# Light into the darkness: Unifying the known and unknown coding sequence space in microbiome analyses

Chiara Vanni[1,2], Matthew S. Schechter[1,3], Silvia G. Acinas[4], Albert Barberán[5], Pier Luigi Buttigieg[6], Emilio O. Casamayor[7], Tom O. Delmont[8], Carlos M. Duarte[9], A. Murat Eren[3,10], Robert D. Finn[11], Renzo Kottmann[1], Alex Mitchell[11], Pablo Sanchez[4], Kimmo Siren[12], Martin Steinegger[13,14], Frank Oliver Glöckner[15,16,2], Antonio Fernandez-Guerra[1,17] *

## Affiliations

1 Microbial Genomics and Bioinformatics Research Group, Max Planck Institute for Marine Microbiology, Celsiusstraße 1, 28359, Bremen, Germany

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

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

4 Department of Marine Biology and Oceanography, Institut de Ciènces del Mar, CSIC, Barcelona, Spain.

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

6 Alfred Wegener Institute, Helmholtz Centre for Polar and Marine Research, Am Handelshafen 12, 27570 Bremerhaven, Germany

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

8 Génomique Métabolique, Genoscope, Institut François Jacob, CEA, CNRS, Univ Evry, Université Paris-Saclay, 91057 Evry, France

9 Red Sea Research Centre (RSRC) and Computational Bioscience Research Center (CBRC), King Abdullah University of Science and Technology, Thuwal 23955, Saudi Arabia

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

11 European Molecular Biology Laboratory, European Bioinformatics Institute (EMBL-EBI), Wellcome Genome Campus, Hinxton, Cambridge CB10 1SD, UK

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

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

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

15 University of Bremen, MARUM, Leobener Str. 8, 28359 Bremen, Germany
Life Sciences and Chemistry, Campus Ring 1, 28759 Bremen, Germany

16 Computing Center, Helmholtz Center for Polar and Marine Research, Am Handelshafen 12, 27570 Bremerhaven, Germany

38    17 Lundbeck GeoGenetics Centre, The Globe Institute, University of Copenhagen, 1350
39    Copenhagen, Denmark
40
41    *Corresponding author: Antonio Fernandez-Guerra, antonio.fernandez-guerra@sund.ku.dk
42

# Abstract

Bridging the gap between the known and the unknown coding sequence space is one of the biggest challenges in molecular biology today. This challenge is especially extreme in microbiome analyses where between 40% and 60% of the coding sequences detected are of unknown function, and ignoring this fraction limits our understanding of microbial systems. Discarding the uncharacterized fraction is not an option anymore. Here, we present an in-depth exploration of the microbial unknown fraction through the lenses of a conceptual framework and a computational workflow we developed to unify the microbial known and unknown coding sequence space. Our approach partitions the coding sequence space in gene clusters and contextualizes them with genomic and environmental information. We analyzed 415,971,742 genes predicted from 1,749 metagenomes and 28,941 bacterial and archaeal genomes, putting into perspective the extent of the unknown fraction, its diversity, and its relevance in a genomic and environmental context. With the identification of a target gene of unknown function for antibiotic resistance, we demonstrate how a contextualized unknown coding sequence space provides a robust framework for the generation 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 how microbes shape life on Earth[1–7], thus bringing about a golden age of microbial genomics. An ever-increasing number of genomes and metagenomes are unlocking uncharted regions of microbial diversity[1,8,9], providing new perspectives on the evolution of life[10,11]. However, our rapidly growing inventories of new genes have a glaring issue: between 40% and 60% cannot be assigned to a known function[12–15]. Current analytical approaches for genomic and metagenomic data[16–20] generally do not include this uncharacterized fraction in downstream analyses, constraining their results to conserved pathways and housekeeping functions[17]. This inability to handle shades of the unknown is an immense impediment to realizing the potential for discovery of microbial genomics and microbiology at large[12,21].

Predicting function from traditional single sequence similarity appears to have yielded all it can[22–24], thus several groups have attempted to resolve gene function by other means. Such efforts include combining biochemistry and crystallography[25]; using environmental co-occurrence[26]; by grouping those genes into evolutionarily related families[27–30]; using remote homologies[31,32]; or more recently using deep learning approaches[33,34]. In 2018, Price et al.[13] developed a high-throughput experimental pipeline that provides mutant phenotypes for thousands of bacterial

3

75    genes of unknown function being one of the most promising methods to tackle the unknown.

76    Despite their promise, experimental methods are labor-intensive and require novel computational

77    methods that could bridge the existing gap between the known and unknown coding sequence

78    space (CDS-space).

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

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

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

82    resource to get a better understanding of the unknown functional fraction and boost the current

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

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

85    that (1) the extent of the unknown fraction is smaller than expected, (2) that the diversity of gene

86    clusters in the unknown fraction is higher than in the known fraction, and (3) that the unknown

87    fraction is phylogenetically more conserved and is predominantly lineage-specific at the species

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

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

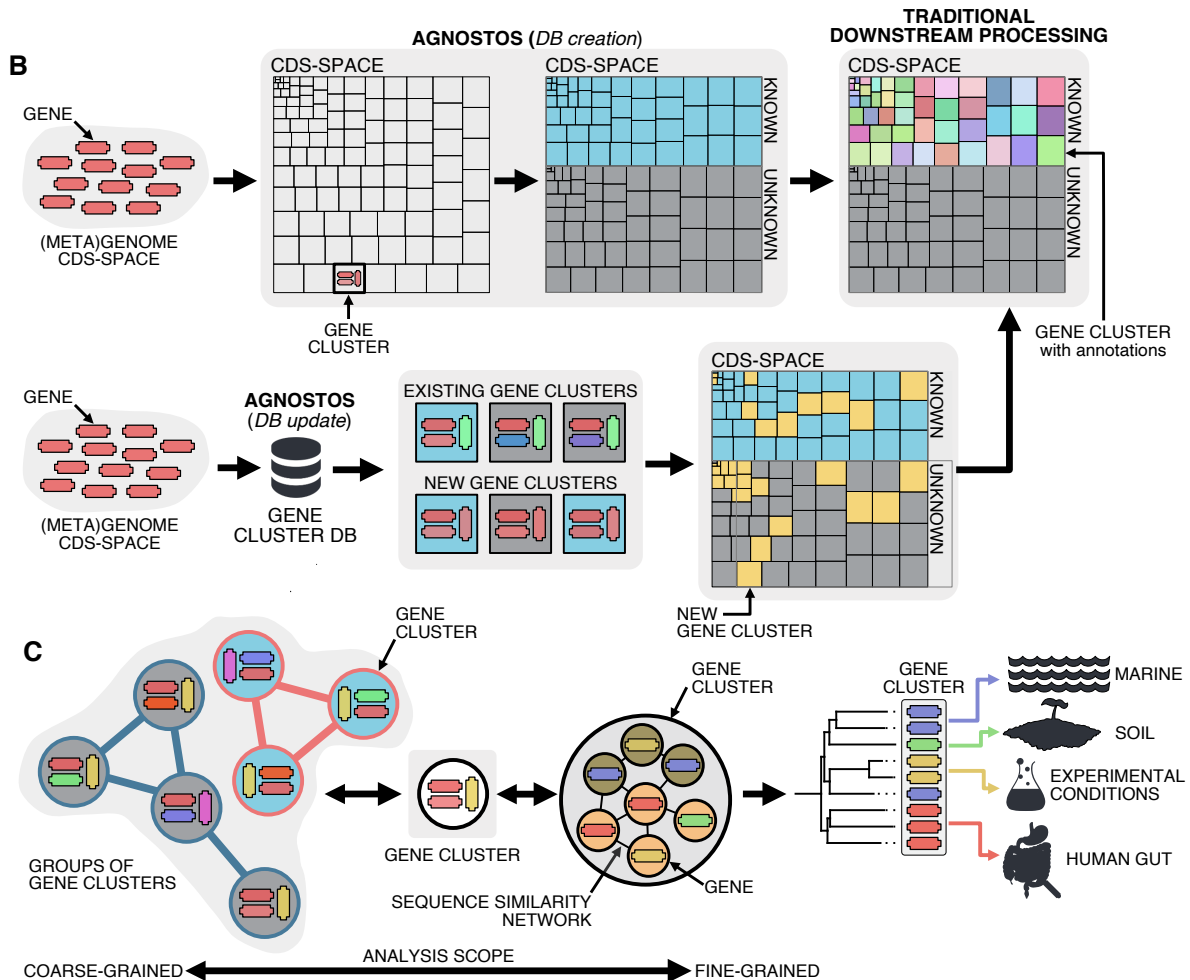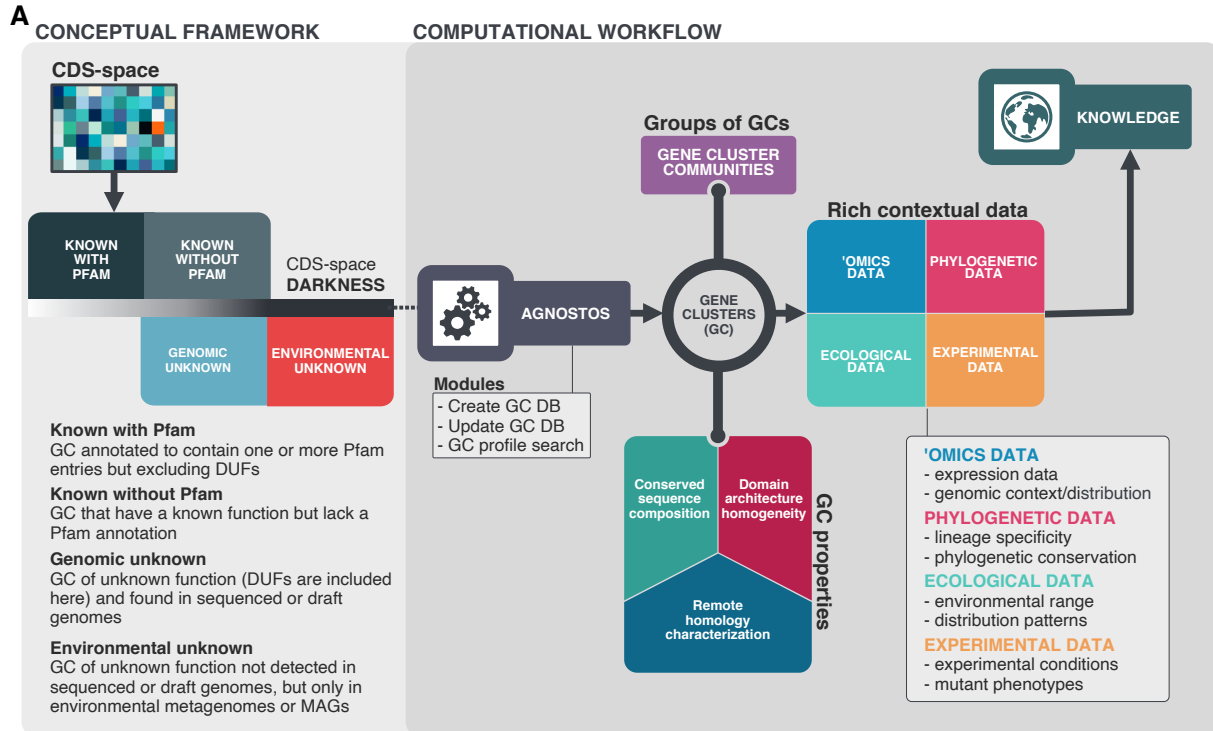90    hypothesis-driven molecular investigations.

91

# Results

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

We created the conceptual and technical foundations to unify the known and unknown CDS-space and provide a practical solution to one of the most significant ongoing challenges in microbiome analyses. 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[35]. 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). GCs produce a structured CDS-space reducing its complexity (Fig. 1B), are independent of the known and unknown fraction, are conserved across environments and organisms, and 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**

KNOWN WITH PFAM | KNOWN WITHOUT PFAM

CDS-space **DARKNESS**

GENOMIC UNKNOWN | ENVIRONMENTAL UNKNOWN

**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

**Groups of GCs**
GENE CLUSTER COMMUNITIES

AGNOSTOS

GENE CLUSTERS (GC)

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

**Rich contextual data**

'OMICS DATA | PHYLOGENETIC DATA
ECOLOGICAL DATA | EXPERIMENTAL DATA

KNOWLEDGE

**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**

GENE

(META)GENOME CDS-SPACE

**AGNOSTOS (**DB creation**)**

CDS-SPACE

GENE CLUSTER

CDS-SPACE

KNOWN | UNKNOWN

**TRADITIONAL DOWNSTREAM PROCESSING**

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**

GROUPS OF GENE CLUSTERS

GENE CLUSTER

GENE CLUSTER

SEQUENCE SIMILARITY NETWORK

GENE CLUSTER

GENE

GENE CLUSTER

MARINE

SOIL

EXPERIMENTAL CONDITIONS

HUMAN GUT

COARSE-GRAINED ← ANALYSIS SCOPE → FINE-GRAINED

113

6

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

132

133 Driven by the concepts defined in the conceptual framework, we developed AGNOSTOS, a

134 computational workflow that infers, validates, refines, and classifies GCs in the four proposed

135 categories (Fig. 1A; Fig. 1B; Supp. Fig 1). AGNOSTOS provides two operational modules (*DB*

136 *creation* and *DB update*) to produce GCs with a highly conserved intra-homogeneous structure

137 (Fig. 1B), both in terms of sequence similarity and domain architecture homogeneity; it exhausts

138 any existing homology to known genes and provides a proper delimitation of the unknown CDS-

139 space before classifying each GC in one of the four categories. In the last step, we decorate each

140 GC with a rich collection of contextual data that we compile from different sources, or that we

141 generate by analyzing the GC contents in different contexts (Fig. 1A). For each GC, we also offer

142 several products that can be used for analytical purposes like improved representative

143 sequences, consensus sequences, sequence profiles for MMseqs2[36] and HHblits[37], or the GC

144 members as a sequence similarity network (see Online Methods). To complement the collection,

145 we also provide a subset of what we define as *high-quality* GCs. The defining criteria are (1) the

146 representative is a complete gene and (2) more than one-third of genes within a GC are complete

147 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 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[38] 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% of them), which agrees with recent findings on novel small proteins from metagenomes[39].

The KWP category contains the largest proportion of incomplete ORFs (Supp. Table 3), impeding the detection and assignment of Pfam domains. But it also incorporates sequences with an unusual amino acid composition that has homology to proteins with high levels of disorder in the DPD database[40] and that have characteristic functions of the intrinsically disordered proteins[41] (IDP) like cellular processes and signaling as predicted by eggNOG annotations (Supp. Table 4). As part of the workflow, each GC is complemented with a rich set of information, as shown in Fig 1A (Supp. Table 5; Supp Note 6).

# Beyond the twilight zone, communities of gene clusters

The method we developed to group GCs in gene cluster communities (GCCs) (Fig. 2A) reduced the final collection of GCs by 87%, producing 673,601 GCCs (Fig. 2B; Supp. Note 7). We validated the ability of our approach in capturing remote homologies between related GCs using two well-known gene families present in our environmental datasets, proteorhodopsins[42] and bacterial ribosomal proteins[43]. In our dataset, 64 GCs (12,184 genes) and 3 GCCs (Supp Note 8) contained sequences classified as proteorhodopsin (PR). One *Known* GCC contained 99% of the PR annotated genes (Fig. 2C), with the only exception of 85 genes taxonomically annotated as viral and assigned to the *PR Supercluster I*[44] enclosed in two GU communities (five GU gene clusters; Supp Note 8). For the ribosomal proteins, the results were not so satisfactory. We identified 1,843 GCs (781,579 genes) and 98 GCCs. The number of GCCs is larger compared to the expected number of ribosomal proteins families (16) used for validation. When we use *high-quality* GCs (Supp. Note 8), we get closer to the expected number of GCCs (Fig. 2D). With this subset, we identified 26 GCCs and 145 GCs (1,687 genes). The cross-validation of our method against the approach used in Méheust et al.[43] (Supp. Note 9) confirms the intrinsic complexity of analyzing metagenomic data. Both approaches showed a high agreement in the GCCs identified (Supp. Table 9-1). Still, our method inferred fewer GCCs for each of the ribosomal protein families (Supplementary Figure 9-3), coping better with the nuisances of a metagenomic setup, like incomplete genes (Supp. Table 6).

**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 take into account 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

# A smaller but highly diverse unknown coding sequence space

211    Combining clustering and remote homology searches reduce the extent of the unknown CDS-

212    space compared to the traditional genomic and metagenomic analysis approaches (Fig. 3A). Our

213    workflow recruited as much as 71% of genes in human-related metagenomic samples and 65%

214    of the genes in marine metagenomes into the known CDS-space. In both human and marine

215    microbiomes, the genomic unknown fraction showed a similar proportion of genes (21%, Fig. 3A).

216    The number of genes corresponding to EU gene clusters is higher in marine metagenomes; in

217    total, 12% of the genes are part of this GC category. We obtained a comparable result when we

218    evaluated the genes from the GTDB_r86, 75% of bacterial and 64% of archaeal genes were part

219    of the known CDS-space. Archaeal genomes contained more unknowns than those from Bacteria,

220    where 30% of the genes are classified as genomic unknowns in Archaea, and only 20% in

221    Bacteria (Fig. 3A; Supp. Table 7). We observed a similar trend when we evaluated the number of

222    amino acids belonging to the known and unknown CDS-space. From the 90,128,659,316 amino

223    acids analyzed, the majority of the amino acids in metagenomes (74%) and in GTDB_r86 (80%)

224    are in the known CDS-space (Fig. 3B; Supp. Table 7). In both cases, approximately 40% of the

225    amino acids in the known CDS-space were part of a Pfam domain (Fig. 3B; Supp. Table 7). The

226    proportion of amino acids in the unknown CDS-space ranged from the 22% in metagenomes and

227    15% in GTDB_r86. In both cases, only 2% of the amino acids in the unknown CDS-space were

228    covered by a Pfam domain.

229    To evaluate the coverage of our dataset, we calculated the accumulation rates of GCs and GCCs.

230    For the metagenomic dataset we used 1,264 metagenomes (18,566,675 GCs and 282,580

231    GCCs) and for the genomic dataset 28,941 genomes (9,586,109 GCs and 496,930 GCCs). The

232    rate of accumulation of unknown GCs was three times higher than the known (2 times for the

233    genomic), and both cases were far from reaching a plateau (Fig. 3C). This is not the case for the

234    GCC accumulation curves (Supp Fig 4B), where they reached a plateau. The rate of accumulation

235    is largely determined by the number of singletons, and especially singletons from EUs (Supp note

236    11 and Supp Fig 5). While the accumulation rate of known GCs between marine and human

237    metagenomes is almost identical, there are striking differences for the unknown GCs (Fig. 3D).

238    These differences are maintained even when we remove the virus-enriched samples from the

11

239     marine metagenomes (Supp Fig 4A). Although the marine metagenomes include a large variety

240     of environments, from coastal to the deep sea, the known space remains quite constrained.

241     Despite only including marine and human metagenomes in our database, our coverage to other

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

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

244     by the screening of MGnify[16] (release 2018_09; 11 biomes; 843,535,6116 proteins) where we

245     assigned 74% of the MGnify proteins into one of our categories (Supplementary Fig. 6).

246

**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 represent the mean of 1000 permutations and shaded in grey the standard deviation. Non-abundant singleton clusters were excluded from the accumulation curves calculation.
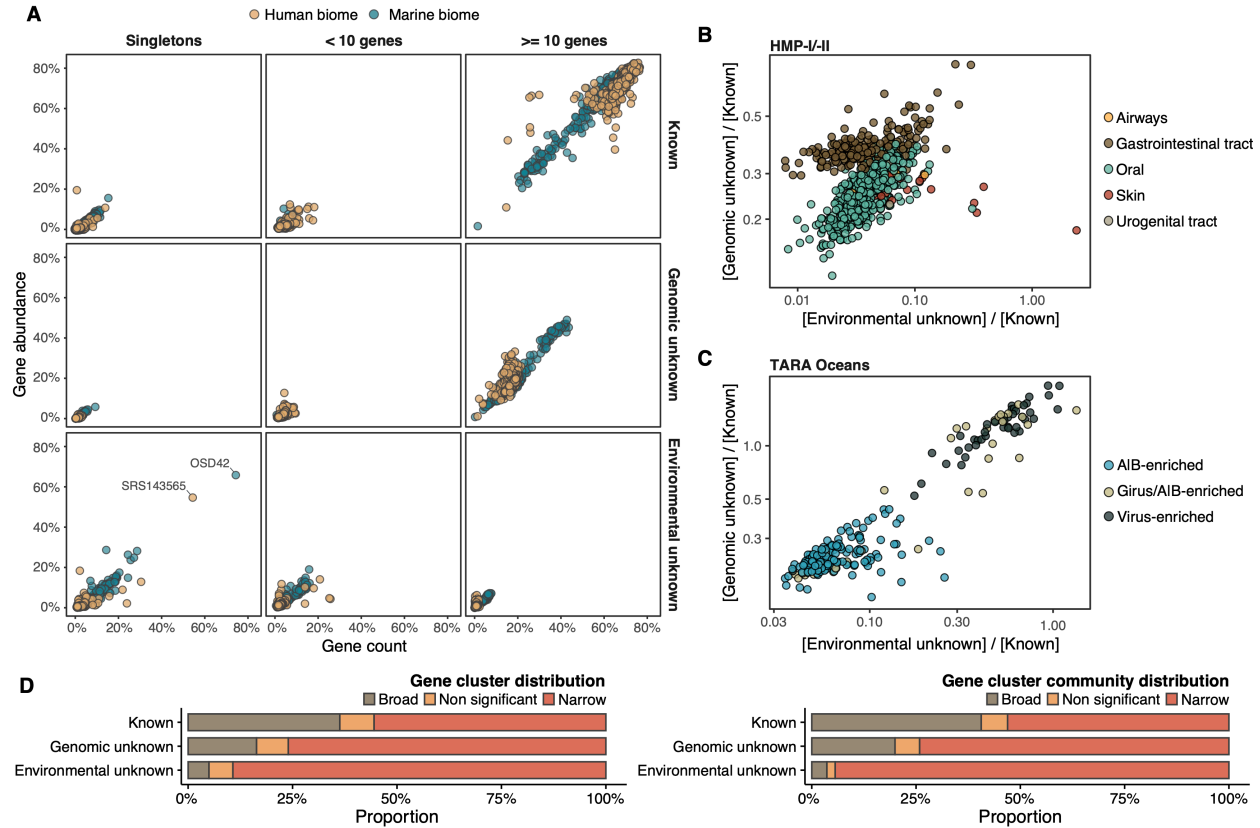
# Revealing the importance of the unknown coding sequence space in marine and human environments

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 samples, 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[45]. The HMP metagenome (SRS143565) corresponds to a human sample from the right cubital fossa from a healthy female subject. To understand the unusual composition of this metagenome, 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, especially enriched with genomic unknowns. The HMP metagenomes with the largest ratio of unknowns are those samples identified to contain crAssphages[46,47] and HPV viruses[48] (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.

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 for the adaptation to specific environmental conditions. But the pool of broadly distributed environmental unknowns is the most interesting result. We identified traces of potential ubiquitous organisms left

14

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

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

283



**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) compared to 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.
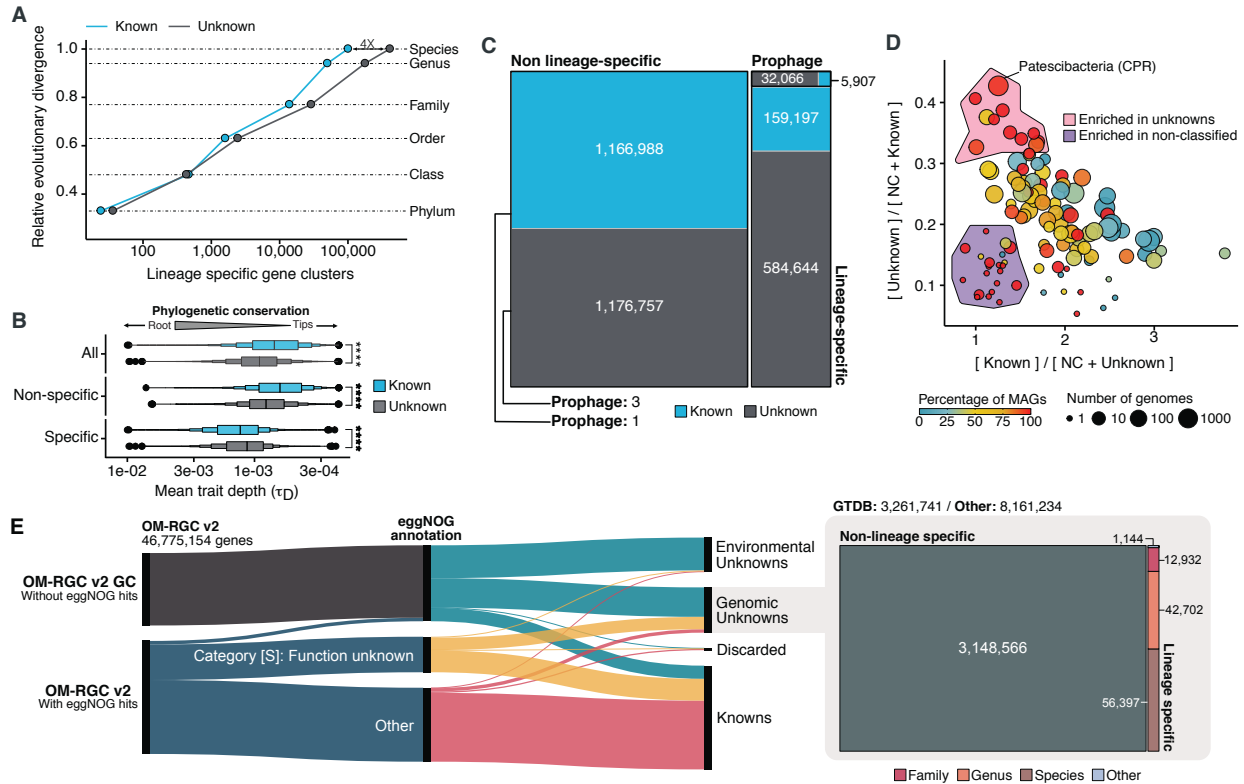
296

# The genomic unknown coding sequence space is lineage-specific

298 We already showed that the unknown CDS-space is habitat-specific and might be relevant for

299 organism adaptation. With the inclusion of the genomes from GTDB_r86, we have accessed a

15

300 phylogenomic framework to assess how conserved and exclusive is a GC within a lineage
301 (lineage-specifity[49]) and the clade depth where organisms share a GC (phylogenetic
302 conservation[50]). We identified 781,814 lineage-specific GCs and 464,923 phylogenetically
303 conserved (P < 0.05) GCs in Bacteria (Supp. Table 10; Supp. Note 13 for Archaea). The number
304 of lineage-specific GCs increases with the Relative Evolutionary Distance[11] (Fig. 5A) and
305 differences between the known and the unknown fraction start to be evident at the Family level.
306 The unknown GCs are more phylogenetically conserved than the known (Fig. 5B, p < 0.0001),
307 revealing the importance of the genome's uncharacterized fraction. However, this is not the case
308 for the lineage-specific and phylogenetically conserved GCs, where the unknown GCs are less
309 phylogenetically conserved (Fig. 5B), agreeing with the large number of lineage-specific GCs at
310 Genus and Species level. To discard the possibility that the lineage-specific GCs of unknown
311 function have a viral origin, we screened all GTDB_r86 genomes for prophages. We only found
312 37,163 lineage-specific GCs in prophage genomic regions, being 86% of them GCs of unknown
313 function. After unveiling the potential relevance of the GCs of unknown function in bacterial
314 genomes, we identified phyla in GTDB_r86 enriched with these types of clusters. A clear pattern
315 emerged when we partitioned the phyla based on the ratio of known to unknown GCs and vice
316 versa (Fig. 5D), the phyla with a larger number of MAGs are enriched in GCs of unknown function
317 Figure 5D. Phyla with a high proportion of non-classified GCs (those discarded during the
318 validation steps) contain a small number of genomes and are primarily composed of MAGs. These
319 groups of phyla highly enriched in unknowns and represented mainly by MAGs include newly
320 described phyla such as *Cand.* Riflebacteria and *Cand.* Patescibacteria[9,51,52], both with the largest
321 unknown to known ratio.
322 We demonstrate the possibility to bridge genomic and metagenomic data and simultaneously
323 unify the known and unknown CDS-space by integrating the new Ocean Microbial Reference
324 Gene Catalog[53] (OM-RGC v2) in our database. We assigned 26,170,875 genes to known GCs,
325 11,422,975 to genomic unknowns, 8,661,221 to environmental unknown and 520,083 were
326 discarded. From the 11,422,975 genes classified as genomic unknowns, we could associate
327 3,261,741 to a GTDB_r86 genome and we identified 113,175 as lineage-specific. The alluvial plot
328 in Fig. 5E depicts the new organization of the OM-RGC v2 after being integrated into our
329 framework, and how we can provide context to the two original types of unknowns in the OM-
330 RGC (those annotated as category S in eggNOG[54] and those without known homologs in the
331 eggNOG database[53]) that can lead to potential experimental targets at the organism level to
332 complement the metatranscriptomic approach proposed by Salazar et al[53].
333
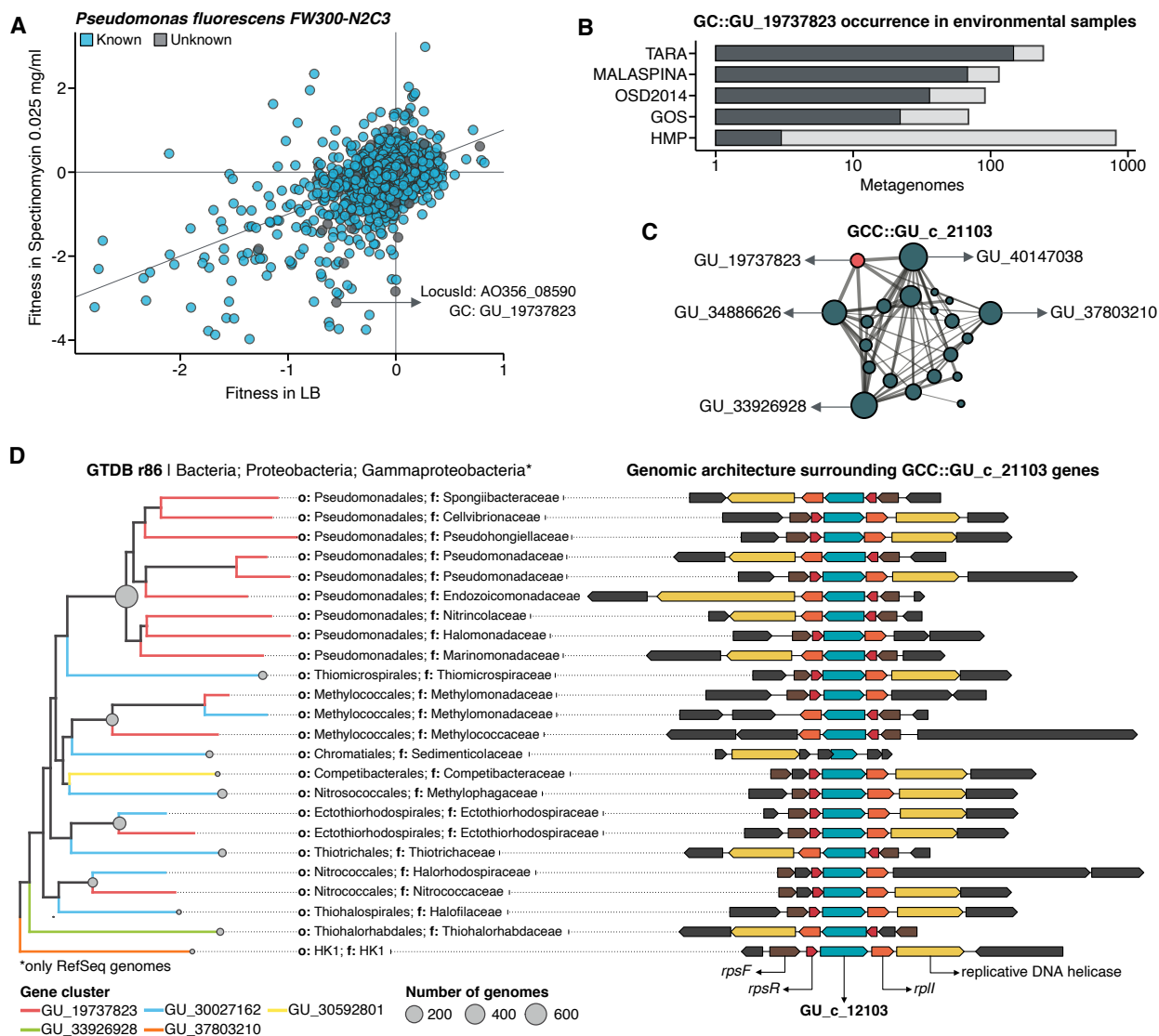
**Figure 5:** Phylogenomic exploration of the unknown coding sequence space. (A) Distribution of the lineage-specific GCs by taxonomic level. Lineage-specific unknown GCs are more abundant in the lower taxonomic levels (genus, species). (B) Phylogenetic conservation of the known and unknown coding sequence space in 27,372 bacterial genomes from GTDB_r86. We observe differences in the conservation between the known and the unknown coding sequence space for lineage- and non-lineage specific GCs (paired Wilcoxon rank-sum test; all p-values < 0.0001). (C) The majority of the lineage-specific clusters are part of the unknown coding sequence space, being a small proportion found in prophages present in the GTDB_r86 genomes. (D) Known and unknown coding sequence space of the 27,732 GTDB_r86 bacterial genomes grouped by bacterial phyla. Phyla are partitioned based on the ratio of known to unknown GCs and vice versa. Phyla enriched in MAGs have higher proportions in GCs of unknown function. Phyla with a high proportion of non-classified clusters (NC; discarded during the validation steps) tend to contain a small number of genomes. (E) The left side of the alluvial plot shows the uncharacterized (OM-RGC v2 GC) and characterized (OM-RGC v2) fraction of the gene catalog. The functional annotation is based on the eggNOG annotations provided by Salazar et al.[53]. The right side of the alluvial plot shows the new organization of the OM-RGC v2 coding sequence space based on the approach described in this study. The treemap in the right links the metagenomic and genomic space adding context to the unknown fraction of the OM-RGC v2

## 353 Augmenting experimental data through a structured coding
## 354 sequence space

355   We selected one of the experimental conditions tested in Price et al.[13] to demonstrate the potential

356   of our approach to augment experimental data. We compared the fitness values in plain rich

357   medium with added Spectinomycin dihydrochloride pentahydrate to the fitness in plain rich

358   medium (LB) in *Pseudomonas fluorescens FW300-N2C3* (Fig. 6A). This antibiotic inhibits protein

359   synthesis and elongation by binding to the bacterial 30S ribosomal subunit and interferes with the

360   peptidyl tRNA translocation. We identified the gene with locus id AO356_08590 that presents a

361   strong phenotype (fitness = -3.1; t = -9.1) and has no known function. This gene belongs to the

362   genomic unknown GC GU_19737823. We can track this GC into the environment and explore

363   the occurrence in the different samples we have in our database. As expected, the GC is mostly

364   found in non-human metagenomes (Fig. 6B) as *Pseudomonas* are common inhabitants of soil

365   and water environments[55].   However, finding this GC also in human-related samples is very

366   interesting, due to the potential association of *P. fluorescens* and human disease where Crohn's

367   disease patients develop serum antibodies to this microbe[56]. We can add another layer of

368   information to the selected GC by looking at the associated remote homologs in the GCC

369   GU_c_21103 (Fig. 6C). We identified all the genes in the GTDB_r86 genomes that belong to the

370   GCC GU_c_21103 (Supp. Table 11) and explored their genomic neighborhoods. All members

371   from GU_c_21103 are constrained to the class *Gammaproteobacteria*, and interestingly

372   GU_19737823 is mostly exclusive to the order *Pseudomonadales*. The gene order in the different

373   genomes analyzed is highly conserved, finding GU_19737823 after the *rpsF*::*rpsR* operon and

374   before *rpll*. *rpsF* and *rpsR* encode for 30S ribosomal proteins, the prime target of spectinomycin.

375   The combination of the experimental evidence and the associated data inferred by our approach

376   provides strong support to generate the hypothesis that the gene AO356_08590 might be involved

377   in the resistance to spectinomycin.

378

379

**Figure 6:** Augmenting experimental data with GCs of unknown function. (A) We used the fitness values from the experiments from Price et al.[13] 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 *rplI*. *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. Bubble plot depicts the number of genomes with a gene that belongs to GU_c_21103.

19

# Discussion

We present a new conceptual framework and computational workflow to unify the known and unknown CDS-space in microbial analyses. Using this framework, we performed an in-depth exploration of the microbial unknown CDS-space. We demonstrated that we could link the unknown fraction of metagenomic studies to specific genomes and provide a powerful tool for hypothesis generation. During the last years, the microbiome community has established a standard operating procedure[17] for analyzing metagenomes that we can briefly summarize into (1) assembly, (2) gene prediction, (3) gene catalog inference, (4) binning, and (5) characterization. Thanks to recent computational developments[36,57], we envisioned an alternative to this workflow where we can maximize the information used when analyzing genomic and metagenomic data. In addition, we provide a mechanism to reconcile top-down and bottom-up approaches, thanks to the well-structured CDS-space proposed by our framework. AGNOSTOS can create environmental- and organism-specific variations of a seed GC database. Then, it integrates the predicted genes from new genomes and metagenomes and dynamically creates and classifies new GCs with those genes not integrated during the initial step (Fig. 1B). Afterward, the potential functions of the known GCs can be carefully characterized by incorporating them into the traditional workflows.

One of the most appealing characteristics of our approach is that the GCs provide unified groups of homologous genes across environments and organisms indifferently if they belong to the known or unknown CDS-space, and we can contextualize the unknown fraction using this genomic and environmental information. Our combination of partitioning and contextualization features a smaller unknown CDS-space than we expected. On average, for our genomic and metagenomic data, only 30% of the genes fall in the unknown fraction. One hypothesis to reconcile this surprising finding is that until recently, the methodologies to identify remotely homologous sequences in large datasets were computationally prohibitive. New methods[36,37], like the ones used in AGNOSTOS, are enabling large scale distant homology searches. Still, one has to apply conservative measures to control the trade-off between specificity and sensitivity to avoid overclassification.

We found that the majority of the coding sequence space at gene and amino acid is known, both in genomes and metagenomes. However, it presents a high diversity as shown in the GC accumulation curves highlighting the vast remaining untapped microbial fraction and its potential importance for niche adaptation owing to its narrow ecological distribution. In a genomic context, the unknown fraction is predominantly species' lineage-specific and phylogenetically more

20

429    conserved than the known fraction, supporting the signal observed in the environmental data and

430    emphasizing that the unknown fraction should not be ignored. We also ruled out the effect of

431    prophages, strengthening the hypothesis that the lineage-specific GCs of unknown function might

432    be associated with the mechanisms of microbial diversification and niche adaptation as a result

433    of the constant diversification of gene families and the survival of new gene lineages[58,59]. It is

434    worth noting that we need to explore further the unknown fraction to identify new potential protein

435    domains. Only 10% of the unknown CDS-space amino acids are part of a Pfam domain (DUF and

436    others); this contrasts with the numbers observed in the known CDS-space, where Pfam domains

437    include 50% of the amino acids.

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

439    (42% of the genomes in GTDB_r86), but they are also enriching genes of unknown function in

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

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

442    interest due to their unusual biology[9]. We provide a collection of 54,343 lineage-specific GCs of

443    unknown function at different taxonomic level resolutions (Supp. Table 12; Supp. Note 14), which

444    will be a valuable resource for the advancement of knowledge in the CPR research efforts.

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

446    that remains untapped to analysis and discovery. With the identification of a potential target gene

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

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

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

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

451    sequence assembly has been an active area of research[60], this has not been the case for gene

452    prediction methods[60], which are becoming outdated[61] and cannot cope with the current amount

453    of data. Alternatives like protein-level assembly[62] combined with the exploration of the assembly

454    graphs' neighborhoods[63] become very attractive for our purposes. In any case, we still face the

455    challenge of discriminating between real and artifactual singletons[64]. At the moment, there are no

456    methods available to provide a plausible solution and, at the same time, being scalable. We devise

457    a potential solution in the recent developments in unsupervised deep learning methods where

458    they use large corpora of proteins to define a language model *embedding* for protein sequences[65].

459    These models could be applied to predict *embeddings* in singletons, which could be clustered or

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

461    than expected. We follow a conservative approach to avoid mixing multidomain proteins in GCs

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

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

464 to the use of remote homologies. Although the inference of GCCs is using very sensitive methods

465 to compare profile HMMs, low sequence diversity in GCs can limit its effectiveness. Our approach

466 is affected by the presence and propagation of contamination in reference databases, a significant

467 problem in 'omics[66,67]. In our case, we only use Pfam as a source for annotation owing to its high-

468 quality and manual curation process. The categorization process of our GCs depends on the

469 information from other databases, and to minimize the potential impact of contamination, we apply

470 methods that weight the annotations of the identified homologs to discriminate if a GC belongs to

471 the known or unknown CDS-space. We foresee the integration of our approach to assist in the

472 manual curation process and increase the quality of the recovered MAGs[68].

473 The work presented here should incentivize the scientific community to build a collective effort to

474 define the different levels of unknown[69] where clear guidelines and protocols should be

475 established. Our work proves that the integration of the unknown fraction is possible and aims to

476 provide a new brighter future for microbiome analyses.

477

# 478 Material and methods

## 479 Genomic and metagenomic dataset

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

481 ocean microbiome: Tara Oceans expedition (TARA)[2], Malaspina expedition[70], Ocean Sampling

482 Day (OSD)[3], and Global Ocean Sampling Expedition (GOS)[71]. We complemented this set with

483 1,246 metagenomes obtained from the Human Microbiome Project (HMP) phase I and II[72]. We

484 used the assemblies provided by TARA, Malaspina, OSD and HMP projects and the long Sanger

485 reads from GOS[73]. A total of 156M (156,422,969) contigs and 12.8M long-reads were collected

486 (Supp. Table 6).

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

488 archaeal) from the Genome Taxonomy Database[11] (GTDB) Release 03-RS86 (19th August

489 2018).

## 490 Computational workflow development

491 We implemented a computation workflow based on Snakemake[74] for the easy processing of large

492 datasets in a reproducible manner. The workflow provides three different strategies to analyze

493   the data. The module *DB-creation* creates the gene cluster database, validates and partitions the

494   gene clusters (GCs) in the main functional categories. The module *DB-update* allows the

495   integration of new sequences (either at the contig or predicted gene level) in the existing gene

496   cluster database. In addition, the workflow has a *profile-search* function to quickly screen samples

497   using the gene cluster PSSM profiles in the database.

## Metagenomic and genomic gene prediction

499   We used Prodigal (v2.6.3)[75] in metagenomic mode to predict the genes from the metagenomic

500   dataset. For the genomic dataset, we used the gene predictions provided by Annotree[76], since

501   they were obtained, consistently, with Prodigal v2.6.3. We identified potential spurious genes

502   using the *AntiFam* database[77]. Furthermore, we screened for '*shadow' genes* using the procedure

503   described in Yooseph et al.[78].

## PFAM annotation

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

506   (version: 3.1b2)[79] in combination with the Pfam database v31[80]. We kept the matches exceeding

507   the internal gathering threshold and presenting an independent e-value < 1e-5 and coverage >

508   0.4. In addition, we took into account multi-domain annotations, and we removed overlapping

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

## Determination of the gene clusters

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

512   MMseqs2 software[57] ("*--cov-mode 0 -c 0.8 --min-seq-id 0.3"*). We discarded from downstream

513   analyses the singletons and clusters with a size below a threshold identified after applying a

514   broken-stick model[81]. We integrated the genomic data into the metagenomic cluster database

515   using the "DB-update" module of the workflow. This module uses the *clusterupdate* module of

516   MMseqs2[36], with the same parameters used for the metagenomic clustering.

## Quality-screening of gene clusters

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

519   methodologies to validate their cluster sequence composition and functional annotation

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

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

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

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

524    alignment methods implemented in *PARASAIL*[82] (version 2.1.5). We trimmed the SSN using a

525    custom algorithm[83,84] that removes edges while maintaining the network structural integrity and

526    obtaining the smallest connected graph formed by a single component. Finally, the new cluster

527    representative was identified as the most central node of the trimmed SSN by the eigenvector

528    centrality algorithm, as implemented in igraph[85]. After this step, we built a multiple sequence

529    alignment for each cluster using *FAMSA*[86] (version 1.1). Then, we screened each cluster-MSA

530    for non-homologous sequences to the new cluster representative. Owing to computational

531    limitations, we used two different approaches to evaluate the cluster-MSAs. We used *LEON-BIS*[87]

532    for the clusters with a size ranging from 10 to 1,000 genes and OD-SEQ[88] for the clusters with

533    more than 1,000 genes. In the end, we applied a broken-stick model[81] to determine the threshold

534    to discard a cluster.

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

536    the consistency of intra-cluster Pfam annotations, we applied a combination of w-shingling[89] and

537    Jaccard similarity. We used w-shingling (k-shingle = 2) to group consecutive domain annotations

538    as a single object. We measured the homogeneity of the *shingle sets* (sets of domains) between

539    genes using the Jaccard similarity and reported the median similarity value for each cluster.

540    Moreover, we took into consideration the Clan membership of the Pfam domains and that a gene

541    might contain N-, C- and M-terminal domains for the functional homogeneity validation. We

542    discarded clusters with a median similarity < 1.

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

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

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

546    2).


## 547   Remote homology classification of gene clusters

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

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

550    GCs, we inferred a consensus protein domain architecture (DA) (an ordered combination of

551    protein domains) for each annotated gene cluster. To identify each gene cluster consensus DA,

552    we created directed acyclic graphs connecting the Pfam domains based on their topological order

553    on the genes using *igraph*[85]. We collapsed the repetitions of the same domain. Then we used the

554     gene completeness as a positive-weighting value for the selection of the cluster consensus DA.

555     Within this step, we divided the GCs into "Knowns" (Known) if annotated to at least one Pfam

556     domains of known function (DKFs) and "Genomic unknowns" (GU) if annotated entirely to Pfam

557     domains of unknown function (DUFs).

558     We aligned the sequences of the non-annotated GCs with FAMSA[86] and obtained cluster

559     consensus sequences with the *hhconsensus* program from *HH-SUITE*[37]. We used the cluster

560     consensus sequences to perform a nested search against the UniRef90 database (release

561     2017_11)[90] and NCBI *nr* database (release 2017_12)[91] to retrieve non-Pfam annotations with

562     *MMSeqs2*[36] ("*-e 1e-05 --cov-mode 2 -c 0.6*"). We kept the hits within 60% of the Log(best-e-value)

563     and searched the annotations for any of the terms commonly used to define proteins of unknown

564     function (Supp. Table 12). We used a quorum majority voting approach to decide if a gene cluster

565     would be classified as *Genomic Unknown* or *Known without Pfams* based on the annotations

566     retrieved. We searched the consensus sequences without any homologs in the UniRef90

567     database against NCBI *nr*. We applied the same approach and criteria described for the first

568     search. Ultimately, we classified as *Environmental Unknown* those GCs whose consensus

569     sequences did not align with any of the NCBI *nr* entries.

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

571     and sensitivity for the remote homology searches such as (1) a modification of the algorithm

572     described in Hingamp et al.[92] to get a confident group of homologs to determine if a query protein

573     is known or unknown by a quorum majority voting approach (Supp Note 3); (2) strict parameters

574     in terms of iterations, bidirectional coverage and probability thresholds for the HHblits alignments

575     to minimize the inclusion of non-homologous sequences; and (3) avoid providing annotations for

576     our gene clusters, as we believe that annotation should be a careful process done on a smaller

577     scale and with experimental context.

# Gene cluster remote homology refinement

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

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

581     HMM/HMM alignment method *HHblits*[93]. We created the HMM profiles with the *hhmake* program

582     from the *HH-SUITE*[37]. We only accepted those hits with an *HHblits-probability* ≥ 90% and we re-

583     classified them following the same majority vote approach as previously described. The clusters

584     with no hits remained as the refined set of EUs. We applied a similar refinement approach to the

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

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

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

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

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

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

# Gene cluster characterization

591

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

593    (version: b43de8b7559a3b45c8e5e9e02cb3023dd339231a), which allows computing the lowest

594    common ancestor through the implementation of the 2bLCA protocol [92]. We searched all cluster

595    genes against UniProtKB (release of January 2018) [94] using the following parameters "*-e 1e-05 -*

596    *-cov-mode 0 -c 0.6"*. We parsed the results to keep only the hits within 60% of the log10(best-e-

597    value). To retrieve the taxonomic lineages, we used the R package *CHNOSZ*[95]. We measured

598    the intra-cluster taxonomic admixture by applying the *entropy.empirical()* function from the *entropy*

599    R package[96]. This function estimates the Shannon entropy based on the different taxonomic

600    annotation frequencies. For each cluster, we also retrieved the cluster consensus taxonomic

601    annotation, which we defined as the taxonomic annotation of the majority of the genes in the

602    cluster.

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

604    Dark Proteome Database (DPD)[40] as reference. We searched the cluster genes against the DPD,

605    applying the MMseqs2 search program[36] with "*-e 1e-20 --cov-mode 0 -c 0.6*". For each cluster,

606    we then retrieved the mean and the median level of darkness, based on the gene DPD

607    annotations.

# High-quality clusters

608

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

610    their representatives. We identified the minimum required percentage of complete genes per

611    cluster by a broken-stick model[81] applied to the percentage distribution. Then, we selected the

612    GCs found above the threshold and with a complete representative.

# A set of non-redundant domain architectures

613

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

615    account the large proportion of fragmented genes in the metagenomic dataset and that could

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

617    architecture, we took into account their topological order in the genes. To reduce the number of

618    comparisons, we calculated the pairwise string cosine distance (q-gram = 3) between domain

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

620    collapsed a fragmented domain architecture to the larger one when it contained less than 75% of

621    complete genes.

## Inference of gene cluster communities

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

624    combined an all-vs-all HMM gene cluster comparison with Markov Cluster Algorithm (MCL)[97]

625    community identification. We started performing the inference on the Known GCs to use the Pfam

626    DAs as constraints. We aligned the gene cluster HMMs using HHblits[93] (-n 2 -Z 10000000 -B

627    10000000 -e 1) and we built a homology graph using the cluster pairs with probability ≥ 50% and

628    bidirectional coverage > 60%. We used the ratio between HHblits-bitscore and aligned-columns

629    as the edge weights (Supp. Note 9). We used MCL[97] (v. 12-068) to identify the communities

630    present in the graph. We developed an iterative method to determine the optimal MCL inflation

631    parameter that tries to maximize the relationship of five intra-/inter-community properties: (1) the

632    proportion of MCL communities with one single DA, based on the consensus DAs of the cluster

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

634    communities with a PFAM clan entropy equal to 0; (4) the intra-community HHblits-score/Aligned-

635    columns score (normalized by the maximum value); and (5) the number of MCL communities,

636    which should, in the end, reflect the number of non-redundant DAs. We iterated through values

637    ranging from 1.2 to 3.0, with incremental steps of 0.1. During the inference process, some of the

638    GCs became orphans in the graph. We applied a three-step approach to assigning a community

639    membership to these GCs. First, we used less stringent conditions (probability ≥ 50% and

640    coverage >= 40%) to find homologs in the already existing GCCs. Then, we ran a second iteration

641    to find secondary relationships between the newly assigned GCs and the missing ones. Lastly,

642    we created new communities with the remaining GCs. We repeated the whole process with the

643    other categories (KWP, GU and EU), applying the optimal inflation value found for the Known (2.2

644    for metagenomic and 2.5 for genomic data).

## Gene cluster communities validation

646    We tested the biological significance of the GCCs using the phylogeny of proteorhodopsin[44] (PR).

647    We used the proteorhodopsin HMM profiles[42] to screen the marine metagenomic datasets using

27

648    *hmmsearch* (version 3.1b2)[79]. We kept the hits with a coverage > 0.4 and e-value <= 1e-5. We

649    removed identical duplicates from the sequences assigned to PR with CD-HIT[98] (v4.6) and

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

651    into the MicRhode[44] PR tree first, we optimized the initial tree parameters and branch lengths with

652    RAxML (v8.2.12)[99]. We used PaPaRA (v2.5)[100] to incrementally align the query PR sequences

653    against the MicRhode PR reference alignment and *pplacer*[101] (v1.1.alpha19-0-g807f6f3) to place

654    the sequences into the tree. Finally, we assigned the query PR sequences to the MicRhode PR

655    Superclusters based on the phylogenetic placement. We further investigated the GCs annotated

656    as viral (196 genes, 14 GC) comparing them to the six newly discovered viral PRs[102] using

657    Parasail[82] (-a sg_stats_scan_sse2_128_16 -t 8 -c 1 -x). As an additional evaluation, we

658    investigated the distributions of standard GCCs and HQ GCCs within ribosomal protein families.

659    We obtained the ribosomal proteins used for the analysis combining the set of 16 ribosomal

660    proteins from Méheust et al.[43] and those contained in the collection of bacterial single-copy genes

661    of Anvi'o[103]. Also, for the ribosomal proteins, we compared the outcome of our method to the one

662    proposed by Méheust et al.[43] (Supp. Note 9).

663    # Metagenomic sample selection for downstream analyses

664    For the subsequent ecological analyses we selected those metagenomes with a number of genes

665    larger or equal to the first quartile of the distribution of all the metagenomic gene counts. (Supp.

666    Table 13).

667    # Gene cluster abundance profiles in genomes and metagenomes

668    We estimated abundance profiles for the metagenomic cluster categories using the read coverage

669    to each predicted gene as a proxy for abundance. We calculated the coverage by mapping the

670    reads against the assembly contigs using the *bwa-mem* algorithm from *BWA mapper*[104]. Then,

671    we used *BEDTOOLS*[105], to find the intersection of the gene coordinates to the assemblies, and

672    normalize the per-base coverage by the length of the gene. We calculated the cluster abundance

673    in a sample as the sum of the cluster gene abundances in that sample, and the cluster category

674    abundance in a sample as the sum of the cluster abundances. We obtained the proportions of

675    the different gene cluster categories applying a total-sum-scaling normalization. For the genomic

676    abundance profiles, we used the number of genes in the genomes and normalized by the total

677    gene counts per genome.

28

## 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 that got reclassified as good-quality clusters after integrating the GTDB data to minimize the impact of potential spurious singletons. To complement those analyses, we evaluated the coverage of our dataset by searching seven different state-of-the-art databases against our set of metagenomic GC HMM profiles (Supp. Note 12).

## Occurrence of gene clusters in the environment

We used 1,264 metagenomes from the TARA Oceans, MALASPINA Expedition, OSD2014 and HMP-I/II to explore the properties of the unknown CDS-space in the environment. We applied the Levins Niche Breadth (NB) index[106] to investigate the GCs and GCCs environmental distributions. We removed the GCs and cluster communities with a mean relative abundance < 1e-5. We followed a divide-and-conquer strategy to avoid the computational burden of generating the null-models to test the significance of the distributions owing to the large number of metagenomes and GCs. First, we grouped similar samples based on the gene cluster content using the Bray-Curtis dissimilarity[107] in combination with the *Dynamic Tree Cut*[108] R package. We created 100 random datasets picking up one random sample from each group. For each of the 100 random datasets, we created 100 random abundance matrices using the *nullmodel* function of the *quasiswap* count method[109]. Then we calculated the *observed* NB and obtained the 2.5% and 97.5% quantiles based on the randomized sets. We compared the observed and quantile values for each gene cluster and defined it to have a *Narrow distribution* when the *observed* was smaller than the 2.5% quantile and to have a *Broad distribution* when it was larger than the 97.5% quantile. Otherwise, we classified the cluster as *Non-significant*[110]. We used a majority voting approach to get a consensus distribution 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[76] to identify the lineage-specific GCs and to which rank they are specific. Following similar criteria to the ones used in Mendler et al.[76], 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*[50] function implemented in the R package *castor*[50]. 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)[53] into our cluster database using the same procedure as for the genomic data. We evaluated the uncharacterized fraction and the genes classified into the eggNOG[54] category S within the context of our database.

## Augmenting experimental data

We searched the 37,684 genes of unknown function associated with mutant phenotypes from Price et al.[13] against our gene cluster profiles. We kept the hits with e-value ≤ 1e-20 and a query coverage > 60%. Then we filtered the results to keep the hits within 90% of the Log(best-e-value), and we used a majority vote function to retrieve the consensus category for each hit. Lastly, we

735    selected the best-hits based on the smallest e-value and the largest query and target coverage

736    values. We used the fitness values from the RB-TnSeq experiments from Price et al. to identify

737    genes of unknown function that are important for fitness under certain experimental conditions.

## Code and data availability

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

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

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

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

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

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

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

746

747

# Acknowledgments

764     Consolider-Ingenio program (ref. CSD2008-00077). The authors thank Johannes Söding and

765     Alex Bateman for helpful discussions.

766

# Author Contributions

768     CV, MSS and AF-G performed the analyses and wrote the computational workflow. MS assisted

769     with the clustering and remote homology searches. KS helped with the identification of prophages

770     in genomic sequences. PLB and AB provided feedback and assisted with the ecological analyses.

771     RDF and AM provided feedback and information on the MGnify and Pfam databases. CMD, PS

772     and SGA provided the Malaspina metagenomes. TOD and AME analyzed data in the context of

773     metagenome-assembled genomes. AF-G conceived the study and supervised the work. CV and

774     AF-G wrote the manuscript. All authors read, edited and approved the final manuscript.

775

776

# Competing Interests

778     The authors declare no competing interests.

779

780

# References

1. Hug, L. A. *et al.* A new view of the tree of life. *Nat Microbiol* **1**, 16048 (2016).

2. Sunagawa, S. *et al.* Ocean plankton. Structure and function of the global ocean microbiome. *Science* **348**, 1261359 (2015).

3. Kopf, A. *et al.* The ocean sampling day consortium. *Gigascience* **4**, 27 (2015).

4. Almeida, A. *et al.* A new genomic blueprint of the human gut microbiota. *Nature* **568**, 499–504 (2019).

5. Pasolli, E. *et al.* Extensive Unexplored Human Microbiome Diversity Revealed by Over 150,000 Genomes from Metagenomes Spanning Age, Geography, and Lifestyle. *Cell* **176**, 649-662.e20 (2019).

6. Pachiadaki, M. G. *et al.* Charting the Complexity of the Marine Microbiome through Single-Cell Genomics. *Cell* **179**, 1623-1635.e11 (2019).

7. Cross, K. L. *et al.* Targeted isolation and cultivation of uncultivated bacteria by reverse genomics. *Nat. Biotechnol.* **37**, 1314–1321 (2019).

8. Eloe-Fadrosh, E. A. *et al.* Global metagenomic survey reveals a new bacterial candidate phylum in geothermal springs. *Nat. Commun.* **7**, 10476 (2016).

9. Brown, C. T. *et al.* Unusual biology across a group comprising more than 15% of domain Bacteria. *Nature* **523**, 208–211 (2015).

10. Spang, A. *et al.* Complex archaea that bridge the gap between prokaryotes and eukaryotes. *Nature* **521**, 173–179 (2015).

11. Parks, D. H. *et al.* A standardized bacterial taxonomy based on genome phylogeny substantially revises the tree of life. *Nat. Biotechnol.* **36**, 996–1004 (2018).

12. Bernard, G., Pathmanathan, J. S., Lannes, R., Lopez, P. & Bapteste, E. Microbial Dark Matter Investigations: How Microbial Studies Transform Biological Knowledge and

808    Empirically Sketch a Logic of Scientific Discovery. *Genome Biol. Evol.* **10**, 707–715 (2018).

809   13. Price, M. N. *et al.* Mutant phenotypes for thousands of bacterial genes of unknown function.

810    *Nature* **557**, 503–509 (2018).

811   14. Carradec, Q. *et al.* A global ocean atlas of eukaryotic genes. *Nat. Commun.* **9**, 373 (2018).

812   15. Almeida, A. *et al.* A unified catalog of 204,938 reference genomes from the human gut

813    microbiome. *Nat. Biotechnol.* (2020) doi:10.1038/s41587-020-0603-3.

814   16. Mitchell, A. L. *et al.* MGnify: the microbiome analysis resource in 2020. *Nucleic Acids Res.*

815    **48**, D570–D578 (2020).

816   17. Quince, C., Walker, A. W., Simpson, J. T., Loman, N. J. & Segata, N. Shotgun

817    metagenomics, from sampling to analysis. *Nat. Biotechnol.* **35**, 833–844 (2017).

818   18. Franzosa, E. A. *et al.* Species-level functional profiling of metagenomes and

819    metatranscriptomes. *Nat. Methods* **15**, 962–968 (2018).

820   19. Huerta-Cepas, J. *et al.* Fast Genome-Wide Functional Annotation through Orthology

821    Assignment by eggNOG-Mapper. *Mol. Biol. Evol.* **34**, 2115–2122 (2017).

822   20. Chen, I.-M. A. *et al.* IMG/M v.5.0: an integrated data management and comparative

823    analysis system for microbial genomes and microbiomes. *Nucleic Acids Res.* **47**, D666–

824    D677 (2019).

825   21. Hanson, A. D., Pribat, A., Waller, J. C. & Crécy-Lagard, V. de. 'Unknown'proteins and

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

827    *J* **425**, 1–11 (2010).

828   22. Arnold, F. H. Design by Directed Evolution. *Acc. Chem. Res.* **31**, 125–131 (1998).

829   23. Brandenberg, O. F., Fasan, R. & Arnold, F. H. Exploiting and engineering hemoproteins for

830    abiological carbene and nitrene transfer reactions. *Curr. Opin. Biotechnol.* **47**, 102–111

831    (2017).

832   24. Arnold, F. H. Directed Evolution: Bringing New Chemistry to Life. *Angew. Chem. Int. Ed*

833    *Engl.* **57**, 4143–4148 (2018).

834    25. Jaroszewski, L. *et al.* Exploration of uncharted regions of the protein universe. *PLoS Biol.* **7**,

835          (2009).

836    26. Buttigieg, L. P. *et al.* Ecogenomic Perspectives on Domains of Unknown Function:

837          Correlation-Based Exploration of Marine Metagenomes. *PLoS One* **8**, (2013).

838    27. Yooseph, S. *et al.* The Sorcerer II global ocean sampling expedition: Expanding the

839          universe of protein families. *PLoS Biol.* **5**, 0432–0466 (2007).

840    28. Wyman, S. K., Avila-Herrera, A., Nayfach, S. & Pollard, K. S. A most wanted list of

841          conserved microbial protein families with no known domains. *PLoS One* **13**, e0205749

842          (2018).

843    29. Brum, J. R. *et al.* Illuminating structural proteins in viral "dark matter" with metaproteomics.

844          *Proc. Natl. Acad. Sci. U. S. A.* **113**, 2436–2441 (2016).

845    30. Bateman, A., Coggill, P. & Finn, D. R. DUFs: Families in search of function. *Acta

846          Crystallogr. Sect. F Struct. Biol. Cryst. Commun.* **66**, 1148–1152 (2010).

847    31. Lobb, B., Kurtz, D. A., Moreno-Hagelsieb, G. & Doxey, A. C. Remote homology and the

848          functions of metagenomic dark matter. *Front. Genet.* **6**, 1–12 (2015).

849    32. Bitard-Feildel, T. & Callebaut, I. Exploring the dark foldable proteome by considering

850          hydrophobic amino acids topology. *Sci. Rep.* **7**, 41425 (2017).

851    33. Bileschi, M. L. *et al.* Using Deep Learning to Annotate the Protein Universe. *bioRxiv* 626507

852          (2019) doi:10.1101/626507.

853    34. Liu, X. L. Deep Recurrent Neural Network for Protein Function Prediction from Sequence.

854          *bioRxiv* 103994 (2017) doi:10.1101/103994.

855    35. Rost, B. Twilight zone of protein sequence alignments. *Protein Eng. Des. Sel.* **12**, 85–94

856          (1999).

857    36. Steinegger, M. & Söding, J. MMseqs2 enables sensitive protein sequence searching for the

858          analysis of massive data sets. *Nat. Biotechnol.* **advance on**, (2017).

859    37. Steinegger, M. *et al.* HH-suite3 for fast remote homology detection and deep protein

bioRxiv preprint doi: https://doi.org/10.1101/2020.06.30.180448; this version posted August 11, 2020. 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.

860    annotation. *BMC Bioinformatics* **20**, 473 (2019).

861    38. Skewes-Cox, P., Sharpton, T. J., Pollard, K. S. & DeRisi, J. L. Profile hidden Markov

862    models for the detection of viruses within metagenomic sequence data. *PLoS One* **9**,

863    e105067 (2014).

864    39. Sberro, H. *et al.* Large-Scale Analyses of Human Microbiomes Reveal Thousands of Small,

865    Novel Genes. *Cell* **178**, 1245-1259.e14 (2019).

866    40. Perdigão, N., Rosa, A. C. & O'Donoghue, S. I. The Dark Proteome Database. *BioData Min.*

867    **10**, 1–11 (2017).

868    41. Habchi, J., Tompa, P., Longhi, S. & Uversky, V. N. Introducing protein intrinsic disorder.

869    *Chem. Rev.* **114**, 6561–6588 (2014).

870    42. Olson, D. K., Yoshizawa, S., Boeuf, D., Iwasaki, W. & DeLong, E. F. Proteorhodopsin

871    variability and distribution in the North Pacific Subtropical Gyre. *ISME J.* **12**, 1047–1060

872    (2018).

873    43. Méheust, R., Burstein, D., Castelle, C. J. & Banfield, J. F. The distinction of CPR bacteria

874    from other bacteria based on protein family content. *Nat. Commun.* **10**, 4173 (2019).

875    44. Boeuf, D., Audic, S., Brillet-Guéguen, L., Caron, C. & Jeanthon, C. MicRhoDE: a curated

876    database for the analysis of microbial rhodopsin diversity and evolution. *Database* **2015**,

877    (2015).

878    45. La Cono, V. *et al.* Partaking of Archaea to biogeochemical cycling in oxygen-deficient

879    zones of meromictic saline Lake Faro (Messina, Italy). *Environ. Microbiol.* **15**, 1717–1733

880    (2013).

881    46. Edwards, R. A. *et al.* Global phylogeography and ancient evolution of the widespread

882    human gut virus crAssphage. *Nat Microbiol* **4**, 1727–1736 (2019).

883    47. Dubinkina, V. B., Ischenko, D. S., Ulyantsev, V. I., Tyakht, A. V. & Alexeev, D. G.

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

885    *Bioinformatics* vol. 17 (2016).

886   48.  Ma, Y. *et al.* Human papillomavirus community in healthy persons, defined by

887         metagenomics analysis of human microbiome project shotgun sequencing data sets. *J.*

888         *Virol.* **88**, 4786–4797 (2014).

889   49.  Mendler, K. *et al.* AnnoTree: visualization and exploration of a functionally annotated

890         microbial tree of life. *Nucleic Acids Res.* **47**, 4442–4448 (2019).

891   50.  Martiny, A. C., Treseder, K. & Pusch, G. Phylogenetic conservatism of functional traits in

892         microorganisms. *ISME J.* **7**, 830–838 (2013).

893   51.  Rinke, C. *et al.* Insights into the phylogeny and coding potential of microbial dark matter.

894         *Nature* **499**, 431–437 (2013).

895   52.  Anantharaman, K. *et al.* Expanded diversity of microbial groups that shape the dissimilatory

896         sulfur cycle. *ISME J.* **12**, 1715–1728 (2018).

897   53.  Salazar, G. *et al.* Gene Expression Changes and Community Turnover Differentially Shape

898         the Global Ocean Metatranscriptome. *Cell* **179**, 1068-1083.e21 (2019).

899   54.  Huerta-Cepas, J. *et al.* eggNOG 5.0: a hierarchical, functionally and phylogenetically

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

901         *Res.* **47**, D309–D314 (2019).

902   55.  Heffernan, B., Murphy, C. D. & Casey, E. Comparison of planktonic and biofilm cultures of

903         Pseudomonas fluorescens DSM 8341 cells grown on fluoroacetate. *Appl. Environ.*

904         *Microbiol.* **75**, 2899–2907 (2009).

905   56.  Scales, B. S., Dickson, R. P., LiPuma, J. J. & Huffnagle, G. B. Microbiology, genomics, and

906         clinical significance of the Pseudomonas fluorescens species complex, an unappreciated

907         colonizer of humans. *Clin. Microbiol. Rev.* **27**, 927–948 (2014).

908   57.  Steinegger, M. & Söding, J. Clustering huge protein sequence sets in linear time. *Nat.*

909         *Commun.* **9**, 2542 (2018).

910   58.  Francino, M. P. The ecology of bacterial genes and the survival of the new. *Int. J. Evol.*

911         *Biol.* **2012**, 394026 (2012).

912    59. Muller, E. E. L. Determining Microbial Niche Breadth in the Environment for Better

913        Ecosystem Fate Predictions. *mSystems* **4**, (2019).

914    60. Roumpeka, D. D., Wallace, R. J., Escalettes, F., Fotheringham, I. & Watson, M. A Review

915        of Bioinformatics Tools for Bio-Prospecting from Metagenomic Sequence Data. *Front.*

916        *Genet.* **8**, 23 (2017).

917    61. Ivanova, N. N. *et al.* Stop codon reassignments in the wild. *Science* **344**, 909–913 (2014).

918    62. Steinegger, M., Mirdita, M. & Söding, J. Protein-level assembly increases protein sequence

919        recovery from metagenomic samples manyfold. *Nat. Methods* **16**, 603–606 (2019).

920    63. Titus Brown, C. *et al.* Exploring neighborhoods in large metagenome assembly graphs

921        reveals hidden sequence diversity. *bioRxiv* 462788 (2018) doi:10.1101/462788.

922    64. Höps, W., Jeffryes, M. & Bateman, A. Gene Unprediction with Spurio: A tool to identify

923        spurious protein sequences. *F1000Res.* **7**, 261 (2018).

924    65. Heinzinger, M. *et al.* Modeling aspects of the language of life through transfer-learning

925        protein sequences. *BMC Bioinformatics* **20**, 723 (2019).

926    66. Breitwieser, F. P., Pertea, M., Zimin, A. & Salzberg, S. L. Human contamination in bacterial

927        genomes has created thousands of spurious proteins. *Genome Res.* (2019)

928        doi:10.1101/gr.245373.118.

929    67. Steinegger, M. & Salzberg, S. L. Terminating contamination: large-scale search identifies

930        more than 2,000,000 contaminated entries in GenBank. *Genome Biol.* **21**, 115 (2020).

931    68. Chen, L.-X., Anantharaman, K., Shaiber, A., Eren, A. M. & Banfield, J. F. Accurate and

932        complete genomes from metagenomes. *Genome Res.* **30**, 315–333 (2020).

933    69. Thomas, A. M. & Segata, N. Multiple levels of the unknown in microbiome research. *BMC*

934        *Biol.* **17**, 48 (2019).

935    70. Duarte, C. M. Seafaring in the 21St Century: The Malaspina 2010 Circumnavigation

936        Expedition. *Limnol. Oceanog. Bull.* **24**, 11–14 (2015).

937    71. Rusch, D. B. *et al.* The Sorcerer II Global Ocean Sampling Expedition: Northwest Atlantic

938       through Eastern Tropical Pacific. *PLoS Biol.* **5**, 1–34 (2007).

939    72. Lloyd-Price, J. *et al.* Strains, functions and dynamics in the expanded Human Microbiome

940       Project. *Nature* **550**, 61–66 (2017).

941    73. Sanger, F., Nicklen, S. & Coulson, A. R. DNA sequencing with chain-terminating inhibitors.

942       *Proc. Natl. Acad. Sci. U. S. A.* **74**, 5463–5467 (1977).

943    74. Köster, J. Reproducible data analysis with Snakemake. *F1000Res.* **7**, (2018).

944    75. Hyatt, D. *et al.* Prodigal: prokaryotic gene recognition and translation initiation site

945       identification. *BMC Bioinformatics* **11**, 119–119 (2010).

946    76. Mendler, K. *et al.* AnnoTree: visualization and exploration of a functionally annotated

947       microbial tree of life. *Nucleic Acids Res.* **47**, 4442–4448 (2019).

948    77. Eberhardt, R. Y. *et al.* AntiFam: a tool to help identify spurious ORFs in protein annotation.

949       *Database* **2012**, bas003–bas003 (2012).

950    78. Yooseph, S., Li, W. & Sutton, G. Gene identification and protein classification in microbial

951       metagenomic sequence data via incremental clustering. *BMC Bioinformatics* **9**, 1–13

952       (2008).

953    79. Finn, R. D., Clements, J. & Eddy, S. R. HMMER web server: interactive sequence similarity

954       searching. *Nucleic Acids Res.* **39**, W29–W37 (2011).

955    80. Finn, R. D. *et al.* The Pfam protein families database: towards a more sustainable future.

956       *Nucleic Acids Res.* **44**, D279–D285 (2016).

957    81. Bennett, K. D. Determination of the number of zones in a biostratigraphical sequence. *New*

958       *Phytol.* **132**, 155–170 (1996).

959    82. Daily, J. Parasail: SIMD C library for global, semi-global, and local pairwise sequence

960       alignments. *BMC Bioinformatics* **17**, 81–81 (2016).

961    83. Žure, M., Fernandez-Guerra, A., Munn, C. B. & Harder, J. Geographic distribution at

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

963       sediments. *ISME J.* **11**, 478–489 (2017).

964    84. Chafee, M. *et al.* Recurrent patterns of microdiversity in a temperate coastal marine

965         environment. *ISME J.* **12**, 237–252 (2018).

966    85. Csardi, G. & Nepusz, T. The igraph software package for complex network research.

967         *InterJournal* vol. Complex Systems 1695 (2006).

968    86. Deorowicz, S., Debudaj-Grabysz, A. & Gudyś, A. FAMSA: Fast and accurate multiple

969         sequence alignment of huge protein families. *Sci. Rep.* **6**, 33964–33964 (2016).

970    87. Vanhoutreve, R. *et al.* LEON-BIS: multiple alignment evaluation of sequence neighbours

971         using a Bayesian inference system. *BMC Bioinformatics* **17**, 271–271 (2016).

972    88. Jehl, P., Sievers, F. & Higgins, D. G. OD-seq: outlier detection in multiple sequence

973         alignments. *BMC Bioinformatics* **16**, 269–269 (2015).

974    89. Broder, A. Z. On the resemblance and containment of documents. in *Proceedings.*

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

976         1997).

977    90. The UniProt Consortium. UniProt: the universal protein knowledgebase. *Nucleic Acids Res.*

978         **45**, D158–D169 (2017).

979    91. NCBI Resource Coordinators. Database resources of the National Center for Biotechnology

980         Information. *Nucleic Acids Res.* **46**, D8–D13 (2018).

981    92. Hingamp, P. *et al.* Exploring nucleo-cytoplasmic large DNA viruses in Tara Oceans

982         microbial metagenomes. *ISME J.* **7**, 1678–1695 (2013).

983    93. Remmert, M., Biegert, A., Hauser, A. & Söding, J. HHblits: Lightning-fast iterative protein

984         sequence searching by HMM-HMM alignment. *Nat. Methods* **9**, 173–175 (2012).

985    94. UniProt Consortium, T. UniProt: the universal protein knowledgebase. *Nucleic Acids Res.*

986         **46**, 2699 (2018).

987    95. Dick, J. M. Calculation of the relative metastabilities of proteins using the CHNOSZ

988         software package. *Geochem. Trans.* **9**, 10 (2008).

989    96. Hausser, J. & Strimmer, K. Entropy inference and the James-Stein estimator, with

990      application to nonlinear gene association networks. *arXiv [stat.ML]* (2008).

991    97. van Dongen, S. & Abreu-Goodger, C. Using MCL to Extract Clusters from Networks. in

992      *Bacterial Molecular Networks: Methods and Protocols* (eds. van Helden, J., Toussaint, A. &

993      Thieffry, D.) 281–295 (Springer New York, 2012).

994    98. Li, W. & Godzik, A. Cd-hit: a fast program for clustering and comparing large sets of protein

995      or nucleotide sequences. *Bioinformatics* **22**, 1658–1659 (2006).

996    99. Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large

997      phylogenies. *Bioinformatics* **30**, 1312–1313 (2014).

998    100. Berger, S. A. & Stamatakis, A. PaPaRa 2.0: a vectorized algorithm for probabilistic

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

1000      *http://sco. h-its. org/exelixis/publica tions. html. Exelixis-RRDR-2012-2015* (2012).

1001    101. Matsen, F. A., Kodner, R. B. & Armbrust, E. V. pplacer: linear time maximum-likelihood and

1002      Bayesian phylogenetic placement of sequences onto a fixed reference tree. *BMC*

1003      *Bioinformatics* **11**, 538 (2010).

1004    102. Needham, D. M. *et al.* A distinct lineage of giant viruses brings a rhodopsin photosystem to

1005      unicellular marine predators. *Proc. Natl. Acad. Sci. U. S. A.* **116**, 20574–20583 (2019).

1006    103. Murat Eren, A. *et al.* Anvi'o: an advanced analysis and visualization platform for 'omics

1007      data. *PeerJ* **3**, e1319 (2015).

1008    104. Li, H. & Durbin, R. Fast and accurate long-read alignment with Burrows–Wheeler transform.

1009      *Bioinformatics* **26**, 589–595 (2010).

1010    105. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic

1011      features. *Bioinformatics* **26**, 841–842 (2010).

1012    106. Levins, R. THE STRATEGY OF MODEL BUILDING IN POPULATION BIOLOGY. *Am. Sci.*

1013      **54**, 421–431 (1966).

1014    107. Bray, J. R., Roger Bray, J. & Curtis, J. T. An Ordination of the Upland Forest Communities

1015      of Southern Wisconsin. *Ecological Monographs* vol. 27 325–349 (1957).

1016    108. Langfelder, P., Zhang, B. & Horvath, S. Defining clusters from a hierarchical cluster tree:

1017        the Dynamic Tree Cut package for R. *Bioinformatics* **24**, 719–720 (2008).

1018    109. Miklós, I. & Podani, J. RANDOMIZATION OF PRESENCE–ABSENCE MATRICES:

1019        COMMENTS AND NEW ALGORITHMS. *Ecology* vol. 85 86–92 (2004).

1020    110. Salazar, G. *et al.* Particle-association lifestyle is a phylogenetically conserved trait in

1021        bathypelagic prokaryotes. *Mol. Ecol.* **24**, 5692–5706 (2015).