



**HAL**  
open science

# Beyond the limits of the unassigned protist microbiome: inferring large-scale spatio-temporal patterns of Syndiniales marine parasites

Iris Rizos, Pavla Debeljak, Thomas Finet, Dylan Klein, Sakina-Dorothee  
Ayata, Fabrice Not, Lucie Bittner

## ► To cite this version:

Iris Rizos, Pavla Debeljak, Thomas Finet, Dylan Klein, Sakina-Dorothee Ayata, et al.. Beyond the limits of the unassigned protist microbiome: inferring large-scale spatio-temporal patterns of Syndiniales marine parasites. *ISME Communications*, 2023, 3 (1), pp.16. 10.1038/s43705-022-00203-7 . hal-04230262

**HAL Id: hal-04230262**

**<https://hal.science/hal-04230262v1>**

Submitted on 5 Oct 2023

**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.

## ARTICLE OPEN



# Beyond the limits of the unassigned protist microbiome: inferring large-scale spatio-temporal patterns of Syndiniales marine parasites

Iris Rizos<sup>1,2✉</sup>, Pavla Debeljak<sup>1</sup>, Thomas Finet<sup>1</sup>, Dylan Klein<sup>1</sup>, Sakina-Dorothee Ayata<sup>3</sup>, Fabrice Not<sup>2</sup> and Lucie Bittner<sup>1,4</sup>

© The Author(s) 2023

Marine protists are major components of the oceanic microbiome that remain largely unrepresented in culture collections and genomic reference databases. The exploration of this uncharted protist diversity in oceanic communities relies essentially on studying genetic markers from the environment as taxonomic barcodes. Here we report that across 6 large scale spatio-temporal planktonic surveys, half of the genetic barcodes remain taxonomically unassigned at the genus level, preventing a fine ecological understanding for numerous protist lineages. Among them, parasitic Syndiniales (Dinoflagellata) appear as the least described protist group. We have developed a computational workflow, integrating diverse 18S rDNA gene metabarcoding datasets, in order to infer large-scale ecological patterns at 100% similarity of the genetic marker, overcoming the limitation of taxonomic assignment. From a spatial perspective, we identified 2171 unassigned clusters, i.e., Syndiniales sequences with 100% similarity, exclusively shared between the Tropical/Subtropical Ocean and the Mediterranean Sea among all Syndiniales orders and 25 ubiquitous clusters shared within all the studied marine regions. From a temporal perspective, over 3 time-series, we highlighted 39 unassigned clusters that follow rhythmic patterns of recurrence and are the best indicators of parasite community's variation. These clusters withhold potential as ecosystem change indicators, mirroring their associated host community responses. Our results underline the importance of Syndiniales in structuring planktonic communities through space and time, raising questions regarding host-parasite association specificity and the trophic mode of persistent Syndiniales, while providing an innovative framework for prioritizing unassigned protist taxa for further description.

ISME Communications; <https://doi.org/10.1038/s43705-022-00203-7>

## INTRODUCTION

The advances in high-throughput sequencing technologies have provided a new perspective to microbial diversity at a global scale. Studying the DNA of environmental microbial communities (i.e., microbiome) allowed to overcome the limit of non-cultivability and provided access to an unprecedented large quantity of high resolution genetic information [1–4]. In silico downstream analysis of genetic big-data shed light on a new challenge in microbial ecology: exploring the unassigned microbiome [5, 6]. From the point of view of environmental genomics, the unassigned microbiome encompasses all genetic sequences that cannot be annotated with referenced biological information as they have no match in databases at a functional [6–9] and/or taxonomic level [5, 10, 11]. Recent studies have pointed out that unassigned sequences contribute to 25–58% of microbial communities' diversity observed across a variety of aquatic and soil ecosystems [4, 5, 12]. In the marine realm, large scale sequencing studies have revealed that the unassigned microbiome represents half of the functional diversity (including samples enriched in viruses, prokaryotes and protists) [13]. In terms of taxonomic diversity, the unassigned protist microbiome, defined as taxa with V9 regions of

the 18S rDNA marker having a sequence similarity <80% with reference sequences, represents ~30% at the supergroup level [14].

In marine metabarcoding studies, Syndiniales (a clade of marine alveolates, MALVs [15]) represent an ubiquitous and hyperdiverse lineage of protistan endoparasites [16–18]. Syndiniales are distributed worldwide from tropical and temperate zones [14, 19] to both arctic and antarctic poles [20, 21]. Their unexpected contribution to protist community composition has been revealed by metabarcoding studies both in open sea and coastal environments, with Syndiniales being the third most abundant lineage of the circumglobal Tara Oceans expedition [14] and representing up to 11% of community's abundance in fjordic-bays [21] and 28% at a North-Atlantic river estuary [22]. Accumulating observations and correlations of metabarcoding data support that Syndiniales are opportunistically infecting a wide spectrum of hosts, including other protists (dinoflagellates, ciliates, radiolarians) but also metazoans (e.g., crustaceans) [22, 23]. Their wide abundance and distribution confers them global ecological importance for microbial food webs and biogeochemical cycling, by regulating host populations [22, 24, 25] and supplying the microbial loop with organic

<sup>1</sup>Institut de Systématique, Evolution, Biodiversité (ISYEB), Muséum National d'Histoire Naturelle, CNRS, Sorbonne Université, EPHE, Université des Antilles, Paris, France. <sup>2</sup>Sorbonne Université, CNRS, AD2M-UMR7144 Station Biologique de Roscoff, 29680 Roscoff, France. <sup>3</sup>Sorbonne Université, Laboratoire d'Océanographie et du Climat: Expérimentation et Analyses Numériques (LOCEAN, SU/CNRS/IRD/MNHN), 75252 Paris Cedex 05, France. <sup>4</sup>Institut Universitaire de France, Paris, France. ✉email: iris.rizos@sb-roscoff.fr

Received: 27 July 2022 Revised: 15 November 2022 Accepted: 16 November 2022

Published online: 28 February 2023

matter [26]. Syndiniales are notoriously known for having an ecological impact at the entire host population level, by being responsible for the collapse of red tide dinoflagellate blooms [27] and mass mortality of fish larvae [28]. Yet, the great majority of Syndiniales remain uncultivable and show a high degree of divergence in genomic sequences [29]. A recent study in an estuary revealed the existence of at least 8 cryptic Syndiniales species, among which 6 could be differentiated by the V4 region of the 18S marker by a 100% sequence similarity threshold [17]. Five groups of Syndiniales have been described (MALVs I-V) [19]. The core Syndiniales diversity is encompassed within groups II and I, while the majority of in vitro studies have been conducted on the genus *Amoebophrya* (MALV II) [17, 29, 30]. Their complex lifestyle, small size (0.2–20 µm) and lack of distinctive morphological features makes Syndiniales' description a laborious process relying on designing specific probes for in situ hybridization [15, 24, 25, 31]. Thus, Syndiniales diversity still remains a blackbox in protistology [22, 25, 32], rendering the ecological understanding of these widespread microorganisms below the order level presently beyond reach [17].

In this study, we explored marine planktonic protist communities at a wide spatio-temporal scale, in order to: (i) quantify the taxonomically unassigned sequences and reveal protist lineages for which there is a major scarcity of taxonomic references, (ii) highlight unassigned protist diversity shared between contrasted marine environments and (iii) identify unassigned taxa which are ecologically relevant and recurrent, that should be prioritized for further characterization. We integrated 12 years of data and 155 different sampling locations from 6 environmental metabarcoding datasets, combining 3 coastal time-series (ASTAN, BBMO, SOLA), 1 European coastal Sea sampling project (BioMarKs) and 2 oceanographic campaigns (Malaspina, MOOSE). As a study case, we focused our analyses on the parasite group of Syndiniales and, by clustering the gathered metabarcodes in a Sequence Similarity Network (SSN), we revealed novel ecological patterns of Syndiniales at a taxonomic resolution of 100% similarity between V4 regions of the 18S rDNA marker.

## RESULTS

### Diversity and abundance of taxonomically unassigned protists: the uncharted territory of Syndiniales

Among the 343,165 metabarcodes (i.e., sequences representative of read clusters cf. Materials and Methods) we considered in our study (Supplementary Table S1), those that were taxonomically unassigned at a given taxonomic rank (i.e., without any match with reference sequences under 80% of sequence similarity) according to the PR2 or SILVA reference databases were considered as unassigned at this taxonomic rank (Supplementary Fig. S1A). Unassigned metabarcodes occurred in every sampled region and at every taxonomic rank, from kingdom to species (Fig. 1A, Supplementary Fig. S2). Both the relative abundance and number of unassigned metabarcodes increased from high to low taxonomic ranks contributing respectively to an average of 0.03% and 0.28% of the whole protist community at the kingdom rank and to 69.35% and 82.67% at the species rank (Fig. 1A, Supplementary Fig. S3B). At kingdom level, 628 metabarcodes remained unassigned among which 87.70% originated from bathypelagic samples (2150–4000 m) of the Malaspina expedition (Supplementary Fig. S4). The biggest increase in unassigned metabarcode proportion was observed from family to genus level for which 71.14% and 58.95% of metabarcodes were unassigned in relative number and relative abundance respectively (increase of 35% in unassigned metabarcodes). Overall, at the lowest taxonomic levels of our global dataset, i.e., genus and species, the proportion of unassigned metabarcodes was similar and represented more than half of the metabarcodes that could not be assigned to any referenced protist taxon (Fig. 1A). The study of

unassigned sequences was thus conducted from the viewpoint of the genus taxonomic level.

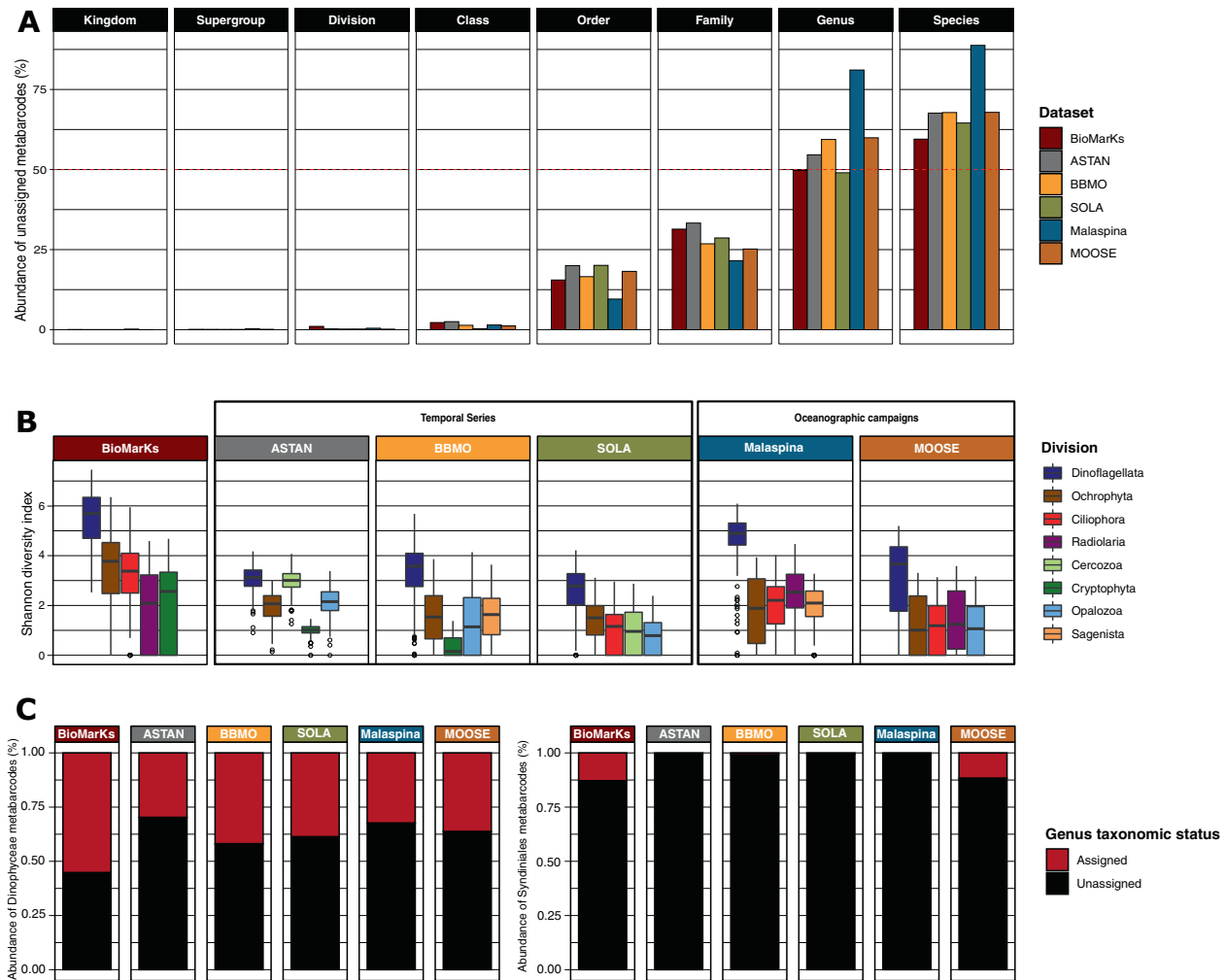
Across protist divisions, a higher diversity index was obtained for unassigned metabarcodes belonging to Dinoflagellata for all datasets (Fig. 1B). Overall, 54% of unassigned metabarcodes in relative number and 63% in relative abundance belonged to Dinoflagellata (Supplementary Fig. S4A). Among other protist divisions lacking taxonomic assignment at the genus level were Ochrophyta (all datasets), Ciliophora (BioMarKs, SOLA, Malaspina, MOOSE), Radiolaria (BioMarKs, Malaspina, MOOSE), Cercozoa (ASTAN, SOLA), Cryptophyta (BioMarKs, ASTAN, BBMO), Opalozoa (ASTAN, BBMO, SOLA, MOOSE) and Sagenista (BioMarKs, BBMO, Malaspina) (Fig. 5A). A higher diversity index was obtained for unassigned sequences, compared to assigned sequences, for the divisions Opalozoa, Sagenista and Cercozoa (Supplementary Fig. S5B). Thus, when studying only assigned genera of the latter protist divisions, their diversity could be largely underestimated.

Dinoflagellata metabarcodes represent 52% of our global dataset (179,615 metabarcodes). Among unassigned Dinoflagellata, Dinophyceae and Syndiniales were the two dominant classes and Syndiniales represented 66% and 48% of metabarcodes in terms of number and abundance respectively (Supplementary Fig. S6A). Within these two classes, the proportion of unassigned metabarcodes at the genus level was 2-fold higher for Syndiniales, with 98% and 95% of metabarcodes unassigned in terms of relative number and abundance (Fig. 1C, Supplementary Fig. S6B). Only 4 species of Syndiniales had a taxonomic assignment (0.01% of total metabarcodes and 0.53% of Syndiniales metabarcodes). Syndiniales metabarcodes unassigned at genus level represented 21% of our global dataset (72,789 metabarcodes). Given the contribution and overwhelming majority of unassigned Syndiniales in our dataset, we decided to focus the rest of our study on this lineage.

### Shared patterns of unassigned Syndiniales diversity between sunlit mediterranean and tropical waters

To investigate the spatio-temporal distribution of Syndiniales at genus level we built connected components (CCs), i.e. clusters of metabarcodes with 100% sequence identity and a minimum of 80% coverage. We consider the CCs as a proxy for clustering metabarcodes of the same Syndiniales genera or at least as pragmatic units to deal with Syndiniales molecular diversity across multiple datasets. After clustering, our global dataset contained 4317 Syndiniales CCs (30% of all CCs) out of which 4245 CCs were unassigned at the genus level (98% of Syndiniales CCs) (Supplementary Fig. S7A). These unassigned CCs belonged to 5 orders of Syndiniales: Dino-Group-I to III, Dino-Group-V [19] and an "Unknown" order (rank not assigned). Out of the unassigned Syndiniales CCs, 58% (2478 CCs) were exclusively shared within 2 sea regions, being mainly the Tropical/Subtropical Ocean and the Mediterranean Sea (which both include samples at depth >1000 m), regrouping 51% of the unassigned Syndiniales CCs (2171 CCs) (Fig. 2, light blue and orange). Unassigned CCs endemic to one region represented 23% of Syndiniales CCs (961 CCs) and were mostly found at the surface of the Tropical/Subtropical ocean (Fig. 2, light blue), while 12% of CCs (518 CCs) were shared between 3 regions (Fig. 2) and 7% CCs (288 CCs) were shared between more than 3 regions (Fig. 2). All studied sea regions shared 25 ubiquitous unassigned Syndiniales CCs, among which 14 CCs belonged to the Dino-Group-II Syndiniales order (Fig. 2).

Among Syndiniales orders, Dino-Group-II and Group-I were the most represented in our dataset (2954 CCs, 70%; 1056 CCs, 25% of unassigned Syndiniales CCs respectively (Supplementary Fig. S7B)) and their distribution was mostly restricted to the Subtropical Ocean and Mediterranean Sea (Fig. 2A, B). Dino-Group-III (212 CCs; 5% of unassigned Syndiniales CCs (Supplementary Fig. S7B)) had the widest distribution, including diversity shared between relatively more different pairs of regions, with some patterns being unique to this order, i.e. CCs exclusively



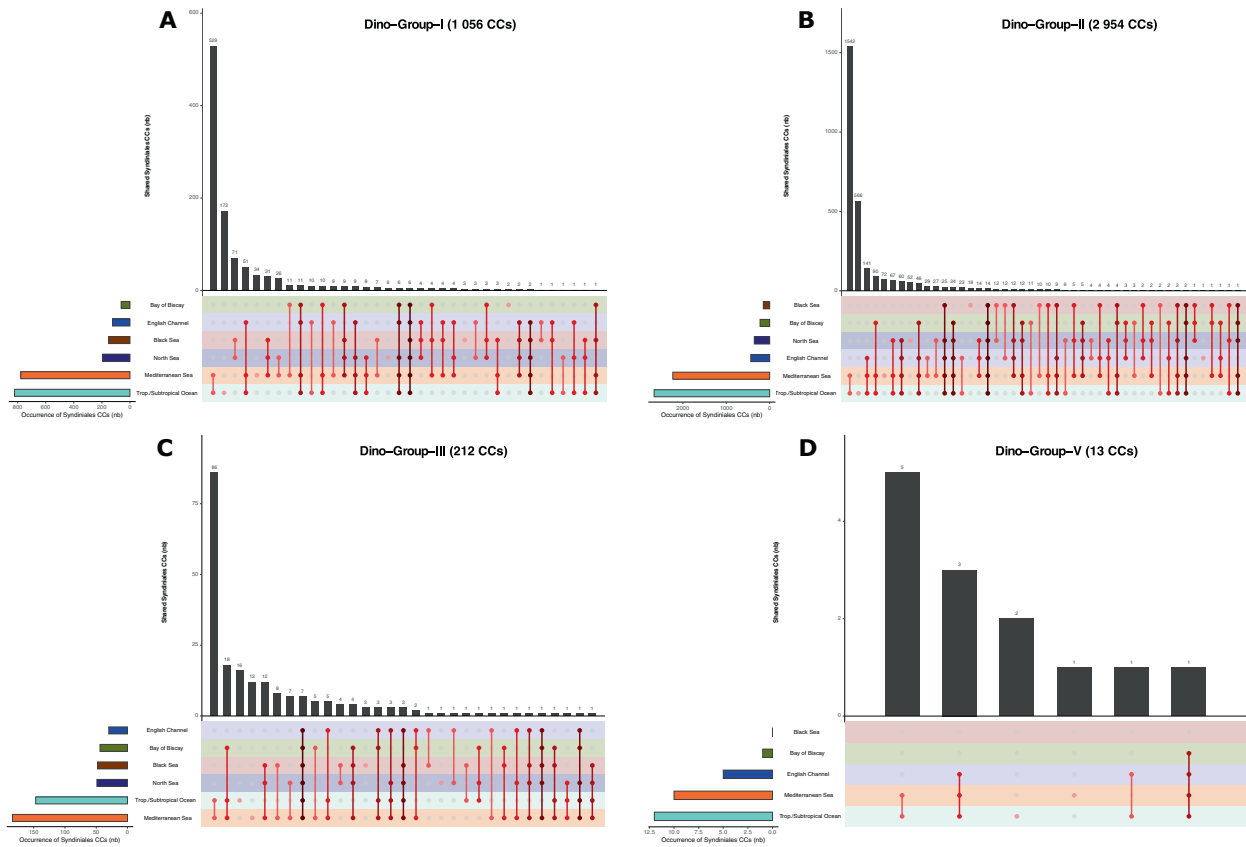
**Fig. 1** Relative abundance and diversity of unassigned metabarcodes. **A** Relative abundance of unassigned metabarcodes at each taxonomic level from kingdom to species. Colors represent the 6 studied datasets. The horizontal red dashed line marks 50% of the dataset in terms of relative abundance. **B** Shannon Weiner diversity index calculated at genus level within major protist divisions in each dataset. Only metabarcodes unassigned at genus level are selected. Colors indicate the protist divisions that represent >50% of unassigned metabarcodes at genus level in each dataset (Supplementary Fig. S5). **C** Relative abundance of assigned and unassigned metabarcodes within the class Dinophyceae (left) and Syndiniales (right) found in each dataset. Colors indicate the taxonomic status (Assigned/Unassigned) of metabarcodes at genus level.

shared between the Bay of Biscay and the Mediterranean Sea and between the Black Sea and the Mediterranean Sea (Fig. 2C). DinoGroup-V included 13 CCs (0.3% of unassigned Syndiniales CCs (Supplementary Fig. S7B)) and included CCs exclusively shared between the English Channel and the Tropical/Subtropical Ocean (Fig. 2D). The Unknown Syndiniales order included 10 CCs (0.2% of unassigned Syndiniales CCs (Supplementary Fig. S7B)) and was found in 3 sea regions: Mediterranean Sea, Tropical/Subtropical Ocean (main pattern for this order, 9 CCs) and North Sea (1 CC shared between the 3 mentioned sea regions) (Fig. 2D).

Since for all Syndiniales orders 50% of unassigned CCs were found to be exclusively common to mediterranean and tropical regions we further explored how this pattern was distributed across the water column. Among the 2171 CCs exclusively shared between mediterranean and tropical waters, Syndiniales communities were the most similar in the photic zone with 63% of CCs common between DCM (Deep Chlorophyll Maximum) layers and ~30% common between surface and DCM reciprocally (34% CCs common between Tropical/Subtropical Ocean DCM and Mediterranean Sea surface; 32% CCs common between mediterranean DCM and Tropical/Subtropical Ocean, n.b. percentages are indicative of major trends and not proportion as combinations

are not exclusive) (Fig. 3A). Notably, a pattern of shared Syndiniales CCs was also found between bathypelagic samples from the Mediterranean Sea and samples from the photic zone of the Tropical/Subtropical Ocean (29% CCs) (Fig. 3A).

In order to test if these shared diversity patterns can be explained by similar physicochemical conditions, we explored the abundance variation of unassigned Syndiniales CCs in the mediterranean and tropical waters in an RDA using the physicochemical parameters as explanatory variables. The variation of physicochemical parameters explained ~10% (11.2% in first 6 RDA dimensions) of the abundance variation of unassigned Syndiniales CCs (Fig. 3B). Based on this result, two communities of Syndiniales could be distinguished according to the first two dimensions of the RDA: a deep water (>200 m) community associated with colder and more eutrophic conditions (Fig. 3B, left) and a photic (surface and DCM) community associated with warmer and more oligotrophic conditions (Fig. 3B, right). In the RDA space associated with the photic zone, mediterranean and tropical samples partly overlap and correspond to warmer and less salty waters, hence providing an environmental basis for the observed Syndiniales pattern in these two marine environments (Fig. 3A).



**Fig. 2 Geographical distribution patterns of Connected Components (CCs) among Syndiniales orders unassigned at genus taxonomic level.** The number of CCs (y axis and above each bar) within each Syndiniales order (Dino-Group-I, II, III and V; **A**, **B**, **C** and **D** respectively) is represented according to their patterns of occurrence across the 6 sea regions defined in our metadata (Supplementary Fig. S2). Results are presented in ascending order across the x axis. The 4245 CCs containing only unassigned sequences at genus level within each order of Syndiniales were selected for this analysis. The non-assigned Syndiniales order “Unknown”, composed of 10 CCs (Supplementary Fig. S7B), is not depicted as it follows only 2 patterns: 9 CCs shared between the Mediterranean Sea and Trop./Subtropical Ocean and 1 CC shared between the Mediterranean Sea, Trop./Subtropical Ocean and North Sea. Syndiniales order Dino-Group-IV contained only assigned sequences (at genus level) and was thus not included in the plot. The occurrence color code indicates the different sea regions, while the red gradient indicates the geographical distribution range of CCs (i.e., light red: CC found in 1 sea region; bordeaux: CC found across the 6 sea regions). (N.b. The number of sampled stations is variable between sea regions with a higher number of stations sampled in Subtropical ocean (122 stations) and Mediterranean sea (35 stations). The other 4 sea regions are represented by samplings at a single station. Also Subtropical ocean and Mediterranean sea include samplings located between 200 m and 4000 m deep. North sea and Black sea samples are from surface, DCM and anoxic layers. The English Channel and Bay of Biscay include only surface samples (Supplementary Table S1).

To further investigate this hypothesis, we compared the community composition of protist divisions known to be major hosts for Syndiniales, between the marine regions of our global dataset. For Dinophyceae, Radiolaria and Ciliophora, the jaccard dissimilarity index was the lowest between Mediterranean Sea and Tropical/Subtropical Ocean compared to community comparisons between the other sea regions (Supplementary Table S2). This was also the case for Syndiniales, further supporting the results illustrated above (Fig. 3A). Neither of the remaining protist divisions found as having an important contribution to our dataset (Ochrophyta, Cercozoa, Cryptophyta, Opalozoa) showed the same tendency apart from Sagenista (Supplementary Fig. S4, Supplementary Table S2).

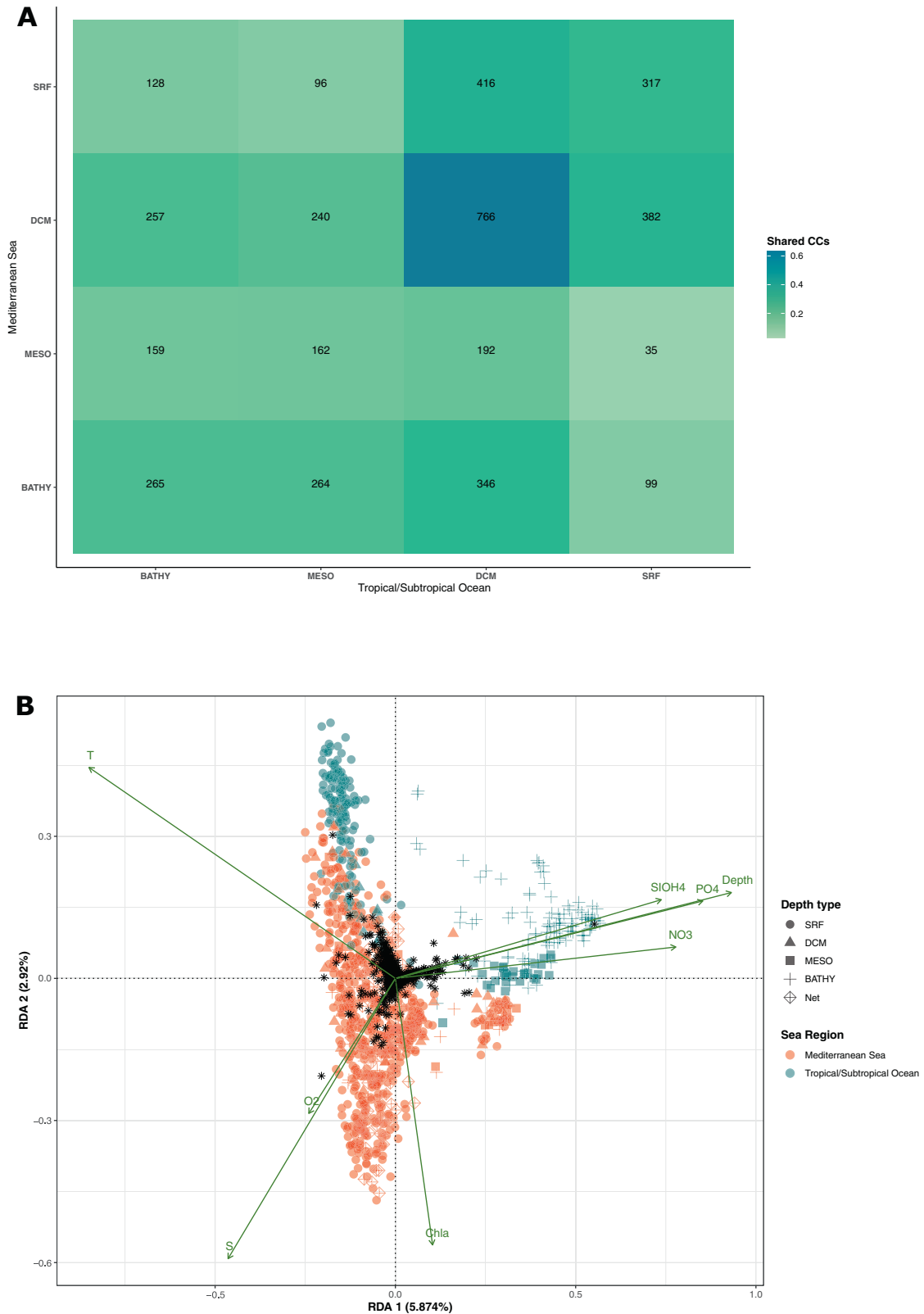
### Rhythmic ecological indicators among unassigned Syndiniales community

Temporal aspects of the Syndiniales community were studied across the three time-series (ASTAN, BBMO, SOLA) in our dataset (Fig. 4). Unassigned Syndiniales clusters did not indicate any clear seasonal preference based on monthly abundance for any of the time-series (Supplementary Fig. S8). The correlation of CCs to the overall Syndiniales community dynamics and their

rhythmicity was computed with two methods. The Escouffier’s equivalent vectors selected the CCs that are the best indicators of community abundance variation according to a PCA and the Lomb-Scargle periodogram algorithm detected if CCs follow rhythmic patterns of occurrence across time. In the studied time-series, 75% of the Syndiniales community response to environmental variation was described by 45 CCs at ASTAN, 36 CCs at BBMO and 17 CCs at SOLA (Supplementary Table S3). These community indicator CCs were all unassigned at the genus level. Rhythmic occurrence among Syndiniales CCs was found to be more prevalent in the Western Channel with 208 rhythmic Syndiniales CCs found at ASTAN, 118 CCs found at BBMO and 15 CCs found at SOLA (Supplementary Table S4). Some of the unassigned Syndiniales CCs were found to be both community indicators and rhythmic throughout the time-series: 27 CCs at ASTAN, 7 CCs at BBMO and 5 CCs at SOLA (Supplementary Table S5). The average recurrence period of these clusters was ~1.5 years at ASTAN and BBMO ~1 year at SOLA (Supplementary Table S6). We identified two rhythmic indicator CCs shared between the time-series of the English Channel (i.e., ASTAN) and Mediterranean Sea (Fig. 4): CC\_unknown\_154, shared with BBMO, and CC\_unknown\_183, shared with SOLA (recurrence periods are

indicated in Supplementary Table S6). One indicator CC, CC\_unknown\_126, was found to be shared between all the studied time series (Fig. 4) with quicker recurrence periods in the Mediterranean Sea (Supplementary Table S6). All other rhythmic indicator CCs were specific to each time-series. CC\_unknown\_126 was the CC with the highest monthly relative abundance at BBMO

and SOLA, while having the 4th highest monthly relative abundance at ASTAN. The seasonal prevalence for the majority of rhythmic indicator CCs was up to 3 seasons (Fig. 4, Supplementary Table S7). Rhythmic indicators with a 4 season prevalence occurred (45 CCs across the three time-series) and were more numerous at the Western Channel (41 CCs)



**Fig. 3 Similarity in Syndiniales genera communities between Mediterranean Sea and Subtropical Ocean. A** Proportion of Syndiniales CCs unassigned at genus level and shared between the Mediterranean Sea (y axis) and the Subtropical ocean (x axis) per depth layer (SRF for surface, DCM for Deep Chlorophyll Maximum, MESO for mesopelagic layer (>200–1000 m), BATHY for bathypelagic layer (>1000–4000 m)). The percentages illustrate major trends and not proportions (i.e., sums of percentages exceed 100% as combinations of shared CCs are not exclusive and some CCs are present in multiple depth layers). The number of samples from each depth layer for the Mediterranean Sea are: SRF = 571, DCM = 88, MESO = 97 and BATHY = 46 and for the Subtropical ocean: SRF = 136, DCM = 13, MESO = 30 and BATHY = 110. The number of CCs found in each depth layer is: SRF = 1620, DCM = 1221, MESO = 449 and BATHY = 518 for the Mediterranean Sea and SRF = 1726, DCM = 943, MESO = 611 and BATHY = 1281 for the Subtropical ocean. **B** Redundancy Analysis (RDA) for Mediterranean Sea and Subtropical Ocean data. The variation of abundance in unassigned Syndiniales CCs (black stars) is correlated to the variation of physicochemical parameters (green arrows). The most pertinent environmental parameters allowing to differentiate the studied marine regions were selected (cf. Materials and Methods: Spatiotemporal patterns of metabarcodes and CCs). The samples are represented by different colors for Mediterranean Sea (orange) and Subtropical Ocean (blue). The shapes indicate the depth layer: dot; SRF, triangle; DCM, square; MESO, cross; BATHY and square/cross; NET (vertical profile samples (0–500 m)). The dimensions of the input abundance matrix are: 4037 CCs and 1055 samples (768 samples for the Mediterranean Sea and 287 for the Tropical/Subtropical Ocean). The global RDA (cf. Materials and Methods: Spatiotemporal patterns of metabarcodes and CCs) was statistically significant at 0.005% and the first 2 axes of the RDA with the selected explanatory variables (shown below) were significant at 0.01%.

(Supplementary Table S6). The shared indicator CC\_unknown\_126 maintained a high seasonal prevalence occurring at 3 seasons in the Mediterranean Sea (i.e., BBMO and SOLA) and 4 seasons in the English Channel (i.e., ASTAN) (Fig. 4).

## DISCUSSION

### What are we missing from eukaryotic diversity with metabarcoding?

In environmental genomics investigations, the 18S rDNA marker sequence constitutes the gold standard for the exploration of eukaryotic diversity in environmental communities, shedding light on uncultivable and rare taxa [16, 30]. Yet, by integrating different metabarcoding datasets, we report that in the marine realm half of protist sequences cannot be taxonomically assigned at the genus level (57% of sequences in our dataset) and these unassigned protist taxa represent 36% to 82% of the protist community in terms of abundance across 6 diverse marine environments. Few metabarcoding studies have quantified unassigned protist diversity. In Tara Oceans, unassigned protist diversity revealed with the V9 region of the 18S rDNA marker at the supergroup level was found to be <3% of total reads [14] when referring to unassigned sequences as marker sequences with <80% identity with reference sequences. Here, with the V4 region of the 18S rDNA marker we find that unassigned protist sequences represent in abundance <1% at the supergroup level. We also report that at the genus level unassigned sequences are not rare among the protist community as they represent in abundance >45% of metabarcodes in each studied dataset and up to 80% of metabarcodes for the Malaspina expedition dataset.

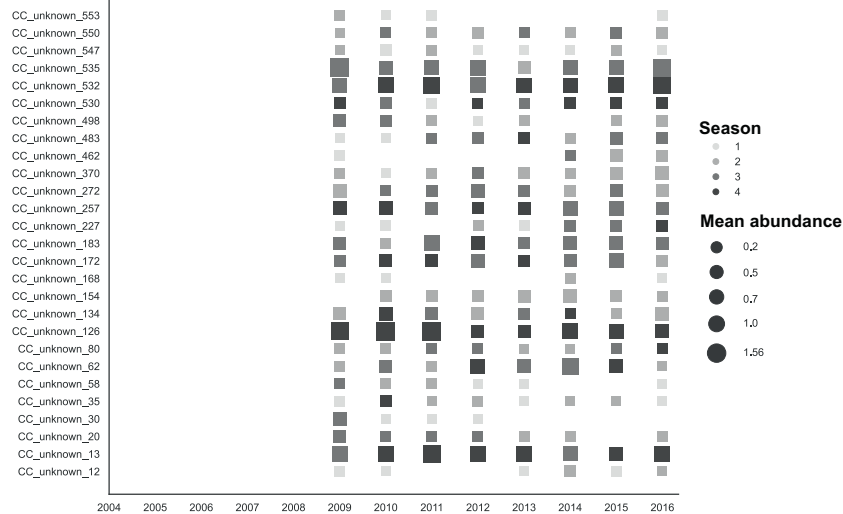
Our results confirm the current biased view of eukaryotic diversity, mostly focusing on multicellular and cultivable taxa, neglecting >70% of eukaryote diversity, including key lineages for the evolution of life and to understand ecosystems functioning [33–35]. This missing picture can be addressed, for metabarcoding studies, in the context of sample acquisition but also data acquisition in reference databases. Some oceanic regions are more sampled than others, i.e., coastal locations compared to deep/open-sea environments [36]. Moreover, the maintenance and update of reference databases is a laborious but critical process whose pace is difficult to synchronize with the generation of an ever-increasing amount of environmental sequences [37]. Metabarcoding assessments of the diversity also depends on the choice of ribosomal marker genes. In our study, the largest proportion of unassigned protist diversity was found at low taxonomic levels, a trend that has also been observed for prokaryotes [6]. Universal ribosomal markers such as 16S rDNA and 18S rDNA can have a distinct taxonomic resolution depending on the lineage considered and within each lineage [33, 38], for instance in order to describe diatom diversity a threshold >95% similarity of the V9 regions of the 18S rDNA gene with reference

sequences delimits some genera (e.g., *Undatella*) while a threshold of <90% is sufficient for assigning some other genera (e.g., *Synedropsis*) [39]. The taxonomic resolution challenge of barcoding markers is particularly relevant for rapidly evolving lineages, like predicted by evolutionary theory for parasites [40]. Studies on life history traits of multicellular parasites have demonstrated their quick adaptive plasticity, being involved in an evolutionary arms race with their host [40, 41]. Parasites are the most abundant component in many eukaryotic communities investigated through metabarcoding approaches, whether using high throughput sequencing technologies such as Illumina in tropical soils [42], subtropical marine ecosystems [43] and polar regions [20, 21], or low throughput cloning-sequencing methods in a lacustrine ecosystem [44]. In our study, parasitic Dinoflagellates (*Syndiniales*) represented 22% of metabarcodes and only 0.4% (1537 metabarcodes) could be assigned to a referenced genus, being the major contributor to the unassigned marine protist microbiome.

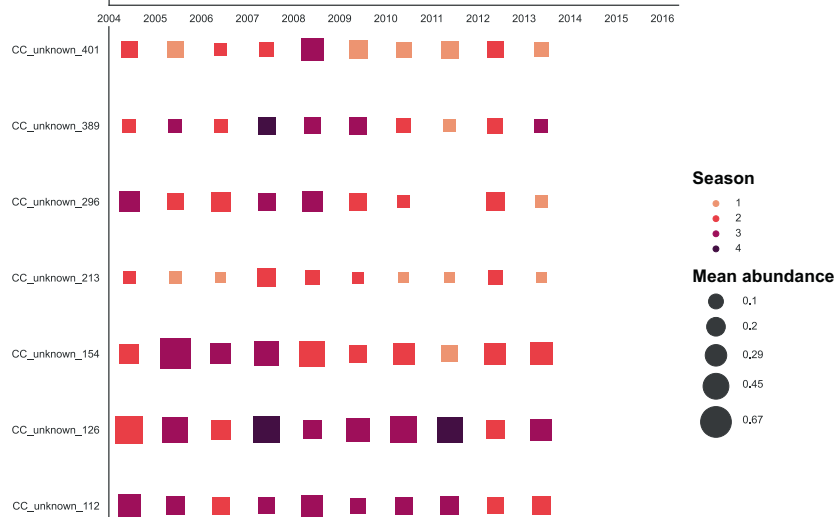
### Perspectives on *Syndiniales* biogeography

When studying the geographical distribution patterns of *Syndiniales*, we found CCs of 100%-similar sequences shared between disconnected marine regions included along a latitudinal gradient from the North Sea to the South Subtropical Atlantic, Indian, and Pacific Oceans. Similar results have been reported by Clarke et al. (2019) regarding a *Syndiniales* Group I OTU with identical V4 regions of the 18S rDNA marker retrieved in surface samples along a Southern Ocean transect near sea-ice edge and seven different Northern Hemisphere coastal locations including tropical/subtropical zones. The inferred putative 18S rDNA marker V9 region of this abundant *Syndiniales* was present in every station of the Tara Oceans voyage, including Mediterranean samples [20]. This suggests that closely related parasites can infect a wide range of hosts [20], which could also be the case for the shared *Syndiniales* CCs in our study. Our results indicated 50% (2171 CCs) of the *Syndiniales* community in common between tropical/subtropical waters and the Mediterranean basin in the euphotic zone. For these ecosystems, a convergent selection of host-parasite pairs in distant but physicochemically similar oligotrophic environments could be another hypothesis. The statistical analyses we conducted, reported physicochemical similarities in surface waters of these marine environments, while the composition of potential *Syndiniales* hosts was more similar between the Tropical/Subtropical Ocean and the Mediterranean Sea. Similar host communities could also explain the pattern of shared *Syndiniales* CCs between bathypelagic tropical/subtropical and photic Mediterranean layers. For example, *Syndiniales* host lineages of some Dinoflagellata and Radiolaria are known to be able to transit into a cyst stage during their life cycle and sediment in the water column [30, 45]. Hence, the presence of identical *Syndiniales* sequences across different water layers, that hold

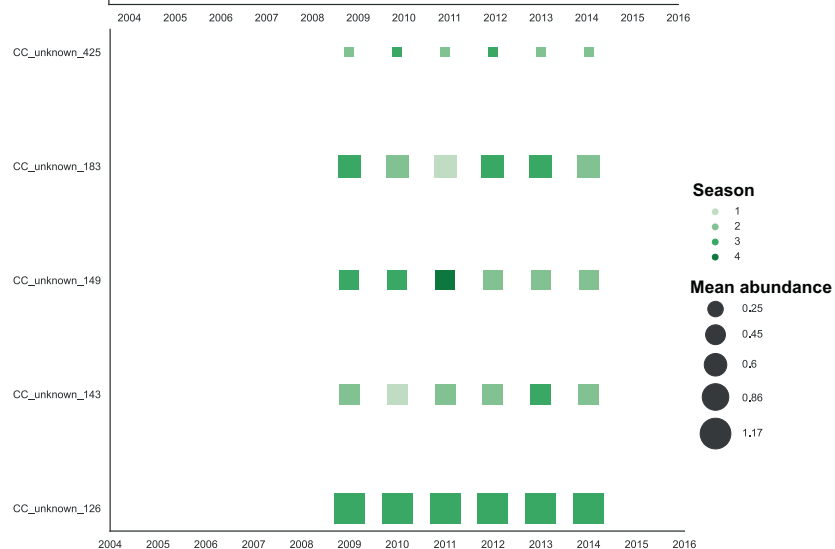
## A.ASTAN



## B.BBMO



## C.SOLA



**Fig. 4 Annual seasonal prevalence and abundance of rhythmic indicator Syndiniales CCs.** The occurrence of CCs selected by the Escouffier's equivalent vectors and Lomb-Scargle Periodogram methods was studied across each time-series: **A** ASTAN (top); **B** BBMO (middle); **C** SOLA (bottom). Relative abundance was computed per year as an average value of each month and is represented by square size. Colors indicate the seasonal prevalence of the CC throughout each year and the color gradient indicates the prevalence extent (i.e., 1 season prevalence indicated by the lightest color and 4 seasons indicated by the darkest color of the gradient). A CC is considered prevalent if it is present at least once during each season. Taxonomically unassigned CCs at genus level are indicated by "unknown" in the CC ids (y axis).



similar host communities could indicate the presence of Syndiniales within different stages of the same hosts [22, 30]. Syndiniales diversity patterns should be interpreted with caution, as the degree of genomic divergence within this lineage is high [29]. The geographical patterns we inferred with the 18S rDNA V4 region reveal genetic proximity between distant Syndiniales communities. Yet, their species composition cannot be resolved solely based on the 18S rDNA V4 region and should be defined by the combination of distinct genetic markers (e.g., V4 and V9 regions of 18S, ITS or COI) [17, 33, 46, 47]. Complementary studies need to be done comparing open sea to coastal regions and a lineage-specific primer should be designed for Syndiniales, as has been done with other parasite lineages like Perkinsea [42, 48] and Microsporidia [43, 49]. Our results open up perspectives for exploring host-parasite comparative biogeography patterns through co-occurrence networks in order to elucidate the describe host-parasite associations at a global scale and encourage the definition of host-ranges among parasites at low taxonomic resolution [50–52].

### Perspectives on Syndiniales temporal dynamics

By studying temporal patterns of Syndiniales diversity across 3 time-series we highlighted a small number of CCs that are recurrent over time, persistent through seasons and some indicators of parasite community variation (Supplementary Table S6). We can hypothesize that the recurrence of these taxa could be associated with rhythmic host patterns like annual blooms, as parasites can respond quickly to elevated host density [22, 26]. The seasonally persistent taxa we found, could indicate a generalist and opportunistic parasite behavior, infecting the hosts that are available during each season, while surviving in spore form during low host densities [23, 53]. Flexible host-parasite associations have already been described in coastal estuaries using co-occurrence networks [22]. Alternatively, parasites cannot persist below a critical host threshold [30], which questions the trophic mode of the persistent Syndiniales in our study. Up to date only parasitic and parasitoid Syndiniales have been described [50]. Nevertheless, parasitism is a mode of symbiosis along a parasite-mutualist continuum and transitions from one mode to the other should not be excluded [54]. Moreover, the novel closest group to MALVs, Eleutherids, has been recently described to be composed of free-living protists [55].

### Syndiniales as potential indicators of ecosystem change?

Our analysis also highlighted Syndiniales CCs that were both recurrent over time and good indicators of parasite community abundance variation. These Syndiniales CC hold the potential for monitoring changes in environmental microbial communities, reflecting shifts not only among the Syndiniales communities but also mirroring their associated host community. The absence of these Syndiniales CC could, for instance, indicate a shift in microbial community composition during or after an environmental perturbation. In marine environments, multicellular parasites (e.g., trematodes) have been employed as bioindicators of host physiology in response to accumulating pollution for environmental monitoring [56]. The diversity of frog and fish endoparasites was shown to reflect their surrounding ecological conditions. Selecting endoparasite taxa that are sensitive to environmental perturbation, like trematodes according to host landscape anthropogenization [56] is crucial for a potential bioindicator. In that respect, our analysis throughout a 6–10 years of abundance information and metadata suggests that dinoflagellate parasites could be used for marine habitat monitoring as it has been done with diatoms, ciliates and foraminifera [57]. Behind the blackbox of Syndiniales taxonomy could be hidden a promising global ecosystem change indicator; thanks to their worldwide distribution [14], abundance [22], quick response time

to host community shifts [22] and intimate implication in marine food webs [26].

### SSNs as integrative tools to prioritize unassigned protist taxa

In this integrative study we have used a sequence similarity network to explore the ecology of the main components of the unassigned protist microbiome by combining 6 metabarcoding datasets. SSNs are relevant and efficient analytical tools for addressing the unassigned microbiome challenge as they allow studying simultaneously large datasets, in order to categorize and prioritize unassigned sequences. They have been recently employed among prokaryotes for surveying the coding part of genomes and metagenomes [6] and taxonomy across extreme aquatic environments [12]. By exploring the biogeography of these sequences we can reveal core taxa shared across ecosystems [12, 58]. Here we have explored both biogeographical and temporal patterns of protists at the species level without requiring a reference taxonomic match. Our FAIR (Findable, Accessible, Interoperable and Reusable) computational workflow that allows to integrate data from heterogeneous ecosystem sampling protocols, such as coastal time-series and open sea campaigns and can be applied to any targeted protist group of any metabarcoding dataset, of the same marker gene, for example originating from the metaPR2 database [59] and Ocean Barcode Atlas [60]. The taxa identified by the network could then be specifically targeted for in situ hybridization [44] and isolation for single-cell omics [33]. Other approaches to reduce the unassigned taxonomic load encompass long-read sequencing [16], sequencing multiple metabarcoding markers [32] and combining metabarcoding and microscopy [36]. The unassigned microbiome holds an unexplored potential of novel taxa and functions that will surely challenge the current view of microbial ecology in the ocean and beyond [5, 33, 61].

## MATERIALS AND METHODS

### Gathering and homogenization of metabarcoding datasets

Metabarcoding datasets of pre-processed and clustered 18S rDNA marker sequences containing the variable region V4 and originating from 6 distinct sampling projects were gathered. The datasets include three temporal series of bimensual samplings at a single station: ASTAN in Roscoff, English Channel, France (8 years of data), BBMO in Blanes Bay, Mediterranean Sea, Spain (10 years of data) and SOLA in Banyuls-sur-Mer, Mediterranean Sea, France (9 years of data) (Supplementary Fig. S1D); two oceanographic campaigns of punctual samplings across 148 locations: Malaspina Expedition (122 stations, circumglobal Tropical/Subtropical Ocean) and MOOSE (26 stations, Mediterranean Sea, 10.18142/235, campaigns 2017 (10.17600/17001500) and 2018 (10.17600/18000442)) (Supplementary Fig. S1B, C) [62]; and one European project of punctual samplings at 6 marine coastal stations: BioMarKs project (samples from: Oslo, Norway; Roscoff, France; Varna, Bulgaria; Gijon, Spain; Barcelona, Spain; Naples, Italy). Sequencing was done with Illumina MiSeq technology, except for the BioMarKs project sequenced by 454 pyrosequencing. Each metabarcoding dataset contained the abundance tables of reads clean-processed and inferred into ASVs (OTUs for BioMarKs) and their taxonomic affiliation (details in Supplementary Table S1). The initial global dataset contained 539,546 metabarcodes. For homogenization purposes, the same two filtering conditions were applied independently to each of the 6 datasets (Supplementary Fig. S1A, Step 1): removal of sequences corresponding to metazoans, terrestrial plants (Streptophyta) and macroalgae (Floriophyceae, Bangiophyceae, Phaeophyceae, and Ulvophyceae); removal of sequences having less than 80% identity with reference databases. The latter threshold was chosen according to the original preprocessing of the datasets: the MOOSE dataset had beforehand implemented a minimum identity threshold of 80% and Malaspina and BBMO of 95%. A 95% filter was considered too stringent, as too many unknown sequences of interest might be removed, a 80% threshold was applied to the global dataset for homogenization. The global abundance table resulting from the homogenization workflow involved at this stage 343,165 metabarcodes, and each sample was normalized by total read number and scaled from 0 to 1.

To account for variations in the taxonomic assignment procedure (assignment tools, database versions) across datasets, a new taxonomic assignment (Supplementary Fig. S1A, Step 2) was performed on the global set of metabarcodes with the PR2 database (version 4.12.0, released on 08.08.2019, <https://pr2-database.org>; blast parameters: -evalue 0.01 -max\_target\_seqs 15, [63]). Only the best hit (best e-value) of each alignment was kept. These new assignments were filtered again for multicellular taxa and only sequences with a length greater to 200 bp were kept (Supplementary Fig. S1A, Step 3). The PR2 database includes 8 taxonomic ranks: kingdom, supergroup, division, class, order, family, genus, and species. To avoid prokaryotic contamination at the kingdom level, an assignment was performed using the SILVA database (<https://www.arb-silva.de/>, version 138) implemented in the DADA2 algorithm [64]. 3874 prokaryotic metabarcodes were removed out of the 4519 unassigned sequences at the kingdom level. The taxonomic ranks that were left unassigned were marked as “Unknown” and the taxonomy of the sequence was considered unassigned at this given rank. Unassigned ranks located between attributed ranks were regarded as gaps in the taxonomic hierarchy and not as unassigned ranks. The diversity and abundance of unassigned sequences were explored on Rstudio (R version 4.1.1, [65]), using the packages: ‘data.table’, ‘vegan’, ‘ggplot2’, ‘ggsci’ and ‘gridExtra’.

### Homogenization and analysis of environmental data

Our global dataset included 1531 samples (ASTAN: 374, BBMO: 327, SOLA: 154, Malaspina: 289, MOOSE: 272) (Supplementary Table S1). The metadata and environmental information associated with the studied samples were retrieved from the initial studies [66–72] and supplemented with public oceanographic databases (cf. additional information in the next paragraph). The information contained 14 metadata variables: name of the campaign, sampled region, station (for oceanographic campaigns), sequencing technology, sampling date, year, month, season, depth (m), depth type (surface (depth ≤ 5 m), deep maximum chlorophyll (DCM), mesopelagic zone (depth ≥ 200 m), bathypelagic zone (depth ≥ 1000 m)), sampled size fraction, latitude, longitude). The 3 temporal series datasets (ASTAN, SOLA, BBMO) were sampled only at surface, BioMarks dataset was sampled at surface and DCM, while the 2 oceanographic campaigns (MOOSE and Malaspina) were sampled at surface, DCM, mesopelagic and bathypelagic zones (up to 2000 m depth for MOOSE and 4000 m for Malaspina). The sampled size fractions are: 0–0.2 μm, 0.2–3 μm, 0.2–0.8 μm, 0.8–3 μm, 0.8–20 μm, 3–20 μm, 20–2 000 μm. The information contained as well 10 environmental variables: temperature (°C), salinity (PSU), pH, concentrations of oxygen (ml/L), nitrate (μmol/L), nitrite (μmol/L), ammonium (μmol/L), phosphate (μmol/L), silicate (μmol/L) and chlorophyll-*a* (μg/L). For ASTAN and BioMarks datasets, when in situ environmental variables were missing, metadata were retrieved from public oceanographic databases (SOMLIT database (<https://www.somlit.fr>); World Ocean Database (<https://www.ncei.noaa.gov/access/world-ocean-database-select/dbsearch.html>), SeaDataNet (<https://cdi.seadatanet.org/search>)). No additional information could be retrieved for 2 locations (Varna and Gijon). The environmental data and metadata were explored on Rstudio (R version 4.1.1), using the packages: ‘maps’, ‘tidyverse’, ‘sp’, ‘reshape2’, ‘tidyr’, ‘ade4’, ‘factoextra’ (Principal Component Analysis), ‘ggplot2’, ‘ggsci’ and ‘gridExtra’.

### Sequence Similarity Network as a framework for heterogeneous datasets comparison

The 343,165 metabarcodes were aligned against each other with the following options: e-value <1e–4; >80% coverage for both subject and query (except for the alignments involving SOLA sequences (maximum sequence length = 230 bp compared to a mean of 430 bp for other datasets) in which case the coverage threshold was applied only to the SOLA sequence in order to avoid a misrepresentation of SOLA sequences in our analysis). Self-hits and reciprocal hits (same query-subject pair) were discarded. The filtered blast output (2,942,982 alignments) was used to cluster sequences by similarity in a Sequence Similarity Network (SSN), with ‘igraph’ R package (version 1.2.6, <https://igraph.org/r/>, [73]). The sequences (i.e., the network nodes) were labeled according to metadata and taxonomic affiliation. The sequences were clustered into Connected Components (CCs) by setting an identity threshold of 100% sequence similarity, and CCs involving less than 6 sequences were removed (this number of 6 was chosen in order to enable the representativity of all 6 datasets in small CCs). The taxonomic homogeneity of CCs in the network was evaluated for known sequences at the genus level, and if only a single genus assignment was found this name was extrapolated to the other

nodes of the CC even if these ones were of unknown genera. Thus, CCs were considered here as a proxy for studying taxonomic diversity at the genus level. The final network was composed of 12,619 CCs.

### Spatio-temporal patterns of metabarcodes and CCs

CCs including only Syndiniales sequences unassigned at genus level were extracted from the network (4245 CCs; 33.6% of network and 47.6% of unassigned network CCs at genus level, Supplementary Figs. S7A and S8). The distribution of clusters across marine environments and time was explored with R functions that were coded to extract the sequence attributes related to sampling data in each CC (location, dataset, depth, season month). A Redundancy Analysis (RDA) was performed on the abundance matrix of Syndiniales CCs using the metadata for Tropical/Subtropical Ocean and Mediterranean Sea samples as explanatory variables. ANOVA tests were run to assess the robustness of the global RDA (all environmental variables included) and of the first two dimensions of the RDA with selected environmental variables. Both the RDA and ANOVA were run via the *vegan* package. Potential Syndiniales host communities were compared with the Jaccard dissimilarity index of the based on the Bray–Curtis compositional dissimilarity of abundances [74]. Jaccard index was computed with the *vegdist* function of the ‘vegan’ package, according to the formula:  $2B/(1+B)$ , where B is Bray–Curtis dissimilarity. The temporal patterns of Syndiniales among each Time Series (ASTAN, BBMO, SOLA) were explored for both assigned and unassigned genera clusters (4317 CCs; 34.2% of network, Supplementary Fig. S7B). Diversity indexes (species richness (S), Shannon’s diversity (H) and reverse Pielou index (J), using the *vegan* package) and statistical metrics (mean abundance per month) were computed. The Escoufier’s equivalent vector method was applied on CCs present at least 5 times across each time series. This method was run with the package *pastecs* and sorted clusters according to their correlation to a principal component analysis (PCA) [75]. The cumulated correlation level chosen was 75% in order to avoid retrieving clusters with negligible correlation (100% would result in retrieving the whole dataset). The rhythmicity of CCs across time was computed by the Lomb–Scargle Periodogram (LSP) [72] via the *lomb* package. Each CC was associated with a PNmax value, a *p*-value and a rhythmicity period (in days). The LSP method was applied according to Lambert et al., 2019 and is particularly well suited for our time-series data, as it allows us to detect the periodic patterns in unevenly sampled data. The PNmax is the decision variable corresponding to the peak normalized power, and CCs were considered rhythmic for a PNmax >10 (i.e., *p*-value <0.01). Graphical representations were plotted on Rstudio (R version 4.1.1) and Python (v3.8, package ‘seaborn’).

### DATA AVAILABILITY

Scripts, data and Rmarkdown files necessary to run all the analyses included in this work are publicly available on the github page [https://github.com/IrisRizos/Unassigned\\_Protists\\_SSN](https://github.com/IrisRizos/Unassigned_Protists_SSN).

### REFERENCES

- Forster D, Bittner L, Karkar S, Dunthorn M, Romac S, Audic S, et al. Testing ecological theories with sequence similarity networks: marine ciliates exhibit similar geographic dispersal patterns as multicellular organisms. *BMC Biol.* 2015;13:16. <http://bmcbiol.biomedcentral.com/articles/10.1186/s12915-015-0125-5>.
- Galperin MY, Koonin EV. From complete genome sequence to ‘complete’ understanding? *Trends Biotechnol.* 2010;28:398–406. <https://linkinghub.elsevier.com/retrieve/pii/S0167779910000892>.
- Modha S, Robertson DL, Hughes J, Orton RJ. Quantifying and cataloguing unknown sequences within human microbiomes. *mSystems.* 2022;7:e01468–21. <https://journals.asm.org/doi/10.1128/msystems.01468-21>.
- Wyman SK, Avila-Herrera A, Nayfach S, Pollard KS. A most wanted list of conserved microbial protein families with no known domains. *PLoS One.* 2018;13:e0205749. <https://dx.plos.org/10.1371/journal.pone.0205749>.
- Bernard G, Pathmanathan JS, Lannes R, Lopez P, Baptiste E. Microbial dark matter investigations: how microbial studies transform biological knowledge and empirically sketch a logic of scientific discovery. *Genome Biol Evol.* 2018;10:707–15. <https://academic.oup.com/gbe/article/10/3/707/4840377>.
- Vanni C, Schechter MS, Acinas SG, Barberán A, Buttigieg PL, Casamayor EO, et al. Unifying the known and unknown microbial biological knowledge space. *eLife.* 2022;11:e67667. <https://elifesciences.org/articles/67667>.
- Jaroszewski L, Li Z, Krishna SS, Bakolitsa C, Woolley J, Deacon AM, et al. Exploration of uncharted regions of the protein universe. *PLoS Biol.* 2009;7:e1000205. <https://dx.plos.org/10.1371/journal.pbio.1000205>.

8. Meng A, Corre E, Probert I, Gutierrez-Rodriguez A, Siano R, Annamale A, et al. Analysis of the genomic basis of functional diversity in dinoflagellates using a transcriptome-based sequence similarity network. *Mol Ecol*. 2018;27:2365–80. <https://onlinelibrary.wiley.com/doi/10.1111/mec.14579>.
9. Meng A, Marchet C, Corre E, Peterlongo P, Alberti A, Da Silva C, et al. A de novo approach to disentangle partner identity and function in holobiont systems. *Microbiome*. 2018;6:105. <https://microbiomejournal.biomedcentral.com/articles/10.1186/s40168-018-0481-9>.
10. Carradec Q, Pelletier E, Da Silva C, Alberti A, Seeleuthner Y, Tara Oceans Coordinators et al. A global ocean atlas of eukaryotic genes. *Nat Commun*. 2018;9:373. <http://www.nature.com/articles/s41467-017-02342-1>.
11. Ramond P, Sourisseau M, Simon N, Romac S, Schmitt S, Rigaut-Jalabert F, et al. Coupling between taxonomic and functional diversity in protistan coastal communities: functional diversity of marine protists. *Environ Microbiol*. 2019;21:730–49. <http://doi.wiley.com/10.1111/1462-2920.14537>.
12. Zamkovaya T, Foster JS, de Crécy-Lagard V, Conesa A. A network approach to elucidate and prioritize microbial dark matter in microbial communities. *ISME J*. 2021;15:228–44. <https://www.nature.com/articles/s41396-020-00777-x>.
13. Sunagawa S, Coelho LP, Chaffron S, Kultima JR, Labadie K, Salazar G, et al. Structure and function of the global ocean microbiome. *Science*. 2015;348:1261359–1261359. <https://www.sciencemag.org/lookup/doi/10.1126/science.1261359>.
14. de Vargas C, Audic S, Henry N, Decelle J, Mahe F, Logares R, et al. Eukaryotic plankton diversity in the sunlit ocean. *Science*. 2015;348:1261605–1261605. <https://www.sciencemag.org/lookup/doi/10.1126/science.1261605>.
15. Strassert JFH, Karnkowska A, Hehenberger E, del Campo J, Kolisko M, Okamoto N, et al. Single cell genomics of uncultured marine alveolates shows paraphyly of basal dinoflagellates. *ISME J*. 2018;12:304–8. <http://www.nature.com/articles/ismej2017167>.
16. Burki F, Sandin MM, Jamy M. Diversity and ecology of protists revealed by metabarcoding. *Curr Biol*. 2021;31:R1267–80. <https://linkinghub.elsevier.com/retrieve/pii/S0960982221010563>.
17. Cai R, Kayal E, Alves-de-Souza C, Bigeard E, Corre E, Jeanthon C, et al. Cryptic species in the parasitic Amoebophrya species complex revealed by a polyphasic approach. *Sci Rep*. 2020;10:2531. <http://www.nature.com/articles/s41598-020-59524-z>.
18. Singer D, Sepey CVW, Lentendu G, Dunthorn M, Bass D, Belbahri L, et al. Protist taxonomic and functional diversity in soil, freshwater and marine ecosystems. *Environ Int*. 2021;146:106262. <https://linkinghub.elsevier.com/retrieve/pii/S0160412020322170>.
19. Guillou L, Viprey M, Chambouvet A, Welsh RM, Kirkham AR, Massana R, et al. Widespread occurrence and genetic diversity of marine parasitoids belonging to Syndiniales (Alveolata). *Environ Microbiol*. 2008;10:3349–65. <https://onlinelibrary.wiley.com/doi/10.1111/j.1462-2920.2008.01731.x>.
20. Clarke LJ, Bestley S, Bissett A, Deagle BE. A globally distributed Syndiniales parasite dominates the Southern Ocean micro-eukaryote community near the sea-ice edge. *ISME J*. 2019;13:734–7. <http://www.nature.com/articles/s41396-018-0306-7>.
21. Cleary AC, Durbin EG. Unexpected prevalence of parasite 18S rDNA sequences in winter among Antarctic marine protists. *J Plankton Res*. 2016;38:401–17. <https://academic.oup.com/plankt/article-lookup/doi/10.1093/plankt/fbw005>.
22. Anderson SR, Harvey EL. Temporal variability and ecological interactions of parasitic marine syndiniales in coastal protist communities. *mSphere*. 2020;5. <https://journals.asm.org/doi/10.1128/mSphere.00209-20>.
23. Käse L, Metfies K, Neuhaus S, Boersma M, Wiltshire KH, Kraberg AC. Host-parasitoid associations in marine planktonic time series: can metabarcoding help reveal them? Amato A, editor. *PLoS One*. 2021;16:e0244817. <https://dx.plos.org/10.1371/journal.pone.0244817>.
24. Jephcott TG, Alves-de-Souza C, Gleason FH, van Ogtrop FF, Sime-Ngando T, Karpov SA, et al. Ecological impacts of parasitic chytrids, syndiniales and perkinsids on populations of marine photosynthetic dinoflagellates. *Fungal Ecology*. 2016; <https://linkinghub.elsevier.com/retrieve/pii/S175450481500032X>.
25. Siano R, Alves-de-Souza C, Foulon E, Bendif EM, Simon N, Guillou L, et al. Distribution and host diversity of Amoebophryidae parasites across oligotrophic waters of the Mediterranean Sea. *Biogeosciences*. 2011;8:267–78. <https://bg.copernicus.org/articles/8/267/2011/>.
26. Moran MA, Ferrer-González FX, Fu H, Nowinski B, Olofsson M, Powers MA, et al. The Ocean's labile DOC supply chain. *Limnol Oceanogr*. 2022;ino.12053. <https://onlinelibrary.wiley.com/doi/10.1002/lno.12053>.
27. Chambouvet A, Morin P, Marie D, Guillou L. Control of toxic marine dinoflagellate blooms by serial parasitic killers. *Science*. 2008;322:1254–7. <https://www.science.org/doi/10.1126/science.1164387>.
28. Shadrin AM, Kholodova MV, Pavlov DS. Geographic distribution and molecular genetic identification of the parasite of the genus Ichthyodinium causing mass mortality of fish eggs and larvae in coastal waters of Vietnam. *Dokl Biol Sci*. 2010;432:220–3. <http://link.springer.com/10.1134/S0012496610030154>.
29. Farhat S, Le P, Kayal E, Noel B, Bigeard E, Corre E, et al. Rapid protein evolution, organelle reductions, and invasive intronic elements in the marine aerobic parasite dinoflagellate Amoebophrya spp. *BMC Biol*. 2021;19:1. <https://bmcbiol.biomedcentral.com/articles/10.1186/s12915-020-00927-9>.
30. Chambouvet A, Alves-de-Souza C, Cuffe V, Marie D, Karpov S, Guillou L. Interplay between the parasite Amoebophrya sp. (Alveolata) and the cyst formation of the red tide Dinoflagellate Scrippsiella trochoidea. *Protist*. 2011;162:637–49. <https://linkinghub.elsevier.com/retrieve/pii/S1434461011000022>.
31. Okamura B, Hartigan A, Naldoni J. Extensive uncharted biodiversity: the parasite dimension. *integrative and comparative biology*. 2018. <https://academic.oup.com/icb/advance-article/doi/10.1093/icb/icy039/5026008>.
32. Rohde K. Ecology and Biogeography, Future Perspectives: Example Marine Parasites. *Geoinfor Geostat Overview*. 2016;4; [http://www.scitechnol.com/peer-review/ecology-and-biogeography-future-perspectives-example-marine-parasites-wRny.php?article\\_id=4869](http://www.scitechnol.com/peer-review/ecology-and-biogeography-future-perspectives-example-marine-parasites-wRny.php?article_id=4869).
33. Pawlowski J, Audic S, Adl S, Bass D, Belbahri L, Berney C, et al. CBOL Protist Working Group: barcoding eukaryotic richness beyond the animal, plant, and fungal kingdoms. *PLoS Biol*. 2012;10:e1001419. <https://dx.plos.org/10.1371/journal.pbio.1001419>.
34. del Campo J, Sieracki ME, Molestina R, Keeling P, Massana R, Ruiz-Trillo I. The others: our biased perspective of eukaryotic genomes. *Trends Ecol Evol*. 2014;29:252–9. <https://linkinghub.elsevier.com/retrieve/pii/S0169534714000640>.
35. Sibbald SJ, Archibald JM. More protist genomes needed. *Nat Ecol Evol*. 2017;1:0145. <http://www.nature.com/articles/s41559-017-0145>.
36. Egge E, Elferink S, Vaulot D, John U, Bratbak G, Larsen A, et al. An 18S V4 rRNA metabarcoding dataset of protist diversity in the Atlantic inflow to the Arctic Ocean, through the year and down to 1000 m depth. *Earth Syst Sci Data*. 2021;13:4913–28. <https://essd.copernicus.org/articles/13/4913/2021/>.
37. Mugnai F, Megléc E, Abbiati M, Bavestrello G, Bertasi F, Bo M, et al. Are well-studied marine biodiversity hotspots still blackspots for animal barcoding? *Global Ecol Conserv*. 2021;32:e01909. <https://linkinghub.elsevier.com/retrieve/pii/S2351989421004595>.
38. Bittner L, Gobet A, Audic S, Romac S, Egge ES, Santini S, et al. Diversity patterns of uncultured Haptophytes unravelled by pyrosequencing in Naples Bay. *Mol Ecol*. 2013;22:87–101. <https://onlinelibrary.wiley.com/doi/10.1111/mec.12108>.
39. Malviya S, Scalco E, Audic S, Vincent F, Veluchamy A, Poullain J, et al. Insights into global diatom distribution and diversity in the world's ocean. *Proc Natl Acad Sci USA*. 2016;113. <https://pnas.org/doi/full/10.1073/pnas.1509523113>.
40. Kochin BF, Bull JJ, Antia R. Parasite evolution and life history theory. *PLoS Biol*. 2010;8:e1000524. <https://dx.plos.org/10.1371/journal.pbio.1000524>.
41. Sheath DJ, Dick JTA, Dickey JWE, Guo Z, Andreou D, Britton JR. Winning the arms race: host–parasite shared evolutionary history reduces infection risks in fish final hosts. *Biol Lett*. 2018;14:20180363. <https://royalsocietypublishing.org/doi/10.1098/rsbl.2018.0363>.
42. Mahé F, de Vargas C, Bass D, Czech L, Stamatakis A, Lara E, et al. Parasites dominate hyperdiverse soil protist communities in Neotropical rainforests. *Nat Ecol Evol*. 2017;1:0091. <http://www.nature.com/articles/s41559-017-0091>.
43. Blanco-Bercial L, Parsons R, Bolaños L, Johnson R, Giovannoni S, Curry R. The protist community mirrors seasonality and mesoscale hydrographic features in the oligotrophic Sargasso Sea. 2022. <https://www.authorea.com/users/453879/articles/551657-the-protist-community-mirrors-seasonality-and-mesoscale-hydrographic-features-in-the-oligotrophic-sargasso-sea?commit=ba32b47ec0dfbf4865eb448dd0b5dd27d5f8cd15>.
44. Lepère C, Domaizon I, Debroas D. Unexpected importance of potential parasites in the composition of the freshwater small-eukaryote community. *Appl Environ Microbiol*. 2008;74:2940–9. <https://journals.asm.org/doi/10.1128/AEM.01156-07>.
45. Decelle J, Martin P, Paborstava K, Pond DW, Tarling G, Mahé F, et al. Diversity, ecology and biogeochemistry of cyst-forming acantharia (radiolaria) in the Oceans. *PLoS One*. 2013;8:e53598. <https://dx.plos.org/10.1371/journal.pone.0053598>.
46. Stern RF, Horak A, Andrew RL, Coffroth MA, Andersen RA, Küpper FC, et al. Environmental barcoding reveals massive dinoflagellate diversity in marine environments. *PLoS One*. 2010;5:e13991. <https://dx.plos.org/10.1371/journal.pone.0013991>.
47. Stoeck T, Bass D, Nebel M, Christen R, Jones MDM, Breiner HW, et al. Multiple marker parallel tag environmental DNA sequencing reveals a highly complex eukaryotic community in marine anoxic water. *Mol Eco*. 2010;19:21–31. <http://doi.wiley.com/10.1111/j.1365-294X.2009.04480.x>.
48. Chambouvet A, Gower DJ, Jirků M, Yabsley MJ, Davis AK, Leonard G, et al. Cryptic infection of a broad taxonomic and geographic diversity of tadpoles by Perkinsae protists. *Proc Natl Acad Sci USA*. 2015;112. <https://pnas.org/doi/full/10.1073/pnas.1500163112>.
49. Chauvet M, Debroas D, Moné A, Dubuffet A, Lepère C. Temporal variations of Microsporidia diversity and discovery of new host–parasite interactions in a lake ecosystem. *Environ Microbiol*. 2022;1462-2920.15950. <https://onlinelibrary.wiley.com/doi/10.1111/1462-2920.15950>.
50. Bjorbaek MFM, Evenstad A, Røseg LL, Krabberød AK, Logares R. The planktonic protist interactome: where do we stand after a century of research? *ISME J*. 2020;14:544–59. <http://www.nature.com/articles/s41396-019-0542-5>.

51. Dallas TA, Han BA, Nunn CL, Park AW, Stephens PR, Drake JM. Host traits associated with species roles in parasite sharing networks. *Oikos*. 2019;128:23–32. <https://onlinelibrary.wiley.com/doi/10.1111/oik.05602>.
52. Lima-Mendez G, Faust K, Henry N, Decelle J, Colin S, Carcillo F, et al. Determinants of community structure in the global plankton interactome. *Science*. 2015;348:1262073. <https://www.science.org/doi/10.1126/science.1262073>.
53. Hayashi A, Crombie A, Lacey E, Richardson A, Vuong D, Piggott A, et al. Aspergillus Sydowii marine fungal bloom in Australian coastal waters, its metabolites and potential impact on symbiodinium dinoflagellates. *Marine Drugs*. 2016;14:59. <http://www.mdpi.com/1660-3397/14/3/59>.
54. Drew GC, Stevens EJ, King KC. Microbial evolution and transitions along the parasite–mutualist continuum. *Nat Rev Microbiol*. 2021;19:623–38. <https://www.nature.com/articles/s41579-021-00550-7>.
55. Hehenberger E, Tikhonenkov D, Cooney E, Jacko-Reynolds V, Irwin N, Keeling P. Free-living relatives of highly abundant unicellular marine parasites elucidate plastid loss. 2022. <https://www.researchsquare.com/article/rs-1472581/v1>.
56. Sures B, Nachev M, Selbach C, Marcogliese DJ. Parasite responses to pollution: what we know and where we go in ‘Environmental Parasitology’. *Parasites Vectors*. 2017;10:65. <http://parasitesandvectors.biomedcentral.com/articles/10.1186/s13071-017-2001-3>.
57. Payne RJ. Seven reasons why protists make useful bioindicators. *Acta Protozoologica*. 2013;52:105–13. <https://doi.org/10.4467/16890027AP.13.0011.1108>
58. Vaulot D, Sim CVH, Ong D, Teo B, Biver C, Jamy M, et al. metaPR 2: a database of eukaryotic 18S rRNA metabarcodes with an emphasis on protists. *Mol Ecol Resour*. 2022;17:55–0998.13674. <https://onlinelibrary.wiley.com/doi/10.1111/1755-0998.13674>.
59. Vernet C, Henry N, Lecubin J, Vargas C, Hingamp P, Lescot M. The Ocean barcode atlas: a web service to explore the biodiversity and biogeography of marine organisms. *Mol Ecol Resour*. 2021;21:1347–58. <https://onlinelibrary.wiley.com/doi/10.1111/1755-0998.13322>.
60. Chust G, Vogt M, Benedetti F, Nakov T, Villéger S, Aubert A, et al. Mare incognitum: a glimpse into future plankton diversity and ecology research. *Front Mar Sci*. 2017;4. <http://journal.frontiersin.org/article/10.3389/fmars.2017.00068/full>.
61. Guillou L, Bachar D, Audic S, Bass D, Berney C, Bittner L, et al. The Protist Ribosomal Reference database (PR2): a catalog of unicellular eukaryote Small Sub-Unit rRNA sequences with curated taxonomy. *Nucl Acids Res*. 2012;41:D597–604. <http://academic.oup.com/nar/article/41/D1/D597/1064851/The-Protist-Ribosomal-Reference-database-PR2-a>.
62. Coppola L, Raimbault P, Mortier L, Testor P. Monitoring the Environment in the Northwestern Mediterranean Sea. *Eos* 2019;100. <https://doi.org/10.1029/2019EO125951>.
63. Callahan BJ, McMurdie PJ, Holmes SP. Exact sequence variants should replace operational taxonomic units in marker-gene data analysis. *ISME J*. 2017;11:2639–43; <http://www.nature.com/articles/ismej2017119>.
64. RStudio Team. RStudio: integrated development for R. Boston, MA: RStudio, PBC; 2020. <http://www.rstudio.com/>.
65. Caracciolo M, Rigaut-Jalabert F, Romac S, Mahé F, Forsans S, Gac J, et al. Seasonal dynamics of marine protist communities in tidally mixed coastal waters. *Mol Ecol*. 2022;31:3761–83. <https://onlinelibrary.wiley.com/doi/10.1111/mec.16539>.
66. Giner CR, Balagué V, Krabberød AK, Ferrera I, Reñé A, Garcés E, et al. Quantifying long-term recurrence in planktonic microbial eukaryotes. *Mol Ecol*. 2019;28:923–35. <https://onlinelibrary.wiley.com/doi/10.1111/mec.14929>.
67. Giner CR, Pernice MC, Balagué V, Duarte CM, Gasol JM, Logares R, et al. Marked changes in diversity and relative activity of picoeukaryotes with depth in the world ocean. *ISME J*. 2020;14:437–49. <http://www.nature.com/articles/s41396-019-0506-9>.
68. Lambert S, Tragin M, Lozano JC, Ghiglione JF, Vaulot D, Bouget FY, et al. Rhythmicity of coastal marine picoeukaryotes, bacteria and archaea despite irregular environmental perturbations. *ISME J*. 2019;13:388–401. <http://www.nature.com/articles/s41396-018-0281-z>.
69. Logares R, Deutschmann IM, Junger PC, Giner CR, Krabberød AK, Schmidt TSB, et al. Disentangling the mechanisms shaping the surface microbial microbiota. *Microbiome*. 2020;8:55. <https://microbiomejournal.biomedcentral.com/articles/10.1186/s40168-020-00827-8>.
70. Pernice MC, Giner CR, Logares R, Perera-Bel J, Acinas SG, Duarte CM, et al. Large variability of bathypelagic microbial eukaryotic communities across the world’s oceans. *ISME J*. 2016;10:945–58. <http://www.nature.com/articles/ismej2015170>.
71. Massana R, Gobet A, Audic S, Bass D, Bittner L, Boutte C, et al. Marine protist diversity in European coastal waters and sediments as revealed by high-throughput sequencing: protist diversity in European coastal areas. *Environ Microbiol*. 2015;17:4035–49. <http://doi.wiley.com/10.1111/1462-2920.12955>.
72. Csardi G, Nepusz T. The igraph software package for complex network research. *InterJ Complex Syst*. 2006;1695:1–9.
73. Bray JR, Curtis JT. An ordination of the upland forest communities of Southern Wisconsin. *Ecol Monographs*. 1957;27:325–49. <https://onlinelibrary.wiley.com/doi/10.2307/1942268>.
74. Robert P, Escoufier Y. A unifying tool for linear multivariate statistical methods: the RV-coefficient. *J R Stat Soc Ser C Appl Stat*. 1976;25:257–65. <https://doi.org/10.2307/2347233>
75. Ruf T. The lomb-scargle periodogram in biological rhythm research: analysis of incomplete and unequally spaced time-series. *Biol Rhythm Res*. 1999;30:178–201. <https://www.tandfonline.com/doi/full/10.1076/brhm.30.2.178.1422>.

## ACKNOWLEDGEMENTS

This work was supported by the *Institut des Sciences du Calcul et des Données* (ISCD) of Sorbonne University via funding from the project FORMAL (From Observing to Modeling ocean Life, <https://iscd.sorbonne-universite.fr/research/sponsored-junior-teams/formal-2/>). The authors thank the researchers that provided the datasets and metadata for this work: N Simon and M Caracciolo (Ecology of Marine Plankton (ECOMAP), Station Biologique de Roscoff, France) for the ASTAN time series, M Mendez-Sandin (Systematic Biology, Dept. of Organismal Biology Uppsala University, Sweden) for the MOOSE campaign, R Logares (Institut de Ciències del Mar (ICM), Barcelona, Spain) for the Malaspina campaign, Blanes Bay Observatory (BBMO) time series and the BioMarks project, P Galand (Laboratory of Microbial Oceanography, Banyuls-sur-Mer, France) for the SOLA time series. The authors acknowledge the MOOSE program (Mediterranean Ocean Observing System for the Environment) coordinated by CNRS-INSU and the Research Infrastructure ILICO (CNRS-IFREMER). All analyses of this work were performed remotely on the AbiMs cluster of the marine station of Roscoff (<http://abims.sb-roscoff.fr>). LB acknowledges the Institut Universitaire de France for her 5-year nomination as Junior Member (2020–2025).

## AUTHOR CONTRIBUTIONS

LB, FN and SDA conceived the research; IR and TF analyzed and plotted graphically the data and conducted statistical analysis (c.f. first/second and third result sections respectively); LB, PD, SDA and DK contributed to the analysis by providing feedback; PD and DK contributed to the bioinformatic workflow by providing feedback and scripts; IR wrote the paper with improvement suggestions from all co-authors. All authors gave final approval for publication.

## COMPETING INTERESTS

The authors declare no competing interests.

## ADDITIONAL INFORMATION

**Supplementary information** The online version contains supplementary material available at <https://doi.org/10.1038/s43705-022-00203-7>.

**Correspondence** and requests for materials should be addressed to Iris Rizos.

**Reprints and permission information** is available at <http://www.nature.com/reprints>

**Publisher’s note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article’s Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article’s Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2023