Berlin Heidelberg, Berlin, Heidelberg, 2008), pages 159–178 (cited on page 69).
- [191] P. Hine, “The ecology of *Bonamia* and decline of bivalve molluscs”, *New Zealand Journal of Ecology* **20**, 109–116 (1996) (cited on page 79).
- [192] S. Culloty and M. Mulcahy, “*Bonamia ostreae* in the Native Oyster *Ostrea edulis*”, *Marine Environment and Health Series* **29**, 1–36 (2007) (cited on page 79).
- [193] C. Audemard et al., “*Bonamia exitiosa* transmission among, and incidence in, Asian oyster *Crassostrea ariakensis* under warm euhaline conditions”, *Diseases of Aquatic Organisms* **110**, 143–50 (2014) (cited on page 79).

- [194] J. Andrews, “Epizootiology of diseases of oysters (*Crassostrea virginica*), and parasites of associated organisms in eastern North America”, *Helgolander Meeresuntersuchungen* **37**, 149–166 (1984) (cited on page 79).
- [195] H. H. Haskin and J. D. Andrews, “Uncertainties and speculations about the life cycle of the eastern oyster pathogen *Haplosporidium nelsoni* (MSX)”, *Special Publication (American Fisheries Society)* **18** (1988) (cited on page 79).
- [196] E. N. Powell et al., “Modeling the MSX parasite in eastern oyster (*Crassostrea virginica*) populations. III. Regional application and the problem of transmission”, *Journal of Shellfish Research* **18**, 517–537 (1999) (cited on pages 79, 88).
- [197] R. Fletcher, “Newton-Like Methods”, in *Practical Methods of Optimization* (John Wiley & Sons, Ltd, 2013) Chap. 3, pages 44–79 (cited on page 80).
- [198] J. Bezanson et al., “Julia: A fresh approach to numerical computing”, *SIAM Review* **59**, 65–98 (2017) (cited on pages 80, 134, 135, 157, 333).
- [199] C. Rackauckas and Q. Nie, “DifferentialEquations.jl – A Performant and Feature-Rich Ecosystem for Solving Differential Equations in Julia”, *The Journal of Open Research Software* **5**, 15 (2017) (cited on pages 80, 135, 333).
- [200] P. K. Molnár et al., “Thermal Performance Curves and the Metabolic Theory of Ecology—A Practical Guide to Models and Experiments for Parasitologists”, *Journal of Parasitology* **103**, 423–439 (2017) (cited on page 83).
- [201] J. R. Rohr and J. M. Cohen, “Understanding how temperature shifts could impact infectious diseases”, *PLoS Biology* **18**, e3000938 (2020) (cited on page 84).
- [202] M. Vurro et al., “Emerging infectious diseases of crop plants in developing countries: impact on agriculture and socio-economic consequences”, *Food Security* **2**, 113–132 (2010) (cited on page 88).
- [203] F. M. Tomley and M. W. Shirley, “Livestock infectious diseases and zoonoses”, *Philosophical transactions of the Royal Society of London. Series B, Biological sciences* **364**, 2637–2642 (2009) (cited on page 88).
- [204] F. Pernet et al., “Infectious diseases in oyster aquaculture require a new integrated approach”, *Philosophical transactions of the Royal Society of London. Series B, Biological sciences* **371**, 20150213 (2016) (cited on page 88).
- [205] C. Salata et al., “Coronaviruses: a paradigm of new emerging zoonotic diseases”, *Pathogens and Disease* **77**, ftaa006 (2020) (cited on page 88).

- [206] D. M. Morens et al., “The challenge of emerging and re-emerging infectious diseases”, *Nature* **430**, 242–249 (2004) (cited on page 88).
- [207] A. A. Cunningham et al., “One Health, emerging infectious diseases and wildlife: two decades of progress?”, *Philosophical Transactions of the Royal Society B: Biological Sciences* **372**, 20160167 (2017) (cited on page 88).
- [208] A. A. Aguirre and G. M. Tabor, “Global Factors Driving Emerging Infectious Diseases”, *Annals of the New York Academy of Sciences* **1149**, 1–3 (2008) (cited on page 88).
- [209] M. E. Eisenlord et al., “Ochre star mortality during the 2014 wasting disease epizootic: role of population size structure and temperature”, *Philosophical Transactions of the Royal Society B: Biological Sciences* **371**, 20150212 (2016) (cited on page 88).
- [210] K. Jones et al., “A review of fibropapillomatosis in Green turtles (*Chelonia mydas*)”, *The Veterinary Journal* **212**, 48–57 (2016) (cited on page 88).
- [211] X. Guo and S. E. Ford, “Infectious diseases of marine molluscs and host responses as revealed by genomic tools”, *Philosophical Transactions of the Royal Society B: Biological Sciences* **371**, 20150206 (2016) (cited on page 88).
- [212] B. Petton et al., “The Pacific Oyster Mortality Syndrome, a Polymicrobial and Multifactorial Disease: State of Knowledge and Future Directions”, *Frontiers in Immunology* **12**, 52 (2021) (cited on page 88).
- [213] F. Pernet et al., “Determination of risk factors for herpesvirus outbreak in oysters using a broad-scale spatial epidemiology framework”, *Scientific Reports* **8**, 10869 (2018) (cited on page 88).
- [214] J. B. McLaughlin et al., “Outbreak of *Vibrio parahaemolyticus* Gastroenteritis Associated with Alaskan Oysters”, *New England Journal of Medicine* **353**, 1463–1470 (2005) (cited on page 88).
- [215] À. Giménez-Romero et al., “Modelling parasite-produced marine diseases: The case of the mass mortality event of *Pinna nobilis*”, *Ecological Modelling* **459**, 109705 (2021) (cited on pages 88–90, 92–94, 96, 101, 102, 111, 265, 335).
- [216] P. C. Cross et al., “Utility of  $R_0$  as a predictor of disease invasion in structured populations”, *Journal of The Royal Society Interface* **4**, 315–324 (2007) (cited on page 89).
- [217] J. Li et al., “The failure of  $R_0$ ”, *Computational and Mathematical Methods in Medicine* **2011**, 527610 (2011) (cited on page 89).

- [218] S. Riley et al., “Five challenges for spatial epidemic models”, [Epidemics](#) **10**, 68–71 (2015) (cited on page 89).
- [219] C. A. Gilligan and F. van den Bosch, “Epidemiological Models for Invasion and Persistence of Pathogens”, [Annual Review of Phytopathology](#) **46**, 385–418 (2008) (cited on page 89).
- [220] V. Grimm and S. F. Railsback, *Individual-based Modeling and Ecology* (Princeton University Press, Princeton (NJ), 2005) (cited on pages 89, 146).
- [221] B. Breckling, “Individual-Based Modelling Potentials and Limitations”, [The Scientific World Journal](#) **2**, 684985 (2002) (cited on page 89).
- [222] À. Giménez-Romero, “Spatial effects in parasite-produced marine diseases”, [GitHub Repository](#) (2022) (cited on page 92).
- [223] F. Brauer, “Compartmental Models in Epidemiology”, in *Mathematical Epidemiology*, edited by F. Brauer et al. (Springer Berlin Heidelberg, Berlin, Heidelberg, 2008), pages 19–79 (cited on pages 96, 110).
- [224] M. J. Keeling and B. T. Grenfell, “Individual-based Perspectives on  $R_0$ ”, [Journal of Theoretical Biology](#) **203**, 51–61 (2000) (cited on page 96).
- [225] J. A. N. Filipe and M. M. Maule, “Analytical methods for predicting the behaviour of population models with general spatial interactions”, [Mathematical Biosciences](#) **183**, 15–35 (2003) (cited on page 96).
- [226] J. A. N. Filipe and M. M. Maule, “Effects of dispersal mechanisms on spatio-temporal development of epidemics”, [Journal of Theoretical Biology](#) **226**, 125–141 (2004) (cited on page 96).
- [227] Y. F. Suprunenko et al., “Analytical approximation for invasion and endemic thresholds, and the optimal control of epidemics in spatially explicit individual-based models”, [Journal of The Royal Society Interface](#) **18**, 20200966 (2021) (cited on pages 96, 97).
- [228] E. Bertuzzo et al., “On spatially explicit models of cholera epidemics”, [Journal of The Royal Society Interface](#) **7**, 321–333 (2010) (cited on page 99).
- [229] T. S. Athni et al., “The influence of vector-borne disease on human history: socio-ecological mechanisms”, [Ecology Letters](#) **24**, 829–846 (2021) (cited on page 110).
- [230] S. K. Schumacher and J. I. Campbell, “Travel Medicine”, in *Urgent Care Medicine Secrets*, edited by R. P. Olympia et al. (Elsevier, 2018) Chap. 56, pages 352–357 (cited on page 110).
- [231] “Fact sheets of vector-borne diseases”, [World Health Organization](#) (2018) (cited on page 110).

- [232] W. Huang et al., “Bacterial Vector-Borne Plant Diseases: Unanswered Questions and Future Directions”, *Molecular Plant* **13**, 1379–1393 (2020) (cited on page 110).
- [233] C. Bragard et al., “Status and Prospects of Plant Virus Control Through Interference with Vector Transmission”, *Annual Review of Phytopathology* **51**, 177–201 (2013) (cited on page 110).
- [234] K. P. Tumber et al., “Pierce’s disease costs California 104 million per year”, *California Agriculture* **68**, 20–29 (2014) (cited on pages 110, 152, 197).
- [235] K. Schneider et al., “Impact of *Xylella fastidiosa* subspecies *pauca* in European olives”, *Proceedings of the National Academy of Sciences* **117**, 9250–9259 (2020) (cited on pages 110, 153, 154, 173, 178, 186, 197).
- [236] E. P. Rybicki, “A Top Ten list for economically important plant viruses”, *Archives of Virology* **160**, 17–20 (2015) (cited on page 110).
- [237] P. van den Driessche, “Reproduction numbers of infectious disease models”, *Infectious Disease Modelling* **2**, 288–303 (2017) (cited on page 110).
- [238] I. G. Laukó, “Stability of disease free sets in epidemic models”, *Mathematical and Computer Modelling* **43**, 1357–1366 (2006) (cited on page 110).
- [239] J. C. Kamgang and G. Sallet, “Computation of threshold conditions for epidemiological models and global stability of the disease-free equilibrium (DFE)”, *Mathematical Biosciences* **213**, 1–12 (2008) (cited on page 110).
- [240] F. Van den Bosch and M. J. Jeger, “The basic reproduction number of vector-borne plant virus epidemics”, *Virus Research* **241**, 196–202 (2017) (cited on page 111).
- [241] R. Garms et al., “Studies on the reinvasion of the *Onchocerciasis* Control Programme in the Volta River Basin by *Simulium damnosum* sI with emphasis on the south-western areas”, *Tropenmedizin und Parasitologie* **30**, 345–362 (1979) (cited on page 111).
- [242] J. Rocklöv and R. Dubrow, “Climate change: an enduring challenge for vector-borne disease prevention and control”, *Nature Immunology* **21**, 479–483 (2020) (cited on pages 111, 178).
- [243] G. Chowell, “Fitting dynamic models to epidemic outbreaks with quantified uncertainty: A primer for parameter uncertainty, identifiability, and forecasts”, *Infectious Disease Modelling* **2**, 379–398 (2017) (cited on pages 111, 144).

- [244] K. Roosa and G. Chowell, “Assessing parameter identifiability in compartmental dynamic models using a computational approach: application to infectious disease transmission models”, *Theoretical Biology and Medical Modelling* **16**, 1 (2019) (cited on pages 111, 144).
- [245] Y.-H. Kao and M. C. Eisenberg, “Practical unidentifiability of a simple vector-borne disease model: Implications for parameter estimation and intervention assessment”, *Epidemics* **25**, 89–100 (2018) (cited on page 111).
- [246] H.-M. Wei et al., “An epidemic model of a vector-borne disease with direct transmission and time delay”, *Journal of Mathematical Analysis and Applications* **342**, 895–908 (2008) (cited on page 115).
- [247] A. A. Lashari and G. Zaman, “Global dynamics of vector-borne diseases with horizontal transmission in host population”, *Computers & Mathematics with Applications* **61**, 745–754 (2011) (cited on page 115).
- [248] N. Shah and G. Jyoti, “SEIR Model and Simulation for Vector Borne Diseases”, *Applied Mathematics* **4**, 13–17 (2013) (cited on page 115).
- [249] S. Zhao et al., “Modelling the effective reproduction number of vector-borne diseases: the yellow fever outbreak in Luanda, Angola 2015–2016 as an example”, *PeerJ* **8**, e8601–e8601 (2020) (cited on page 115).
- [250] L. Esteva and C. Vargas, “Analysis of a dengue disease transmission model”, *Mathematical Biosciences* **150**, 131–151 (1998) (cited on page 115).
- [251] H. R. Thieme, “Convergence results and a Poincaré-Bendixson trichotomy for asymptotically autonomous differential equations”, *Journal of Mathematical Biology* **30**, 755–763 (1992) (cited on page 115).
- [252] C. Castillo-Chavez and H. R. Thieme, “Asymptotically autonomous epidemic models”, in *Mathematical Population Dynamics: Analysis of Heterogeneity, Vol. I, Theory of Epidemics*, edited by O. Arino et al. (Wuerz publishing (Winnipeg), 1995), pages 33–50 (cited on page 115).
- [253] C. L. Wesley and L. J. Allen, “The basic reproduction number in epidemic models with periodic demographics”, *Journal of Biological Dynamics* **3**, 116–129 (2009) (cited on page 121).
- [254] N. Bacaër and S. Guernaoui, “The epidemic threshold of vector-borne diseases with seasonality”, *Journal of Mathematical Biology* **53**, 421–436 (2006) (cited on pages 121, 268).
- [255] Y. H. Chew et al., “Mathematical Models Light Up Plant Signaling”, *The Plant Cell* **26**, 5–20 (2014) (cited on page 128).

- [256] M. J. Jeger et al., “A model for analysing plant-virus transmission characteristics and epidemic development”, *Mathematical Medicine and Biology: A Journal of the IMA* **15**, 1–18 (1998) (cited on page 128).
- [257] M. J. Jeger et al., “Epidemiology of insect-transmitted plant viruses: modelling disease dynamics and control interventions”, *Physiological Entomology* **29**, 291–304 (2004) (cited on page 128).
- [258] L. V. Madden et al., “A Theoretical Assessment of the Effects of Vector-Virus Transmission Mechanism on Plant Virus Disease Epidemics”, *Phytopathology* **90**, 576–594 (2000) (cited on page 128).
- [259] M. J. Jeger and C. Bragard, “The Epidemiology of *Xylella fastidiosa*; A Perspective on Current Knowledge and Framework to Investigate Plant Host–Vector–Pathogen Interactions”, *Phytopathology* **109**, 200–209 (2019) (cited on pages 128, 129, 154, 173, 178, 196).
- [260] C. Chiyaka et al., “Modeling huanglongbing transmission within a citrus tree”, *Proceedings of the National Academy of Sciences* **109**, 12213–12218 (2012) (cited on page 128).
- [261] D. L. Hopkins and A. H. Purcell, “*Xylella fastidiosa*: Cause of Pierce’s Disease of Grapevine and Other Emergent Diseases”, *Plant Disease* **86**, 1056–1066 (2002) (cited on pages 128, 152, 153, 161, 164, 196).
- [262] M. Vanhove et al., “Genomic Diversity and Recombination among *Xylella fastidiosa* Subspecies”, *Applied and Environmental Microbiology* **85**, e02972–18 (2019) (cited on pages 128, 152, 159, 175).
- [263] M. Saponari et al., “*Xylella fastidiosa* in Olive in Apulia: Where We Stand”, *Phytopathology* **109**, 175–186 (2019) (cited on page 128).
- [264] S. Soubeyrand et al., “Inferring pathogen dynamics from temporal count data: the emergence of *Xylella fastidiosa* in France is probably not recent”, *New Phytologist* **219**, 824–836 (2018) (cited on pages 128, 129).
- [265] E. Moralejo et al., “Phylogenetic inference enables reconstruction of a long-overlooked outbreak of almond leaf scorch disease (*Xylella fastidiosa*) in Europe”, *Communications Biology* **3**, 560 (2020) (cited on pages 128, 133, 134, 152, 153, 205, 343, 356).
- [266] D. Cornara et al., “Spittlebugs as vectors of *Xylella fastidiosa* in olive orchards in Italy”, *Journal of Pest Science* **90**, 521–530 (2017) (cited on pages 128, 132).
- [267] J. López-Mercadal et al., “Mechanical management of weeds drops nymphal density of *Xylella fastidiosa* vectors”, *bioRxiv* (2022) (cited on pages 128, 145).

- [268] E. Moralejo et al., “Insights into the epidemiology of Pierce’s disease in vineyards of Mallorca, Spain”, *Plant Pathology* **68**, 1458–1471 (2019) (cited on pages 128, 152–154, 161, 166, 175, 191, 205, 343, 344, 356).
- [269] N. Bodino et al., “Phenology, seasonal abundance and stage-structure of spittlebug (Hemiptera: Aphrophoridae) populations in olive groves in Italy”, *Scientific reports* **9**, 1–17 (2019) (cited on page 129).
- [270] S. M. Chmiel and M. C. Wilson, “Estimation of the Lower and Upper Developmental Threshold Temperatures and Duration of the Nymphal Stages of the Meadow Spittlebug, *Philaenus spumarius*”, *Environmental Entomology* **8**, 682–685 (1979) (cited on page 129).
- [271] D. Cornara et al., “Natural areas as reservoir of candidate vectors of *Xylella fastidiosa*”, *Bulletin of Insectology* **74**, 173–180 (2021) (cited on page 129).
- [272] J. H. Freitag, “Host range of the Pierce’s disease virus of grapes as determined by insect transmission”, *Phytopathology* **41**, 920–934 (1951) (cited on pages 129, 130).
- [273] A. H. Purcell and A. Finlay, “Evidence for noncirculative transmission of Pierce’s disease bacterium by sharpshooter leafhoppers”, *Phytopathology* **69**, 393–395 (1979) (cited on pages 129, 130).
- [274] R. P. P. Almeida and L. Nunney, “How do plant diseases caused by *Xylella fastidiosa* emerge?”, *Plant Disease* **99**, 1457–1467 (1987) (cited on pages 129, 152).
- [275] S. M. White et al., “Modelling the spread and control of *Xylella fastidiosa* in the early stages of invasion in Apulia, Italy”, *Biological Invasions* **19**, 1825–1837 (2017) (cited on page 129).
- [276] C. Abboud et al., “Dating and localizing an invasion from post-introduction data and a coupled reaction–diffusion–absorption model”, *Journal of Mathematical Biology* **79**, 765–789 (2019) (cited on page 129).
- [277] M. P. Daugherty and R. P. P. Almeida, “Understanding How an Invasive Vector Drives Pierce’s Disease Epidemics: Seasonality and Vine-to-Vine Spread”, *Phytopathology* **109**, 277–285 (2019) (cited on pages 129, 153, 175).
- [278] S. M. White et al., “Estimating the epidemiology of emerging *Xylella fastidiosa* outbreaks in olives”, *Plant Pathology* **69**, 1403–1413 (2020) (cited on pages 129, 134, 157).
- [279] M. Brunetti et al., “A mathematical model for *Xylella fastidiosa* epidemics in the Mediterranean regions. Promoting good agronomic practices for their effective control.”, *Ecological Modelling* **432**, 109204 (2020) (cited on page 129).

- [280] À. Giménez-Romero et al., “Global predictions for the risk of establishment of Pierce’s disease of grapevines”, *Communications Biology* **5**, 1389 (2022) (cited on pages 129, 178, 179, 182, 186, 190, 191, 197, 198, 205, 206, 268, 269, 385).
- [281] N. Bodino et al., “Temporal dynamics of the transmission of *Xylella fastidiosa* subsp. *pauca* by *Philaenus spumarius* to olive plants”, *Entomologia Generalis* **41**, 463–480 (2021) (cited on pages 129, 139).
- [282] V. Cavalieri et al., “Transmission of *Xylella fastidiosa* Subspecies *Pauca* Sequence Type 53 by Different Insect Species”, *Insects* **10**, 324 (2019) (cited on pages 130, 138, 139).
- [283] B. L. Teviotdale and J. H. Connell, “Almond Leaf Scorch”, *ANR University of California*, 8106 (2003) (cited on pages 130, 134).
- [284] J. F. Stevenson et al., “Grapevine Susceptibility to Pierce’s Disease II: Progression of Anatomical Symptoms”, *American Journal of Enology and Viticulture* **55**, 238–245 (2004) (cited on page 130).
- [285] A. Fierro et al., “A lattice model to manage the vector and the infection of the *Xylella fastidiosa* on olive trees”, *Scientific Reports* **9**, 8723 (2019) (cited on pages 130, 134).
- [286] S. Antonatos et al., “Seasonal Appearance, Abundance, and Host Preference of *Philaenus spumarius* and *Neophilaenus campestris* (Hemiptera: Aphrophoridae) in Olive Groves in Greece”, *Environmental Entomology* **50**, 1474–1482 (2021) (cited on page 132).
- [287] D. J. Beal et al., “Seasonal Abundance and Infectivity of *Philaenus spumarius* (Hemiptera: Aphrophoridae), a Vector of *Xylella fastidiosa* in California Vineyards”, *Environmental Entomology* **50**, 467–476 (2021) (cited on pages 132, 152, 176).
- [288] N. Bacaër, “Approximation of the Basic Reproduction Number  $R_0$  for Vector-Borne Diseases with a Periodic Vector Population”, *Bulletin of Mathematical Biology* **69**, 1067–1091 (2007) (cited on page 133).
- [289] À. Giménez-Romero et al., “Vector-borne diseases with nonstationary vector populations: The case of growing and decaying populations”, *Phys. Rev. E* **106**, 054402 (2022) (cited on pages 133, 143, 191, 267, 268, 340, 341, 363).
- [290] P. Virtanen et al., “SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python”, *Nature Methods* **17**, 261–272 (2020) (cited on page 134).
- [291] À. Giménez-Romero, “A compartmental model for *Xylella fastidiosa* diseases”, *GitHub Repository* (2022) (cited on pages 134, 226).

- [292] M. D. Homan and A. Gelman, “The No-U-Turn Sampler: Adaptively Setting Path Lengths in Hamiltonian Monte Carlo”, *J. Mach. Learn. Res.* **15**, 1593–1623 (2014) (cited on page 134).
- [293] H. Ge et al., “Turing: a language for flexible probabilistic inference”, in *International Conference on Artificial Intelligence and Statistics, AISTATS 2018, 9-11 April 2018, Playa Blanca, Lanzarote, Canary Islands, Spain* (2018), pages 1682–1690 (cited on page 134).
- [294] A. Saltelli et al., *Sensitivity Analysis in Practice: A Guide to Assessing Scientific Models* (Halsted Press, USA, 2004) (cited on pages 135, 333).
- [295] I. Sobol, “Global sensitivity indices for nonlinear mathematical models and their Monte Carlo estimates”, *Mathematics and Computers in Simulation* **55**, 271–280 (2001) (cited on pages 135, 333).
- [296] D. Cornara et al., “Transmission of *Xylella fastidiosa* by naturally infected *Philaenus spumarius* (Hemiptera, Aphrophoridae) to different host plants”, *Journal of Applied Entomology* **141**, 80–87 (2017) (cited on pages 138, 139).
- [297] N. J. Cunniffe et al., “Thirteen challenges in modelling plant diseases”, *Epidemics* **10**, 6–10 (2015) (cited on page 146).
- [298] M. J. Jeger et al., “Plant virus epidemiology: Applications and prospects for mathematical modeling and analysis to improve understanding and disease control”, *Plant Disease* **102**, 837–854 (2018) (cited on page 146).
- [299] M. Carvajal-Yepes et al., “A global surveillance system for crop diseases: Global preparedness minimizes the risk to food supplies”, *Science* **364**, 1237–1239 (2019) (cited on page 152).
- [300] H. A. Mooney and E. E. Cleland, “The evolutionary impact of invasive species”, *Proceedings of the National Academy of the USA* **98**, 5446–5451 (2001) (cited on page 152).
- [301] D. Pimentel et al., “Environmental and Economic Costs of Nonindigenous Species in the United States”, *BioScience* **50**, 53–65 (2000) (cited on page 152).
- [302] N. Spence et al., “How the global threat of pests and diseases impacts plants, people, and the planet”, *Plants People Planet* **2**, 5–13 (2020) (cited on page 152).
- [303] S. G. Borkar, *History of Plant Pathology* (WPI Publishing, Boca Raton, FL, USA, 2017) (cited on pages 152, 176).

- [304] M. T. Brewer and M. G. Milgroom, “Phylogeography and population structure of the grape powdery mildew fungus, *Erysiphe necator*, from diverse *Vitis* species”, *BMC Evolutionary Biology* **10**, 268 (2010) (cited on pages 152, 176).
- [305] M. Rouxel et al., “Geographic distribution of cryptic species of *Plasmopara viticola* causing downy mildew on wild and cultivated grape in Eastern North America”, *Phytopathology* **104**, 694–701 (2014) (cited on pages 152, 176).
- [306] J. Tello et al., “Major outbreaks in the nineteenth century shaped grape *Phylloxera* contemporary genetic structure in Europe”, *Scientific Reports* **9**, 17540 (2019) (cited on pages 152, 176).
- [307] M. Gomila et al., “Draft genome resources of two strains of *Xylella fastidiosa* XYL1732/17 and XYL2055/17 isolated from Mallorca vineyards”, *Phytopathology* **109**, 222–224 (2019) (cited on pages 152, 343).
- [308] C. C. Su et al., “Pierce’s disease of Grapevines in Taiwan: Isolation, Cultivation and Pathogenicity of *Xylella fastidiosa*”, *Journal of Phytopathology* **161**, 389–396 (2013) (cited on pages 152, 166, 344).
- [309] M. J. Davis et al., “Pierce’s disease of grapevines: Isolation of the causal bacterium”, *Science* **199**, 75–77 (1978) (cited on page 152).
- [310] R. P. P. Almeida et al., “Addressing the New Global Threat of *Xylella fastidiosa*”, *Phytopathology* **109**, 172–174 (2019) (cited on page 152).
- [311] A. Sicard et al., “*Xylella fastidiosa*: Insights into an Emerging Plant Pathogen”, *Annual Review of Phytopathology* **56**, 181–202 (2018) (cited on page 152).
- [312] L. Nunney et al., “An Experimental Test of the Host-Plant Range of Nonrecombinant Strains of North American *Xylella fastidiosa* subsp. *multiplex*”, *Phytopathology* **109**, 294–300 (2019) (cited on page 152).
- [313] N. Denancé et al., “Several subspecies and sequence types are associated with the emergence of *Xylella fastidiosa* in natural settings in France”, *Plant Pathology* **66**, 1054–1064 (2017) (cited on page 152).
- [314] D. Olmo et al., “First Detection of *Xylella fastidiosa* Infecting Cherry (*Prunus avium*) and *Polygala myrtifolia* Plants, in Mallorca Island, Spain”, *Plant Disease* **101**, 1820–1820 (2017) (cited on page 152).
- [315] R. P. P. Almeida, “*Xylella fastidiosa* Vector Transmission Biology”, in *Vector-Mediated Transmission of Plant Pathogens*, edited by J. K. Brown (APS Publications, St Paul, MN, USA, 2016) Chap. 12, pages 165–173 (cited on page 152).

- [316] L. M. Overall and E. J. Rebeck, “Insect vectors and current management strategies for diseases caused by *Xylella fastidiosa* in the Southern United States”, *Journal of Integrated Pest Management* **8**, 1–12 (2017) (cited on page 152).
- [317] H. H. P. Severin, “Spittle-insect vectors of Pierce’s disease virus: II. Life history and virus transmission”, *Hilgardia* **19**, 357–382 (1950) (cited on page 152).
- [318] D. Cornara et al., “Transmission of *Xylella fastidiosa* to Grapevine by the Meadow Spittlebug”, *Phytopathology* **106**, 1285–1290 (2016) (cited on page 152).
- [319] D. Cornara et al., “An overview on the worldwide vectors of *Xylella fastidiosa*”, *Entomologia Generalis* **39**, 157–181 (2019) (cited on page 153).
- [320] M. Godefroid et al., “Climate tolerances of *Philaenus spumarius* should be considered in risk assessment of disease outbreaks related to *Xylella fastidiosa*”, *Journal of Pest Science* **1**, 1 (2021) (cited on pages 153, 154, 160, 161, 171, 172, 175, 176, 359).
- [321] A. Castillo et al., “Allopatric Plant Pathogen Population Divergence following Disease Emergence”, *Applied and Environmental Microbiology* **87**, <https://doi.org/10.1128/AEM.02095-203> (2021) (cited on page 153).
- [322] A. H. Purcell, “Paradigms: Examples from the Bacterium *Xylella fastidiosa*”, *Annual Review of Phytopathology* **51**, 339–356 (2013) (cited on pages 153, 161, 357).
- [323] H. Feil and A. H. Purcell, “Temperature-Dependent Growth and Survival of *Xylella fastidiosa* in Vitro and in Potted Grapevines”, *Plant Disease* **85**, 1230–1234 (2001) (cited on pages 153, 155, 156, 161–163, 173, 349, 350, 352).
- [324] J. H. Lieth et al., “Modeling cold curing of Pierce’s disease in *Vitis vinifera* ‘Pinot Noir’ and ‘Cabernet Sauvignon’ grapevines in California”, *Phytopathology* **101**, 1492–1500 (2011) (cited on pages 153, 156, 157, 159, 163, 349).
- [325] A. H. Purcell, “Almond Leaf Scorch: Leafhopper and Spittlebug Vectors<sup>12</sup>”, *Journal of Economic Entomology* **73**, 834–838 (1980) (cited on page 153).
- [326] O. Anas et al., “The effect of warming winter temperatures on the severity of Pierce’s disease in the Appalachian mountains and Piedmont of the southeastern United States”, *Plant Health Progress* **9**, 1–17 (2008) (cited on pages 153, 156, 159, 163).

- [327] H. Feil et al., “Effects of Date of Inoculation on the Within-Plant Movement of *Xylella fastidiosa* and Persistence of Pierce’s Disease Within Field Grapevines”, *Phytopathology* **93**, 244–251 (2003) (cited on pages 153, 156, 163).
- [328] B. R. Gruber and M. P. Daugherty, “The biology of xylem fluid-feeding insect vectors of *Xylella fastidiosa* and their relation to disease epidemiology”, *Plant Pathology* **62**, 194–204 (2012) (cited on pages 153, 175).
- [329] C. Bragard et al., “Update of the Scientific Opinion on the risks to plant health posed by *Xylella fastidiosa* in the EU territory”, *EFSA Journal* **17**, 5655 (2019) (cited on pages 153, 154, 156, 171, 174, 191).
- [330] M. Godefroid et al., “*Xylella fastidiosa*: climate suitability of European continent”, *Scientific Reports* **9**, 8844 (2019) (cited on pages 153, 171, 178).
- [331] M. S. Hoddle, “The potential adventive geographic range of glassy-winged sharpshooter, *Homalodisca coagulata* and the grape pathogen *Xylella fastidiosa*: implications for California and other grape growing regions of the world”, *Crop Protection* **23**, 691–699 (2004) (cited on page 153).
- [332] L. Bosso et al., “Shedding light on the effects of climate change on the potential distribution of *Xylella fastidiosa* in the Mediterranean basin”, *Biological Invasions* **23**, 1759–1768 (2016) (cited on page 153).
- [333] D. P. Bebber et al., “Crop pests and pathogens move polewards in a warming world”, *Nature Climate Change* **3**, 985–988 (2013) (cited on pages 153, 196).
- [334] S. M. Coakley et al., “Climate change and plant disease management”, *Annual Review of Phytopathology* **37**, 399–426 (1999) (cited on pages 153, 173).
- [335] H. Scherm and A. H. C. van Bruggen, “Global Warming and Nonlinear Growth: How Important are Changes in Average Temperature?”, *Phytopathology* **84**, 1380–1384 (1994) (cited on pages 153, 173, 174, 178, 196, 204).
- [336] J. E. Truscott and C. A. Gilligan, “Response of a deterministic epidemiological system to a stochastically varying environment”, *Proceedings of the National Academy of the USA* **100**, 9067–9072 (2003) (cited on page 153).
- [337] D. P. Bebber et al., “Modelling coffee leaf rust risk in Colombia with climate reanalysis data”, *Philosophical Transactions of the Royal Society B: Biological Sciences* **371**, 20150458 (2016) (cited on page 154).
- [338] A. Giménez-Romero, “Pierce’s Disease Establishment Risk Dashboard”, *Webpage* (2021) (cited on pages 154, 167, 169, 173, 175, 176, 189, 190).

- [339] G. S. McMaster and W. Wilhelm, “Growing degree-days: one equation, two interpretations”, *Agricultural and Forest Meteorology* **87**, 291–300 (1997) (cited on pages 155, 161).
- [340] W. Yan and L. A. Hunt, “An Equation for Modelling the Temperature Response of Plants using only the Cardinal Temperatures”, *Annals of Botany* **84**, 607–614 (1999) (cited on pages 155, 161, 352).
- [341] R. D. Magarey et al., “A Simple Generic Infection Model for Foliar Fungal Plant Pathogens”, *Phytopathology* **95**, 92–100 (2005) (cited on pages 155, 161).
- [342] A. H. Purcell, “Spatial patterns of Pierce’s disease in the Napa Valley”, *American Journal of Enology and Viticulture* **25**, 162–167 (1974) (cited on page 156).
- [343] A. H. Purcell et al., “Vector preference and inoculation efficiency as components of resistance to Pierce’s disease in European grape cultivars”, *Phytopathology* **71**, 429–435 (1981) (cited on page 156).
- [344] T. M. Therneau, “A Package for Survival Analysis in R”, *R package version 3.7-0* (2022) (cited on page 156).
- [345] J. Muñoz Sabater, “ERA5-Land hourly data from 1950 to present”, *Copernicus Climate Change Service (C3S) Climate Data Store (CDS)* (2019) (cited on pages 157, 197, 206).
- [346] A. Weech, “GRIB.jl”, *GitHub Repository* (2021) (cited on page 157).
- [347] L. M. Overall and E. J. Rebek, “Seasonal Abundance and Natural Inoculativity of Insect Vectors of *Xylella fastidiosa* in Oklahoma Tree Nurseries and Vineyards”, *Journal of Economic Entomology Pest Management*, 1–10 (2015) (cited on page 159).
- [348] D. Hail et al., “Detection and analysis of the bacterium, *Xylella fastidiosa*, in glassy-winged sharpshooter, *Homalodisca vitripennis*, populations in Texas”, *Journal of Insect Science* **10**, 168 (2010) (cited on page 159).
- [349] A. K. Wallingford et al., “Expansion of the Range of Pierce’s Disease in Virginia”, *Plant Health Progress* **8**, 42 (2007) (cited on page 159).
- [350] A. L. Myers et al., “Pierce’s disease of grapevines: Identification of the primary vectors in North Carolina”, *Phytopathology* **97**, 1440–1450 (2007) (cited on page 159).
- [351] R. Albibi et al., “RAPD fingerprinting *Xylella fastidiosa* Pierce’s disease strains isolated from a vineyard in North Florida”, *FEMS Microbiology Letters* **165**, 347–352 (1998) (cited on page 159).

- [352] A. Jiménez-Valverde, “Insights into the area under the receiver operating characteristic curve (AUC) as a discrimination measure in species distribution modelling”, *Global Ecology and Biogeography* **21**, 498–507 (2012) (cited on page 159).
- [353] European Environment Agency (EEA), “CORINE Land Cover”, *Copernicus Land Monitoring Service 2018* (2018) (cited on page 160).
- [354] L. Bütikofer et al., “The problem of scale in predicting biological responses to climate”, *Global Change Biology* **26**, 6657–6666 (2020) (cited on page 161).
- [355] S. Fry, R. Milholland, et al., “Multiplication and translocation of *Xylella fastidiosa* in petioles and stems of grapevine resistant, tolerant, and susceptible to Pierce’s disease.”, *Phytopathology* **80**, 61–65 (1990) (cited on pages 161, 173).
- [356] L. C. Galvez et al., “The Threat of Pierce’s disease to Midwest Wine and Table Grapes”, *APSnet Features* (2010) (cited on page 164).
- [357] C. Caminade et al., “Global risk model for vector-borne transmission of Zika virus reveals the role of El Niño 2015”, *Proceedings of the National Academy of the USA* **114**, 119–124 (2017) (cited on page 173).
- [358] B. Berisha et al., “Isolation of Peirce’s disease bacteria from grapevines in Europe”, *European Journal of Plant Pathology* **104**, 427–433 (1998) (cited on page 174).
- [359] R. Karban and M. Huntzinger, “Decline of meadow spittlebugs, a previously abundant insect, along the California coast”, *Bull Ecol Soc Am* **8**, 1–3 (2018) (cited on page 175).
- [360] H. Feng and M. Zhang, “Global land moisture trends: drier in dry and wetter in wet over land”, *Scientific Reports* **5**, 18018 (2016) (cited on page 175).
- [361] C. D. Harvell et al., “Climate Warming and Disease Risks for Terrestrial and Marine Biota”, *Science* **296**, 2158–2162 (2002) (cited on pages 178, 196).
- [362] M. Delgado-Baquerizo et al., “The proportion of soil-borne pathogens increases with warming at the global scale”, *Nature Climate Change* **10**, 550–554 (2020) (cited on pages 178, 196).
- [363] J. Dudney et al., “Nonlinear shifts in infectious rust disease due to climate change”, *Nature Communications* **12**, 5102 (2021) (cited on pages 178, 196, 204).

- [364] T. M. Chaloner et al., “Plant pathogen infection risk tracks global crop yields under climate change”, *Nature Climate Change* **11**, 710–715 (2021) (cited on page 178).
- [365] B. K. Singh et al., “Climate change impacts on plant pathogens, food security and paths forward”, *Nature Reviews Microbiology* **21**, 640–656 (2023) (cited on page 178).
- [366] M. Bergot et al., “Simulation of potential range expansion of oak disease caused by *Phytophthora cinnamomi* under climate change”, *Global Change Biology* **10**, 1539–1552 (2004) (cited on page 178).
- [367] I. B. Pangga et al., “Pathogen dynamics in a crop canopy and their evolution under changing climate”, *Plant Pathology* **60**, 70–81 (2011) (cited on page 178).
- [368] D. P. Bebber, “Climate change effects on Black Sigatoka disease of banana”, *Philosophical Transactions of the Royal Society B* **374**, 20180269 (2019) (cited on page 178).
- [369] P. Juroszek and A. von Tiedemann, “Linking plant disease models to climate change scenarios to project future risks of crop diseases: a review”, *Journal of Plant Diseases and Protection* **122**, 3–15 (2015) (cited on page 178).
- [370] K. A. Garrett et al., “Complexity in climate-change impacts: an analytical framework for effects mediated by plant disease”, *Plant Pathology* **60**, 15–30 (2011) (cited on pages 178, 196).
- [371] M. Godefroid et al., “Climate tolerances of *Philaenus spumarius* should be considered in risk assessment of disease outbreaks related to *Xylella fastidiosa*”, *Journal of Pest Science* **95**, 855–868 (2022) (cited on pages 178, 182–184, 197).
- [372] L. Bosso et al., “Shedding light on the effects of climate change on the potential distribution of *Xylella fastidiosa* in the Mediterranean basin”, *Biological Invasions* **18**, 1759–1768 (2016) (cited on pages 178, 197).
- [373] M. Godefroid et al., “Forecasting future range shifts of *Xylella fastidiosa* under climate change”, *Plant Pathology* **71**, 1839–1848 (2022) (cited on pages 178, 197, 206).
- [374] D. Jacob et al., “Regional climate downscaling over Europe: perspectives from the EURO-CORDEX community”, *Regional Environmental Change* **20**, 51 (2020) (cited on pages 179, 184).
- [375] R. C. Cornes et al., “An Ensemble Version of the E-OBS Temperature and Precipitation Data Sets”, *Journal of Geophysical Research: Atmospheres* **123**, 9391–9409 (2018) (cited on page 179).

- [376] J. Muñoz-Sabater et al., “ERA5-Land: a state-of-the-art global reanalysis dataset for land applications”, *Earth System Science Data* **13**, 4349–4383 (2021) (cited on pages 179, 198, 206).
- [377] F. Giorgi and W. J. Gutowski, “Regional dynamical downscaling and the CORDEX initiative”, *Annual Review of Environment and Resources* **40**, 467–490 (2015) (cited on page 179).
- [378] K. E. Taylor et al., “An overview of CMIP5 and the experiment design”, *Bulletin of the American Meteorological Society* **93**, 485–498 (2011) (cited on page 180).
- [379] J. Diez-Sierra et al., “Consistency of the regional response to global warming levels from CMIP5 and CORDEX projections”, *Climate Dynamics* **61**, 4047–4060 (2023) (cited on page 180).
- [380] M. Iturbide et al., “Repository supporting the implementation of FAIR principles in the IPCC-WGI Atlas”, *Zenodo* (cited on page 180).
- [381] À. Giménez-Romero, “Pierce’s Disease Global Risk Predictions”, *GitHub Repository* (2022) (cited on page 182).
- [382] S. J. Phillips et al., “Maximum entropy modeling of species geographic distributions”, *Ecological Modelling* **190**, 231–259 (2006) (cited on page 182).
- [383] *What is GBIF?*, <https://www.gbif.org/what-is-gbif> (cited on pages 182, 206).
- [384] J. Garcia Molinos and C. Brown, “VoCC”, *Zenodo* (2019) (cited on page 184).
- [385] J. García Molinos et al., “VoCC: An R package for calculating the velocity of climate change and related climatic metrics”, *Methods in Ecology and Evolution* **10**, 2195–2202 (2019) (cited on page 184).
- [386] C. J. Willmott and J. J. Feddema, “A More Rational Climatic Moisture Index”, *The Professional Geographer* **44**, 84–88 (1992) (cited on page 184).
- [387] S. Candiago et al., “A geospatial inventory of regulatory information for wine protected designations of origin in Europe”, *Scientific Data* **9**, 394 (2022) (cited on page 189).
- [388] R. C. Venette et al., “Pest Risk Maps for Invasive Alien Species: A Roadmap for Improvement”, *BioScience* **60**, 349–362 (2010) (cited on page 191).
- [389] L. Hannah et al., “Climate change, wine, and conservation”, *Proceedings of the National Academy of Sciences* **110**, 6907–6912 (2013) (cited on page 192).
- [390] M. Moriondo et al., “Projected shifts of wine regions in response to climate change”, *Climatic change* **119**, 825–839 (2013) (cited on page 192).

- [391] T. Fellmann et al., “Major challenges of integrating agriculture into climate change mitigation policy frameworks”, *Mitigation and Adaptation Strategies for Global Change* **23**, 451–468 (2018) (cited on page 193).
- [392] K. D. Lafferty, “The ecology of climate change and infectious diseases”, *Ecology* **90**, 888–900 (2009) (cited on page 196).
- [393] D. P. Bebber et al., “The global spread of crop pests and pathogens”, *Global Ecology and Biogeography* **23**, 1398–1407 (2014) (cited on page 196).
- [394] H. N. Fones et al., “Threats to global food security from emerging fungal and oomycete crop pathogens”, *Nature Food* **1**, 332–342 (2020) (cited on page 196).
- [395] J. B. Ristaino et al., “The persistent threat of emerging plant disease pandemics to global food security”, *Proceedings of the National Academy of Sciences* **118**, e2022239118 (2021) (cited on page 196).
- [396] S. Skendžić et al., “The Impact of Climate Change on Agricultural Insect Pests”, *Insects* **12**, 440 (2021) (cited on page 196).
- [397] A. Ortiz-Bobea et al., “Anthropogenic climate change has slowed global agricultural productivity growth”, *Nature Climate Change* **11**, 306–312 (2021) (cited on page 196).
- [398] S. A. Levin, “The Problem of Pattern and Scale in Ecology: The Robert H. MacArthur Award Lecture”, *Ecology* **73**, 1943–1967 (1992) (cited on page 196).
- [399] C. Navarro-Racines et al., “High-resolution and bias-corrected CMIP5 projections for climate change impact assessments”, *Scientific Data* **7**, 7 (2020) (cited on page 196).
- [400] D. M. Christiansen et al., “High-resolution data are necessary to understand the effects of climate on plant population dynamics of a forest herb”, *Ecology* **105**, e4191 (2024) (cited on page 196).
- [401] U. A. Abdulwahab et al., “Choice of climate data affects the performance and interpretation of species distribution models”, *Ecological Modelling* **471**, 110042 (2022) (cited on page 196).
- [402] N. Dubos et al., “Choice of climate data influences predictions for current and future global invasion risks for two *Phelsuma* geckos”, *Biological Invasions* **25**, 2929–2948 (2023) (cited on page 196).
- [403] (EFSA) et al., “Update of the *Xylella* spp. host plant database – systematic literature search up to 31 December 2022”, *EFSA Journal* **21**, e08061 (2023) (cited on page 196).

- [404] S. Lindow, “Money Matters: Fueling Rapid Recent Insight Into *Xylella fastidiosa*—An Important and Expanding Global Pathogen”, *Phytopathology* **109**, 210–212 (2019) (cited on page 197).
- [405] C. Sabelli, “Deadly olive tree pathogen came by road and rail”, *Nature Italy*, [10.1038/d43978-023-00118-4](https://doi.org/10.1038/d43978-023-00118-4) (2023) (cited on page 197).
- [406] D. N. Karger et al., “Climatologies at high resolution for the earth’s land surface areas”, *Scientific Data* **4**, 170122 (2017) (cited on pages 197, 198, 206).
- [407] À. Giménez-Romero et al., “Global warming significantly increases the risk of Pierce’s disease epidemics in European vineyards”, *Scientific Reports* **14**, 9648 (2024) (cited on pages 201, 206, 269).
- [408] A. Menzel et al., “European phenological response to climate change matches the warming pattern”, *Global Change Biology* **12**, 1969–1976 (2006) (cited on page 204).
- [409] D. N. Karger et al., “Climatologies at high resolution for the earth’s land surface areas”, *EnviDat*, [10.16904/envidat.228](https://doi.org/10.16904/envidat.228) (2021) (cited on pages 205, 206).
- [410] C.-C. Su et al., “Pierce’s Disease of Grapevines in Taiwan: Isolation, Cultivation and Pathogenicity of *Xylella fastidiosa*”, *Journal of Phytopathology* **161**, 389–396 (2013) (cited on page 205).
- [411] M. Gomila et al., “Draft genome resources of two strains of *Xylella fastidiosa* XYL1732/17 and XYL2055/17 isolated from Mallorca vineyards”, *Phytopathology* **109**, 222–224 (2019) (cited on page 205).
- [412] T. Loureiro et al., “*Xylella fastidiosa*: A Glimpse of the Portuguese Situation”, *Microbiology Research* **14**, 1568–1588 (2023) (cited on page 205).
- [413] P. Friedlingstein et al., “Global Carbon Budget 2021”, *Earth System Science Data Discussions* **2021**, 1–191 (2021) (cited on page 212).
- [414] K. Caldeira and M. E. Wickett, “Anthropogenic carbon and ocean pH”, *Nature* **425**, 365–365 (2003) (cited on page 212).
- [415] S. C. Doney et al., “Ocean acidification: the other CO<sub>2</sub> problem”, *Annual Review of Marine Science* **1**, 169–192 (2009) (cited on page 212).
- [416] G. E. Nilsson et al., “Near-future carbon dioxide levels alter fish behaviour by interfering with neurotransmitter function”, *Nature Climate Change* **2**, 201–204 (2012) (cited on page 212).
- [417] S. Zunino et al., “Impact of ocean acidification on ecosystem functioning and services in habitat-forming species and marine ecosystems”, *Ecosystems* **24**, 1561–1575 (2021) (cited on page 212).

- [418] F. Giorgi, “Climate change hot-spots”, *Geophysical Research Letters* **33**, L08707 (2006) (cited on page 212).
- [419] J. P. Bethoux et al., “The Mediterranean Sea: a miniature ocean for climatic and environmental studies and a key for the climatic functioning of the North Atlantic”, *Progress in Oceanography* **44**, 131–146 (1999) (cited on page 212).
- [420] C. N. Bianchi and C. Morri, “Marine biodiversity of the Mediterranean Sea: Situation, problems and prospects for future research”, *Marine Pollution Bulletin* **40**, 367–376 (2000) (cited on page 212).
- [421] F. Micheli et al., “Cumulative human impacts on Mediterranean and Black Sea marine ecosystems: assessing current pressures and opportunities”, *PLoS One* **8**, e79889 (2013) (cited on page 212).
- [422] M. Vargas-Yáñez et al., “Warming trends and decadal variability in the Western Mediterranean shelf”, *Global and Planetary Change* **63**, 177–184 (2008) (cited on page 212).
- [423] M. Vargas-Yáñez et al., “Climate change in the Western Mediterranean sea 1900–2008”, *Journal of Marine Systems* **82**, 171–176 (2010) (cited on page 212).
- [424] V. Masson-Delmotte et al., “Climate Change 2021: The Physical Science Basis. Contribution of Working Group I to the Sixth Assessment Report of the Intergovernmental Panel on Climate Change”, in *Sixth Assessment Report of the Intergovernmental Panel on Climate Change* (Cambridge University Press, 2021) (cited on page 212).
- [425] J. García-Lafuente et al., “Hotter and Weaker Mediterranean Outflow as a Response to Basin-Wide Alterations”, *Frontiers in Marine Science* **8** (2021) (cited on page 213).
- [426] M. Álvarez et al., “The CO<sub>2</sub> system in the Mediterranean Sea: A basin wide perspective”, *Ocean Science* **10**, 69–92 (2014) (cited on page 213).
- [427] A. E. R. Hassoun et al., “Acidification of the Mediterranean Sea from anthropogenic carbon penetration”, *Deep-Sea Research Part I: Oceanographic Research Papers* **102**, 1–15 (2015) (cited on page 213).
- [428] J. Palmiéri et al., “Simulated anthropogenic CO<sub>2</sub> storage and acidification of the Mediterranean Sea”, *Biogeosciences* **12**, 781–802 (2015) (cited on pages 213, 221).
- [429] S. Flecha et al., “Trends of pH decrease in the Mediterranean Sea through high frequency observational data: Indication of ocean acidification in the basin”, *Scientific Reports* **5**, 1–8 (2015) (cited on pages 213, 214, 221).

- [430] S. Flecha et al., “Decadal acidification in Atlantic and Mediterranean water masses exchanging at the Strait of Gibraltar”, *Scientific Reports* **9**, 1–11 (2019) (cited on page 213).
- [431] L. Kapsenberg et al., “Coastal ocean acidification and increasing total alkalinity in the northwestern Mediterranean Sea”, *Ocean Science* **13**, 411–426 (2017) (cited on pages 213, 221, 222).
- [432] K. M. Yao et al., “Time variability of the north-western Mediterranean Sea pH over 1995–2011”, *Marine Environmental Research* **116**, 51–60 (2016) (cited on pages 213, 221).
- [433] EEA, “State and pressures of the marine and coastal Mediterranean environment”, *European Environment Agency*, 1–44 (1999) (cited on page 213).
- [434] C. J. Crossland et al., “The coastal zone—a domain of global interactions”, in *Coastal fluxes in the Anthropocene* (Springer, 2005), pages 1–37 (cited on page 213).
- [435] A. V. Borges and N. Gypens, “Carbonate chemistry in the coastal zone responds more strongly to eutrophication than to ocean acidification”, *Limnology and Oceanography* **55**, 346–353 (2010) (cited on page 213).
- [436] J. Carstensen and C. M. Duarte, “Drivers of pH Variability in Coastal Ecosystems”, *Environmental Science and Technology* **53**, 4020–4029 (2019) (cited on pages 213, 214).
- [437] N. R. Bates et al., “A Time-Series View of Changing Surface Ocean Chemistry Due to Ocean Uptake of Anthropogenic CO<sub>2</sub> and Ocean Acidification”, *Oceanography* **27**, 126–141 (2014) (cited on pages 213, 221).
- [438] G. E. Hofmann et al., “High-Frequency Dynamics of Ocean pH: A Multi-Ecosystem Comparison”, *PLOS ONE* **6**, 1–11 (2011) (cited on page 213).
- [439] C. M. Duarte et al., “Is Ocean Acidification an Open-Ocean Syndrome? Understanding Anthropogenic Impacts on Seawater pH”, *Estuaries and Coasts* **36**, 221–236 (2013) (cited on pages 213, 214).
- [440] J. M. Mercado and F. J. L. Gordillo, “Inorganic carbon acquisition in algal communities: are the laboratory data relevant to the natural ecosystems?”, *Photosynthesis Research* **109**, 257 (2011) (cited on page 213).
- [441] D. Krause-Jensen et al., “Macroalgae contribute to nested mosaics of pH variability in a subarctic fjord”, *Biogeosciences* **12**, 4895–4911 (2015) (cited on page 213).
- [442] S. Goffredo and Z. Dubinsky, *The Mediterranean Sea: Its history and present challenges* (Springer Science & Business Media, 2013) (cited on page 213).

- [443] K. Murphy et al., “World distribution, diversity and endemism of aquatic macrophytes”, *Aquatic Botany* **158**, 103127 (2019) (cited on page 213).
- [444] I. E. Hendriks et al., “Photosynthetic activity buffers ocean acidification in seagrass meadows”, *Biogeosciences* **11**, 333–346 (2014) (cited on page 214).
- [445] A. M. Ricart et al., “Coast-wide evidence of low pH amelioration by seagrass ecosystems”, *Global Change Biology* **27**, 2580–2591 (2021) (cited on page 214).
- [446] J. Newton et al., “Global ocean acidification observing network: requirements and governance plan”, *Archimer* (2015) (cited on page 214).
- [447] H. Hewamalage et al., “Recurrent Neural Networks for Time Series Forecasting: Current status and future directions”, *International Journal of Forecasting* **37**, 388–427 (2021) (cited on page 214).
- [448] Y. Huang et al., “Reconstructing coupled time series in climate systems using three kinds of machine-learning methods”, *Earth System Dynamics* **11**, 835–853 (2020) (cited on page 214).
- [449] M. Fourier et al., “A Regional Neural Network Approach to Estimate Water-Column Nutrient Concentrations and Carbonate System Variables in the Mediterranean Sea: CANYON-MED”, *Frontiers in Marine Science* **7**, 620 (2020) (cited on pages 215, 217).
- [450] T. Friedrich and A. Oschlies, “Basin-scale pCO<sub>2</sub> maps estimated from ARGO gfloat data: A model study”, *Journal of Geophysical Research: Oceans* **114**, 1–9 (2009) (cited on page 215).
- [451] H. C. Bittig et al., “An Alternative to Static Climatologies: Robust Estimation of Open Ocean CO<sub>2</sub> Variables and Nutrient Concentrations From T, S, and O<sub>2</sub> Data Using Bayesian Neural Networks”, *Frontiers in Marine Science* **5**, 328 (2018) (cited on pages 215, 220).
- [452] P. Landschützer et al., “A neural network-based estimate of the seasonal to inter-annual variability of the Atlantic Ocean carbon sink”, *Biogeosciences* **10**, 7793–7815 (2013) (cited on page 215).
- [453] D. Broullón et al., “A global monthly climatology of total alkalinity: a neural network approach”, *Earth System Science Data* **11**, 1109–1127 (2019) (cited on pages 215, 220).
- [454] D. Broullón et al., “Weekly reconstruction of pH and total alkalinity in an upwelling-dominated coastal ecosystem through neural networks (A<sub>T</sub>pH – NN): The case of Ría de Vigo (NW Spain) between 1992 and 2019”, *Biogeosciences Discussions*, 1–36 (2021) (cited on pages 215, 217, 220).

- [455] S. Contractor and M. Roughan, “Efficacy of Feedforward and LSTM Neural Networks at Predicting and Gap Filling Coastal Ocean Timeseries: Oxygen, Nutrients, and Temperature”, *Frontiers in Marine Science* **8**, 368 (2021) (cited on pages 215, 220).
- [456] M. P. Seidel et al., “A sensor for in situ indicator-based measurements of seawater pH”, *Marine Chemistry* **109**, 18–28 (2008) (cited on page 220).
- [457] L. Gregor et al., “A comparative assessment of the uncertainties of global surface ocean CO<sub>2</sub> estimates using a machine-learning ensemble (CSIR-ML6 version 2019a)–have we hit the wall?”, *Geoscientific Model Development* **12**, 5113–5136 (2019) (cited on page 220).
- [458] N. Lefèvre et al., “A comparison of multiple regression and neural network techniques for mapping in situ pCO<sub>2</sub> data”, *Tellus B: Chemical and Physical Meteorology* **57**, 375–384 (2005) (cited on page 220).
- [459] X. Li et al., “A Neural Network-Based Analysis of the Seasonal Variability of Surface Total Alkalinity on the East China Sea Shelf”, *Frontiers in Marine Science* **7**, 219 (2020) (cited on page 220).
- [460] R. Sauzède et al., “Estimates of water-column nutrient concentrations and carbonate system parameters in the global ocean: a novel approach based on neural networks”, *Frontiers in Marine Science*, 128 (2017) (cited on page 220).
- [461] A. Velo et al., “Total alkalinity estimation using MLR and neural network techniques”, *Journal of Marine Systems* **111**, 11–18 (2013) (cited on page 220).
- [462] A. Graves and J. Schmidhuber, “Framewise phoneme classification with bidirectional LSTM and other neural network architectures”, *Neural networks* **18**, 602–610 (2005) (cited on page 220).
- [463] K. Lee et al., “Roles of marginal seas in absorbing and storing fossil fuel CO<sub>2</sub>”, *Energy & Environmental Science* **4**, 1133–1146 (2011) (cited on page 221).
- [464] A. Schneider et al., “High anthropogenic carbon content in the eastern Mediterranean”, *Journal of Geophysical Research: Oceans* **115**, C12050 (2010) (cited on page 221).
- [465] F. F. Pérez et al., “Contrasting drivers and trends of ocean acidification in the subarctic Atlantic”, *Scientific Reports* **11**, 1–16 (2021) (cited on page 221).
- [466] G. Cossarini et al., “Spatiotemporal variability of alkalinity in the Mediterranean Sea”, *Biogeosciences* **12**, 1647–1658 (2015) (cited on page 222).

- [467] D. Koopmans et al., “High net primary production of Mediterranean seagrass (*Posidonia oceanica*) meadows determined with aquatic eddy covariance”, *Frontiers in Marine Science* **7**, 118 (2020) (cited on page 222).
- [468] C. Barrón et al., “Organic carbon metabolism and carbonate dynamics in a Mediterranean seagrass (*Posidonia oceanica*), meadow”, *Estuaries and Coasts* **29**, 417–426 (2006) (cited on page 222).
- [469] W. Champenois and A. V. Borges, “Seasonal and interannual variations of community metabolism rates of a *Posidonia oceanica* seagrass meadow”, *Limnology and Oceanography* **57**, 347–361 (2012) (cited on page 222).
- [470] P. Rivaro et al., “Distributions of carbonate properties along the water column in the Mediterranean Sea: Spatial and temporal variations”, *Marine Chemistry* **121**, 236–245 (2010) (cited on page 222).
- [471] A. E. R. Hassoun et al., “Modeling of the total alkalinity and the total inorganic carbon in the Mediterranean Sea”, *Journal of Water Resources and Ocean Science* **4**, 24–32 (2015) (cited on page 222).
- [472] E. Gemayel et al., “Climatological variations of total alkalinity and total dissolved inorganic carbon in the Mediterranean Sea surface waters”, *Earth System Dynamics* **6**, 789–800 (2015) (cited on pages 222, 225).
- [473] A. Schneider et al., “Alkalinity of the Mediterranean Sea”, *Geophysical Research Letters* **34**, L15608 (2007) (cited on page 222).
- [474] C. Millot, “Circulation in the Western Mediterranean Sea”, *Journal of Marine Systems* **20**, 423–442 (1999) (cited on page 222).
- [475] R. Pawlowicz, “M\_Map: A mapping package for MATLAB, version 1.4 m”, *Computer software*, UBC EOAS (2020) (cited on page 223).
- [476] F. Gazeau et al., “Whole-system metabolism and CO<sub>2</sub> fluxes in a Mediterranean Bay dominated by seagrass beds (Palma Bay, NW Mediterranean)”, *Biogeosciences* **2**, 43–60 (2005) (cited on page 223).
- [477] N. Marbà et al., “Effectiveness of protection of seagrass (*Posidonia oceanica*) populations in Cabrera National Park (Spain)”, *Environmental Conservation* **29**, 509–518 (2002) (cited on page 223).
- [478] J. Tintoré and B. Casas Pérez, “Buoy Bahía de Palma Physicochemical parameters of sea water data”, *Balearic Islands Coastal Observing and Forecasting System*, SOCIB (2022) (cited on page 224).
- [479] B. B. Benson and D. Krause, “The concentration and isotopic fractionation of oxygen dissolved in freshwater and seawater in equilibrium with the atmosphere”, *Limnology and Oceanography* **29**, 620–632 (1984) (cited on page 224).

- [480] T. D. Clayton and R. H. Byrne, “Spectrophotometric seawater pH measurements: total hydrogen ion concentration scale calibration of m-cresol purple and at-sea results”, *Deep-Sea Research Part I* **40**, 2115–2129 (1993) (cited on page 224).
- [481] A. G. Dickson et al., *Guide to best practices for ocean CO<sub>2</sub> measurements* (North Pacific Marine Science Organization, 2007) (cited on page 225).
- [482] J. D. Sharp et al., “CO2SYSv3 for MATLAB”, *Zenodo* (2020) (cited on page 225).
- [483] C. Mehrbach et al., “Measurement of the apparent dissociation constants of carbonic acid in seawater at atmospheric pressure”, *Limnology and Oceanography* **18**, 897–907 (1973) (cited on page 225).
- [484] A. Dickson and F. J. Millero, “A comparison of the equilibrium constants for the dissociation of carbonic acid in seawater media”, *Deep Sea Research Part A. Oceanographic Research Papers* **34**, 1733–1743 (1987) (cited on page 225).
- [485] A. G. Dickson, “Standard potential of the reaction:  $\text{AgCl (s)} + 12\text{H}_2 \text{(g)} = \text{Ag (s)} + \text{HCl (aq)}$ , and the standard acidity constant of the ion  $\text{HSO}_4^-$  in synthetic sea water from 273.15 to 318.15 K”, *The Journal of Chemical Thermodynamics* **22**, 113–127 (1990) (cited on page 225).
- [486] R. J. Woosley, “Evaluation of the temperature dependence of dissociation constants for the marine carbon system using pH and certified reference materials”, *Marine Chemistry* **229**, 103914 (2021) (cited on page 225).
- [487] R. Weiss and B. Price, “Nitrous oxide solubility in water and seawater”, *Marine Chemistry* **8**, 347–359 (1980) (cited on page 225).
- [488] E. Dlugokencky et al., “Atmospheric Nitrous Oxide Dry Air Mole Fractions from the NOAA GML Carbon Cycle Cooperative Global Air Sampling Network, 1997-2020”, *NOAA* (2021) (cited on page 225).
- [489] R. B. Cleveland et al., “STL: A seasonal-trend decomposition procedure based on loess”, *J. Off. Stat* **6**, 3–73 (1990) (cited on page 225).
- [490] À. Giménez-Romero, “Coastal pH variability reconstructed through neural networks: the coastal Balearic Sea case study”, *GitHub Repository* (2021) (cited on page 227).
- [491] M. F. Kallesøe et al., *Linking Coastal Ecosystems and Human Well-Being: Learning from conceptual frameworks and empirical results* (Colombo: Ecosystems and Livelihoods Group, Asia, IUCN, 2008) (cited on page 230).

- [492] C. M. Duarte et al., “The role of coastal plant communities for climate change mitigation and adaptation”, *Nature Climate Change* **3**, 961–968 (2013) (cited on page 230).
- [493] P. Macreadie et al., “Quantifying and modelling the carbon sequestration capacity of seagrass meadows – A critical assessment”, *Marine Pollution Bulletin* **83**, 430–439 (2014) (cited on page 230).
- [494] G. Jorda et al., “Ocean warming compresses the three-dimensional habitat of marine life”, *Nature Ecology and Evolution* **4**, 109–114 (2020) (cited on page 230).
- [495] N. J. Waltham et al., “UN Decade on Ecosystem Restoration 2021–2030—What Chance for Success in Restoring Coastal Ecosystems?”, *Frontiers in Marine Science* **7**, 71 (2022) (cited on page 230).
- [496] “United Nations Conference on Sustainable Development, Rio+20”, *Convention on Biological Diversity* (2012) (cited on page 230).
- [497] P. J. Mumby and A. J. Edwards, “Mapping marine environments with IKONOS imagery: enhanced spatial resolution can deliver greater thematic accuracy”, *Remote Sensing of Environment* **82**, 248–257 (2002) (cited on page 230).
- [498] D. Mishra et al., “Benthic Habitat Mapping in Tropical Marine Environments Using QuickBird Multispectral Data”, *Photogrammetric Engineering & Remote Sensing* **72**, 1037–1048 (2006) (cited on page 230).
- [499] A. Le Quilleuc et al., “Very High-Resolution Satellite-Derived Bathymetry and Habitat Mapping Using Pleiades-1 and ICESat-2”, *Remote Sensing* **14**, 133 (2022) (cited on page 230).
- [500] Allen Coral Atlas, “Imagery, maps and monitoring of the world’s tropical coral reefs”, *Zenodo* (2022) (cited on pages 230, 250, 251, 259).
- [501] C. Zhang et al., “Object-based benthic habitat mapping in the Florida Keys from hyperspectral imagery”, *Estuarine, Coastal and Shelf Science* **134**, 88–97 (2013) (cited on page 231).
- [502] J. J. Senecal et al., “Efficient Convolutional Neural Networks for Multi-Spectral Image Classification”, in *2019 International Joint Conference on Neural Networks (IJCNN)* (2019), pages 1–8 (cited on page 231).
- [503] P. Wicaksono et al., “Benthic Habitat Mapping Model and Cross Validation Using Machine-Learning Classification Algorithms”, *Remote Sensing* **11**, 1279 (2019) (cited on page 231).

- [504] P. Gudžius et al., “Deep learning-based object recognition in multispectral satellite imagery for real-time applications”, *Machine Vision and Applications* **32**, 98 (2021) (cited on page 231).
- [505] S. Chand and B. Bollard, “Detecting the Spatial Variability of Seagrass Meadows and Their Consequences on Associated Macrofauna Benthic Activity Using Novel Drone Technology”, *Remote Sensing* **14**, 160 (2022) (cited on page 231).
- [506] E.-i. Jeon et al., “Semantic segmentation of seagrass habitat from drone imagery based on deep learning: A comparative study”, *Ecological Informatics* **66**, 101430 (2021) (cited on page 231).
- [507] D. Traganos et al., “Towards Global-Scale Seagrass Mapping and Monitoring Using Sentinel-2 on Google Earth Engine: The Case Study of the Aegean and Ionian Seas”, *Remote Sensing* **10**, 1227 (2018) (cited on page 231).
- [508] A. Ariasari et al., “Random forest classification and regression for seagrass mapping using PlanetScope image in Labuan Bajo, East Nusa Tenggara”, in *Sixth International Symposium on LAPAN-IPB Satellite*, Vol. 11372, edited by Y. Setiawan et al. (International Society for Optics and Photonics, 2019), 113721Q (cited on page 231).
- [509] N. T. Ha et al., “A Comparative Assessment of Ensemble-Based Machine Learning and Maximum Likelihood Methods for Mapping Seagrass Using Sentinel-2 Imagery in Tauranga Harbor, New Zealand”, *Remote Sensing* **12**, 355 (2020) (cited on page 231).
- [510] A. Mederos-Barrera et al., “Seagrass mapping using high resolution multispectral satellite imagery: A comparison of water column correction models”, *International Journal of Applied Earth Observation and Geoinformation* **113**, 102990 (2022) (cited on page 231).
- [511] M. M. Coffey et al., “Performance across WorldView-2 and RapidEye for reproducible seagrass mapping”, *Remote Sensing of Environment* **250**, 112036 (2020) (cited on page 231).
- [512] J. Marcello et al., “Seabed mapping in coastal shallow waters using high resolution multispectral and hyperspectral imagery”, *Remote Sensing* **10**, 1208 (2018) (cited on page 231).
- [513] D. Traganos and P. Reinartz, “Machine learning-based retrieval of benthic reflectance and *Posidonia oceanica* seagrass extent using a semi-analytical inversion of Sentinel-2 satellite data”, *International Journal of Remote Sensing* **39**, 9428–9452 (2018) (cited on page 231).

- [514] D. Poursanidis et al., “On the use of Sentinel-2 for coastal habitat mapping and satellite-derived bathymetry estimation using downscaled coastal aerosol band”, *International Journal of Applied Earth Observation and Geoinformation* **80**, 58–70 (2019) (cited on page 231).
- [515] J. E. Duffy et al., “Toward a coordinated global observing system for seagrasses and marine macroalgae”, *Frontiers in Marine Science* **6**, 317 (2019) (cited on page 231).
- [516] K. A. Islam et al., “Semi-Supervised Adversarial Domain Adaptation for Seagrass Detection in Multispectral Images”, in *2019 IEEE International Conference on Data Mining (ICDM)* (2019), pages 1120–1125 (cited on page 231).
- [517] D. Poursanidis et al., “Mapping coastal marine habitats and delineating the deep limits of the Neptune’s seagrass meadows using very high resolution Earth observation data”, *International Journal of Remote Sensing* **39**, 8670–8687 (2018) (cited on page 231).
- [518] D. Traganos et al., “Towards global-scale seagrass mapping and monitoring using Sentinel-2 on Google Earth Engine: The case study of the Aegean and Ionian Seas”, *Remote Sensing* **10**, 1–14 (2018) (cited on page 231).
- [519] D. Traganos and P. Reinartz, “Mapping Mediterranean seagrasses with Sentinel-2 imagery”, *Marine Pollution Bulletin* **134**, 197–209 (2018) (cited on page 231).
- [520] T. Bakirman and M. U. Gumusay, “Assessment of machine learning methods for seagrass classification in the mediterranean”, *Baltic Journal of Modern Computing* **8**, 315–326 (2020) (cited on page 231).
- [521] A. Kellaris et al., “Using low-cost drones to monitor heterogeneous submerged seaweed habitats: A case study in the Azores”, *Aquatic Conservation: Marine and Freshwater Ecosystems* **29**, 1909–1922 (2019) (cited on page 231).
- [522] M. Chowdhury et al., “AI-driven remote sensing enhances Mediterranean seagrass monitoring and conservation to combat climate change and anthropogenic impacts”, *Scientific Reports* **14**, 8360 (2024) (cited on page 231).
- [523] M. D. Zeiler and R. Fergus, “Visualizing and Understanding Convolutional Networks”, in *Computer Vision – ECCV 2014*, edited by D. Fleet et al. (2014), pages 818–833 (cited on page 231).
- [524] F. Milletari et al., “V-Net: Fully Convolutional Neural Networks for Volumetric Medical Image Segmentation”, in *2016 Fourth International Conference on 3D Vision (3DV)* (2016), pages 565–571 (cited on page 231).

- [525] Planet Team, *Planet Application Program Interface: In Space for Life on Earth* (San Francisco, CA, 2017) (cited on pages 232, 242).
- [526] J. Dai et al., “Instance-aware semantic segmentation via multi-task network cascades”, in [Proceedings of the IEEE conference on computer vision and pattern recognition](#) (2016), pages 3150–3158 (cited on page 234).
- [527] À. Giménez-Romero, “CAMELE: Consensus for Automated Marine Ecosystem Labelling and Evaluation”, [Zenodo \(2024\)](#) (cited on page 237).
- [528] A. Giménez-Romero, “CAMELE dashboard”, [Webpage \(2024\)](#) (cited on page 239).
- [529] Government of the Balearic Islands, *Atlas Posidonia, 2000-2019* (cited on page 242).
- [530] L. del Valle Villalonga et al., “*Posidonia oceanica* Cartography and Evolution of the Balearic Sea (Western Mediterranean)”, [Remote Sensing](#) **15**, 5748 (2023) (cited on page 242).
- [531] B. Martín Míguez et al., “The European Marine Observation and Data Network (EMODnet): Visions and Roles of the Gateway to Marine Data in Europe”, [Frontiers In Marine Science](#) **6**, 24 (2019) (cited on page 243).
- [532] O. Ronneberger et al., “U-Net: Convolutional Networks for Biomedical Image Segmentation”, in [Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015](#), edited by N. Navab et al. (2015), pages 234–241 (cited on pages 244, 404).
- [533] A. Chaurasia and E. Culurciello, “LinkNet: Exploiting encoder representations for efficient semantic segmentation”, in [2017 IEEE Visual Communications and Image Processing \(VCIP\)](#) (2017) (cited on pages 244, 404).
- [534] T.-Y. Lin et al., “Feature Pyramid Networks for Object Detection”, in [2017 IEEE Conference on Computer Vision and Pattern Recognition \(CVPR\)](#) (2017), pages 936–944 (cited on pages 244, 404).
- [535] H. Zhao et al., “Pyramid Scene Parsing Network”, in [2017 IEEE Conference on Computer Vision and Pattern Recognition \(CVPR\)](#) (2017), pages 6230–6239 (cited on pages 244, 405).
- [536] P. Iakubovskii, “Segmentation Models”, [GitHub Repository \(2019\)](#) (cited on page 244).
- [537] M. A. Rahman and Y. Wang, “Optimizing Intersection-Over-Union in Deep Neural Networks for Image Segmentation”, in [Advances in Visual Computing](#), edited by G. Bebis et al. (2016), pages 234–244 (cited on page 245).

- [538] T. Sorensen, “A method of establishing groups of equal amplitude in plant sociology based on similarity of species content and its application to analyses of the vegetation on Danish commons”, *Biologiske skrifter* **5**, 1–34 (1948) (cited on page 245).
- [539] L. R. Dice, “Measures of the amount of ecologic association between species”, *Ecology* **26**, 297–302 (1945) (cited on page 245).
- [540] F. Moberg and C. Folke, “Ecological goods and services of coral reef ecosystems”, *Ecological Economics* **29**, 215–233 (1999) (cited on page 250).
- [541] A. W. Droxler and S. J. Jorry, “The Origin of Modern Atolls: Challenging Darwin’s Deeply Ingrained Theory”, *Annual Review of Marine Science* **13**, 537–573 (2021) (cited on page 250).
- [542] T. P. Scoffin and J. E. Dixon, “The distribution and structure of coral reefs: one hundred years since Darwin”, *Biological Journal of the Linnean Society* **20**, 11–38 (1983) (cited on page 250).
- [543] S. J. Purkis et al., “The Statistics of Natural Shapes in Modern Coral Reef Landscapes”, *The Journal of Geology* **115**, 493–508 (2007) (cited on page 250).
- [544] R. H. Bradbury and R. E. Reichelt, “Fractal dimension of a coral reef at ecological scales”, *Marine Ecology Progress Series* **10**, 169–171 (1983) (cited on page 250).
- [545] D. G. Zawada and J. C. Brock, “A Multiscale Analysis of Coral Reef Topographic Complexity Using Lidar-Derived Bathymetry”, *Journal of Coastal Research* **2009**, 6–15 (2009) (cited on page 250).
- [546] L. Alvarez-Filip et al., “Flattening of Caribbean coral reefs: region-wide declines in architectural complexity”, *Proceedings of the Royal Society B: Biological Sciences* **276**, 3019–3025 (2009) (cited on page 250).
- [547] Y.-M. Bozec et al., “The dynamics of architectural complexity on coral reefs under climate change”, *Global Change Biology* **21**, 223–235 (2015) (cited on page 250).
- [548] D. Sous et al., “On the small-scale fractal geometrical structure of a living coral reef barrier”, *Earth Surface Processes and Landforms* **45**, 3042–3054 (2020) (cited on pages 250, 258).
- [549] T. D. Eddy et al., “Global Decline in Capacity of Coral Reefs to Provide Ecosystem Services”, *One Earth* **4**, 1278–1285 (2021) (cited on page 250).

- [550] N. L. Bindoff et al., “Changing Ocean, Marine Ecosystems, and Dependent Communities”, in *IPCC Special Report on the Ocean and Cryosphere in a Changing Climate*, edited by H.-O. Pörtner et al. (Cambridge University Press, Cambridge, UK, 2019) Chap. 5, pages 447–587 (cited on page 250).
- [551] C. 15, “Kunming-Montreal Global Biodiversity Framework”, *Convention on Biological Diversity (2022)* (cited on page 250).
- [552] S. O’Connor et al., *NOAA Coral Reef Information System, National Centers for Environmental Information*, Silver Spring, MD, 2p, 2020 (cited on page 250).
- [553] R. Carlton et al., *Global Reef Expedition Final Report*, Vol. 15 (Khaled bin Sultan Living Oceans Foundation, Annapolis, MD, 2021) (cited on page 250).
- [554] UNEP-WCMC et al., *Global distribution of warm-water coral reefs, compiled from multiple sources (listed in “Coral\_Source.mdb”), and including IMaRS-USF and IRD (2005), IMaRS-USF (2005) and Spalding et al. (2001)*, Cambridge (UK): UNEP World Conservation Monitoring Centre, 2010 (cited on page 250).
- [555] À. Giménez-Romero et al., “Universal spatial properties of coral reefs”, *Zenodo (2023)* (cited on page 251).
- [556] T. Mori et al., “Common power laws for cities and spatial fractal structures”, *Proceedings of the National Academy of Sciences* **117**, 6469–6475 (2020) (cited on page 252).
- [557] C. M. A. Pinto et al., “Double power laws, fractals and self-similarity”, *Applied Mathematical Modelling* **38**, 4019–4026 (2014) (cited on page 252).
- [558] D. A. Seekell et al., “A fractal-based approach to lake size-distributions”, *Geophysical Research Letters* **40**, 517–521 (2013) (cited on page 252).
- [559] C. M. Sorensen and G. M. Wang, “Size distribution effect on the power law regime of the structure factor of fractal aggregates”, *Phys. Rev. E* **60**, 7143–7148 (1999) (cited on page 252).
- [560] B. Vidondo et al., “Some aspects of the analysis of size spectra in aquatic ecology”, *Limnology and Oceanography* **42**, 184–192 (1997) (cited on page 252).
- [561] B. B. Mandelbrot, *The fractal geometry of nature*, 3rd edition (W. H. Freeman and Comp., New York, 1983) (cited on pages 252, 260).
- [562] S. Lovejoy, “Area-Perimeter Relation for Rain and Cloud Areas”, *Science* **216**, 185–187 (1982) (cited on page 252).

- [563] E. E. George et al., “Space-filling and benthic competition on coral reefs”, *PeerJ* **9**, e11213 (2021) (cited on pages 255, 258).
- [564] C. M. Sorensen and C. Oh, “Divine proportion shape preservation and the fractal nature of cluster-cluster aggregates”, *Phys. Rev. E* **58**, 7545–7548 (1998) (cited on page 255).
- [565] T. Nakamura and T. Nakamori, “A geochemical model for coral reef formation”, *Coral Reefs* **26**, 741–755 (2007) (cited on page 257).
- [566] S. Mistr and D. Bercovici, “A Theoretical Model of Pattern Formation in Coral Reefs”, *Ecosystems* **6**, 0061–0074 (2003) (cited on pages 257, 258).
- [567] H. Bosscher and W. Schlager, “Computer simulation of reef growth”, *Sedimentology* **39**, 503–512 (1992) (cited on page 257).
- [568] R. V. Solé and S. C. Manrubia, “Are rainforests self-organized in a critical state?”, *Journal of Theoretical Biology* **173**, 31–40 (1995) (cited on page 257).
- [569] G. B. West et al., “Scaling in biology: patterns and processes, causes and consequences”, in *Scaling in biology*, edited by J. H. Brown and G. B. West (Oxford University Press, Oxford, UK, 2000), pages 1–24 (cited on page 257).
- [570] J. H. Brown et al., “The fractal nature of nature: power laws, ecological complexity and biodiversity”, *Philosophical Transactions of the Royal Society of London. Series B: Biological Sciences* **357**, 619–626 (2002) (cited on pages 257, 258).
- [571] J. Chave and S. Levin, “Scale and Scaling in Ecological and Economic Systems”, *Environmental and Resource Economics* **26**, 527–557 (2003) (cited on page 257).
- [572] P. A. Marquet et al., “Scaling and power-laws in ecological systems”, *Journal of Experimental Biology* **208**, 1749–1769 (2005) (cited on page 257).
- [573] Á. Corral and Á. González, “Power Law Size Distributions in Geoscience Revisited”, *Earth and Space Science* **6**, 673–697 (2019) (cited on page 257).
- [574] D. Marković and C. Gros, “Power laws and self-organized criticality in theory and nature”, *Physics Reports* **536**, 41–74 (2014) (cited on page 257).
- [575] P. Bak et al., “Self-organized criticality”, *Phys. Rev. A* **38**, 364–374 (1988) (cited on page 257).
- [576] J. M. Carlson and J. Doyle, “Highly optimized tolerance: A mechanism for power laws in designed systems”, *Phys. Rev. E* **60**, 1412–1427 (1999) (cited on page 257).

- [577] J. M. Carlson and J. Doyle, “Highly Optimized Tolerance: Robustness and Design in Complex Systems”, *Phys. Rev. Lett.* **84**, 2529–2532 (2000) (cited on page 257).
- [578] M. E. J. Newman and K. Sneppen, “Avalanches, scaling, and coherent noise”, *Phys. Rev. E* **54**, 6226–6231 (1996) (cited on page 257).
- [579] T. M. Scanlon et al., “Positive feedbacks promote power-law clustering of Kalahari vegetation”, *Nature* **449**, 209–212 (2007) (cited on page 257).
- [580] S. Kéfi et al., “Spatial vegetation patterns and imminent desertification in Mediterranean arid ecosystems”, *Nature* **449**, 213–217 (2007) (cited on page 257).
- [581] K. G. Wilson, “Problems in Physics with Many Scales of Length”, *Scientific American* **241**, 158–179 (1979) (cited on page 257).
- [582] J. J. Binney et al., *The Theory of Critical Phenomena: An Introduction to the Renormalization Group* (Clarendon Press, Oxford (UK), 1992) (cited on page 257).
- [583] M. A. Muñoz, “Colloquium: Criticality and dynamical scaling in living systems”, *Review of Modern Physics* **90**, 031001 (2018) (cited on page 258).
- [584] P. E. Schmid, “Fractal Properties of Habitat and Patch Structure in Benthic Ecosystems”, in , Vol. 30, edited by A. Fitter and D. Raffaelli, *Advances in Ecological Research* (Academic Press, 1999), pages 339–401 (cited on page 258).
- [585] D. M. Kennedy and C. D. Woodroffe, “Fringing reef growth and morphology: a review”, *Earth-Science Reviews* **57**, 255–277 (2002) (cited on page 258).
- [586] M. Cross and H. Greenside, *Pattern formation and Dynamics in Nonequilibrium Systems* (Cambridge University Press, Cambridge, UK, 2009) (cited on page 258).
- [587] S. Leutenegger et al., “STR: a simple and efficient algorithm for R-tree packing”, in *Proceedings 13th International Conference on Data Engineering* (1997), pages 497–506 (cited on pages 259, 260).
- [588] S. Gillies et al., “Shapely: manipulation and analysis of geometric objects”, *GitHub Repository* (2007) (cited on pages 259, 260).
- [589] À. Giménez-Romero, “Coral Reef Analysis”, *GitHub Repository* (2022) (cited on pages 259, 260).
- [590] J. V. den Bossche et al., “Geopandas”, *Zenodo* (2023) (cited on page 259).
- [591] J. Alstott et al., “powerlaw: A Python Package for Analysis of Heavy-Tailed Distributions”, *PLoS ONE* **9**, 1–11 (2014) (cited on page 260).

- [592] A. Clauset et al., “Power-Law Distributions in Empirical Data”, *SIAM Review* **51**, 661–703 (2009) (cited on page 260).
- [593] Y. Chen, “A set of formulae on fractal dimension relations and its application to urban form”, *Chaos, Solitons & Fractals* **54**, 150–158 (2013) (cited on page 260).
- [594] M. Altman, “Traditional Districting Principles: Judicial Myths vs. Reality”, *Social Science History* **22**, 159–200 (1998) (cited on page 261).
- [595] S. J. Rey and L. Anselin, “PySAL: A Python Library of Spatial Analytical Methods”, in *Handbook of Applied Spatial Analysis: Software Tools, Methods and Applications*, edited by M. M. Fischer and A. Getis (Springer Berlin Heidelberg, Berlin, Heidelberg, 2010), pages 175–193 (cited on page 261).
- [596] À. Giménez-Romero et al., “Spatial effects in parasite-induced marine diseases of immobile hosts”, *Royal Society Open Science* **9**, 212023 (2022) (cited on page 266).
- [597] À. Giménez-Romero et al., “A Compartmental Model for *Xylella fastidiosa* Diseases with Explicit Vector Seasonal Dynamics”, *Phytopathology*® **113**, 1686–1696 (2023) (cited on page 267).
- [598] À. Giménez-Romero et al., “High-resolution climate data reveals increased risk of Pierce’s Disease for grapevines worldwide”, *bioRxiv* (2024) (cited on page 269).
- [599] S. Flecha et al., “pH trends and seasonal cycle in the coastal Balearic Sea reconstructed through machine learning”, *Scientific Reports* **12**, 12956 (2022) (cited on page 270).
- [600] À. Giménez-Romero et al., “Mapping the distribution of seagrass meadows from space with deep convolutional neural networks”, *bioRxiv* (2024) (cited on page 270).
- [601] O. Diekmann and J. A. P. Heesterbeek, *Mathematical Epidemiology of Infectious Diseases. Model Building, Analysis and Interpretation* (John Wiley and Sons, Chichester (UK), 2000) (cited on page 332).
- [602] J. Cariboni et al., “The Role of Sensitivity Analysis in Ecological Modelling”, *Ecological Modelling* **203**, 167–182 (2007) (cited on page 333).
- [603] J. F. Gillooly et al., “Effects of Size and Temperature on Metabolic Rate”, *Science* **293**, 2248–2251 (2001) (cited on pages 333, 334).
- [604] J. R. Coelho and F. S. Bezerra, “The effects of temperature change on the infection rate of *Biomphalaria glabrata* with *Schistosoma mansoni*”, in *Memórias do Instituto Oswaldo Cruz* **101**, 223–224 (2006) (cited on page 334).

- 
- [605] L. L. M. Shapiro et al., “Quantifying the effects of temperature on mosquito and parasite traits that determine the transmission potential of human malaria”, *PLOS Biology* **15**, e2003489 (2017) (cited on page 334).
- [606] EPPO, “PM 7/24-3 *Xylella fastidiosa*”, *EPPO Bulletin. European and Mediterranean Plant Protection Organisation* **48**, 175–218 (2018) (cited on page 343).
- [607] R. D. C. Team, “R: A language and environment for statistical computing”, 2017 (cited on page 344).
- [608] D. Bates et al., “Fitting Linear Mixed-Effects Models Using lme4”, *Journal of Statistical Software* **67**, 1–48 (2015) (cited on page 344).
- [609] M. Daugherty et al., “Severe pruning of infected grapevines has limited efficacy for managing Pierce’s disease”, *American Journal of Enology and Viticulture* **69**, 289–294 (2018) (cited on page 349).
- [610] R. M. Maier, “Chapter 3 - Bacterial Growth”, in *Environmental Microbiology (Second Edition)*, edited by R. M. Maier et al., Second Edition (Academic Press, San Diego, 2009), pages 37–54 (cited on page 352).
- [611] F. Pedregosa et al., “Scikit-learn: Machine Learning in Python”, *Journal of Machine Learning Research* **12**, 2825–2830 (2011) (cited on page 361).



## A. Analysis of the SIRP model

### A.1 Finding a conserved quantity for the SIRP model

Starting with the SIRP model,

$$\begin{aligned}\dot{S} &= -\bar{\beta}PS \\ \dot{I} &= \bar{\beta}PS - \gamma I \\ \dot{R} &= \gamma I \\ \dot{P} &= \lambda I - \bar{\beta}PS - \mu P ,\end{aligned}\tag{A.1}$$

from the  $\dot{S}$  equation,  $P$  can be written as follows,

$$P = -\frac{1}{\bar{\beta}} \frac{\dot{S}}{S},\tag{A.2}$$

and summing up the equations for  $\dot{S}$  and  $\dot{I}$  the following relation for  $I$  is obtained

$$I = -(\dot{S} + \dot{I})/\gamma.\tag{A.3}$$

Replacing Eq. (A.2), Eq. (A.3) and the differential equation for  $\dot{S}$  in the 4th differential equation in Eq. (A.1) one obtains,

$$\dot{P} = -\frac{\lambda}{\gamma}(\dot{S} + \dot{I}) + \dot{S} + \frac{\mu}{\bar{\beta}} \cdot \frac{\dot{S}}{S}\tag{A.4}$$

As  $\dot{S}/S = d(\ln S)/dt$ , all terms in the previous equation are exact differentials with respect to time, and the equation can be integrated yielding,

$$P + \frac{\lambda}{\gamma}(S+I) - S - \frac{\mu}{\beta} \ln S = C \quad (\text{A.5})$$

with the integration constant  $C$ , that is a conserved quantity, i.e., it takes the same value at one time of the dynamical evolution of the system.  $C$  is related to the initial conditions by,

$$\begin{aligned} C &= P(0) + \frac{\lambda}{\gamma}(S(0) + I(0)) - \frac{\mu}{\beta} \ln S(0) - S(0) = \\ &P(0) + \frac{\lambda}{\gamma}(N - R(0)) - \frac{\mu}{\beta} \ln S(0) - S(0) \end{aligned} \quad (\text{A.6})$$

It is possible to use Eq. (A.5)-Eq. (A.6) to express one of variables as a function of the others, for example the parasite concentration  $P$  as,

$$P(S, I) = P(0) - \frac{\lambda}{\gamma}(S+I - N + R(0)) + \frac{\mu}{\beta} \ln \frac{S}{S(0)} + S - S(0), \quad (\text{A.7})$$

or equivalently as,

$$P(S, R) = P(0) + \frac{\lambda}{\gamma} \left[ R - R(0) + \frac{\mu\gamma}{\beta\lambda} \ln \frac{S}{S(0)} \right] + S - S(0) \quad (\text{A.8})$$

From Eq. (A.5), it is easy to show that the SIP model of Ref. [182], that differs from the SIRP model in that the fourth equation is simplified to  $\dot{P} = \lambda I - \mu P$ , has as exact conserved quantity,

$$P + \frac{\lambda}{\gamma}(S+I) - \frac{\mu}{\beta} \ln S = \mathcal{C} \quad (\text{A.9})$$

as the extra term in the SIRP model  $-\bar{\beta}SP$  is equal to  $\dot{S}$  from the first equation Eq. (A.1).

The SIR model has a conserved quantity [49], that in the case of Eq. (3.9) takes the form,

$$I + S - \frac{\gamma}{\beta'} \ln S = C. \quad (\text{A.10})$$

Rewriting Eq. (A.5) in the alternative form,

$$\frac{\gamma}{\lambda}P + \left(1 - \frac{\gamma}{\lambda}\right)S + I - \frac{\mu\gamma}{\lambda\beta} \ln S = C' \quad (\text{A.11})$$

it can be seen that if  $\lambda \gg \gamma$  Eq. (A.11) reduces to Eq. (A.10), remembering that in Eq. (3.9)  $\beta' = \lambda\bar{\beta}/\mu$ . The assumptions used to arrive to Eq. (A.10) in Section 3.2.3.3 where  $\mu \gg (\gamma, \bar{\beta})$ , and taking into account the expression for  $R_0$  Eq. (3.3), that  $\lambda \gtrsim \mu$  is most plausible to keep  $R_0$  above the epidemic threshold ( $R_0 > 1$ ).

## A.2 Stability analysis of the fixed points of the SIRP model

Here we will assume the initial fixed point of our SIRP model, with  $I(0) = P(0) = 0$  right before the introduction of the infection, either through  $I$  or  $P$ . We will assume that  $R(0) = 0$ , so that  $S(0) = N$ . To study the linear stability of the model we need to write the Jacobian, that takes the form,

$$J = \begin{pmatrix} -\bar{\beta}P & 0 & 0 & \bar{\beta}S \\ \bar{\beta}P & -\gamma & 0 & \bar{\beta}S \\ 0 & \gamma & 0 & 0 \\ -\bar{\beta}P & \lambda & 0 & (\bar{\beta}S - \mu) \end{pmatrix} \quad (\text{A.12})$$

and obtain the eigenvalues for both fixed points, where we have already used the standard incidence,  $\bar{\beta} = \beta/N$ , from the evidence of the validation with experiments. For the pre-epidemic fixed point, the Jacobian becomes,

$$\begin{pmatrix} 0 & 0 & 0 & \bar{\beta}S(0) \\ 0 & -\gamma & 0 & \bar{\beta}S(0) \\ 0 & \gamma & 0 & 0 \\ 0 & \lambda & 0 & (\bar{\beta}S(0) - \mu) \end{pmatrix} \quad (\text{A.13})$$

Matrix Eq. (A.13) has two null (0) eigenvalues and a pair of eigenvalues given by,

$$\Lambda_{1,2} = -\frac{1}{2}(\gamma + \mu + \bar{\beta}S(0) \pm \sqrt{\gamma^2 + \mu^2 + (\bar{\beta}S(0))^2 + 2\mu\bar{\beta}S(0) - 2\gamma\mu - 2\gamma\bar{\beta}S(0) + 4\lambda\bar{\beta}S(0)}) \quad (\text{A.14})$$

from which one can determine that the fixed point is unstable whenever

$$\lambda\bar{\beta}S(0) > \gamma(\mu + \bar{\beta}S(0)) \quad (\text{A.15})$$

and stable if the inequality is reversed. It can be easily shown that Eq. (A.15) is equivalent to  $R_0 > 1$ , with  $R_0$  given by Eq. (3.3).

The final point of the epidemic,  $S(\infty)$ , can be found by solving the transcendental equation,

$$\left(\frac{\lambda}{\gamma} - 1\right)S(\infty) - \frac{\mu}{\bar{\beta}}\ln(S(\infty)) = C \quad (\text{A.16})$$

where  $C$  is determined from the initial conditions (Eq. (A.6)) and  $I(\infty) = P(\infty) = 0$ . (Eq. (A.16)) has two roots, where  $S(\infty)$  represents the smallest one.

### A.3 Calculation of $R_0$ using the Next Generation Matrix method

Here we will use the Next Generation Matrix (NGM) method, explained in Section 2.1.1, to calculate the basic reproduction number  $R_0$  for the SIRP model. The NGM method is based on the decomposition of the Jacobian matrix of the system of differential equations into a transmission matrix  $T$  and a transition matrix  $\Sigma$ , where the transmission matrix contains the terms that contribute to the transmission of the disease, while the transition matrix contains the terms that imply transitions between compartments [601]. The basic reproduction number  $R_0$  is then given by the dominant eigenvalue of the matrix  $K = -T\Sigma^{-1}$ .

In the case of the SIRP model the decomposition is applied to the  $2 \times 2$  Jacobian corresponding to the dynamical evolution of the  $(I, P)$  infectious compartments, being the decomposition,

$$J = \begin{pmatrix} -\gamma & \bar{\beta}S_0 \\ \lambda & -(\bar{\beta}S_0 + \mu) \end{pmatrix} \quad T = \begin{pmatrix} 0 & \bar{\beta}S_0 \\ 0 & 0 \end{pmatrix} \quad \Sigma = \begin{pmatrix} -\gamma & 0 \\ \lambda & -(\bar{\beta}S_0 + \mu) \end{pmatrix}$$

where the  $\bar{\beta}PS$  term in  $\dot{I}$  is the only one that contributes to the transmission matrix, as it is the only process involving infection, while all the other terms in the dynamical equations of  $\dot{I}$  and  $\dot{P}$  imply transitions (to another compartment, like  $I \rightarrow R$  or birth and death of  $P$ ).

Then, the next generation matrix is given by,

$$K = -T\Sigma^{-1} = \begin{pmatrix} \frac{\lambda\beta S_0}{\gamma(\beta S_0 + \mu)} & \frac{\beta S_0}{\beta S_0 + \mu} \\ 0 & 0 \end{pmatrix} \implies R_0 = \frac{\lambda\beta S_0}{\gamma(\beta S_0 + \mu)},$$

This result coincides with the expectation that  $R_0$  should correspond to the number of hosts infected in a single generation by the appearance of an infected host in a completely susceptible population. This can be obtained from the number of parasites produced by an infected host,  $\lambda$ , times the time in which the infected host is alive producing parasites,  $1/\gamma$ , multiplied by the number of infected hosts produced per parasite,  $\beta S_0$ , times the time the parasite is alive available to infect,  $1/(\mu + \beta S_0)$ , taking into account that parasites are inactivated at a rate  $\mu$  and also die when infecting at a rate  $\beta S_0$ , where this result assumes that the susceptible population does not change from its initial value  $S_0$ .

### A.4 Sensitivity Analysis

One particular way to analyse the local sensitivity (LSA) of a given model function,  $F(\vec{p})$ , for each of the parameters that conform it,  $p_i$ , is through the

normalised sensitivity indexes [602],

$$\Omega_{p_i}^F = \frac{\partial F}{\partial p_i} \frac{p_i}{F} \Big|_{p_i=p^0} . \quad (\text{A.17})$$

where the partial derivatives in Eq. (A.17) are determined analytically in our case.

GSA works by studying the influence of a large domain of parameter space in the final state of the epidemic and in the epidemic peak. In our case this will be achieved by means of a variance based analysis, known as Sobol method [295]. This particular method provides information not only on how a particular parameter alone influences the model outputs (as happens with LSA), but also on the influence of its interactions with other parameters. This information is organised in what are known as Sobol indices, that have been implemented within the Julia high-level programming language [198] using the DifferentialEquations.jl package [199], and in particular through its subpackage DiffEqSensitivity.jl. This implementation allows the user to sample the parameter space using QuasiMonteCarlo methods and thus obtain confidence intervals (CI) for the sensitivity indices, which are directly related to the committed statistical error.

The total order indices are a measure of the total variance of the output quantity caused by variations of the input parameter and its interactions. First order (or “main effect”) indices are a measure of the contribution to the output variance given by the variation of the parameter alone, but averaged over variations in other input parameters. Second order indices take into account first order interactions between parameters. Further indices can be obtained, describing the influence of higher-order interactions between parameters, but these are not going to be considered. More detailed information about sensitivity analysis can be found in [294].

## A.5 General rate change with temperature

In [603] the metabolic rate of a wide variety of organisms was studied, showing that the change in the metabolic rate with temperature was similar among them. In particular, the natural logarithm of the metabolic rate linearly depends on the inverse of absolute temperature,

$$\log(R(T)) = a \cdot \left( \frac{100}{T} \right) + b \quad (\text{A.18})$$

and for all the analysed organisms they found that  $a$  lies between  $-5$  and  $-10$  and  $b$  between  $14$  and  $30$ . From their analysis, we can compute the change in

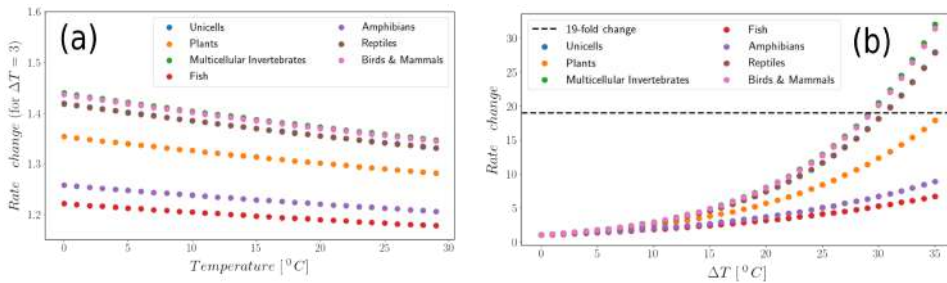
the rate for a given increase of temperature,

$$\frac{R(T + \Delta T)}{R(T)} = \frac{\exp(a \cdot 100 / (T + \Delta T) + b)}{\exp(a \cdot 100 / T + b)} = \exp\left(a \cdot \frac{-100}{T + \Delta T} \cdot \frac{\Delta T}{T}\right). \quad (\text{A.19})$$

Substituting  $T = 287\text{K}$  and  $\Delta T = 3\text{K}$ , that correspond to our available data (cf. Section 3.4) in Eq. (A.19), using both the upper and lower limit of  $a$ , we obtain that the expected increase in the effective transmission rate is between 1.2 to 1.4. This is far from the 19-fold increase that we obtained with the mass action hypothesis in Section 3.4 while it is in good agreement with either the 1.92 ratio we obtained for  $\bar{\beta}$  with the reduction of Section 3.2.3.2 or the 1.43 ratio obtained with the fast-slow approximation of Section 3.2.3.3, both obtained using the standard incidence choice.

Fig. A.1(a) shows the change in the rate with an increase of  $3^\circ\text{C}$  for different base temperatures and for all the organisms analysed in [603], and using their fit. Note that for all temperatures between  $0^\circ\text{C}$  and  $30^\circ\text{C}$  the rate change lies between 1.2 and 1.45. Fig. A.1(b) shows the change in the rate for different temperature increases, with a base temperature of  $T = 287\text{K}$ . Note that in order to obtain a 19-fold increase the temperature change should be at least of  $30^\circ\text{C}$ <sup>1</sup>. The temperature dependence of metabolic rates has been reported in the context of epidemic parameters [604, 605]

The behavior of the metabolic rates re-analysed here has been also found experimentally in epidemic contexts such as [604, 605], i.e., the increase of the rates with temperature fulfill the ranges shown here.



**Figure A.1:** Graphical representation of change in the rate (in ordinates) for different reference temperatures (in abscissae) for: (a) a temperature increase of  $3^\circ\text{C}$ ; (b) a temperature increase of  $14^\circ\text{C}$ . The black dotted line in (b) corresponds to a 19-fold increase in the rate.

<sup>1</sup>A temperature change of  $30^\circ\text{C}$  could fall outside the range in which the study of [603] is valid. We just stress that a 19-fold rate change is unlikely for the case of a  $3^\circ\text{C}$  that correspond to the 2 data sets that we compare in this section.

## A.6 Derivation of the non-spatial equation for $R_\infty$

The model described by the ODE system in Eq. (4.2) has a conserved quantity  $\mathcal{C}$  given by [215].

$$\mathcal{C} = P + \frac{\lambda}{\gamma} (S + I) - S - \frac{\mu}{\beta} \ln S \quad (\text{A.20})$$

At  $t = \infty$  the system reaches an absorbing state completely determined by  $S(\infty)$ , as  $P(\infty) = I(\infty) = 0$  and  $N = S(\infty) + R(\infty)$ . Thus, from Eq. (A.20) we have

$$S(\infty) \left( \frac{\lambda}{\gamma} - 1 \right) - \frac{\mu}{\beta} \ln(S(\infty)) = \mathcal{C}_0 \quad (\text{A.21})$$

The transcendental equation Eq. (A.21) can be solved by means of the Lambert's W function,

$$S(\infty) = -\frac{\mu\gamma}{\beta(\lambda-\gamma)} W_0 \left( -\frac{\beta(\lambda-\gamma)}{\mu\gamma} \exp(-\beta\mathcal{C}_0/\mu) \right) \quad (\text{A.22})$$

which can be simplified to

$$S(\infty) = -\frac{S(0)}{\xi} W_0 \left( -\xi \exp \left( -\frac{\beta}{\mu} C \right) \right), \quad (\text{A.23})$$

with  $\xi = S(0) \frac{\beta(\lambda-\gamma)}{\mu\gamma}$  and  $C = P(0) + \frac{\lambda}{\gamma} (S(0) + I(0)) - S(0)$ .

Finally, the absorbing state fulfils the condition  $N = S(\infty) + R(\infty)$  so that the final number of dead individuals can be expressed as

$$R(\infty) = N + \frac{S(0)}{\xi} W_0 \left( -\xi \exp \left( -\frac{\beta}{\mu} C \right) \right). \quad (\text{A.24})$$



## B. The SIR-V model and its application to Xf diseases

### B.1 Calculation of $R_0$ from standard methods

The standard methods of calculation of  $R_0$  are based in the linear stability analysis of the disease-free equilibrium, either directly, through the linear analysis of the fixed point, that yields the stability condition from which  $R_0$  can be obtained, or using the Next Generation Method (NGM) [68] that provides directly  $R_0$  by solving a suitable linear problem. Customarily these methods are applied to a pre-pandemic disease-free equilibrium, but as there is no such state in the case of non-stationary populations, here a similar approach is applied to a post-pandemic or asymptotic disease-free equilibrium.

#### Linear stability analysis

In order to perform the linear stability analysis of the fixed point ( $I_H = I_V = 0$ ) we first need to compute the Jacobian matrix,  $J$ ,

$$J = \begin{pmatrix} -\beta \frac{I_V}{N_H} & 0 & 0 & -\beta \frac{S_H}{N_H} \\ \beta \frac{I_V}{N_H} & -\gamma & 0 & \beta \frac{S_H}{N_H} \\ 0 & -\alpha \frac{S_V}{N_H} & -\alpha \frac{I_H}{N_H} - \mu & 0 \\ 0 & \alpha \frac{S_V}{N_H} & \alpha \frac{I_H}{N_H} & -\mu \end{pmatrix} \quad (\text{B.1})$$

Then, we evaluate the Jacobian at the fixed point (or disease free equilibrium, DFE), yielding

$$J|_{DFE} = \begin{pmatrix} 0 & 0 & 0 & -\beta \\ 0 & -\gamma & 0 & \beta \\ 0 & -\alpha \frac{C}{N_H} \frac{\delta}{\mu} & -\mu & 0 \\ 0 & \alpha \frac{C}{N_H} \frac{\delta}{\mu} & 0 & -\mu \end{pmatrix} \quad (\text{B.2})$$

where  $S_H = N_H$  has been considered.

The eigenvalues of Eq. (B.2) are,

$$\begin{aligned} \lambda_0 &= 0 \\ \lambda_\mu &= -\mu \\ \lambda_\pm &= -\frac{(\gamma + \mu)}{2} \pm \frac{1}{2} \sqrt{(\gamma - \mu)^2 + 4\beta\alpha \frac{C}{N_H} \frac{\delta}{\mu}} \end{aligned} \quad (\text{B.3})$$

It is straightforward to see that all eigenvalues are real and the stability of the disease-free equilibrium is determined by the sign of the eigenvalues.  $\lambda_\mu = -\mu < 0$  as  $\mu$  is defined positive, so in order to discuss the stability of this fixed point, we need to study the  $\lambda_\pm$  eigenvalues.  $\lambda_-$  is always negative, but  $\lambda_+$  changes sign depending on the values of the parameters. The threshold condition  $\lambda_+ = 0$  leads to:

$$\lambda_+ = 0 \Rightarrow \frac{\beta\alpha}{\gamma\mu} \frac{C}{N_H} \frac{\delta}{\mu} = 1 \quad (\text{B.4})$$

So, for  $\frac{\beta\alpha}{\gamma\mu} \frac{C}{N_H} \frac{\delta}{\mu} < 1 \Rightarrow \lambda_+ < 0$  the fixed point is stable and for  $\frac{\beta\alpha}{\gamma\mu} \frac{C}{N_H} \frac{\delta}{\mu} > 1 \Rightarrow \lambda_+ > 0$  a perturbation will grow in the direction of the eigenvector associated to  $\lambda_+$ . Thus, this threshold defines the basic reproduction number,

$$R_0 = \frac{\beta\alpha}{\gamma\mu} \frac{C}{N_H} \frac{\delta}{\mu} \quad (\text{B.5})$$

If instead of  $S_H = N_H$  one considers any initial condition of hosts,  $S_H(0)$ , the basic reproduction number is given by,

$$R_0 = \frac{\beta\alpha}{\gamma\mu} \frac{C}{N_H} \frac{\delta}{\mu} \frac{S_H(0)}{N_H} \quad (\text{B.6})$$

### Next Generation Matrix method

The previous result can also be obtained by means of the NGM method, which is explained in detail in [68]. Basically the method is based in decomposing

the Jacobian in the form  $J = T + \Sigma$ , where  $T$  is the *transmission part*, that describes the production of new infections, and  $\Sigma$  the *transition part*, that describes changes of state (including death). Then, it can be proved [68] that the *basic reproduction number*  $R_0$  is given by the spectral radius (i.e. the largest eigenvalue) of the (next generation) matrix  $K = -T\Sigma^{-1}$ .

$$K = -T\Sigma^{-1} = \begin{pmatrix} \frac{\beta\alpha}{\gamma\mu} \frac{C}{N_H} \frac{\delta}{\mu} & \frac{\beta}{\mu} \\ 0 & 0 \end{pmatrix} \quad (\text{B.7})$$

with,

$$T = \begin{pmatrix} 0 & \beta \frac{N_H}{N_H} \\ 0 & 0 \end{pmatrix}, \quad \Sigma = \begin{pmatrix} -\gamma & 0 \\ \alpha \frac{C}{N_H} \frac{\delta}{\mu} & -\mu \end{pmatrix}$$

$$\text{and } -\Sigma^{-1} = \begin{pmatrix} \frac{1}{\gamma} & 0 \\ \frac{\alpha}{\gamma\mu} \frac{C}{N_H} \frac{\delta}{\mu} & \frac{1}{\mu} \end{pmatrix}$$

The basic reproduction number is the spectral radius of this matrix so:

$$\begin{aligned} \det(K - \sigma\mathbb{I}) = 0 &\implies \begin{vmatrix} \frac{\beta\alpha}{\gamma\mu} \frac{C}{N_H} \frac{\delta}{\mu} - \sigma & \frac{\beta}{\mu} \\ 0 & -\sigma \end{vmatrix} = \\ &= (-\sigma) \left( \frac{\beta\alpha}{\gamma\mu} \frac{C}{N_H} \frac{\delta}{\mu} - \sigma \right) = 0. \end{aligned}$$

Solving for  $\sigma$  one obtains the solutions,

$$\sigma_1 = \frac{\beta\alpha}{\gamma\mu} \frac{C}{N_H}; \quad \sigma_2 = 0 \quad (\text{B.8})$$

Therefore, the basic reproduction number is

$$R_0 = \frac{\beta\alpha}{\gamma\mu} \frac{C}{N_H} \frac{\delta}{\mu} \quad (\text{B.9})$$

If instead of  $S_H = N_H$  one considers any initial condition of hosts,  $S_H(0)$ , the basic reproduction number is given by,

$$R_0 = \frac{\beta\alpha}{\gamma\mu} \frac{C}{N_H} \frac{\delta}{\mu} \frac{S_H(0)}{N_H} \quad (\text{B.10})$$

## B.2 Calculation of $R_0$ for non-stationary vector populations

We extend the computation of  $R_0$  in the case of non-stationary and non-periodic vector populations by following the natural definition of *basic reproductive number*. Thus,  $R_0$  is computed by averaging the number of secondary infections

produced by an infected individual along one generation, that is equivalent to averaging the instantaneous definition of  $R_0$ , namely  $R_0^i$ , over one generation,

$$\overline{R_0} = \langle R_0^i(t) \rangle \Big|_0^{t_g} = \frac{R_0}{N_v^*} \langle N_v(t) \rangle \Big|_0^{t_g} = \frac{R_0}{N_v^*} \frac{1}{t_g} \int_0^{t_g} N_v(t) dt, \quad (\text{B.11})$$

where the integral in Eq. (B.17) is solved as

$$\begin{aligned} \int_0^{t_g} N_v(t) dt &= \left[ N_v^* t - \frac{1}{\mu} (N_v(0) - N_v^*) e^{-\mu t} \right]_0^{t_g} = \\ &N_v^* t_g - \frac{1}{\mu} (N_v(0) - N_v^*) [e^{-\mu t_g} - 1]. \end{aligned} \quad (\text{B.12})$$

Thus, the basic reproduction number for non-stationary vector populations is given by

$$\overline{R_0} = \frac{R_0}{N_v^*} \left\{ N_v^* - \frac{1}{\mu t_g} [N_v(0) - N_v^*] [e^{-\mu t_g} - 1] \right\}, \quad (\text{B.13})$$

where the generation time,  $t_g$ , is Eq. (5.8). Eq. (B.19) can be rewritten as,

$$\overline{R_0} = \langle R_0^i(t) \rangle \Big|_0^{t_g} = R_0 \left[ 1 - \frac{1}{\tau} (f - 1) (e^{-\tau} - 1) \right] = R_0 \cdot \mathcal{F}, \quad (\text{B.14})$$

where  $\tau = 1 + \mu/\gamma$  and  $\mathcal{F}$  is the expression in brackets, which accounts for the effect of the decaying vector population on the stationary  $R_0$ .

In our approach, a generation is defined as the time elapsed in the following sequence of processes: 1) A host individual becomes infected; 2) The infected host passes the disease to a susceptible vector; 3) The infected vector dies. Basically, the time elapsed from the first to the last process is the time in which new infections can be produced, i.e.,  $t_g$  Eq. (5.8).

### B.3 Determination of $R_0$ for Xf diseases

The handicap of determining the basic reproductive number of the model Eq. (6.2) is that the pre-pandemic fixed point given by  $I_H = I_v = 0$  and  $S_H = S_H(0)$  is not a fixed point of the system of differential equations, because vector population decays, so that the standard methods to compute  $R_0$  such as the Next Generation Matrix [68, 289] do not apply. In [289] a method was suggested to determine the basic reproductive number in the case of compartmental models of vector-borne transmitted diseases in which the vector population grows or decays. It consists in averaging the instantaneous basic reproductive number over the time of a generation.

To proceed we consider that  $I_H = I_v = 0$ ,  $S_H = S_H(0)$  is indeed a fixed point of the system. Then, the basic reproductive number could be determined, e.g.,

as shown in [69]. First, an infectious host infects vectors at a rate  $\beta S_H(0)/N_H$  for a time  $1/\gamma$ . This produces  $\beta S_H(0)/\gamma N_H$  infected vectors. The second stage is that these infectious vectors infect hosts at a rate  $\alpha N_v(0)/N_H$  for a time  $1/\mu$ , producing  $\alpha N_v/\mu N_H$  infectious hosts per vector. The net result of these two stages is

$$\tilde{R}_0 = \frac{\alpha\beta}{\mu\gamma} \frac{S_H(0)}{N_H^2} N_v(0) = R_0^* \cdot N_v(0) . \quad (\text{B.15})$$

This result coincides with the value of  $R_0$  obtained using the standard NGM method, that can be applied in this case because we are assuming that we use a nongeneric initial condition that sits at the fixed point of the model.

In practice, our initial condition will never be a fixed point of the model, and, as mentioned above, we will obtain an approximate basic reproductive number, to which we will refer as  $R_0$  using the method suggested in [289], that consists in calculating the *average* number of secondary infections produced by an infectious host in *one generation*. One first defines an instantaneous basic reproductive number,

$$R_0^{(i)}(t) = \frac{\beta\alpha}{\mu\gamma} \frac{S_H(0)}{N_H^2} N_v(t) = R_0^* N_v(t) , \quad (\text{B.16})$$

from which the average is simply computed as

$$R_0 = \left\langle R_0^{(i)}(t) \right\rangle \Big|_0^\tau = R_0^* \langle N_v(t) \rangle \Big|_0^\tau = R_0^* \frac{1}{\tau} \int_0^\tau N_v(t) dt . \quad (\text{B.17})$$

In our model, the time-dependent vector population can be obtained from Eq. (6.2),

$$\dot{N}_v = \dot{S}_v + \dot{I}_v = -\mu N_v \implies N_v(t) = N_v(0) e^{-\mu t} , \quad (\text{B.18})$$

and introducing this expression for  $N_v(t)$  in Eq. (B.17) the integral can be solved

$$R_0 = \frac{\beta\alpha S_H(0)}{\mu\gamma N_H^2} \frac{N_v(0)}{\mu\tau} (1 - e^{-\mu\tau}) = R_0^* \frac{N_v(0)}{\mu\tau} (1 - e^{-\mu\tau}) , \quad (\text{B.19})$$

that is an approximated expression to the basic reproductive number for our model, in which the vector population is nonstationary, where, in Eq. (B.16) and Eq. (B.19) it has been defined,  $R_0^* = (\beta\alpha S_H(0))/(\mu\gamma N_H^2)$ .

Note that in our model one generation correspond to one year and that  $N_v(0)$  is reset every year.



## C. Modeling Pierce's disease risk

### C.1 Inoculation tests on European grapevine varieties

**Plants.** Grapevine saplings were annually supplied from a nursery in mainland Spain (Viveros Villanueva Vides, SL), consisting of one-year-old rootstocks grafted in winter with dormant grapevine cultivars, and grown in 20-L plastic pots with a standard potting mix. Fifty-seven rootstock-scion cultivar combinations were used in the inoculation assay (Table C.1). Potted plants were randomly distributed in 12-plant rows along an insect-proof tunnel exposed to air temperature and daily drip-irrigated to field capacity, fortnightly sprinkled with a slow-release fertiliser and treated with insecticides and fungicides when needed until the end of the experiment. Two weeks before the onset of the inoculation assay, leaf samples of all plants were collected and tested for the presence of Xf through qPCR as described elsewhere [268].

**Isolates and inoculation.** We used for the inoculation experiment two isolates of Xf. subsp. *fastidiosa* (ST1) recovered from grapevines: XLY 2055/17 (GenBank WGS: QTJS01) and XYL2177/18 (JAAGVM01) [265, 307]. In the third-year assay, we included an isolate of Xf subsp. *multiplex* ST81 XYL1981/18 (JAAGV1) to test whether other strains in Majorca could cause PD as well. Isolates were grown on BYCE medium at 28°C for 7-10 days, following EPPO protocols [606]. Cells were collected by scraping the colonies and suspending them in 1.5 ml Eppendorf tubes each with 1 ml of phosphate-buffered saline (PBS) solution until obtaining a turbid ( $10^8 - 10^9$  cell/ml) suspension. Plants were mechanically inoculated by pin-prick inoculation [130] with slight modifi-

cations. A 10- $\mu$ l drop of the bacterial suspension was pipetted on the leaf axil and punctured five times with an entomological needle. Eight-nine replicates per scion-rootstock combination were inoculated with the bacterial suspension and four-three plants per cultivar with a drop of PBS as a control at the end of May. Inoculation was repeated two weeks thereafter by piercing the next leaf axil above that previously inoculated [268].

**Disease score.** Disease severity was rated by counting the number of symptomatic leaves eight weeks post-inoculation (WPI) and then biweekly until the 16th week. A disease index was calculated according to Su et al. [308]. To determine the basipetal and systemic movement of Xf<sub>PD</sub>, we counted the number of symptomatic leaves below the point of inoculation from the same stems or any stem below at 12 WPI. Symptomatic and asymptomatic plants were tested by qPCR for Xf infection at 12 WPI taking the petiole of the second and fifth leaf above the point of inoculation. On the 14th week, five leaves per plant of all inoculated plants were used for Xf<sub>PD</sub> isolation, as described below. Those plants for which the qPCR was negative and Xf<sub>PD</sub> could not be isolated were treated as not infected.

**Data analysis.** All statistical analyses on disease scores were carried out using R. 3.5.2 version software [607]. We used the functions `glm` and `glmer` in the R package `lme4` [608] for fitting Generalized Linear Models and Generalized Linear Mixed Models (GLMMs) in the analysis of disease incidence and severity in the inoculation assays. In all tests, we modeled the response variable (i) disease incidence with the binomial error (logit-link function) and (ii) disease severity with the Poisson error (log-link function). A within-subject (repeated measures) factorial design was performed to evaluate differences in disease severity over time among different cultivar-rootstock combinations. Cultivars-rootstock and time were treated as fixed factors and plant subjects as a random effect. Rootstock and time were analyzed as fixed factors and plant subjects as a random effect. Controls were excluded from the analysis, as lesions did not develop on them. For each cultivar, we calculated the area under the disease progress curve (AUDPC) from weeks-post-inoculation and disease index using the package `Agricolae`. To test whether genotypes within the ST1-grapevine population vary in virulence, we included a second strain (XYL2177/18) in the 2019 inoculation experiment.

**Varietal response to Xf.** In our three-year inoculation tests, we included a representative number of local and international varieties (Table C.1). In total, among 886 inoculated-grapevine plants comprising 36 varieties in 57 unique

combinations (scion-rootstock), 86.1% ( $n = 764$ ) of them developed typical PD symptoms at 16 WPI. In contrast, none of the grapevine plants inoculated with the strain XYL 1981/17, ST81 subsp. *multiplex* presented symptoms. The results of the pathogenicity tests on European grape varieties are shown in [Table C.1](#).

Table C.1: **Summary of the inoculation tests on grapevine varieties ranked from most to less susceptible in the disease index.** Thirty-six local, regional and international varieties were screened in combination with eight rootstocks. The number of symptomatic leaves was counted 16 weeks after inoculation and infections were confirmed by qPCR. DI: disease index; AUDCP: area under the disease progress curve.

Scion	Rootstock	Nº leaves	DI	AUDCP	% Positive	Year
Gorgollassa	R110	24.5 ± 8.8	5.00	31.29	100	2018
Sauvignon Blanc	R110	16.5 ± 3.3	5.00	28.37	100	2018
Tempranillo	SO4	15.7 ± 4.0	5.00	37.22	100	2019
Garnacha tintorera	R110	19.0 ± 7.4	5.00	32.17	94.44	2019
Tempranillo	41B	13.8 ± 2.7	4.89	13.43	0.00	2020
Syrah	R140	13.6 ± 3.0	4.86	29.46	98.21	2019
Tempranillo Blanco	R110	16.0 ± 4.9	4.83	41.61	94.44	2019
Chardonnay	R110	16.6 ± 5.6	4.83	26.39	94.44	2019
Bobal	R110	15.2 ± 6.8	4.78	25.44	94.44	2019
Prensal	161/49	15.7 ± 5.0	4.75	21.37	100	2018
Viura	SO4	16.7 ± 4.7	4.75	26.25	100	2018
Garnacha tintorera	P1103	15.8 ± 6.8	4.72	35.94	94.44	2019
Graciano	R140	14.3 ± 5.6	4.61	22.44	88.89	2019
Airen	R110	15.2 ± 5.9	4.61	23.06	94.44	2019
Mandó	R110	11 ± 3.6	4.56	21.22	100	2018
Tempranillo	R110 RJ43	11.9 ± 3.0	4.56	31.17	100	2019
Tempranillo	R110 RJ78	11.7 ± 3.8	4.50	34.33	94.44	2019
Tempranillo	41B	12.4 ± 4.7	4.44	26.06	61.11	2019

Continued on next page

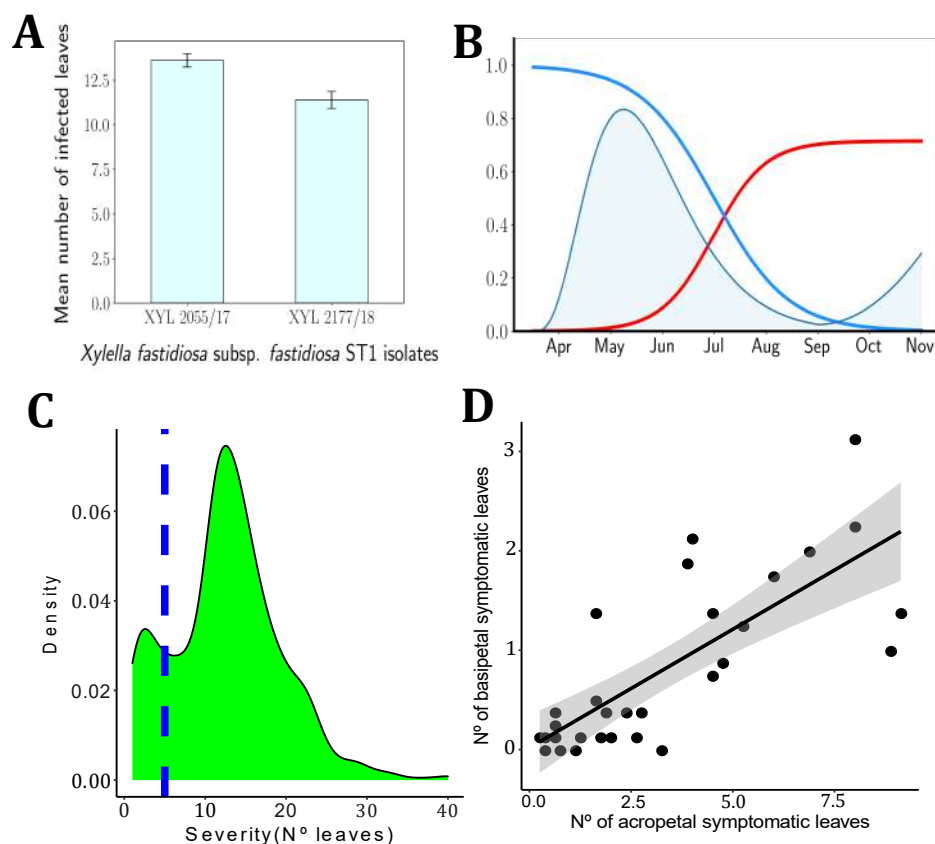
Table C.1: Inoculation tests on different grapevine varieties (Continued)

Garnacha	R110	11.1 ± 3.4	4.44	25.00	83.33	2019
Viura	P1103	20.3 ± 8.4	4.37	25.875	87.5	2018
Malvasia	R110	13.8 ± 6.5	4.33	26.71	71.73	2019
Tempranillo	R110 RJ43	12.2 ± 5.9	4.33	24.23	88.89	2020
Pedro Ximenez	R140	10.0 ± 4.33	4.33	16.87	0.00	2020
Hondarrabi Beltza	196-17	11.2 ± 4.6	4.22	19.90	87.5	2019
Tempranillo	P1103	12.7 ± 6.3	4.17	30.62	62.5	2018
Albariño	R110	9.5 ± 4.3	4.06	24.50	66.67	2019
Garnacha tintorera	P1103	13.1 ± 8.0	4.0	16.85	88.89	2020
Tempranillo	R110	11.7 ± 7.0	3.94	26.72	77.78	2019
Merlot	R110	16.4 ± 13.1	3.75	24.87	75	2018
Manto Negro	R110	11.9 ± 3.1	3.75	16.37		2018
Macabeo (Viura)	R110	13.0 ± 9.6	3.72	19.54	70.83	2019
Pinot Noir	R110	10.3 ± 7.1	3.67	13.42	44.44	2020
Tempranillo Blanco	R110	11.2 ± 8.4	3.56	17.31	55.56	2020
Syrah	R110	9.0 ± 5.6	3.50	12.37	100	2018
Verdejo	R110	9.5 ± 6.6	3.50	13.89	72.22	2019
Airen	R110	8.7 ± 5.3	3.44	15.98		2020
Monastrell	R110	12.3 ± 8.8	3.39	20.33	67.78	2019
Mencia	R110	12.4 ± 9.8	3.39	18.28	61.11	2019
Cabernet	R110	8.2 ± 6.0	3.33	11.56	72.22	2019
Tempranillo	SO4	8.0 ± 4.9	3.33	20.44	44.44	2020
Garnacha tintorera	R110	9.2 ± 7.6	3.33	15.68	55.56	2020
Chardonnay	R110	7.0 ± 3.6	3.33	8.45	66.67	2020
Viura	R140	10.5 ± 10.0	3.20	14.10	70	2018
Pedro Ximénez	R110	7.5 ± 4.9	3.17	12.67	66.67	2019
Garnacha	R110	9.2 ± 8.2	3.11	16.77	33.33	2020
Pinot Noir	R110	8.4 ± 5.0	3.00	11.28	61.11	2019

Continued on next page

Table C.1: Inoculation tests on different grapevine varieties (Continued)

Cabernet Sauvignon	R110	6.6 ± 4.5	2.89	7.22	77.78	2020
Syrah	R140	12.2 ± 10.0	2.89	15.44	55.56	2018
Hondarrabi zuri	SO4	6.9 ± 6.3	2.78	12.78	61.11	2019
Chardonnay	R110	8.4 ± 6.1	2.62	13.50	75	2018
Graciano	R140	6.7 ± 5.7	2.56	8.12	55.56	2020
Tempranillo	R140	8.1 ± 8.8	2.50	19.50	50	2018
Tempranillo	R110 RJ78	5.3 ± 5.2	2.44	11.88	55.56	2020
Prensal	R110	5.4 ± 5.5	2.37	9.25		2018
Tempranillo	SO4	9.7 ± 11.2	2.37	11.87	50	2018
Giró Ros	161/49	6.1 ± 7.6	2.29	7.71		2018
Giró Negre	R110	5.2 ± 5.7	2.25	8.87	62.5	2018
Viognier	R110	4.7 ± 5.0	2.25	6.50	75	2018
Callet	R110	7.9 ± 9.8	2.25	8.87	37.5	2018
Tempranillo	R110	3.9 ± 3.6	2.11	9.22	0.00	2020
Hondarrabi beltza	196-17	3.1 ± 1.8	1.89	5.89	11.11	2020
Tempranillo	R110	4.0 ± 4.9	1.75	7.62	25	2018
Argamussa	R110	3.2 ± 4.5	1.62	4.25		2018
Tempranillo	41B	5.6 ± 11.1	1.25	7.37	25	2018
Albariño	R110	2.1 ± 2.1	1.22	5.23	33.33	2020
Vinater Blanc	R110	2.0 ± 3.5	1.00	3.75	25	2018
Hondarrabi zuri	SO4	1.6 ± 1.3	1	5.78	0.00	2020
Cabernet	R110	2.9 ± 6.7	0.87	3.00	25	2018
Syrah	41B	1.7 ± 3.9	0.87	3.75	12.5	2018
Esperó de Gall	R110	4.0 ± 10.9	0.75	4.62	12.5	2018
Sauvignon Blanc	SO4	0.5 ± 0.5	0.50	2.50	0	2018
Giró Ros	R110	0.6 ± 1.9	0.37	1.50	0	2018
Mancés	R110	0.4 ± 0.7	0.25	1.25		2018



**Figure C.1: Factors influencing *Xf-Philaenus spumarius-Vitis vinifera* pathosystem.** (a) Virulence differences between *Xf* subsp. *fastidiosa* isolates on grapevines. Bars represent the mean number of symptomatic infected leaves four months after inoculating. Both isolates XYL2055/17 ( $n = 316$  inoculated plants) and XYL2177/18 ( $n = 260$ ) were collected from vineyards on Majorca. Scores were pooled among the 21 varieties inoculated; (b) conceptual graph of the population dynamics of *P. spumarius* on vineyards in Majorca and the effect on winter curing. Blue: density function of *P. spumarius*; red line: proportion of *P. spumarius* carrying *Xf*; blue line: proportion of plants recovering according to the time they are infected; (c) bimodal density function of the number of symptomatic leaves. The blue dash line marks 5 symptomatic leaves; and (d) correlation between the upward and downward movement of  $Xf_{PD}$  within the canes from the inoculation point. Each point depicts the mean distance travelled in both directions by the bacteria.

Overall, European *V. vinifera* varieties exhibited significant differences in

their susceptibility to  $Xf_{PD}$ , which could imply differences in risk of PD establishment at the regional scale (Table C.1). When compared between grape major phenotypic groups, red grape varieties were 1.45 times more prone to  $Xf_{PD}$  infection than white grape varieties ( $\chi^2 = 41.58$ ,  $df = 1$ ,  $P = 1.072 \times 10^{-10}$ ), while symptoms were 36.7% more severe in red grapes than in white grape cultivars ( $\chi^2 = 554.54$ ,  $df = 1$ ,  $P = 2.2 \times 10^{-16}$ ). In addition, we probed whether  $Xf_{PD}$  strains isolated from grapevines in Majorca differ in their virulence pooled across all grapevine varieties, finding significant differences in virulence ( $\chi^2 = 68.73$ ,  $df = 1$ ,  $P = 2.2 \times 10^{-16}$ ) and infectiveness ( $\chi^2 = 8.07$ ,  $df = 1$ ,  $P = 0.0045$ ) (Fig. C.1).



**Figure C.2: Experimental setup.** Greenhouse facilities and general view of its interior and the arrangement of the vine plants. The metallic structure is covered with an anti-thrips mesh.

Early-season  $Xf$ -infections on grapevines are considered to be more likely to survive the following year than late-season infections [323, 324]. By contrast, varieties developing symptoms, later on, may affect pathogen acquisition efficiency by vectors and thus decrease the rate of disease transmission. We found a positive correlation ( $F_{1,28} = 39.58$ ,  $P < 0.001$ ;  $R^2 = 0.57$ ) between the number of symptomatic leaves formed above the point of inoculation and those formed below (Fig. C.1). This acropetal/basipetal ratio of infected leaves is indicative of systemic movement of the pathogen and of a greater probability that infections on vines showing a lower number of symptomatic leaves will be more likely eliminated by winter pruning or by low temperatures [609]. As

a result, we assumed in our model that Xf-infected plants that develop fewer symptomatic leaves at the end of 16 weeks of incubation will contribute less to the spread of the disease within vineyards (Fig. C.1).

## C.2 Modeling climate suitability for PD

### C.2.1 Modified Growing Degree Days (MGDD) from Arrhenius Equation

Feil and Purcell estimated Xf growth rate as a function of temperature,  $\sigma(T)$ , using Arrhenius' Law, i.e.,  $\ln K \sim -1/T$  (see Fig. 3 in [323]). The overall dependence on  $T$  is nonmonotonic with two different types of behavior: i)  $k$  grows with  $T$  until a maximum value is attained at  $T = 28^\circ\text{C} = 301.15\text{K}$ ; and ii)  $k$  decreases beyond the maximum. Growth is zero beyond the lowest and highest threshold temperatures.

The mathematical form of the Arrhenius' Law dependence between the growth rate  $k$  and the absolute temperature  $T$  reads as follows,

$$k = A \exp(-E/T) , \quad (\text{C.1})$$

where  $A$  is a pre-exponential factor and  $E$  an activation energy in units of the Boltzmann constant  $k_B$ . The original use of this equation is for the rate constant of a chemical reaction that increases monotonically with  $T$ , and so  $E > 0$ . To fit the non-monotonic whole growth behavior of Xf, we considered two Arrhenius functions with opposite signs in the activation rate,

$$k = A_1 \exp(-E_1/T) + A_2 \exp(+E_2/T) , \quad (\text{C.2})$$

where  $E_1 > 0$  and  $E_2 < 0$ .

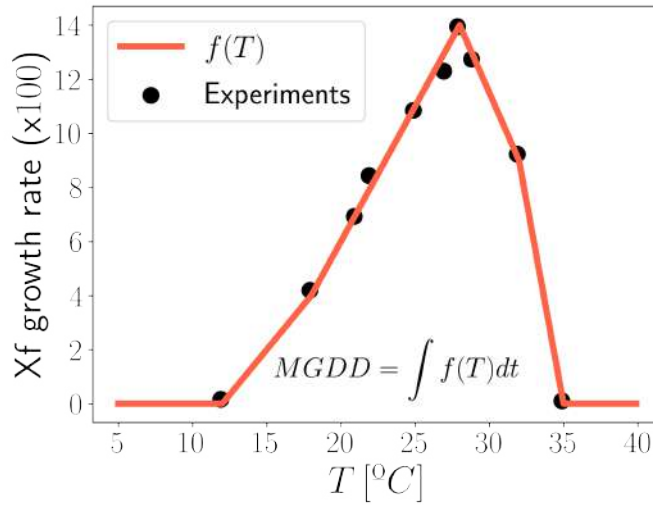
Now let us denote by  $t$ , the temperature in Celsius,  $t = T - 273.15$ . Importantly, within the bacterial temperature growth range ( $10\text{-}36^\circ\text{C}$ ) in the Arrhenius equation Eq. (C.2),  $t$  is quite small respect to  $b = 273.15$ , the absolute (Kelvin) temperature corresponding to  $0^\circ\text{C}$  in  $T = 273.15 + t$ . The two exponents in Eq. (C.2) can be approximated now as,

$$\begin{aligned} k &= A \exp\left(-\frac{E}{b+t}\right) = A \exp\left(-\frac{E}{b(1+t/b)}\right) \approx A \exp\left[-\frac{E}{b}\left(1-\frac{t}{b}\right)\right] = \\ &A \exp\left(-\frac{E}{b}\right) \exp\left(\frac{E}{b^2}t\right) \approx A \exp\left[-\frac{E}{b}\right] \left(1 + \frac{E}{b^2}t\right) = \\ &A \exp\left(-\frac{E}{b}\right) + A \frac{E}{b^2} \exp\left(-\frac{E}{b}\right) t = B + Ct , \quad (\text{C.3}) \end{aligned}$$

where we assume that  $t/b = t/273.15 \ll 1$  and  $(Et/(b^2)) = (Et)/(273.15^2) \ll 1$ , whereas  $B$  and  $C$  are constants. In particular,  $C > 0$  if  $E > 0$  fits the region

before the maximum in which the growth rate increases, while  $C < 0$  if  $E < 0$  fits the region after the maximum where  $k$  decreases. The positive/negative sign stems from the coefficient of the linear term in  $t$ ,  $E/b^2$ .

Each exponential in Eq. (C.2) can be expressed with a simple straight line, valid in our temperature range. This approach can be extended by adding more exponential terms in Eq. (C.2) to further improve the fit with a multilinear dependence between  $Xf_{PD}$  growth rate and temperature, obtaining a function proportional to  $Xf_{PD}$  growth rate  $F(T) = C \cdot \sigma(T)$  (see Fig. C.3).



**Figure C.3: Relationship between MGDD and temperature.** Contribution to the  $MGDD$  resulting from the fitting to the data in (1). The original Arrhenius plot,  $\log k$  vs.  $1/T$  in kelvin was converted to a linear dependence in Celsius temperature  $t$  (cf. Appendix C.2.1)

Now, this multilinear fit to the  $Xf_{PD}$  growth rate can be used to redefine the classical Growing Degree-Days (GDD) metric into the new Modified Growing Degree Days (MGDD).  $GDDs$  are computed as the integral of a particular function of temperature

$$GDD(t) = \int_{t_0}^t f(T(t))dt \tag{C.4}$$

where  $f(T)$  is defined as

$$f(T(t)) = \begin{cases} T(t) - T_{base} & \text{if } T \geq T_{base} \\ 0 & \text{if } T < T_{base} \end{cases} \tag{C.5}$$

Considering different slopes relating  $X_{f_{PD}}$  growth rate and temperature at different temperature intervals, as shown in Fig. C.3, we modified this particular function to now account for  $X_{f_{PD}}$  growth,

$$MGDD(t) = \int_{t_0}^t F(T(t))dt = C \cdot \int_{t_0}^t \sigma(T(t))dt \quad (C.6)$$

We wish to stress that the use of a multilinear form to represent MGDDs Fig. C.3 stems from the fundamental temperature dependence of the kinetics of bacterial growth as described by the Arrhenius equation, and is not an arbitrary simplified representation. Moreover, the MGDD function is fitted using the whole set of data published in [323], and is not simply based on the knowledge of the cardinal temperatures, as customarily done when writing smooth interpolating functions with the sole input of the cardinal temperatures (see, e.g., [340]).

## C.2.2 Relation between MGDD and within-plant bacterial population

The usual growth cycle of bacteria consists of several phases (lag, exponential, stationary and death phase), being of most interest to environmental microbiologists the interval between the lag and the onset of the stationary phase [610]. During the exponential phase, the rate of increase of cells is proportional to the number of cells present at any particular time. Thus, the evolution of the bacterial population,  $N$ , over time is given by the following differential equation,

$$\frac{dN}{dt} = \sigma N \implies N(t) = N_0 \cdot \exp(\sigma t) , \quad (C.7)$$

where  $\sigma$  is the specific growth rate constant.

As shown in the previous section, the growth rate of  $X_f$  has specific temperature dependence,  $\sigma(T)$ . In our study, temperature varies over time, so we can write the growth rate as a time-dependent quantity,  $\sigma(T(t))$ . With this, the evolution of the bacterial population will be given by

$$N(t) = N_0 \exp\left(\int_{t_0}^{t_f} \sigma(T(t))dt\right) . \quad (C.8)$$

Recalling Eq. (C.6) we can write the previous equation as

$$N(t) = \frac{N_0}{C} \cdot \exp(MGDD(t)) = C' \cdot N_0 \exp(MGDD(t)) \quad (C.9)$$

Indeed, the same can be done considering the logistic differential equation (that includes the stationary phase),

$$\frac{dN}{dt} = \sigma(T(t)) \cdot N \cdot \left(1 - \frac{N}{K}\right) \quad (C.10)$$

whose solution is

$$N(t) = \frac{K}{1 + C \exp\left(-\int_{t_0}^{t_f} \sigma(T(t))\right)} \quad (\text{C.11})$$

and using Eq. (C.6) it can be rewritten as

$$N(t) = \frac{K}{1 + C' \exp(-MGDD(t))} \quad (\text{C.12})$$

and thereby the bacterial population after a given time  $t$  is related to the *MGDD* by Eq. (C.12).

**Note:** We are assuming a correspondence between *in vitro* and *in planta* growth rates of Xf.

### C.2.3 Epidemiological and theoretical basis

A standard SIR model was considered as a basis to assess the risk of PD outbreaks worldwide (see Appendix C.4 for an analytical derivation of the relation between a vector-borne disease model and a standard SIR model). The model is represented by the following three equations,

$$\begin{aligned} \dot{S} &= -\beta SI/N \\ \dot{I} &= \beta SI/N - \gamma I \\ \dot{R} &= \gamma I, \end{aligned} \quad (\text{C.13})$$

where  $S$  is the susceptible host population,  $I$  is the infected population,  $R$  is the dead population and the total population  $N$  is conserved,  $S + I + R = N$ , as hosts die only when they contract the disease. The transmission of the disease from infected hosts to susceptible ones is mediated by the *transmission rate*  $\beta$  while the death of infected individuals is regulated by the *mortality rate*  $\gamma$ .

Analyzing the non-trivial fixed point,  $\vec{x} = (N, 0, 0)$ , it can be proved the existence of an epidemic threshold. As  $S$  is a monotonically decreasing function, which implies  $S(t) < S_0$ , one can write the following relation,

$$\frac{dI}{dt} = I(\beta S/N - \gamma) \leq I(\beta N/N - \gamma) = \gamma I(\beta/\gamma - 1) = \gamma I(R_0 - 1), \quad (\text{C.14})$$

where  $R_0 = \beta/\gamma$ . Thus,  $R_0 < 1$  implies  $dI/dt < 0 \forall t$  and  $I_0 > I(t)$  as  $t \rightarrow \infty$ , basically meaning that the epidemic dies out, while for  $R_0 > 1$ ,  $I(t)$  grows initially until  $S(t_c) = \gamma/\beta$ , at which  $\dot{I}(t_c) = 0$ , and the epidemic starts waning out.  $R_0$  corresponds to the so-called *basic reproduction number* and measures the number of secondary infections given by a primary infection in a fully susceptible population.

We wish to model the risk of PD establishment in a susceptible (healthy) population. For this, we characterized the maximum growth rate of the epidemic, when  $S(t) \sim S(0)$ . Thus, the growth is well approximated under these conditions with the (linearised) differential equation,

$$dI/dt = \beta SI - \gamma I \approx \gamma I(\beta N/\gamma - 1) = \gamma I(R_0 - 1) . \quad (\text{C.15})$$

where we have assumed the initial conditions,  $S_0 \approx N$ ,  $I(0) \approx 0$  and  $R(0) = 0$ . This linear differential equation can be integrated exactly,

$$I(t) = I(0) \exp(\gamma(R_0 - 1)t) . \quad (\text{C.16})$$

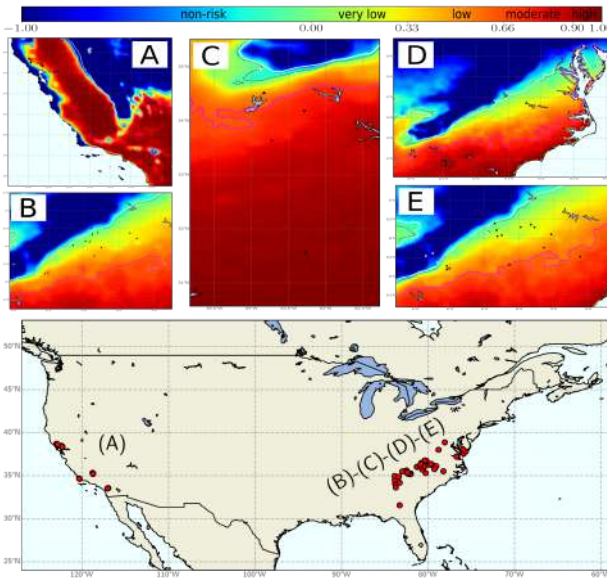
As explained in the main text, to account for the effect of temperature in the epidemic process we modify the previous expression as follows

$$I(t) = I(0) \exp(\gamma(R_0 - 1)t) \cdot \mathcal{F}(MGDD(t)) \cdot \mathcal{G}(CDD(t)) = I(0) \exp(\gamma(R_0 - 1)t) \cdot \Pi(t) , \quad (\text{C.17})$$

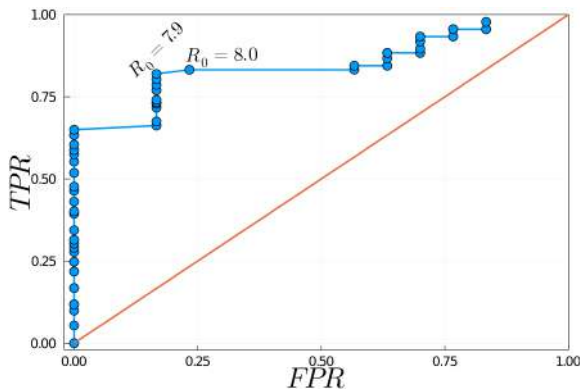
where  $\Pi(t) = \mathcal{F}(MGDD(t)) \cdot \mathcal{G}(CDD(t))$  is the cumulative probability of chronic infection that depends on temperature.

### C.2.4 Model validation

We ran several simulations of the model Eq. (7.7) with  $R_0$  values between 1 and 14 to validate PD spatio-temporal distribution in the US. We found  $R_0 = 8$  as the optimal parameter for maximizing the area under a ROC curve (Fig. C.5), returning an accuracy of more than 80%, except for 2006, due to data obtained from an area at the transient-risk zone (Fig. C.4). The ROC curve is a graphical representation of the trade-off between the true positive rate (TPR) and the false positive rate (FPR) for every possible cut-off of the model. The area under the ROC curve is a measure of the model's accuracy, with a value of 1 indicating perfect accuracy and 0.5 indicating no better than random accuracy. The ROC curve for  $R_0 = 8$  shows that the model is accurate in predicting the presence of PD in the US, with an area under the curve of 0.85.



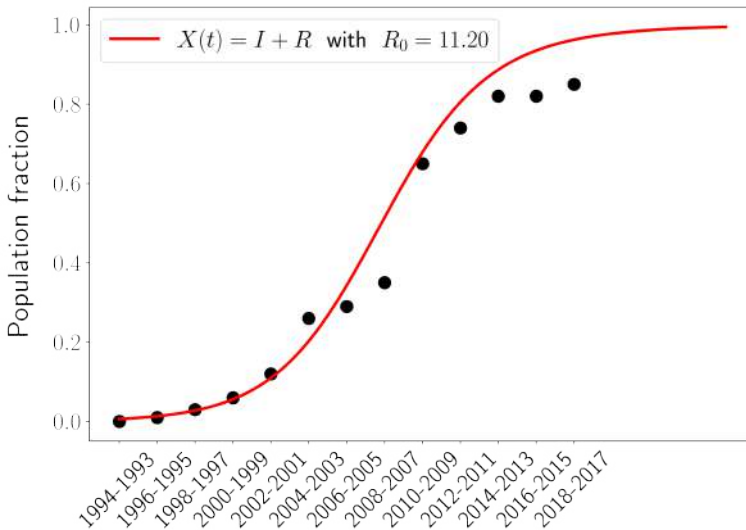
**Figure C.4: Model validation for an  $R_0 = 8$  scenario with presence/absence data (black/white stars) of PD in the United States.** Panel (A) corresponds to data from California in 2015 while the other panels show data from 2002, 2005, 2006 and 2001 (respectively) in the east of the United States. The last panel clarifies the validation zones previously mentioned.



**Figure C.5: ROC curve illustrating the model validation procedure with spatio-temporal data from PD distribution in the US.** TPR is the *true positive rate* and FPR the *false positive rate*. The model accuracy reaches its optimum in  $R_0 \approx 8$  by maximizing the true positive rate and minimizing the false positive rate. The spatio-temporal PD distribution in the US was obtained from data collected from publications between 2001 and 2015.

### C.2.5 Determination of $R_0$ for Europe

Unlike the validation of our model based on the distribution of PD in the USA, there are no spatio-temporal data on PD outbreaks available in Europe to estimate  $R_0$ . One way to solve this problem is to use data on the incidence of almond leaf scorch disease in Majorca to fit a SIR model and obtain an approximation of  $R_0$ . The initial date of introduction and progression of the almond leaf scorch epidemics is well characterized, and both diseases are transmitted by *P. spumarius* [265, 268].



**Figure C.6: Fitting a SIR model to the progress of the almond leaf scorch disease in the Balearic Islands from 1993 onward.** The best match was obtained with  $R_0 = 11.2$ . Points represent an estimate of the proportion of infected trees (incidence) from dendrochronological analysis and detection of Xf DNA in growth rings by qPCR. The incidence of ALSD in 2012 and 2017 was independently validated by field and Google Map Street View image observations.

Using  $\gamma = 1/14\text{years}^{-1}$  as the mortality rate [265], the best fit was provided by  $\beta_{\text{opt}} = 0.8$ , giving rise to  $R_0 = 11.2$  (Fig. C.6), which is in good agreement with the order of magnitude of  $R_0 = 8$  in the United States (Fig. C.5). To find a proper scenario for PD in Europe, we considered a constant transmission rate  $\beta_{\text{opt}}^1$  and applied an average mortality rate  $\gamma \sim 1/5\text{years}^{-1}$  of PD infected vines

<sup>1</sup>As the vector that transmits both ALSD and PD is the same (*Philaenus spumarius*) we considered that the transmission rate should not vary much.

[322], which gives rise to  $R_0 = 4$ . Finally, we used the information on the climate suitability for the vector in Majorca ( $\approx 0.8$  on average) to determine a baseline scenario for Europe,  $R_0 = 4/0.8 = 5$ . Thus, we can argue that  $R_0 = 5$  is a good proxy for modelling the establishment of PD in Europe. The use of  $R_0 = 5$  is furthermore corroborated by the reasonability of the predictions obtained.

### C.2.6 Simulation details

To assess the risk of  $Xf$  establishment in vineyards, we performed spatio-temporal simulations for the world's largest wine-growing areas. The cell size of the abstract grid was determined by the resolution of the data collected from ERA5-Land,  $0.1^\circ \times 0.1^\circ$ , so the spatial resolution is approximately 9km in the latitudes of the Mediterranean basin. A small initial infected-plant population was introduced annually into each cell assuming that if the conditions are or become favorable the disease will propagate locally. We chose  $I_i(0) = 1$  to rescale the results to any initial population size, and implemented Eq. (C.17) in each cell. Simulation time was discretized in years and computed in two steps incorporating summer ( $F(MGDD)$ ) and winter ( $F'(CDD)$ ) periods. To implement Eq. (C.17) we took into account that the  $MGDD$  and  $CDD$  differ at each time step; thereby it required to convert Eq. (C.16) into a mathematical map. The equation can be expressed as,

$$I(t) = I(0) \exp(\gamma(R_0 - 1)t) = I(0) [\exp(\gamma(R_0 - 1))]^t, \quad (\text{C.18})$$

where  $t$  is the discrete-time in years, so that

$$I(t - 1) = I(0) [\exp(\gamma(R_0 - 1))]^{t-1}, \quad (\text{C.19})$$

and, thus,

$$I(t) = I(t - 1) \exp(\gamma(R_0 - 1)). \quad (\text{C.20})$$

The discretized form of Eq. (C.17) is then

$$I(t_i) = I(t_{i-1}) \exp(\gamma(R_0 - 1)) \cdot F(MGDD(t_i)) \cdot F'(CDD(t_i)). \quad (\text{C.21})$$

A risk index was created to represent the relative velocity of PD local exponential propagation,

$$r(\tau) = \max \left\{ \frac{\log(I(\tau)/I(0))}{\gamma(R_0 - 1)\tau}, -1 \right\}, \quad (\text{C.22})$$

where  $\tau$  is the simulated time,  $R_0$  is the basic reproduction number and  $I(0)$  the initial condition (initial number of infected plants). The index ranges from -1 to 1 as the maximum risk value always occurs under optimal climatic conditions

( $F(MGDD) = F(CDD) = 1$ ) and thus  $I(\tau) = I(0) \exp((R_0 - 1)\tau)$ . The minimum risk was intentionally cut off at -1 to use a symmetric scale, as otherwise, the logarithmic scale is unbounded.

The numerator of the risk index defined in Eq. (C.22) is formally similar to the definition of Lyapunov exponents (LEs), which characterize predictability in chaotic systems (the denominator normalizes this quantity to its maximum value). This is not surprising because both the risk of Eq. (C.22) and the growth of perturbations in chaotic systems correspond to an exponential process. Following this analogy, we would expect a growing exponential process in the risk of the establishment if  $r > 0$ , while a decreasing exponential that goes to 0 would denote no risk if  $r < 0$ . However, Lyapunov exponents are (normally) calculated for autonomous (i.e. unforced, and so *steady*) dynamical systems, while Eq. (C.21) has 2 forcing terms (i.e., is non-autonomous). The result is a non-exponential behavior found when  $|r|$  is small. So beyond the expected regions with growing exponential and decreasing exponential behavior, we find a transition zone, where the system is oscillatory and not exponential, as neither growth in more auspicious years for  $X_{f_{PD}}$  or decrease in less auspicious ones prevails, and neither of the growing or decreasing pure exponential behaviors manifests.

We define the borderlines of this transition region by  $I(\tau) \leq 10 \cdot I(0)$  in the southern boundary and  $I(\tau) \geq 0.05 \cdot I(0)$  in the northern one when  $\tau = 40$  years. Basically, for the upper boundary, we assume that if an initial infection is multiplied by 10 after 40 years, then the exponential growth would be unstoppable. Conversely, if an initial introduction of infected individuals decays more than 95% of its original value after 40 years, we then assume that the exponential decay would continue and clearly PD cannot be established. Since  $\tau, \gamma$  are fixed, the limits of the transition zones depend on  $R_0$  and it is given by the risk index instead of the number of infected plants as follows,

$$\begin{aligned} \text{Upper limit: } r_{\text{trans}}^{\text{max}} &= \frac{\log(10)}{\gamma(R_0 - 1)\tau} \\ \text{Lower limit: } r_{\text{trans}}^{\text{min}} &= \frac{\log(0.05)}{\gamma(R_0 - 1)\tau} \end{aligned} \quad (\text{C.23})$$

For instance, with  $\gamma = 0.2 \text{ years}^{-1}$ ,  $\tau = 39$  years and  $R_0 = 5$  (values used for Europe) the transition zones are delimited by  $-0.09 < r(\tau) < 0.075$ . So, the model outputs can be associated with the following behaviors:

1. Epidemic-risk zones:  $r(\tau) > r_{\text{trans}}^{\text{max}}$ . The risk index  $r_j(\tau)$  is ranked as high ( $r_j(\tau) > 0.9$ ), moderate (0.9-0.66), low (0.66-0.33) and very low (0.33- $r_{\text{trans}}^{\text{max}}$ ).

## 2. Transition-risk zones:

$r_{\text{trans}}^{\min} < r(\tau) < r_{\text{trans}}^{\max}$ . In this zone the incidence,  $I(t)$ , predicted by the model does not grow clearly, but neither it does disappear, and incidence oscillates. This region is expected to be very sensitive to changes induced by climate change, and transit to epidemic-risk zones with low growth rates.

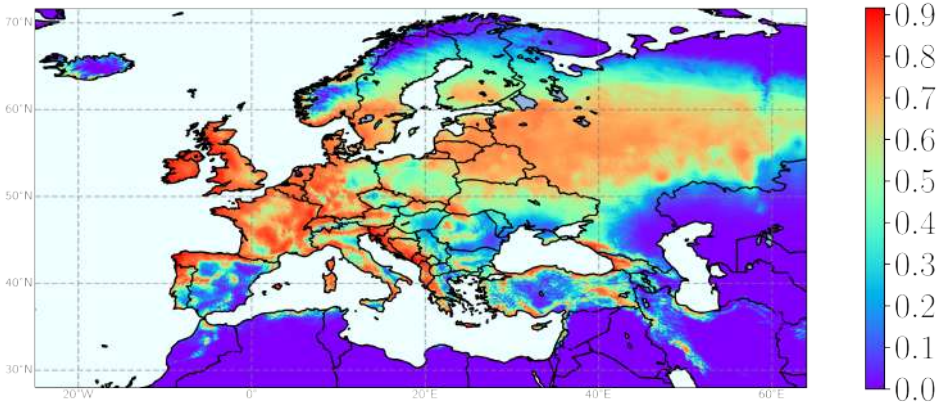
3. Non-risk zone:  $r(\tau) < r_{\text{trans}}^{\min}$ . Incidence decrease exponentially due to the combined effect of the *MGDD* and *CDD*, or to the (low) vector abundance in the case of predictions for Europe. Cells in this region with  $r_j$  not far from  $-0.1$  could become transitional due to the effect of climate change.

### C.2.7 Vector distribution influence

Information on the climatic suitability of the vector *P. spumarius* [320] was used to modulate the value of the basic reproduction number (Fig. C.7). We assumed a linear dependence of  $\beta$ , the transmission rate, with the vector climatic suitability resulting in each of the model cells,

$$R_0(x) = \frac{\beta v(x)}{\gamma} = R_0 \cdot v(x) , \quad (\text{C.24})$$

where  $x$  illustrates the space dependence.

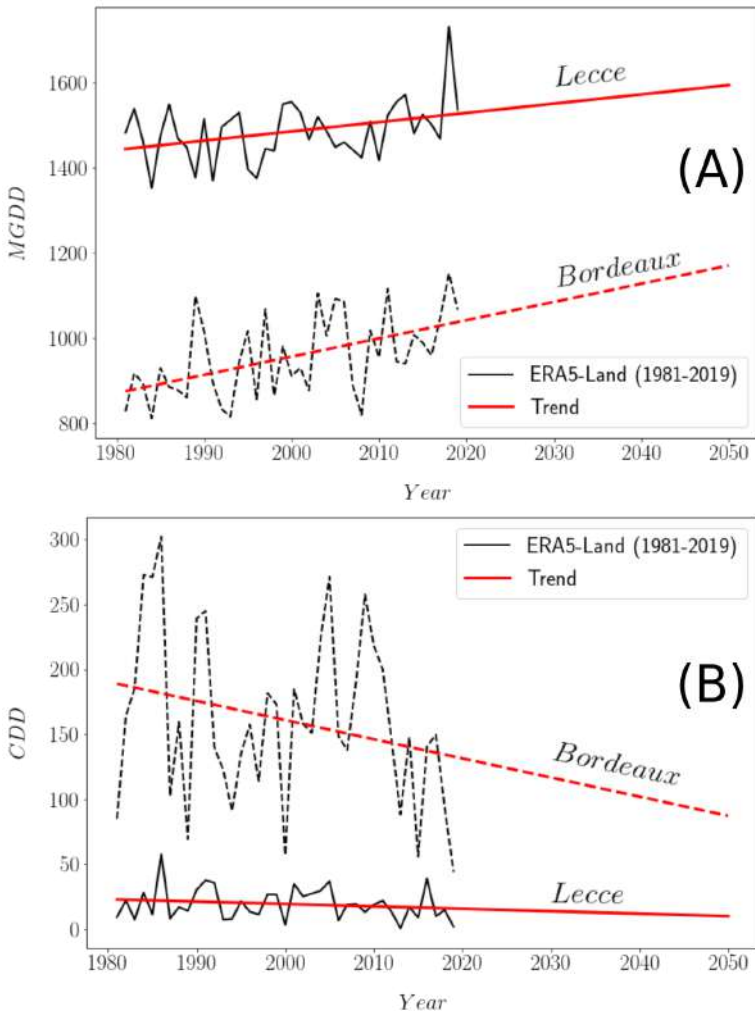


**Figure C.7: Average climatic suitability for *Philaenus spumarius* in Europe.** The map shows the climatic suitability of the vector estimated from a generalized additive model of insect distribution and the correlation of two bioclimatic descriptors, a climatic humidity index for the period of 8 coldest months of the year and the average maximum temperature in spring.

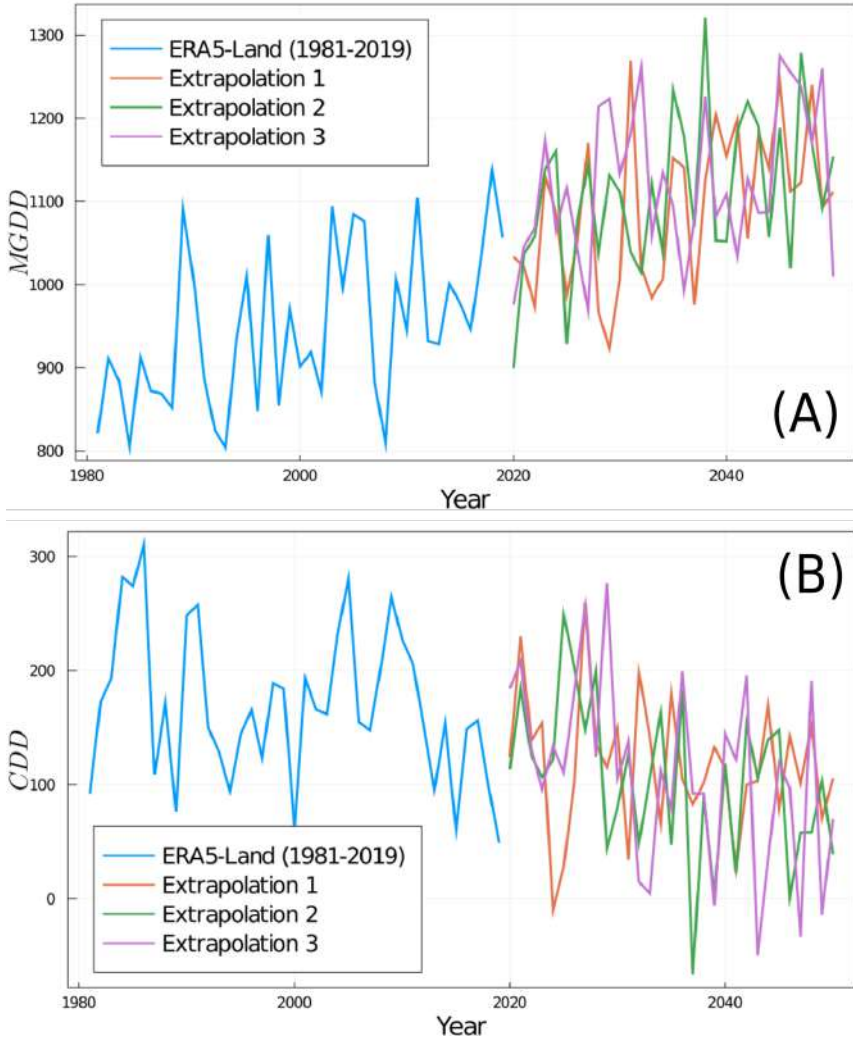
In Appendix C.4 we show an analytical derivation of the linear dependence between  $R_0$  and the vector population (i.e., the number of vectors). Then,

assuming that climatic suitability (i.e., probability of presence) is directly related to the number of vectors we obtain the linear scaling between  $R_0$  and climatic suitability for vectors.

### C.3 Future risk extrapolation



**Figure C.8: Determination of *MGDD* and *CDD* metric trends and future projections for two different European regions.** The *MGDD* (A) and *CDD* (B) trends show steeper slopes in the temperate climate of Bordeaux than in the Mediterranean climate of Lecce.



**Figure C.9: Interannual climatic variability extrapolations of  $MGDD$  (A) and  $CDD$  (B) for Bordeaux.** A linear model was fitted using Sklearn's LinearRegression module in Python and the interannual climatic variability was included as a Gaussian noise distribution by calculating the mean and fluctuations of the variance around  $MGDD$  and  $CDD$  trends.

To project PD risk in a climate change scenario, historical  $CDD$  and  $MGDD$  data were calculated to generate annual time series for each location recorded in the data set. To obtain the time trend of the variables in each pixel, a linear model was fitted using Sklearn's LinearRegression module in Python [611]. The interannual climatic variability was also included as a Gaussian noise distribution

by calculating the mean and fluctuations of the variance around the trend of the *MGDD* and *CDD* metrics for any record in the data set. We show in Fig. C.8 the determination of the trend of the metrics *MGDD* and *CDD* for Lecce and Bordeaux. Fig. C.9 shows three realizations to extrapolate the *MGDD* and *CDD* metrics for Bordeaux after applying Gaussian noise to the trend. This risk extrapolation to 2050 implies a linear extrapolation of past *MGDD* and *CDD* tendencies. Note that because *MGDD* and *CDD* functions are nonlinear this is just a rough approximation to the future risk, as nonlinearities could play a major role in a climate change scenario.

## C.4 Mathematical justifications

We show how a linear scaling between the vector population and the basic reproduction number can be obtained from a vector-borne disease model. Moreover, a SIR model can be derived from the same vector-borne disease model (under some assumptions).

In a model defined according to the following processes,

$$S_H + I_V \xrightarrow{\beta} I_H + I_V \quad I_H \xrightarrow{\gamma} R_H \quad S_V + I_H \xrightarrow{\alpha} I_V + I_H \quad S_V \xrightarrow{\mu} \emptyset \quad I_V \xrightarrow{\mu} \emptyset, \quad (\text{C.25})$$

where the birth of new susceptible vectors is described as a source term, the host-vector compartmental model can be written as,

$$\begin{aligned} \dot{S}_H &= -\beta S_H I_V / N_H \\ \dot{I}_H &= \beta S_H I_V / N_H - \gamma I_H \\ \dot{R}_H &= \gamma I_H \\ \dot{S}_V &= \delta C - \alpha S_V I_H / N_H - \mu S_V \\ \dot{I}_V &= \alpha S_V I_H / N_H - \mu I_V, \end{aligned} \quad (\text{C.26})$$

when a standard incidence [187] is considered.

The model describes the infection of susceptible hosts ( $S_H$ ) at a rate  $\beta$  through their interaction with infected vectors ( $I_V$ ), while susceptible vectors ( $S_V$ ) are infected at a rate  $\alpha$  through their interaction with infected hosts ( $I_H$ ). Infected hosts exit the infected compartment at a rate  $\gamma$ , while infected vectors stay infected for the rest of their life since they are not affected by the pathogen. The model assumes that vectors die naturally (or disappear from the population by some mechanism) at a rate  $\mu$  and are born (appear) at a constant rate  $\delta$ , being susceptible. The constant term  $C$  sets the scale of the stationary value of the vector population.

### C.4.1 Linear scaling of $R_0$ with vector population

The standard methods of calculation of  $R_0$  are based on the linear stability analysis of the disease-free equilibrium, either directly, through the linear analysis of the fixed point that yields the stability condition from which  $R_0$  can be obtained, or using the Next Generation Method (NGM) [68] that provides directly  $R_0$  by solving a suitable linear problem. The disease-free equilibrium of the model (the fixed point) is given by  $I_H = I_v = 0$  yielding  $\dot{S}_v = 0 \implies S_v = \delta C / \mu = N_v^*$ , where  $N_v^*$  is the stationary value of the vector population.

As shown in [289], both methods yield the following relation for the basic reproduction number,

$$R_0 = \frac{\beta \alpha}{\gamma \mu} \frac{C}{N_H} \frac{\delta}{\mu} \frac{S_H(0)}{N_H} = \frac{\beta \alpha N_v^*}{\gamma \mu} \frac{S_H(0)}{N_H}, \quad (\text{C.27})$$

in which the basic reproduction number scales linearly with the vector population.

### C.4.2 Reduction to a SIR model

In a timescale where the vector population changes faster than the host population (a good approximation for  $X_{\text{fPD}}$ -related diseases), the former will almost instantaneously reach the stationary value. Thus, if  $1/\mu \ll 1/\gamma$ , or equivalently if  $\mu \gg \gamma$ , we can rewrite the time derivative of the vector infected population as

$$\varepsilon \dot{I}_v = \frac{\alpha}{\mu} S_v \frac{I_H}{N_H} - I_v, \quad (\text{C.28})$$

with  $\varepsilon = 1/\mu$  being a small parameter. Then,  $\dot{I}_v$  can be neglected and the infected vector population can be obtained from the relationship,

$$I_v \approx \frac{\alpha S_v I_H}{\mu N_H}. \quad (\text{C.29})$$

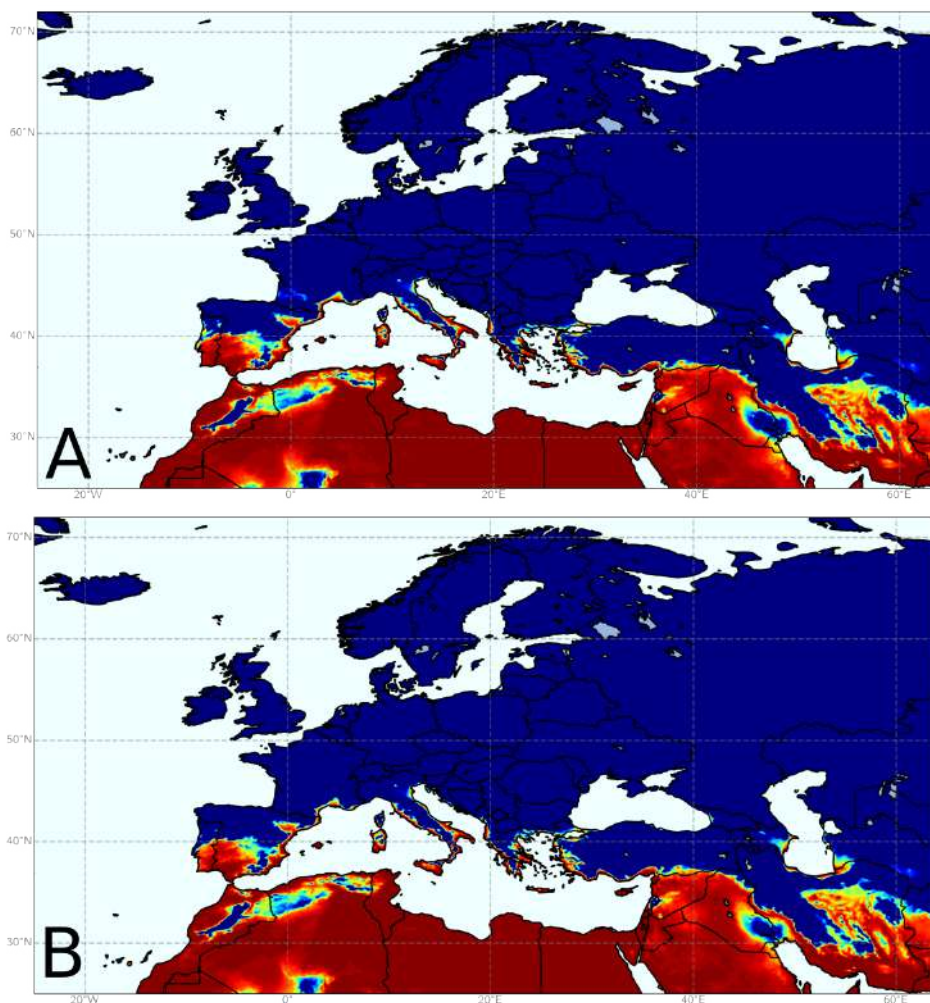
Furthermore, if  $\lambda N_H \gg I_H$  (which is indeed plausible in this limit) the model can be written as a SIR model with constant coefficients,

$$\begin{aligned} \dot{S}_H &= -\beta_{eff} \frac{S_H I_H}{N_H} \\ \dot{I}_H &= \beta_{eff} \frac{S_H I_H}{N_H} - \gamma I_H \\ \dot{R}_H &= \gamma I_H, \end{aligned} \quad (\text{C.30})$$

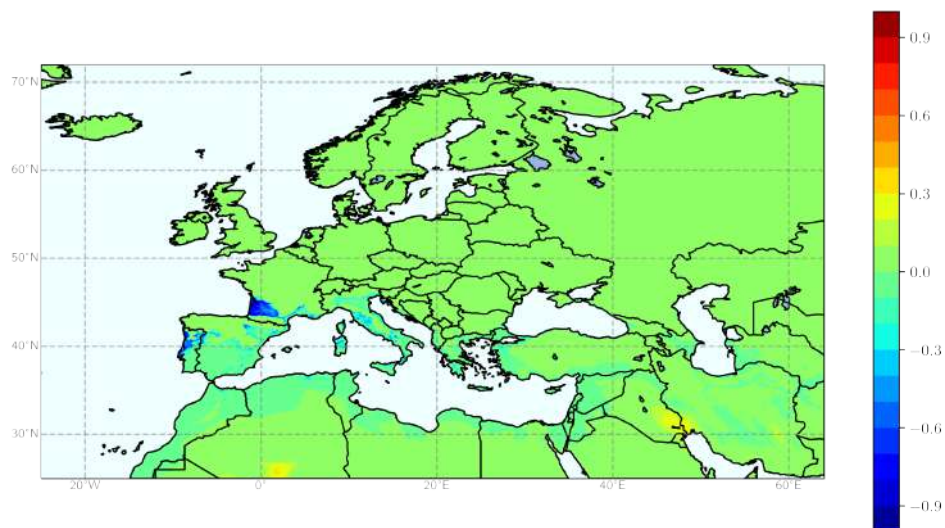
where  $\beta_{eff} = \frac{\beta'}{\lambda} = \frac{\beta \alpha N_v^*}{\mu N_H}$ .

Note that in the SIR model reduction, the effective  $\beta_{eff}$  coefficient depends linearly on the vector population  $N_v^*$ .

## C.5 Comparison of MGDD calculations from different models

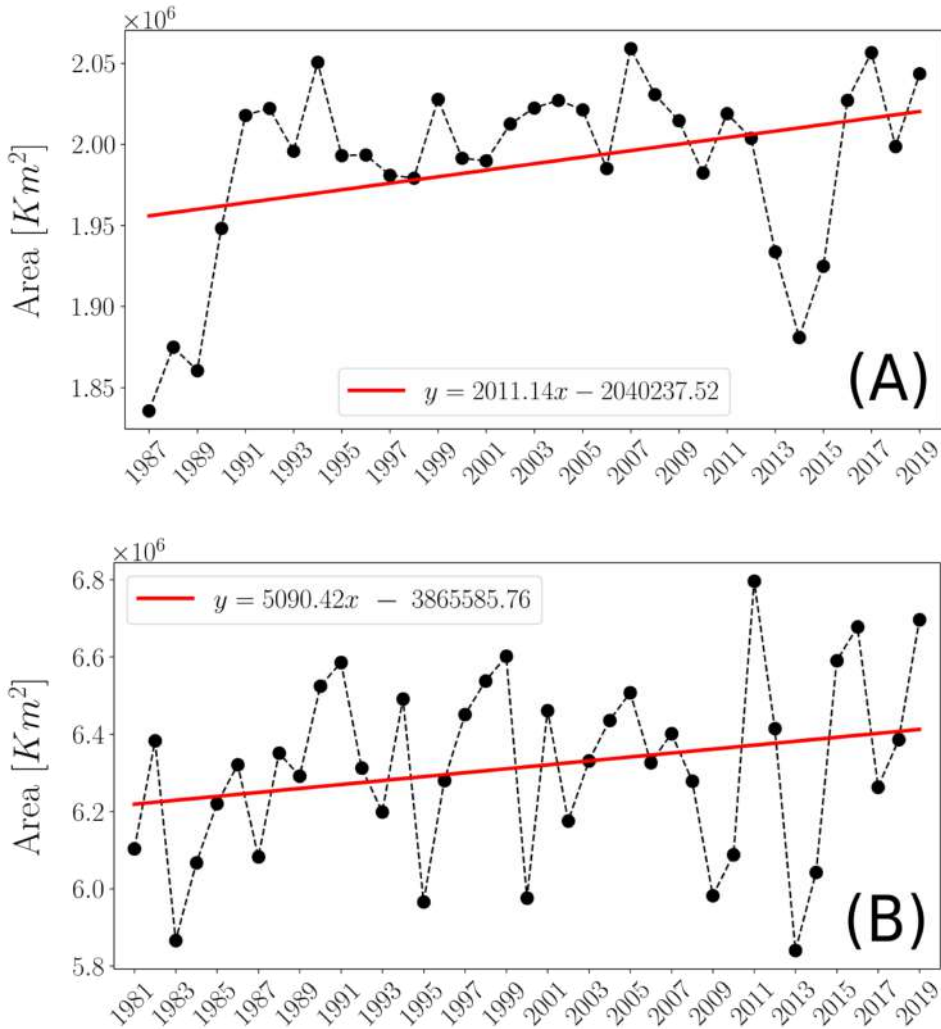


**Figure C.10:** Risk index computed with (A) MGDD from the Arrhenius-based fit and (B) MGDD from the beta function fit.

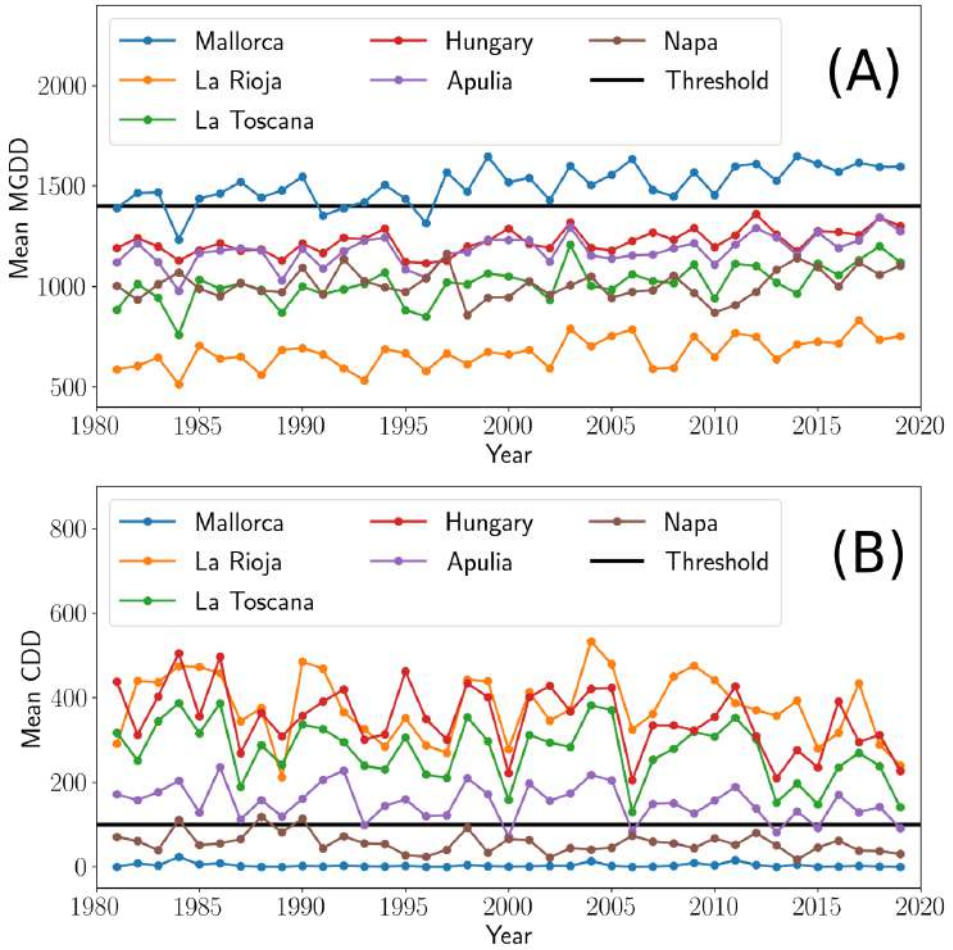


**Figure C.11:** Difference in risk index when computed using MGDD calculated from the Arrhenius-based fit or the beta function fit.

## C.6 Analysis of MGDD and CDD time series



**Figure C.12: Trends in the risk-epidemic zones during the 1981-2019 period (A) and the areas encompassed below the  $CDD < 314$  line (B) comprising land areas between  $103^\circ W$  and  $70^\circ W$  of the United States**



**Figure C.13: Trends in *MGDD* (A) and *CDD* (B) values and oscillations during 1981-2019 for seven wine regions with different climates from Europe and the US. *MGDDs* show a slight upward trend and lesser oscillations than the *CDDs*.**

## C.7 Detailed analysis of PD risk

Table C.2: **PD risk areas in Europe after running the model under a  $R_0 = 5$  scenario and a homogeneous spatial vector distribution.** The epidemic-risk zones are classified according to the relative disease growth rates defined by the risk index, as very low, low, moderate, and high growth rates. The total risk refers to the sum of the epidemic-risk zones

Country	No risk (km <sup>2</sup> )	Trans. (km <sup>2</sup> )	Very low (km <sup>2</sup> )	Low (km <sup>2</sup> )	Med. (km <sup>2</sup> )	High (km <sup>2</sup> )	Total risk (km <sup>2</sup> )	Total surf. (km <sup>2</sup> )	Risk (%)
Russia	4.22e6	1.57e4	1.70e3	0.0	0.0	0.0	1.70e3	4.24e6	0.0
Norway	3.25e5	0.0	0.0	0.0	0.0	0.0	0.0	3.25e5	0.0
France	4.51e5	3.62e4	4.05e4	8.33e3	7.02e3	1.63e3	5.74e4	5.44e5	10.6
Sweden	4.42e5	0.0	0.0	0.0	0.0	0.0	0.0	4.42e5	0.0
Belarus	2.08e5	0.0	0.0	0.0	0.0	0.0	0.0	2.08e5	0.0
Ukraine	5.69e5	0.0	0.0	0.0	0.0	0.0	0.0	5.69e5	0.0
Poland	3.12e5	0.0	0.0	0.0	0.0	0.0	0.0	3.12e5	0.0
Austria	8.33e4	0.0	0.0	0.0	0.0	0.0	0.0	8.33e4	0.0
Hungary	9.23e4	0.0	0.0	0.0	0.0	0.0	0.0	9.23e4	0.0
Moldova	3.21e4	0.0	0.0	0.0	0.0	0.0	0.0	3.21e4	0.0
Romania	2.33e5	1.93e3	0.0	0.0	0.0	0.0	0.0	2.35e5	0.0
Lithuania	6.46e4	0.0	0.0	0.0	0.0	0.0	0.0	6.46e4	0.0
Latvia	6.42e4	0.0	0.0	0.0	0.0	0.0	0.0	6.42e4	0.0
Estonia	4.55e4	0.0	0.0	0.0	0.0	0.0	0.0	4.55e4	0.0
Germany	3.55e5	0.0	0.0	0.0	0.0	0.0	0.0	3.55e5	0.0
Bulgaria	1.02e5	8.93e3	1.28e3	0.0	0.0	0.0	1.28e3	1.12e5	1.1
Greece	3.77e4	1.72e4	1.76e4	1.05e4	1.46e4	3.21e4	7.47e4	1.30e5	57.7
Albania	1.89e4	2.42e3	2.42e3	2.05e3	2.24e3	2.14e3	8.84e3	3.02e4	29.3
Croatia	4.57e4	2.13e3	3.27e3	1.68e3	1.06e3	2.70e2	6.27e3	5.41e4	11.6
Switzer.	4.62e4	0.0	0.0	0.0	0.0	0.0	0.0	4.62e4	0.0
Luxemb.	2.71e3	0.0	0.0	0.0	0.0	0.0	0.0	2.71e3	0.0
Belgium	3.05e4	0.0	0.0	0.0	0.0	0.0	0.0	3.05e4	0.0
Nether.	3.69e4	0.0	0.0	0.0	0.0	0.0	0.0	3.69e4	0.0
Portugal	1.42e4	1.42e4	1.38e4	5.61e3	2.48e4	1.58e4	6.00e4	8.84e4	67.9
Spain	2.04e5	4.16e4	5.15e4	5.32e4	6.81e4	8.52e4	2.58e5	5.04e5	51.2
Ireland	6.82e4	0.0	0.0	0.0	0.0	0.0	0.0	6.82e4	0.0
Italy	1.13e5	3.95e4	4.65e4	2.21e4	3.04e4	4.88e4	1.48e5	2.99e5	49.3

Continued on next page

Table C.2: PD risk areas in Europe with an homogeneous spatial vector distribution (Continued)

Denmark	4.23e4	0.0	0.0	0.0	0.0	0.0	0.0	4.23e4	0.0
UK	2.43e5	0.0	0.0	0.0	0.0	0.0	0.0	2.43e5	0.0
Iceland	1.07e5	0.0	0.0	0.0	0.0	0.0	0.0	1.07e5	0.0
Slovenia	2.02e4	2.58e2	8.64	0.0	0.0	0.0	8.64	2.06e4	0.4
Finland	3.29e5	0.0	0.0	0.0	0.0	0.0	0.0	3.29e5	0.0
Slovakia	4.81e4	0.0	0.0	0.0	0.0	0.0	0.0	4.81e4	0.0
Czechia	8.08e4	0.0	0.0	0.0	0.0	0.0	0.0	8.08e4	0.0
Bosnia	4.96e4	4.50e2	0.0	0.0	0.0	0.0	0.0	5.00e4	0.0
Macedonia	2.28e4	2.03e3	9.27	0.0	0.0	0.0	9.27	2.50e4	0.4
Serbia	7.60e4	0.0	0.0	0.0	0.0	0.0	0.0	7.60e4	0.0
Monteneg.	1.14e4	3.64e2	4.55e2	7.30e2	3.66e2	0.0	1.55e3	1.33e4	11.7
Kosovo	1.12e4	0.0	0.0	0.0	0.0	0.0	0.0	1.12e4	0.0
Cyprus	4.04e2	0.0	0.0	0.0	1.01e2	4.85e3	4.95e3	5.35e3	92.5
Czech Rep.	7.85e4	0.0	0.0	0.0	0.0	0.0	0.0	7.85e4	0.0
Malta	0.0	0.0	0.0	0.0	0.0	1.99e2	1.99e2	1.99e2	100.0
<b>TOTAL (%)</b>	9.21	1.8	1.8	1.0	1.5	1.9	6.1		

Table C.3: PD risk areas in the United States after running the model under a  $R_0 = 8$  scenario and using a homogeneous spatial vector distribution. The epidemic-risk zones are classified according to the relative disease growth rates defined by the risk index, as very low, low, moderate and high growth rates. The total risk refers to the sum of the epidemic-risk zones

State	No risk (km <sup>2</sup> )	Trans. (km <sup>2</sup> )	Very low (km <sup>2</sup> )	Low (km <sup>2</sup> )	Med. (km <sup>2</sup> )	High (km <sup>2</sup> )	Total risk (km <sup>2</sup> )	Total surf. (km <sup>2</sup> )	High Risk (%)
Maine	8.46e4	0.00	0.00	0.00	0.00	0.00	0.00	8.46e4	0.00
Massachus.	2.08e4	0.00	0.00	0.00	0.00	0.00	0.00	2.08e4	0.00
Michigan	1.50e5	0.00	0.00	0.00	0.00	0.00	0.00	1.50e5	0.00
Montana	3.74e5	0.00	0.00	0.00	0.00	0.00	0.00	3.74e5	0.00
Nevada	2.39e5	1.07e4	7.19e3	6.92e3	1.03e4	7.98e3	3.24e4	2.82e5	2.8

Continued on next page

Table C.3: PD risk areas in the US (Continued)

New Jersey	1.60e4	3.80e3	2.86e2	0.00	0.00	0.00	2.86e2	2.01e4	0.00
New York	1.28e5	0.00	0.00	0.00	0.00	0.00	0.00	1.28e5	0.00
North Carolina	1.56e4	5.71e3	1.63e4	4.39e4	4.02e4	7.50e3	1.08e5	1.29e5	5.8
Ohio	1.08e5	0.00	0.00	0.00	0.00	0.00	0.00	1.08e5	0.00
Pennsylvania	1.15e5	0.00	0.00	0.00	0.00	0.00	0.00	1.15e5	0.00
Rhode Island	2.67e3	0.00	0.00	0.00	0.00	0.00	0.00	2.67e3	0.00
Tennessee	3.38e3	1.24e4	6.59e4	2.87e4	0.00	0.00	9.46e4	1.10e5	0.00
Texas	3.90e3	3.36e4	3.58e4	5.26e4	1.97e5	3.62e5	6.47e5	6.84e5	5.28
Utah	2.12e5	3.63e3	1.18e3	2.95e2	0.00	0.00	1.47e3	2.17e5	0.00
Washington	1.75e5	1.71e2	0.00	0.00	0.00	0.00	0.00	1.75e5	0.00
Wisconsin	1.44e5	0.00	0.00	0.00	0.00	0.00	0.00	1.44e5	0.00
Puerto Rico	1.40e3	0.00	0.00	0.00	0.00	7.72e3	7.72e3	9.13e3	84.6
Maryland	1.46e4	5.46e3	7.04e3	2.91e2	0.00	0.00	7.34e3	2.74e4	0.00
Alabama	3.19e2	0.00	0.00	2.25e4	4.19e4	6.96e4	1.34e5	1.34e5	51.8
Alaska	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Arizona	8.96e4	3.59e4	1.68e4	1.69e4	3.76e4	9.83e4	1.70e5	2.95e5	33.3
Arkansas	0.00	8.36e3	3.27e4	4.68e4	4.60e4	0.00	1.26e5	1.34e5	0.00
California	1.31e5	1.62e4	2.23e4	2.48e4	6.15e4	1.54e5	2.62e5	4.09e5	37.5
Colorado	2.72e5	0.00	0.00	0.00	0.00	0.00	0.00	2.72e5	0.00
Connecticut	1.33e4	0.00	0.00	0.00	0.00	0.00	0.00	1.33e4	0.00
Delaware	9.50e2	2.10e3	2.31e3	0.00	0.00	0.00	2.31e3	5.36e3	0.00
District of Columbia	0.00	9.59	0.00	0.00	0.00	0.00	0.00	9.59	0.00
Florida	7.16e3	0.00	0.00	0.00	0.00	1.43e5	1.43e5	1.50e5	95.2
Georgia	5.25e2	2.02e2	3.74e3	1.24e4	3.58e4	9.89e4	1.51e5	1.52e5	65.3
Hawaii	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Idaho	2.18e5	0.00	0.00	0.00	0.00	0.00	0.00	2.18e5	0.00
Illinois	1.37e5	7.72e3	0.00	0.00	0.00	0.00	0.00	1.45e5	0.00
Indiana	9.17e4	6.81e2	0.00	0.00	0.00	0.00	0.00	9.24e4	0.00
Iowa	1.47e5	0.00	0.00	0.00	0.00	0.00	0.00	1.47e5	0.00
Kansas	2.03e5	1.43e4	0.00	0.00	0.00	0.00	0.00	2.17e5	0.00

Continued on next page

Table C.3: PD risk areas in the US (Continued)

Kentucky	3.48e4	5.89e4	1.03e4	0.00	0.00	0.00	1.03e4	1.04e5	0.00
Louisiana	6.74e3	0.00	0.00	0.00	7.97e3	1.08e5	1.16e5	1.23e5	88.0
Minnesota	2.16e5	0.00	0.00	0.00	0.00	0.00	0.00	2.16e5	0.00
Mississippi	4.25e2	0.00	0.00	1.50e4	5.42e4	5.26e4	1.22e5	1.22e5	43.0
Missouri	1.38e5	3.62e4	7.83e3	0.00	0.00	0.00	7.83e3	1.82e5	0.00
Nebraska	1.94e5	0.00	0.00	0.00	0.00	0.00	0.00	1.94e5	0.00
New Hampshire	2.38e4	0.00	0.00	0.00	0.00	0.00	0.00	2.38e4	0.00
New Mexico	1.57e5	3.18e4	4.13e4	4.10e4	3.95e4	0.0	1.22e5	3.11e5	0.00
North Dakota	1.80e5	0.0	0.0	0.0	0.0	0.0	0.0	1.80e5	0.0
Oklahoma	5.03e3	4.56e4	5.51e4	5.80e4	1.54e4	0.0	1.29e5	1.79e5	0.0
Oregon	2.50e5	6.01e2	0.0	0.0	0.0	0.0	0.0	2.50e5	0.0
South Carolina	6.21e2	0.0	5.04e2	6.16e3	3.90e4	3.41e4	7.97e4	8.04e4	42.4
South Dakota	2.01e5	0.0	0.0	0.0	0.0	0.0	0.0	2.01e5	0.0
Vermont	2.50e4	0.0	0.0	0.0	0.0	0.0	0.0	2.50e4	0.0
Virginia	4.42e4	1.70e4	2.87e4	1.48e4	7.88e2	0.0	4.42e4	1.05e5	0.0
West Virginia	6.24e4	3.89e2	0.0	0.0	0.0	0.0	0.0	6.28e4	0.0
Wyoming	2.53e5	0.0	0.0	0.0	0.0	0.0	0.0	2.53e5	0.0
<b>TOTAL (%)</b>	<b>63.10</b>	<b>4.50</b>	<b>4.60</b>	<b>5.00</b>	<b>8.10</b>	<b>14.70</b>	<b>32.40</b>		

Table C.4: **Potential distribution of PD in other world winegrowing regions.** In most areas of China and Australia *Vitis vinifera* is not cultivated and epidemic-risk zones with high growth rate correspond mainly to tropical areas in China, Australia, South Africa and Argentina

Country	No risk (km <sup>2</sup> )	Trans. (km <sup>2</sup> )	Very low (km <sup>2</sup> )	Low (km <sup>2</sup> )	Med. (km <sup>2</sup> )	High (km <sup>2</sup> )	Total risk (km <sup>2</sup> )	Total surf. (km <sup>2</sup> )
<i>China</i>	6.78e6	3.10e5	2.11e5	3.45e5	4.13e5	9.92e5	1.96e6	9.05e6

Continued on next page

Table C.4: Potential distribution of PD in other world winegrowing regions  
(Continued)

<i>Australia</i>	5.05e5	3.05e5	4.41e5	1.33e6	2.72e6	2.38e6	6.87e6	7.68e6
<i>South Africa</i>	2.17e5	1.84e5	1.20e5	1.52e5	2.79e5	2.64e5	8.15e5	1.22e6
<i>Argent.</i>	9.92e5	1.48e5	9.41e4	2.19e5	3.74e5	9.46e5	1.63e6	2.77e6
<i>Chile</i>	7.04e5	5.66e4	2.11e4	2.00e4	8.83e3	9.13e2	5.09e4	8.11e5

Table C.5: **Predicted PD risk areas for the US in 2050 considering a  $R_0 = 8$  scenario and a homogeneous spatial vector distribution.** The epidemic-risk zones are classified according to the relative disease growth rates defined by the risk index, as very low, low, moderate and high growth rates. The total risk refers to the sum of the epidemic-risk zones

State	No risk (km <sup>2</sup> )	Trans. (km <sup>2</sup> )	Very low (km <sup>2</sup> )	Low (km <sup>2</sup> )	Med. (km <sup>2</sup> )	High (km <sup>2</sup> )	Total risk (km <sup>2</sup> )	Total surf. (km <sup>2</sup> )	High Risk (%)
Maine	8.46e4	0.0	0.0	0.0	0.0	0.0	0.0	8.46e4	0.00
Massachus.	2.08e4	0.0	0.0	0.0	0.0	0.0	0.0	2.08e4	0.00
Michigan	1.50e5	0.0	0.0	0.0	0.0	0.0	0.0	1.50e5	0.00
Montana	3.74e5	0.0	0.0	0.0	0.0	0.0	0.0	3.74e5	0.00
Nevada	2.30e5	1.23e4	8.29e3	7.79e3	1.02e4	1.33e4	3.96e4	2.82e5	4.73
New Jersey	1.08e4	7.00e3	2.29e3	0.0	0.0	0.0	2.29e3	2.00e4	0.00
New York	1.26e5	1.12e3	0.0	0.0	0.0	0.0	0.0	1.28e5	0.00
North Carolina	1.03e4	4.40e3	8.61e3	2.81e4	6.58e4	1.20e4	1.15e5	1.29e5	9.26
Ohio	1.07e5	2.89e2	0.0	0.0	0.0	0.0	0.0	1.08e5	0.00
Pennsylvania	1.15e5	9.45	0.0	0.0	0.0	0.0	0.0	1.15e5	0.00
Rhode Island	2.67e3	0.0	0.0	0.0	0.0	0.0	0.0	2.67e3	0.00
Tennessee	2.29e3	1.09e3	2.38e4	7.97e4	3.53e3	0.0	1.07e5	1.10e5	0.00
Texas	2.71e3	0.0	5.21e4	4.66e4	1.62e5	4.21e5	6.82e5	6.84e5	6.15
Utah	2.10e5	2.65e3	3.33e3	7.86e2	9.84	0.0	4.22e3	2.17e5	0.00
Washington	1.66e5	8.24e3	1.37e3	0.0	0.0	0.0	1.37e3	1.75e5	0.00
Wisconsin	1.44e5	0.0	0.0	0.0	0.0	0.0	0.0	1.44e5	0.00

Continued on next page

Table C.5: Extrapolated PD risk areas in the US (Continued)

Puerto Rico	1.40e3	0.0	0.0	0.0	0.0	7.72e3	7.72e3	9.13e3	8.46
Maryland	1.16e4	4.00e3	8.65e3	3.10e3	0.0	0.0	1.17e4	2.74e4	0.00
Alabama	3.19e2	0.0	0.0	3.64e3	4.45e4	8.60e4	1.34e5	1.34e5	6.39
Alaska	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.00
Arizona	7.53e4	2.77e4	2.58e4	1.44e4	2.69e4	1.25e5	1.92e5	2.95e5	4.24
Arkansas	0.0	0.0	1.75e4	4.15e4	5.57e4	1.92e4	1.34e5	1.34e5	1.44
California	1.24e5	9.57e3	1.79e4	2.17e4	4.03e4	1.96e5	2.76e5	4.09e5	4.79
Colorado	2.55e5	1.52e4	2.06e3	0.0	0.0	0.0	2.06e3	2.72e5	0.00
Connecticut	1.33e4	0.0	0.0	0.0	0.0	0.0	0.0	1.33e4	0.00
Delaware	3.81e2	8.54e2	3.65e3	4.81e2	0.0	0.0	4.13e3	5.36e3	0.00
District of Columbia	0.0	9.59	0.0	0.0	0.0	0.0	0.0	9.59	0.00
Florida	7.16e3	0.0	0.0	0.0	0.0	1.43e5	1.43e5	1.50e5	9.52
Georgia	5.26e2	0.0	1.01e3	4.96e3	2.81e4	1.17e5	1.51e5	1.52e5	7.72
Hawaii	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.00
Idaho	2.13e5	4.48e3	0.0	0.0	0.0	0.0	0.0	2.18e5	0.00
Illinois	1.19e5	1.87e4	6.85e3	0.0	0.0	0.0	6.85e3	1.45e5	0.00
Indiana	8.41e4	8.33e3	0.0	0.0	0.0	0.0	0.0	9.24e4	0.00
Iowa	1.47e5	0.0	0.0	0.0	0.0	0.0	0.0	1.47e5	0.00
Kansas	1.51e5	4.56e4	2.07e4	0.0	0.0	0.0	2.07e4	2.17e5	0.00
Kentucky	1.54e4	3.53e4	5.22e4	1.09e3	0.0	0.0	5.33e4	1.04e5	0.00
Louisiana	6.74e3	0.0	0.0	0.0	0.0	1.16e5	1.16e5	1.23e5	9.45
Minnesota	2.16e5	0.0	0.0	0.0	0.0	0.0	0.0	2.16e5	0.00
Mississippi	4.25e2	0.0	0.0	5.05e2	4.31e4	7.83e4	1.22e5	1.22e5	6.40
Missouri	8.54e4	5.61e4	3.80e4	2.39e3	0.0	0.0	4.04e4	1.82e5	0.00
Nebraska	1.94e5	0.0	0.0	0.0	0.0	0.0	0.0	1.94e5	0.00
New Hampshire	2.38e4	0.0	0.0	0.0	0.0	0.0	0.0	2.38e4	0.00
New Mexico	1.36e5	2.39e4	4.30e4	4.45e4	5.13e4	1.18e4	1.51e5	3.11e5	3.80
North Dakota	1.80e5	0.0	0.0	0.0	0.0	0.0	0.0	1.80e5	0.00
Oklahoma	0.0	5.92e2	7.25e4	6.90e4	3.71e4	0.0	1.79e5	1.79e5	0.00
Oregon	2.46e5	2.59e3	1.89e3	0.0	0.0	0.0	1.89e3	2.51e5	0.00

Continued on next page

Table C.5: Extrapolated PD risk areas in the US (Continued)

South Carolina	6.21e2	0.0	0.0	1.01e3	2.03e4	5.85e4	7.97e4	8.04e4	7.28
South Dakota	2.01e5	0.0	0.0	0.0	0.0	0.0	0.0	2.01e5	0.00
Vermont	2.50e4	0.0	0.0	0.0	0.0	0.0	0.0	2.50e4	0.00
Virginia	3.45e4	1.19e4	2.66e4	3.06e4	1.87e3	0.0	5.91e4	1.05e5	0.00
West Virginia	5.62e4	6.59e3	0.0	0.0	0.0	0.0	0.0	6.28e4	0.00
Wyoming	2.53e5	0.0	0.0	0.0	0.0	0.0	0.0	2.53e5	0.00
<b>TOTAL (%)</b>	59.6	4.0	5.6	5.2	7.6	18.1	36.5		

Table C.6: **Predicted PD risk areas in Europe in 2050 after running the model under a  $R_0 = 5$  scenario and a homogeneous spatial vector distribution.** The epidemic-risk zones are classified according to the relative disease growth rates defined by the risk index, as very low, low, moderate, and high growth rates. The total risk refers to the sum of the epidemic-risk zones

Country	No risk (km <sup>2</sup> )	Trans. (km <sup>2</sup> )	Very low (km <sup>2</sup> )	Low (km <sup>2</sup> )	Med. (km <sup>2</sup> )	High (km <sup>2</sup> )	Total risk (km <sup>2</sup> )	Total surf. (km <sup>2</sup> )	Risk (%)
Russia	4.14e6	6.00e4	2.67e4	1.04e4	1.08e3	0.0	3.82e4	4.24e6	0.9
Norway	3.25e5	0.0	0.0	0.0	0.0	0.0	0.0	3.25e5	0.0
France	3.60e5	6.26e4	4.45e4	5.36e4	1.24e4	1.14e4	1.22e5	5.44e5	22.4
Sweden	4.42e5	0.0	0.0	0.0	0.0	0.0	0.0	4.42e5	0.0
Belarus	2.08e5	0.0	0.0	0.0	0.0	0.0	0.0	2.08e5	0.0
Ukraine	5.55e5	1.37e4	6.87e2	0.0	0.0	0.0	6.87e2	5.69e5	0.1
Poland	3.12e5	0.0	0.0	0.0	0.0	0.0	0.0	3.12e5	0.0
Austria	8.18e4	1.49e3	0.0	0.0	0.0	0.0	0.0	8.33e4	0.0
Hungary	3.96e4	5.27e4	0.0	0.0	0.0	0.0	0.0	9.23e4	0.0
Moldova	3.21e4	0.0	0.0	0.0	0.0	0.0	0.0	3.21e4	0.0
Romania	2.06e5	2.65e4	2.19e3	2.62e2	0.0	0.0	2.46e3	2.35e5	1.0
Lithuania	6.46e4	0.0	0.0	0.0	0.0	0.0	0.0	6.46e4	0.0
Latvia	6.42e4	0.0	0.0	0.0	0.0	0.0	0.0	6.42e4	0.0

Continued on next page

Table C.6: Extrapolated PD risk areas in Europe in 2050 with a homogeneous vector spatial distribution (Continued)

Estonia	4.55e4	0.0	0.0	0.0	0.0	0.0	0.0	4.55e4	0.0
Germany	3.53e5	1.85e3	0.0	0.0	0.0	0.0	0.0	3.55e5	0.0
Bulgaria	7.92e4	2.24e4	8.39e3	1.82e3	9.14	0.0	1.03e4	1.12e5	9.2
Greece	2.61e4	9.74e3	1.62e4	1.70e4	1.63e4	4.42e4	9.38e4	1.30e5	72.4
Albania	1.56e4	2.60e3	2.99e3	2.32e3	2.88e3	3.73e3	1.19e4	3.02e4	39.5
Croatia	2.26e4	2.19e4	2.30e3	2.66e3	2.82e3	1.77e3	9.55e3		
Switzer.	4.61e4	8.48	0.0	0.0	0.0	0.0	0.0	4.62e4	0.0
Luxemb.	2.71e3	0.0	0.0	0.0	0.0	0.0	0.0	2.71e3	0.0
Belgium	3.05e4	0.0	0.0	0.0	0.0	0.0	0.0	3.05e4	0.0
Nether.	3.69e4	0.0	0.0	0.0	0.0	0.0	0.0	3.69e4	0.0
Portugal	8.02e3	5.73e3	1.08e4	1.29e4	1.03e4	4.07e4	7.47e4	8.84e4	84.5
Spain	1.65e5	4.09e4	3.79e4	5.43e4	6.90e4	1.37e5	2.98e5	5.04e5	59.2
Ireland	6.82e4	0.0	0.0	0.0	0.0	0.0	0.0	6.82e4	0.0
Italy	8.37e4	1.65e4	3.15e4	5.15e4	3.89e4	7.77e4	1.99e5	3.00e5	66.6
Denmark	4.23e4	0.0	0.0	0.0	0.0	0.0	0.0	4.23e4	0.0
UK	2.43e5	0.0	0.0	0.0	0.0	0.0	0.0	2.43e5	0.0
Iceland	1.07e5	0.0	0.0	0.0	0.0	0.0	0.0	1.07e5	0.0
Slovenia	1.92e4	5.97e2	4.29e2	2.58e2	8.64	0.0	7.73e2	2.06e4	3.8
Finland	3.29e5	0.0	0.0	0.0	0.0	0.0	0.0	3.29e5	0.0
Slovakia	4.55e4	2.64e3	0.0	0.0	0.0	0.0	0.0	4.81e4	0.0
Czechia	8.08e4	0.0	0.0	0.0	0.0	0.0	0.0	8.08e4	0.0
Bosnia	4.45e4	4.74e3	7.19e2	0.0	0.0	0.0	7.19e2	5.00e4	1.4
Macedonia	1.73e4	4.23e3	3.23e3	1.85e2	0.0	0.0	3.41e3	2.50e4	13.7
Serbia	4.76e4	2.84e4	0.0	0.0	0.0	0.0	0.0	7.60e4	0.0
Monteneg.	1.09e4	1.82e2	3.63e2	4.55e2	9.11e2	4.57e2	2.19e3	1.33e4	16.5
Kosovo	9.71e3	1.45e3	0.0	0.0	0.0	0.0	0.0	1.12e4	0.0
Cyprus	4.04e2	0.0	0.0	0.0	0.0	4.95e3	4.95e3	5.35e3	92.5
Czech Rep.	7.85e4	0.0	0.0	0.0	0.0	0.0	0.0	7.85e4	0.0
Malta	0.0	0.0	0.0	0.0	0.0	1.99e2	1.99e2	1.99e2	100
<b>TOTAL (%)</b>	<b>87.6</b>	<b>3.8</b>	<b>1.9</b>	<b>2.0</b>	<b>1.5</b>	<b>3.2</b>	<b>8.6</b>		

Table C.7: **PD risk areas in Europe after running the model under a  $R_0 = 5$  scenario and a spatial heterogeneous vector distribution (climatic suitability)**. The epidemic-risk zones are classified according to the relative disease growth rates defined by the risk index, as very low, low, moderate and high growth rates. The total risk refers to the sum of the epidemic-risk zones

Country	No risk (km <sup>2</sup> )	Trans. (km <sup>2</sup> )	Very low (km <sup>2</sup> )	Low (km <sup>2</sup> )	Med. (km <sup>2</sup> )	High (km <sup>2</sup> )	Total risk (km <sup>2</sup> )	Total surf. (km <sup>2</sup> )	Risk (%)
Russia	4.23e6	3.07e3	3.56e2	0.0	0.0	0.0	3.56e2	4.24e6	0.01
Norway	3.25e5	0.0	0.0	0.0	0.0	0.0	0.0	3.25e5	0.0
France	4.66e5	5.79e4	1.32e4	6.31e3	9.15e2	0.0	2.05e4	5.44e5	3.76
Sweden	4.42e5	0.0	0.0	0.0	0.0	0.0	0.0	4.42e5	0.0
Belarus	2.08e5	0.0	0.0	0.0	0.0	0.0	0.0	2.08e5	0.0
Ukraine	5.69e5	0.0	0.0	0.0	0.0	0.0	0.0	5.69e5	0.0
Poland	3.12e5	0.0	0.0	0.0	0.0	0.0	0.0	3.12e5	0.0
Austria	8.33e4	0.0	0.0	0.0	0.0	0.0	0.0	8.33e4	0.0
Hungary	9.23e4	0.0	0.0	0.0	0.0	0.0	0.0	9.23e4	0.0
Moldova	3.21e4	0.0	0.0	0.0	0.0	0.0	0.0	3.21e4	0.0
Romania	2.35e5	0.0	0.0	0.0	0.0	0.0	0.0	2.35e5	0.0
Lithuania	6.46e4	0.0	0.0	0.0	0.0	0.0	0.0	6.46e4	0.0
Latvia	6.42e4	0.0	0.0	0.0	0.0	0.0	0.0	6.42e4	0.0
Estonia	4.55e4	0.0	0.0	0.0	0.0	0.0	0.0	4.55e4	0.0
Germany	3.55e5	0.0	0.0	0.0	0.0	0.0	0.0	3.55e5	0.0
Bulgaria	1.11e5	1.09e3	0.0	0.0	0.0	0.0	0.0	1.12e5	0.0
Greece	5.43e4	2.76e4	14845.01	71e4	1.58e4	0.0	4.77e4	1.30e5	36.8
Albania	1.97e4	2.89e3	2.88e3	4.56e3	9.48	0.0	7.54e3	3.02e4	25
Croatia	4.63e4	2.65e3	3.09e3	1.59e3	4.47e2	0.0	5.13e3	5.41e4	9.48
Switzer.	4.62e4	0.0	0.0	0.0	0.0	0.0	0.0	4.62e4	0.0
Luxemb.	2.71e3	0.0	0.0	0.0	0.0	0.0	0.0	2.71e3	0.0
Belgium	3.05e4	0.0	0.0	0.0	0.0	0.0	0.0	3.05e4	0.0
Nether.	3.69e4	0.0	0.0	0.0	0.0	0.0	0.0	3.69e4	0.0
Portugal	1.95e4	1.77e4	4.50e4	6.17e3	0.0	0.0	5.12e4	8.84e4	57.9
Spain	2.63e5	1.59e5	6.65e4	1.21e4	3.98e3	0.0	8.26e4	5.04e5	16.4
Ireland	6.82e4	0.0	0.0	0.0	0.0	0.0	0.0	6.82e4	0.0

Continued on next page

Table C.7: PD risk areas in Europe with a heterogeneous vector spatial distribution (Continued)

Italy	1.34e5	7.30e4	4.13e4	4.04e4	1.13e4	0.0	9.29e4	2.99e5	31.0
Denmark	4.23e4	0.0	0.0	0.0	0.0	0.0	0.0	4.23e4	0.0
UK	2.43e5	0.0	0.0	0.0	0.0	0.0	0.0	2.43e5	0.0
Iceland	1.07e5	0.0	0.0	0.0	0.0	0.0	0.0	1.07e5	0.0
Slovenia	2.02e4	2.58e2	8.64	0.0	0.0	0.0	8.64	2.06e4	0.42
Finland	3.29e5	0.0	0.0	0.0	0.0	0.0	0.0	3.29e5	0.0
Slovakia	4.81e4	0.0	0.0	0.0	0.0	0.0	0.0	4.81e4	0.0
Czechia	8.08e4	0.0	0.0	0.0	0.0	0.0	0.0	8.08e4	0.0
Bosnia	4.98e4	1.80e2	0.0	0.0	0.0	0.0	0.0	5.00e4	0.0
Macedonia	2.50e4	0.0	0.0	0.0	0.0	0.0	0.0	2.50e4	0.0
Serbia	7.60e4	0.0	0.0	0.0	0.0	0.0	7.60e4	0.0	
Monteneg.	1.15e4	3.64e2	7.29e2	7.31e2	0.0	0.0	1.46e3	1.33e4	11.0
Kosovo	1.12e4	0.0	0.0	0.0	0.0	0.0	0.0	1.12e4	0.0
Cyprus	4.04e2	0.0	1.41e3	2.22e3	1.31e3	0.0	4.95e3	5.35e3	92.4
Czech Rep.	7.85e4	0.0	0.0	0.0	0.0	0.0	0.0	7.85e4	0.0
Malta	0.0	0.0	0.0	0.0	1.99e2	0.0	1.99e2	1.99e2	100
<b>TOTAL (%)</b>	93.5	3.40	1.90	0.91	0.3	0.00	3.10		

Table C.8: **Surface of European vineyards in risk of PD given by the intersection of (Corine-Land-Cover) and the projected model in the ERA5-land data under a  $R_0 = 5$  scenario with the layer of vector climatic suitability.** The epidemic-risk zones are classified according to the relative disease growth rates defined by the risk index, as very low (0.1-0.33), low (0.33-0.66), moderate (0.66-0.9) and high exponential growth rates ( $> 90$ ). The total risk refers to the sum of the epidemic-risk zones.

Country	No risk (Ha)	Trans. (Ha)	Very low (Ha)	Low (Ha)	Med. (Ha)	High (Ha)	Risk (%)
Albania	1.64e3	4.42e2	1.12e2	1.37e3	0.00	0.00	4.16
Austria	6.65e4	0.00	0.00	0.00	0.00	0.00	0.00

Continued on next page

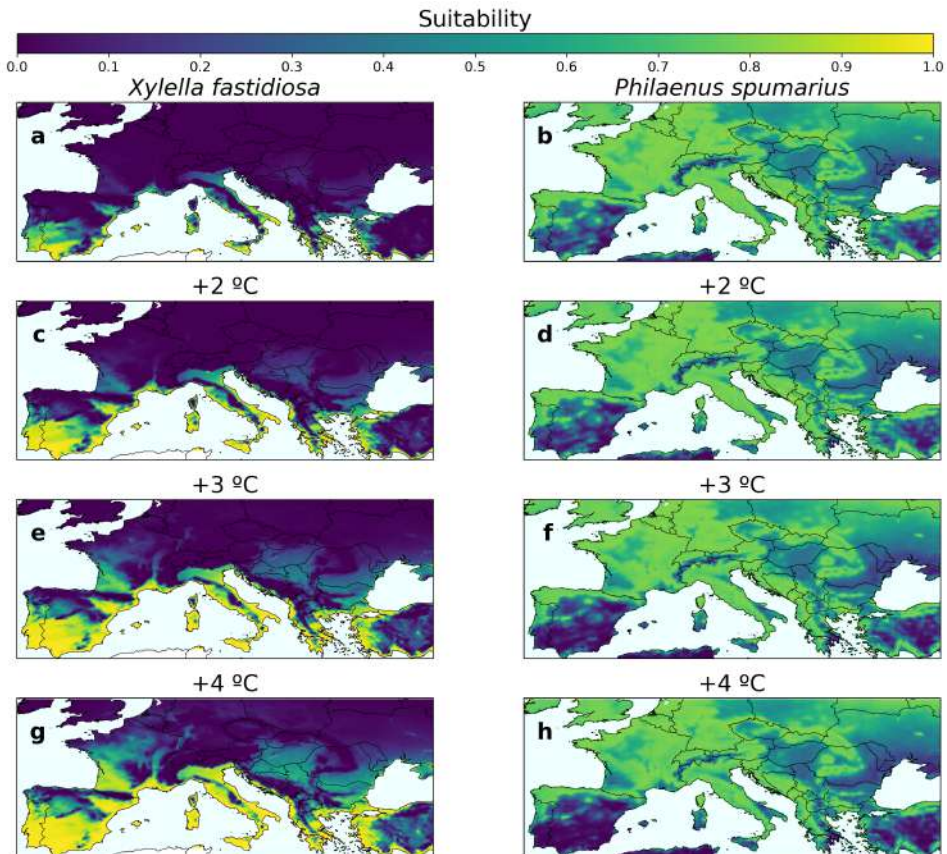
Table C.8: Surface of European vineyards at risk (Continued)

Bulgaria	1.13e5	3.20e3	0.00	0.00	0.00	0.00	0.00
Switzer.	1.45e4	0.00	0.00	0.00	0.00	0.00	0.00
Cyprus	0.00	0.00	1.95e2	9.54e3	4.41e3	0.00	100
Czech Rep.	1.69e4	0.00	0.00	0.00	0.00	0.00	0.00
Germany	1.30e5	0.00	0.00	0.00	0.00	0.00	0.00
Greece	1.38e4	2.06e4	1.35e4	8.54e3	2.43e4	0.00	5.74
Spain	2.84e5	6.97e5	6.41e4	5.17e3	1.21e3	0.00	6.71
France	3.71e5	4.07e5	2.82e5	6.53e4	3.52e3	0.00	3.11
Croatia	1.62e4	1.52e3	3.22e3	2.58e3	1.34e3	0.00	2.87
Hungary	1.01e5	0.00	0.00	0.00	0.00	0.00	0.00
Italy	1.60e5	1.80e5	8.37e4	1.16e5	8.01e4	0.00	4.51
Luxemb.	1.63e3	0.00	0.00	0.00	0.00	0.00	0.00
Monteneg.	0.00	1.90	2.63e3	2.29e2	0.00	0.00	100
Macedonia	2.77e4	0.00	0.00	0.00	0.00	0.00	0.00
Malta	2.57	0.00	0.00	0.00	2.74	0.00	100
Portugal	4.37e4	7.42e4	9.13e4	2.07e3	0.00	0.00	4.42
Romania	2.25e5	0.00	0.00	0.00	0.00	0.00	0.00
Serbia	8.61e3	0.00	0.00	0.00	0.00	0.00	0.00
Slovenia	2.60e4	1.17e3	8.41e2	0.00	0.00	0.00	3.00
Slovakia	2.06e4	0.00	0.00	0.00	0.00	0.00	0.00
<b>TOTAL (%)</b>	42.1	35.6	13.9	5.4	3.0	0.00	

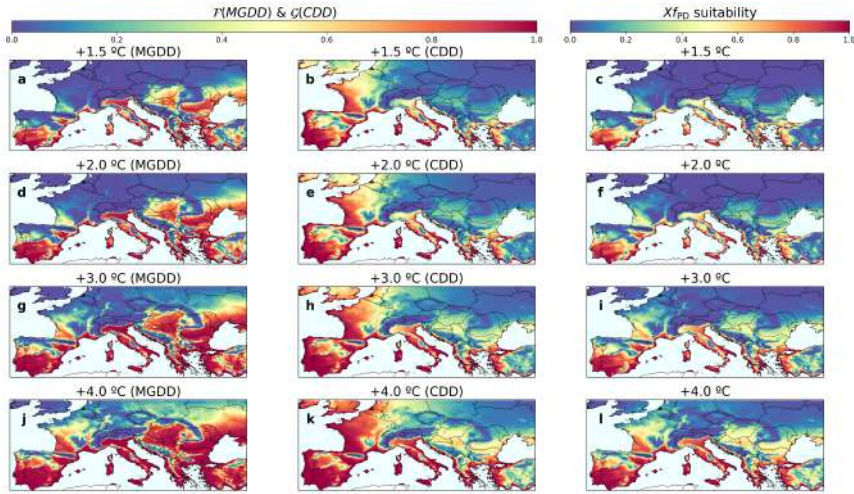
Table C.9: **Predicted PD risk in 2050 in European vineyards (Corine-Land-Cover) considering a  $R_0 = 5$  scenario and the vector climatic suitability.** The epidemic-risk zones are classified according to the relative disease growth rates defined by the risk index, as very low, low, moderate, and high growth rates. The total risk refers to the sum of the epidemic-risk zones.

Country	No risk (Ha)	Trans. (Ha)	Very low (Ha)	Low (Ha)	Med. (Ha)	High (Ha)	Risk (%)
Albania	7.67e2	9.22e2	4.76e2	1.35e3	4.58	0.00	5.26
Austria	6.65e4	0.00	0.00	0.00	0.00	0.00	0.00
Bulgaria	1.10e5	5.73e3	5.27e2	0.00	0.00	0.00	0.45
Switzer.	1.45e4	0.00	0.00	0.00	0.00	0.00	0.00
Cyprus	0.00	0.00	8.49e2	1.23e4	9.64e2	0.00	100
Czech Rep.	1.69e4	0.00	0.00	0.00	0.00	0.00	0.00
Germany	1.30e5	0.00	0.00	0.00	0.00	0.00	0.00
Greece	1.34e4	1.31e4	2.45e4	7.91e3	2.19e4	0.00	6.72
Spain	3.65e5	6.18e5	6.04e4	6.00e3	1.21e3	0.00	6.44
France	2.45e5	1.98e5	5.10e5	1.71e5	5.71e3	0.00	6.08
Croatia	1.59e4	2.43e2	1.88e3	6.20e3	7.01e2	0.00	3.53
Hungary	1.01e5	0.00	0.00	0.00	0.00	0.00	0.00
Italy	9.28e4	1.31e5	2.29e5	1.55e5	1.26e4	0.00	6.39
Luxemb.	1.63e3	0.00	0.00	0.00	0.00	0.00	0.00
Monteneg.	0.00	0.00	1.90	2.86e3	0.00	0.00	100
Macedonia	2.65e4	1.20e3	0.00	0.00	0.00	0.00	0.00
Malta	2.57	0.00	0.00	0.00	2.74	0.00	100
Portugal	1.94e4	7.67e4	1.01e5	1.39e4	0.00	0.00	5.45
Romania	2.25e5	0.00	0.00	0.00	0.00	0.00	0.00
Serbia	8.61e3	0.00	0.00	0.00	0.00	0.00	0.00
Slovenia	2.25e4	1.61e3	3.07e3	8.41e2	0.00	0.00	1.40
Slovakia	2.06e4	0.00	0.00	0.00	0.00	0.00	0.00
<b>Tot. (%)</b>	38.40	26.90	23.90	9.70	1.10	0.00	

## C.8 Present and future climate suitability for Xf & Ps

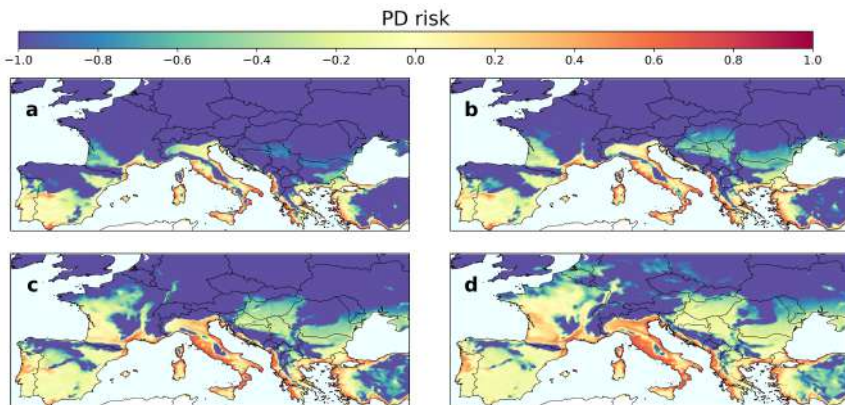


**Figure C.14:** Climate suitability for  $Xf_{PD}$  and *P. spumarius* under the current scenario and different climate projections. (a,b) Current scenario. (c,d) +2 °C climate projection. (e,f) +3 °C climate projection. (g,h) +4 °C climate projection.



**Figure C.15:**  $\mathcal{F}(MGDD)$ ,  $\mathcal{G}(CDD)$  and suitability of  $Xf_{PD}$  ( $\mathcal{F}(MGDD) \cdot \mathcal{G}(CDD)$ ) under different climate projections. (a-c) +1.5 °C climate projection. (d-f) +2 °C climate projection. (g-i) +3 °C climate projection. (j-l) +4 °C climate projection.

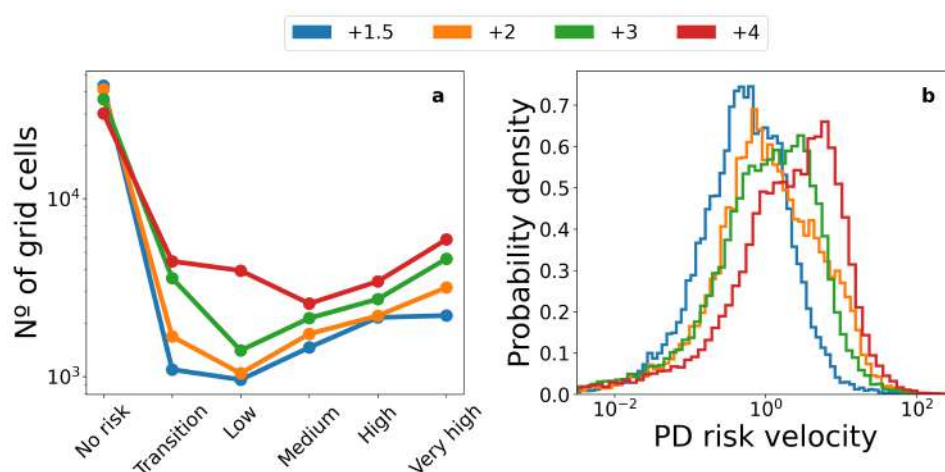
## C.9 Future Pierce's disease risk projections under climate change



**Figure C.16:** Future risk of PD epidemics under different climate projections. (a) +1.5 °C climate projection. (b) +2 °C climate projection. (c) +3 °C climate projection. (d) +4 °C climate projection.

**Table C.10:** Some risk velocity statistics for each climate projection.

Scenario	+1.5 °C	+2 °C	+3 °C	+4 °C
Risk velocity > 5 km/year (%)	6.34	33.44	23.56	52.48
Max risk velocity (km/year)	181	691	393	667
Mean risk velocity (km/year)	0.78	1.9	3.84	5.15



**Figure C.17: Risk and risk velocity shifts as function of the climate projections.** (a) Number of grid cells at each risk level across the different climate projections. (b) Histograms of PD risk velocities in each climate projection. A shift toward higher velocities can be clearly observed.

## C.10 Risk analysis within European wine PDOs

An extended analysis on the potential impact of PD in European PDOs can be easily performed using our PD establishment risk [webpage](#), where the following insights can be contrasted in an interactive way.

Some French wine regions consistently remain risk-free in all scenarios, such as Alsace, Jura, and Champagne, while others, like Bourgogne (including Beaujolais) and the Loire Valley, only reach the transitional risk region in the warmest scenarios. Conversely, the Rhône Valley, Southwest, Languedoc, Roussillon, Provence, and Bordeaux face increasing risk in warmer scenarios. Within the Rhône region, Condrieu PDO stays out of risk in the northern part, while Côte-Rotie reaches low risk only in the +4 °C scenario. Further south, Châteauneuf-du-Pape reaches low risk from the +2 °C scenario, and Costières

de Nîmes becomes low risk from the +1.5 °C scenario, escalating to medium risk from the +2 °C scenario. Regions directly influenced by the Mediterranean, like Cassis and Bandol in Provence or Muscat de Frontignan in Languedoc, experience higher risk, with some coastal areas reaching medium-high risk levels. In Provence, the risk becomes medium in the warmest scenarios, while in Languedoc-Roussillon, it typically decreases to low risk from the +2 or +3 °C scenarios (except for Muscat DOPs). In Bordeaux, certain right bank PDOs (Pomerol, Lalande de Pomerol) become low risk in the +3 °C scenario, while Saint Emilion achieves the same in the +4 °C scenario. On the left bank, major DOPs such as Margaux, Saint Julien, Pauillac, and Saint Estèphe reach low risk in the +4 °C scenario, along with Pessac-Léognan in the region below the city of Bordeaux. Graves and Sauternes, however, reach at most the transitional risk level in the +3 °C scenario.

In Spain, contrasting patterns are observed. The northwestern part, including Rioja and Ribera del Duero DOPs, consistently remains risk-free in all four scenarios, except for the coastal Rías Baixas DOP, which becomes low risk in the +4 °C scenario. In the southern region of Andalusia, a decrease in risk level is observed with increasing warming, as seen in the cases of the Jerez/Sherry PDO (labeled as Manzanilla) and Sierras de Málaga, which have low risk in the +1.5 °C and +2 °C scenarios and become transitional at +3 °C and +4 °C, respectively. This decrease in risk is associated with decreased vector suitability as temperatures rise. In the more continental DOPs of Aragón and Valencia (Utiel-Requena), the risk remains non-existent, while those closer to the Mediterranean exhibit increased risk levels. Penedès becomes low risk in the +2 °C scenario, and the coastal DOPs of Alella and Empordà clearly increase to medium risk at the +2 °C scenario from a low risk level at the +1.5 °C scenario. Other more continental DOPs, such as Priorat, become transitional at most. A similar pattern is observed in central Spain, with northern DOPs remaining risk-free and more southern ones (Méntrida, Ribera del Guadiana) becoming transitional. Notably, the large La Mancha DOP remains risk-free, again due to decreased vector suitability.

**Table C.11: Percentage of land surface, vineyard surface and PDOs at risk in Europe under different climate projections.**

Scenario	Land surface (%)				Vineyard surface (%)				PDOs (%)			
	+1.5 °C	+2 °C	+3 °C	+4 °C	+1.5 °C	+2 °C	+3 °C	+4 °C	+1.5 °C	+2 °C	+3 °C	+4 °C
No risk	99.4	99.15	98.61	96.14	68.89	61.3	45.54	37.64	72.7	64.88	49.17	33.89
Transition	0.28	0.39	0.64	1.99	12.44	14.35	21.86	22.0	9.13	9.39	15.28	18.79
Low	0.12	0.18	0.32	0.82	8.01	10.26	14.58	18.51	7.11	12.64	15.45	20.19
Medium	0.15	0.22	0.34	0.87	9.33	12.1	15.93	19.96	9.13	10.01	15.63	22.3
High	0.04	0.06	0.09	0.18	1.34	1.99	2.09	1.89	1.93	3.07	4.48	4.83
Total risk (r>0.1)	0.32	0.46	0.75	1.87	18.67	24.35	32.6	40.35	18.17	25.72	35.56	47.32

**Table C.12: Percentage of land surface, vineyard surface and PDOs at risk per different European countries under different climate projections.**

Scenario	Total surface (%)				Vineyard surface (%)				PDO (%)			
	+1.5 °C	+2 °C	+3 °C	+4 °C	+1.5 °C	+2 °C	+3 °C	+4 °C	+1.5 °C	+2 °C	+3 °C	+4 °C
Austria	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Belgium	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Bulgaria	0.01	0.03	0.11	3.43	0.25	0.41	2.74	4.88	0.22	0.56	1.91	3.27
Croatia	0.24	0.38	0.64	14.82	31.53	42.6	45.89	47.65	31.27	23.75	27.93	32.47
France	0.53	0.84	2.32	7.25	24.21	35.75	56.43	80.0	13.37	15.62	22.73	41.65
Germany	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Greece	1.62	2.33	3.32	53.37	59.61	65.08	70.36	66.69	66.38	71.19	72.61	74.68
Italy	3.44	4.98	8.39	13.8	57.49	65.4	77.34	81.79	45.77	58.34	74.71	82.65
Portugal	4.66	8.36	30.05	30.22	12.7	28.14	33.77	42.71	19.54	27.9	35.43	36.82
Slovenia	0.01	0.01	0.05	0.14	1.26	4.5	12.99	21.3	2.86	3.73	6.05	11.14
Slovakia	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Spain	6.98	7.43	6.33	6.34	3.77	3.99	4.55	4.42	12.74	13.41	10.6	10.83

In Portugal, a similar pattern of decreasing risk is observed in the southern regions of Algarve and Setúbal. However, there is a clear increase in risk in the central DOPs of Bairrada (reaching low-risk in the +2°C scenario and medium-risk from +3°C), Dão DOP (becoming low-risk from the +3°C scenario), and Vinho Verde DOP, which exhibits the same pattern as the Rías Baixas region in Galicia, becoming low risk in the +4°C scenario. The Óbidos DOP also becomes low risk from the +2°C scenario. The Carcavelos and Colares DOPs, with their strong maritime influence, maintain or reach a medium risk level, while other DOPs like Alentejo, Beira Interior, and Trás-os-Montes are not at risk. The Douro DOP transitions from no risk to the transitional level, but this change may be influenced by the relatively low spatial resolution of our study, as the wine-producing area is relatively narrow (less than 0.1°).

In Italy, there is an overall increase in risk observed across the country. However, the increase is more limited in Piedmont, where DOPs such as Barolo, Alba, Barbera, and Langhe only reach low-risk levels at 4°C. Similarly, in Veneto and Friuli, they become low risk from the 3°C scenario, while Alto Adige DOP remains non-risk. In Tuscany, there is a considerable higher increase in risk. Coastal regions like Maremma and Bolgheri reach medium and high risk at 2°C, respectively, having positive risk already in the 1.5°C scenario and the island of Elba already is already at high risk in the base scenario. More interior tuscan DOPs like Chianti and Brunello di Montalcino become low-risk from the +2°C scenario. In Umbria, the risk increases rapidly, with Colli Perugini DOP reaching medium-risk at the +3°C scenario, and Torgiano and Orvieto transitioning from low-risk at +2°C to medium risk at +3°C. In the Apulia region, there are distinct differences between the lower tip of the Apulian peninsula, where the risk levels remain consistently high, and the more continental north, where DOPs like Gravina in Bari province become low-risk from the +2°C scenario,

and San Severo DOP in Foggia province consistently stays at low-risk across all scenarios. In the islands, the risk decreases with increasing warming scenarios in some zones of Sicily, like Marsala, while others stay at risk and not changing the level, and a nonmonotonic behavior is found in Etna (decreases at 2°C and increases again at 3°C), while Sardinian DOPs maintain their risk level, typically medium, like Malvasia di Bosa, Vermentino di Gallura and Vernaccia di Oristano.

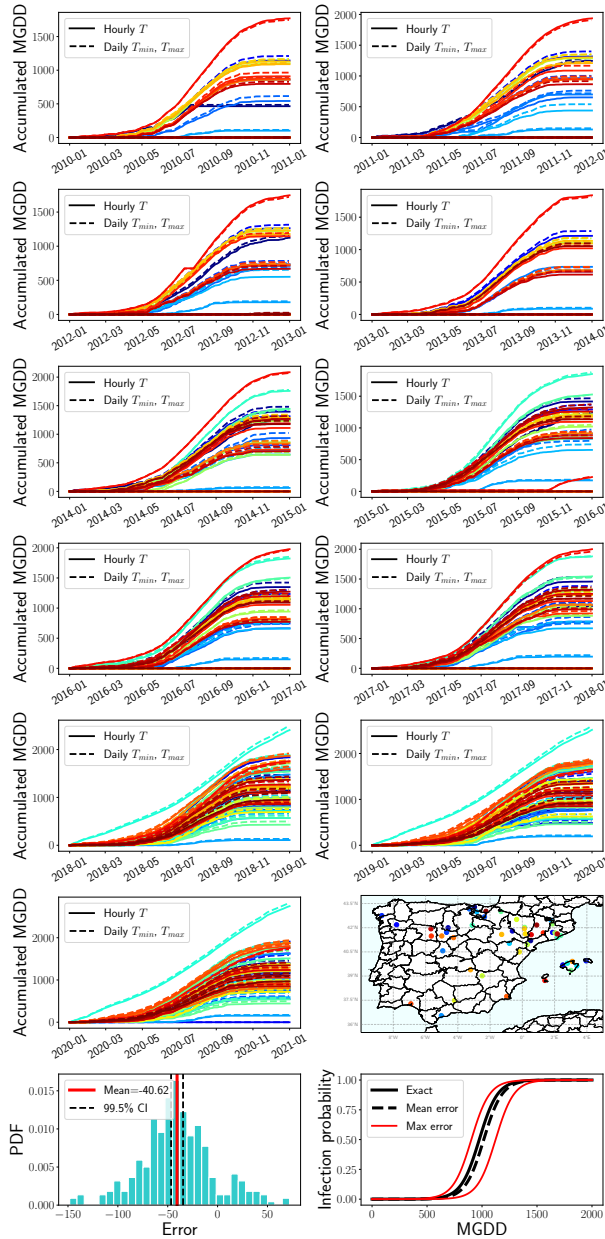
## C.11 MGDD and CDD approximation

MGDD and CDD metrics were originally defined using hourly temperature data from ERA-5 land dataset (Chapter 7, [280]). However, the E-OBS and CORDEX datasets used for the climate projections only provide daily granularity. To overcome this limitation we use a basic sinusoidal extrapolation relating maximum and minimum daily temperature to hourly temperatures,

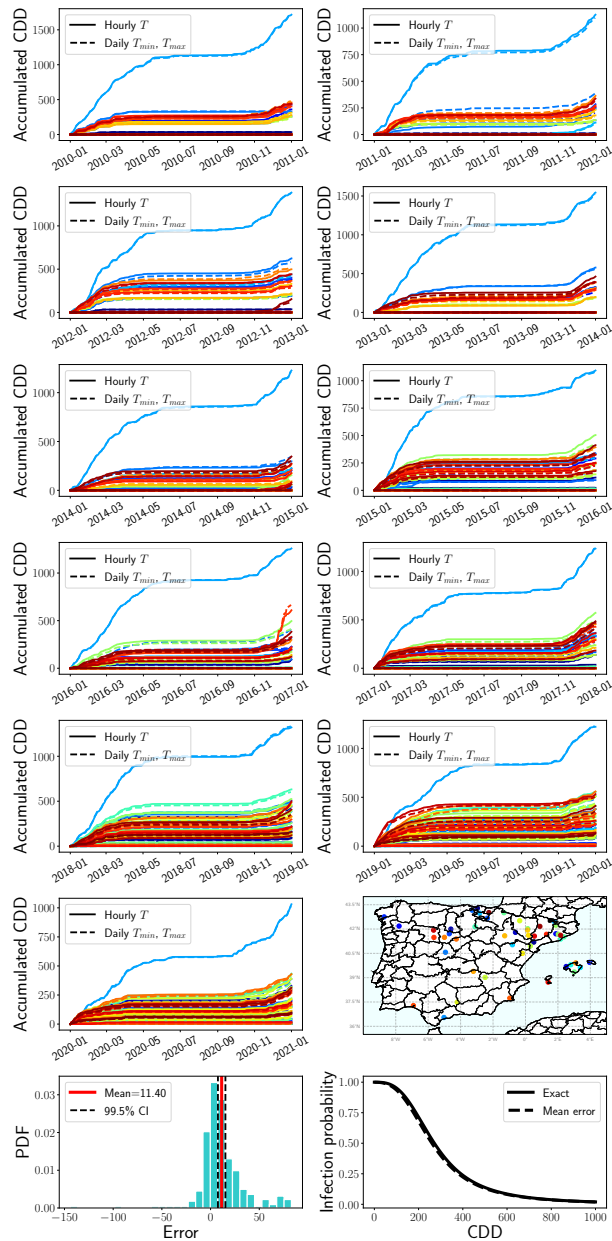
$$T_h = \frac{T_{max} + T_{min}}{2} + \frac{T_{max} - T_{min}}{2} \sin(w \cdot h) , \quad (C.31)$$

with  $w = 2\pi/24$  and  $h$  ranging from 0 to 23 (as MGDD and CDD are cumulative sums, the phase is irrelevant for the approximation).

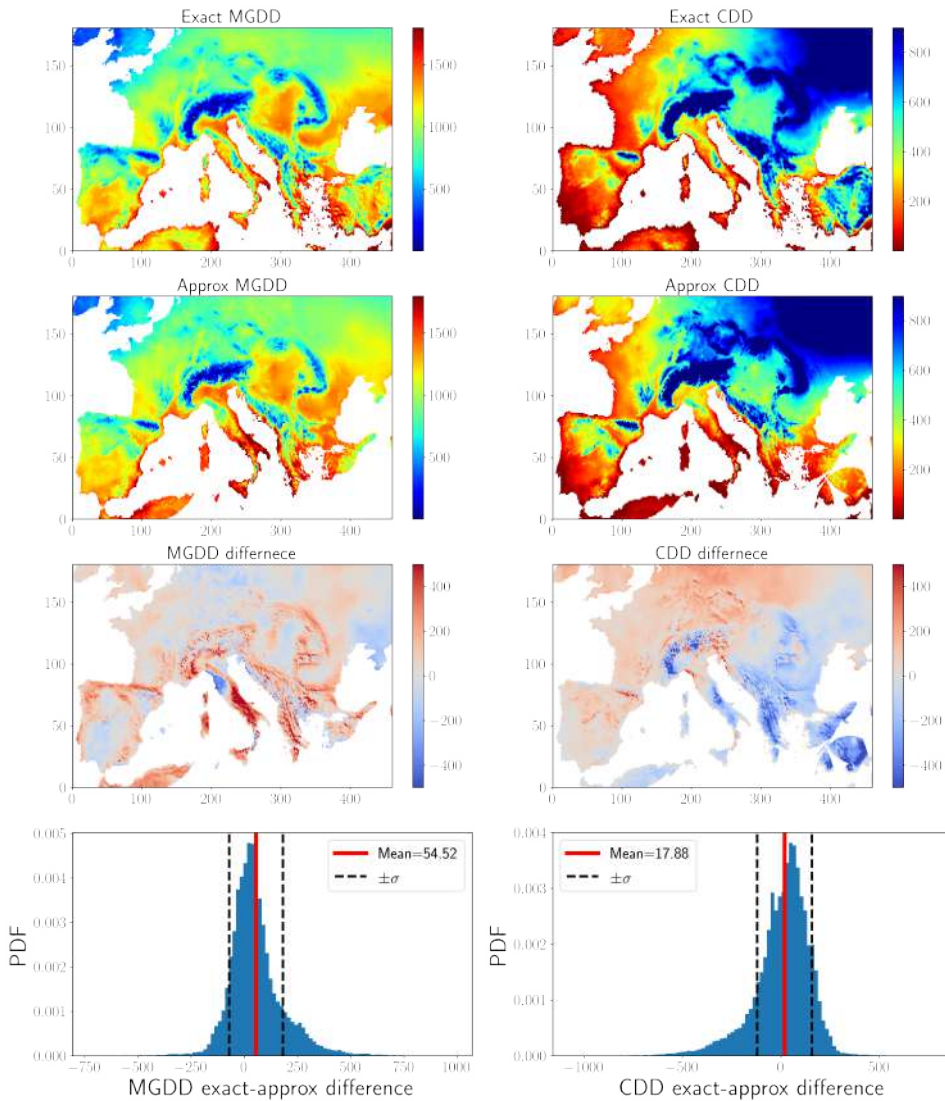
The approximation was validated with data from the national meteorological agency in Spain (AEMET). using hourly temperature data obtained from 50 meteorological stations in the period 2010-2020. To this end, we computed MGDD and CDD using both the full hourly data and only the daily maximum and minimum temperatures, in the latter case using Eq. (C.31). The resulting MGDD and CDD estimates using both approaches were very similar (Figs. C.18 and C.19). Because the temporal resolution of the E-OBS and ERA-5 land data sets are different and are acquired using different methodologies, we evaluated the possible divergence between the MGDD and CDD estimated from each dataset. The results calculated with both data agreed, showing a mean difference of 54 and 17 units for MGDD and CDD, respectively, with a standard deviation of about 200 units for both metrics (Fig. C.20).



**Figure C.18:** Comparison of MGDD accumulation computed with hourly mean temperature data (solid line) and daily maximum and minimum values (dashed line) using data from several meteorological stations in Spain and different years. The last row shows the distribution of errors and mean error in MGDD and the mean and maximum potential errors in infection probability.



**Figure C.19:** Comparison of CDD accumulation computed with hourly mean temperature data (solid line) and daily maximum and minimum values (dashed line) using data from several meteorological stations in Spain and different years. The last row shows the distribution of errors and mean error in CDD and the mean and maximum potential errors in infection probability.

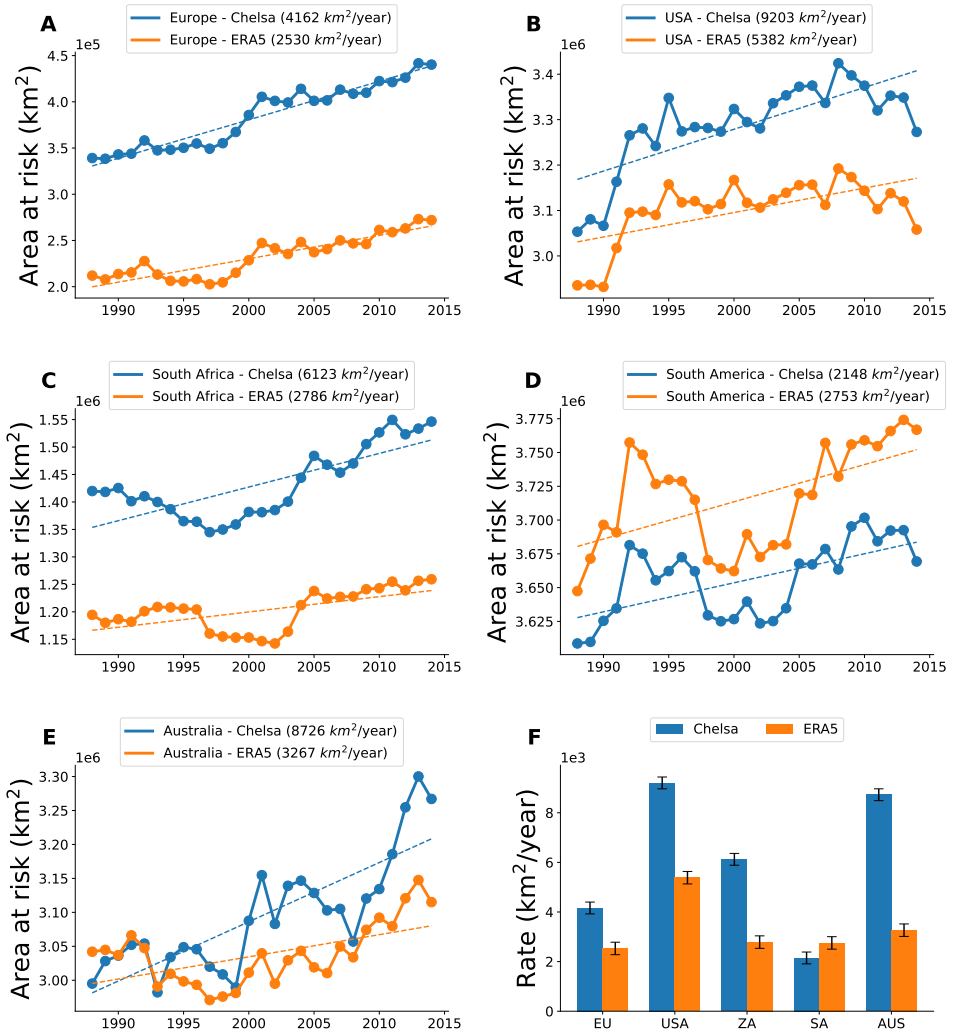


**Figure C.20:** Differences in annual accumulated MGDD and CDD when using ERA5-Land dataset with hourly mean temperature and EOBS dataset with maximum and minimum daily temperatures.

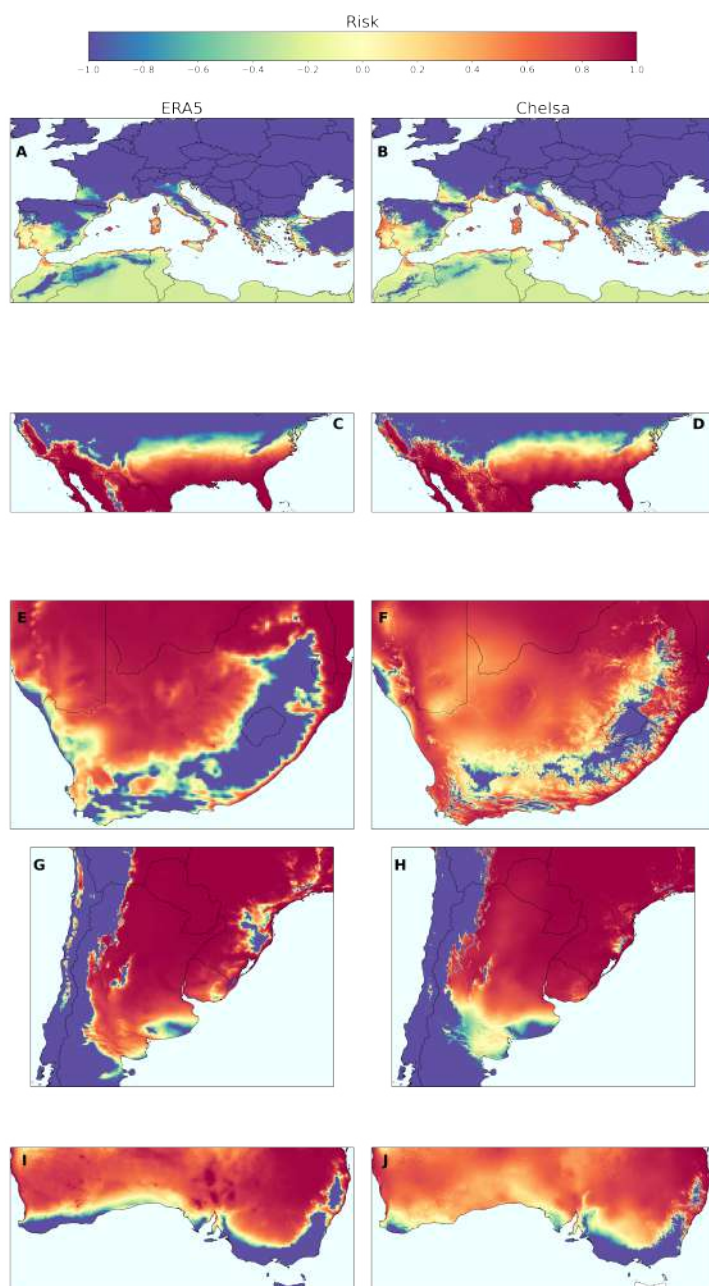
## C.12 The effect of spatial resolution on risk projections

We compared the risk indices obtained with ERA5 (10 km resolution) and CHELSA (1 km resolution) datasets in different viticulture areas. The results show that the risk indices obtained with CHELSA are generally higher than those obtained with ERA5, specially in river valleys and coastal areas, where

the temperature is more influenced by the local topography. This is particularly evident in the Mediterranean basin or South Africa (Fig. C.22). We also compared the projected risk increase rate in different viticulture areas, showing that the risk increase rate is generally higher in CHELSA than in ERA5, with the exception of South America, where the risk increase rate is similar in both datasets (Fig. C.21).

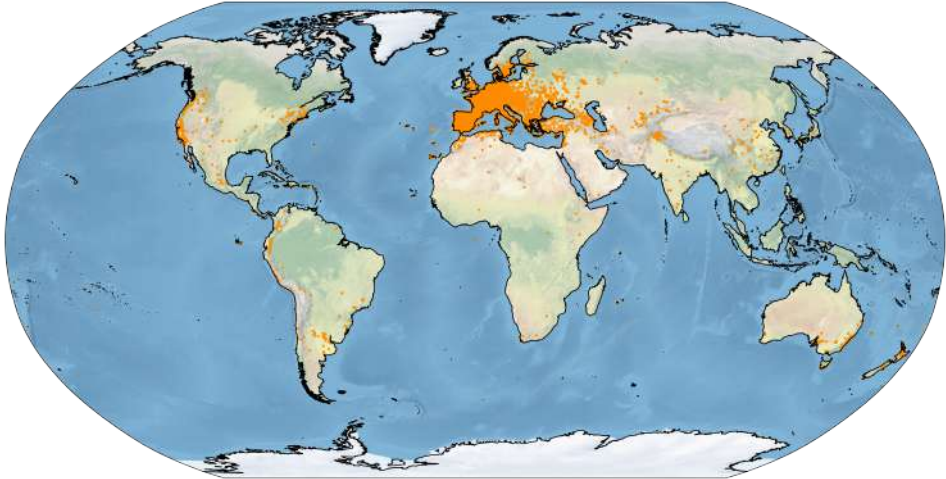


**Figure C.21:** Difference in projected risk in increase rate based on CHELSA (high-resolution, 1 km) and ERA5 (mid-resolution, 10 km) datasets in global viticulture areas. (A) Europe (B) United States (C) South Africa (D) South America (E) Australia.



**Figure C.22:** Comparison of risk indices obtained with ERA5 (mid-resolution – 10 km, left column) and CHELSA (high-resolution – 1 km, right column) datasets in Europe (A-B), United States (C-D), South Africa (E-F), South America (G-H) and Australia (I-J).

## C.13 *Vitis vinifera* global distribution



**Figure C.23:** Presence locations of *Vitis vinifera* obtained from GBIF.



## D. Nonlinear time-series analysis and reconstruction

### D.1 Seasonal adjusted fits for pH and temperature

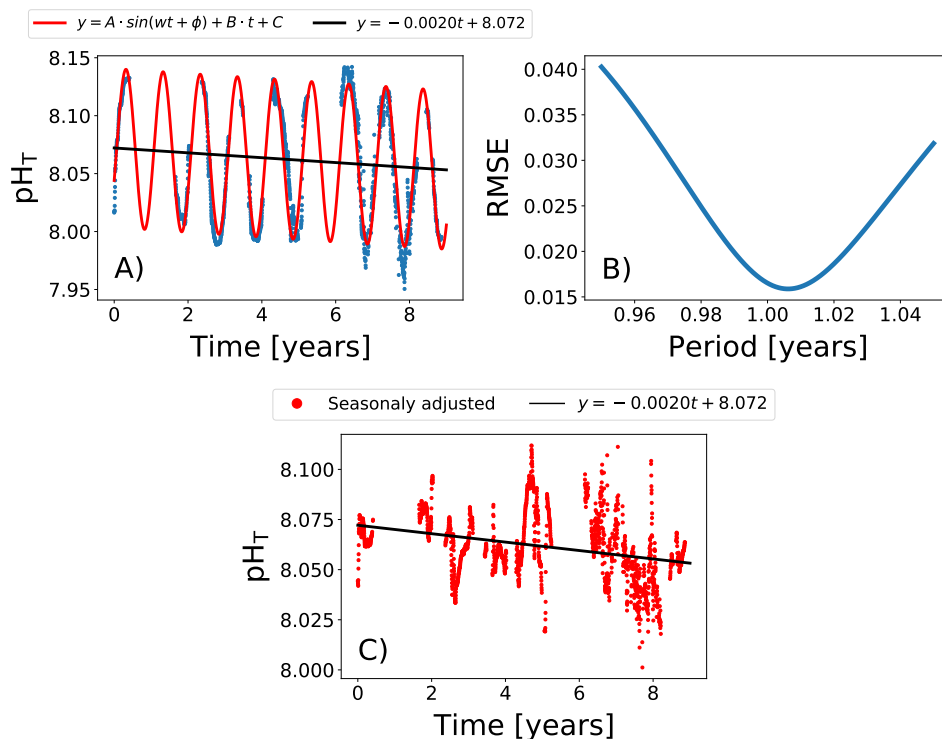
Linear regression is not suitable for analyzing seasonal data due to its assumption of constant variance and independence of observations, which are often violated in seasonal time series. Seasonal data typically exhibit patterns that repeat at regular intervals over time, leading to violations of the independence assumption. Moreover, the variance of the data may not be constant across different seasons, violating the assumption of constant variance. These violations can result in biased parameter estimates and inaccurate predictions when using linear regression models to analyze seasonal data. To account for seasonality in data, more appropriate methods such as seasonal decomposition or time series models like SARIMA (Seasonal Autoregressive Integrated Moving Average) should be employed, which explicitly capture the seasonal patterns and dependencies present in the data. However, the presence of missing data make them impractical to use in this study.

Another possibility to obtain the long-term trend of the data is to remove the seasonal component from the time series. This can be done by fitting a sinusoidal function to the data and subtracting the fit from the original time series. Indeed, one can directly fit a sinusoidal function with a linear component to the data, which can be expressed as:

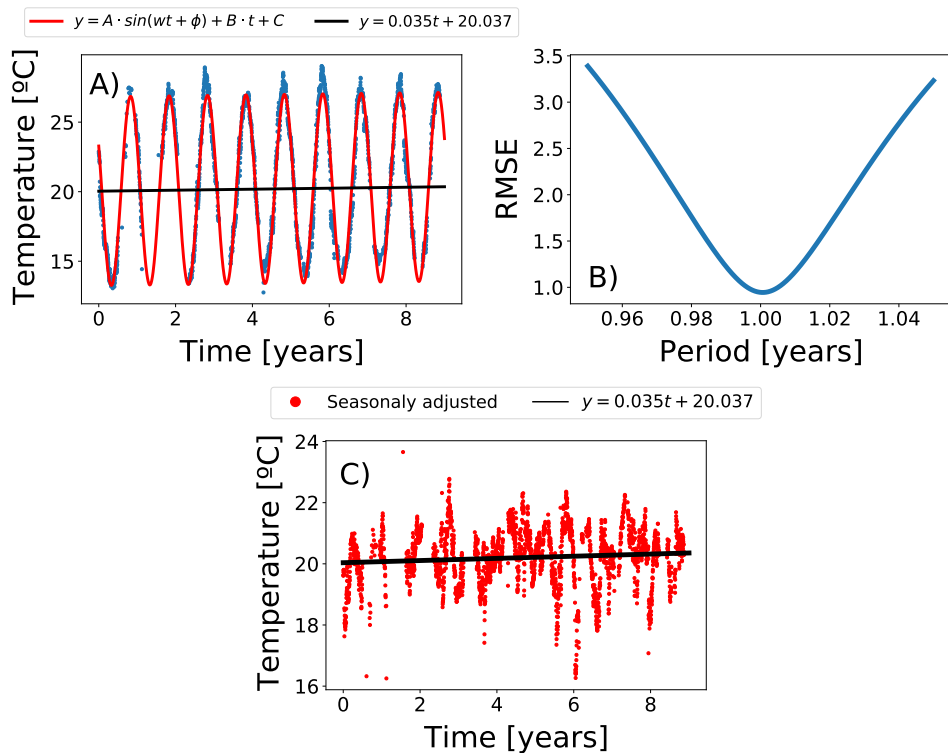
$$f(t) = A \sin(\omega t + \phi) + Bt + C \quad (\text{D.1})$$

where  $A$  is the amplitudes of the sine function,  $\omega$  is the angular frequency,

$\phi$  is the phase shift,  $B$  is the slope of the linear component and  $C$  is the offset. The angular frequency is related to the period of the function by  $\omega = 2\pi/T$ , where  $T$  is the period. The optimal values for the parameters were found by minimizing the mean squared error between the seasonally adjusted data and the fit. The optimal period was found by minimizing the mean squared error between the seasonally adjusted data and the fit for different periods. The optimal period was found to be  $T = 1.006$  for pH and  $T = 1.0$  for temperature. The optimal fits are shown in Figs. D.1 and D.2.

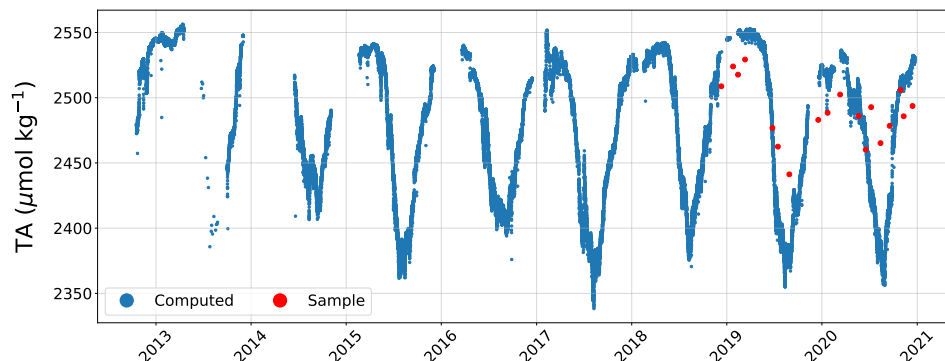


**Figure D.1: Seasonally adjusted fit to reconstructed and measured pH data.** A) Full fit of Eq. (1) with  $A = 0.0686 \pm 0.0005$ ,  $B = -0.0020 \pm 0.0002$ ,  $\phi = -6.704 \pm 0.008$ ,  $C = 8.0721 \pm 0.0008$ . B) Optimal period ( $T = 1.006$ ,  $\omega = 2\pi/T$ ) found to fit the data. C) Linear regression to the seasonal adjusted data.



**Figure D.2: Seasonally adjusted fit to reconstructed and measured temperature data.** A) Full fit of Eq. (1) with  $A = 6.792 \pm 0.029$ ,  $B = -0.0353 \pm 0.0080$ ,  $\phi = -3.6396 \pm 0.0040$ ,  $C = 20.037 \pm 0.043$ . B) Optimal period ( $T = 1.0$ ,  $\omega = 2\pi/T$ ) found to fit the data. C) Linear regression to the seasonal adjusted data.

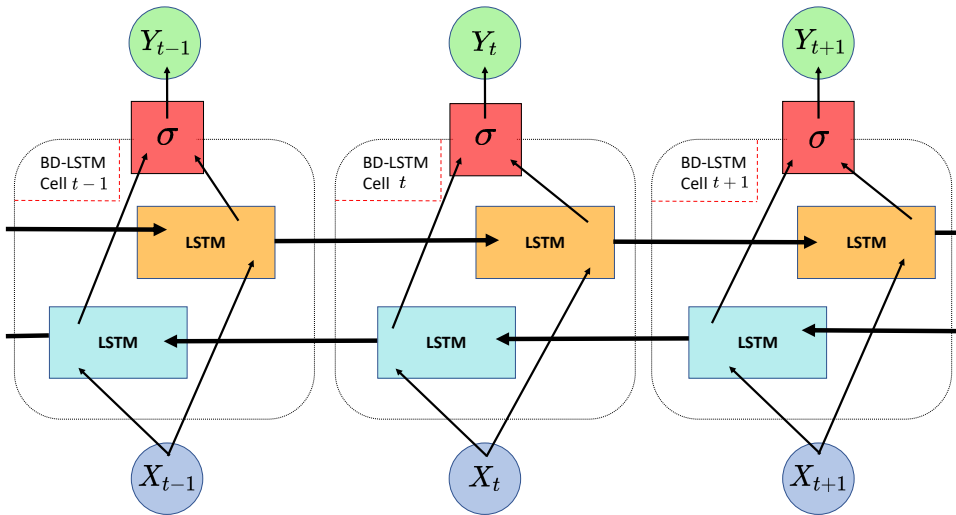
## D.2 Total alkalinity in the Bay of Palma



**Figure D.3: Total alkalinity in the Bay of Palma.** Total alkalinity in  $\mu\text{mol kg}^{-1}$  calculated values (blue dots) and concentrations obtained from samples (red dots) in the Bay of Palma.

## D.3 Bidirectional Long-Short Term Memory neural network

Recurrent Neural Networks (RNNs) are a class of artificial neural networks in which node connections arise along a temporal sequence, i.e., previous values in a time series are linked to current values. In simple words, RNNs predict a point of the time series using past information. Simple Recurrent Neural Networks (SRNNs) are the straightforward extension of Feedforward Neural Networks, in which past information and learned knowledge is encoded in the network as state vectors. SRNNs suffer from the so-called vanishing gradient problem, i.e., distant parts of the time series do not play a role in the training process. Thus, SRNNs are not capable to learn long-term dependencies. Long-Short Term Memory (LSTM) neural networks overcome this limitation by implementing three gates to update and control the cell state (forget gate, input gate, output gate), thus allowing to keep long-term dependencies [106]. Bidirectional Long-Short Term Memory (BD-LSTM) neural networks are able to encode both past and future information by implementing two LSTM layers flowing in opposite time directions. The forward layer preserve past information while the backwards layer preserves future information. Thus, using the two hidden states combined BD-LSTM are able in any point in time to preserve information from both past and future. A schematic representation of the BD-LSTM neural network is shown in (Fig. D.4).



**Figure D.4: Scheme for the Bidirectional-LSTM Neural Network.** The network receives as input a tensor of shape  $(\text{batch\_size}, \text{window\_size}, N_{\text{features}})$ , where  $\text{batch\_size}$  is the number of examples to train per iteration,  $\text{batch\_size}$  is the number of past and future points considered and  $N_{\text{features}}$  is the number of features used to predict the target series. First, the input is linearly transformed to match the number of cells of the network (3 in the figure). Then, the input is transformed by the backward (blue) and forward (orange) LSTM layers, which propagates information backwards and forwards, respectively. Finally, the values are transformed through an activation gate. In our model, this values will be ultimately transformed via a dense layer to a unique output.



# E. Mapping marine habitats with deep learning

## E.1 Satellite imagery and ground truth data

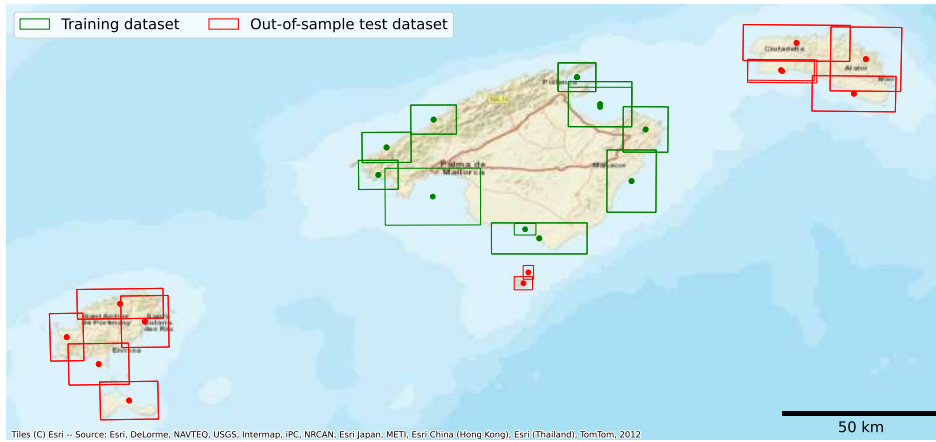
### E.1.1 Study region covering

We obtained a total of 60 Planetscope satellite images covering the coast of the Balearic Islands for the years 2020 to 2023. From those, the 20 images covering the island of Mallorca for years 2020-2022 were used to conform the training dataset, while the images from the other islands for the years 2020-2022 were included in the out-of-sample test set (Fig. E.1). The remaining images from the whole region in the year 2023 were left as final test. The metadata of the satellite images used in the study is shown in Table E.1.

Here we must emphasize that we refer to this test set as “out-of-sample” because the images conforming the set are relatively distant to the ones in the training set. Furthermore, the sub-classes conforming the 4 major ecological classes are not exactly the same in each island (see below). This, in principle, hinders the robustness and generalization ability of deep learning models. Perhaps this is one of the reasons for which a general and robust model for habitat mapping in the Mediterranean Sea has been hitherto lacking. Thus, with this strict separation of training and out-of-sample test set we aim to comprehensively study the robustness and generalization power of our deep learning models.

Because we will finally train our model with all available data (from all regions), the images from 2023 will conform the final test set, to ensure the robustness of the model to changes in environmental conditions at the moment

of image acquisition and image metadata.



**Figure E.1: Coverage of the training and out-of-sample test datasets in the Balearic Islands.** Each image was obtained for different years from 2020 to 2023, depending on the availability.

Table E.1: Metadata of the satellite images used in the study.

Name	Satellite azimuth	Sun azimuth	Sun elevation
Formentera_18_july_2021	99.4	113.5	57.3
Formentera_24_july_2022	97.7	127.8	63
Formentera_3_august_2023	110	128.8	60.5
South_Oeste_Menorca_29_July_2022	100.6	116.1	54.1
South_Oeste_Menorca_09_July_2021	182.4	114.2	58.1
Sur_Oeste_Menorca_9_august_2023	101.5	121.5	53.3
Sur_Ibiza_20_july_2022	101.1	121.9	61.4
Sur_Ibiza_18_july_2021	99.5	113.7	57.3
Sur_Ibiza_26_august_2023	101.1	139.3	55.5
Es_trenc_25_July_2022	270.7	113.8	54.6
Es_Trenc_15_july_2023	112.8	112.6	56.7
Norte_Menorca_29_July_2021	104.5	136.3	63
Norte_Menorca_20_July_2022	176.1	118.6	58.1

Continued on next page

Table E.1: Metadata of the satellite images used in the study. (Continued)

Norte_Menorca_23_june_2023	102	123.6	64.5
Este_Menorca_19_July_2021	271.4	115.9	56.7
Este_Menorca_15_July_2022	268	113.3	56.1
Este_Menorca_23_june_2023	266.2	112.5	58.8
CalaPi_CalaFiguera_29_july_2021	100.9	117.6	55.7
CalaPi_CalaFiguera_23_july_2022	101.8	123.8	60.8
CalaPi_CalaFiguera_12_july_2023	155.9	112.8	57.5
PortoColom_CalaMillor_27_july_2021	278	117.4	55.9
PortoColom_CalaMillor_22_july_2022	101.4	123.9	60.8
PortoColom_CalaMillor_13_july_2023	101.7	113.1	57.5
Sur_Este_menorca_02_July_2021	101.2	130.5	66.7
Sur_Este_Menorca_30_July_2022	275.5	116.4	53.9
Sur_Este_Menorca_9_august_2023	273.9	119.9	52.4
Palma_7_july_2022	102.1	120.9	62.8
Palma_29_june_2023	101.3	123.1	64.8
Alcudia_25_July_2022	101	125.1	60.3
Alcudia_21_May_2020	103.7	119.7	58.7
Alcudia_22_July_2021	101.2	133.1	64.1
Alcudia_15_july_2023	111.9	113.2	56.6
Capdepera_20_july_2021	269.6	115.5	56.7
Capdepera_22_july_2022	101.4	123.9	60.8
Capdepera_31_july_2023	111.6	116.2	53.8
Este_Ibiza_18_july_2021	99.6	113.8	57.3
Este_Ibiza_24_july_2022	98.1	128.6	62.8
Este_Ibiza_14_july_2023	277.5	124	63.3
Pollença_6_July_2021	278.2	113.9	58.4
Pollença_23_May_2020	277.1	119.7	58.8
Pollença_21_July_2022	101.2	124.2	60.8
Pollença_27_june_2023	106.9	123.5	64.4
Oeste_Ibiza_18_july_2021	101.6	113.8	57.2
Oeste_Ibiza_14_july_2022	101.3	120.8	62.3
Oeste_Ibiza_19_july_2023	277.3	112.6	56.1

Continued on next page

Table E.1: Metadata of the satellite images used in the study. (Continued)

Banyalbufar_Soller_23_july_2021	176.3	116.8	56.5
Banyalbufar_Soller_17_july_2022	110.5	123	61.5
Banyalbufar_Soller_12_july_2023	101.5	124.5	63.2
Dragonera_Banyalbufar_29_july_2022	102.8	120.4	56.9
Dragonera_Banyalbufar_10_july_2021	110.5	113.7	58.1
Dragonera_Banyalbufar_16_july_2023	101.1	125.8	63.1
Norte_Ibiza_18_july_2021	101.6	113.8	57.2
Norte_Ibiza_18_july_2022	277.9	112.7	56.2
Norte_Ibiza_27_junio_2023	102	110.9	58.9
South_Cabrera_23_july_2022	101.8	123.3	60.9
South_Cabrera_29_june_2023	101.3	110.7	58.7
North_Cabrera_19_july_2022	101.7	112	55.5
North_Cabrera_14_july_2023	111.3	124	63.3

### E.1.2 Ground truth dataset composition

The original seabed cartography contains a total of 28 different classes, which were aggregated into 4 major ecological groups or habitat types (*Posidonia oceanica*, Other green plants, Brown algae & rocks and Sandy bottoms) based on feature similarity and ecological function (Table E.2). Although these habitat classes are present in the whole Mediterranean Sea, the specific composition of underlying sub-classes can vary among the different islands (e.g., one particular species of algae might be present in only one island, such as *Zostera noltii*).

## E.2 Dataset creation

Satellite imagery, along with habitat and bathymetry data, were integrated to create comprehensive training and testing datasets for our model. Initially, the Near-Infrared (NIR) band was utilized to eliminate land pixels through a clustering algorithm, namely K-means, as this band is not able to penetrate water beyond 1 or 2 meters. Subsequently, the NIR band was replaced with bathymetry information. These processed satellite images served as the primary input data for our model. The ground truth dataset, or labels, consisted of raster files mirroring the satellite images, with single-band values indicating the benthic class for each pixel. Construction of this dataset involved associating each pixel in the processed satellite images with a corresponding benthic class based on aggregated habitat data. Pixels lacking a class assignment were masked out in

**Table E.2:** Ecological Categories and Subcategories in the ground truth habitat data.

Category	Subcategory	Area (km <sup>2</sup> )	Presence zone	
<b>Posidonia oceanica</b>	Posidonia oceanica	538.61	Mallorca, Menorca, Ibiza, Formentera	
	Barrier reef of Posidonia oceanica	0.49	Mallorca, Menorca	
	Posidonia oceanica on stone with sand	20.33	Mallorca, Menorca	
	Mixed Posidonia oceanica with dead rhizome	0.10	Ibiza	
	Meadows of Posidonia oceanica on dead mat (rhizome)	5.23	Mallorca, Menorca	
<b>Other Green Plants</b>	Algae photophilic on stone with Posidonia oceanica	6.63	Mallorca, Menorca	
	Caulerpa prolifera	0.82	Mallorca, Menorca, Ibiza, Formentera	
	Meadows of phanerogams and green rhizomatous algae	4.87	Mallorca, Menorca	
	Fine sands with Cymodocea nodosa	1.95	Mallorca, Menorca, Ibiza, Formentera	
	Cymodocea nodosa	1.86	Mallorca, Menorca, Ibiza, Formentera	
	Zostera noltii	0.01	Menorca	
	Cymodocea nodosa and Zostera noltii	0.04	Menorca	
	Mixed meadows of Cymodocea nodosa and Caulerpa prolifera	4.62	Mallorca, Menorca, Ibiza	
	Muddy bays with red algae (Alisidium corralinum, Rytiphlaea tinctoria)	0.01	Menorca	
	<b>Sandy Bottoms</b>	Coarse sands	250.19	Ibiza, Formentera
Soft or sedimentary substrate		584.71	Mallorca, Menorca	
Mud		0.03	Ibiza	
Fine sands		74.36	Ibiza, Formentera	
Muddy detrital bottom		11.52	Ibiza	
Leptometra phalangium fields in bathyal bottoms of platform edge		208.85	Menorca	
Bathyal bottoms of platform edge with Gryphus vitreus		91.93	Menorca	
Medium sands		42.41	Ibiza, Formentera	
Muddy detrital bottoms infralittoral and circalittoral		589.81	Mallorca, Menorca, Formentera	
<b>Brown Algae and Rocks</b>		Rocky bottoms with photophilic algae and sands	36.49	Mallorca, Menorca
		Rocky bottoms dominated by sciafilic and hemisciafilic algae	21.17	Mallorca, Menorca, Ibiza, Formentera
		Cliffs, walls, and rocky slopes of the deep sea	13.96	Menorca
	Rocky bottoms with photophilic algae	14.35	Ibiza, Formentera	
	Rocky bottoms with photophilic algae and sands	15.58	Mallorca, Menorca	

both the satellite and label data. Similarly, pixels already masked in the satellite image were masked in the label image to ensure consistency.

Finally, patches of 256 × 256 pixels were created from each satellite and label image, forming the final dataset comprising up to 19369 patches. To study model performance in a real-case scenario, we trained our model with only data from the island of Mallorca (8488 patches, from which 1698 were used as validation set), while left as out-of-sample test set the data from the islands of Menorca, Ibiza, Formentera and Cabrera (7942 patches). We specifically call our test set “out-of-sample” test set to highlight the non-traditional way in which we test our model, which allow to test the extrapolation power and robustness of our model for real-case scenarios.

## E.3 Deep learning models

The performance of deep learning models can differ from one another due to its different architectures. Certain models can perform better than others in some situations, like segmenting specific classes or in images taken on different environmental conditions. Thus, we explored a selection of state-of-the-art deep learning models for semantic image segmentation such as UNET, Linknet, FPN and PSPNet (see below). Each model is formed by Convolutional Neural Network (CNN) blocks, designed to address specific challenges in semantic segmentation tasks. The specific architecture of the CNN blocks are usually

referred to as the “backbone” of the model. We tested 10 different backbone models for each deep learning model, leading to the training and evaluation of 40 models.

### E.3.1 UNET

UNET [532] is a popular architecture for image segmentation. It utilizes an encoder-decoder structure to capture high-level semantic information and detailed spatial features. The key feature of UNET is the inclusion of skip connections, enabling the model to propagate information from earlier layers to the corresponding decoder layers. This integration of skip connections allows the model to leverage both global context from the encoder and fine-grained spatial details from earlier layers, resulting in improved performance and accuracy. UNET’s design aims to learn hierarchical representations while preserving spatial information, making it effective for various computer vision tasks.

### E.3.2 Linknet

Linknet [533] follows an encoder-decoder structure and focuses on achieving a balance between accuracy and computational efficiency. Instead of traditional skip connections, Linknet incorporates “link blocks” to facilitate information flow between corresponding layers in the encoder and decoder paths. These link blocks consist of a shortcut connection and a residual connection, preserving and propagating important information during the upsampling process. Additionally, Linknet employs batch normalization and ReLU activation to enhance training convergence. The combination of skip connections, link blocks, and optimization techniques in Linknet results in improved information flow, enhanced spatial details, and reduced computational complexity.

### E.3.3 FPN

FPN (Feature Pyramid Network) [534] is an architecture was introduced to tackle object detection and semantic segmentation across different scales. It addresses the challenge of capturing multi-scale information by constructing a feature pyramid with varying levels of feature maps. The architecture comprises a bottom-up pathway that extracts high-level semantic features using a CNN, such as ResNet or VGG, and a top-down pathway that generates feature maps by upsampling and merging information from higher-resolution levels. FPN incorporates lateral connections to combine low-level and high-level features, facilitating the fusion of fine-grained spatial details and high-level semantic information. The resulting feature pyramid enables effective detection and classification of objects of different sizes, making it well-suited multi-scale analysis classification problems.

### E.3.4 PSPNet

PSPNet (Pyramid Scene Parsing Network) [535] is an architecture that utilizes a pyramid pooling module to capture contextual information at multiple scales. By dividing the input image into regions of varying sizes and aggregating global contextual information within each region, PSPNet improves the model's understanding of the scene and enhances object classification accuracy. The architecture comprises a convolutional neural network (CNN) backbone followed by the pyramid pooling module. This module performs pooling operations at different levels and spatial resolutions, capturing multi-scale context. The pooled features from each level are concatenated and passed to subsequent layers for classification. This enables PSPNet to effectively capture both local and global contextual information. The incorporation of the pyramid pooling mechanism enhances the model's comprehension of objects within the scene, ensuring robustness to variations in object scale and size.

### E.3.5 Backbones

- **ResNet34**: A variant of ResNet with 34 layers, introducing residual learning to mitigate the vanishing gradient problem.
- **ResNet152**: A deeper variant of ResNet with 152 layers, capable of capturing more complex features.
- **SeResNet152 (SE-ResNet152)**: Based on ResNet, This architecture incorporates a Squeeze-and-Excitation block for adaptive recalibration of channel importance.
- **ResNeXt101**: While the ResNet model makes use of many smaller paths, ResNeXt substitutes "groups" for this function. There are several parallel pathways in these groupings, and distinct features are learned via each path.
- **SeResNeXt101 (SE-ResNeXt101)**: Combines ResNeXt architecture with Squeeze-and-Excitation for enhanced feature representation.
- **DenseNet201**: A Dense Convolutional Network connecting each layer in a feed-forward fashion for improved parameter efficiency.
- **InceptionV3**: Part of the Inception family, using parallel convolutional operations for features at various scales.
- **InceptionResNetV2**: Extends InceptionV3 with residual connections for improved training convergence.
- **EfficientNetB7**: Part of the EfficientNet family, balancing accuracy and efficiency for state-of-the-art performance using "Compound Scaling" methods.
- **MobileNetV2**: Optimized for mobile and edge devices, using depthwise separable convolutions for efficiency, to create deep neural networks that

are lightweight and have minimal latency for embedded and mobile devices.

## E.4 Performance Metrics

The evaluation of the deep learning models in our study involves the use of various performance metrics to assess their effectiveness in segmenting seagrass habitats. These metrics provide insights into the models' accuracy, precision, recall, F1 score, Cohen's kappa and intersection over union (IoU). Each metric serves a specific purpose in evaluating different aspects of model performance.

### E.4.1 Accuracy

Accuracy measures the overall correctness of the model's predictions. It is calculated as the ratio of correctly predicted pixels to the total number of pixels in the dataset.

$$\text{Accuracy} = \frac{\text{True Positives} + \text{True Negatives}}{\text{Total Pixels}}$$

### E.4.2 Precision

Precision is the ratio of true positive predictions to the total predicted positives, indicating how well the model performs when it predicts a certain class. Higher precision values imply fewer false positives.

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

### E.4.3 Recall

Recall, also known as sensitivity or true positive rate, measures the ability of the model to capture all instances of a given class. It is calculated as the ratio of true positive predictions to the total actual positives.

$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

### E.4.4 F1 Score

The F1 score is the harmonic mean of precision and recall. It provides a balanced measure of a model's performance, especially when dealing with imbalanced datasets. A higher F1 score indicates better overall performance.

$$\text{F1 Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} = \frac{2 \times \text{TP}}{2 \times \text{TP} + \text{FN} + \text{FP}}$$

### E.4.5 Cohen's Kappa

Cohen's Kappa is a statistical measure of inter-rater agreement for categorical items. It is generally thought to be a more robust measure than simple percent agreement calculation, as Kappa takes into account the possibility of the agreement occurring by chance.

$$\kappa = \frac{p_o - p_e}{1 - p_e}$$

where  $p_o$  is the relative observed agreement among raters, and  $p_e$  is the hypothetical probability of chance agreement, using the observed data to calculate the probabilities of each observer randomly seeing each category.

### E.4.6 Intersection over Union (IoU)

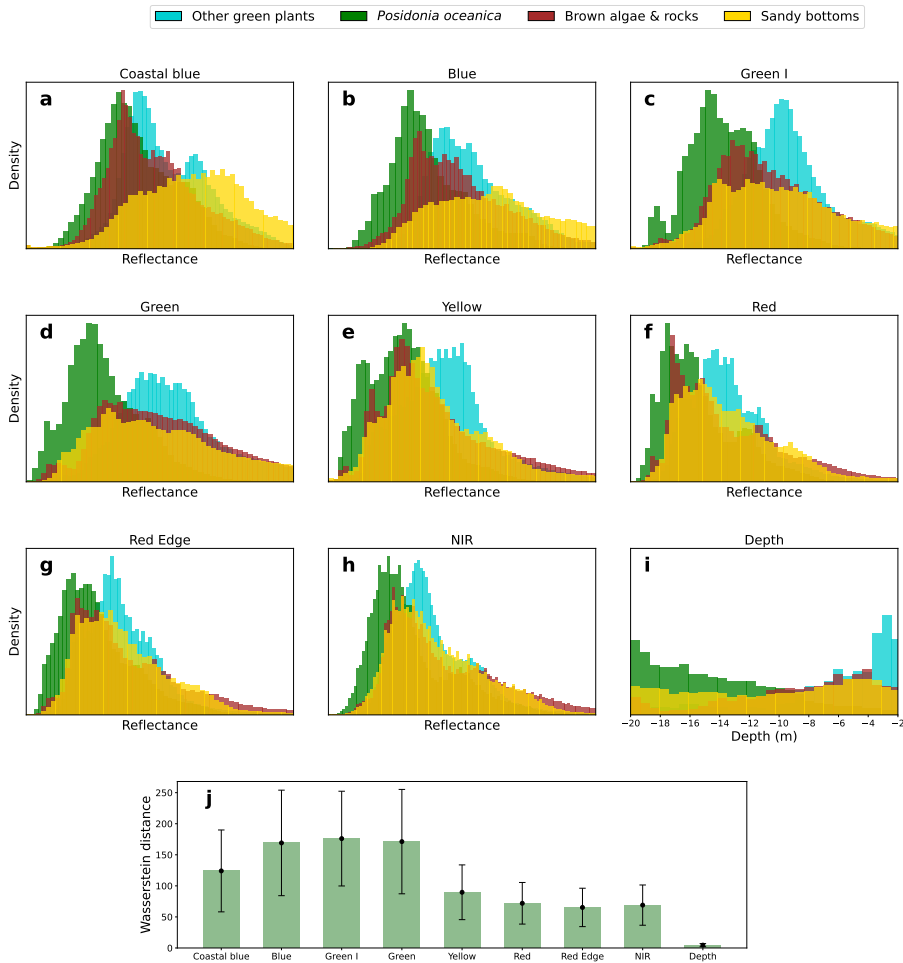
IoU, also known as Jaccard Index, measures the overlap between the predicted and true positive pixels. It is calculated as the ratio of the intersection of predicted (P) and label (L) pixels to the union of all pixels. IoU is particularly useful for semantic segmentation tasks, providing insights into the spatial accuracy of predictions.

$$\text{IoU} = \frac{|P \cap L|}{|P \cup L|} = \frac{|P \cap L|}{|P| + |L| - |P \cap L|}$$

## E.5 Spectral reflectance analysis

To gain a little understanding on our problem, we first perform a basic analysis on the spectral reflectance of the habitat classes. The distribution of reflectance values among the spectral bands plays a crucial role in the ability of the model to segment the different classes. When the reflectance values of the different classes do not overlap, the discrimination capability of the machine learning model is enhanced. Conversely, when reflectance values show significant overlap between classes, the model may encounter difficulties in categorizing pixels accurately, as the spectral information becomes less informative.

We computed the distribution of the response values of each habitat class to each of the satellite bands together with the mean Wasserstein distance among the distributions (Fig. E.2). The Blue, Green and Green 1 bands are highlighted as potentially more informative for the model, followed by Coastal Blue and Yellow. Red, Red Edge and NIR are identified as the bands with potentially less discrimination power. This could be expected from the fact that the wavelengths at the extreme ends of the visible spectrum are attenuated faster than those wavelengths in the middle. Finally, we observe that depth data is, by itself, basically uninformative.



**Figure E.2: Distribution of response values of different habitat classes.** Distribution of surface reflectance values for the different processed habitat classes with respect to the different bands available from the satellite imagery and depth.

Of course, this is just a simple, rather linear, analysis of the information provided to the model by the satellite images. In the end, the AI models will try to segment the different classes based on the representation of the response values in an abstract hyperspace. Something interesting is that depth has a very low discrimination power alone, but, we advance, observed that introducing the depth to the model highly improves model predictions. Thus, we can hypothesize that the model is somehow learning and applying a kind of depth reflectance correction to the input data to increase its accuracy.

## E.6 Architecture selection

We trained the 40 deep learning models as specified in the Methods section. After training, the performance on both train and validation sets was compared among the models. Models based on UNET and Linknet clearly outperformed models based on PSPNet and FPN architectures (Table E.3). Despite the similar performance of the models based on both UNET and Linknet architectures, UNET is a more complex architecture than Linknet, which translates into a bigger number of trainable parameters (Table E.4). Thus, models based on UNET architecture are more computationally expensive and, in addition, are more prone to suffer from over-fitting. Thus, we finally selected Linknet as the main architecture for our models.

**Table E.3:** Performance among all architectures for all backbones based on the Intersection over Union metric.

Architecture	UNET	Linknet	PSPNet	FPN
densenet201	91.15	<b>91.51</b>	86.93	87.97
resnet152	<b>92.01</b>	91.95	86.94	90.64
seresnext101	91.52	<b>91.54</b>	85.52	87.88
efficientnetb7	<b>89.27</b>	88.65	85.99	87.19
inceptionv3	89.49	<b>90.59</b>	85.13	89.33
seresnet152	91.76	<b>91.83</b>	86.49	90.48
inceptionresnetv2	91.85	<b>91.91</b>	86.16	89.61
resnext101	92.16	<b>92.5</b>	86.96	89.09
mobilenetv2	<b>89.35</b>	89.16	83.15	89.67
resnet34	<b>90.44</b>	90.14	86.12	90.05

## E.7 Out-of-sample predicting power and robustness

Because the performance of the different backbones can differ from one another, with some being better than others in certain situations, we implemented a pixel-wise consensus algorithm to enhance the robustness and reliability of model predictions. Basically, the results from all models were aggregated and the class label with the highest frequency across all predictions was assigned to the each pixel. Following this algorithm, we tested different aggregation strategies: selecting all the models (voting\_all), selecting only the best 3 or 5 models (voting\_top\_3, voting\_top\_5) and selecting the models that performed better in segmenting each class individually (voting\_specialists).

**Table E.4:** Number of total parameters among all architectures for all backbones (in millions)

Architecture	UNET	Linknet
resnet152	67.31	<b>63.53</b>
mobilenetv2	8.04	<b>4.14</b>
seresnet152	73.95	<b>70.17</b>
inceptionv3	29.93	<b>26.27</b>
densenet201	26.39	<b>22.56</b>
inceptionresnetv2	62.06	<b>57.87</b>
seresnext101	56.07	<b>52.29</b>
efficientnetb7	75.05	<b>72.26</b>
resnet34	24.47	<b>21.65</b>
resnext101	51.29	<b>47.52</b>

We compared the performance of our Linknet models and consensus strategies in both the training and out-of-sample test datasets. Regarding the training set, we observe that some models perform better than others, as expected. At an individual basis, inceptionresnetv2 and resnext101 are the best-performing backbones, while efficientnetb7 and densenet201 are the worst-performing ones (Table E.5). The consensus strategies cluster at the first positions of the ranking, with voting\_top\_3 and voting\_top\_5 occupying the two first positions. Interestingly, we observe that efficientnetb7 is the best-performing backbone in the test set, followed by voting\_all and inceptionresnetv2 (Table E.6). This suggests that efficientnetb7 is one of the most robust models, with higher generalization power. At any rate, we observe that it is closely followed by the voting\_all consensus strategy. Here we must emphasize that while the mask class is considered during training, it is omitted during evaluation and prediction phases to focus solely on the specified classes of interest. Consequently, this contributes to the observed differences in IoU values between Table E.3 and Tables E.5 and E.6.

Overall, we note that the relative performance of the model in the out-of-sample test set is significantly reduced in comparison with the training set, although the absolute performance is still considerable. This, indeed, is expected, as the habitat classes of the out-of-sample test are not formed by exactly the same sub-categories that in the training set. Furthermore, there is a huge environmental variability and intrinsic noise in the overall framework proposed here: the relative position of the satellite with respect to the sun and the earth at the time of each image acquisition, the specific atmospheric and marine conditions, the fact that some species forming the habitat classes are seasonal, etc.

**Table E.5:** Performance metrics for all models based on Linknet architecture in the training dataset

Backbone	IoU	f1	Kappa	Precision	Recall	Accuracy
<b>voting_top_3</b>	<b>89.57</b>	<b>94.01</b>	<b>80.72</b>	<b>94.79</b>	<b>94.37</b>	<b>94.28</b>
voting_top_5	89.01	93.64	80.12	94.59	94.01	93.9
inceptionresnetv2	88.95	93.72	81.35	94.3	93.99	93.9
resnext101	88.46	93.27	78.48	94.22	93.65	93.55
resnet152	88.34	93.23	78.72	94.12	93.62	93.54
voting_all	88.22	93.13	78.63	94.24	93.56	93.44
voting_specialists	88.22	93.21	79.27	94.2	93.57	93.46
mobilenetv2	88.01	93.16	80.74	93.78	93.37	93.28
resnet34	86.43	91.95	75.29	93.12	92.47	92.33
inceptionv3	85.96	91.74	77.56	92.88	92.09	92.04
seresnet152	85.74	91.54	76.36	92.94	91.91	91.79
seresnext101	85.65	91.55	76.2	92.75	91.92	91.77
efficientnetb7	84.56	90.79	72.95	91.86	91.35	91.19
densenet201	82.21	89.04	69.72	90.99	89.67	89.48
Median	88.12	93.16	78.56	93.95	93.47	93.36

**Table E.6:** Performance metrics for all models based on Linknet architecture in the out-of-sample test dataset

Backbone	IoU	f1	Kappa	Precision	Recall	Accuracy
<b>efficientnetb7</b>	<b>62.19</b>	<b>73.05</b>	<b>54.34</b>	76.1	74.78	74.62
<b>voting_all</b>	61.97	72.77	53.7	<b>76.74</b>	<b>75.02</b>	<b>74.88</b>
inceptionresnetv2	61.80	72.69	52.8	75.31	74.87	74.7
voting_top_5	61.35	72.38	53.11	76.29	74.38	74.25
voting_specialists	61.10	71.92	52.2	75.93	74.47	74.37
voting_top_3	60.88	71.92	52.03	75.94	74.08	73.94
mobilenetv2	60.87	72.01	51.74	74.40	74.12	73.97
seresnext101	60.58	71.82	52.02	75.01	73.57	73.35
resnext101	60.46	71.68	51.64	75.39	73.68	73.50
seresnet152	60.14	71.48	52.77	75.58	72.94	72.67
inceptionv3	59.65	70.82	50.66	73.81	72.85	72.65
densenet201	59.2	70.69	51.57	75.37	72.01	71.83
resnet152	57.58	69.16	48.89	73.04	71.00	70.91
resnet34	57.54	69.2	47.38	72.54	71.32	71.16
Median	60.73	71.87	52.03	75.38	73.88	73.72

### E.7.1 Understanding model performance

We then studied the performance of the model on segmenting each of the ecological habitats individually. We observed the emergence of “specialists”: in the training set, inceptionresnet outperforms in segmenting *Posidonia oceanica* meadows and Other green plants, resnext101 is better suited for mapping Sandy bottoms while resnet152 shows enhanced performance for detecting Brown algae & rocks (Table E.7).

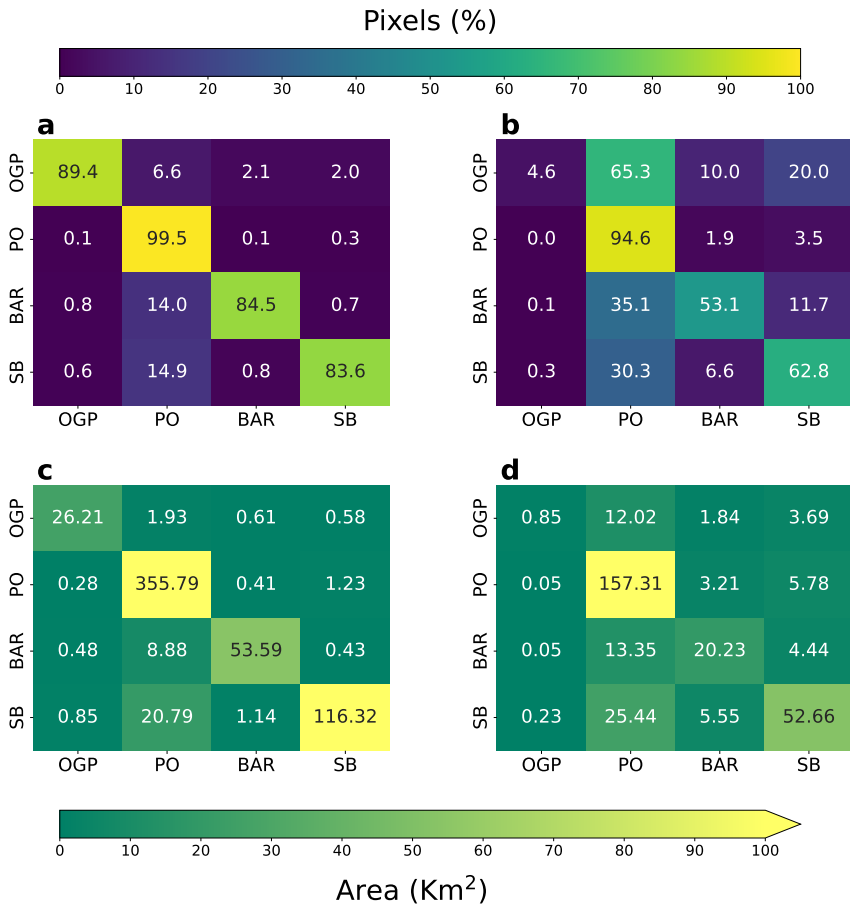
However, the results change when looking at the test dataset. This somehow explains why the voting\_all consensus method outperforms the voting\_smart one. We also note that the performance of the model with respect to Other green plants is drastically reduced in the test set, which undoubtedly affect the overall performance of the model shown in Table E.6.

**Table E.7:** Performance in segmenting each habitat class individually, measured by Intersection over Union, for all models in training and out-of-sample test datasets.

Method	Training dataset				Test dataset			
	PO	OGP	SB	BAR	PO	OGP	SB	BAR
densenet201	88.69	69.9	75.35	66.25	72.48	5.44	55.12	38.24
efficientnetb7	90.26	69.56	77.37	74.66	76.99	3.86	<b>56.31</b>	40.59
inceptionresnetv2	<b>92.61</b>	<b>85.81</b>	83.63	81.17	<b>77.34</b>	3.72	53.85	<b>41.17</b>
inceptionv3	90.8	73.6	80.82	75.57	75.79	3.79	50.99	38.21
mobilenetv2	92.05	83.6	81.49	81.44	76.76	<b>5.49</b>	53.16	37.27
resnet152	91.87	85.61	82.79	<b>81.7</b>	71.54	2.92	51.04	37.92
resnet34	90.8	82.7	80.78	75.84	73.11	3.37	49.71	35.7
resnext101	92.09	85.78	<b>83.65</b>	79.56	75.74	4.39	52.92	39.6
seresnet152	90.87	75.98	79.63	74.45	73.09	4.82	55.65	40.4
seresnext101	90.67	77.45	78.53	76.53	75.61	4.09	53.69	39.49
voting_all	91.97	85.29	82.2	81.54	77.30	<b>4.24</b>	<b>54.50</b>	<b>41.99</b>
voting_top_3	<b>92.79</b>	<b>87.28</b>	<b>84.72</b>	<b>83.01</b>	76.25	4.12	53.11	40.57
voting_top_5	92.5	86.1	83.54	82.53	76.18	4.13	54.44	41.71
voting_specialists	92.01	86.78	82.38	80.21	<b>77.31</b>	4.22	51.88	40.69
Area in (km <sup>2</sup> )	359.96	29.31	139.46	63.82	156.12	18.17	80.59	36.96

To further understand this drastic loss of performance on segmenting the “Other green plants” class, we computed the confusion matrix for the model (voting\_all from now on) predictions (Fig. E.3). In the training dataset there is already a significant difference between the performance in segmenting *Posidonia oceanica* meadows, with 99.5% True Positives (TP) and the other classes, with 89.4% TP for Other green plants, 84.5% TP for Brown algae & rocks

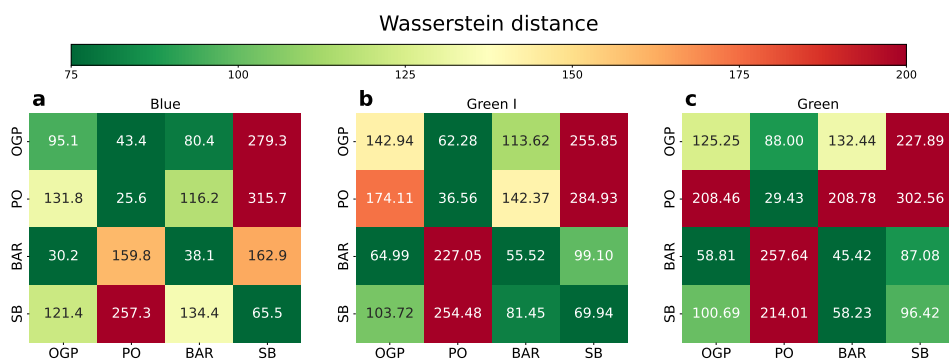
and 83.6% TP for Sandy Bottoms (Fig. E.3 a). At any rate, those differences clearly increase in the out-of-sample test set. While the segmentation of *Posidonia oceanica* meadows keeps a remarkable 94.6% TP, with a still low FP rate, the TP rate for the other classes is reduced (Fig. E.3 b). Specifically, a significant part of the pixels that were categorized by the ground truth data as Other green plants, Brown algae & rocks or Sandy bottoms are being classified by the model as *Posidonia oceanica* meadows. This is the reason for the very low IoU score of the Other green plant class in the testing set and the overall IoU decrease in all classes. In addition, we can also observe that the area of each class is significantly reduced in the test set in comparison with the training set (Fig. E.3c-d, sum over rows).



**Figure E.3: Confusion matrices for predicted classes and their area.** Confusion matrix for train (a) and out-of-sample test (b) dataset. Confusion matrix for the predicted area in train (c) and out-of-sample test (d) datasets.

To better understand this behavior, we compared the response values of each class between train and test datasets. In principle, the distribution of response values for a given class in the test set should be similar to the distribution of that class in the training set. This is indeed the rationale behind the use of correlative models. Of course, these distributions should be similar in a high dimensional space that we are not able to visualize. At any rate, to gain some understanding, we computed the Wasserstein distance between the distribution of response values in train and test for the Blue, Green and Green I bands, i.e., the most influential bands, as previously found.

The results revealed that the distribution of response values for the Other green plants class in the test dataset is much more similar to the *Posidonia oceanica* distribution of response values in the train set than to its own class (Fig. E.4), specially for Green I and Green bands (Fig. E.4 b-c). Then it is not surprising that the model classifies all those samples as *Posidonia oceanica*. Indeed, because of the great seasonality and natural variability of some of the species conforming the Other green plants class, the real habitat class present at the time of the satellite image acquisition could deviate from the initial one present at sonar-based classification. Similarly, there are some mixed sub-classes such as "Photophilic algae on stone with *Posidonia oceanica*" which further increase the probability of the previous argument taking place. That would explain the fact that the distribution of the response values for the "Other green plants" class in the test dataset is more similar to the *Posidonia oceanica* distribution in the training set than to its own class in the training set. So, in summary, it is plausible that some of the model predictions that are supposed to be incorrect, were indeed correct.



**Figure E.4: Wasserstein distance between Train and Test datasets.** Blue (a), Green I (b) and Green (c) bands.

All the analysis on model performance done up to this point were based on

**mean** results over the training and out-of-sample test datasets, which are based on 21 images each. Now we analyze the performance of the model predictions in each image individually (Table E.8). We observe that there is a great variability in model performance within each dataset. In the training set, 60% of the images are segmented with  $\text{IoU} > 90\%$  and 30% with  $\text{IoU} > 80\%$ .

**Table E.8:** Performance metrics for voting\_all model in segmenting individual images from training and testing datasets.

Image Name	IoU	f1	Kappa	Precision	Recall	Accuracy	Area (Km <sup>2</sup> )	Dataset
PortoColom_CalaMillor_22_july_2022	94.77	97.27	94.21	97.31	97.3	97.24	28.22	Train
Es_trenc_25_july_2022	93.56	96.53	87.1	96.81	96.57	96.54	12.65	Train
Capdepera_22_july_2022	93.23	96.43	94.04	96.52	96.47	96.37	15.89	Train
Pollença_6_july_2021	92.78	96.06	86.42	96.3	96.27	96.22	26.02	Train
CalaPi_CalaFiguera_23_july_2022	92.34	95.83	90.89	96.12	96.02	95.96	36.58	Train
Alcudia_25_july_2022	91.97	95.52	79.92	95.92	95.83	95.81	76.18	Train
Capdepera_20_july_2021	91.37	95.36	92.31	95.6	95.43	95.34	15.89	Train
Pollença_23_may_2020	91.12	94.99	81.0	95.44	95.39	95.33	24.7	Train
CalaPi_CalaFiguera_29_july_2021	90.87	94.95	90.16	95.47	95.15	95.09	36.53	Train
Banyalbufar_Soller_23_july_2021	90.38	94.89	90.71	95.07	94.94	94.81	8.88	Train
Dragonera_Toro_23_july_2021	90.06	94.63	90.71	94.82	94.76	94.5	11.93	Train
Dragonera_Toro_22_july_2022	90.05	94.65	90.62	94.84	94.77	94.49	11.87	Train
Pollença_21_july_2022	89.82	94.01	77.29	94.83	94.7	94.65	26.05	Train
PortoColom_CalaMillor_27_july_2021	88.91	93.99	86.49	94.44	94.14	94.08	28.22	Train
Alcudia_22_july_2021	86.06	91.52	63.17	93.14	92.51	92.47	75.69	Train
Alcudia_21_may_2020	84.62	90.44	58.93	92.69	91.34	91.3	73.63	Train
Dragonera_Banyalbufar_10_july_2021	83.97	91.04	85.27	92.17	91.11	90.87	7.09	Train
Formentera_24_july_2022	82.57	89.53	82.24	89.73	89.74	88.55	28.96	Test
Formentera_18_july_2021	82.54	89.52	82.27	89.9	89.84	88.68	28.99	Test
Palma_7_july_2022	82.42	90.03	81.53	91.59	90.1	88.99	57.4	Train
Sur_Oeste_Menorca_09_July_2021	82.12	88.25	58.03	87.84	89.5	89.55	10.87	Test
North_Cabrera_19_july_2022	77.78	86.24	49.14	87.42	85.54	85.54	1.37	Test
Sur_Oeste_Menorca_29_July_2022	75.14	84.15	36.85	88.46	81.84	80.29	16.4	Test
Sur_Ibiza_18_july_2021	70.41	79.92	51.71	80.02	81.73	82.14	21.47	Test
Sur_Ibiza_20_july_2022	70.11	79.64	50.29	79.84	81.55	81.97	22.52	Test
South_Cabrera_23_july_2022	64.89	77.3	47.39	81.76	76.02	76.0	1.45	Test
Banyalbufar_Soller_17_july_2022	64.49	76.47	54.73	84.17	79.5	79.4	8.83	Train
Dragonera_Banyalbufar_29_july_2022	63.14	76.1	62.49	85.53	76.27	75.88	7.25	Train
Este_Ibiza_18_july_2021	62.28	73.4	43.37	77.75	76.52	76.53	13.19	Test
Este_Ibiza_24_july_2022	61.53	72.59	42.15	77.8	75.62	75.67	13.21	Test
Oeste_Ibiza_14_july_2022	59.56	73.26	54.19	76.74	75.17	75.92	7.72	Test
Oeste_Ibiza_18_july_2021	58.9	72.63	53.26	76.26	74.69	75.59	7.75	Test
Sur_Este_Menorca_29_July_2022	58.0	69.34	57.21	67.93	74.34	75.16	16.15	Test
Sur_Este_Menorca_02_July_2021	57.01	68.53	55.82	71.64	73.55	74.38	16.42	Test
Norte_Ibiza_18_july_2022	53.01	68.57	53.45	70.57	69.52	69.21	4.27	Test
Norte_Ibiza_18_july_2021	50.03	66.05	49.13	71.43	66.85	67.13	4.4	Test
Norte_Menorca_20_July_2022	43.32	55.8	43.98	69.5	58.77	58.33	13.85	Test
Este_Menorca_19_July_2021	42.26	56.51	42.5	61.63	60.13	60.21	25.49	Test
Este_Menorca_15_July_2022	38.35	52.87	36.78	63.33	57.53	57.42	25.95	Test
Norte_Menorca_29_July_2021	38.28	51.3	37.52	67.92	54.69	54.37	14.23	Test

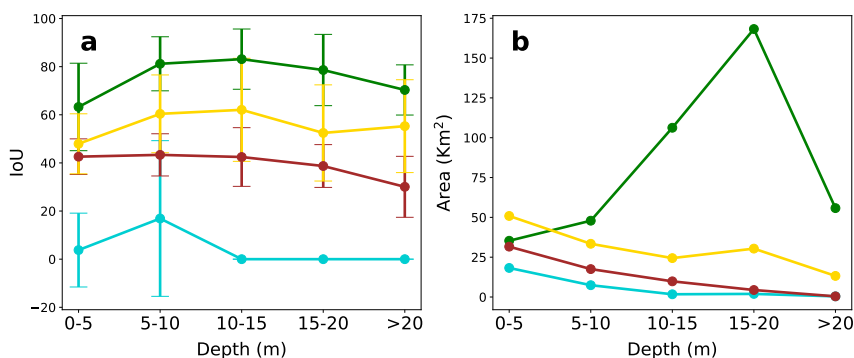
The resting 2 images, however, are segmented with “only” around 64% IoU score, overlapping with the out-of-sample test dataset (Table E.8). In the test set, 15% of the images are segmented with  $\text{IoU} > 80\%$ , 20% with  $\text{IoU} > 70\%$  and 45% with  $\text{IoU} > 50\%$ . The resting 4 images are segmented with an IoU score of around 40%. In any case, the accuracy of all segmentations are always

over 50%, much above the threshold for random segmentation of 25% (recall it is a 4-class classification problem).

These results present two interesting insights. First, the model is able to segment some images in the test set with very high performance (15% of the images). Recall this is not a usual testing set, but we have only trained the model in a given island (Mallorca) and tested it in 3 other islands (Menorca, Ibiza and Formentera), not controlling for neither weather or satellite-position variability. Second, we observe that 2 images in the training test are segmented with rather “low” performance, situating at the middle of the testing set in the performance ranking (Table E.8). Interestingly, those are images situated in the “Serra de Tramuntana”, the mountain range in Mallorca. The sea bottom in this area is characterized by a rocky substrate where depth increases rapidly. Indeed, we can observe that the images that are segmented with higher performance in the test set (Formentera\_24\_july\_2022 to Sur\_Ibiza\_20\_july\_2022) are precisely characterized by more sandy and clear sea bottoms with slowly increasing depths, while the rest are rather rocky with rapidly increasing depth (Table E.8). This insights can help in taking future steps to develop a general model for segmenting sea habitats in the whole Mediterranean sea.

## E.8 Effect of depth on model performance

Finally, we studied the effect of depth in model performance, focusing on the test set.



**Figure E.5: Model performance as function of depth.** (a) Performance of voting\_all model as function of the depth of classified pixels. (b) Area of samples contained in each depth range for the different habitat classes.

We observe that, in general, model performance is rather independent on depth (Fig. E.5 a). Interestingly, we observe that the *Posidonia oceanica* class consistently exhibits the highest IoU values across all depths compared to other

classes. This noteworthy trend may be attributed to the abundance of *Posidonia oceanica* samples (Fig. E.5 b), suggesting that the model’s performance excels in accurately identifying and delineating *Posidonia* in underwater environments.

**Table E.9:** Performance metrics for the final model in segmenting individual images from previous training and testing datasets.

Image Name	IoU	f1	Kappa	Precision	Recall	Accuracy	Area (Km <sup>2</sup> )	Dataset
Formentera_18_july_2021	98.50	99.24	98.50	99.25	99.25	99.25	29.07	Test
Formentera_24_july_2022	98.42	99.20	98.43	99.20	99.20	99.20	29.04	Test
Sur_Oeste_Menorca_09_July_2021	97.35	98.62	93.39	98.65	98.65	98.65	11.01	Test
Sur_Oeste_Menorca_29_July_2022	97.33	98.59	91.86	98.64	98.65	98.65	16.40	Test
Sur_Ibiza_20_july_2022	96.74	98.32	95.88	98.34	98.34	98.34	22.60	Test
Es_trenc_25_july_2022	96.70	98.28	92.34	98.31	98.31	98.31	12.65	Train
Sur_Ibiza_18_july_2021	96.69	98.29	95.99	98.31	98.31	98.31	21.67	Test
Norte_Menorca_20_July_2022	96.46	98.19	97.36	98.20	98.20	98.20	14.13	Test
Este_Menorca_19_July_2021	96.04	97.97	97.05	97.99	97.98	97.98	27.18	Test
CalaPi_CalaFiguera_29_july_2021	96.03	97.94	95.71	97.97	97.97	97.97	36.63	Train
Este_Menorca_15_July_2022	96.03	97.97	97.04	97.99	97.97	97.97	27.04	Test
CalaPi_CalaFiguera_23_july_2022	96.01	97.94	95.70	97.96	97.96	97.96	36.59	Train
PortoColom_CalaMillor_22_july_2022	95.85	97.85	95.50	97.87	97.87	97.87	28.22	Train
Sur_Este_Menorca_29_July_2022	95.80	97.84	96.59	97.85	97.85	97.85	16.54	Test
Sur_Este_Menorca_02_July_2021	95.79	97.83	96.58	97.85	97.84	97.84	16.52	Test
Norte_Menorca_29_July_2021	95.69	97.79	96.76	97.81	97.79	97.79	14.25	Test
Palma_7_july_2022	95.68	97.78	95.95	97.79	97.78	97.78	57.39	Train
PortoColom_CalaMillor_27_july_2021	95.55	97.69	95.15	97.71	97.71	97.71	28.22	Train
Alcudia_25_july_2022	95.26	97.47	89.41	97.57	97.57	97.57	76.19	Train
North_Cabrera_19_july_2022	95.25	97.48	89.80	97.50	97.53	97.53	1.38	Test
Alcudia_21_may_2020	94.88	97.23	87.03	97.35	97.37	97.37	73.63	Train
Alcudia_22_july_2021	94.80	97.21	88.23	97.34	97.33	97.33	75.70	Train
Capdepera_22_july_2022	94.64	97.19	95.32	97.22	97.22	97.22	15.90	Train
Capdepera_20_july_2021	94.58	97.15	95.28	97.19	97.19	97.19	15.92	Train
Este_Ibiza_24_july_2022	94.02	96.88	92.80	96.93	96.91	96.91	13.28	Test
Pollença_23_may_2020	94.01	96.80	89.11	96.89	96.91	96.91	24.73	Train
Este_Ibiza_18_july_2021	93.85	96.79	92.54	96.83	96.82	96.82	13.29	Test
Pollença_21_july_2022	93.41	96.46	88.34	96.59	96.59	96.59	26.05	Train
Oeste_Ibiza_14_july_2022	93.34	96.51	94.01	96.57	96.55	96.55	7.76	Test
Oeste_Ibiza_18_july_2021	93.34	96.51	94.00	96.57	96.56	96.56	7.77	Test
Pollença_6_july_2021	92.71	96.07	87.63	96.19	96.21	96.21	26.05	Train
Banyalbufar_Soller_23_july_2021	92.64	96.14	93.06	96.20	96.17	96.17	8.89	Train
Banyalbufar_Soller_17_july_2022	92.61	96.12	93.02	96.19	96.15	96.15	8.88	Train
Dragonera_Banyalbufar_29_july_2022	92.42	95.91	93.47	96.04	96.04	96.04	7.28	Train
Dragonera_Banyalbufar_10_july_2021	92.35	95.88	93.41	96.01	96.00	96.00	7.25	Train
Norte_Ibiza_18_july_2022	92.04	95.74	93.75	95.91	95.85	95.85	4.40	Test
South_Cabrera_23_july_2022	91.75	95.56	89.64	95.52	95.66	95.66	1.53	Test
Norte_Ibiza_18_july_2021	90.76	95.01	92.67	95.26	95.13	95.13	4.41	Test
Dragonera_Toro_23_july_2021	90.48	94.88	91.19	95.05	94.99	94.99	11.93	Train
Dragonera_Toro_22_july_2022	89.95	94.58	90.50	94.78	94.71	94.71	11.93	Train

## E.9 CAMELE trained with all available data

Finally, we trained the CAMELE model with all available data. The results are shown in Table E.9. We observe that the model is able to segment all images with an IoU score above 90%, with a mean, median and maximum IoU score of 94.64% 95.22% and 98.5%, respectively. This is a remarkable result, considering the high variability in the images. We assessed the model

robustness by predicting on a new set of images from the Balearic Islands in 2023 (Tables E.10 and E.11), not used in the training or out-of-sample test datasets. The model achieved a remarkable mean IoU score of 80% in this new dataset, which is encouraging for future applications of the model.

**Table E.10:** Performance metrics for the final model in segmenting individual images for the year 2023.

Image Name	IoU	f1	Kappa	Precision	Recall	Accuracy	Area (Km <sup>2</sup> )
Formentera_3_august_2023	93.78	96.77	93.76	96.78	96.77	96.77	29.04
Sur_Oeste_Menorca_9_august_2023	93.71	96.55	83.18	96.65	96.7	96.7	16.4
PortoColom_CalaMillor_13_july_2023	92.21	95.85	91.41	95.9	95.89	95.89	28.23
Sur_Ibiza_26_august_2023	90.07	94.56	86.7	94.62	94.69	94.69	22.55
North_Cabrera_14_july_2023	87.8	92.87	69.54	93.3	93.28	93.28	1.36
Pollença_27_june_2023	87.68	93.0	77.37	93.16	93.34	93.34	26.05
Capdepera_31_july_2023	87.3	92.93	88.4	93.2	93.17	93.17	15.88
South_Cabrera_29_june_2023	86.38	92.42	81.93	92.38	92.5	92.5	1.53
Sur_Este_Menorca_9_august_2023	85.02	91.69	87.11	92.36	91.93	91.93	16.57
Norte_Ibiza_27_june_2023	84.32	91.28	87.09	91.95	91.45	91.45	4.42
Alcudia_15_july_2023	81.08	87.67	50.55	90.25	89.27	89.27	76.2
Oeste_Ibiza_19_july_2023	80.96	89.21	80.97	90.04	89.59	89.59	7.74
Es_Trenc_15_july_2023	79.47	86.24	67.01	91.98	86.47	86.47	12.66
Este_Ibiza_14_july_2023	78.24	87.03	69.24	87.84	87.69	87.69	13.27
CalaPi_CalaFiguera_12_july_2023	76.17	85.29	69.68	87.56	86.37	86.37	36.52
Banyalbufar_Soller_12_july_2023	75.25	84.9	72.54	87.29	86.13	86.13	8.85
Dragonera_Banyalbufar_16_july_2023	72.61	83.72	73.37	86.49	83.65	83.65	7.17
Palma_29_june_2023	68.32	80.13	63.45	81.64	81.39	81.39	62.94
Norte_Menorca_23_june_2023	67.65	79.55	69.89	86.56	78.33	78.33	14.16
Este_Menorca_23_june_2023	66.3	79.37	69.65	83.97	79.37	79.37	27.04
Soller_Calobra_18_july_2023	44.43	60.2	40.03	66.68	60.43	60.43	6.0

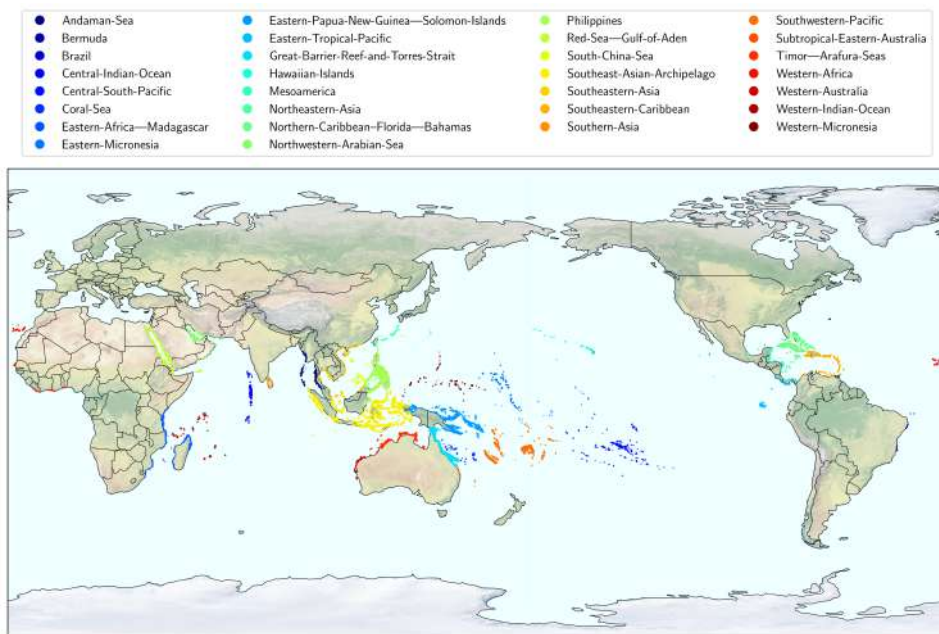
**Table E.11:** Average performance metrics for the final model with 2023 images.

Method	IoU	f1	Kappa	Precision	Recall	Accuracy
voting_top_10	79.92	87.63	71.69	89.46	88.26	88.26
voting_top_5	79.38	87.24	71.14	88.95	87.87	87.87
resnet34	79.29	87.28	71.34	88.56	87.71	87.71
voting_top_smart	79.29	87.2	70.49	89.13	87.83	87.83
voting_top_3	79.27	87.11	70.95	88.88	87.74	87.74
inceptionresnetv2	79.17	87.2	70.96	88.42	87.75	87.75
mobilenetv2	77.76	86.34	69.05	87.46	86.84	86.84
resnext101	77.61	85.9	68.89	87.7	86.58	86.58
seresnet152	77.31	86.01	68.99	87.31	86.51	86.51
inceptionv3	77	85.71	68.53	87.22	86.13	86.13
resnet152	76.81	85.43	67.79	87.3	86.07	86.07
seresnext101	76.52	85.23	67.32	86.95	85.98	85.98
efficientnetb7	75.22	84.42	65.69	85.85	85.18	85.18
densenet201	74.85	84.18	65.63	86.19	84.63	84.63

## **F. Coral reefs' macroecological patterns**

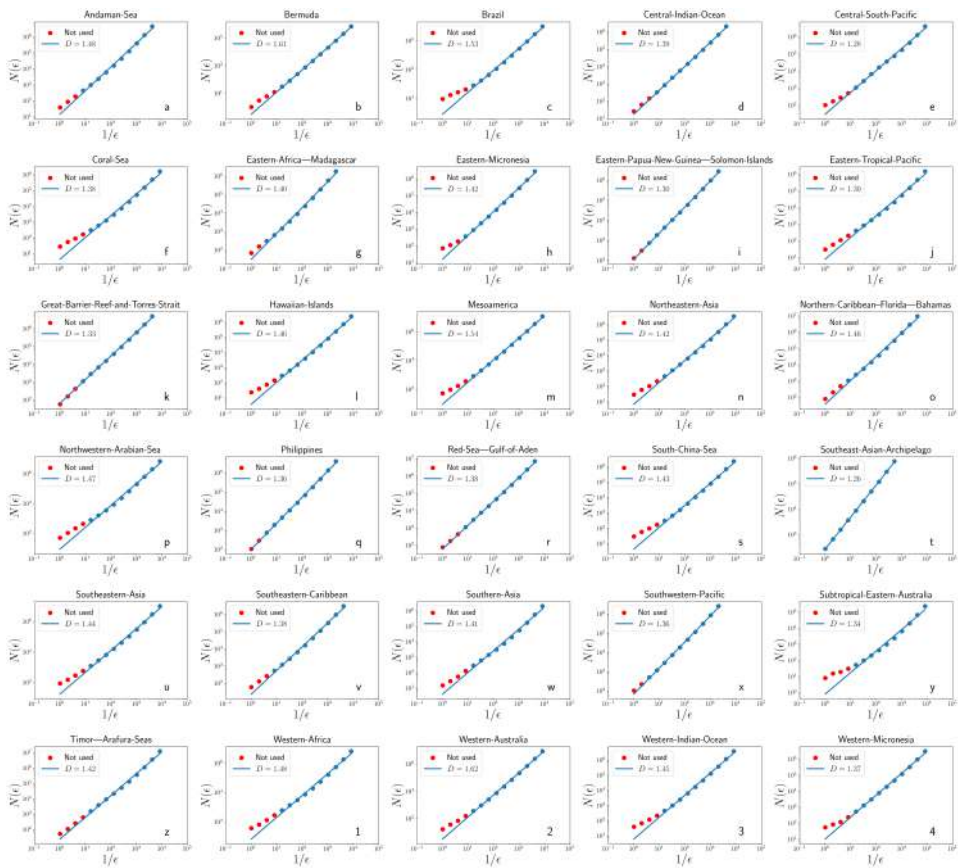
## F.1 Macroecological patterns & fractal dimension

Fig. F.1 shows the location of the 30 coral reef provinces considered in this study. The provinces are distributed across the world's oceans, with the largest number of reefs located in the Southeast Asian Archipelago and the smallest in the Subtropical Eastern Australia province (Table F.1).



**Figure F.1: Geographic location of the studied coral reefs.** Individual reefs described in this study are plotted in a world Atlas, with a different color for each of the 30 provinces.

Table F.1 shows the number of reefs, total reef area, total reef perimeter, mean reef area, mean reef perimeter, and fractal dimension of each province. The fractal dimension was computed using the box-counting method, as shown in Fig. F.2. The fractal dimension of the provinces ranges from 1.26 to 1.62, with the Bermuda province having the highest fractal dimension and the Central South Pacific province having the lowest. The fractal dimension of the coral reef provinces indicates that the distribution of coral reefs is self-similar at different scales, with the same patterns repeating at different levels of magnification. This self-similarity is a characteristic of fractal objects and is consistent with the idea that coral reefs are complex systems that exhibit scaling behavior across different levels of organization.



**Figure F.2: Computation of the fractal dimensions of each coral province using the box-counting method.**  $N(\epsilon)$  corresponds to the number of boxes of length  $\epsilon$  needed to cover the object. The slope of  $\log(N(\epsilon))$  vs  $1/\epsilon$  approximates the fractal dimension of the object. Blue points were used to perform this computation while red dots were discarded.

**Table F.1:** Number of reefs, total reef area, total reef perimeter, mean reef area, mean reef perimeter and fractal dimension of coral reef provinces.

Province	N° reefs	Total reef area (km <sup>2</sup> )	Total reef perimeter (km)	Mean reef area (ha)	Mean reef perimeter (km)	Fractal dimension
Andaman-Sea	2.46e+04	2458.49	82370.95	9.99	3.35	1.48
Bermuda	3.70e+03	82.51	4619.18	2.23	1.25	1.61
Brazil	5.20e+03	117.61	11040.94	2.26	2.12	1.53
Central-Indian-Ocean	3.11e+04	937.56	48889.55	3.02	1.57	1.39
Central-South-Pacific	3.04e+04	509.03	31947.16	1.68	1.05	1.28
Coral-Sea	5.53e+03	204.5	13578.79	3.7	2.46	1.38
Eastern-Africa—Madagascar	5.83e+04	3168.68	104343.91	5.44	1.79	1.4
Eastern-Micronesia	2.45e+04	1445.35	54452.86	5.91	2.23	1.42
Eastern-Papua-New-Guinea—Solomon-Islands	1.22e+05	3643.85	217711.52	2.98	1.78	1.3
Eastern-Tropical-Pacific	9.46e+03	201.48	10588.55	2.13	1.12	1.3
Great-Barrier-Reef-and-Torres-Strait	1.54e+05	1952.24	127829.58	1.27	0.83	1.33
Hawaiian-Islands	1.33e+04	309.32	17102.41	2.33	1.29	1.46
Mesoamerica	5.08e+04	1644.13	76462.11	3.24	1.51	1.54
Northeastern-Asia	1.44e+04	438.56	20612.65	3.05	1.43	1.42
Northern-Caribbean—Florida—Bahamas	1.12e+05	4475.7	179093.5	4.01	1.61	1.46
Northwestern-Arabian-Sea	3.13e+04	873.91	46168.12	2.79	1.47	1.47
Philippines	1.47e+05	6231.33	263475.43	4.25	1.8	1.36
Red-Sea—Gulf-of-Aden	1.64e+05	2876.87	186036.16	1.76	1.14	1.38
South-China-Sea	1.22e+04	302.18	16809.32	2.48	1.38	1.43
Southeast-Asian-Archipelago	2.59e+05	8941.66	426626.12	3.45	1.65	1.26
Southeastern-Asia	3.60e+04	1490.35	51462.3	4.14	1.43	1.44
Southeastern-Caribbean	3.67e+04	1343.08	71308.15	3.66	1.94	1.38
Southern-Asia	6.78e+03	289.07	13465.78	4.27	1.99	1.41
Southwestern-Pacific	9.54e+04	3670.4	164888.22	3.85	1.73	1.36
Subtropical-Eastern-Australia	5.66e+02	31.9	914.1	5.64	1.62	1.34
Timor—Arafura-Seas	5.02e+04	1743.77	77897.78	3.47	1.55	1.42
Western-Africa	2.09e+04	960.09	29485.25	4.6	1.41	1.48
Western-Australia	2.36e+04	1174.42	49788.67	4.97	2.11	1.62
Western-Indian-Ocean	2.36e+04	502.17	29976.67	2.13	1.27	1.45
Western-Micronesia	1.47e+04	403.21	24713.6	2.73	1.68	1.37

## F.2 Coral reef size distribution

Table F.2 shows the statistical comparison of the power-law fit of the coral reef size distribution to other plausible distributions such as Exponential, Stretched exponential, Lognormal, Lognormal positive, or Truncated power-law. A negative ratio implies that the compared distribution is more plausible than the assumed power-law, while a positive ratio indicates that the assumed power-law distribution is a better fit. The p-value measures the statistical significance of the result of each comparison, e.g.,  $p < 0.05$  is assumed to be significant. The results show that the power-law (or truncated power-law) distribution is a better fit for most of the provinces.

**Table F.2:** Statistical comparison of power-law fit of the coral reef size distribution to other plausible distributions.

		Exponential	Stretched Exponential	Lognormal	Lognormal Positive	Truncated Power Law
Andaman-Sea	p-value	1.24e-23	2.93e-03	6.21e-01	2.70e-02	6.03e-02
	Ratio	7958.48	24.72	-0.40	13.52	-1.76
Bermuda	p-value	4.79e-05	7.81e-03	8.64e-01	1.51e-02	9.96e-01
	Ratio	2768.75	14.95	-0.01	12.44	0.00
Brazil	p-value	6.00e-31	1.73e-04	4.95e-01	3.31e-04	2.44e-01
	Ratio	3824.34	22.83	-0.38	19.27	-0.68
Central-Indian-Ocean	p-value	5.49e-151	1.65e-02	1.47e-08	6.27e-04	0.00e+00
	Ratio	17443.86	-29.98	-41.40	-35.42	-63.45
Central-South-Pacific	p-value	5.02e-62	8.63e-02	1.87e-01	2.72e-01	6.31e-07
	Ratio	3649.17	10.17	-1.98	5.02	-12.41
Coral-Sea	p-value	1.98e-98	5.67e-04	5.65e-01	5.47e-04	1.50e-07
	Ratio	4541.59	20.22	-0.31	17.95	-13.79
Eastern-Africa—Madagascar	p-value	9.13e-91	1.72e-25	6.28e-02	7.04e-20	7.31e-03
	Ratio	40940.17	201.41	0.33	143.46	-3.60
Eastern-Micronesia	p-value	1.27e-80	4.21e-05	9.93e-01	1.46e-03	9.47e-03
	Ratio	5767.13	28.85	0.00	16.84	-3.37
Eastern-Papua-New-Guinea—Solomon-Islands	p-value	6.01e-09	3.65e-01	1.62e-01	1.62e-01	6.97e-03
	Ratio	231.98	-2.48	-2.68	-2.68	-3.64
Eastern-Tropical-Pacific	p-value	1.57e-13	9.32e-01	1.42e-01	3.59e-01	1.19e-02
	Ratio	1655.07	-0.42	-3.42	-2.98	-3.16
Great-Barrier-Reef-and-Torres-Strait	p-value	1.11e-29	1.98e-01	1.14e-04	1.98e-01	9.99e-10
	Ratio	16943.45	20.29	-23.03	-13.15	-18.66
Hawaiian-Islands	p-value	1.19e-12	1.29e-01	1.31e-01	4.65e-01	4.88e-03
	Ratio	6041.76	12.22	-3.46	4.49	-3.96
Mesoamerica	p-value	3.43e-75	2.80e-09	5.84e-01	1.04e-06	9.24e-06
	Ratio	16356.74	71.41	-0.40	46.78	-9.83
Northeastern-Asia	p-value	3.38e-71	4.70e-04	1.68e-06	3.30e-06	0.00e+00
	Ratio	6360.31	-29.02	-29.21	-29.19	-41.23
Northern-Caribbean—Florida—Bahamas	p-value	6.32e-69	3.44e-14	8.90e-01	7.65e-09	9.09e-03
	Ratio	29091.42	135.67	-0.04	72.55	-3.40
Northwestern-Arabian-Sea	p-value	3.91e-78	2.60e-17	8.19e-01	9.45e-16	4.98e-02
	Ratio	48828.62	183.48	0.06	162.94	-1.92
Philippines	p-value	1.39e-286	1.26e-13	1.82e-05	8.97e-07	0.00e+00
	Ratio	71195.32	188.54	-24.29	97.32	-77.33
Red-Sea—Gulf-of-Aden	p-value	4.17e-42	2.35e-01	9.45e-04	2.20e-01	2.37e-06
	Ratio	15470.31	17.44	-18.59	-11.52	-11.13
South-China-Sea	p-value	1.61e-76	3.76e-06	8.31e-01	4.34e-05	6.40e-04
	Ratio	7335.20	36.79	-0.06	28.19	-5.83
Southeast-Asian-Archipelago	p-value	0.00e+00	8.93e-02	3.49e-30	1.11e-12	0.00e+00
	Ratio	88730.10	-54.99	-182.78	-158.78	-215.45
Southeastern-Asia	p-value	1.59e-78	2.63e-01	1.50e-02	8.40e-01	4.98e-11
	Ratio	6729.44	9.15	-6.99	-1.16	-21.59
Southeastern-Caribbean	p-value	3.46e-67	1.01e-04	2.21e-06	2.21e-06	0.00e+00
	Ratio	4468.22	-29.16	-28.16	-28.16	-40.67
Southern-Asia	p-value	8.59e-26	1.00e-01	5.96e-01	2.37e-01	5.18e-02
	Ratio	2295.95	7.27	-0.39	4.18	-1.89
Southwestern-Pacific	p-value	0.00e+00	2.38e-46	5.12e-01	1.72e-43	0.00e+00
	Ratio	103608.60	387.77	-0.57	360.89	-43.43
Subtropical-Eastern-Australia	p-value	8.79e-11	9.90e-01	3.70e-01	9.48e-01	4.65e-02
	Ratio	1047.13	0.03	-1.08	0.17	-1.98
Timor—Arafura-Seas	p-value	1.70e-97	1.89e-07	2.10e-01	2.46e-05	3.22e-07
	Ratio	20686.45	70.02	-2.18	46.00	-13.06
Western-Africa	p-value	1.45e-14	2.66e-09	2.64e-01	8.27e-07	9.70e-01
	Ratio	19217.88	97.00	0.10	58.74	-0.00
Western-Australia	p-value	2.29e-15	6.49e-07	6.90e-01	1.80e-05	2.90e-01
	Ratio	14180.60	57.94	0.04	38.93	-0.56
Western-Indian-Ocean	p-value	1.57e-43	3.21e-11	7.65e-01	7.13e-10	5.82e-02
	Ratio	22068.17	98.19	0.05	79.29	-1.79
Western-Micronesia	p-value	1.12e-41	3.45e-03	1.00e-03	1.00e-03	4.14e-12
	Ratio	1432.37	-13.38	-12.08	-12.08	-24.03

Table F.3 shows the results of the power-law fit of the coral reef size distribution.  $\alpha$  is the exponent of the distribution,  $D$  is the Kolmogorov-Smirnov distance (fit error),  $\sigma$  is the standard error of the exponent,  $x_{\max}$  is the optimal maximum value for the power-law fit, and  $x_{\min}$  is the optimal minimum value for the power-law fit.

**Table F.3:** Results of the power-law fit of the coral reef size distribution.  $\alpha$  is the exponent of the distribution,  $D$  is the Kolmogorov-Smirnov distance (fit error),  $\sigma$  is the standard error of the exponent,  $x_{\max}$  is the optimal maximum value for the power-law fit and  $x_{\min}$  is the optimal minimum value for the power-law fit

	$\alpha$	$D$	$\sigma$	$x_{\max}$	$x_{\min}$
<b>Eastern-Micronesia</b>	1.80	1.60e-02	1.23e-02	3.58e+07	1.76e+04
<b>Southeast-Asian-Archipelago</b>	1.78	1.22e-02	2.72e-03	6.93e+07	6.95e+03
<b>Coral-Sea</b>	1.65	2.48e-02	1.16e-02	6.34e+06	2.05e+03
<b>Northwestern-Arabian-Sea</b>	1.83	8.60e-03	5.11e-03	6.11e+07	1.25e+03
<b>Western-Indian-Ocean</b>	1.86	8.58e-03	7.21e-03	2.97e+07	1.88e+03
<b>Southeastern-Asia</b>	1.77	8.28e-03	9.57e-03	4.29e+07	1.47e+04
<b>Northern-Caribbean-Florida-Bahamas</b>	1.86	5.43e-03	6.15e-03	2.35e+08	1.21e+04
<b>Hawaiian-Islands</b>	1.84	8.31e-03	1.19e-02	4.34e+07	4.12e+03
<b>Great-Barrier-Reef-and-Torres-Strait</b>	1.98	8.10e-03	6.68e-03	6.27e+07	1.01e+04
<b>Red-Sea-Gulf-of-Aden</b>	1.95	1.02e-02	7.27e-03	6.92e+07	1.63e+04
<b>Mesoamerica</b>	1.83	7.59e-03	7.28e-03	7.67e+07	7.70e+03
<b>Subtropical-Eastern-Australia</b>	1.61	2.62e-02	2.55e-02	8.35e+06	1.01e+03
<b>Southeastern-Caribbean</b>	1.81	1.80e-02	1.03e-02	2.63e+07	2.00e+04
<b>Central-South-Pacific</b>	1.86	1.24e-02	1.29e-02	8.58e+06	1.02e+04
<b>Western-Africa</b>	1.93	1.18e-02	1.00e-02	2.40e+08	3.64e+03
<b>Bermuda</b>	2.06	1.95e-02	2.77e-02	3.44e+07	3.18e+03
<b>Andaman-Sea</b>	1.78	8.19e-03	1.14e-02	3.16e+08	1.89e+04
<b>Western-Australia</b>	1.83	1.00e-02	9.50e-03	2.42e+08	5.51e+03
<b>Central-Indian-Ocean</b>	1.71	1.27e-02	5.83e-03	1.37e+07	3.11e+03
<b>Timor-Arafura-Seas</b>	1.79	7.21e-03	6.47e-03	5.41e+07	5.88e+03
<b>Brazil</b>	1.82	1.39e-02	1.66e-02	9.54e+06	2.27e+03
<b>Eastern-Tropical-Pacific</b>	1.93	1.64e-02	2.04e-02	1.02e+07	1.06e+04
<b>South-China-Sea</b>	1.82	1.47e-02	1.06e-02	1.02e+07	3.21e+03
<b>Northeastern-Asia</b>	1.72	1.86e-02	8.95e-03	1.27e+07	4.10e+03
<b>Eastern-Africa-Madagascar</b>	1.79	1.23e-02	5.36e-03	2.07e+08	4.70e+03
<b>Western-Micronesia</b>	1.82	1.89e-02	1.65e-02	7.28e+06	1.70e+04
<b>Southern-Asia</b>	1.78	1.44e-02	1.90e-02	1.99e+07	8.45e+03
<b>Philippines</b>	1.75	5.20e-03	3.32e-03	1.17e+08	5.41e+03
<b>Southwestern-Pacific</b>	1.74	7.58e-03	3.01e-03	5.55e+07	2.05e+03
<b>Eastern-Papua-New-Guinea-Solomon-Islands</b>	2.41	2.20e-02	4.28e-02	2.92e+07	5.38e+05

