



HAL
open science

The defensome of complex bacterial communities

Angelina Beavogui, Auriane Lacroix, Nicolas Wiart, Julie Poulain, Tom O Delmont, Lucas Paoli, Patrick Wincker, Pedro H. Oliveira

► **To cite this version:**

Angelina Beavogui, Auriane Lacroix, Nicolas Wiart, Julie Poulain, Tom O Delmont, et al.. The defensome of complex bacterial communities. 2024. hal-04406189

HAL Id: hal-04406189

<https://hal.science/hal-04406189>

Preprint submitted on 19 Jan 2024

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 Attribution - NonCommercial - NoDerivatives 4.0 International License

The defensome of complex bacterial communities

Angelina Beavogui¹, Auriane Lacroix¹, Nicolas Wiart², Julie Poulain^{1,3}, Tom O. Delmont^{1,3},
Lucas Paoli^{4,5}, Patrick Wincker^{1,3}, Pedro H. Oliveira^{1,#}

¹Génomique Métabolique, Genoscope, Institut François Jacob, Commissariat à l'Energie Atomique (CEA), CNRS, Université Evry, Université Paris-Saclay, 2 Rue Gaston Crémieux, 91057 Evry, France.

²Genoscope, Institut François Jacob, CEA, Université Paris-Saclay, 2 Rue Gaston Crémieux, 91057 Evry, France.

³Research Federation for the Study of Global Ocean Systems Ecology and Evolution, FR2022 / Tara GOsee, Paris, France.

⁴Department of Biology, Institute of Microbiology and Swiss Institute of Bioinformatics, ETH Zürich, Zürich 8093, Switzerland

⁵Institut Pasteur, Université Paris Cité, INSERM U1284, Molecular Diversity of Microbes lab, Paris, France

To whom correspondence should be addressed. Pedro H. Oliveira (pcoutool@genoscope.cns.fr)

Keywords: defense systems; phage-bacteria arms race; metagenome assembled genomes; defense islands; environmental defensomes

1 **Abstract**

2 Bacteria have developed various defense mechanisms to avoid infection and killing in
3 response to the fast evolution and turnover of viruses and other genetic parasites. Such pan-
4 immune system (or *defensome*) encompasses a growing number of defense lines that include
5 well-studied innate and adaptive systems such as restriction-modification, CRISPR-Cas and
6 abortive infection, but also newly found ones whose mechanisms are still poorly understood.
7 While the abundance and distribution of defense systems is well-known in complete and
8 culturable genomes, there is a void in our understanding of their diversity and richness in
9 complex microbial communities. Here we performed a large-scale in-depth analysis of the
10 defensomes of 7,759 high-quality bacterial population genomes reconstructed from soil,
11 marine, and human gut environments. We observed a wide variation in the frequency and
12 nature of the defensome among large phyla, which correlated with lifestyle, genome size,
13 habitat, and geographic background. The defensome's genetic mobility, its clustering in
14 defense islands, and genetic variability was found to be system-specific and shaped by the
15 bacterial environment. Hence, our results provide a detailed picture of the multiple immune
16 barriers present in environmentally distinct bacterial communities and set the stage for
17 subsequent identification of novel and ingenious strategies of diversification among
18 uncultivated microbes.

1 Introduction

2 Bacteria are under constant threat of infection by a variety of genetic parasites such as
3 bacteriophages (henceforth called phages)¹. As a result of this strong selective pressure, they
4 have evolved multiple sophisticated defense mechanisms capable of regulating the flux of
5 genetic information spread by mobile genetic elements (MGEs) via horizontal gene transfer
6 (HGT)²⁻⁴. The complete set of a bacterial defense systems' repertoire can be designated as
7 their *defensome*. Several bacterial defense systems have been discovered and extensively
8 discussed in the literature, revealing two major groupings based on their components and
9 modes of action: innate (non-specific) and adaptive immune systems^{5,6}. Typical examples of
10 innate immunity include prevention of phage adsorption⁷, restriction-modification (R-M)
11 systems that use methylation to recognize self from non-self-DNA⁸, and abortive infection
12 (Abi), in which the infected cell commits suicide before the invading phage can complete its
13 replication cycle⁹. Recent efforts to *de-novo* identify microbial defense systems resulted in the
14 discovery of several additional innate immune mechanisms with a wide range of genetic
15 architectures^{3,4}, highlighting the strong selective pressure imposed by genetic parasites on
16 microbial communities. Adaptive immune systems, on the other hand, are so far exclusively
17 represented by clustered, regularly interspaced short palindromic repeats (CRISPR)-Cas, a
18 family of defense systems that provides acquired immunity through the acquisition of short
19 DNA sequences from MGEs that are incorporated into the host genome as spacers¹⁰. Large-
20 scale efforts for defense system mapping have been recently propelled by the development
21 of bioinformatic tools such as DefenseFinder¹¹ and PADLOC¹² that rely on a profuse collection
22 of HMM profiles and specific decision rules for each known defense system. Such mapping
23 has been mainly conducted in bacterial species from reference genome databases (e.g., NCBI
24 RefSeq) that are known to overrepresent acute / common human pathogens and organisms
25 that can largely be cultivated in laboratory¹¹⁻¹³. While extremely insightful, such studies provide
26 a limited snapshot of the bacterial defensome, as they miss the uncharted fraction of
27 environmental microbial diversity that remains uncultured.

1 The current global Earth microbiome has been estimated at approximately 5×10^{30} prokaryotic
2 cells¹⁴ scattered throughout a wide range of environments, including deep oceanic and
3 continental subsurfaces, upper oceanic sediment, soil, and oceans as the most densely
4 populated cases. In many environments, 99% of microbes are yet uncultured¹⁵, while cultured
5 representatives belong overwhelmingly to the phyla Actinobacteria, Bacteroidetes, Firmicutes,
6 and Proteobacteria. For nearly 4 billion years, bacteriophages have co-evolved with bacteria,
7 with recent estimates pointing to the presence of $\sim 10^{31}$ viral particles in the biosphere¹⁶, and
8 up to 10^{23} infection events per second taking place just in the global ocean¹⁷.

9 During the last decade, extensive progress in high-throughput sequencing technologies and
10 computational methods enabled culture-independent genome-resolved metagenomics to
11 recover draft or complete metagenome-assembled genomes (MAGs)^{18–20}. The latter have
12 advanced our understanding on the diversity, abundance, and functional potential of
13 microbiota and phageome composition and corresponding ratios across different
14 environments. A healthy adult human gut for example, is a reservoir for $\sim 4 \times 10^{13}$ bacterial
15 cells (mostly Firmicutes and Bacteroidetes)²¹, and low (10^{-3} -1) virus-to-prokaryote ratios
16 (VPRs)²². In contrast, marine ecosystems typically show larger VPRs (between 8×10^{-3} - 2.15
17 $\times 10^3$, mean of 21.9), followed by soil environments which show the largest ratios (between 2
18 $\times 10^{-3}$ - 8.2×10^3 , mean of 704) (reviewed in ²³). We hypothesize that the strong VPR dynamics
19 across temporal and spatial scales is likely to profoundly shape the defensome arsenal across
20 biomes.

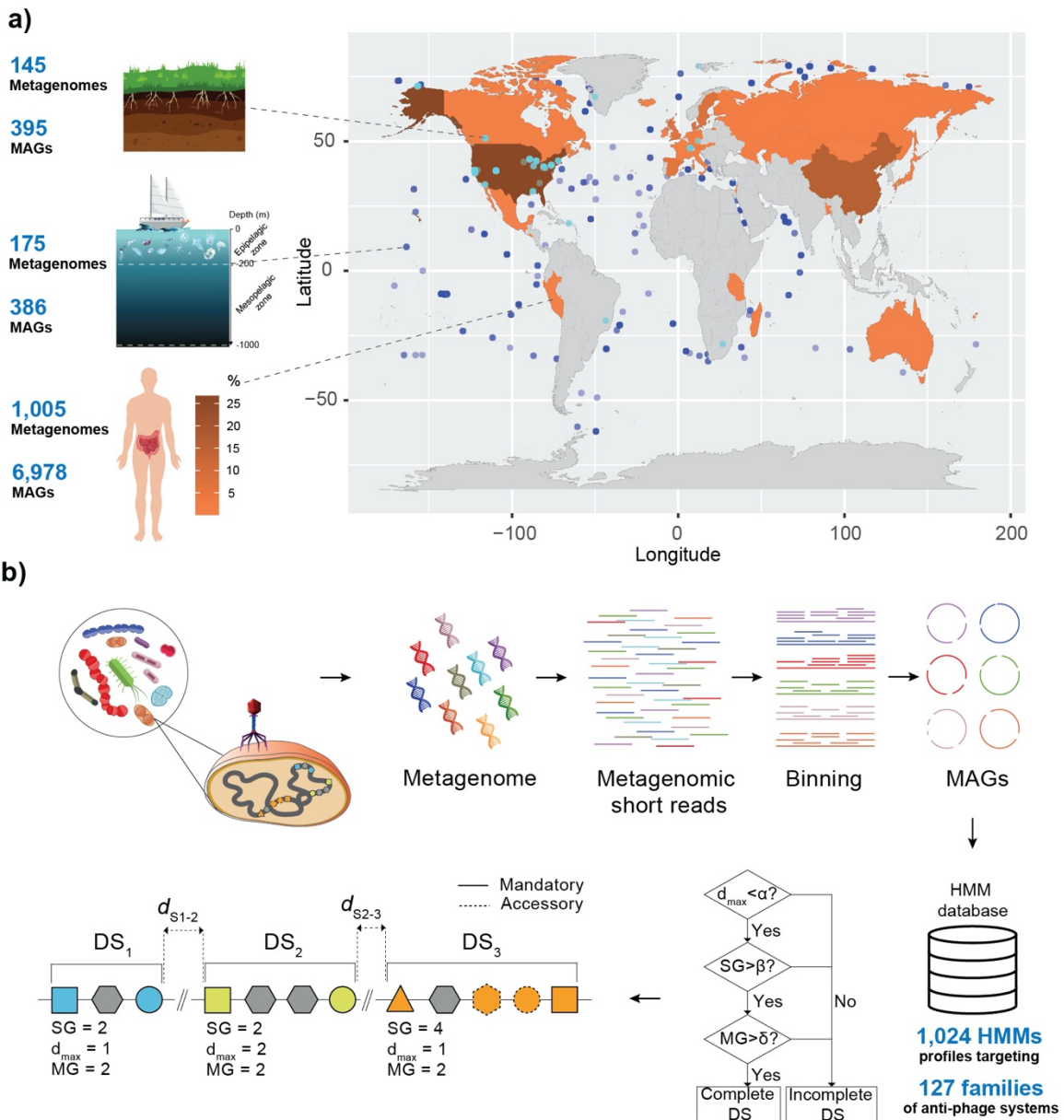
21 In this study, we conducted a large-scale in-depth investigation on the abundance, distribution,
22 and diversity of the defensome in complex bacterial communities from three key
23 environments: soil, marine, and the human gut. We tested the association between
24 defensome and different mechanisms of genetic mobility, the former's colocalization in
25 defense islands, and assessed the mutational landscape of high-frequency single nucleotide
26 polymorphisms (SNPs) and insertion-deletions (indels) across defense gene families. These
27 results provide a unique view of the interplay between microbial communities and their phage
28 invaders, and will pave the way to the identification of hitherto unknown defense systems and

- 1 / or other phage-resistance mechanisms across the enormous diversity of yet-uncultivated
- 2 microbial populations.

1 Results

2 Abundance and distribution of defensomes in bacterial MAGs

3 We performed a defensome mapping across a large dataset of 7,759 high-quality ($\geq 90\%$
 4 completeness, $\leq 5\%$ contamination/redundancy, see Methods) soil, marine, and human gut
 5 MAGs^{24–26} (**Fig. 1a, Supplementary Tables 1-4, Supplementary Fig. 1**). For this purpose,
 6 we used the DefenseFinder pipeline¹¹, which relies on a comprehensive collection of hidden
 7 Markov models (HMM) protein families and genetic organization rules targeting all major
 8 defense system families described in the literature (**Fig. 1b**).



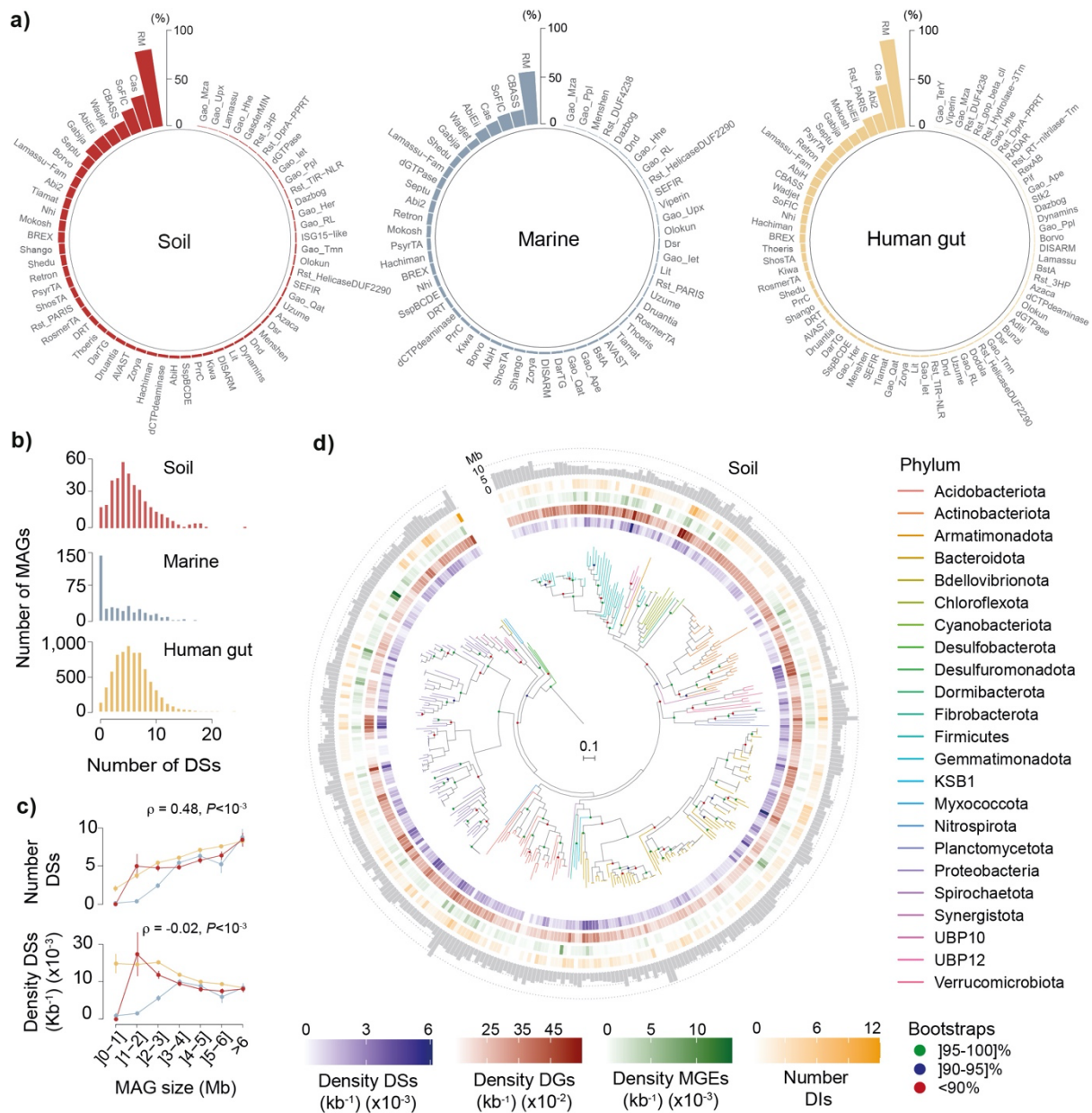
1 **Fig. 1.** Defensome analysis. (a) Our analyses focused on 7,759 high-quality near-complete MAGs recovered from
2 three distinct ecosystems: soil, marine, and human gut^{24–26}. The geographical distribution of soil and marine sample
3 collection sites is shown in the world map, as well as the percentage of human gut samples recovered from each
4 country (shown as colored heatmap). Our dataset includes at least 385 Genera (corresponding to a total of 7,593
5 MAGs) and 25 Classes (corresponding to a total of 93 MAGs) not previously covered in a recent study focusing on
6 the defensome of the NCBI RefSeq prokaryotic database¹¹. (b) A collection of 1,024 HMM profiles targeting 127
7 families of anti-MGE defense systems from DefenseFinder, was used to query the entire MAG dataset. Briefly, this
8 was performed by means of genetic organization rules allowing for two types of genetic components: "mandatory"
9 and "accessory" (as described previously¹¹). Given the wide diversity of genetic organization of anti-MGE systems,
10 rules were written differently for different types of systems. Shown is an example of a genomic region containing
11 three defense systems (DS1-DS3), respectively characterized by a total sum of genes (SG), a maximum distance
12 between defense genes (d_{\max}), and a given number of mandatory genes (MG), which will allow disentangling
13 between complete or incomplete defense systems based on established thresholds (α , β , δ).

14

15 Throughout this manuscript we will refer to *complete* anti-MGE defense systems as those
16 whose currently described genetic organization has been experimentally shown to confer anti-
17 MGE activity. Such concept of defense system *completeness* is expected to evolve in the
18 future (in particular for the recently described cryptic large multigenic systems), as more
19 details will emerge regarding their functional intra-operability. In the case of defense genes,
20 they can either belong to complete defense systems, or classify as *solitary*, i.e., those often
21 shuttled by HGT or arising from genetic erosion of complete defense systems. Of note, the
22 solitary nature of defense genes does not necessarily preclude its functional activity or even
23 implication in anti-MGE defense roles, as it has been previously shown for solitary bacterial
24 methyltransferases (MTases)²⁷.

25 In this study we found 43,263 defense systems and 764,507 defense genes pertaining to a
26 total of 70 defense families across our full MAG dataset (**Supplementary Tables 3, 4**). The
27 relative distribution of defense systems differed considerably across environments, with R–M,
28 CRISPR-Cas and the SoFIC AMPylase being the most predominant (**Fig. 2a**). When the
29 distribution of total defense genes was represented instead, we observed multiple solitary
30 genes / incomplete systems (e.g., Gabija, Gao_Qat / Gao_Mza, or Dodola) consistently
31 present across most MAGs (**Supplementary Fig. 2**). The latter suggests either non-defensive
32 roles or genetic erosion of complete systems similarly to previous observations in complete
33 genomes^{13,27}. While defense system distribution across soil and human gut MAGs followed a

1 typical binomial distribution (with most genomes encoding between 3-4 defense systems), that
 2 observed in genomes from marine environments was geometric-like, with most MAGs (~65%)
 3 showing a limited defensome (**Fig. 2b**). Such observations are in agreement with recent
 4 observations describing a 10^3 times lower effective rate of HGT in marine bacteria compared
 5 with gut bacteria, with soil bacteria occupying an intermediate position between the former
 6 two²⁸.



7
 8 **Fig. 2.** Abundance and distribution of defense systems in MAGs. (a) Percentage of soil, marine and human gut
 9 MAGs harboring each family of defense system. (b) Distribution of number of defense systems (DSs, per MAG)
 10 across environments. (c) Variation of number and density (per MAG and per kb) of defense systems (DSs) with
 11 MAG size (Mb) for each biome. Error bars represent standard deviations of the mean. (d) Phylogenetic

1 representation of 373 soil MAGs, their corresponding phyla, density (per kb) of defense systems (DSs, purple),
2 defense genes (DGs, red), MGEs (green), and number of defense islands (DIs, yellow). Distribution of MAG sizes
3 (Mb) are shown as outer layer barplots. All data corresponds to analyses performed in assemblies with values of
4 $N_{50} \geq 100$ kb.

5

6 Similarly to what has been described for R–M systems¹³, we observed positive correlations
7 between the total number of defense systems and MAG size and concomitant negative
8 correlations between the density of defense systems and size (**Fig. 2c**). Such trends can be
9 explained by the fact that bacteria with larger genomes typically engage in more HGT^{13,2}, thus
10 requiring a more abundant and diverse defensive arsenal. No qualitative differences were
11 observed when the analyses shown in **Figs. 2a-c** were performed using MAG assemblies
12 having values of $N_{50} \geq 200$ and 300 kb to control for the effect of contiguity (**Supplementary**
13 **Fig. 3**).

14 The density of defense systems (per MAG and per kb) differed widely among clades, from
15 none (largely in intracellular bacteria and obligatory endosymbionts) to more than 8×10^{-3} in
16 *Phascolarctobacterium* (human gut) and $\sim 1.5 \times 10^{-2}$ in *Elsteraceae* (soil) and UBA9040
17 (marine) environments (**Fig. 2d, Supplementary Fig. 4a, Supplementary Tables 3, 5**). No
18 MAG was entirely devoid of defense genes, with maximum densities (per MAG and per kb)
19 $\sim 8.5 \times 10^{-2}$ across the different biomes (**Fig. 2d, Supplementary Fig. 4a, Supplementary**
20 **Table 3**). When defense systems were split according to its mechanism of action, R–M, Abi,
21 and potential Abi systems were the most prevalent across biomes (**Supplementary Fig. 4b,**
22 **Supplementary Table 6**), similarly to recent observations²⁹.

23 Apart from MAG size, the abundance of defense genes was expected to depend on
24 phylogenetic depth, as deeper lineages accumulate more events of HGT exchanges
25 presumably leading to defensible buildup. We ran stepwise linear regression analyses to
26 assess the role of these variables in explaining the variance of the defensible
27 (**Supplementary Table 7**). These showed that MAG size had the strongest direct effect on
28 defensible abundance, and that phylogenetic depth had a significant but less important
29 explanatory role.

1 We found in our dataset multiple occurrences of ligand binding WYL domains and protein
2 interaction CARD-like domains (**Supplementary Figs. 4c-e**), with previously demonstrated
3 regulatory activity of phage defense systems, namely BREX, CRISPR-Cas, CBASS and
4 gasdermins³⁰⁻³³. Interestingly, we found here a significant colocalization between these
5 domains and multiple defense genes belonging to additional families involved in regulated cell
6 death, such as Lamassu, RosmerTA, and Rst_PARIS. Very few WYL and CARD-like domains
7 were found in genes from marine MAGs (< 0.75% of the dataset), in agreement with the latter's
8 more limited defensome. The patterns of colocalization differed across genomes recovered
9 from the soil and human gut environments (**Supplementary Figs. 4d-e**). For example, WYL
10 preferentially colocalized with CBASS and RosmerTA, respectively in soil and human gut
11 environments. We also found in the Bacteroidetes bacterium UBA1952, instances of an
12 operon with some similarity to the recently described *Pedobacter rhizosphaerae* CARD-
13 encoding defense system³³. In particular, UBA1952 codes for a VapC-like nuclease of the PIN
14 domain superfamily presumably operating as effector, and the SMC-like RecN with ATPase
15 function (**Supplementary Fig. 4c**).

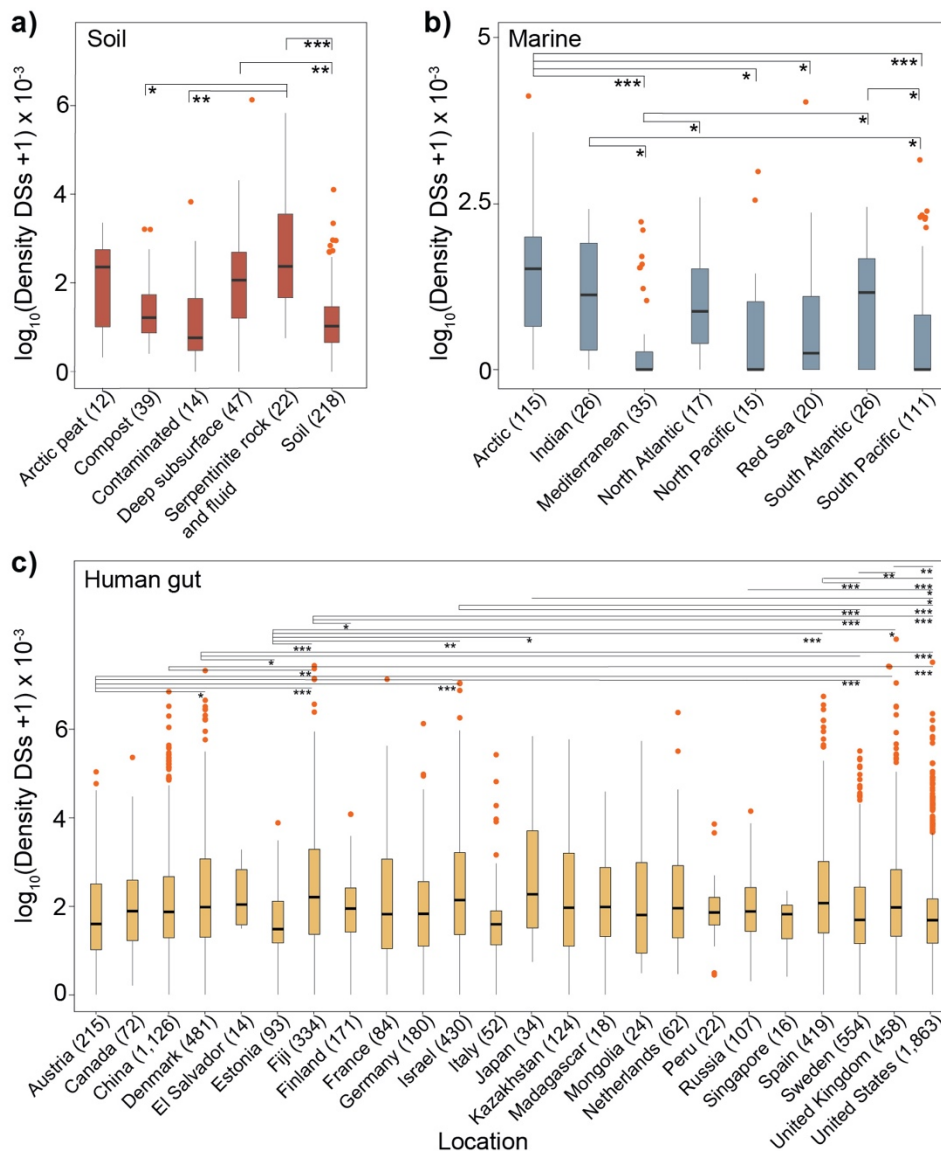
16 Hence, bacterial MAGs possess a diverse repertoire of defense systems (being defense
17 genes essentially ubiquitous), and the patterns of their distribution are very diverse and
18 dependent on genome size and taxonomy. Moreover, defense genes pertaining to systems
19 typically implicated in regulated cell death mechanisms preferentially colocalize with WYL and
20 CARD-like domains and change according to environment.

21

22 **The interplay between defensome repertoire and bacterial biogeography**

23 Fluctuations in microbial community composition are a function of a large ensemble of diverse
24 biotic and abiotic drivers. Factors such as pH (and other physicochemical parameters),
25 temperature, nutrient availability, or pollution can fundamentally reshape the spatiotemporal
26 dynamics of soil / marine bacterial and viral communities³⁴⁻³⁷. In parallel, multiple variables
27 such as host lifestyle, nutritional needs, genetics, age, medication, urbanization, and the

1 impact of westernization are known to significantly impact the human gut microbiome and
 2 virome^{38,39}. Concurrent with this dynamic interplay between environmental filtering and phage-
 3 bacteria antagonistic and/or mutualistic coevolutionary interactions, one expects concomitant
 4 changes in defensome composition. This prompted us to examine how the defensome's
 5 abundance and diversity correlated with bacterial biogeography. The top five most
 6 represented Classes in our dataset for each environment are Gammaproteobacteria (soil,
 7 marine, human gut), Alphaproteobacteria (soil, marine), Dehalococcoidia (marine),
 8 Bacteroidia (soil, human gut), and Clostridia (human gut) (**Supplementary Table 2**). Such
 9 different patterns in species richness, and relative phylogenetic diversity across environments,
 10 are expected to impact genetic flux and concomitantly, defensome profiles.



1 **Fig. 3.** Defensome variation across different ecological and geographical backgrounds. Defense system (DS)
2 density (per MAG per kb) across distinct ecological (soil, marine) (*a*, *b*) and geographical (human gut) (*c*) contexts.
3 Error bars represent standard deviations and Mann–Whitney–Wilcoxon test *P* values are indicated. **P* < 0.05;
4 ***P* < 10⁻²; ****P* < 10⁻³. Number of MAGs analyzed are shown in parentheses.
5

6 In soil environments, the highest and lowest densities of defense systems were respectively
7 observed in MAGs recovered from serpentine-hosted ecosystems and contaminated or
8 regular soils (**Fig. 3a**). These observations are consistent with the fact that serpentine
9 environments are among the most challenging niches on Earth, characterized by low cellular
10 abundances, limited microbial diversity, high VPRs^{40,41}, and consequently, the likely need for
11 additional anti-MGE systems. Conversely, contaminated and regular surface soils impose a
12 type of environmental stress (namely chemical and UV radiation) that is expected to push
13 phage–bacterium interaction from parasitism to mutualism^{42–44}. The latter should provide
14 bacterial hosts with diversified competitiveness and environmental adaptability while allowing
15 prophages to avoid direct exposure to the stressor. Interestingly, while R–M, CRISPR-Cas
16 and SoFIC were prevalent in MAGs recovered from almost all types of soil, arctic peat (richer
17 in Bacteroidales) stands out as an outlier with a high abundance ($\geq 50\%$) of AbiEii and BREX
18 (**Supplementary Fig. 5a, Supplementary Table 8**). While it remains unclear which processes
19 drive the overrepresentation of these particular defense families in MAGs recovered from
20 arctic peat, the latter could be explained by the cell's need for a second layer of resistance
21 under conditions of high VPRs (see below), or eventually to enforce cooperation between
22 individuals, or even with MGEs^{45,46}.

23 In marine MAGs we observed prevalence of R–Ms, but also of the abortive infection system
24 CBASS and the less-known standalone SoFIC. The highest defense system densities were
25 found in MAGs originating from the arctic ocean (**Fig. 3b**). Such increased defensive repertoire
26 fits previous observations describing high VPRs and virus-to-bacteria contact rates in sea ice
27 brine compared to seawater^{47,48}. Following our observations for ice peat soil, we also found a
28 particularly high abundance (~28%) of the AbiEii system in arctic ocean MAGs
29 (**Supplementary Fig. 5b, Supplementary Table 8**). The overall low defensome abundance

1 and diversity in the Mediterranean Sea can be due to the latter's conditions of seasonal
2 oligotrophic conditions, higher temperature (>13°C), and lower concentrations of inorganic
3 nutrients N and P compared to waters of similar depth in open oceans, leading to very low
4 VPRs⁴⁹.

5 To what concerns human gut MAGs, the difference in amplitude in defense system densities
6 across different countries is more subtle (albeit often significant) and harder-to-interpret
7 compared to other environments. While there is a strong trend in the literature supporting a
8 gradual reduction in microbial diversity (and subsequent disruption of metavirome profiles)
9 concomitant with westernization⁵⁰, the latter did not translate into a clear cut geographical
10 trend in regards to the defensome (**Fig. 3c, Supplementary Fig. 5c**).

11 When defense systems were split according to its mechanism of action, their variation in
12 density across distinct ecological and geographical backgrounds was kept qualitatively the
13 same, at least for the most abundant mechanisms (R-M, Abi, and potential Abi systems)
14 (**Supplementary Fig. 6**).

15 Hence, not only the microbiome but also its defensome is dramatically shaped by different
16 ecological and geographical constraints. Higher densities of defense systems were found in
17 MAGs recovered from particularly challenging biomes such as serpentine soils or the arctic
18 itself, in line with the high VPRs described in such environments.

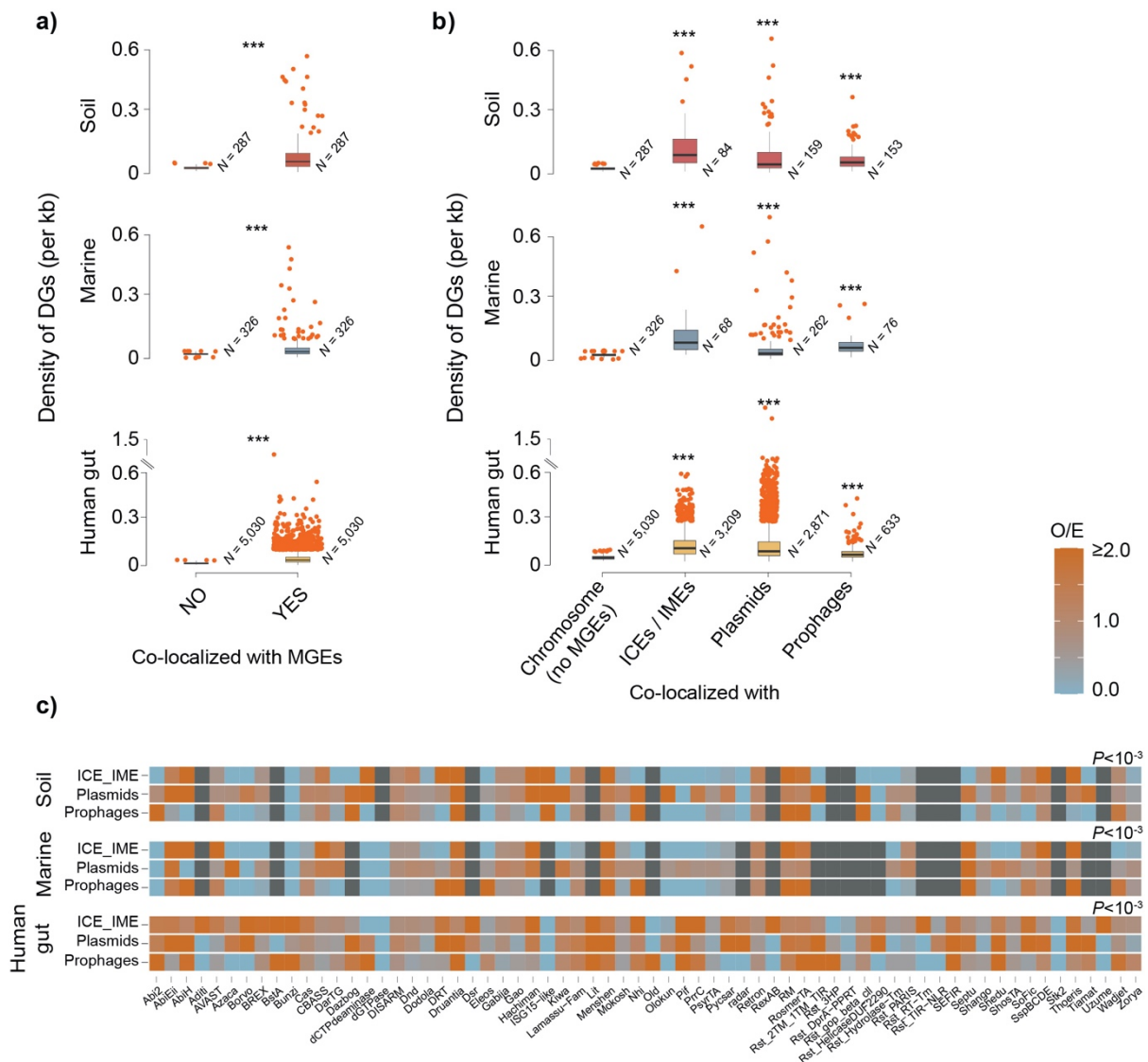
19

20 **The genetic mobility of bacterial MAG defensomes**

21 Cellular defense genes typically propagate by HGT, in a process frequently mediated by
22 MGEs. Physical colocalization between defense genes and MGEs allows for an efficient
23 strategy to modulate and / or resolve potential conflicts in the interactions between the host
24 and the MGE itself. In this context, a growing number of MGE-encoded defense systems or
25 defense genes have been described in several bacteria, particularly involving the most well
26 studied ones (R-Ms, Abi, CRISPR-Cas) and major families of MGEs (phages, plasmids,
27 integrons, ICEs / IMEs)^{13,51,52}. Yet, there is a paucity of data on the genetic mobility of the

1 defensome in complex bacterial environmental communities. We consistently observed more
2 defense genes in MGEs than in chromosomes (excluding MGEs), irrespective of the
3 environment (**Fig. 4a, Supplementary Fig. 7a, Supplementary Tables 9, 10**). This is in line
4 with current evidence that MGE-encoded defense systems protect their host cells as a side-
5 effect of their action to protect the MGE from other MGEs⁵¹. When MGEs were split according
6 to family (excluding integrons which are rare in the human gut microbiota⁵³), there was a slight
7 trend for higher colocalization of defense genes with ICEs / IMEs irrespective of the
8 environment (**Fig. 4b, Methods**), in agreement with recent observations⁵². When integrons
9 were included for comparison, they showed the highest colocalization densities with defense
10 genes in the human gut (**Supplementary Fig. 7b**), a result that should be taken cautiously
11 given its low statistical power.

12 A further split of defense genes according to their corresponding family, allowed us to evaluate
13 the former's over- or underrepresentation across MGE classes (**Fig. 4c, Supplementary Fig.**
14 **7c**). The results put into evidence a few curious aspects of defensome mobility. The first is
15 that irrespectively of the environment, plasmids generally carry a higher than expected by
16 random chance number of defense genes across a large breadth of defense families when
17 compared to other MGE classes. This observation aligns with the fact that plasmids typically
18 allow for a high genetic plasticity and can sustain large gene exchange networks throughout
19 phylogenetically diverse communities⁵⁴.



1
2
3
4
5
6
7
8
9

Fig. 4. The genetic mobility of the defensome. (a) Box plots of the genomic colocalization of defense genes and MGEs. Number of MAGs analyzed are shown as *N* values. (b) Box plots of the genomic colocalization of defense genes with plasmids, prophages, and ICEs / IMEs. Error bars represent standard deviations and Mann–Whitney–Wilcoxon test *P* values are indicated. (c) Heatmap of observed / expected (O / E) ratios of colocalization between genes belonging to distinct defense families and MGEs. Expected values were obtained by multiplying the total number of genes pertaining to a given defense family by the fraction of defense genes of that family assigned to each MGE. * $P < 0.05$; ** $P < 10^{-2}$; *** $P < 10^{-3}$, χ^2 -test. Dark gray squares represent absence of colocalization.

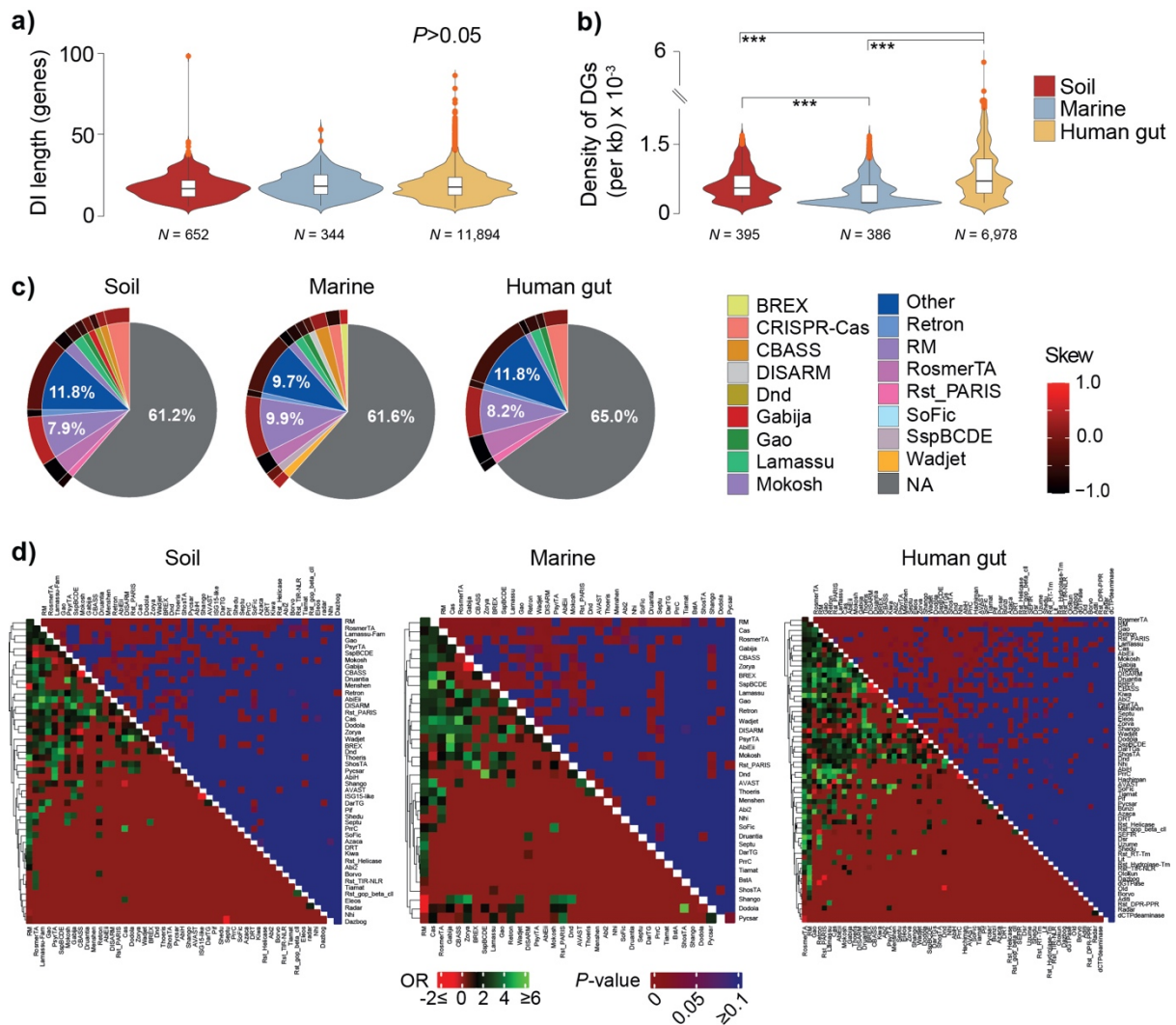
10 The second aspect relates to the highly heterogeneous landscape of combinations defense
11 family / MGE class across multiple environments. This reflects the dynamic interplay between
12 a multitude of parameters, including the density and phylogenetic composition of host cells
13 and MGEs present in the community, habitat structure, and environmental pressures. These
14 results also suggest that certain defense genes / systems favor different classes of MGEs for

1 their shuttling, in a likely dynamic and multilayered interplay with shifting allegiances. Overall,
2 these data shows that a wide range of defense families is carried by MGEs presumably
3 favoring their selfish spread, and that different associations defense family / MGE class are
4 favored across distinct biomes.

5

6 **Encoded functional potential of defense islands and defensome colocalization**

7 Defense genes are typically carried in MGEs by HGT. The former may allow the MGE to be
8 kept in the host by promoting addiction, but on the reverse side of the coin, may carry
9 beneficial traits capable of positive epistatic interactions with the resident host functions. To
10 conciliate these two scenarios, defense genes tend to cluster in so-called defense islands,
11 i.e., high-turnover sinks of genetic diversity, that may serve as catalysts of novel defensive
12 strategies. Therefore, we queried the abundance of such islands and their content. We found
13 12,890 defense islands in 6,217 MAGs (**Supplementary Fig. 8a, Supplementary Table 11a,**
14 **Methods**), with a similar size distribution across environments (median ~ 17 genes) (**Fig. 5a**),
15 suggesting that there is an optimal size range for these defense sinks. Defense gene density
16 in defense islands was significantly lower in marine environments, followed by soil and human
17 gut (**Fig. 5b**). The latter is in line with the above observations on a limited defensome in marine
18 MAGs when compared with other environments. Defense islands' anti-MGE content was very
19 diverse (**Fig. 5c**), with several defense families being overrepresented compared to regions
20 outside defense islands (e.g., Hachiman, R-M, Thoeris) while others being underrepresented
21 (e.g., PsyrTA, ShosTA, Zorya) (**Supplementary Fig. 8b**).



1
2 **Fig. 5.** The MAG defense island repertoire. (a) Defense island (DI) length distribution (given in genes) in soil,
3 marine, and human gut MAGs. (b) Density distribution of defense genes (DGs) (per MAG per kb) across each
4 environment. (c) Pie-plots of the relative abundance (%) of gene content in defense islands. Colored slices
5 correspond to defense genes (those with a relative abundance < 1% were merged as 'Other'), and gray slices (NA)
6 correspond to genes not classified as defensive by DefenseFinder. The outer layer corresponds to the skew ratio
7 between genes belonging to complete and incomplete systems given by
8 $\frac{\#genes\ belonging\ to\ complete\ systems - \#genes\ belonging\ to\ incomplete\ systems}{\#genes\ belonging\ to\ complete\ systems + \#genes\ belonging\ to\ incomplete\ systems}$. (d) Defense families' odds ratio (OR) of
9 colocalization in defense islands (bottom heatmaps) and associated Fisher's exact test P value (upper heatmaps)
10 for the three environments. To eliminate the confounding (inflating) effect of colocalized genes pertaining to the
11 same system, we only considered solitary genes or those pertaining to independent defense systems distanced of
12 5 genes or less.

13

14 This bias for certain defense families to locate in defense islands, suggests either positive
15 epistatic interactions with vicinal genes, or a preferential shuttling by a particular family of
16 MGEs. Despite its diversity, defense families are largely similar across environments and

1 skewed towards incomplete systems, pointing as expected, towards a high gene turnover at
2 defense islands (**Fig. 5c**). Interestingly, the large majority (~63%) of defense islands' gene
3 content was not predicted to have a defensive role. A COG classification of such 'non-
4 defensive' genes revealed a high prevalence of functions linked to replication / recombination
5 / repair and transcription (**Supplementary Fig. 8c**). The latter can be at least partially
6 explained by the fact that defense genes are often shuttled by MGEs, which rely on such
7 functions for target selectivity, insertion, and excision. The above COG categories and the
8 most abundant defense families (R–M and CRISPR-Cas for soil / human gut; R–M, CBASS
9 and RosmerTA for marine biomes) remained unchanged even when considering defense
10 systems (instead of genes) as the main counting unit in the definition of defense islands (see
11 **Methods**) (**Supplementary Fig. 8d, Supplementary Table 11b**).

12 Since MGEs have different distribution patterns, we quantified the frequency of colocalization
13 of defensome families (≤ 5 genes apart) in defense islands compared to regions outside the
14 latter (**Fig. 5d, Supplementary Table 12, Methods**). In line with their abundance, frequent
15 shuttling by MGEs and defensive role, R–Ms significantly colocalized with most other defense
16 families in defense islands irrespectively of the environment. Inversely, R–Ms showed a
17 preference to colocalize with genes pertaining to Menshen, Shango and Dodola families
18 outside defense islands. Interestingly, and despite their general underrepresentation in
19 defense islands (**Supplementary Fig. 8b**), genes pertaining to families such as PsyrTA and
20 Zorya showed significant colocalization with other defense families inside defense islands.
21 Conversely, defense families significantly overrepresented in defense islands (e.g.,
22 Hachiman) (**Supplementary Fig. 8b**), rarely colocalized with other families. Upon splitting our
23 dataset according to biogeographical zones, and despite the subsequent decrease in
24 statistical power, the colocalization trends of the most abundant defense families still hold
25 qualitatively (**Supplementary Fig. 9, Supplementary Table 13**). These observations point for
26 the possibility of previously unappreciated epistatic interactions between selected families of
27 defense genes / systems in defense islands.

1 Hence, we found approximately 11% of the defensome concentrated in defense islands, an
2 environment-dependent highly heterogeneous distribution of defense families in such regions,
3 a large proportion of ‘non-defensive’ functions, and a significant colocalization of a subset of
4 families of defense.

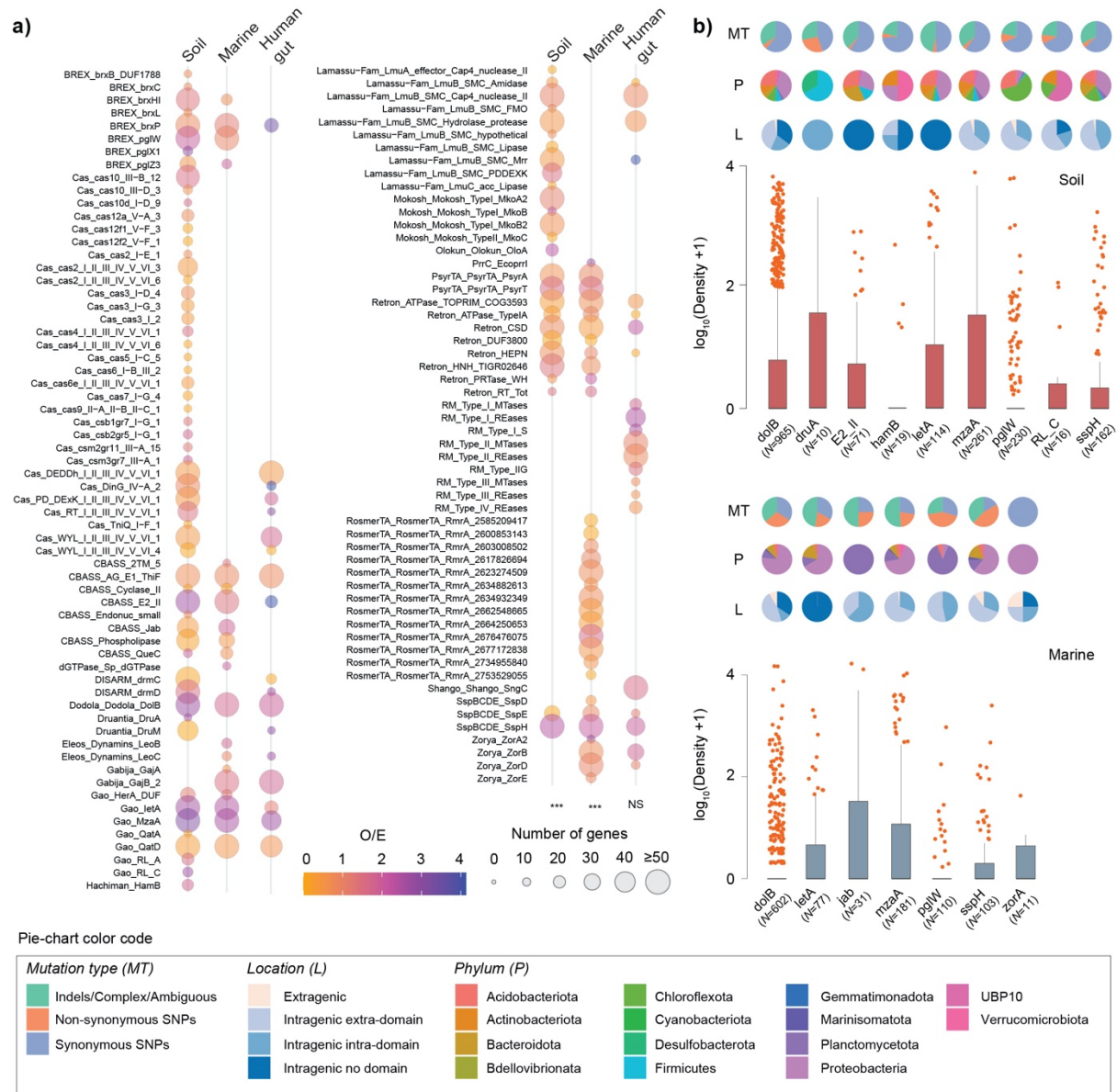
6 The genetic variability of the defensome

7 The coevolutionary dynamics between defenses and counter-defenses contributes to an
8 endless process of genetic diversification and evolution of sequence specificity, that can take
9 place through point mutations, recombination, gene duplications, replication slippage, or
10 transposition⁵⁵. Such panoply of diversification processes has been particularly studied in well-
11 described innate immune systems like R–Ms, and can take the form of, for example, target
12 recognition domain swapping in Type I *hsdS* subunits, or phase variability of Type III *modH*
13 genes. However, there is a void in our current understanding on the extent to which differences
14 in selection strength act across distinct defensome gene families. To this end, we performed
15 metagenome read recruitment over defensome genes, assessed the frequency and type of
16 short variants found, and used this information to pinpoint consistently fast or slow evolving
17 genes across environments (see **Methods** for further details).

18 We observed multiple defense gene families with higher-than-expected values of SNP + indel
19 density across multiple biomes (**Fig. 6a**, only defense families for which at least one defense
20 gene showing a O / E ratio ≥ 1.5 per environment are shown, **Supplementary Table 14**).
21 Genes such as *doiB*, *mzaA*, and *sspH* were among this ‘high-mutation frequency’ subset
22 irrespective of the environment, while others like *druA*, *zorA*, or *letA* were environment-
23 specific. The results were qualitatively similar when all defense families were included
24 (**Supplementary Fig. 10**). The range of SNP + indel densities differed considerably across
25 defense gene families (**Fig. 6b**). Mutation types were also profoundly affected by the
26 environment (and thus by population structure). For example, indels and non-synonymous
27 SNPs were consistently more abundant in marine than in soil MAGs, even when comparing

1 across same defense families (e.g., *dolB*, *mzaA* or *sspH*) (**Fig. 6b**). While most variants found
2 was intragenic, *sspH* and particularly *zorA* had as much as 25% of variants located in the first
3 200 bp upstream the annotated start codons, suggesting potential regulatory effects. The rapid
4 turnover of defense gene repertoires in bacteria, many of which in MGEs, can be followed by
5 selection for the former's conservation or loss in a cell. To investigate the action of natural
6 selection on the defensome gene families showing the highest frequency of variants, we
7 computed the ratio of nonsynonymous over synonymous substitution rates (dN/dS) for the
8 pools of orthologous defense genes within our MAG dataset. Similar to previous observations
9 for CRISPR-Cas and R-M gene families^{56,13}, all defense genes analyzed were found to be
10 under strong purifying selection (dN/dS<<1; **Supplementary Fig. 11a, Supplementary Table**
11 **15**). The preferential purge of nonsynonymous mutations by natural selection contributes to
12 maintain the defensive functions of these genes and can be reconciled with a scenario of high
13 turnover, if the selection pressure on the system fluctuates in time, i.e., if these genes alternate
14 periods of strong purifying selection and periods of relaxed selection (e.g., as a result of
15 competition with other defense systems, or during strong selection for HGT). Interestingly,
16 despite their overall negative selection, we observed relatively high levels of divergence and
17 positive selection in certain portions of their sequences (**Supplementary Fig. 11b**). The latter
18 matched, for example, PFAM domains with predicted AAA+ ATPase activity (PF07724 /
19 PF10431 in *DolB*, and to a less extent PF00004 / PF17862 in *letA*), an *ftsH*-like extracellular
20 domain (PF06480 in *letA*), and a Sigma70-like non-essential domain (PF04546 in *MzaA*).

21
22
23
24
25



1
2 **Fig. 6.** The genetic variability of the defensome. (a) 90 metagenomes (30 for each environment) having a broad
3 representativity in terms of sampling sites (soil and marine) and countries (human gut), as well as in terms of
4 presence of most defense families previously identified by DefenseFinder were selected. Shown in circles are the
5 observed / expected (O / E) ratios of number of defense gene families harboring high-frequency SNPs + indels
6 ($\geq 25\%$ at the variant position) in their gene body (including 200 bp upstream the start codon). Expected values
7 were obtained by multiplying the total number of genes pertaining to a given defense family by the fraction of
8 defense genes of that family harboring high-frequency alleles. Circle radius corresponds to the total number of
9 defense genes analyzed per family. To ease visualization, we limited the figure to only defense families for which
10 at least one defense gene showed an O / E ratio ≥ 1.5 per environment. The complete representation is shown in
11 **Supplementary Fig. 10.** * $P < 0.05$; ** $P < 10^{-2}$; *** $P < 10^{-3}$, χ^2 -test. NS: not significant. (b) Density distribution of
12 SNPs + indels for a selection of defense gene families showing the highest O / E values in (a). Information on
13 mutation type, location and phylum are indicated in pie-plots. The number of genes analyzed is shown in
14 parentheses.

1 Discussion

2 In this study, we present a large-scale analysis of the abundance and diversity of defensomes
3 of genomes/species from complex microbial communities and three representative biomes:
4 soil, marine and human gut. Our results on the quantification of the defensome in marine
5 MAGs lend support to a scenario of a limited defense arsenal in this environment (**Fig. 2b**).
6 The latter can be accounted by a variety of potential explanations namely: *i*) the fact that
7 oligotrophic open oceans typically show an overrepresentation of clades characterized by
8 heavily streamlined genomes⁵⁷ (e.g., Dadabacteria, Chloroflexota) (**Supplementary Fig 1**,
9 **Supplementary Table 2**), and thus, more likely to opt for more transient defense systems and
10 little metabolic plasticity to better cope with the limiting environment of the surface ocean; *ii*)
11 the dominantly planktonic lifestyle and low cell-density in the marine environment (at least for
12 the free-living fractions accounted for in our MAG dataset) which may in itself, or through a
13 reduced frequency of HGT, contribute to a more limited anti-MGE arsenal; *iii*) the fact that the
14 large majority of HMMs currently available to detect defense systems were initially developed
15 on the basis of genetic data that overrepresents not only cultivable bacteria, but also lineages
16 (e.g., *Escherichia*, *Bacillus*, *Pseudomonas*) that are more distantly related to those that make
17 up the global ocean microbiome (**Supplementary Table 3**). On a broader view, our results
18 qualitatively match those recently obtained for RefSeq genomes in terms of the most abundant
19 systems (R–Ms, CRISPR-Cas) and overall diversity of families identified¹¹. The enhanced
20 granularity offered by our cross-environment comparison revealed a few curious differences
21 at the level of preferential ‘second line’ defense families. One of such differences concerns
22 SoFIC and CBASS which are present in roughly 20% of soil and marine MAGs (mainly in
23 Chitinophagales and Caulobacterales), but considerably less predominant (~8%) in human
24 gut MAGs (mainly in Verrucomicrobiales, Enterobacterales, and Bacteroidales) (**Fig. 2a**).
25 Inversely, the abortive infection system Rst_PARIS is present in 20% of human gut MAGs
26 (mainly in Bacteroidales) but is virtually absent in soil or marine environments (**Fig. 2a**). While
27 R–Ms (and to a lesser extent CRISPR-Cas) are largely ubiquitous, our results are supportive

1 of a 'second line' of defense systems (SoFIC, CBASS, Rst_Paris, etc) that is also mostly non-
2 species-specific, differentially favored across distinct environments, and privileged by
3 distinctive strategies of genetic mobility (**Fig. 4c**, see below). As we move down the ladder of
4 defense system abundance, we face an increasing variety of cryptic, presumably highly
5 specialized, and more species / population-specific systems. By further splitting our dataset
6 into sub-environments or by geographic location, we observed significant differences in
7 defense system abundance (**Fig. 3**). And while the increased densities observed at serpentine
8 systems and across the arctic ocean can be explained by the extreme conditions experienced
9 at such environments and a subsequent phage-bacteria imbalance, the more subtle variations
10 in defense system density in human gut MAGs across multiple countries and the panoply of
11 confounding variables associated, preclude the identification of more explanatory scenarios.
12 Higher densities of defense genes were consistently observed in (or in the close vicinity) of
13 MGEs compared to those found in the chromosome (excluding MGEs) (**Fig. 4a**). Such
14 colocalization facilitates the rapid acquisition and / or diversification of the defense to
15 provide resistance against multiple other MGEs. It was recently suggested that the carrying of
16 certain defense systems by MGEs by a given bacterial host, may not always relate with the
17 latter's need for protection, but instead on the best interest of the MGE itself in order to
18 overcome or displace antagonistic MGEs⁵¹. Our observation of a complex and heterogeneous
19 distribution of defense gene families across different classes of MGEs supports such
20 hypothesis and suggests an exploitation of MGEs by defense genes / systems for purposes
21 other than host defense. It ultimately highlights the need to better understand the molecular,
22 and evolutionary interactions between the threesome host-phage-mobilome.

23 Genes acquired by HGT, and MGEs in particular, tend to integrate in a small number of
24 chromosome hotspots to decrease the fitness cost of their integration. Successive rounds of
25 integration / excision / partial deletion of MGEs when accompanied by the co-option of defense
26 genes / systems may result in the formation of defense islands. While initially thought that the
27 latter were merely "genomic junkyards" in which the defense genes that are frequently
28 acquired via HGT accumulate because insertion in these regions is unlikely to be deleterious,

1 it has now become clear that there is a specific selective advantage in such clustering of
2 genes, such as functional cooperation between different defensive modules and generation
3 of novel functions. When compared across environments, defense islands did not show
4 significant differences in terms of size, relative abundances of major defense families, or at
5 the topmost abundant COG functional categories for genes classified as 'non-defensive' (**Fig.**
6 **5a, c, Supplementary Fig. 8c**). While many of these genes seem to encode factors involved
7 in genetic mobility, others have hitherto unknown functions. In this line, an interesting next
8 step would be to build upon our precise delimitation of defense islands in such a large and
9 phylogenetically diverse MAG dataset and use a previously developed colocalization
10 framework³ to leverage the identification of novel defense systems. A significant
11 overrepresentation of several defense families (e.g., Hachiman, R–M, Thoeris) was observed
12 in defense islands (compared to non-island regions). Yet, for certain of these families, such
13 overabundance was not synonymous of a higher likelihood of colocalization with the
14 remainder of the defensome (and vice-versa). These observations point for the possibility of
15 previously unappreciated epistatic interactions or increased probability of functional
16 diversification for a selected subset of families of defense genes / systems in defense islands.
17 In this regard, the extent to which non-canonical HGT mechanisms (e.g., gene transfer agents,
18 nanotubes, membrane vesicles) and MGE-independent mechanisms of diversification (e.g.,
19 homologous recombination) respectively shape the movement of defense genes and the
20 evolution of defense islands remains unclear.

21 Under the Red Queen evolutionary dynamics, the coevolution between opposing hosts and
22 parasites portrays evolution as a never-ending evolutionary arms-race between defense and
23 counter-defense strategies. Such antagonistic coevolution pervades evolutionary change
24 through multiple ingenious strategies, including: *i*) point mutations in phage DNA recognition
25 sites to reduce the likelihood of restriction by R–M systems⁵⁸; *ii*) phase-variation / inversions /
26 point mutations in MTases, REases or S modules leading to altered R–M systems'
27 specificity^{55,59}; *iii*) ON/OFF switch in CRISPR immunity through mutations in *cas* genes⁶⁰;
28 among others. Thus, not only turnover and recombination, but also rapid sequence evolution

1 of certain defense genes / systems through mutation are key factors shaping the host-parasite
2 evolutionary trajectory. Such diversification strategies are a function of the size and the
3 diversity of the defensome gene pool in a bacterial population, and will shape how the latter
4 remains evolutionarily responsive to temporally or spatially variable selection imposed by
5 phages. Different defense genes are expected to evolve at different rates. For example,
6 significant differences in purifying selection have been described across different Types of R-
7 M REases and MTases, highlighting distinct signatures of adaptive evolution¹³. To gain a
8 birds-eye-view on potentially coexisting subpopulations bearing substantial defense gene-
9 level diversity, we built upon a metagenome read recruitment approach. This allowed us to
10 identify a subset of defense genes having a higher-than-expected frequency of SNPs + indels,
11 globally evolving under strong purifying selection, and a heterogeneous landscape of mutation
12 types profoundly affected by the environment (and thus by population structure). Whereas for
13 some of these genes we can point out determinants capable of explaining such observations
14 – namely the presence of domains known for their predisposition to genetic variation (e.g., the
15 motility-associated *motA* domain⁶¹ in *zorA*, or the *ftsk* translocase domain⁶² in *sspH*) –, the
16 lack of substantial functional and mechanistic insights on the remaining ones (and on their
17 systems) precludes further meaningful ascertainties.

18 It is important to appreciate that our computational analysis is challenged by a few difficult-to-
19 control confounding variables and limitations that are worth discussing. The first, concerns the
20 imbalance in our dataset between the number of samples recovered from each biome, as well
21 as their geographic distribution. While the number of soil and marine MAGs analyzed was
22 essentially the same, human gut MAGs were roughly one order of magnitude greater. From
23 the geographic standpoint, marine samples have a global representation, but soil and human
24 gut microbiome data are greatly biased towards the US and China. These observations
25 highlight a critical need for thorough geographic sampling, more global representation of
26 participants in microbiome studies, and a fairer access to genomics resources especially in
27 resource-poor countries. A second confounding variable, likely more relevant, concerns the
28 fact that MAG binning methods using short reads tend to miss certain low-abundance or

1 difficult to resolve MGE families. The fact that defense genes are often carried or colocalize
2 with MGEs, necessarily indicates that our results *i*) may have a bias in the ratio of defense
3 genes inside versus outside the mobilome, and *ii*) are most likely a partial underestimated
4 picture of the real defensome abundance. Future inclusion of long-read data will enable
5 reference-quality genome reconstruction from metagenomes, and further improve our
6 observations. Third, our observations are not representative of all bacterial communities and
7 are likely influenced by characteristics of the sampled environments. Still, the stringent dataset
8 filtering used in our study in terms of MAG completeness and N50 (with associated controls
9 shown in **Supplementary Fig. 1**), together with previous demonstration on the accuracy of
10 MAG size estimates (that are part of our dataset) compared with reference genomes²⁶, makes
11 us have good reasons to think that our analyses constitute a reasonable proxy of the defense
12 landscape diversity carried by such populations, and of the complex interplay underlying their
13 interactions at the intra- and inter-environment level. Lastly, while this study provides novel
14 and intriguing insights into the defensome co-localization, it does not address the specific
15 mechanisms and interactions between different systems, nor the interplay with phage counter-
16 defense strategies^{63,64}.

17 The efforts recently undertaken to identify novel defense mechanisms in typically easily
18 cultivable bacteria must now be followed by initiatives to expand the search to uncultivated
19 microbes in complex microbial communities, to understand how such mechanisms collaborate
20 or antagonize with one another, how they co-opt or are co-opted by MGEs, and how they are
21 shaped by the surrounding environment. Our work provides a first steppingstone in such a
22 direction.

23

1 **Methods**

2 **Data**

3 In this study we built upon a large dataset of 7,759 high-quality soil, marine, and human gut
4 MAGs^{24–26} (**Supplementary Table 1**). These MAGs were filtered on the basis of the Minimum
5 Information about a Metagenome Assembled Genome (MIMAG) standard ($\geq 90\%$
6 completeness, $\leq 5\%$ contamination/redundancy, ≥ 18 tRNA genes and presence of at least
7 one class of 5S, 16S and 23S rRNA genes). When not clearly stated in the original study, we
8 performed identification of rRNA genes using both Infernal⁶⁵ v1.1.4 (options: `-Z 1000 --`
9 `hmmonly --cut_ga --noali -tblout`) and RNAmmer⁶⁶ v1.2 (options: `-S bac -m tsu,ssu,lsu -h -f -`
10 `gff`) (**Supplementary Table 3**). Since defense systems are often *i*) multigenic and *ii*) clustered
11 in defense islands, we further selected for highly contiguous MAGs to more accurately reflect
12 the defensome abundance and distribution. In particular, we selected assemblies having
13 values of $N_{50} \geq 100$ kb (corresponding to at least the top 99.5 % best assemblies), and
14 repeated the analyses for $N_{50} \geq 200$ and 300 kb (chosen upon visual inspection of the density
15 distribution) (**Supplementary Fig. 1**) to account for the effect of contiguity in our observations.
16 MAG annotation was performed with PROKKA⁶⁷ v1.14.5 (default parameters).

17

18 **Identification of anti-MGE defense genes, systems and islands**

19 MAGs were queried for anti-MGE defense genes / systems using DefenseFinder¹¹ v1.0.8
20 (option: `--preserve-raw`). The current version of this tool allows for the screening of 1,024
21 genes pertaining to 127 families of anti-MGE defense systems. Defense islands were defined
22 as arrays of defense genes (or defense systems) separated from one another by 10 genes or
23 less and with a minimum of 5 genes pertaining to at least 3 different defense families.
24 Functional annotation of 'non-defensive' genes was performed with eggNOG-mapper⁶⁸ v.2.1.9
25 (default parameters). To test for colocalization of defense families in defense islands, we
26 computed their odds ratio and associated Fisher's exact test P value. For this purpose, we
27 considered all colocalized defense genes distanced of 5 genes or less both inside and outside

1 defense islands. Genes belonging to the same defense system are necessarily colocalized,
2 so we deliberately eliminated such hits to avoid inflating same system colocalization
3 frequencies. To determine the presence of putative defense system regulators harboring WYL
4 or Caspase Recruitment Domains (CARD), all MAG proteomes were scanned against the
5 Pfam-HMMs PF13280 (WYL) and PF00619 (CARD) using HMMER3⁶⁹ and a cut-off *e*-value
6 of 0.01.

7

8 **Identification of mobile genetic elements**

9 Classification of contigs as belonging to chromosomes or plasmids was performed using
10 PlasClass⁷⁰ v.0.1.1 and PlasFlow⁷¹ v.1.1 (both with default parameters). Plasmid hits were
11 selected as those with a score greater than or equal to 0.7. Integrons were identified using
12 IntegronFinder⁷² v.2.0.1 (option --local_max). Prophages were detected with Virsorter2⁷³
13 v.2.2.3 (options --include-groups dsDNAphage,ssDNA --min-length 5000 --min-score 0.5).
14 Despite recent evidence for phage satellites carrying defense systems⁷⁴, we deliberately
15 excluded them from our analyses, mainly due to the very few examples of experimentally
16 validated satellites (particularly in non-cultivable bacteria), which precludes the development
17 of robust detection tools and an accurate evaluation of their classification. Integrative
18 Conjugative Elements (ICEs) and Integrative Mobilizable Elements (IMEs) were detected with
19 ICEfinder⁷⁵ v.2.6.32-696.10.2.el6.x86_64 (default parameters). All MGE hits matching multiple
20 families were not considered in the analyses (~2.8% of the total MGE dataset detected). While
21 MGE carriage by other MGEs (e.g., integrons by plasmids) is indeed expected, we deliberately
22 eliminated such hits to avoid the confounding effects of their co-occurrence on the defensome
23 analyses.

24

25 **Phylogenetic analyses**

26 For phylogenetic tree construction we took for each MAG a concatenate of 15 ribosomal
27 proteins (L2, L3, L4, L5, L6, L14, L16, L18, L22, L24, S3, S8, S10, S17, S19), aligned them

1 with MAFFT⁷⁶ v7.490 (options: --maxiterate 1000 -globalpair) (soil, marine) or Muscle⁷⁷ v.5.1
2 (option: -super5) (human gut), and trimmed poorly aligned regions with BMGE⁷⁸ v2.0 (option:
3 -t AA). To avoid plotting poorly supported branches, MAGs harboring less than 50% of the
4 abovementioned ribosomal list were omitted from the phylogenetic representations (> 95 %
5 had the expected number of proteins across the three environments). The trees were
6 computed by maximum likelihood with RaxML⁷⁹ v8.2.12 (options: raxmlHPC-PTHREADS-
7 AVX -f a -m PROTGAMMAAUTO -N autoMRE -p 12345 -x 12345) (soil, marine) or iqtree2⁸⁰
8 v2.2.6 (options: -nt 56 -cmax 15 -bb 1000 -alrt 1000 -m TESTNEW -safe) (human gut)
9 (**Supplementary Table 5**). The phylogenetic depth was defined as the average root-to-tip
10 distance, and was computed as the diagonal mean of the phylogenetic variance–covariance
11 matrix of each tree, using the vcv.phylo function in the R package “ape”.

12

13 **Variant analysis of the defensome**

14 To evaluate which defense gene families are preferential targets for increased genetic
15 diversity (SNPs + indels), we selected 90 metagenomes (30 for each environment with similar
16 sequencing depth) having a broad representativity in terms of sampling sites (soil and marine)
17 and countries (human gut), as well as in terms of presence of most defense families that were
18 characteristic to each environment. Fragment recruitment was performed by mapping
19 metagenomic reads from each sample against the ensemble of defense genes (including 200
20 bp upstream of the start codon) pertaining to the previously selected 90 metagenomes using
21 BWA-MEM⁸¹ v.0.7.17 (default parameters). Genetic variants were identified from aligned
22 reads with FreeBayes⁸² v1.1.0 (options: freebayes-parallel -p 1 -P 0 -C 1 -F 0.025 --min-
23 repeat-entropy 1.0 -q 13 -m 60 --strict-vcf --f) and a subsequent filtering step was performed
24 to select only genes (including upstream regions) containing variants having a minimum
25 frequency of 25% supported by at least 10 reads. A minimum of 10 genes per defense family
26 per environment was considered in the analysis. Alignments were visualized using IGV
27 v.2.14.1. Finally, SNPGenie⁸³ v1.0 (options: --vcfformat=4 --snpreport --fastafilename --gtffilename --

1 outdir) was used for variant classification. For each environment, we computed the observed
2 / expected (O / E) ratio of defense genes harboring high-frequency alleles across all defense
3 families. Expected values were obtained by multiplying the total number of genes pertaining
4 to a given defense family by the fraction of defense genes of that family harboring high-
5 frequency alleles.

6

7 **Analysis of substitution rates**

8 All-against-all BLASTP searches were performed on the sets of defense genes scanned in
9 the genomes (default settings, e -value $<10^{-3}$). Clustering was performed using the SILIX
10 package⁸⁴ v.1.3.0 using a minimum identity threshold of 80% and default values for the
11 remaining parameters. Singletons were eliminated from our data set. The remaining protein
12 sequences (putative orthologs) were reverse-translated to the corresponding DNA sequences
13 using PAL2NAL⁸⁵ v14. Pairwise rates of non-synonymous substitutions (dN), synonymous
14 substitutions (dS) and ω (dN/dS) were computed using the KaKs_Calculator⁸⁶ v.2.0
15 implementing the Yang-Nielsen⁸⁷ and Nei-Gojobori⁸⁸ methods. Estimations yielding dS > 1
16 (corresponding to situations of substitution saturation and representing 0.2 % of the total data)
17 were discarded to improve the quality of estimation of ω .

18

19 **Statistical and graphical analyses of data**

20 All statistical and graphical analyses were conducted using R v.4.3.1. Geographical
21 representation of metagenome sampling locations was generated using the *mapdata*
22 package. Visualization of genomic contexts was performed with the package *gggenes*.
23 Colocalization heatmaps were created using the ComplexHeatmap package. Stepwise linear
24 regression analyses were performed by using the *step* function from the *stats* package.

25

1 **Data availability**

2 All data supporting the findings of this study are available within the article and its
3 supplementary files. Source data are provided with this paper.

4

5 **Code availability**

6 Wrapper scripts supporting all key analyses of this work are publicly available
7 at <https://github.com/oliveira-lab/Defensome>.

1 References

- 2 1. Haudiquet, M., De Sousa, J. M., Touchon, M. & Rocha, E. P. C. Selfish, promiscuous and sometimes useful:
3 how mobile genetic elements drive horizontal gene transfer in microbial populations. *Philos. Trans. R. Soc. B*
4 *Biol. Sci.* **377**, 20210234 (2022).
- 5 2. Oliveira, P. H., Touchon, M. & Rocha, E. P. C. Regulation of genetic flux between bacteria by restriction–
6 modification systems. *Proc. Natl. Acad. Sci.* **113**, 5658–5663 (2016).
- 7 3. Doron, S. *et al.* Systematic discovery of antiphage defense systems in the microbial pangenome. *Science* **359**,
8 eaar4120 (2018).
- 9 4. Millman, A. *et al.* An expanded arsenal of immune systems that protect bacteria from phages. *Cell Host Microbe*
10 **30**, 1556–1569.e5 (2022).
- 11 5. Bernheim, A. & Sorek, R. The pan-immune system of bacteria: antiviral defence as a community resource. *Nat.*
12 *Rev. Microbiol.* **18**, 113–119 (2020).
- 13 6. Mayo-Muñoz, D., Pinilla-Redondo, R., Birkholz, N. & Fineran, P. C. A host of armor: Prokaryotic immune
14 strategies against mobile genetic elements. *Cell Rep.* **42**, 112672 (2023).
- 15 7. Rostøl, J. T. & Marraffini, L. (Ph)ighting phages: How bacteria resist their parasites. *Cell Host Microbe* **25**, 184–
16 194 (2019).
- 17 8. Kobayashi, I. Behavior of restriction-modification systems as selfish mobile elements and their impact on
18 genome evolution. *Nucleic Acids Res.* **29**, 3742–3756 (2001).
- 19 9. Lopatina, A., Tal, N. & Sorek, R. Abortive infection: Bacterial suicide as an antiviral immune strategy. *Annu.*
20 *Rev. Virol.* **7**, 371–384 (2020).
- 21 10. Koonin, E. V. & Makarova, K. S. Origins and evolution of CRISPR-Cas systems. *Philos. Trans. R. Soc. B Biol.*
22 *Sci.* **374**, 20180087 (2019).
- 23 11. Tesson, F. *et al.* Systematic and quantitative view of the antiviral arsenal of prokaryotes. *Nat. Commun.* **13**,
24 2561 (2022).
- 25 12. Payne, L. J. *et al.* PADLOC: a web server for the identification of antiviral defence systems in microbial
26 genomes. *Nucleic Acids Res.* **50**, W541–W550 (2022).
- 27 13. Oliveira, P. H., Touchon, M. & Rocha, E. P. C. The interplay of restriction-modification systems with mobile
28 genetic elements and their prokaryotic hosts. *Nucleic Acids Res.* **42**, 10618–10631 (2014).
- 29 14. Whitman, W. B., Coleman, D. C. & Wiebe, W. J. Prokaryotes: The unseen majority. *Proc. Natl. Acad. Sci.* **95**,
30 6578–6583 (1998).
- 31 15. Hofer, U. The majority is uncultured. *Nat. Rev. Microbiol.* **16**, 716–717 (2018).
- 32 16. Mushegian, A. R. Are there 10^{31} virus particles on Earth, or more, or fewer? *J. Bacteriol.* **202**, (2020).
- 33 17. Suttle, C. A. Marine viruses — major players in the global ecosystem. *Nat. Rev. Microbiol.* **5**, 801–812 (2007).
- 34 18. Tyson, G. W. *et al.* Community structure and metabolism through reconstruction of microbial genomes from the
35 environment. *Nature* **428**, 37–43 (2004).
- 36 19. Alneberg, J. *et al.* Binning metagenomic contigs by coverage and composition. *Nat. Methods* **11**, 1144–1146
37 (2014).
- 38 20. Chen, L.-X., Anantharaman, K., Shaiber, A., Eren, A. M. & Banfield, J. F. Accurate and complete genomes from
39 metagenomes. *Genome Res.* **30**, 315–333 (2020).
- 40 21. Sender, R., Fuchs, S. & Milo, R. Revised Estimates for the Number of Human and Bacteria Cells in the Body.
41 *PLOS Biol.* **14**, e1002533 (2016).
- 42 22. Sausset, R., Petit, M. A., Gaboriau-Routhiau, V. & De Paepe, M. New insights into intestinal phages. *Mucosal*
43 *Immunol.* **13**, 205–215 (2020).
- 44 23. Parikka, K. J., Le Romancer, M., Wauters, N. & Jacquet, S. Deciphering the virus-to-prokaryote ratio (VPR):
45 insights into virus-host relationships in a variety of ecosystems. *Biol. Rev.* **92**, 1081–1100 (2017).
- 46 24. Nayfach, S. *et al.* A genomic catalog of Earth’s microbiomes. *Nat. Biotechnol.* **39**, 499–509 (2021).
- 47 25. Almeida, A. *et al.* A unified catalog of 204,938 reference genomes from the human gut microbiome. *Nat.*
48 *Biotechnol.* **39**, 105–114 (2021).
- 49 26. Paoli, L. *et al.* Biosynthetic potential of the global ocean microbiome. *Nature* **607**, 111–118 (2022).
- 50 27. Oliveira, P. H. & Fang, G. Conserved DNA methyltransferases: A window into fundamental mechanisms of
51 epigenetic regulation in bacteria. *Trends Microbiol.* **29**, 28–40 (2021).
- 52 28. Sheinman, M. *et al.* Identical sequences found in distant genomes reveal frequent horizontal transfer across
53 the bacterial domain. *eLife* **10**, e62719 (2021).
- 54 29. Rousset, F. & Sorek, R. The evolutionary success of regulated cell death in bacterial immunity. *Curr. Opin.*
55 *Microbiol.* **74**, 102312 (2023).
- 56 30. Makarova, K. S., Anantharaman, V., Grishin, N. V., Koonin, E. V. & Aravind, L. CARF and WYL domains: ligand-
57 binding regulators of prokaryotic defense systems. *Front. Genet.* **5**, (2014).
- 58 31. Blankenchip, C. L. *et al.* Control of bacterial immune signaling by a WYL domain transcription factor. *Nucleic*
59 *Acids Res.* **50**, 5239–5250 (2022).
- 60 32. Picton, D. M. *et al.* A widespread family of WYL-domain transcriptional regulators co-localizes with diverse
61 phage defence systems and islands. *Nucleic Acids Res.* **50**, 5191–5207 (2022).
- 62 33. Wein, T. *et al.* CARD-like domains mediate anti-phage defense in bacterial gasdermin systems.
63 <http://biorxiv.org/lookup/doi/10.1101/2023.05.28.542683> (2023) doi:10.1101/2023.05.28.542683.

- 1 34. Nogales, B., Lanfranconi, M. P., Piña-Villalonga, J. M. & Bosch, R. Anthropogenic perturbations in marine
2 microbial communities. *FEMS Microbiol. Rev.* **35**, 275–298 (2011).
- 3 35. Mojica, K. D. A. & Brussaard, C. P. D. Factors affecting virus dynamics and microbial host - virus interactions
4 in marine environments. *FEMS Microbiol. Ecol.* **89**, 495–515 (2014).
- 5 36. Williamson, K. E., Fuhrmann, J. J., Wommack, K. E. & Radosevich, M. Viruses in soil ecosystems: an unknown
6 quantity within an unexplored territory. *Annu. Rev. Virol.* **4**, 201–219 (2017).
- 7 37. Mishra, A., Singh, L. & Singh, D. Unboxing the black box—one step forward to understand the soil microbiome:
8 A systematic review. *Microb. Ecol.* **85**, 669–683 (2023).
- 9 38. Martiny, J. B. H. *et al.* Microbial biogeography: putting microorganisms on the map. *Nat. Rev. Microbiol.* **4**, 102–
10 112 (2006).
- 11 39. Hanson, C. A., Fuhrman, J. A., Horner-Devine, M. C. & Martiny, J. B. H. Beyond biogeographic patterns:
12 processes shaping the microbial landscape. *Nat. Rev. Microbiol.* **10**, 497–506 (2012).
- 13 40. Twing, K. I. *et al.* Serpentinization-influenced groundwater harbors extremely low diversity microbial
14 communities adapted to high pH. *Front. Microbiol.* **8**, (2017).
- 15 41. Thomas, E., Anderson, R. E., Li, V., Rogan, L. J. & Huber, J. A. Diverse viruses in deep-sea hydrothermal vent
16 fluids have restricted dispersal across ocean basins. *mSystems* **6**, e00068-21 (2021).
- 17 42. Drew, G. C., Stevens, E. J. & King, K. C. Microbial evolution and transitions along the parasite–mutualist
18 continuum. *Nat. Rev. Microbiol.* **19**, 623–638 (2021).
- 19 43. Huang, D. *et al.* Enhanced mutualistic symbiosis between soil phages and bacteria with elevated chromium-
20 induced environmental stress. *Microbiome* **9**, 150 (2021).
- 21 44. Tang, X. *et al.* Lysogenic bacteriophages encoding arsenic resistance determinants promote bacterial
22 community adaptation to arsenic toxicity. *ISME J.* **17**, 1104–1115 (2023).
- 23 45. Dy, R. L., Przybilski, R., Semeijn, K., Salmond, G. P. C. & Fineran, P. C. A widespread bacteriophage abortive
24 infection system functions through a Type IV toxin–antitoxin mechanism. *Nucleic Acids Res.* **42**, 4590–4605
25 (2014).
- 26 46. Trubl, G. *et al.* Active virus–host interactions at sub-freezing temperatures in Arctic peat soil. *Microbiome* **9**, 208
27 (2021).
- 28 47. Wells, L. E. & Deming, J. W. Modelled and measured dynamics of viruses in Arctic winter sea-ice brines.
29 *Environ. Microbiol.* **8**, 1115–1121 (2006).
- 30 48. Collins, R. E. & Deming, J. W. Abundant dissolved genetic material in Arctic sea ice Part II: Viral dynamics
31 during autumn freeze-up. *Polar Biol.* **34**, 1831–1841 (2011).
- 32 49. Danovaro, R., Manini, E. & Dell’Anno, A. Higher abundance of bacteria than of viruses in deep Mediterranean
33 sediments. *Appl. Environ. Microbiol.* **68**, 1468–1472 (2002).
- 34 50. Segata, N. Gut microbiome: westernization and the disappearance of intestinal diversity. *Curr. Biol.* **25**, R611–
35 R613 (2015).
- 36 51. Rocha, E. P. C. & Bikard, D. Microbial defenses against mobile genetic elements and viruses: Who defends
37 whom from what? *PLOS Biol.* **20**, e3001514 (2022).
- 38 52. Botelho, J. Defense systems are pervasive across chromosomally integrated mobile genetic elements and are
39 inversely correlated to virulence and antimicrobial resistance. *Nucleic Acids Res.* **51**, 4385–4397 (2023).
- 40 53. Buongiorno Pereira, M. *et al.* A comprehensive survey of integron-associated genes present in metagenomes.
41 *BMC Genomics* **21**, 495 (2020).
- 42 54. Redondo-Salvo, S. *et al.* Pathways for horizontal gene transfer in bacteria revealed by a global map of their
43 plasmids. *Nat. Commun.* **11**, 3602 (2020).
- 44 55. Saravanan, M., Vasu, K. & Nagaraja, V. Evolution of sequence specificity in a restriction endonuclease by a
45 point mutation. *Proc. Natl. Acad. Sci.* **105**, 10344–10347 (2008).
- 46 56. Takeuchi, N., Wolf, Y. I., Makarova, K. S. & Koonin, E. V. Nature and intensity of selection pressure on CRISPR-
47 associated genes. *J. Bacteriol.* **194**, 1216–1225 (2012).
- 48 57. Giovannoni, S. J., Cameron Thrash, J. & Temperton, B. Implications of streamlining theory for microbial ecology.
49 *ISME J.* **8**, 1553–1565 (2014).
- 50 58. Pleška, M. & Guet, C. C. Effects of mutations in phage restriction sites during escape from restriction–
51 modification. *Biol. Lett.* **13**, 20170646 (2017).
- 52 59. Atack, J. M., Guo, C., Yang, L., Zhou, Y. & Jennings, M. P. DNA sequence repeats identify numerous Type I
53 restriction-modification systems that are potential epigenetic regulators controlling phase-variable regulons;
54 phasevarions. *FASEB J.* **34**, 1038–1051 (2020).
- 55 60. Watson, B. N. J., Steens, J. A., Staals, R. H. J., Westra, E. R. & Van Houte, S. Coevolution between bacterial
56 CRISPR-Cas systems and their bacteriophages. *Cell Host Microbe* **29**, 715–725 (2021).
- 57 61. Mohawk, K. L., Poly, F., Sahl, J. W., Rasko, D. A. & Guerry, P. High frequency, spontaneous *motA* mutations
58 in *Campylobacter jejuni* strain 81-176. *PLoS ONE* **9**, e88043 (2014).
- 59 62. Diez, A. A., Farewell, A., Nannmark, U. & Nyström, T. A mutation in the *ftsK* gene of *Escherichia coli* affects
60 cell-cell separation, stationary-phase survival, stress adaptation, and expression of the gene encoding the
61 stress protein UspA. *J. Bacteriol.* **179**, 5878–5883 (1997).
- 62 63. Srikant, S., Guegler, C. K. & Laub, M. T. The evolution of a counter-defense mechanism in a virus constrains
63 its host range. *eLife* **11**, e79549 (2022).
- 64 64. Ho, P. *et al.* Bacteriophage antidefense genes that neutralize TIR and STING immune responses. *Cell Rep.* **42**,
65 112305 (2023).
- 66 65. Nawrocki, E. P. & Eddy, S. R. Infernal 1.1: 100-fold faster RNA homology searches. *Bioinformatics* **29**, 2933–
67 2935 (2013).

- 1 66. Lagesen, K. *et al.* RNAMmer: consistent and rapid annotation of ribosomal RNA genes. *Nucleic Acids Res.* **35**,
2 3100–3108 (2007).
- 3 67. Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068–2069 (2014).
- 4 68. Cantalapiedra, C. P., Hernández-Plaza, A., Letunic, I., Bork, P. & Huerta-Cepas, J. eggNOG-mapper v2:
5 Functional annotation, orthology assignments, and domain prediction at the metagenomic scale. *Mol. Biol.*
6 *Evol.* **38**, 5825–5829 (2021).
- 7 69. Eddy, S. R. A probabilistic model of local sequence alignment that simplifies statistical significance estimation.
8 *PLoS Comput. Biol.* **4**, e1000069 (2008).
- 9 70. Pellow, D., Mizrahi, I. & Shamir, R. PlasClass improves plasmid sequence classification. *PLOS Comput. Biol.*
10 **16**, e1007781 (2020).
- 11 71. Krawczyk, P. S., Lipinski, L. & Dziembowski, A. PlasFlow: predicting plasmid sequences in metagenomic data
12 using genome signatures. *Nucleic Acids Res.* **46**, e35–e35 (2018).
- 13 72. Néron, B. *et al.* IntegronFinder 2.0: Identification and analysis of integrons across bacteria, with a focus on
14 antibiotic resistance in *Klebsiella*. *Microorganisms* **10**, 700 (2022).
- 15 73. Guo, J. *et al.* VirSorter2: a multi-classifier, expert-guided approach to detect diverse DNA and RNA viruses.
16 *Microbiome* **9**, 37 (2021).
- 17 74. Rousset, F. *et al.* Phages and their satellites encode hotspots of antiviral systems. *Cell Host Microbe* **30**, 740-
18 753.e5 (2022).
- 19 75. Liu, M. *et al.* ICEberg 2.0: an updated database of bacterial integrative and conjugative elements. *Nucleic Acids*
20 *Res.* **47**, D660–D665 (2019).
- 21 76. Katoh, K. MAFFT: a novel method for rapid multiple sequence alignment based on fast Fourier transform.
22 *Nucleic Acids Res.* **30**, 3059–3066 (2002).
- 23 77. Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids*
24 *Res.* **32**, 1792–1797 (2004).
- 25 78. Criscuolo, A. & Gribaldo, S. BMGE (Block Mapping and Gathering with Entropy): a new software for selection
26 of phylogenetic informative regions from multiple sequence alignments. *BMC Evol. Biol.* **10**, 210 (2010).
- 27 79. Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies.
28 *Bioinformatics* **30**, 1312–1313 (2014).
- 29 80. Minh, B. Q. *et al.* IQ-TREE 2: New models and efficient methods for phylogenetic inference in the genomic era.
30 *Mol. Biol. Evol.* **37**, 1530–1534 (2020).
- 31 81. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. (2013)
32 doi:10.48550/ARXIV.1303.3997.
- 33 82. Garrison, E. & Marth, G. Haplotype-based variant detection from short-read sequencing. (2012)
34 doi:10.48550/ARXIV.1207.3907.
- 35 83. Nelson, C. W., Moncla, L. H. & Hughes, A. L. SNPGenie: estimating evolutionary parameters to detect natural
36 selection using pooled next-generation sequencing data. *Bioinformatics* **31**, 3709–3711 (2015).
- 37 84. Miele, V., Penel, S. & Duret, L. Ultra-fast sequence clustering from similarity networks with SiLiX. *BMC*
38 *Bioinformatics* **12**, 116 (2011).
- 39 85. Suyama, M., Torrents, D. & Bork, P. PAL2NAL: robust conversion of protein sequence alignments into the
40 corresponding codon alignments. *Nucleic Acids Res.* **34**, W609–W612 (2006).
- 41 86. Wang, D., Zhang, Y., Zhang, Z., Zhu, J. & Yu, J. KaKs_Calculator 2.0: A toolkit incorporating gamma-series
42 methods and sliding window strategies. *Genomics Proteomics Bioinformatics* **8**, 77–80 (2010).
- 43 87. Yang, Z. & Nielsen, R. Estimating synonymous and nonsynonymous substitution rates under realistic
44 evolutionary models. *Mol. Biol. Evol.* **17**, 32–43 (2000).
- 45 88. Simple methods for estimating the numbers of synonymous and nonsynonymous nucleotide substitutions. *Mol.*
46 *Biol. Evol.* (1986) doi:10.1093/oxfordjournals.molbev.a040410.
- 47

1 **Acknowledgements**

2 This work was supported by the Genoscope, the Commissariat à l'Énergie Atomique et aux
3 Énergies Alternatives (CEA), France Génomique (ANR-10-INBS-09-08), and by the
4 Interdisciplinary Center MICROBES of the University Paris-Saclay, as part of the France 2030
5 program ANR-11-IDEX-0003. L.P. is supported by a European Molecular Biology
6 Organization Postdoctoral Fellowship (EMBO ALTF 100-2023). We thank Eduardo P. C.
7 Rocha (Institut Pasteur, Paris) for critical reading of the manuscript, and Hadrien Guichard
8 (CEA, Genoscope) for initial efforts in defensome analyses. We acknowledge Dr. João
9 Botelho (CBGP, Universidad Politécnica de Madrid) and the two additional anonymous
10 reviewers for their careful reading and insightful comments and suggestions.

11

12 **Author Contributions**

13 P.H.O supervised the project. A.B. and P.H.O. designed the computational methods. A.B. and
14 A.L. performed most of the computational analyses and developed most of the scripts that
15 support the analyses. A.B., A.L., N.W., J.P., T.O.D., L.P., P.W. and P.H.O. analysed the data.
16 A.B. and P.H.O. wrote the manuscript with additional information inputs from other co-authors.

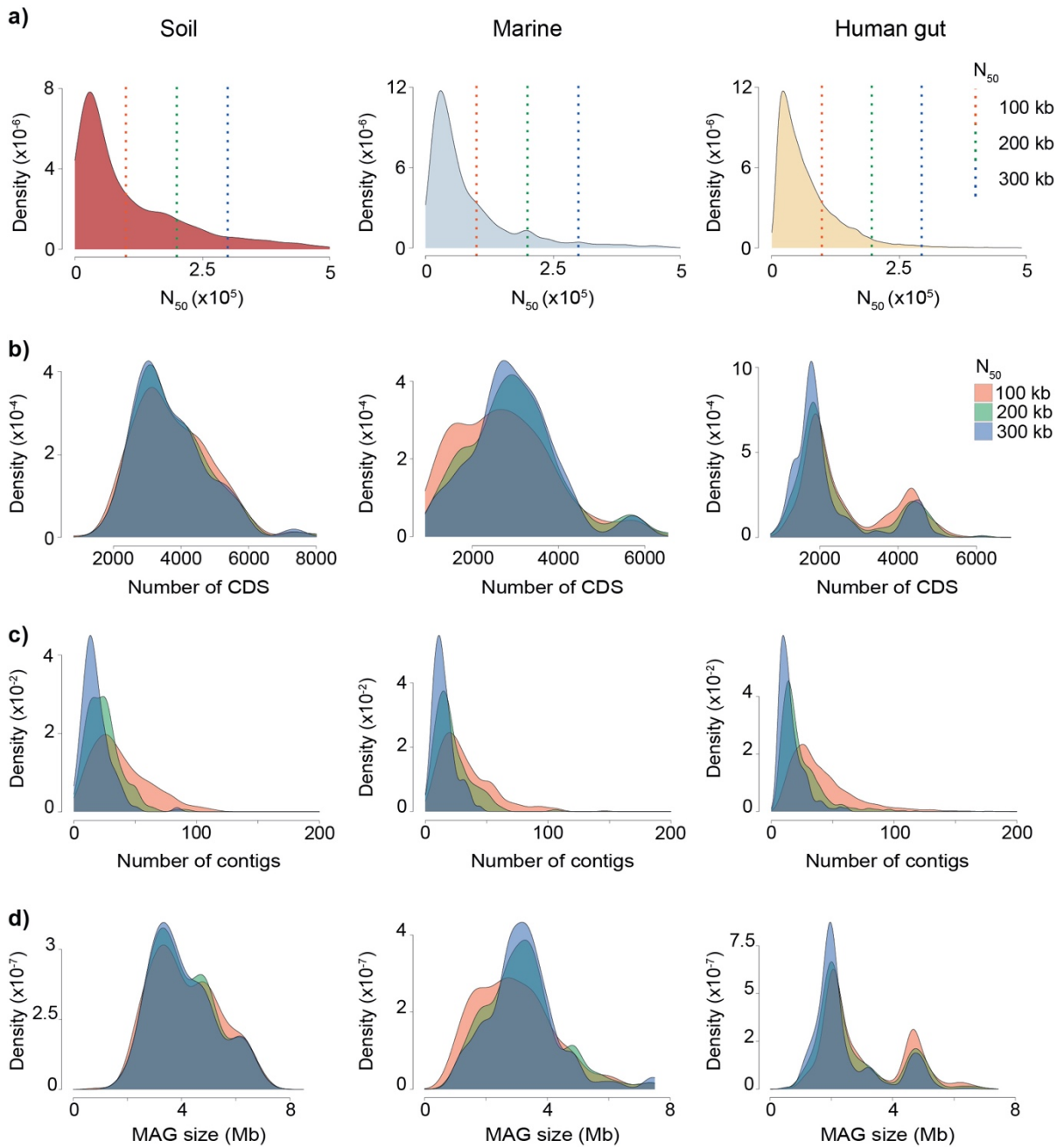
17

18 **Competing interests**

19 The authors declare no competing interests.

20

1 Supplementary Figures

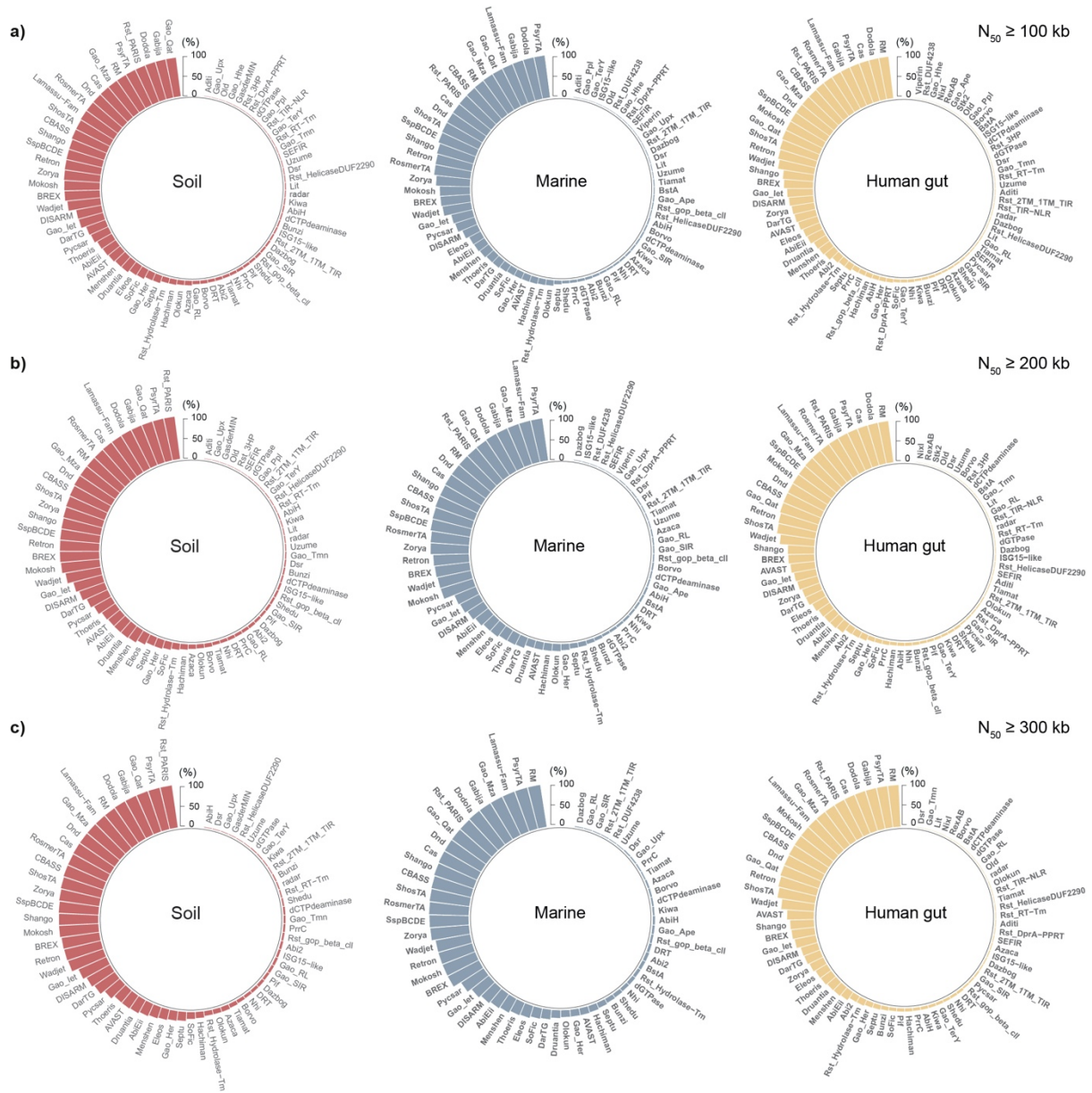


2

3 **Supplementary Fig. 1.** Density distributions of key MAG parameters according to environment: (a) N_{50} ; (b) number

4 of CDS; (c) number of contigs; and (d) size.

1

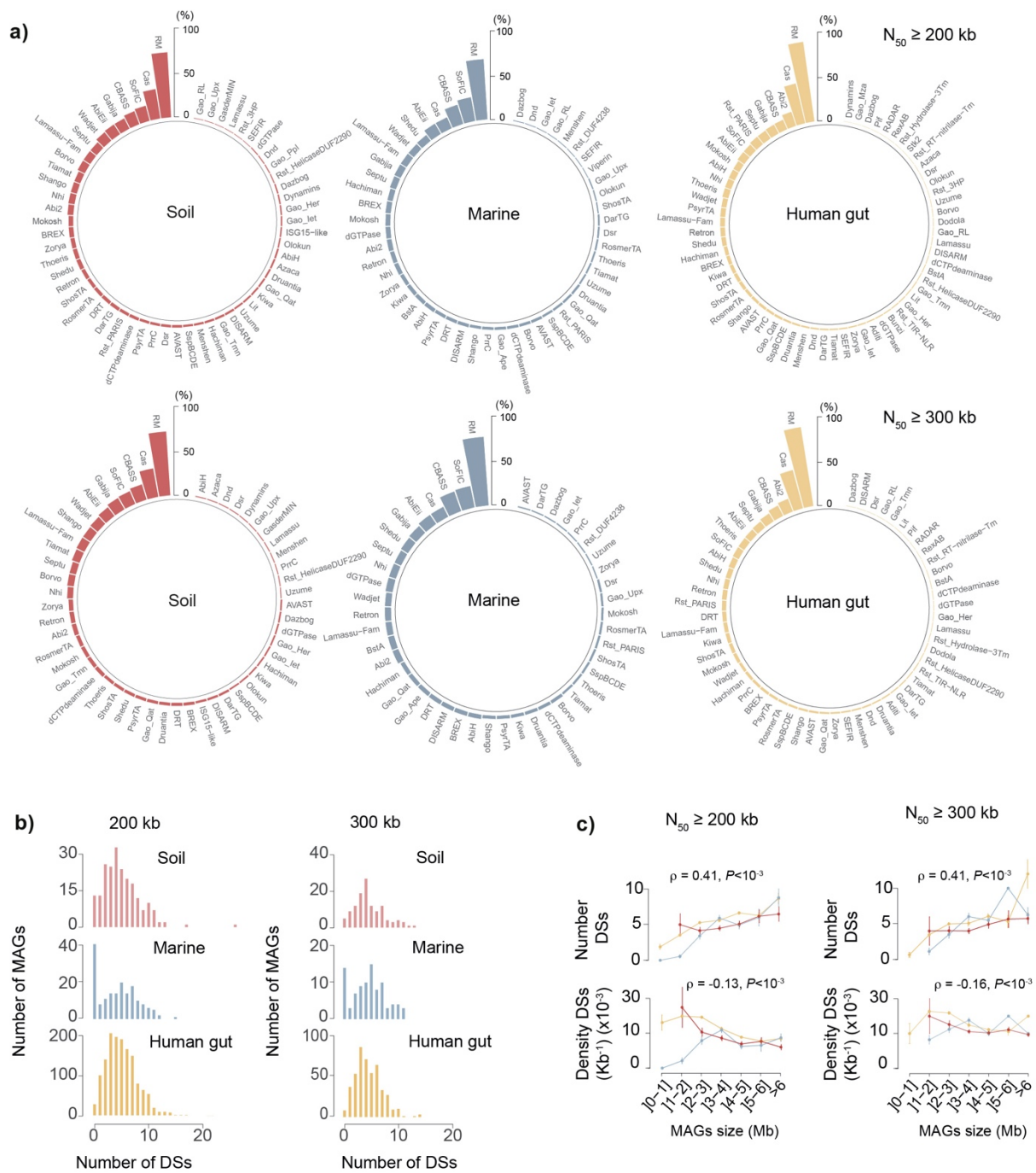


2

3 **Supplementary Fig. 2.** Percentage of soil, marine and human gut MAGs harboring each family of defense genes.

4 Shown are assemblies with values of N₅₀ ≥ 100, 200, and 300 kb.

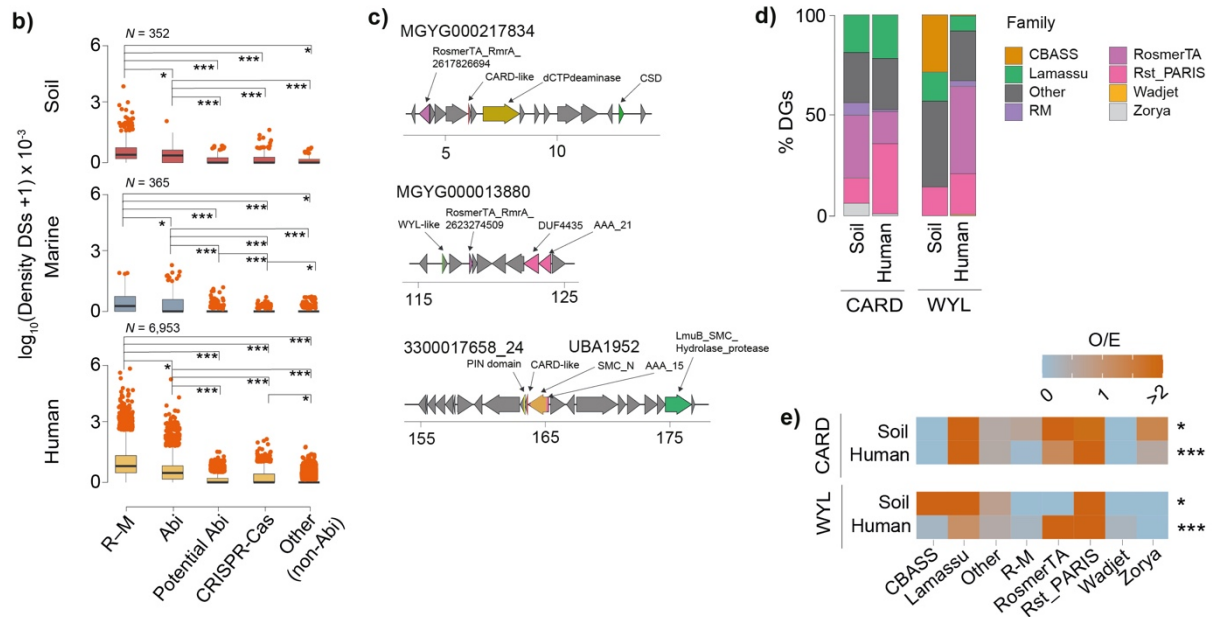
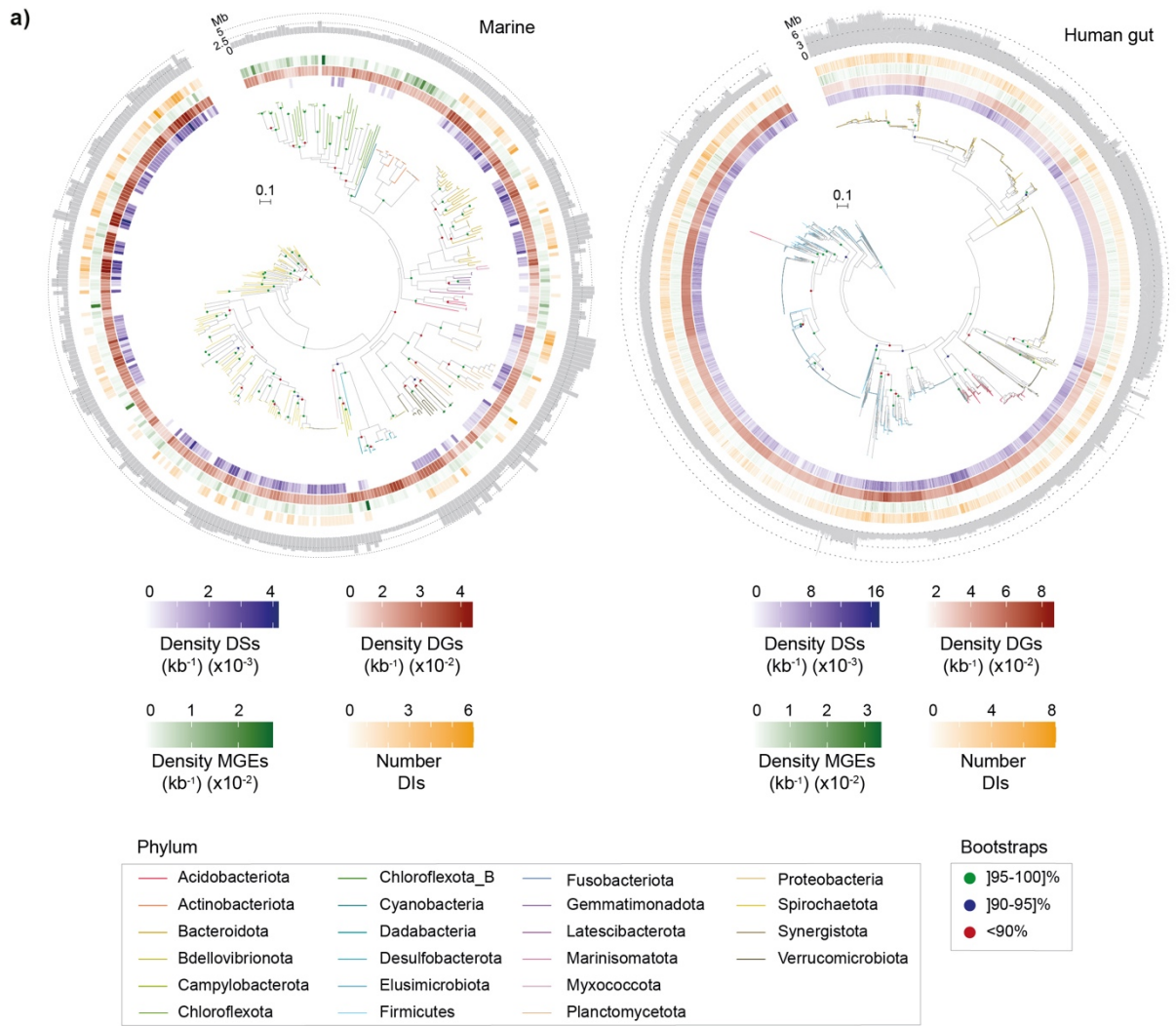
1



2

3 **Supplementary Fig. 3.** Distribution of defense systems in MAGs. (a) Percentage of soil, marine and human gut
4 MAGs harboring each family of defense systems. Shown are assemblies with values of $N_{50} \geq 200$ and 300 kb. (b)
5 Distribution of number of defense systems (DSs, per MAG) across environments for assemblies with values of N_{50}
6 ≥ 200 and 300 kb. (c) Variation of number and density (per kb) of defense systems (DSs) with MAG size (Mb) for
7 assemblies with values of $N_{50} \geq 200$ and 300 kb and for each environment. Error bars represent standard
8 deviations.

1

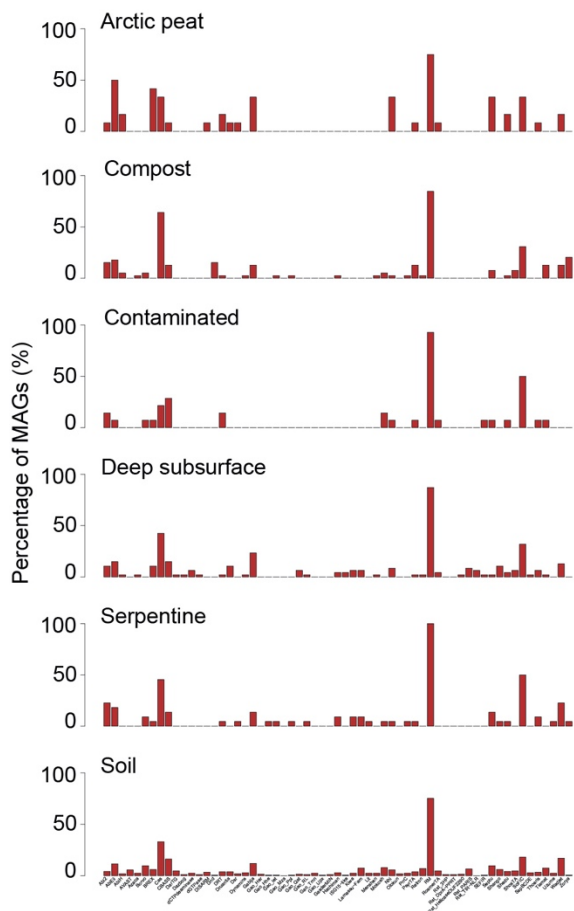


2

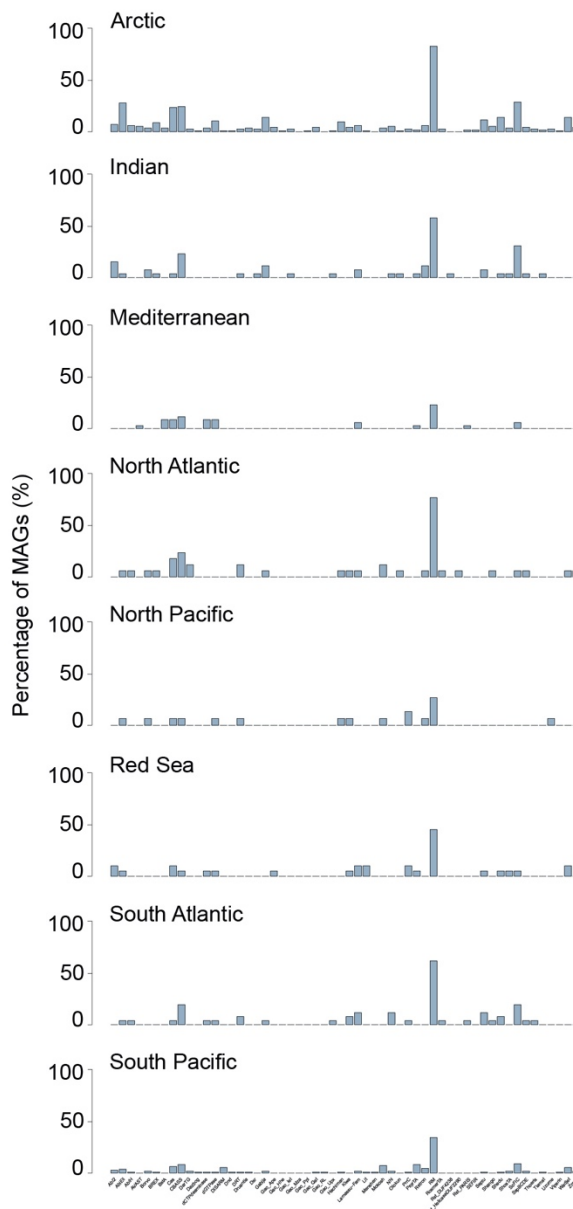
1 **Supplementary Fig. 4.** Phylogenetic distribution of the defensome and association with types of defense
2 mechanism and CARD / WYL domains. (a) Phylogenetic representation of 386 marine and 6,369 human gut MAGs,
3 their corresponding phyla, density (per kb) of defense systems (DSs, purple), defense genes (DGs, red), MGEs
4 (green), and number of defense islands (DIs, yellow). Distribution of MAG sizes (Mb) are shown as outer layer
5 barplots. (b) Defense system (DS) density (per MAG per kb) split per biome and underlying defense mechanism
6 (R–M, Abi, potential Abi, CRISPR-Cas, and other (non-Abi)). Error bars represent standard deviations and Mann–
7 Whitney–Wilcoxon test P values are indicated. $*P < 0.05$; $**P < 10^{-2}$; $***P < 10^{-3}$. (c) Representative instances of
8 colocalized (less than 10 genes apart) WYL and CARD domain-containing genes with defense genes in the human
9 MAGs MGYG000217834 and MGYG000013880. Also shown is a chromosomal region of the Bacteroidetes
10 bacterium UBA1952 sharing similarity with a recently described *Pedobacter rhizosphaerae* CARD-encoding
11 defense system³³. (d) Stacked barplots showing the relative abundance of families of defense genes that colocalize
12 with WYL- and CARD-like domains. (e) Heatmap of observed / expected (O / E) ratios of colocalization between
13 genes belonging to distinct defense families and WYL and CARD-like domains. Expected values were obtained by
14 multiplying the total number of genes pertaining to a given defense family by the fraction of defense genes of that
15 family colocalized with a given domain. P values correspond to the χ^2 -test. To avoid performing analyses with weak
16 statistical power, we omitted the marine MAG dataset due to their low WYL and CARD domain abundance
17 (detected in as much 0.75% of the dataset).

1

a) SOIL

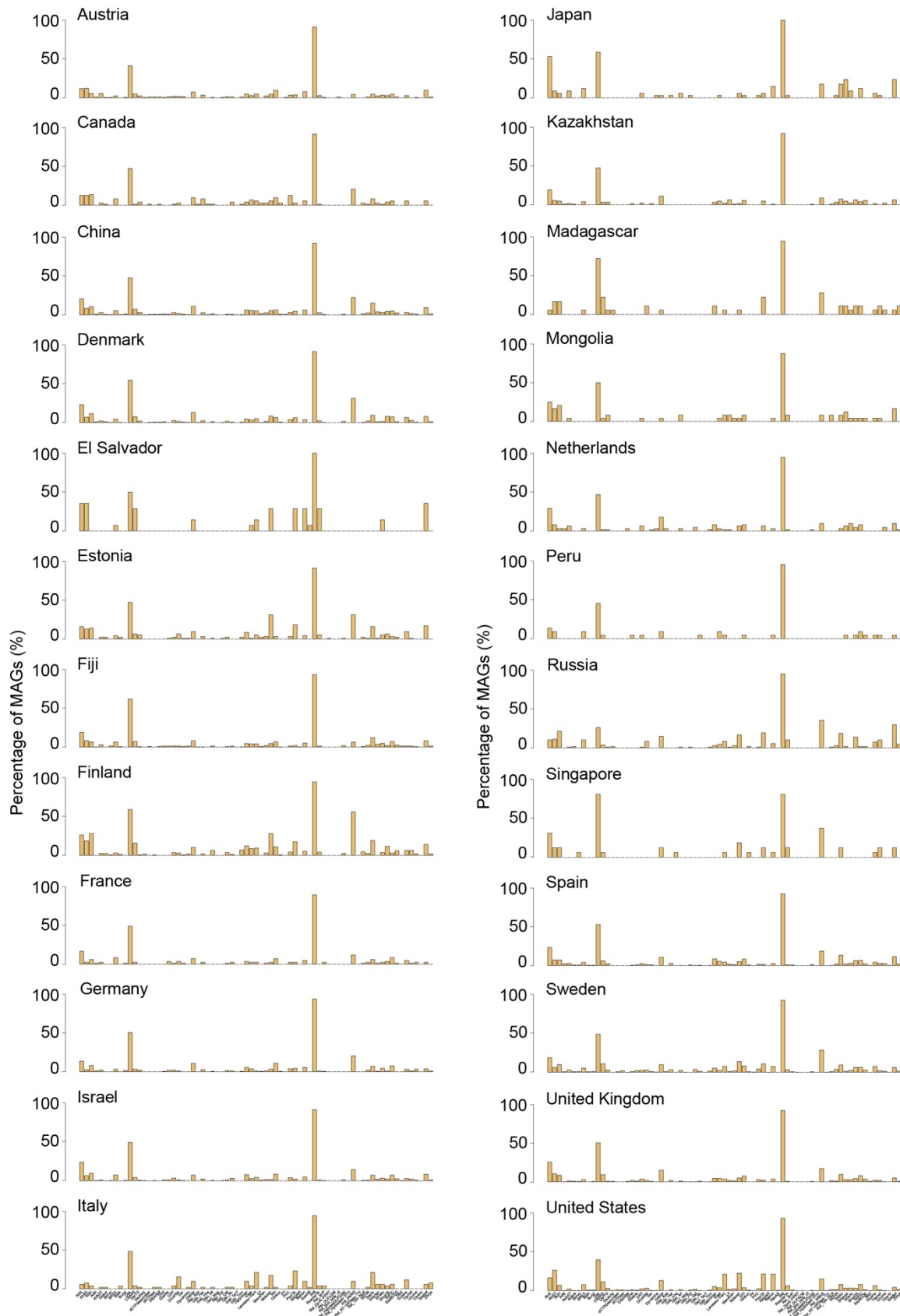


b) MARINE

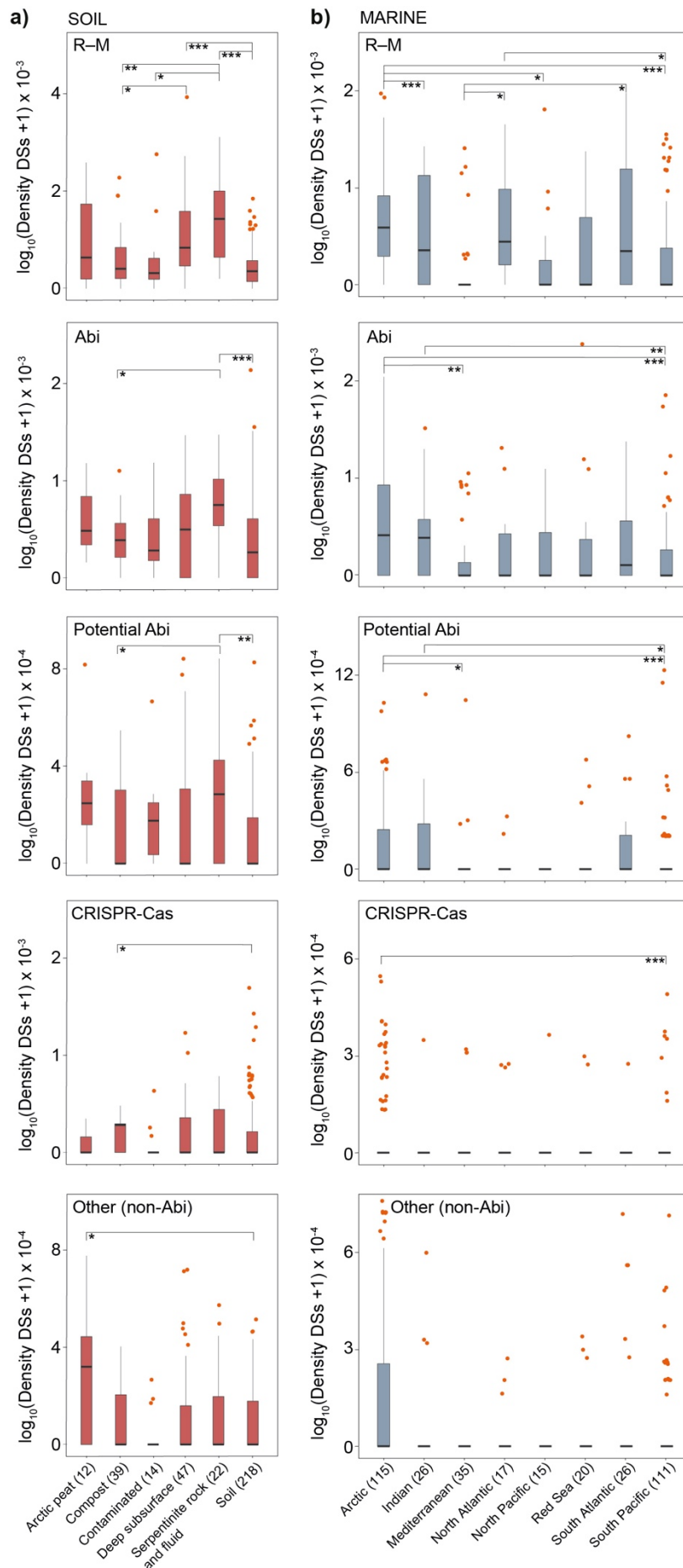


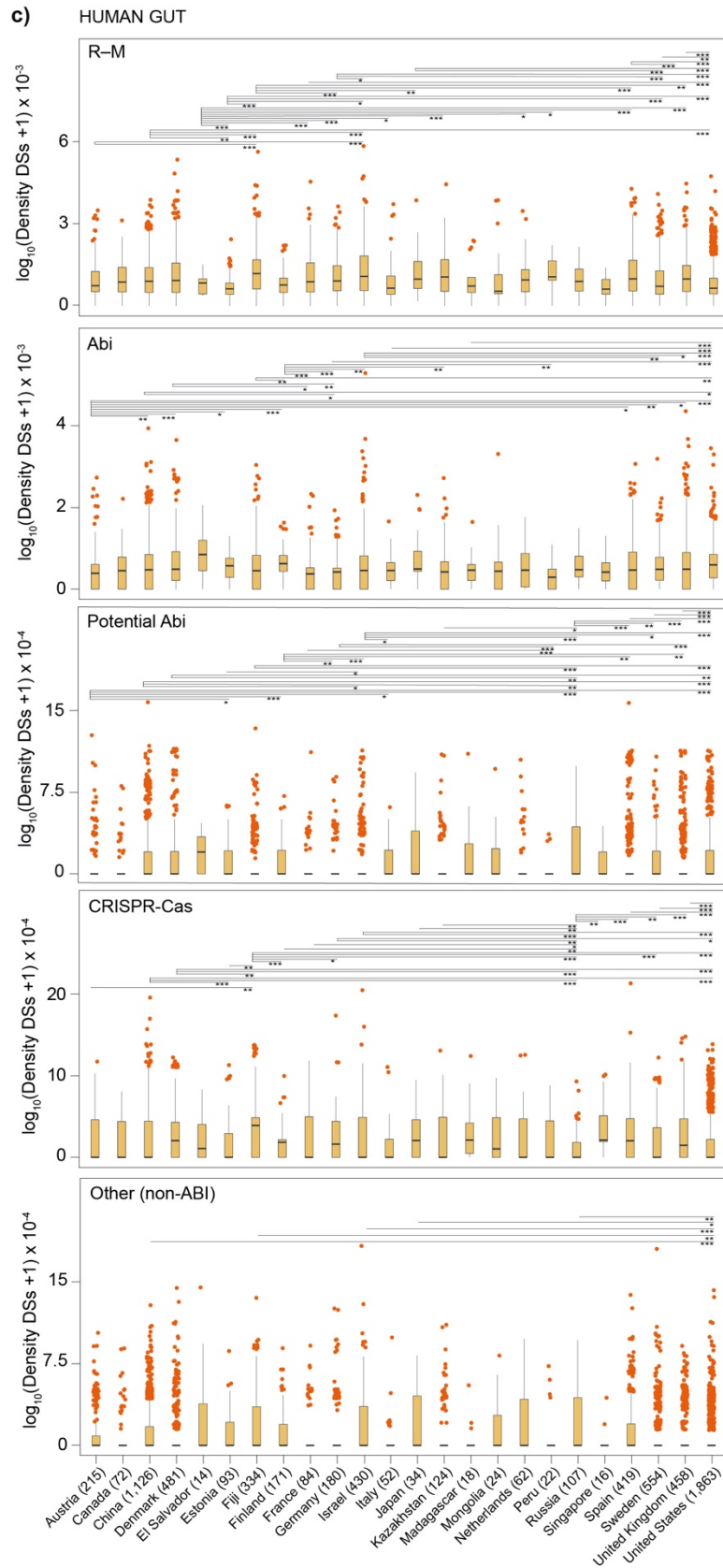
2

c) HUMAN GUT



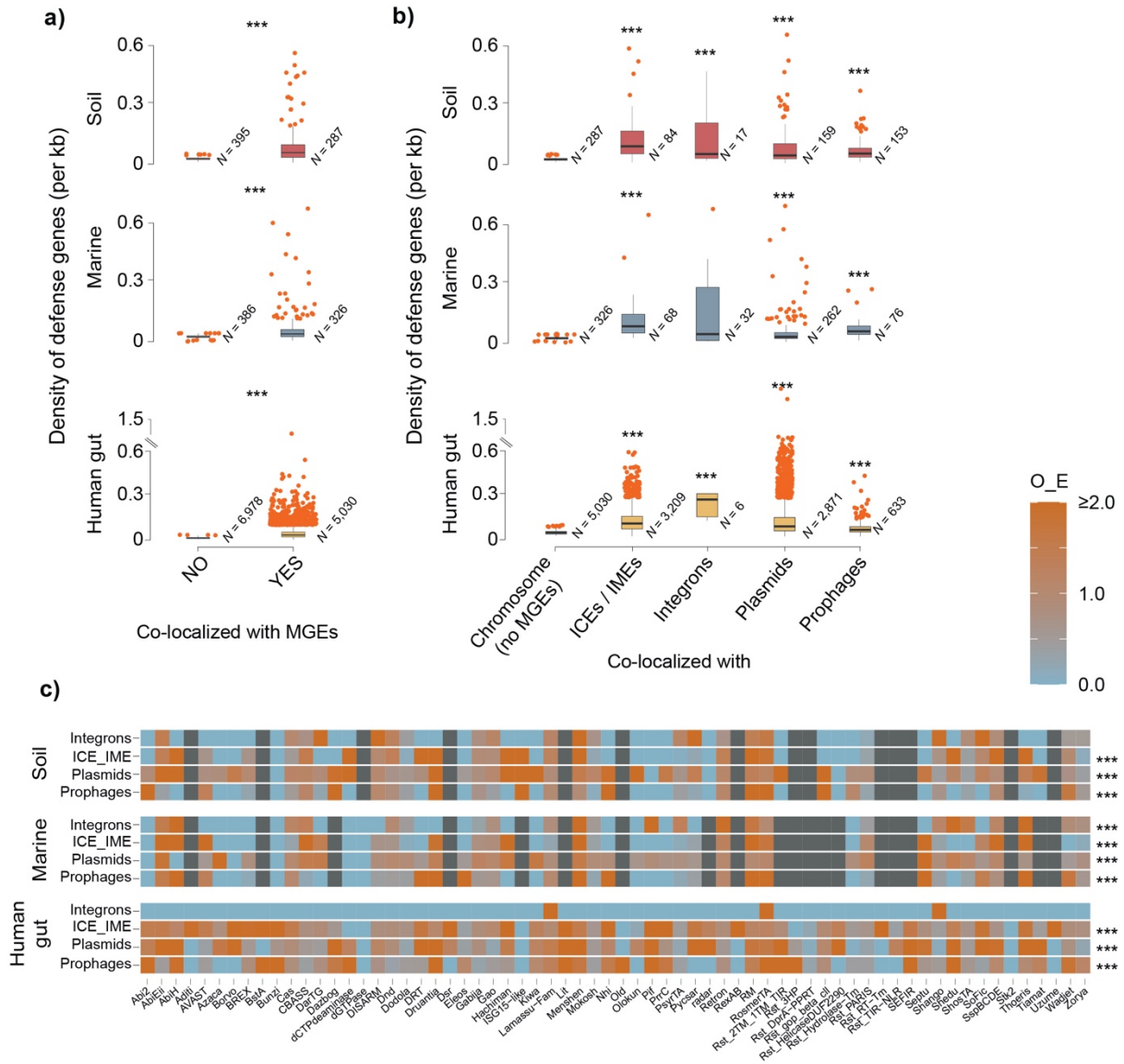
Supplementary Fig. 5. Percentage of (a) soil, (b) marine and (c) human gut MAGs harboring each family of defense system across different ecological and geographical backgrounds.



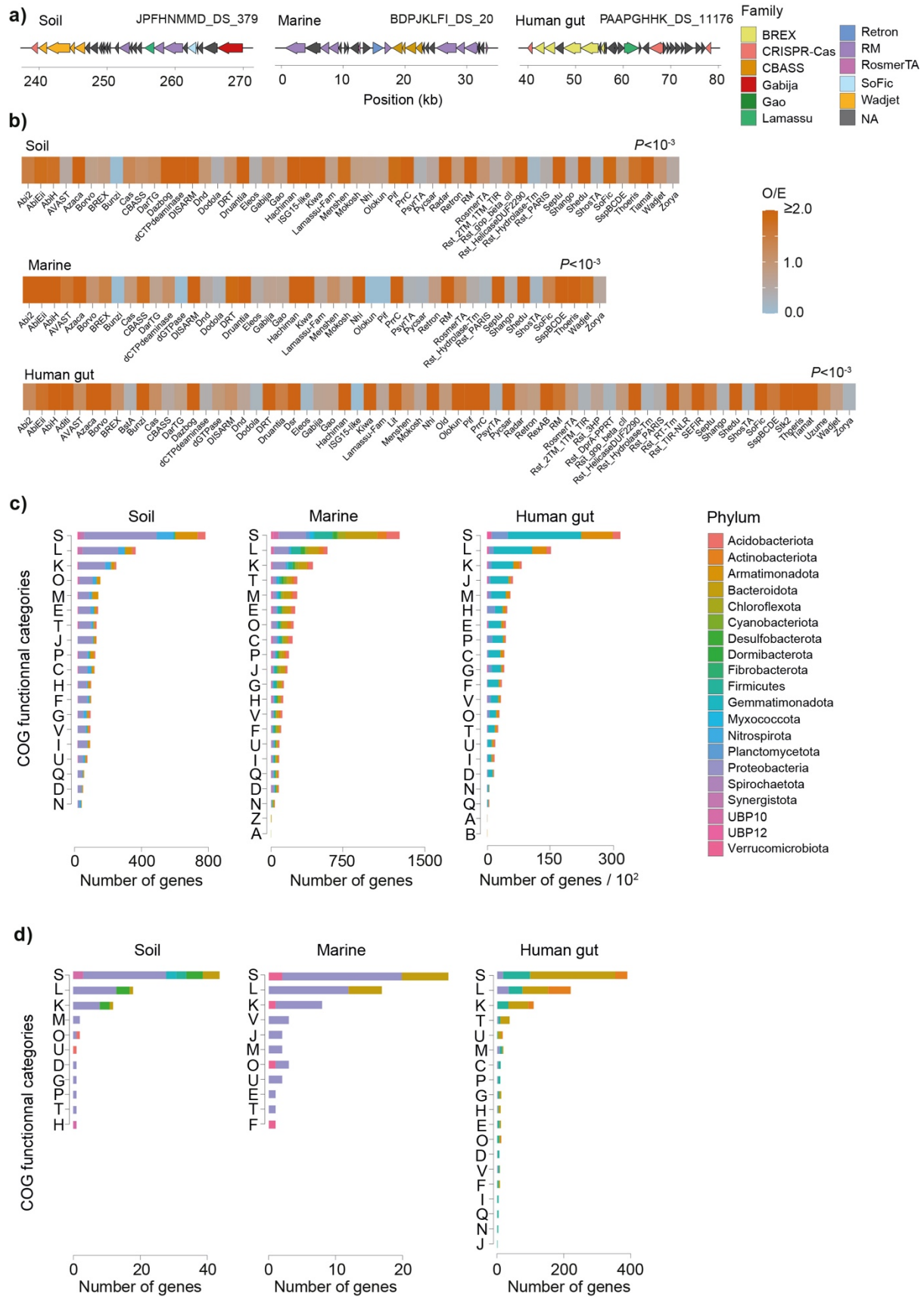


1

2 **Supplementary Fig. 6.** Defense system (DS) density (per MAG per kb) split according to defense mechanism (R–
 3 M, Abi, potential Abi, CRISPR-Cas, and other (non-ABI)), and different ecological (soil, marine) (a, b) and
 4 geographical (human gut) (c) contexts. Error bars represent standard deviations and Mann–Whitney–Wilcoxon test
 5 *P* values are indicated. **P* < 0.05, ***P* < 10⁻², ****P* < 10⁻³.



1
 2 **Supplementary Fig. 7.** The genetic mobility of the defensome. Reproduction of Fig. 4 of the manuscript with the
 3 inclusion of integrations. * $P < 0.05$; ** $P < 10^{-2}$; *** $P < 10^{-3}$, χ^2 -test.



1

2

3

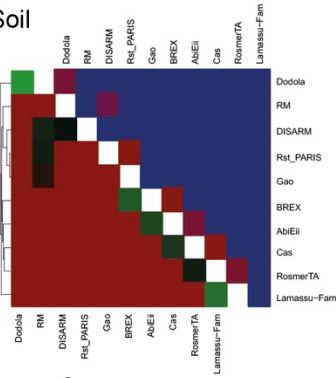
4

Supplementary Fig. 8. The MAG defense island repertoire. (a) Representative examples of defense islands for each environment, illustrating the high diversity of defense families present. (b) Heatmap of observed/expected (O/E) ratios of colocalization between genes belonging to distinct defense families and defense islands. Expected

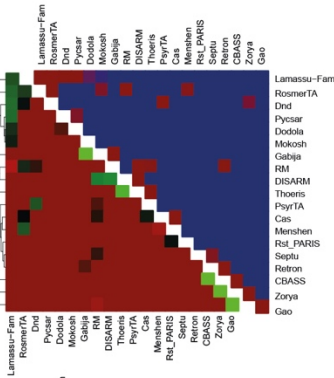
1 values were obtained by multiplying the total number of genes pertaining to a given defense family by the fraction
2 of defense genes of that family assigned to a defense island. *P* values correspond to the χ^2 -test. (c) COG functional
3 annotation assessed by EggNog-mapper of the ensemble of non-defensive genes in defense islands. (d) COG
4 functional annotation of the ensemble of non-defensive genes in defense islands when complete defense systems
5 were used as counting units for the classification of defense islands (see Methods).

a) Soil

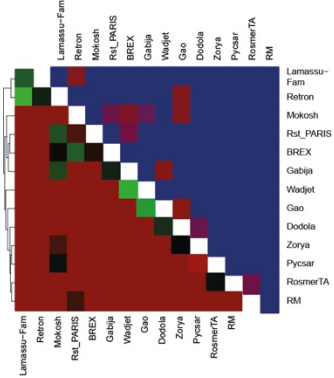
Arctic peat



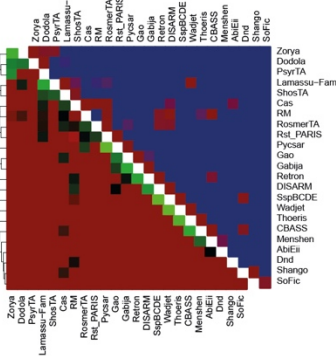
Compost



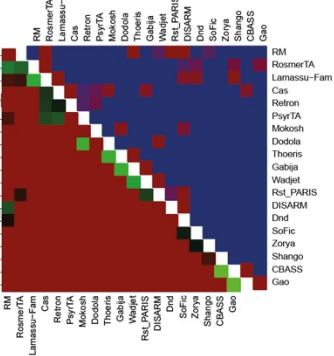
Contaminated



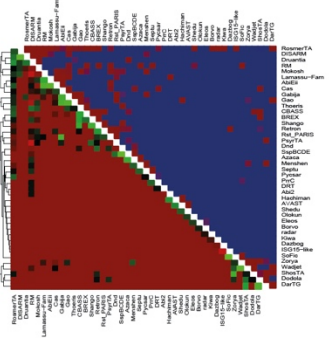
Deep subsurface



Serpentine

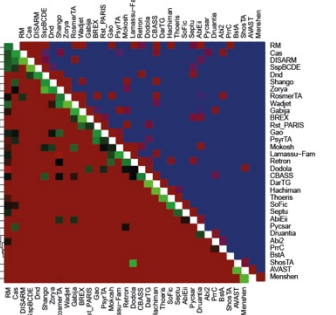


Soil

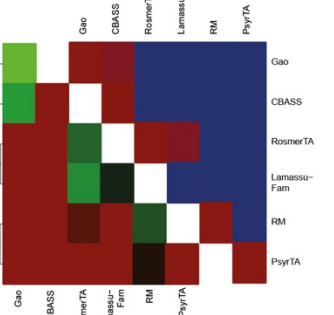


b) Marine

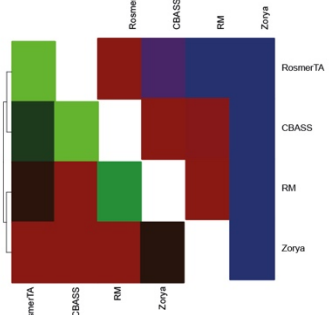
Arctic



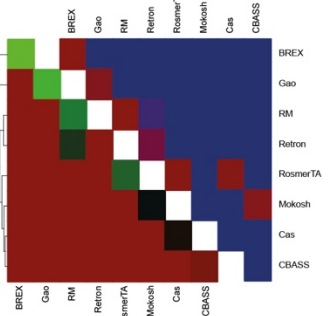
Indian



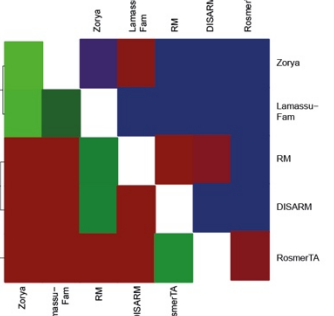
Mediterranean



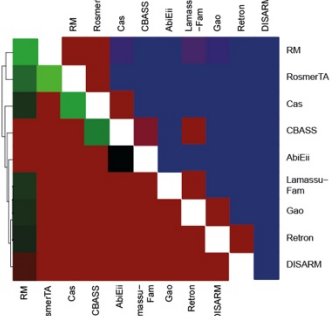
North Atlantic



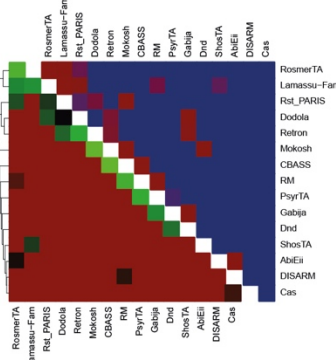
North Pacific



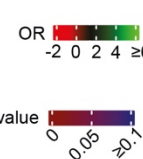
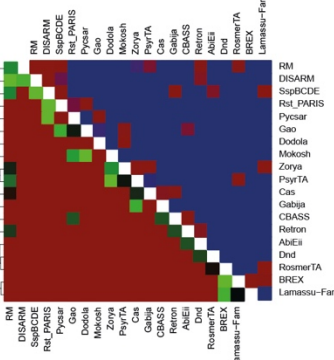
Red Sea



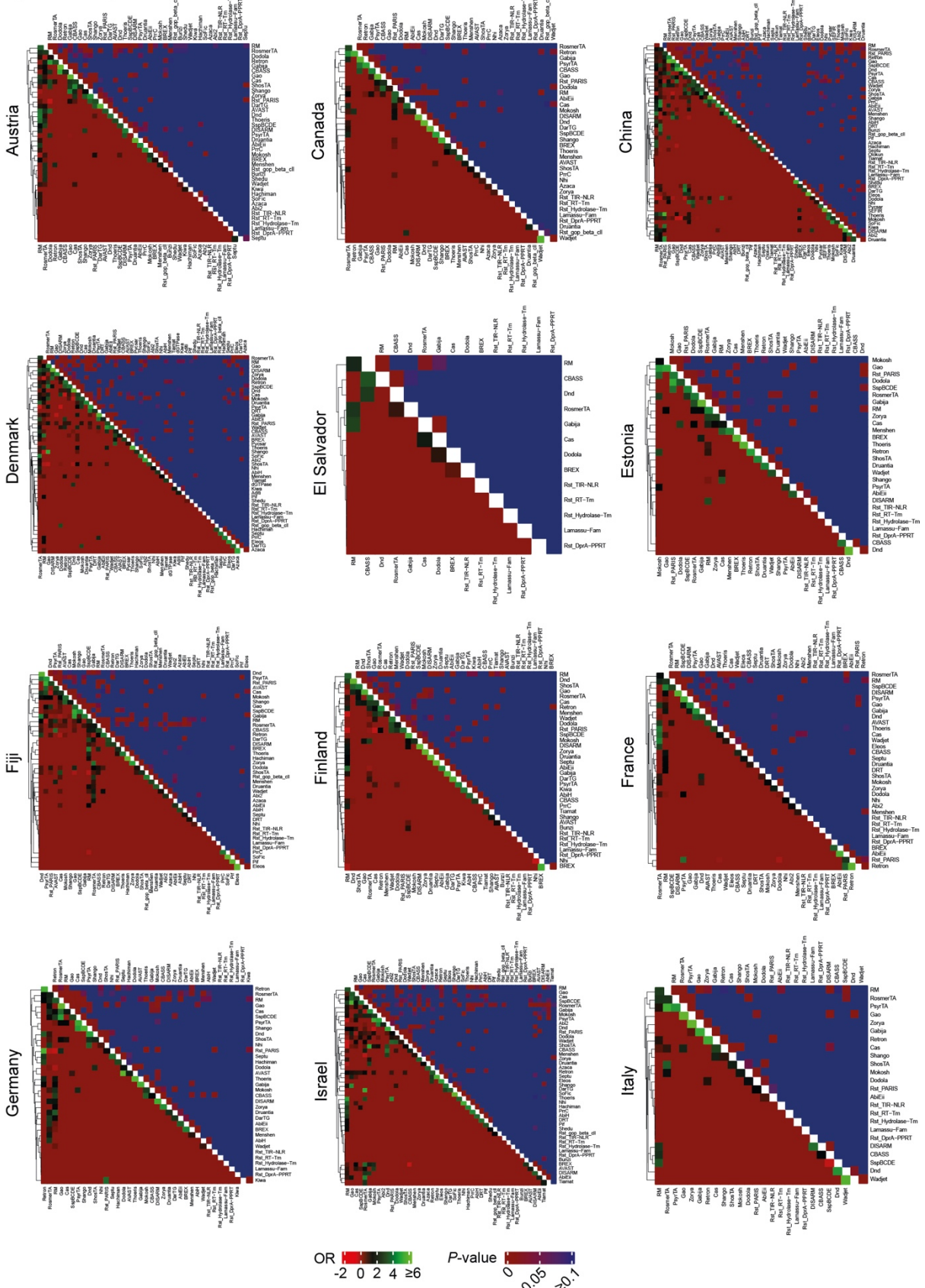
South Atlantic



South Pacific

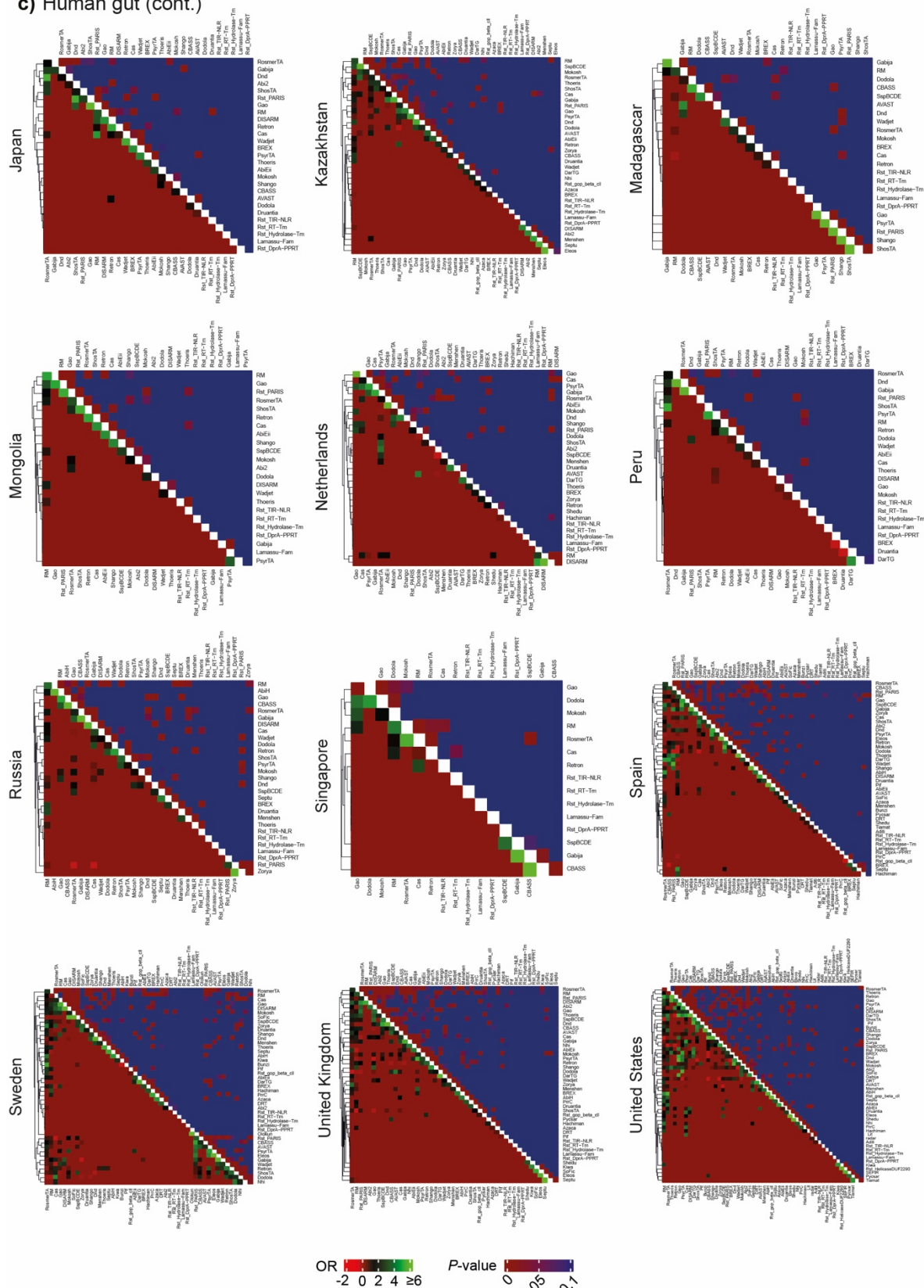


c) Human gut



1
2
3

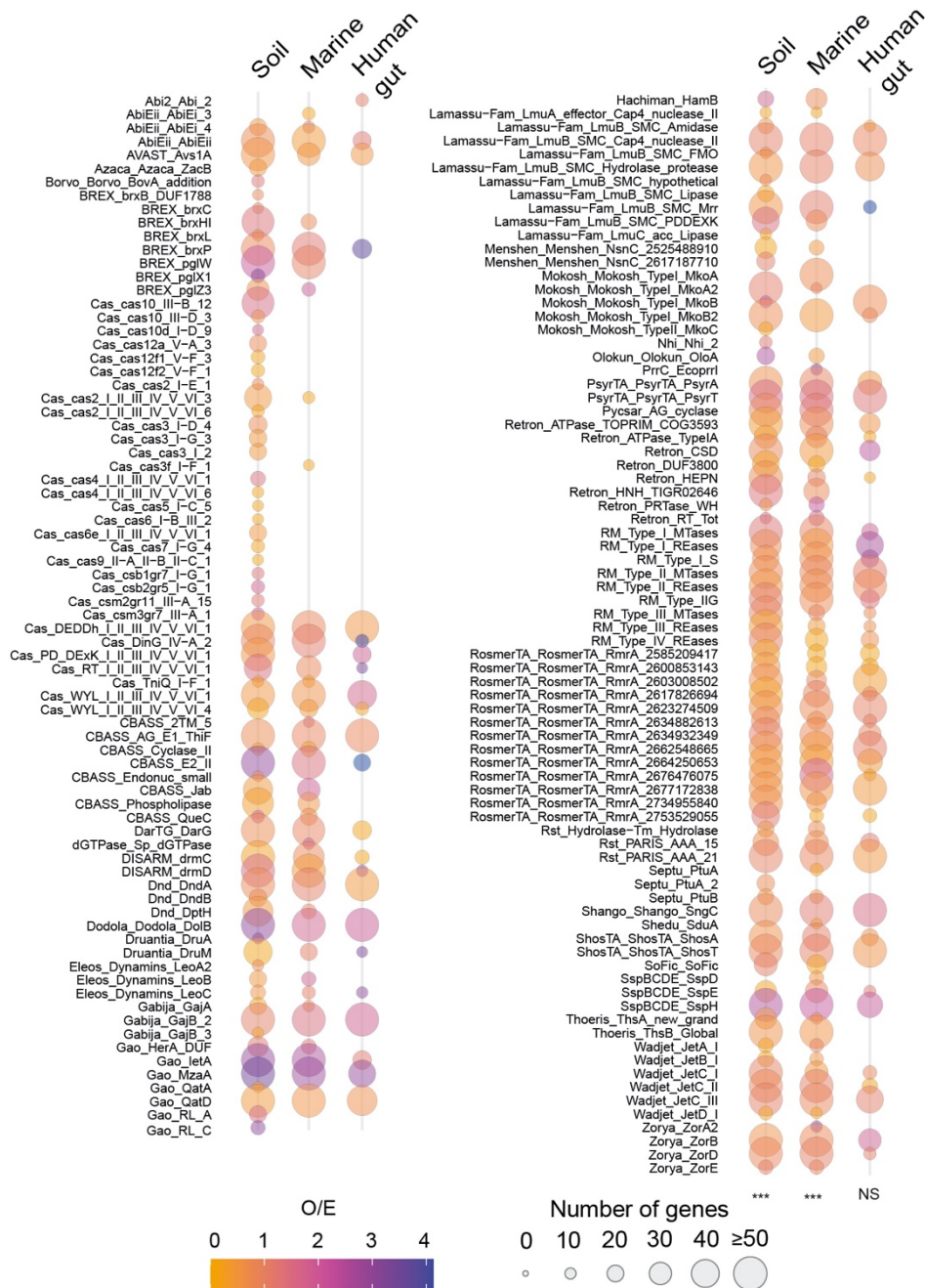
c) Human gut (cont.)



1
2
3
4
5

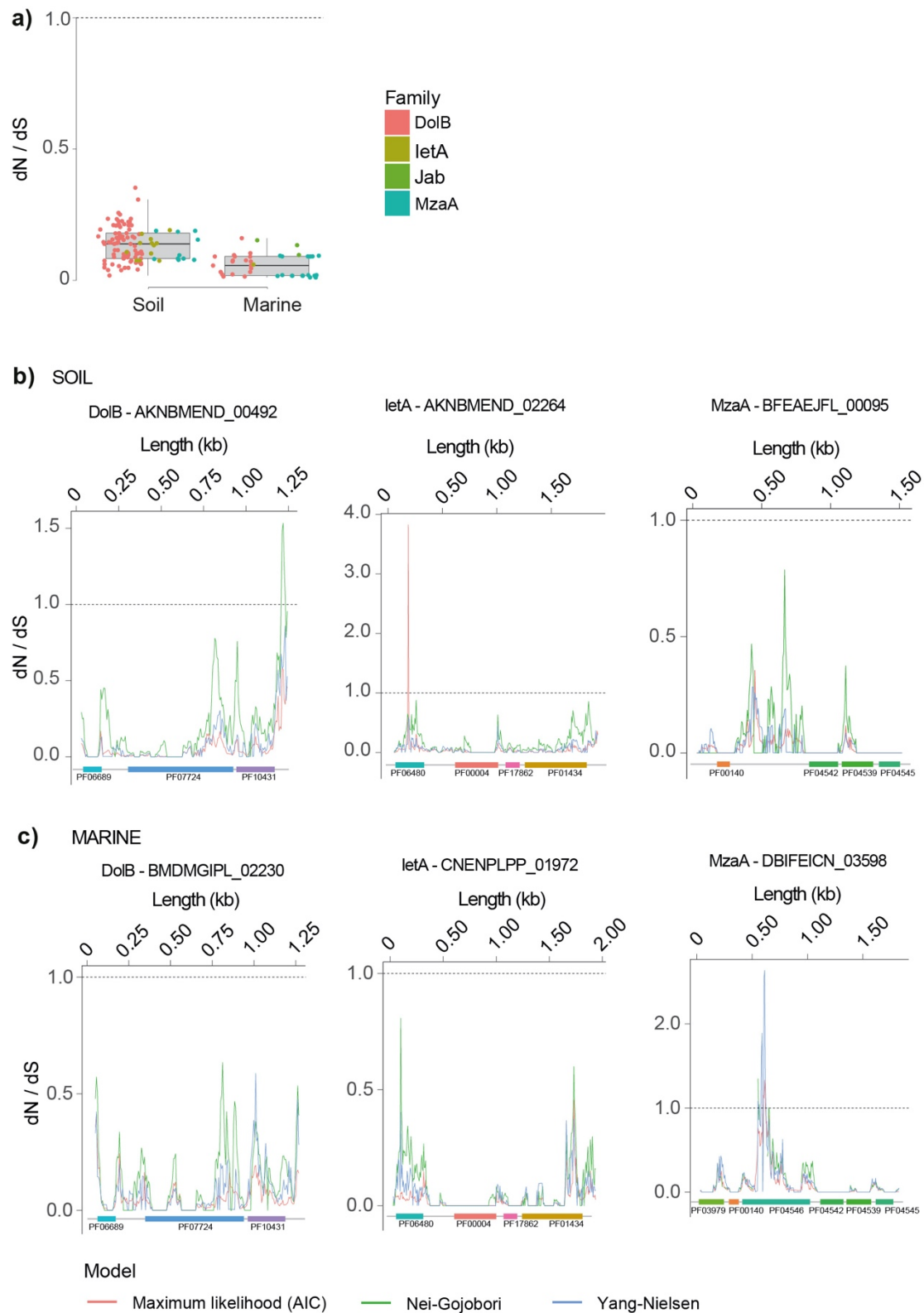
Supplementary Fig. 9. Defense families' odds ratio (OR) of colocalization in defense islands (bottom heatmaps) and associated Fisher's exact test *P* value (upper heatmaps) for different ecological (soil, marine) (a, b) and geographical (human gut) (c) contexts.

1



2

3 **Supplementary Fig. 10.** The genetic variability of the defensome. (a) 90 metagenomes (30 for each environment)
 4 having a broad representativity in terms of sampling sites (soil and marine) and countries (human gut), as well as
 5 in terms of presence of most defense families previously identified by DefenseFinder were selected. Shown in
 6 circles are the observed / expected (O / E) ratios of number of defense gene families harboring high-frequency (\geq
 7 25% of coverage at the variant position) SNPs + indels positions in their gene body (including 200 bp upstream the
 8 start codon). No thresholds on O / E ratio were introduced. Expected values were obtained by multiplying the total
 9 number of genes pertaining to a given defense family by the fraction of defense genes of that family harboring high-
 10 frequency alleles. Circle radius corresponds to the total number of defense genes analyzed per family. All defense
 11 families are represented irrespectively of their O / E ratio span.



1
2
3 **Supplementary Fig. 11.** Evolution of defense genes. (a) Variation in global dN/dS given by the Nei-Gojoberi (NG)
4 and Yang-Nielsen (NY) methods for a selection of defense genes shown to harbor a significantly higher frequency
5 of SNPs + Indels. (b) Across gene profiles of dN/dS given by the Maximum Likelihood (Akaike Information
6 Criterion), Nei-Gojoberi (NG), and Yang-Nielsen (NY) methods for a selection of three defense genes
7 simultaneously present in soil and marine environments. PFAM domains are shown as colored rectangles.