# Unifying the known and unknown microbial coding sequence space

3

4 Chiara Vanni[1,2], Matthew S. Schechter[1,3], Silvia G. Acinas[4], Albert Barberán[5], Pier Luigi

5 Buttigieg[6], Emilio O. Casamayor[7], Tom O. Delmont[8], Carlos M. Duarte[9], A. Murat Eren[3,10],

6 Robert D. Finn[11], Renzo Kottmann[1], Alex Mitchell[11], Pablo Sanchez[4], Kimmo Siren[12], Martin

7 Steinegger[13,14], Frank Oliver Glöckner[15,16,2], Antonio Fernandez-Guerra[1,17] *

8

## Affiliations

10 1 Microbial Genomics and Bioinformatics Research Group, Max Planck Institute for Marine

11 Microbiology, Celsiusstraße 1, 28359, Bremen, Germany

12 2 Jacobs University Bremen, Campus Ring 1, 28759 Bremen, Germany

13 3 Department of Medicine, University of Chicago, Chicago, IL 60637, USA

14 4 Department of Marine Biology and Oceanography, Institut de Ciènces del Mar, CSIC,

15 Barcelona, Spain.

16 5 Department of Environmental Science, University of Arizona, Tucson, 85721 AZ, USA

17 6 Alfred Wegener Institute, Helmholtz Centre for Polar and Marine Research, Am Handelshafen

18 12, 27570 Bremerhaven, Germany

19 7 Center for Advanced Studies of Blanes CEAB-CSIC, Spanish Council for Research, Blanes,

20 Spain

21    8 Génomique Métabolique, Genoscope, Institut François Jacob, CEA, CNRS, Univ Evry,

22    Université Paris-Saclay, 91057 Evry, France

23    9 Red Sea Research Centre (RSRC) and Computational Bioscience Research Center (CBRC),

24    King Abdullah University of Science and Technology, Thuwal 23955, Saudi Arabia

25    10 Josephine Bay Paul Center, Marine Biological Laboratory, Woods Hole, MA 02543, USA

26    11 European Molecular Biology Laboratory, European Bioinformatics Institute (EMBL-EBI),

27    Wellcome Genome Campus, Hinxton, Cambridge CB10 1SD, UK

28    12 Section for Evolutionary Genomics, The GLOBE Institute, University of Copenhagen,

29    Copenhagen, Denmark

30    13 School of Biological Sciences, Seoul National University, Seoul, 08826, South Korea

31    14 Institute of Molecular Biology and Genetics, Seoul National University, Seoul, 08826, South

32    Korea

33    15 University of Bremen, MARUM, Leobener Str. 8, 28359 Bremen, Germany

34    Life Sciences and Chemistry, Campus Ring 1, 28759 Bremen, Germany

35    16 Computing Center, Helmholtz Center for Polar and Marine Research, Am Handelshafen 12,

36    27570 Bremerhaven, Germany

37    17 Lundbeck GeoGenetics Centre, The Globe Institute, University of Copenhagen, 1350

38    Copenhagen, Denmark

39

40    *Corresponding author: Antonio Fernandez-Guerra, antonio.fernandez-guerra@sund.ku.dk

41

2

# Abstract

42

43    Genes of unknown function are among the biggest challenges in molecular biology, especially in

44    microbial systems, where 40%-60% of the predicted genes are unknown. Despite previous

45    attempts, systematic approaches to include the unknown fraction into analytical workflows are

46    still lacking. Here, we propose a conceptual framework and a computational workflow that

47    bridge the known-unknown gap in genomes and metagenomes. We showcase our approach by

48    exploring 415,971,742 genes predicted from 1,749 metagenomes and 28,941 bacterial and

49    archaeal genomes. We quantify the extent of the unknown fraction, its diversity, and its

50    relevance across multiple biomes. Furthermore, we provide a collection of 283,874 lineage-

51    specific genes of unknown function for *Cand*. Patescibacteria, being a significant resource to

52    expand our understanding of their unusual biology. Finally, by identifying a target gene of

53    unknown function for antibiotic resistance, we demonstrate how we can enable the generation

54    of hypotheses that can be used to augment experimental data.

# Introduction

Thousands of isolate, single-cell, and metagenome-assembled genomes are guiding us towards a better understanding of life on Earth (Almeida et al., 2019; Cross et al., 2019; Delmont et al., 2020; Hug et al., 2016; Kopf et al., 2015; Pachiadaki et al., 2019; Pasolli et al., 2019; Sunagawa et al., 2015). At the same time, the ever-increasing number of genomes and metagenomes, unlocking uncharted regions of life's diversity, (Brown et al., 2015; Eloe-Fadrosh et al., 2016; Hug et al., 2016) are providing new perspectives on the evolution of life (Parks et al., 2018; Spang et al., 2015). However, our rapidly growing inventories of new genes have a glaring issue: between 40% and 60% cannot be assigned to a known function (Almeida et al., 2020; Bernard, Pathmanathan, Lannes, Lopez, & Bapteste, 2018; Carradec et al., 2018; Price et al., 2018). Current analytical approaches for genomic and metagenomic data (Chen et al., 2019; Franzosa et al., 2018; Huerta-Cepas et al., 2017; Mitchell et al., 2020; Quince, Walker, Simpson, Loman, & Segata, 2017) generally do not include this uncharacterized fraction in downstream analyses, constraining their results to conserved pathways and housekeeping functions (Quince et al., 2017). This inability to handle the unknown is an immense impediment to realizing the potential for discovery of microbiology and molecular biology at large (Bernard et al., 2018; Hanson, Pribat, Waller, & Crécy-Lagard, 2010). Predicting function from traditional single sequence similarity appears to have yielded all it can (Arnold, 1998, 2018; Brandenberg, Fasan, & Arnold, 2017), thus several groups have attempted to resolve gene function by other means. Such efforts include combining biochemistry and crystallography (Jaroszewski et al., 2009); using environmental co-occurrence (Buttigieg et al., 2013); by grouping those genes into evolutionarily related families (Bateman, Coggill, & Finn, 2010; Brum et al., 2016; Wyman, Avila-Herrera, Nayfach, & Pollard, 2018; Yooseph et al., 2007); using remote homologies (Bitard-Feildel & Callebaut, 2017; Lobb, Kurtz, Moreno-Hagelsieb, & Doxey, 2015); or more recently using deep learning approaches (Bileschi et al., 2019; Liu, 2017). In 2018, Price et al. (Price et
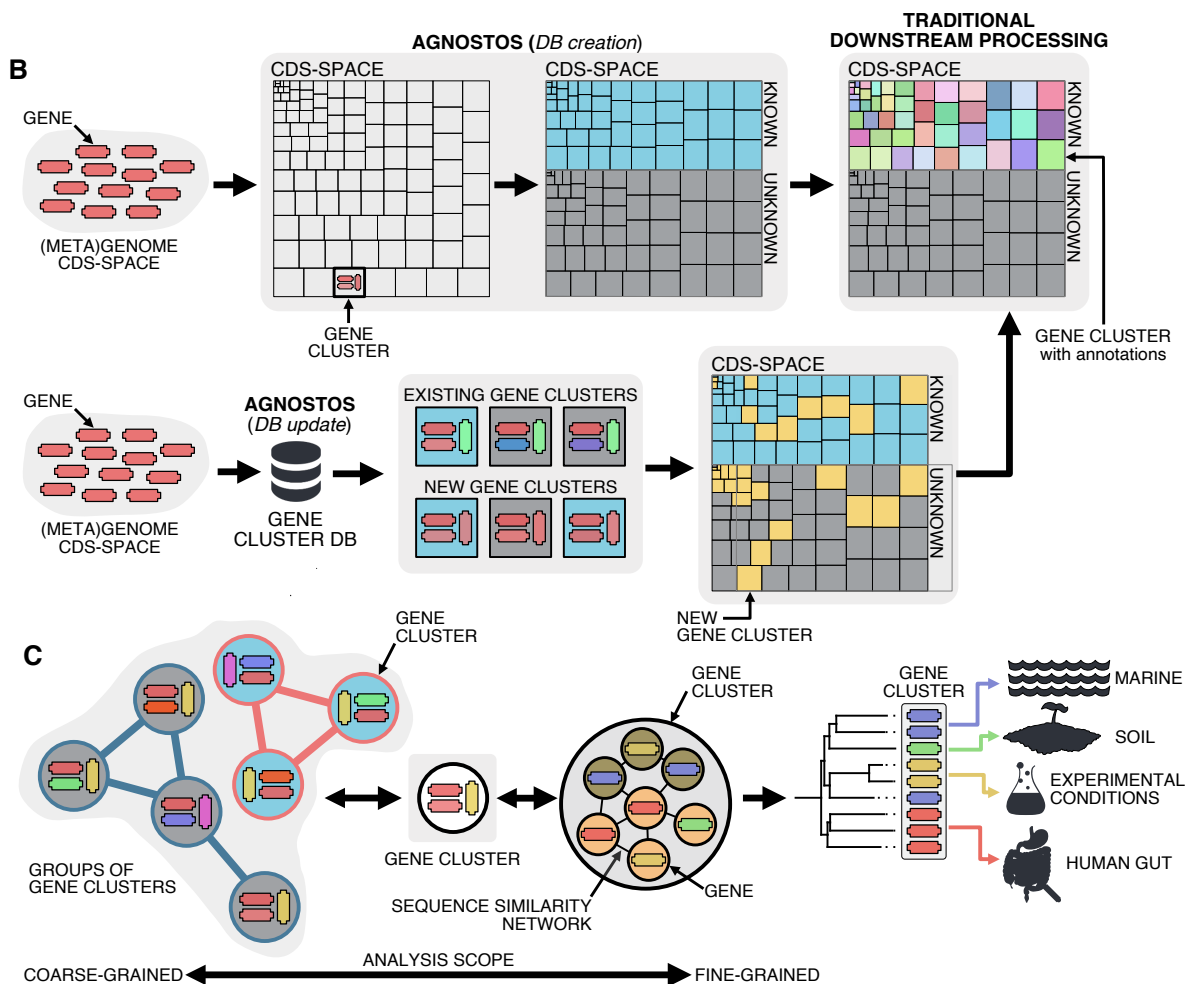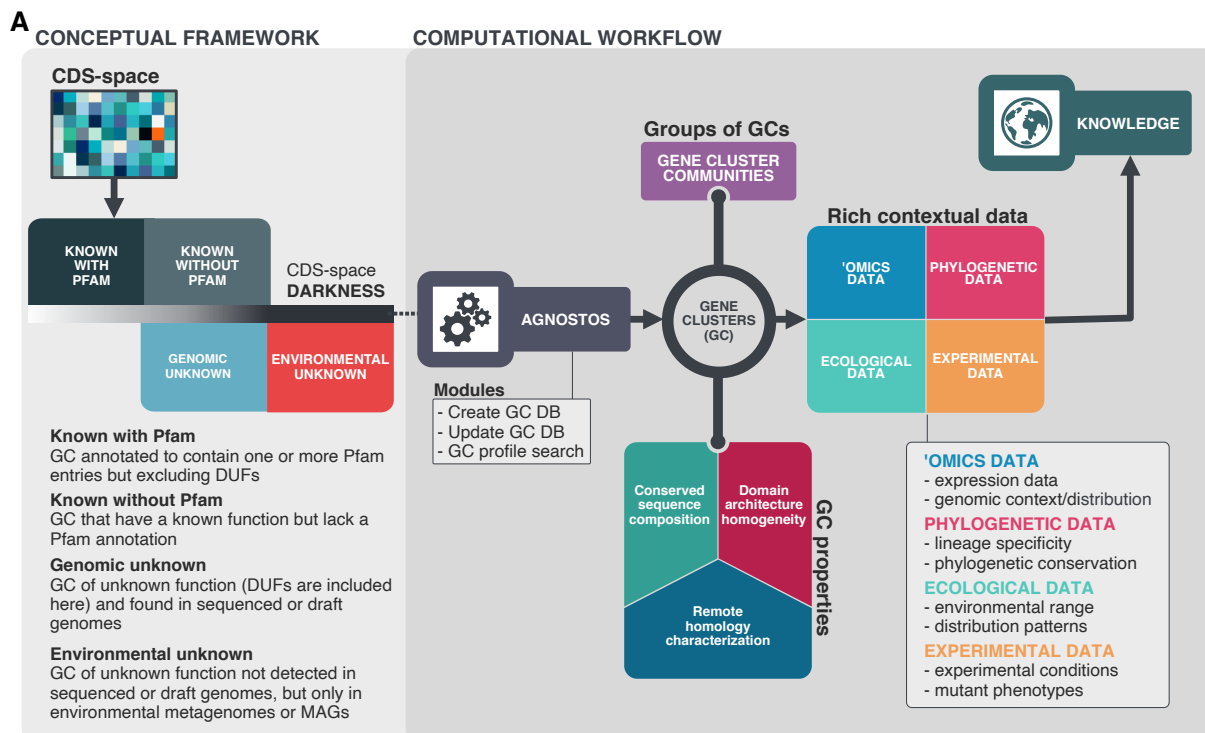
80    al., 2018) developed a high-throughput experimental pipeline that provides mutant phenotypes

81    for thousands of bacterial genes of unknown function being one of the most promising methods

82    to tackle the unknown. Despite their promise, experimental methods are labor-intensive and

83    require novel computational methods that could bridge the existing gap between the known and

84    unknown coding sequence space (CDS-space).

85    Here we present a conceptual framework and a computational workflow that closes the gap

86    between the known and unknown CDS-space by connecting genomic and metagenomic gene

87    clusters. Our approach adds context to vast amounts of unknown biology, providing an

88    invaluable resource to understand the unknown functional fraction better and boost the current

89    methods for its experimental characterization. The application of our approach to 415,971,742

90    genes predicted from 1,749 metagenomes and 28,941 bacterial and archaeal genomes

91    revealed that the unknown fraction (1) is smaller than expected, (2) is exceptionally diverse, and

92    (3) is phylogenetically more conserved and predominantly lineage-specific at the Species level.

93    Finally, we show how we can connect all the outputs produced by our approach to augment the

94    results from experimental data and add context to genes of unknown function through

95    hypothesis-driven molecular investigations.

# Results

## A conceptual framework and a computational workflow to unify the known and the unknown coding sequence space

We created the conceptual and technical foundations to unify the known and unknown CDS-space facilitating the integration of the genes of unknown function into ecological, evolutionary and biotechnological investigations. First, we conceptually partitioned the known and unknown fractions into (1) Known with Pfam annotations (K), (2) Known without Pfam annotations (KWP), (3) Genomic unknown (GU), and (4) Environmental unknown (EU) (Fig. 1A). The framework introduces a subtle change of paradigm compared to traditional approaches where our objective is to provide the best representation of the unknown space. We gear all our efforts towards finding sequences without any evidence of known homologies by pushing the search space beyond the *twilight zone* of sequence similarity (Rost, 1999). With this objective in mind, we use gene clusters (GCs) instead of genes as the fundamental unit to compartmentalize the CDS-space owing to their unique characteristics (Fig. 1B). (1) GCs produce a structured CDS-space reducing its complexity (Fig. 1B), (2) are independent of the known and unknown fraction, (3) are conserved across environments and organisms, and (4) can be used to aggregate information from different sources (Fig. 1A). Moreover, the GCs provide a good compromise in terms of resolution for analytical purposes, and owing to their unique properties, one can perform analyses at different scales. For fine-grained analyses, we can exploit the gene associations within each GC; and for coarse-grained analyses, we can create groups of GCs based on their shared homologies (Fig. 1B).

**A** CONCEPTUAL FRAMEWORK — COMPUTATIONAL WORKFLOW

CDS-space

Groups of GCs — GENE CLUSTER COMMUNITIES

KNOWLEDGE

KNOWN WITH PFAM — KNOWN WITHOUT PFAM — CDS-space DARKNESS

GENOMIC UNKNOWN — ENVIRONMENTAL UNKNOWN

AGNOSTOS → GENE CLUSTERS (GC)

Rich contextual data
'OMICS DATA — PHYLOGENETIC DATA
ECOLOGICAL DATA — EXPERIMENTAL DATA

Modules
- Create GC DB
- Update GC DB
- GC profile search

**Known with Pfam**
GC annotated to contain one or more Pfam entries but excluding DUFs

**Known without Pfam**
GC that have a known function but lack a Pfam annotation

**Genomic unknown**
GC of unknown function (DUFs are included here) and found in sequenced or draft genomes

**Environmental unknown**
GC of unknown function not detected in sequenced or draft genomes, but only in environmental metagenomes or MAGs

GC properties: Conserved sequence composition — Domain architecture homogeneity — Remote homology characterization

**'OMICS DATA**
- expression data
- genomic context/distribution
**PHYLOGENETIC DATA**
- lineage specificity
- phylogenetic conservation
**ECOLOGICAL DATA**
- environmental range
- distribution patterns
**EXPERIMENTAL DATA**
- experimental conditions
- mutant phenotypes

**B**

AGNOSTOS (*DB creation*) — TRADITIONAL DOWNSTREAM PROCESSING

GENE → (META)GENOME CDS-SPACE

CDS-SPACE — GENE CLUSTER

CDS-SPACE — KNOWN / UNKNOWN

CDS-SPACE — KNOWN / UNKNOWN — GENE CLUSTER with annotations

GENE → (META)GENOME CDS-SPACE

AGNOSTOS (*DB update*) — GENE CLUSTER DB

EXISTING GENE CLUSTERS
NEW GENE CLUSTERS

CDS-SPACE — KNOWN / UNKNOWN — NEW GENE CLUSTER

**C**

GENE CLUSTER

GROUPS OF GENE CLUSTERS

GENE CLUSTER — SEQUENCE SIMILARITY NETWORK — GENE

GENE CLUSTER — MARINE / SOIL / EXPERIMENTAL CONDITIONS / HUMAN GUT

COARSE-GRAINED ◄— ANALYSIS SCOPE —► FINE-GRAINED

118

7

**Figure 1:** Conceptual framework to unify the known and unknown CDS-space and integration of the framework in the current analytical workflows (A) Link between the conceptual framework and the computational workflow to partition the CDS-space in the four conceptual categories. AGNOSTOS infers, validates and refines the GCs and combines them in gene cluster communities (GCCs). Then, it classifies them in one of the four conceptual categories based on their level of 'darkness'. Finally, we add context to each GC based on several sources of information, providing a robust framework for generating hypotheses that can be used to augment experimental data. (B) The computational workflow provides two mechanisms to structure the CDS-space using GCs, de novo creation of the GCs (*DB creation*), or integrating the dataset in an existing GC database (*DB update*). The structured CDS-space can then be plugged into traditional analytical workflows to annotate the genes within each GC of the known fraction. With AGNOSTOS, we provide the opportunity to integrate the unknown fraction into the current microbiome analyses easily. C) The versatility of the GCs enables analyses at different scales depending on the scope of our experiments. We can group GCs in gene cluster communities based on their shared homologies to perform coarse-grained analyses. On the other hand, we can design fine-grained analyses using the relationships between the genes in a GC, i.e., detecting network modules in the GC inner sequence similarity network. Additionally, given that GCs are conserved across environments, organisms and experimental conditions give us access to an unprecedented amount of information to design and interpret experimental data.

Driven by the concepts defined in the conceptual framework, we developed AGNOSTOS, a computational workflow that infers, validates, refines, and classifies GCs in the four proposed categories (Fig. 1A; Fig. 1B; Supp. Fig 1). AGNOSTOS provides two operational modules (*DB creation* and *DB update*) to produce GCs with a highly conserved intra-homogeneous structure (Fig. 1B), both in terms of sequence similarity and domain architecture homogeneity; it exhausts any existing homology to known genes and provides a proper delimitation of the unknown CDS-space before classifying each GC in one of the four categories (Methods). In the last step, we decorate each GC with a rich collection of contextual data compiled from different sources or generated by analyzing the GC contents in different contexts (Fig. 1A). For each GC, we also offer several products that can be used for analytical purposes like improved representative sequences, consensus sequences, sequence profiles for MMseqs2 (Steinegger & Soding, 2017) and HHblits (Steinegger, Meier, et al., 2019), or the GC members as a sequence similarity network (Methods). To complement the collection, we also provide a subset of what we define as *high-quality* GCs. The defining criteria are (1) the representative is a complete gene and (2) more than one-third of genes within a GC are complete genes.

## Partitioning and contextualizing the coding sequence space of genomes and metagenomes

We used our approach to explore the unknown CDS-space of 1,749 microbial metagenomes derived from human and marine environments, and 28,941 genomes from the GTDB_r86 (Supp Fig 2A). The initial gene prediction of AGNOSTOS (Supp Fig 1) produced 322,248,552 genes from the environmental dataset and assigned Pfam annotations to 44% of them. Next, it clustered the predicted genes in 32,465,074 GCs. For the downstream processing, we kept 3,003,897 GCs (83% of the original genes) after filtering out any GC that contained less than ten genes (Skewes-Cox, Sharpton, Pollard, & DeRisi, 2014) removing 9,549,853 clusters and 19,911,324 singletons (Supp Fig 2A; Supp. Note 1). The validation process selected 2,940,257 *good-quality* clusters (Fig. 1B; Supp. Table 1; Supp. Note 2), which resulted in 43% of them being members of the unknown CDS-space after the classification and remote homology refinement steps (Supp Fig 2A, Supp. Note 3). We build the link between the environmental and genomic CDS-space by expanding the final collection of GCs with the genes predicted from GTDB_r86 (Supp Fig 2A). Our environmental GCs already included 72% of the genes from GTDB_r86; 22% of them created 2,400,037 new GCs, and the rest 6% resulted in singleton GCs (Supp Fig 2A; Supp. Note 4; Supp. Note 5). The final dataset includes 5,287,759 GCs (Supp Fig 2A), with both datasets sharing only 922,599 GCs (Supp Fig 2B). The addition of the GTDB_r86 genes increased the proportion of GCs in the unknown CDS-space to 54%. As the final step, the workflow generated a subset of 203,217 *high-quality* GCs (Supp. Table 2; Supp. Fig 3). In these *high-quality* clusters, we identified 12,313 clusters potentially encoding for small proteins (<= 50 amino acids). Most of these GCs are unknown (66%), which agrees with recent findings on novel small proteins from metagenomes (Sberro et al., 2019). The KWP category contains the largest proportion of incomplete genes (Supp. Table 3), disrupting the detection and assignment of Pfam domains. But it also incorporates sequences with an unusual amino
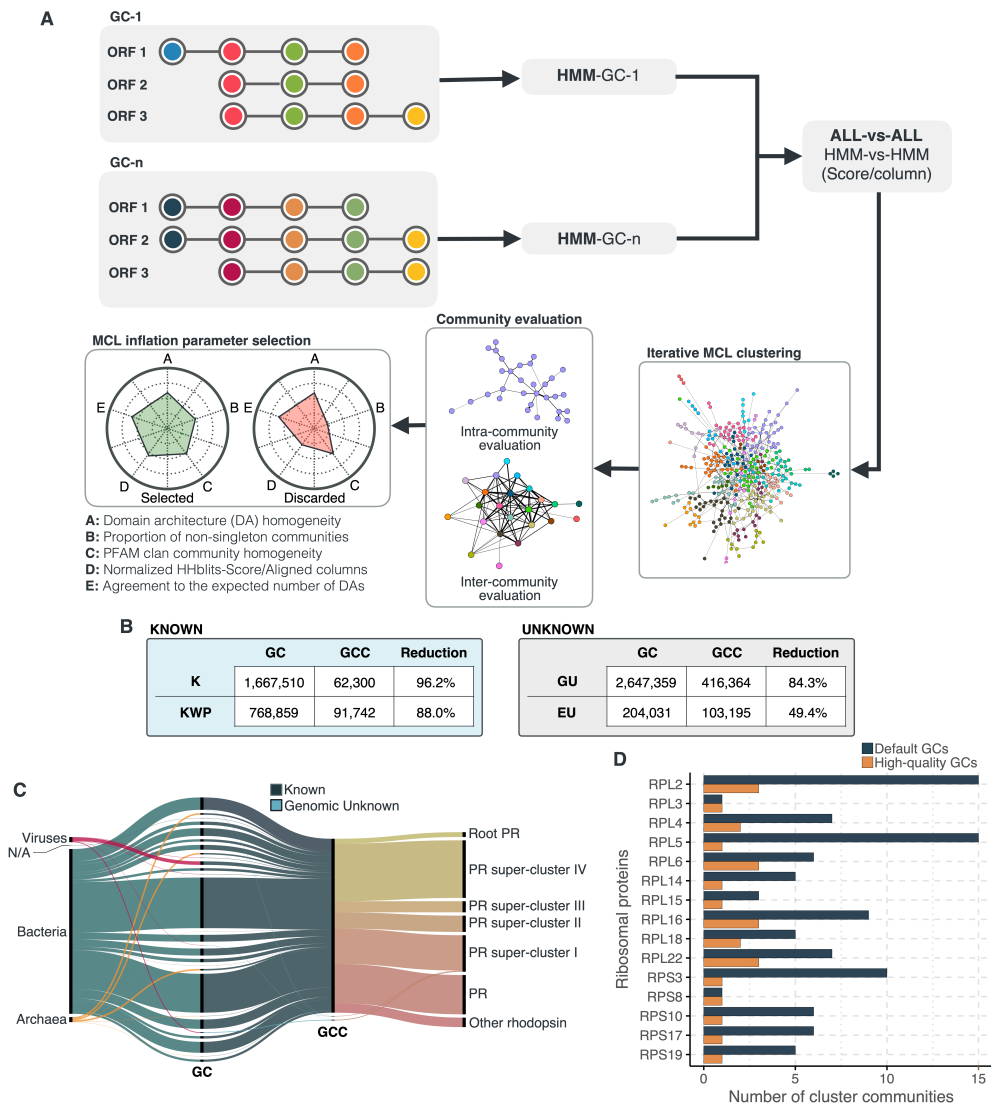
178    acid composition that has homology to proteins with high levels of disorder in the DPD database

179    (Perdigão, Rosa, & O'Donoghue, 2017) and has characteristic functions of intrinsically

180    disordered proteins (Habchi, Tompa, Longhi, & Uversky, 2014) (IDP) like cellular processes and

181    signaling as predicted by eggNOG annotations (Supp. Table 4).

182    As part of the workflow, each GC is complemented with a rich set of information, as shown in

183    Fig 1A (Supp. Table 5; Supp. Note 6).

## Beyond the twilight zone, communities of gene clusters

185    The method we developed to group GCs in gene cluster communities (GCCs) (Fig. 2A) reduced

186    the final collection of GCs by 87%, producing 673,601 GCCs (Methods; Fig. 2B; Supp. Note 7).

187    We validated our approach to capture remote homologies between related GCs using two well-

188    known gene families present in our environmental datasets, proteorhodopsins (Olson,

189    Yoshizawa, Boeuf, Iwasaki, & DeLong, 2018) and bacterial ribosomal proteins (Méheust,

190    Burstein, Castelle, & Banfield, 2019). Our dataset contained 64 GCs (12,184 genes) and 3

191    GCCs (Supp Note 8) classified as proteorhodopsin (PR). One *Known* GCC contained 99% of

192    the PR annotated genes (Fig. 2C), except 85 genes taxonomically annotated as viral and

193    assigned to the *PR Supercluster I* (Boeuf, Audic, Brillet-Guéguen, Caron, & Jeanthon, 2015)

194    within two GU communities (five GU gene clusters; Supp. Note 8). For the ribosomal proteins,

195    the results were not so satisfactory. We identified 1,843 GCs (781,579 genes) and 98 GCCs.

196    The number of GCCs is larger than the expected number of ribosomal protein families (16) used

197    for validation. When we use *high-quality* GCs (Supp. Note 8), we get closer to the expected

198    number of GCCs (Fig. 2D). With this subset, we identified 26 GCCs and 145 GCs (1,687

199    genes). The cross-validation of our method against the approach used in Méheust et al.

200    (Méheust et al., 2019) (Supp. Note 9) confirms the intrinsic complexity of analyzing

201    metagenomic data. Both approaches showed a high agreement in the GCCs identified (Supp.

202    Table 9-1). Still, our method inferred fewer GCCs for each of the ribosomal protein families

203   (Supplementary Figure 9-3), coping better with the nuisances of a metagenomic setup, such as

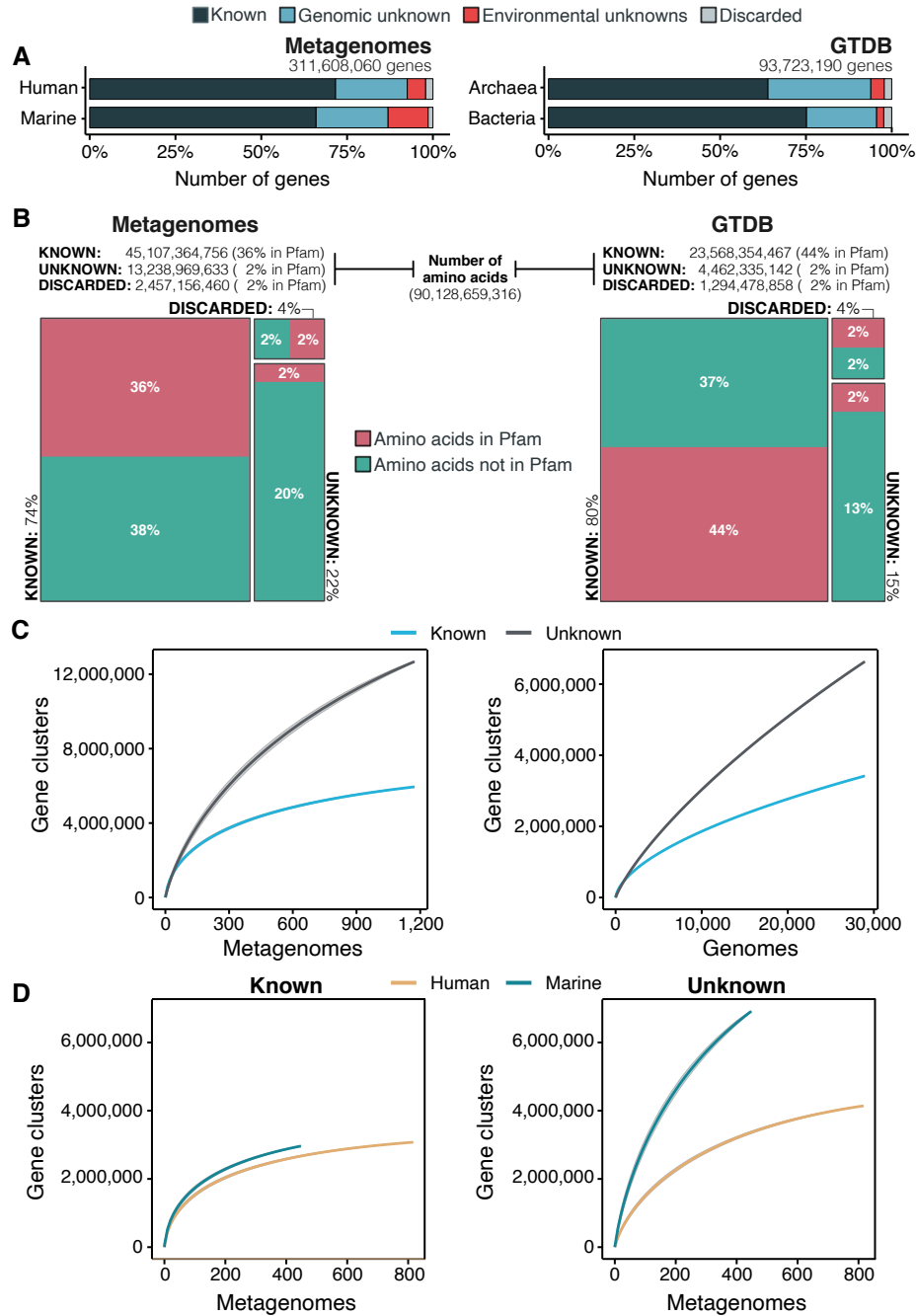204   incomplete genes (Supp. Table 6).

205



206

**Figure 2:** Overview and validation of the workflow to aggregate GCs in communities. (A) We inferred a
gene cluster homology network using the results of an all-vs-all HMM gene cluster comparison with
HHBLITS. The edges of the network are based on the HHblits-score/Aligned-columns. Communities are
identified by an iterative screening of different MCL inflation parameters and evaluated using five different
metrics that consider the inter- and intra-community properties. (B) Comparison of the number of GCs
and GCCs for each of the functional categories. (C) Validation of the GCCs inference based on the
environmental genes annotated as proteorhodopsins. Ribbons in the alluvial plot are genes, and each
stacked bar corresponds (from left to right) to the (1) gene taxonomic classification at the domain level,
(2) GC membership, (3) GCC membership and (4) MicRhoDE operational classification. (D) Validation of
the GCCs inference based on ribosomal proteins based on standard and high-quality GCs.

11

## 217 A smaller but highly diverse unknown coding sequence space

218 Combining clustering and remote homology searches reduces the extent of the unknown CDS-

219 space compared to what is reported by the traditional genomic and metagenomic analysis

220 approaches (Fig. 3A). Our workflow recruited as much as 71% of genes in human-related

221 metagenomic samples and 65% of the genes in marine metagenomes into the known CDS-

222 space. In both human and marine microbiomes, the genomic unknown fraction showed a similar

223 proportion of genes (21%, Fig. 3A). The number of genes corresponding to EU gene clusters is

224 higher in marine metagenomes; 12% of the genes are part of this GC category. We obtained a

225 comparable result when we evaluated the genes from the GTDB_r86, 75% of bacterial and 64%

226 of archaeal genes were part of the known CDS-space. Archaeal genomes contained more

227 unknowns than those from Bacteria, where 30% of the genes are classified as genomic

228 unknowns in Archaea, and only 20% in Bacteria (Fig. 3A; Supp. Table 7). Our approach allows

229 us to go beyond genes, and for the first time, we can provide a detailed characterization of the

230 CDS-space at the amino acid level. From the 90,128,659,316 amino acids analyzed, the

231 majority of the amino acids in metagenomes (74%) and GTDB_r86 (80%) are in the known

232 CDS-space (Fig. 3B; Supp. Table 7). In both cases, approximately 40% of the amino acids in

233 the known CDS-space were part of a Pfam domain (Fig. 3B; Supp. Table 7). The proportion of

234 amino acids in the unknown CDS-space ranged from 22% in metagenomes and 15% in

235 GTDB_r86. Pfam domains covered only 2% of the amino acids in the unknown CDS-space in

236 both cases. To evaluate the differences between the two CDS-space fractions, we calculated

237 the accumulation rates of GCs and GCCs. For the metagenomic dataset we used 1,264

238 metagenomes (18,566,675 GCs and 282,580 GCCs) and for the genomic dataset 28,941

239 genomes (9,586,109 GCs and 496,930 GCCs). The rate of accumulation of unknown GCs was

240 three times higher than the known (2 times for the genomic), and both cases were far from

241 reaching a plateau (Fig. 3C). This is not the case for the GCC accumulation curves (Supp Fig

242    4B), where they reached a plateau. The accumulation rate is largely determined by the number

243    of singletons, especially singletons from EUs (Supp note 11 and Supp. Fig 5). While the

244    accumulation rate of known GCs between marine and human metagenomes is almost identical,

245    there are striking differences for the unknown GCs (Fig. 3D). These differences are maintained

246    even when we remove the virus-enriched samples from the marine metagenomes (Supp Fig

247    4A). Although the marine metagenomes include a large variety of environments, from coastal to

248    the deep sea, the known space remains quite constrained.

249    Despite only including marine and human metagenomes in our database, our coverage of other

250    databases and environments is quite comprehensive, with an overall coverage of 76% (Supp.

251    Note 12). The lowest covered biomes are freshwater, soil and human non-digestive as revealed

252    by the screening of MGnify (Mitchell et al., 2020) (release 2018_09; 11 biomes; 843,535,6116

253    proteins) where we assigned 74% of the MGnify proteins into one of our categories

254    (Supplementary Fig. 6). Furthermore, as a result of this evaluation, we classified 20% of the

255    FunkFams (Wyman et al., 2018) and 44% of the unknowns used by Price et al. (Price et al.,
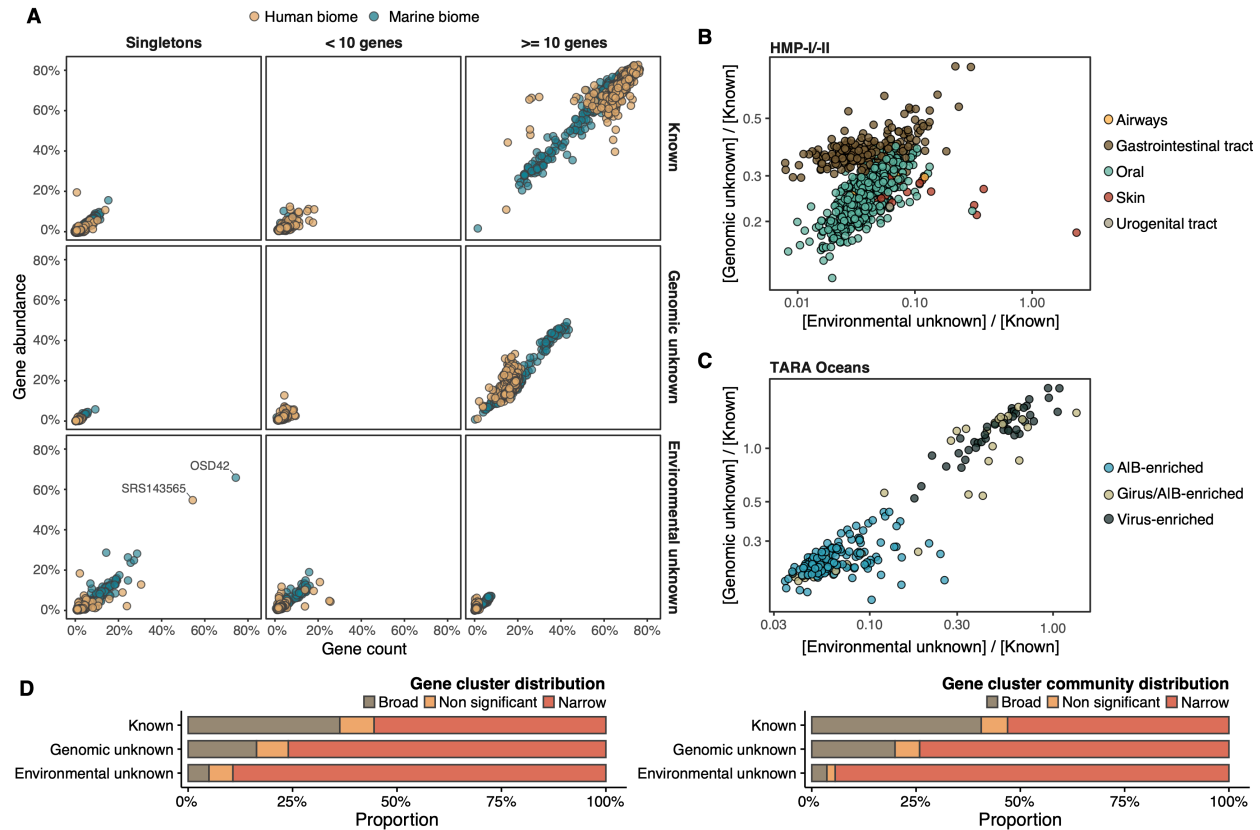
256    2018) to the known fraction (Supp. Table 12-1).

**Figure 3:** The extent of the known and unknown coding sequence space (A) Proportion of genes in the known and unknown. (B) Amino acid distribution in the known and unknown CDS-space. (C) Accumulation curves for the known and unknown CDS-space at the GC- level for the metagenomic and genomic data. from TARA, MALASPINA, OSD2014 and HMP-I/II projects. (D) Collector curves comparing the human and marine biomes. Colored lines represented the mean of 1000 permutations and shaded in grey the standard deviation. Non-abundant singleton clusters were excluded from the accumulation curves calculation.

# The unknown coding sequence space has a limited ecological distribution in human and marine environments

Although the role of the unknown fraction in the environment is still a mystery, the large number of gene counts and abundance observed underlines its inherent ecological relevance (Fig. 4A). In some metagenomes, the genomic unknown fraction can account for more than 40% of the total gene abundance observed (Fig. 4A). The environmental unknown fraction is also relevant in several samples, where singleton GCs are the majority (Fig. 4A). We identified two metagenomes with an unusual composition in terms of environmental unknown singletons. The marine metagenome corresponds to a sample from Lake Faro (OSD42), a meromictic saline with a unique extreme environment where Archaea plays an important role (La Cono et al., 2013). The HMP metagenome (SRS143565) corresponds to a human sample from the right cubital fossa from a healthy female subject. To understand this unusual composition, we should perform further analyses to discard potential technical artifacts like sample contamination. The ratio between the unknown and known GCs revealed that the metagenomes located at the upper left quadrant in Fig. 4B-C are enriched in GCs of unknown function. In human metagenomes, we can distinguish between body sites, with the gastrointestinal tract, where microbial communities are expected to be more diverse and complex, significantly enriched with genomic unknowns. The HMP metagenomes with the largest ratio of unknowns are those samples identified to contain crAssphages (Dubinkina, Ischenko, Ulyantsev, Tyakht, & Alexeev, 2016; Edwards et al., 2019) and HPV viruses (Ma et al., 2014) (Supp. Table 8; Supp. Fig. 7). Consistently, in marine metagenomes (Fig. 4D) we can separate between size fractions, where the highest ratio in genomic and environmental unknowns corresponds to the ones enriched with viruses and giant viruses.

15

**Figure 4:** Distribution of the unknown coding sequence space in the human and marine metagenomes (A) Ratio between the proportion of the number of genes and their estimated abundances per cluster category and biome. Columns represented in the facet depicts three cluster categories based on the size of the clusters. (B) Relationship between the ratio of Genomic unknowns and Environmental unknowns in the HMP-I/II metagenomes. Gastrointestinal tract metagenomes are enriched in Genomic unknown coding sequences compared to the other body sites. (C) Relationship between the ratio of Genomic unknowns and Environmental unknowns in the TARA Oceans metagenomes. Girus and virus enriched metagenomes show a higher proportion of both unknown coding sequences (genomic and environmental) than the Archaea|Bacteria enriched fractions. (D) Environmental distribution of GCs and GCCs based on Levin's niche breadth index. We obtained the significance values after generating 100 null gene cluster abundance matrices using the *quasiswap* algorithm.
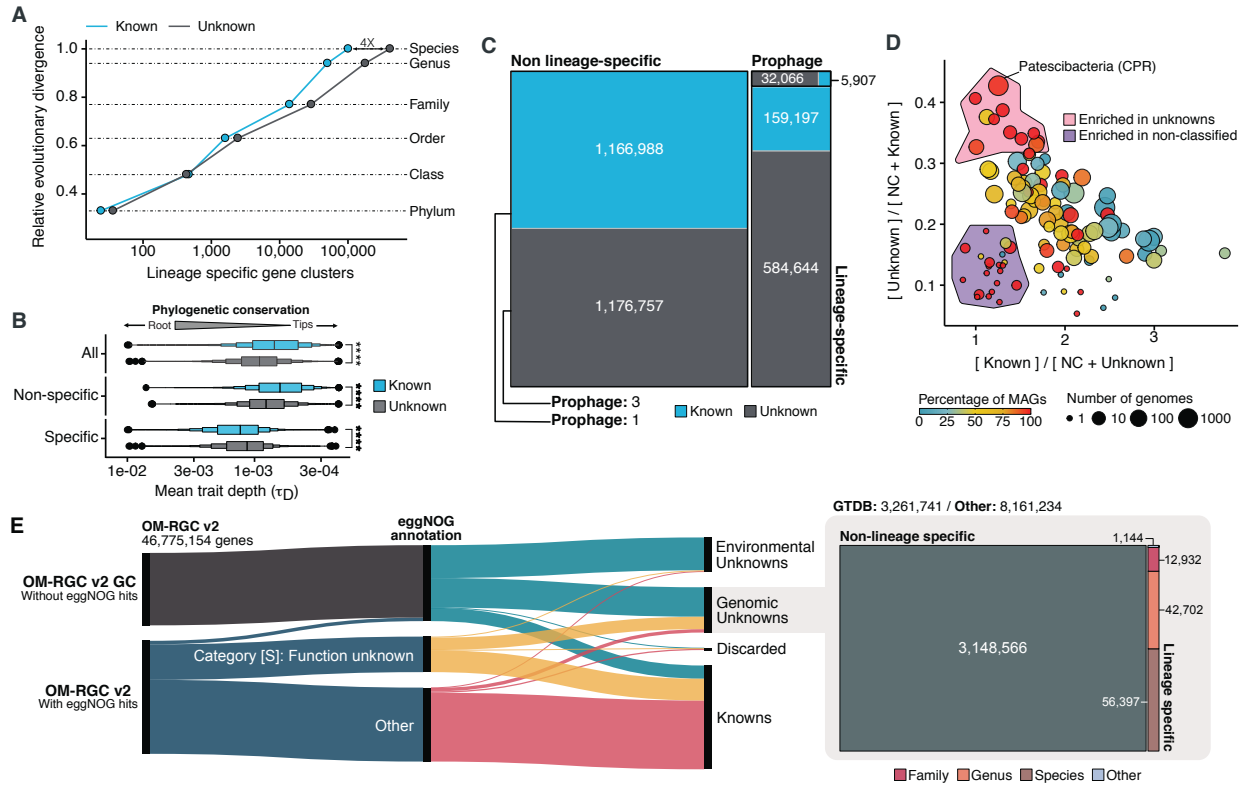
To complement the previous findings, we performed a large-scale analysis to investigate the GC occurrence patterns in the environment. The narrow distribution of the unknown fraction (Fig. 4D) suggests that these GCs might provide a selective advantage and be necessary to adapt to specific environmental conditions. But the pool of broadly distributed environmental unknowns is the most exciting result. We identified traces of potential ubiquitous organisms left

16

306      uncharacterized by traditional approaches, as more than 80% of these GCs cannot be

307      associated with a metagenome-assembled genome (MAG) (Supp Table 9, Supp. Note 10).

## 308 The genomic unknown coding sequence space is lineage-specific

309      With the inclusion of the genomes from GTDB_r86, we have access to a phylogenomic

310      framework that can be used to assess how exclusive is a GC within a lineage (lineage-specifity

311      (Mendler et al., 2019)) and how conserved is a GC across clades (phylogenetic conservation

312      (Martiny, Treseder, & Pusch, 2013)). We identified 781,814 lineage-specific GCs and 464,923

313      phylogenetically conserved ($P < 0.05$) GCs in Bacteria (Supp. Table 10; Supp. Note 13 for

314      Archaea). The number of lineage-specific GCs increases with the Relative Evolutionary

315      Distance (Parks et al., 2018) (Fig. 5A) and differences between the known and the unknown

316      fraction start to be evident at the Family level resulting in 4X more unknown lineage-specific

317      GCs at the Species level. The unknown GCs are more phylogenetically conserved than the

318      known (Fig. 5B, $p < 0.0001$), revealing the importance of the genome's uncharacterized fraction.

319      However, this is not the case for the lineage-specific and phylogenetically conserved GCs,

320      where the unknown GCs are less phylogenetically conserved (Fig. 5B), agreeing with the large

321      number of lineage-specific GCs at Genus and Species level. To discard the possibility that the

322      lineage-specific GCs of unknown function have a viral origin, we screened all GTDB_r86

323      genomes for prophages. We only found 37,163 lineage-specific GCs in prophage genomic

324      regions, being 86% GCs of unknown function. After unveiling the potential relevance of the GCs

325      of unknown function in bacterial genomes, we identified phyla in GTDB_r86 enriched with these

326      types of clusters. A clear pattern emerged when we partitioned the phyla based on the ratio of

327      known to unknown GCs and vice versa (Fig. 5D), the phyla with a larger number of MAGs are

328      enriched in GCs of unknown function Figure 5D. Phyla with a high proportion of non-classified

329      GCs (those discarded during the validation steps) contain a small number of genomes and are

330      primarily composed of MAGs. These groups of phyla highly enriched in unknowns and

17

331     represented mainly by MAGs include newly described phyla such as *Cand.* Riflebacteria and

332     *Cand.* Patescibacteria (Anantharaman et al., 2018; Brown et al., 2015; Rinke et al., 2013), both

333     with the largest unknown to known ratio.

334



335
336
337     **Figure 5:** Phylogenomic exploration of the unknown coding sequence space. (A) Distribution of the
338     lineage-specific GCs by taxonomic level. Lineage-specific unknown GCs are more abundant in the lower
339     taxonomic levels (Genus, Species). (B) Phylogenetic conservation of the known and unknown coding
340     sequence space in 27,372 bacterial genomes from GTDB_r86. We observe differences in the
341     conservation between the known and the unknown coding sequence space for lineage- and non-lineage
342     specific GCs (paired Wilcoxon rank-sum test; all p-values < 0.0001). (C) The majority of the lineage-
343     specific clusters are part of the unknown coding sequence space, being a small proportion found in
344     prophages present in the GTDB_r86 genomes. (D) Known and unknown coding sequence space of the
345     27,732 GTDB_r86 bacterial genomes grouped by bacterial phyla. Phyla are partitioned based on the ratio
346     of known to unknown GCs and vice versa. Phyla enriched in MAGs have higher proportions in GCs of
347     unknown function. Phyla with a high proportion of non-classified clusters (NC; discarded during the
348     validation steps) tend to contain a small number of genomes. (E) The alluvial plot's left side shows the
349     uncharacterized (OM-RGC v2 GC) and characterized (OM-RGC v2) fraction of the gene catalog. The
350     functional annotation is based on the eggNOG annotations provided by Salazar et al.(Salazar et al.,
351     2019). The right side of the alluvial plot shows the new organization of the OM-RGC v2 coding sequence
352     space based on the approach described in this study. The treemap in the right links the metagenomic and
353     genomic space adding context to the unknown fraction of the OM-RGC v2
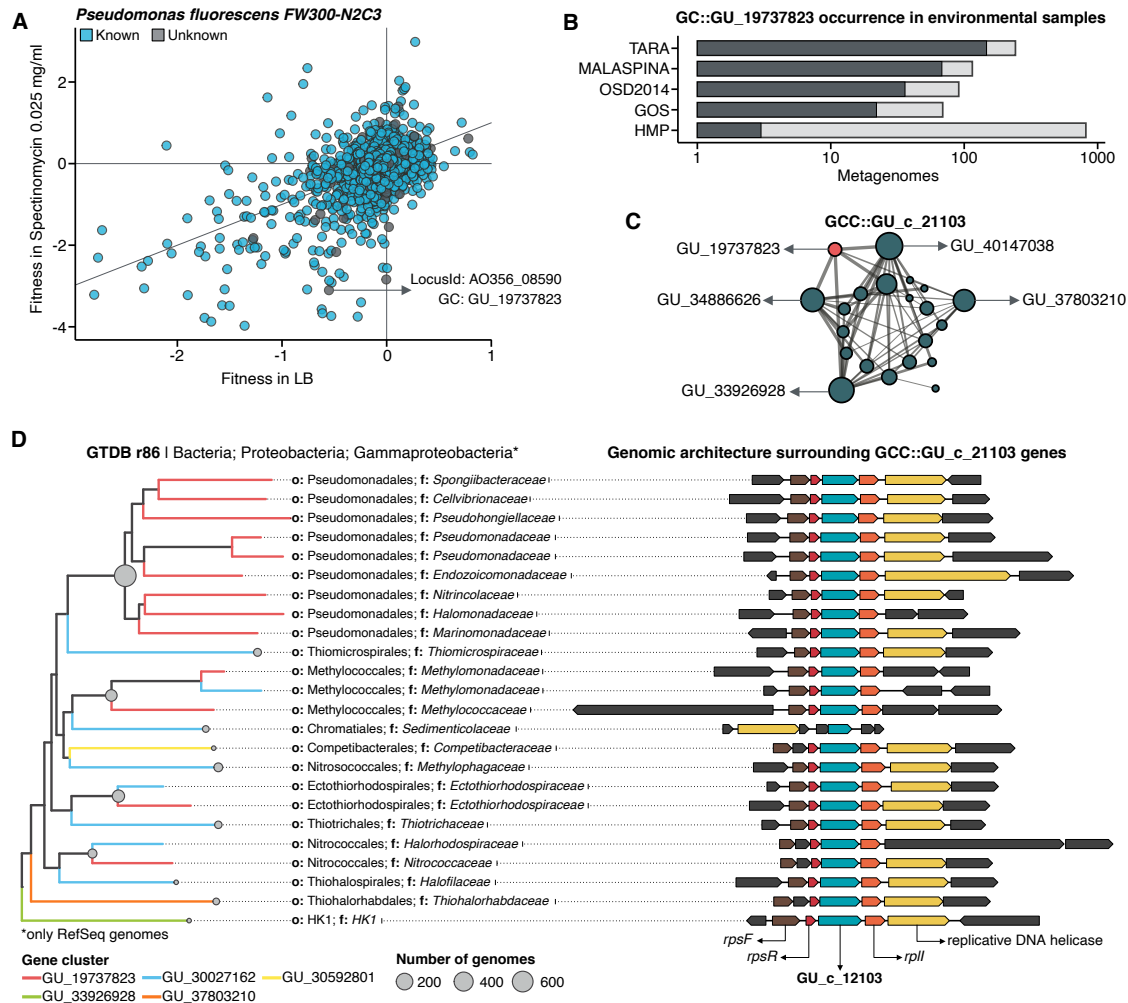
354

355  We demonstrate the possibility of bridging genomic and metagenomic data and simultaneously

356  unifying the known and unknown CDS-space by integrating the new Ocean Microbial Reference

357  Gene Catalog (Salazar et al., 2019) (OM-RGC v2) into our database. We assigned 26,170,875

358  genes to known GCs, 11,422,975 to genomic unknowns, 8,661,221 to environmental unknown

359  and 520,083 were discarded. From the 11,422,975 genes classified as genomic unknowns, we

360  could associate 3,261,741 to a GTDB_r86 genome and we identified 113,175 as lineage-

361  specific. The alluvial plot in Fig. 5E depicts the new organization of the OM-RGC v2 after being

362  integrated into our framework and how we can provide context to the two original types of

363  unknowns in the OM-RGC (those annotated as category S in eggNOG (Huerta-Cepas et al.,

364  2019) and those without known homologs in the eggNOG database (Salazar et al., 2019)) that

365  can lead to potential experimental targets at the organism level to complement the

366  metatranscriptomic approach proposed by Salazar et al. (Salazar et al., 2019).

## 367 A structured coding sequence space augments the interpretation
## 368 of experimental data

369  We selected one of the experimental conditions tested in Price et al. (Price et al., 2018) to

370  demonstrate the potential of our approach to augment experimental data. We compared the

371  fitness values in plain rich medium with added Spectinomycin dihydrochloride pentahydrate to

372  the fitness in plain rich medium (LB) in *Pseudomonas fluorescens FW300-N2C3* (Fig. 6A). This

373  antibiotic inhibits protein synthesis and elongation by binding to the bacterial 30S ribosomal

374  subunit and interferes with the peptidyl tRNA translocation. We identified the gene with locus id

375  AO356_08590 that presents a strong phenotype (fitness = -3.1; t = -9.1) and has no known

376  function. This gene belongs to the genomic unknown GC GU_19737823. We can track this GC

377  into the environment and explore the occurrence in the different samples we have in our

19

378    database. As expected, the GC is mostly found in non-human metagenomes (Fig. 6B) as

379    *Pseudomonas* are common inhabitants of soil and water environments(Heffernan, Murphy, &

380    Casey, 2009). However, finding this GC also in human-related samples is very interesting due

381    to the potential association of *P. fluorescens* and human disease where Crohn's disease

382    patients develop serum antibodies to this microbe (Scales, Dickson, LiPuma, & Huffnagle,

383    2014). We can add another layer of information to the selected GC by looking at the associated

384    remote homologs in the GCC GU_c_21103 (Fig. 6C). We identified all the genes in the

385    GTDB_r86 genomes that belong to the GCC GU_c_21103 (Supp. Table 11) and explored their

386    genomic neighborhoods. All members from GU_c_21103 are constrained to the class

387    *Gammaproteobacteria*, and interestingly GU_19737823 is mostly exclusive to the order

388    *Pseudomonadales*. The gene order in the different genomes analyzed is highly conserved,

389    finding GU_19737823 after the *rpsF*::*rpsR* operon and before *rplI*. *rpsF* and *rpsR* encode for

390    30S ribosomal proteins, the prime target of spectinomycin. The combination of the experimental

391    evidence and the associated data inferred by our approach provides strong support to generate

392    the hypothesis that the gene AO356_08590 might be involved in the resistance to

393    spectinomycin.

394

**Figure 6:** Augmenting experimental data with GCs of unknown function. (A) We used the fitness values from the experiments from Price et al.(Price et al., 2018) to identify genes of unknown function that are important for fitness under certain experimental conditions. The selected gene belongs to the genomic unknown GC GU_19737823 and presents a strong phenotype (fitness = -3.1; t = -9.1) (B) Occurrence of GU_19737823 in the metagenomes used in this study. Darker bars depict the number of metagenomes where the GC is found. (C) GU_19737823 is a member of the GCC GU_c_21103. The network shows the relationships between the different GCs members of the gene cluster community GU_c_21103. The size of the node corresponds to the node degree of each GC. Edge thickness corresponds to the bitscore/column metric. Highlighted in red is GU_19737823. (D) We identified all the genes in the GTDB_r86 genomes that belong to the GCC GU_c_21103 and explored their genomic neighborhoods. GU_c_21103 members were constrained to the class Gammaproteobacteria, and GU_19737823 is mostly exclusive to the order Pseudomonadales. The gene order in the different genomes analyzed is highly conserved, finding GU_19737823 after the *rpsF::rpsR* operon and before *rpII. rpsF* and *rpsR* encode for the *30S ribosomal protein S6* and *30S ribosomal protein S18,* respectively. The GTDB_r86 subtree only shows RefSeq genomes. Branch colors correspond to the different GCs found in GU_c_21103. The bubble plot depicts the number of genomes with a gene that belongs to GU_c_21103.

# Discussion

413

414     We present a new conceptual framework and computational workflow to unify the known and

415     unknown CDS-space. Using this framework, we performed an in-depth exploration of the

416     microbial unknown CDS-space, demonstrating that we can link the unknown fraction of

417     metagenomic studies to specific genomes and provide a powerful tool for hypothesis

418     generation. During the last years, the microbiome community has established a standard

419     operating procedure(Quince et al., 2017) for analyzing metagenomes that we can briefly

420     summarize into (1) assembly, (2) gene prediction, (3) gene catalog inference, (4) binning, and

421     (5) characterization. Thanks to recent computational developments (Steinegger & Soding, 2017;

422     Steinegger & Söding, 2018), we envisioned an alternative to this workflow to maximize the

423     information used when analyzing genomic and metagenomic data. In addition, we provide a

424     mechanism to reconcile top-down and bottom-up approaches, thanks to the well-structured

425     CDS-space proposed by our framework. AGNOSTOS can create environmental- and organism-

426     specific variations of a seed GC database. Then, it integrates the predicted genes from new

427     genomes and metagenomes and dynamically creates and classifies new GCs with those genes

428     not integrated during the initial step (Fig. 1B). Afterward, the potential functions of the known

429     GCs can be carefully characterized by incorporating them into the traditional workflows.

430     One of the most appealing characteristics of our approach is that the GCs provide unified

431     groups of homologous genes across environments and organisms indifferently if they belong to

432     the known or unknown CDS-space, and we can contextualize the unknown fraction using this

433     genomic and environmental information. Our combination of partitioning and contextualization

434     features a smaller unknown CDS-space than we expected. On average, for our genomic and

435     metagenomic data, only 30% of the genes fall in the unknown fraction. One hypothesis to

436     reconcile this surprising finding is that the methodologies to identify remotely homologous

437     sequences in large datasets were computationally prohibitive until recently. New methods

438    (Steinegger, Meier, et al., 2019; Steinegger & Soding, 2017), like the ones used in AGNOSTOS,

439    are enabling large scale distant homology searches. Still, one has to apply conservative

440    measures to control the trade-off between specificity and sensitivity to avoid overclassification.

441    We found that most of the coding sequence space at the gene and amino acid level is known,

442    both in genomes and metagenomes. However, it presents a high diversity, as shown in the GC

443    accumulation curves highlighting the vast remaining untapped microbial fraction and its potential

444    importance for niche adaptation owing to its narrow ecological distribution. In a genomic context

445    and after ruling out the effect of prophages, the unknown fraction is predominantly Species'

446    lineage-specific and phylogenetically more conserved than the known fraction, supporting the

447    signal observed in the environmental data emphasizing that we should not ignore the unknown

448    fraction. It is worth noting that the high diversity observed in the unknowns only represents the

449    20% of the amino acids in the CDS-space we analyzed, and only 10% of this unknown amino

450    acid space is part of a Pfam domain (DUF and others). This contrasts with the numbers

451    observed in the known CDS-space, where Pfam domains include 50% of the amino acids. All

452    this evidence combined strengthens the hypothesis that the genes of unknown function,

453    especially the lineage-specific ones, might be associated with the mechanisms of microbial

454    diversification and niche adaptation due to the constant diversification of gene families and the

455    survival of new gene lineages (Francino, 2012; Muller, 2019).

456    Metagenome-assembled genomes are not only unveiling new regions of the microbial universe

457    (42% of the genomes in GTDB_r86), but they are also enriching genes of unknown function in

458    the tree of life. We investigated the unknown CDS-space of *Cand*. Patescibacteria, more

459    commonly known as Candidate Phyla Radiation (CPR), a phylum that has raised considerable

460    interest due to its unusual biology (Brown et al., 2015). We provide a collection of 54,343

461    lineage-specific GCs of unknown function at different taxonomic level resolutions (Supp. Table

462    12; Supp. Note 14), which will be a valuable resource for the CPR advancement research

463    efforts.

464　Our effort to tackle the unknown provides a pathway to unlock a large pool of likely relevant data

465　that remains untapped to analysis and discovery. By identifying a potential target gene of

466　unknown function for antibiotic resistance, we demonstrate the value of our approach and how it

467　can boost insights from model organism experiments. But severe challenges remain, such as

468　the dependence on the quality of the assemblies and their gene predictions, as shown by the

469　analysis of the ribosomal protein GCCs where many of the recovered genes are incomplete.

470　While sequence assembly has been an active area of research (Roumpeka, Wallace,

471　Escalettes, Fotheringham, & Watson, 2017), this has not been the case for gene prediction

472　methods (Roumpeka et al., 2017), which are becoming outdated(Ivanova et al., 2014) and

473　cannot cope with the current amount of data. Alternatives like protein-level assembly

474　(Steinegger, Mirdita, & Söding, 2019) combined with exploring the assembly graphs'

475　neighborhoods (Titus Brown et al., 2018) become very attractive for our purposes. In any case,

476　we still face the challenge of discriminating between real and artifactual singletons (Höps,

477　Jeffryes, & Bateman, 2018). There are currently no methods available to provide a plausible

478　solution and, at the same time, being scalable. We devise a potential solution in the recent

479　developments in unsupervised deep learning methods where they use large corpora of proteins

480　to define a language model *embedding* for protein sequences (Heinzinger et al., 2019). These

481　models could be applied to predict *embeddings* in singletons, which could be clustered or used

482　to determine their coding potential. Another issue is that we might be creating more GCs than

483　expected. We follow a conservative approach to avoid mixing multi-domain proteins in GCs

484　owing to the fragmented nature of the metagenome assemblies that could result in the split of a

485　GC. However, not only splitting can be a problem, but also lumping unrelated genes or GCs

486　owing to the use of remote homologies. Although the inference of GCCs is using very sensitive

487　methods to compare profile HMMs, low sequence diversity in GCs can limit its effectiveness.

488　Moreover, our approach is affected by the presence and propagation of contamination in

489　reference databases, a significant problem in 'omics (Breitwieser, Pertea, Zimin, & Salzberg,

24

490    2019; Steinegger & Salzberg, 2020). In our case, we only use Pfam as a source for annotation

491    owing to its high-quality and manual curation process. The categorization process of our GCs

492    depends on the information from other databases, and to minimize the potential impact of

493    contamination, we apply methods that weight the annotations of the identified homologs to

494    discriminate if a GC belongs to the known or unknown CDS-space.

495    The results presented here prove that the integration and the analysis of the unknown fraction

496    are possible. We are unveiling a brighter future, not only for microbiome analyses but also for

497    boosting eukaryotic-related studies, thanks to the increasing number of projects, including

498    metatranscriptomic data (Delmont et al., 2020; Vorobev et al., 2020). Furthermore, our work

499    lays the foundations for further developments of clear guidelines and protocols to define the

500    different levels of unknown (Thomas & Segata, 2019) and should encourage the scientific

501    community for a collaborative effort to tackle this challenge.

# Material and methods

## Genomic and metagenomic dataset

504    We used a set of 583 marine metagenomes from four of the major metagenomic surveys of the

505    ocean microbiome: Tara Oceans expedition (TARA) (Sunagawa et al., 2015), Malaspina

506    expedition (Duarte, 2015), Ocean Sampling Day (OSD) (Kopf et al., 2015), and Global Ocean

507    Sampling Expedition (GOS) (Rusch et al., 2007). We complemented this set with 1,246

508    metagenomes obtained from the Human Microbiome Project (HMP) phase I and II (Lloyd-Price

509    et al., 2017). We used the assemblies provided by TARA, Malaspina, OSD and HMP projects

510    and the long Sanger reads from GOS (Sanger, Nicklen, & Coulson, 1977). A total of 156M

511    (156,422,969) contigs and 12.8M long-reads were collected (Supp. Table 6).

25

512   For the genomic dataset, we used the 28,941 prokaryotic genomes (27,372 bacterial and 1,569

513   archaeal) from the Genome Taxonomy Database (Parks et al., 2018) (GTDB) Release 03-RS86

514   (19th August 2018).

# Computational workflow development

516   We implemented a computation workflow based on Snakemake (Köster, 2018) for the easy

517   processing of large datasets in a reproducible manner. The workflow provides three different

518   strategies to analyze the data. The module *DB-creation* creates the gene cluster database,

519   validates and partitions the gene clusters (GCs) in the main functional categories. The module

520   *DB-update* allows the integration of new sequences (either at the contig or predicted gene level)

521   in the existing gene cluster database. In addition, the workflow has a *profile-search* function to

522   quickly screen samples using the gene cluster PSSM profiles in the database.

# Metagenomic and genomic gene prediction

524   We used Prodigal (v2.6.3) (Hyatt et al., 2010) in metagenomic mode to predict the genes from

525   the metagenomic dataset. For the genomic dataset, we used the gene predictions provided by

526   Annotree (Mendler et al., 2019), since they were obtained, consistently, with Prodigal v2.6.3.

527   We identified potential spurious genes using the *AntiFam* database (Eberhardt et al., 2012).

528   Furthermore, we screened for '*shadow*' genes using the procedure described in Yooseph et al.

529   (Yooseph, Li, & Sutton, 2008).

# PFAM annotation

531   We annotated the predicted genes using the *hmmsearch* program from the *HMMER* package

532   (version: 3.1b2) (R. D. Finn, Clements, & Eddy, 2011) in combination with the Pfam database

533   v31 (Robert D. Finn et al., 2016). We kept the matches exceeding the internal gathering

534    threshold and presenting an independent e-value < 1e-5 and coverage > 0.4. In addition, we

535    took into account multi-domain annotations, and we removed overlapping annotations when the

536    overlap is larger than 50%, keeping the ones with the smaller e-value.


# Determination of the gene clusters

538    We clustered the metagenomic predicted genes using the cascaded-clustering workflow of the

539    MMseqs2 software (Steinegger & Söding, 2018) ("*--cov-mode 0 -c 0.8 --min-seq-id 0.3*"). We

540    discarded from downstream analyses the singletons and clusters with a size below a threshold

541    identified after applying a broken-stick model (Bennett, 1996). We integrated the genomic data

542    into the metagenomic cluster database using the "DB-update" module of the workflow. This

543    module uses the *clusterupdate* module of MMseqs2 (Steinegger & Soding, 2017), with the same

544    parameters used for the metagenomic clustering.


# Quality-screening of gene clusters

546    We examined the GCs to ensure their high intra-cluster homogeneity. We applied two

547    methodologies to validate their cluster sequence composition and functional annotation

548    homogeneity. We identified non-homologous sequences inside each cluster combining the

549    identification of a new cluster representative sequence via a sequence similarity network (SSN)

550    analysis, and the investigation of intra-cluster multiple sequence alignments (MSAs), given the

551    new representative. Initially, we generated an SSN for each cluster, using the semi-global

552    alignment methods implemented in *PARASAIL* (Daily, 2016) (version 2.1.5). We trimmed the

553    SSN using a custom algorithm (Chafee et al., 2018; Žure, Fernandez-Guerra, Munn, & Harder,

554    2017) that removes edges while maintaining the network structural integrity and obtaining the

555    smallest connected graph formed by a single component. Finally, the new cluster representative

556    was identified as the most central node of the trimmed SSN by the eigenvector centrality

557    algorithm, as implemented in igraph (Csardi & Nepusz, 2006). After this step, we built a multiple

558 sequence alignment for each cluster using *FAMSA* (Deorowicz, Debudaj-Grabysz, & Gudyś,

559 2016) (version 1.1). Then, we screened each cluster-MSA for non-homologous sequences to

560 the new cluster representative. Owing to computational limitations, we used two different

561 approaches to evaluate the cluster-MSAs. We used *LEON-BIS* (Vanhoutreve et al., 2016) for

562 the clusters with a size ranging from 10 to 1,000 genes and OD-SEQ (Jehl, Sievers, & Higgins,

563 2015) for the clusters with more than 1,000 genes. In the end, we applied a broken-stick model

564 (Bennett, 1996) to determine the threshold to discard a cluster.

565 The predicted genes can have multi-domain annotations in different orders, therefore to validate

566 the consistency of intra-cluster Pfam annotations, we applied a combination of w-shingling

567 (Broder, 1997) and Jaccard similarity. We used w-shingling (k-shingle = 2) to group consecutive

568 domain annotations as a single object. We measured the homogeneity of the *shingle sets* (sets

569 of domains) between genes using the Jaccard similarity and reported the median similarity

570 value for each cluster. Moreover, we took into consideration the Clan membership of the Pfam

571 domains and that a gene might contain N-, C- and M-terminal domains for the functional

572 homogeneity validation. We discarded clusters with a median similarity < 1.

573 After the validation, we refined the gene cluster database removing the clusters identified to be

574 discarded and the clusters containing ≥ 30% *shadow genes*. Lastly, we removed the single

575 shadow, spurious and non-homologous genes from the remaining clusters (Supplementary Note

576 2).

## 577 Remote homology classification of gene clusters

578 To partition the validated GCs into the four main categories, we processed the set of GCs

579 containing Pfam annotated genes and the set of not annotated GCs separately. For the

580 annotated GCs, we inferred a consensus protein domain architecture (DA) (an ordered

581 combination of protein domains) for each annotated gene cluster. To identify each gene cluster

582 consensus DA, we created directed acyclic graphs connecting the Pfam domains based on their

28

583    topological order on the genes using *igraph* (Csardi & Nepusz, 2006). We collapsed the

584    repetitions of the same domain. Then we used the gene completeness as a positive-weighting

585    value for the selection of the cluster consensus DA. Within this step, we divided the GCs into

586    "Knowns" (Known) if annotated to at least one Pfam domains of known function (DKFs) and

587    "Genomic unknowns" (GU) if annotated entirely to Pfam domains of unknown function (DUFs).

588    We aligned the sequences of the non-annotated GCs with FAMSA (Deorowicz et al., 2016) and

589    obtained cluster consensus sequences with the *hhconsensus* program from *HH-SUITE*

590    (Steinegger, Meier, et al., 2019). We used the cluster consensus sequences to perform a

591    nested search against the UniRef90 database (release 2017_11) (The UniProt Consortium,

592    2017) and NCBI *nr* database (release 2017_12) (NCBI Resource Coordinators, 2018) to retrieve

593    non-Pfam annotations with *MMSeqs2* (Steinegger & Soding, 2017) ("*-e 1e-05 --cov-mode 2 -c*

594    *0.6"*). We kept the hits within 60% of the Log(best-e-value) and searched the annotations for

595    any of the terms commonly used to define proteins of unknown function (Supp. Table 12). We

596    used a quorum majority voting approach to decide if a gene cluster would be classified as

597    *Genomic Unknown* or *Known without Pfams* based on the annotations retrieved. We searched

598    the consensus sequences without any homologs in the UniRef90 database against NCBI *nr*. We

599    applied the same approach and criteria described for the first search. Ultimately, we classified

600    as *Environmental Unknown* those GCs whose consensus sequences did not align with any of

601    the NCBI *nr* entries.

602    In addition, we developed some conservative measures to control the trade-off between

603    specificity and sensitivity for the remote homology searches such as (1) a modification of the

604    algorithm described in Hingamp et al. (Hingamp et al., 2013) to get a confident group of

605    homologs to determine if a query protein is known or unknown by a quorum majority voting

606    approach (Supp Note 3); (2) strict parameters in terms of iterations, bidirectional coverage and

607    probability thresholds for the HHblits alignments to minimize the inclusion of non-homologous

608 sequences; and (3) avoid providing annotations for our gene clusters, as we believe that

609 annotation should be a careful process done on a smaller scale and with experimental context.

## 610 Gene cluster remote homology refinement

611 We refined the *Environmental Unknown* GCs to ensure the lack of any characterization by

612 searching for remote homologies in the Uniclust database (release 30_2017_10) using the

613 HMM/HMM alignment method *HHblits* (Remmert, Biegert, Hauser, & Söding, 2012). We created

614 the HMM profiles with the *hhmake* program from the *HH-SUITE* (Steinegger, Meier, et al.,

615 2019). We only accepted those hits with an *HHblits-probability* ≥ 90% and we re-classified them

616 following the same majority vote approach as previously described. The clusters with no hits

617 remained as the refined set of EUs. We applied a similar refinement approach to the KWP

618 clusters to identify GCs with remote homologies to Pfam protein domains. The KWP HMM

619 profiles were searched against the Pfam *HH-SUITE* database (version 31), using *HHblits*. We

620 accepted hits with a probability ≥ 90% and a target coverage > 60% and removed overlapping

621 domains as described earlier. We moved the KWP with remote homologies to known Pfams to

622 the Known set, and those showing remote homologies to Pfam DUFs to the GUs. The clusters

623 with no hits remained as the refined set of KWP.

## 624 Gene cluster characterization

625 To retrieve the taxonomic composition of our clusters we applied the *MMseqs2 taxonomy*

626 program (version: b43de8b7559a3b45c8e5e9e02cb3023dd339231a), which allows computing

627 the lowest common ancestor through the implementation of the 2bLCA protocol (Hingamp et al.,

628 2013). We searched all cluster genes against UniProtKB (release of January 2018) (UniProt

629 Consortium, 2018) using the following parameters "*-e 1e-05 --cov-mode 0 -c 0.6*". We parsed

630 the results to keep only the hits within 60% of the log10(best-e-value). To retrieve the taxonomic

631 lineages, we used the R package *CHNOSZ* (Dick, 2008). We measured the intra-cluster

632    taxonomic admixture by applying the *entropy.empirical()* function from the *entropy* R package

633    (Hausser & Strimmer, 2008). This function estimates the Shannon entropy based on the

634    different taxonomic annotation frequencies. For each cluster, we also retrieved the cluster

635    consensus taxonomic annotation, which we defined as the taxonomic annotation of the majority

636    of the genes in the cluster.

637    In addition to the taxonomy, we evaluated the clusters' level of darkness and disorder using the

638    Dark Proteome Database (DPD) (Perdigão et al., 2017) as reference. We searched the cluster

639    genes against the DPD, applying the MMseqs2 search program (Steinegger & Soding, 2017)

640    with "*-e 1e-20 --cov-mode 0 -c 0.6*". For each cluster, we then retrieved the mean and the

641    median level of darkness, based on the gene DPD annotations.

## High-quality clusters

643    We defined a subset of high-quality clusters based on the completeness of the cluster genes

644    and their representatives. We identified the minimum required percentage of complete genes

645    per cluster by a broken-stick model (Bennett, 1996) applied to the percentage distribution. Then,

646    we selected the GCs found above the threshold and with a complete representative.

## A set of non-redundant domain architectures

648    We estimated the number of potential domain architectures present in the *Known* GCs taking

649    into account the large proportion of fragmented genes in the metagenomic dataset and that

650    could inflate the number of potential domain architectures. To identify fragments of larger

651    domain architecture, we took into account their topological order in the genes. To reduce the

652    number of comparisons, we calculated the pairwise string cosine distance (q-gram = 3) between

653    domain architectures and discarded the pairs that were too divergent (cosine distance ≥ 0.9).

654    We collapsed a fragmented domain architecture to the larger one when it contained less than

655    75% of complete genes.

## Inference of gene cluster communities

656

657 We aggregated distant homologous GCs into GCCs. The community inference approach

658 combined an all-vs-all HMM gene cluster comparison with Markov Cluster Algorithm (MCL) (van

659 Dongen & Abreu-Goodger, 2012) community identification. We started performing the inference

660 on the Known GCs to use the Pfam DAs as constraints. We aligned the gene cluster HMMs

661 using HHblits (Remmert et al., 2012) (-n 2 -Z 10000000 -B 10000000 -e 1) and we built a

662 homology graph using the cluster pairs with probability ≥ 50% and bidirectional coverage > 60%.

663 We used the ratio between HHblits-bitscore and aligned-columns as the edge weights (Supp.

664 Note 9). We used MCL (van Dongen & Abreu-Goodger, 2012) (v. 12-068) to identify the

665 communities present in the graph. We developed an iterative method to determine the optimal

666 MCL inflation parameter that tries to maximize the relationship of five intra-/inter-community

667 properties: (1) the proportion of MCL communities with one single DA, based on the consensus

668 DAs of the cluster members; (2) the ratio of MCL communities with more than one cluster; (3)

669 the proportion of MCL communities with a PFAM clan entropy equal to 0; (4) the intra-

670 community HHblits-score/Aligned-columns score (normalized by the maximum value); and (5)

671 the number of MCL communities, which should, in the end, reflect the number of non-redundant

672 DAs. We iterated through values ranging from 1.2 to 3.0, with incremental steps of 0.1. During

673 the inference process, some of the GCs became orphans in the graph. We applied a three-step

674 approach to assigning a community membership to these GCs. First, we used less stringent

675 conditions (probability ≥ 50% and coverage >= 40%) to find homologs in the already existing

676 GCCs. Then, we ran a second iteration to find secondary relationships between the newly

677 assigned GCs and the missing ones. Lastly, we created new communities with the remaining

678 GCs. We repeated the whole process with the other categories (KWP, GU and EU), applying

679 the optimal inflation value found for the Known (2.2 for metagenomic and 2.5 for genomic data).

32

# Gene cluster communities validation

680

681 We tested the biological significance of the GCCs using the phylogeny of proteorhodopsin

682 (Boeuf et al., 2015) (PR). We used the proteorhodopsin HMM profiles (Olson et al., 2018) to

683 screen the marine metagenomic datasets using *hmmsearch* (version 3.1b2) (R. D. Finn et al.,

684 2011). We kept the hits with a coverage > 0.4 and e-value <= 1e-5. We removed identical

685 duplicates from the sequences assigned to PR with CD-HIT (W. Li & Godzik, 2006) (v4.6) and

686 cleaned from sequences with less than 100 amino acids. To place the identified PR sequences

687 into the MicRhode (Boeuf et al., 2015) PR tree first, we optimized the initial tree parameters and

688 branch lengths with RAxML (v8.2.12) (Stamatakis, 2014). We used PaPaRA (v2.5) (Berger &

689 Stamatakis, 2012) to incrementally align the query PR sequences against the MicRhode PR

690 reference alignment and *pplacer* (Matsen, Kodner, & Armbrust, 2010) (v1.1.alpha19-0-g807f6f3)

691 to place the sequences into the tree. Finally, we assigned the query PR sequences to the

692 MicRhode PR Superclusters based on the phylogenetic placement. We further investigated the

693 GCs annotated as viral (196 genes, 14 GC) comparing them to the six newly discovered viral

694 PRs (Needham et al., 2019) using Parasail (Daily, 2016) (-a sg_stats_scan_sse2_128_16 -t 8 -c

695 1 -x). As an additional evaluation, we investigated the distributions of standard GCCs and HQ

696 GCCs within ribosomal protein families. We obtained the ribosomal proteins used for the

697 analysis combining the set of 16 ribosomal proteins from Méheust et al. (Méheust et al., 2019)

698 and those contained in the collection of bacterial single-copy genes of Anvi'o (Murat Eren et al.,

699 2015). Also, for the ribosomal proteins, we compared the outcome of our method to the one

700 proposed by Méheust et al. (Méheust et al., 2019) (Supp. Note 9).

# Metagenomic sample selection for downstream analyses

For the subsequent ecological analyses, we selected those metagenomes with a number of genes larger or equal to the first quartile of the distribution of all the metagenomic gene counts. (Supp. Table 13).

# Gene cluster abundance profiles in genomes and metagenomes

We estimated abundance profiles for the metagenomic cluster categories using the read coverage to each predicted gene as a proxy for abundance. We calculated the coverage by mapping the reads against the assembly contigs using the *bwa-mem* algorithm from *BWA mapper* (H. Li & Durbin, 2010). Then, we used *BEDTOOLS* (Quinlan & Hall, 2010), to find the intersection of the gene coordinates to the assemblies, and normalize the per-base coverage by the length of the gene. We calculated the cluster abundance in a sample as the sum of the cluster gene abundances in that sample, and the cluster category abundance in a sample as the sum of the cluster abundances. We obtained the proportions of the different gene cluster categories applying a total-sum-scaling normalization. For the genomic abundance profiles, we used the number of genes in the genomes and normalized by the total gene counts per genome.

# Rate of genomic and metagenomic gene clusters accumulation

We calculated the cumulative number of known and unknown GCs as a function of the number of metagenomes and genomes. For each metagenome count, we generated 1000 random sets, and we calculated the number of GCs and GCCs recovered. For this analysis, we used 1,246 HMP metagenomes and 358 marine metagenomes (242 from TARA and 116 from Malaspina). We repeated the same procedure for the genomic dataset. We removed the singletons from the metagenomic dataset with an abundance smaller than the mode abundance of the singletons

34

724 that got reclassified as good-quality clusters after integrating the GTDB data to minimize the

725 impact of potential spurious singletons. To complement those analyses, we evaluated the

726 coverage of our dataset by searching seven different state-of-the-art databases against our set

727 of metagenomic GC HMM profiles (Supp. Note 12).

728

## 729 Occurrence of gene clusters in the environment

730 We used 1,264 metagenomes from the TARA Oceans, MALASPINA Expedition, OSD2014 and

731 HMP-I/II to explore the properties of the unknown CDS-space in the environment. We applied

732 the Levins Niche Breadth (NB) index (Levins, 1966) to investigate the GCs and GCCs

733 environmental distributions. We removed the GCs and cluster communities with a mean relative

734 abundance < 1e-5. We followed a divide-and-conquer strategy to avoid the computational

735 burden of generating the null-models to test the significance of the distributions owing to the

736 large number of metagenomes and GCs. First, we grouped similar samples based on the gene

737 cluster content using the Bray-Curtis dissimilarity(Bray, Roger Bray, & Curtis, 1957) in

738 combination with the *Dynamic Tree Cut* (Langfelder, Zhang, & Horvath, 2008) R package. We

739 created 100 random datasets picking up one random sample from each group. For each of the

740 100 random datasets, we created 100 random abundance matrices using the *nullmodel* function

741 of the *quasiswap* count method (Miklós & Podani, 2004). Then we calculated the *observed* NB

742 and obtained the 2.5% and 97.5% quantiles based on the randomized sets. We compared the

743 observed and quantile values for each gene cluster and defined it to have a *Narrow distribution*

744 when the *observed* was smaller than the 2.5% quantile and to have a *Broad distribution* when it

745 was larger than the 97.5% quantile. Otherwise, we classified the cluster as *Non-significant*

746 (Salazar et al., 2015). We used a majority voting approach to get a consensus distribution

747 classification based on the ten random datasets.

## Identification of prophages in genomic sequences

We used PhageBoost (https://github.com/ku-cbd/PhageBoost/) to find gene regions in the microbial genomes that result in high viral signals against the overall genome signal. We set the following thresholds to consider a region prophage: minimum of 10 genes, maximum 5 gaps, single-gene probability threshold 0.9. We further smoothed the predictions using Parzen rolling windows of 20 periods and looked at the smoothed probability distribution across the genome. We disregarded regions that had a summed smoothed probability less than 0.5, and those regions that did differ from the overall population of the genes in a genome by using Kruskal–Wallis rank test (p-value 0.001).

## Lineage-specific gene clusters

We used the F1-score developed for AnnoTree (Mendler et al., 2019) to identify the lineage-specific GCs and to which rank they are specific. Following similar criteria to the ones used in Mendler et al. (Mendler et al., 2019), we considered a gene cluster to be lineage-specific if it is present in less than half of all genomes and at least 2 with F1-score > 0.95.

## Phylogenetic conservation of gene clusters

We calculated the phylogenetic conservation (τD) of each gene cluster using the *consenTRAIT* (Martiny et al., 2013) function implemented in the R package *castor* (Martiny et al., 2013). We used a paired Wilcoxon rank-sum test to compare the average τD values for lineage-specific and non-specific GCs.

## Evaluation of the OM-RGC v2 uncharacterized fraction

We integrated the 46,775,154 genes from the second version of the TARA Ocean Microbial Reference Gene Catalog (OM-RGC v2) (Salazar et al., 2019) into our cluster database using

36

770    the same procedure as for the genomic data. We evaluated the uncharacterized fraction and the

771    genes classified into the eggNOG (Huerta-Cepas et al., 2019) category S within the context of

772    our database.

# Augmenting RB-TnSeq experimental data

774    We searched the 37,684 genes of unknown function associated with mutant phenotypes from

775    Price et al. (Price et al., 2018) against our gene cluster profiles. We kept the hits with e-value ≤

776    1e-20 and a query coverage > 60%. Then we filtered the results to keep the hits within 90% of

777    the Log(best-e-value), and we used a majority vote function to retrieve the consensus category

778    for each hit. Lastly, we selected the best-hits based on the smallest e-value and the largest

779    query and target coverage values. We used the fitness values from the RB-TnSeq experiments

780    from Price et al. to identify genes of unknown function that are important for fitness under

781    certain experimental conditions.

# Availability of data and materials

783    The code used for the analyses in the manuscript is available at https://github.com/functional-

784    dark-side/functional-dark-side.github.io/tree/master/scripts. The code to recreate the figures is

785    available at https://github.com/functional-dark-side/vanni_et_al-figures. Detailed descriptions of

786    the different methods and results of this manuscript are available at

787    https://dark.metagenomics.eu. The workflow AGNOSTOS is available at

788    https://github.com/functional-dark-side/agnostos-wf, and its database can be downloaded from

789    https://doi.org/10.6084/m9.figshare.12459056.

790

791

# Acknowledgements

# Contributions

CV, MSS and AF-G performed the analyses and wrote the computational workflow. MS assisted with the clustering and remote homology searches. KS helped with the identification of prophages in genomic sequences. PLB and AB provided feedback and assisted with the ecological analyses. RDF and AM provided feedback and information on the MGnify and Pfam

816    databases. CMD, PS and SGA provided the Malaspina metagenomes. TOD and AME analyzed

817    data in the context of metagenome-assembled genomes. AF-G conceived the study and

818    supervised the work. CV and AF-G wrote the manuscript. All authors read, edited and approved

819    the final manuscript.

## Competing Interests

821    The authors declare no competing interests.

# 822 References

823 Almeida, A., Mitchell, A. L., Boland, M., Forster, S. C., Gloor, G. B., Tarkowska, A., … Finn, R.

824         D. (2019). A new genomic blueprint of the human gut microbiota. *Nature*, *568*(7753),

825         499–504.

826 Almeida, A., Nayfach, S., Boland, M., Strozzi, F., Beracochea, M., Shi, Z. J., … Finn, R. D.

827         (2020). A unified catalog of 204,938 reference genomes from the human gut

828         microbiome. *Nature Biotechnology*. doi:10.1038/s41587-020-0603-3

829 Anantharaman, K., Hausmann, B., Jungbluth, S. P., Kantor, R. S., Lavy, A., Warren, L. A., …

830         Banfield, J. F. (2018). Expanded diversity of microbial groups that shape the

831         dissimilatory sulfur cycle. *The ISME Journal*, *12*(7), 1715–1728.

832 Arnold, F. H. (1998). Design by Directed Evolution. *Accounts of Chemical Research*, *31*(3),

833         125–131.

834 Arnold, F. H. (2018). Directed Evolution: Bringing New Chemistry to Life. *Angewandte Chemie* ,

835         *57*(16), 4143–4148.

836 Bateman, A., Coggill, P., & Finn, D. R. (2010). DUFs: Families in search of function. *Acta*

837         *Crystallographica. Section F, Structural Biology and Crystallization Communications*,

838         *66*(10), 1148–1152.

839 Bennett, K. D. (1996). Determination of the number of zones in a biostratigraphical sequence.

840         *The New Phytologist*, *132*(1), 155–170.

841 Berger, S. A., & Stamatakis, A. (2012). PaPaRa 2.0: a vectorized algorithm for probabilistic

842         phylogeny-aware alignment extension. *Heidelberg Institute for Theoretical Studies,*

843         *Http://Sco. h-Its. Org/Exelixis/Publica Tions. Html. Exelixis-RRDR-2012-2015*. Retrieved

844         from https://cme.h-its.org/exelixis/pubs/Exelixis-RRDR-2012-5.pdf

845 Bernard, G., Pathmanathan, J. S., Lannes, R., Lopez, P., & Bapteste, E. (2018). Microbial Dark

846         Matter Investigations: How Microbial Studies Transform Biological Knowledge and

847    Empirically Sketch a Logic of Scientific Discovery. *Genome Biology and Evolution*, *10*(3),

848    707–715.

849 Bileschi, M. L., Belanger, D., Bryant, D., Sanderson, T., Carter, B., Sculley, D., … Colwell, L. J.

850    (2019). Using Deep Learning to Annotate the Protein Universe (p. 626507).

851    doi:10.1101/626507

852 Bitard-Feildel, T., & Callebaut, I. (2017). Exploring the dark foldable proteome by considering

853    hydrophobic amino acids topology. *Scientific Reports*, *7*, 41425.

854 Boeuf, D., Audic, S., Brillet-Guéguen, L., Caron, C., & Jeanthon, C. (2015). MicRhoDE: a

855    curated database for the analysis of microbial rhodopsin diversity and evolution.

856    *Database: The Journal of Biological Databases and Curation*, *2015*.

857    doi:10.1093/database/bav080

858 Brandenberg, O. F., Fasan, R., & Arnold, F. H. (2017). Exploiting and engineering hemoproteins

859    for abiological carbene and nitrene transfer reactions. *Current Opinion in Biotechnology*,

860    *47*, 102–111.

861 Bray, J. R., Roger Bray, J., & Curtis, J. T. (1957). An Ordination of the Upland Forest

862    Communities of Southern Wisconsin. *Ecological Monographs*, Vol. 27, pp. 325–349.

863    doi:10.2307/1942268

864 Breitwieser, F. P., Pertea, M., Zimin, A., & Salzberg, S. L. (2019). Human contamination in

865    bacterial genomes has created thousands of spurious proteins. *Genome Research*.

866    doi:10.1101/gr.245373.118

867 Broder, A. Z. (1997). On the resemblance and containment of documents. *Proceedings.*

868    *Compression and Complexity of SEQUENCES 1997 (Cat. No. 97TB100171)*, 21–29.

869    IEEE.

870 Brown, C. T., Hug, L. A., Thomas, B. C., Sharon, I., Castelle, C. J., Singh, A., … Banfield, J. F.

871    (2015). Unusual biology across a group comprising more than 15% of domain Bacteria.

872    *Nature*, *523*(7559), 208–211.

873    Brum, J. R., Cesar Ignacio-Espinoza, J., Kim, E.-H., Trubl, G., Jones, R. M., Roux, S., …

874        Sullivan, M. B. (2016). Illuminating structural proteins in viral "dark matter" with

875        metaproteomics. *Proceedings of the National Academy of Sciences of the United States*

876        *of America*, *113*(9), 2436–2441.

877    Buttigieg, L. P., Hankeln, W., Kostadinov, I., Kottmann, R., Yilmaz, P., Duhaime, B. M., &

878        Gl??ckner, O. F. (2013). Ecogenomic Perspectives on Domains of Unknown Function:

879        Correlation-Based Exploration of Marine Metagenomes. *PloS One*, *8*(3).

880    Carradec, Q., Pelletier, E., Da Silva, C., Alberti, A., Seeleuthner, Y., Blanc-Mathieu, R., …

881        Wincker, P. (2018). A global ocean atlas of eukaryotic genes. *Nature Communications*,

882        *9*(1), 373.

883    Chafee, M., Fernàndez-Guerra, A., Buttigieg, P. L., Gerdts, G., Eren, A. M., Teeling, H., &

884        Amann, R. I. (2018). Recurrent patterns of microdiversity in a temperate coastal marine

885        environment. *The ISME Journal*, *12*(1), 237–252.

886    Chen, I.-M. A., Chu, K., Palaniappan, K., Pillay, M., Ratner, A., Huang, J., … Kyrpides, N. C.

887        (2019). IMG/M v.5.0: an integrated data management and comparative analysis system

888        for microbial genomes and microbiomes. *Nucleic Acids Research*, *47*(D1), D666–D677.

889    Cross, K. L., Campbell, J. H., Balachandran, M., Campbell, A. G., Cooper, S. J., Griffen, A., …

890        Podar, M. (2019). Targeted isolation and cultivation of uncultivated bacteria by reverse

891        genomics. *Nature Biotechnology*, *37*(11), 1314–1321.

892    Csardi, G., & Nepusz, T. (2006). The igraph software package for complex network research.

893        *InterJournal*, p. 1695. Retrieved from http://igraph.org

894    Daily, J. (2016). Parasail: SIMD C library for global, semi-global, and local pairwise sequence

895        alignments. *BMC Bioinformatics*, *17*(1), 81–81.

896    Delmont, T. O., Gaia, M., Hinsinger, D. D., Fremont, P., Guerra, A. F., Murat Eren, A., … Jaillon,

897        O. (2020). Functional repertoire convergence of distantly related eukaryotic plankton

898        lineages revealed by genome-resolved metagenomics (p. 2020.10.15.341214).

899      doi:10.1101/2020.10.15.341214

900  Deorowicz, S., Debudaj-Grabysz, A., & Gudyś, A. (2016). FAMSA: Fast and accurate multiple

901      sequence alignment of huge protein families. *Scientific Reports*, *6*(1), 33964–33964.

902  Dick, J. M. (2008). Calculation of the relative metastabilities of proteins using the CHNOSZ

903      software package. *Geochemical Transactions*, *9*, 10.

904  Duarte, C. M. (2015). Seafaring in the 21St Century: The Malaspina 2010 Circumnavigation

905      Expedition. *Limnology and Oceanography Bulletin*, *24*(1), 11–14.

906  Dubinkina, V. B., Ischenko, D. S., Ulyantsev, V. I., Tyakht, A. V., & Alexeev, D. G. (2016).

907      Assessment of k-mer spectrum applicability for metagenomic dissimilarity analysis. *BMC*

908      *Bioinformatics*, Vol. 17. doi:10.1186/s12859-015-0875-7

909  Eberhardt, R. Y., Haft, D. H., Punta, M., Martin, M., O'Donovan, C., & Bateman, A. (2012).

910      AntiFam: a tool to help identify spurious ORFs in protein annotation. *Database*, *2012*(0),

911      bas003–bas003.

912  Edwards, R. A., Vega, A. A., Norman, H. M., Ohaeri, M., Levi, K., Dinsdale, E. A., … Dutilh, B.

913      E. (2019). Global phylogeography and ancient evolution of the widespread human gut

914      virus crAssphage. *Nature Microbiology*, *4*(10), 1727–1736.

915  Eloe-Fadrosh, E. A., Paez-Espino, D., Jarett, J., Dunfield, P. F., Hedlund, B. P., Dekas, A. E., …

916      Ivanova, N. N. (2016). Global metagenomic survey reveals a new bacterial candidate

917      phylum in geothermal springs. *Nature Communications*, *7*, 10476.

918  Finn, R. D., Clements, J., & Eddy, S. R. (2011). HMMER web server: interactive sequence

919      similarity searching. *Nucleic Acids Research*, *39*(suppl), W29–W37.

920  Finn, Robert D., Coggill, P., Eberhardt, R. Y., Eddy, S. R., Mistry, J., Mitchell, A. L., … Bateman,

921      A. (2016). The Pfam protein families database: towards a more sustainable future.

922      *Nucleic Acids Research*, *44*(D1), D279–D285.

923  Francino, M. P. (2012). The ecology of bacterial genes and the survival of the new. *International*

924      *Journal of Evolutionary Biology*, *2012*, 394026.

925    Franzosa, E. A., McIver, L. J., Rahnavard, G., Thompson, L. R., Schirmer, M., Weingart, G., …

926          Huttenhower, C. (2018). Species-level functional profiling of metagenomes and

927          metatranscriptomes. *Nature Methods*, *15*(11), 962–968.

928    Habchi, J., Tompa, P., Longhi, S., & Uversky, V. N. (2014). Introducing protein intrinsic disorder.

929          *Chemical Reviews*, *114*(13), 6561–6588.

930    Hanson, A. D., Pribat, A., Waller, J. C., & Crécy-Lagard, V. de. (2010). 'Unknown'proteins and

931          'orphan'enzymes: the missing half of the engineering parts list--and how to find it.

932          *Biochemical Journal*, *425*(1), 1–11.

933    Hausser, J., & Strimmer, K. (2008). Entropy inference and the James-Stein estimator, with

934          application to nonlinear gene association networks. Retrieved from

935          http://arxiv.org/abs/0811.3579

936    Heffernan, B., Murphy, C. D., & Casey, E. (2009). Comparison of planktonic and biofilm cultures

937          of Pseudomonas fluorescens DSM 8341 cells grown on fluoroacetate. *Applied and*

938          *Environmental Microbiology*, *75*(9), 2899–2907.

939    Heinzinger, M., Elnaggar, A., Wang, Y., Dallago, C., Nechaev, D., Matthes, F., & Rost, B.

940          (2019). Modeling aspects of the language of life through transfer-learning protein

941          sequences. *BMC Bioinformatics*, *20*(1), 723.

942    Hingamp, P., Grimsley, N., Acinas, S. G., Clerissi, C., Subirana, L., Poulain, J., … Ogata, H.

943          (2013). Exploring nucleo-cytoplasmic large DNA viruses in Tara Oceans microbial

944          metagenomes. *The ISME Journal*, *7*(9), 1678–1695.

945    Höps, W., Jeffryes, M., & Bateman, A. (2018). Gene Unprediction with Spurio: A tool to identify

946          spurious protein sequences. *F1000Research*, *7*, 261.

947    Huerta-Cepas, J., Forslund, K., Coelho, L. P., Szklarczyk, D., Jensen, L. J., von Mering, C., &

948          Bork, P. (2017). Fast Genome-Wide Functional Annotation through Orthology

949          Assignment by eggNOG-Mapper. *Molecular Biology and Evolution*, *34*(8), 2115–2122.

950    Huerta-Cepas, J., Szklarczyk, D., Heller, D., Hernández-Plaza, A., Forslund, S. K., Cook, H., …

951       Bork, P. (2019). eggNOG 5.0: a hierarchical, functionally and phylogenetically annotated

952       orthology resource based on 5090 organisms and 2502 viruses. *Nucleic Acids*

953       *Research*, *47*(D1), D309–D314.

954   Hug, L. A., Baker, B. J., Anantharaman, K., Brown, C. T., Probst, A. J., Castelle, C. J., …

955       Banfield, J. F. (2016). A new view of the tree of life. *Nature Microbiology*, *1*, 16048.

956   Hyatt, D., Chen, G.-L., LoCascio, P. F., Land, M. L., Larimer, F. W., & Hauser, L. J. (2010).

957       Prodigal: prokaryotic gene recognition and translation initiation site identification. *BMC*

958       *Bioinformatics*, *11*(1), 119–119.

959   Ivanova, N. N., Schwientek, P., Tripp, H. J., Rinke, C., Pati, A., Huntemann, M., … Rubin, E. M.

960       (2014). Stop codon reassignments in the wild. *Science*, *344*(6186), 909–913.

961   Jaroszewski, L., Li, Z., Krishna, S. S., Bakolitsa, C., Wooley, J., Deacon, M. A., … Godzik, A.

962       (2009). Exploration of uncharted regions of the protein universe. *PLoS Biology*, *7*(9).

963   Jehl, P., Sievers, F., & Higgins, D. G. (2015). OD-seq: outlier detection in multiple sequence

964       alignments. *BMC Bioinformatics*, *16*(1), 269–269.

965   Kopf, A., Bicak, M., Kottmann, R., Schnetzer, J., Kostadinov, I., Lehmann, K., … Glöckner, F. O.

966       (2015). The ocean sampling day consortium. *GigaScience*, *4*, 27.

967   Köster, J. (2018). Reproducible data analysis with Snakemake. *F1000Research*, *7*. Retrieved

968       from http://www.dodsc.tu-

969       dortmund.de/cms/Medienpool/files/002_Kolloquium/09_Koester.pdf

970   La Cono, V., La Spada, G., Arcadi, E., Placenti, F., Smedile, F., Ruggeri, G., … Yakimov, M. M.

971       (2013). Partaking of Archaea to biogeochemical cycling in oxygen-deficient zones of

972       meromictic saline Lake Faro (Messina, Italy). *Environmental Microbiology*, *15*(6), 1717–

973       1733.

974   Langfelder, P., Zhang, B., & Horvath, S. (2008). Defining clusters from a hierarchical cluster

975       tree: the Dynamic Tree Cut package for R. *Bioinformatics* , *24*(5), 719–720.

976   Levins, R. (1966). THE STRATEGY OF MODEL BUILDING IN POPULATION BIOLOGY.

bioRxiv preprint doi: https://doi.org/10.1101/2020.06.30.180448; this version posted February 18, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

977     *American Scientist*, *54*(4), 421–431.

978     Li, H., & Durbin, R. (2010). Fast and accurate long-read alignment with Burrows–Wheeler

979         transform. *Bioinformatics* , *26*(5), 589–595.

980     Li, W., & Godzik, A. (2006). Cd-hit: a fast program for clustering and comparing large sets of

981         protein or nucleotide sequences. *Bioinformatics* , *22*(13), 1658–1659.

982     Liu, X. L. (2017). Deep Recurrent Neural Network for Protein Function Prediction from

983         Sequence (p. 103994). doi:10.1101/103994

984     Lloyd-Price, J., Mahurkar, A., Rahnavard, G., Crabtree, J., Orvis, J., Hall, A. B., … Huttenhower,

985         C. (2017). Strains, functions and dynamics in the expanded Human Microbiome Project.

986         *Nature*, *550*(7674), 61–66.

987     Lobb, B., Kurtz, D. A., Moreno-Hagelsieb, G., & Doxey, A. C. (2015). Remote homology and the

988         functions of metagenomic dark matter. *Frontiers in Genetics*, *6*(JUL), 1–12.

989     Ma, Y., Madupu, R., Karaoz, U., Nossa, C. W., Yang, L., Yooseph, S., … Pei, Z. (2014). Human

990         papillomavirus community in healthy persons, defined by metagenomics analysis of

991         human microbiome project shotgun sequencing data sets. *Journal of Virology*, *88*(9),

992         4786–4797.

993     Martiny, A. C., Treseder, K., & Pusch, G. (2013). Phylogenetic conservatism of functional traits

994         in microorganisms. *The ISME Journal*, *7*(4), 830–838.

995     Matsen, F. A., Kodner, R. B., & Armbrust, E. V. (2010). pplacer: linear time maximum-likelihood

996         and Bayesian phylogenetic placement of sequences onto a fixed reference tree. *BMC*

997         *Bioinformatics*, *11*, 538.

998     Méheust, R., Burstein, D., Castelle, C. J., & Banfield, J. F. (2019). The distinction of CPR

999         bacteria from other bacteria based on protein family content. *Nature Communications*,

1000        *10*(1), 4173.

1001    Mendler, K., Chen, H., Parks, D. H., Lobb, B., Hug, L. A., & Doxey, A. C. (2019). AnnoTree:

1002        visualization and exploration of a functionally annotated microbial tree of life. *Nucleic*

1003    *Acids Research*, *47*(9), 4442–4448.

1004    Miklós, I., & Podani, J. (2004). RANDOMIZATION OF PRESENCE–ABSENCE MATRICES:

1005        COMMENTS AND NEW ALGORITHMS. *Ecology*, Vol. 85, pp. 86–92. doi:10.1890/03-

1006        0101

1007    Mitchell, A. L., Almeida, A., Beracochea, M., Boland, M., Burgin, J., Cochrane, G., … Finn, R. D.

1008        (2020). MGnify: the microbiome analysis resource in 2020. *Nucleic Acids Research*,

1009        *48*(D1), D570–D578.

1010    Muller, E. E. L. (2019). Determining Microbial Niche Breadth in the Environment for Better

1011        Ecosystem Fate Predictions. *MSystems*, *4*(3). doi:10.1128/mSystems.00080-19

1012    Murat Eren, A., Esen, Ö. C., Quince, C., Vineis, J. H., Morrison, H. G., Sogin, M. L., & Delmont,

1013        T. O. (2015). Anvi'o: an advanced analysis and visualization platform for 'omics data.

1014        *PeerJ*, *3*, e1319.

1015    NCBI Resource Coordinators. (2018). Database resources of the National Center for

1016        Biotechnology Information. *Nucleic Acids Research*, *46*(D1), D8–D13.

1017    Needham, D. M., Yoshizawa, S., Hosaka, T., Poirier, C., Choi, C. J., Hehenberger, E., …

1018        Worden, A. Z. (2019). A distinct lineage of giant viruses brings a rhodopsin photosystem

1019        to unicellular marine predators. *Proceedings of the National Academy of Sciences of the*

1020        *United States of America*, *116*(41), 20574–20583.

1021    Olson, D. K., Yoshizawa, S., Boeuf, D., Iwasaki, W., & DeLong, E. F. (2018). Proteorhodopsin

1022        variability and distribution in the North Pacific Subtropical Gyre. *The ISME Journal*,

1023        *12*(4), 1047–1060.

1024    Pachiadaki, M. G., Brown, J. M., Brown, J., Bezuidt, O., Berube, P. M., Biller, S. J., …

1025        Stepanauskas, R. (2019). Charting the Complexity of the Marine Microbiome through

1026        Single-Cell Genomics. *Cell*, *179*(7), 1623-1635.e11.

1027    Parks, D. H., Chuvochina, M., Waite, D. W., Rinke, C., Skarshewski, A., Chaumeil, P.-A., &

1028        Hugenholtz, P. (2018). A standardized bacterial taxonomy based on genome phylogeny

1029    substantially revises the tree of life. *Nature Biotechnology*, *36*(10), 996–1004.

1030    Pasolli, E., Asnicar, F., Manara, S., Zolfo, M., Karcher, N., Armanini, F., … Segata, N. (2019).

1031    Extensive Unexplored Human Microbiome Diversity Revealed by Over 150,000

1032    Genomes from Metagenomes Spanning Age, Geography, and Lifestyle. *Cell*, *176*(3),

1033    649-662.e20.

1034    Perdigão, N., Rosa, A. C., & O'Donoghue, S. I. (2017). The Dark Proteome Database. *BioData*

1035    *Mining*, *10*(1), 1–11.

1036    Price, M. N., Wetmore, K. M., Waters, R. J., Callaghan, M., Ray, J., Liu, H., … Deutschbauer, A.

1037    M. (2018). Mutant phenotypes for thousands of bacterial genes of unknown function.

1038    *Nature*, *557*(7706), 503–509.

1039    Quince, C., Walker, A. W., Simpson, J. T., Loman, N. J., & Segata, N. (2017). Shotgun

1040    metagenomics, from sampling to analysis. *Nature Biotechnology*, *35*(9), 833–844.

1041    Quinlan, A. R., & Hall, I. M. (2010). BEDTools: a flexible suite of utilities for comparing genomic

1042    features. *Bioinformatics* , *26*(6), 841–842.

1043    Remmert, M., Biegert, A., Hauser, A., & Söding, J. (2012). HHblits: Lightning-fast iterative

1044    protein sequence searching by HMM-HMM alignment. *Nature Methods*, *9*(2), 173–175.

1045    Rinke, C., Schwientek, P., Sczyrba, A., Ivanova, N. N., Anderson, I. J., Cheng, J.-F., … Woyke,

1046    T. (2013). Insights into the phylogeny and coding potential of microbial dark matter.

1047    *Nature*, *499*(7459), 431–437.

1048    Rost, B. (1999). Twilight zone of protein sequence alignments. *Protein Engineering, Design &*

1049    *Selection: PEDS*, *12*(2), 85–94.

1050    Roumpeka, D. D., Wallace, R. J., Escalettes, F., Fotheringham, I., & Watson, M. (2017). A

1051    Review of Bioinformatics Tools for Bio-Prospecting from Metagenomic Sequence Data.

1052    *Frontiers in Genetics*, *8*, 23.

1053    Rusch, D. B., Halpern, A. L., Sutton, G., Heidelberg, K. B., Williamson, S., Yooseph, S., …

1054    Venter, J. C. (2007). The Sorcerer II Global Ocean Sampling Expedition: Northwest

1055    Atlantic through Eastern Tropical Pacific. *PLoS Biology*, *5*(3), 1–34.

1056  Salazar, G., Cornejo-Castillo, F. M., Borrull, E., Díez-Vives, C., Lara, E., Vaqué, D., … Acinas,

1057    S. G. (2015). Particle-association lifestyle is a phylogenetically conserved trait in

1058    bathypelagic prokaryotes. *Molecular Ecology*, *24*(22), 5692–5706.

1059  Salazar, G., Paoli, L., Alberti, A., Huerta-Cepas, J., Ruscheweyh, H.-J., Cuenca, M., …

1060    Sunagawa, S. (2019). Gene Expression Changes and Community Turnover Differentially

1061    Shape the Global Ocean Metatranscriptome. *Cell*, *179*(5), 1068-1083.e21.

1062  Sanger, F., Nicklen, S., & Coulson, A. R. (1977). DNA sequencing with chain-terminating

1063    inhibitors. *Proceedings of the National Academy of Sciences of the United States of*

1064    *America*, *74*(12), 5463–5467.

1065  Sberro, H., Fremin, B. J., Zlitni, S., Edfors, F., Greenfield, N., Snyder, M. P., … Bhatt, A. S.

1066    (2019). Large-Scale Analyses of Human Microbiomes Reveal Thousands of Small,

1067    Novel Genes. *Cell*, *178*(5), 1245-1259.e14.

1068  Scales, B. S., Dickson, R. P., LiPuma, J. J., & Huffnagle, G. B. (2014). Microbiology, genomics,

1069    and clinical significance of the Pseudomonas fluorescens species complex, an

1070    unappreciated colonizer of humans. *Clinical Microbiology Reviews*, *27*(4), 927–948.

1071  Skewes-Cox, P., Sharpton, T. J., Pollard, K. S., & DeRisi, J. L. (2014). Profile hidden Markov

1072    models for the detection of viruses within metagenomic sequence data. *PloS One*, *9*(8),

1073    e105067.

1074  Spang, A., Saw, J. H., Jørgensen, S. L., Zaremba-Niedzwiedzka, K., Martijn, J., Lind, A. E., …

1075    Ettema, T. J. G. (2015). Complex archaea that bridge the gap between prokaryotes and

1076    eukaryotes. *Nature*, *521*(7551), 173–179.

1077  Stamatakis, A. (2014). RAxML version 8: a tool for phylogenetic analysis and post-analysis of

1078    large phylogenies. *Bioinformatics* , *30*(9), 1312–1313.

1079  Steinegger, M., Meier, M., Mirdita, M., Vöhringer, H., Haunsberger, S. J., & Söding, J. (2019).

1080    HH-suite3 for fast remote homology detection and deep protein annotation. *BMC*

1081    *Bioinformatics*, *20*(1), 473.

1082    Steinegger, M., Mirdita, M., & Söding, J. (2019). Protein-level assembly increases protein

1083        sequence recovery from metagenomic samples manyfold. *Nature Methods*, *16*(7), 603–

1084        606.

1085    Steinegger, M., & Salzberg, S. L. (2020). Terminating contamination: large-scale search

1086        identifies more than 2,000,000 contaminated entries in GenBank. *Genome Biology*,

1087        *21*(1), 115.

1088    Steinegger, M., & Soding, J. (2017). MMseqs2 enables sensitive protein sequence searching for

1089        the analysis of massive data sets. *Nature Biotechnology*, *advance on*.

1090        doi:10.1038/nbt.3988

1091    Steinegger, M., & Söding, J. (2018). Clustering huge protein sequence sets in linear time.

1092        *Nature Communications*, *9*(1), 2542.

1093    Sunagawa, S., Coelho, L. P., Chaffron, S., Kultima, J. R., Labadie, K., Salazar, G., … Bork, P.

1094        (2015). Ocean plankton. Structure and function of the global ocean microbiome.

1095        *Science*, *348*(6237), 1261359.

1096    The UniProt Consortium. (2017). UniProt: the universal protein knowledgebase. *Nucleic Acids*

1097        *Research*, *45*(D1), D158–D169.

1098    Thomas, A. M., & Segata, N. (2019). Multiple levels of the unknown in microbiome research.

1099        *BMC Biology*, *17*(1), 48.

1100    Titus Brown, C., Moritz, D., O'Brien, M. P., Reidl, F., Reiter, T., & Sullivan, B. D. (2018).

1101        Exploring neighborhoods in large metagenome assembly graphs reveals hidden

1102        sequence diversity (p. 462788). doi:10.1101/462788

1103    UniProt Consortium, T. (2018). UniProt: the universal protein knowledgebase. *Nucleic Acids*

1104        *Research*, *46*(5), 2699.

1105    van Dongen, S., & Abreu-Goodger, C. (2012). Using MCL to Extract Clusters from Networks. In

1106        J. van Helden, A. Toussaint, & D. Thieffry (Eds.), *Bacterial Molecular Networks: Methods*

1107    *and Protocols* (pp. 281–295). New York, NY: Springer New York.

1108    Vanhoutreve, R., Kress, A., Legrand, B., Gass, H., Poch, O., & Thompson, J. D. (2016). LEON-

1109        BIS: multiple alignment evaluation of sequence neighbours using a Bayesian inference

1110        system. *BMC Bioinformatics*, *17*(1), 271–271.

1111    Vorobev, A., Dupouy, M., Carradec, Q., Delmont, T., Annamale, A., Wincker, P., & Pelletier, E.

1112        (2020). Transcriptome reconstruction and functional analysis of eukaryotic marine

1113        plankton communities via high-throughput metagenomics and metatranscriptomics.

1114        *Genome Research*. doi:10.1101/gr.253070.119

1115    Wyman, S. K., Avila-Herrera, A., Nayfach, S., & Pollard, K. S. (2018). A most wanted list of

1116        conserved microbial protein families with no known domains. *PloS One*, *13*(10),

1117        e0205749.

1118    Yooseph, S., Li, W., & Sutton, G. (2008). Gene identification and protein classification in

1119        microbial metagenomic sequence data via incremental clustering. *BMC Bioinformatics*,

1120        *9*, 1–13.

1121    Yooseph, S., Sutton, G., Rusch, B. D., Halpern, L. A., Williamson, J. S., Remington, K., …

1122        Venter, C. J. (2007). The Sorcerer II global ocean sampling expedition: Expanding the

1123        universe of protein families. *PLoS Biology*, *5*(3), 0432–0466.

1124    Žure, M., Fernandez-Guerra, A., Munn, C. B., & Harder, J. (2017). Geographic distribution at

1125        subspecies resolution level: closely related Rhodopirellula species in European coastal

1126        sediments. *The ISME Journal*, *11*(2), 478–489.