# Water mass specific genes dominate the Southern Ocean microbiome

**Emile Faure**

emile.faure@univ-brest.fr

Station Biologique de Roscoff, Sorbonne Université

**Jolann Pommellec**

Université de Bretagne Occidentale

**Cyril Noel**

IFREMER

**Alexandre Cormier**

IFREMER

**Lisa-Marie Delpech**

Université de La Rochelle    https://orcid.org/0009-0006-0154-0211

**Murat Eren**

Helmholtz institute for functional marine biodiversity

**Antonio Fernandez-Guerra**

Lundbeck GeoGenetics Centre, The Globe Institute, University of Copenhagen    https://orcid.org/0000-0002-8679-490X

**Chiara Vanni**

Max Planck Institute for Marine Microbiology    https://orcid.org/0000-0002-1124-1147

**Marion Fourquez**

MIO, Aix Marseille Université

**Marie-Noëlle Houssais**

Sorbonne Université

**Corinne Da Silva**

Commissariat à l'Energie Atomique

**Frederick Gavory**

GENOSCOPE, CEA

**Aude Perdereau**

GENOSCOPE, CEA

**Karine Labadie**

GENOSCOPE, CEA

**Patrick Wincker**

Génomique Métabolique, Genoscope, Institut François Jacob, CEA, CNRS, Université d'Evry, Université Paris-Saclay,    https://orcid.org/0000-0001-7562-3454

**Julie Poulain**

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

**Christel Hassler**

  Ecole Polytechnique Fédérale de Lausanne

**Yajuan Lin**

  Texas A&M Corpus Christi    https://orcid.org/0000-0002-9057-9321

**Nicolas Cassar**

  Duke University

**Lois Maignien**

  Department of Biology, University of Southern Denmark, 5230 Odense M, Denmark

---

**Article**

**Keywords:**

**Additional Declarations:** There is **NO** Competing Interest.

---

# Water mass specific genes dominate the Southern Ocean microbiome

*Emile Faure[1,2]\*, Jolann Pommellec[1], Cyril Noel[3], Alexandre Cormier[3], Lisa-Marie Delpech[1,4,5], Murat A Eren[6,7,8,9,10], Antonio Fernandez-Guerra[11,12], Chiara Vanni[13], Marion Fourquez[14], Marie-Noëlle Houssais[15], Corinne da Silva[16], Frederick Gavory[16], Aude Perdereau[16], Karine Labadie[16], Patrick Wincker[17,18], Julie Poulain[17,18], Christel Hassler[19,20,21], Yajuan Lin[22,23,24], Nicolas Cassar[23,24]\*, Loïs Maignien[1,10]\**

*1 - Univ Brest (UBO), CNRS, IFREMER, Laboratoire de Microbiologie des Environnements Extrêmes, Plouzané, 29280, France*

*2 - Sorbonne Université, Centre National de la Recherche Scientifique, UMR 7144 Adaptation and Diversity in the Marine Environment, Station Biologique, Roscoff 29680, France*

*3 - IFREMER, IRSI – Service de Bioinformatique (SeBiMER) Plouzané, 29280 France*

*4 - Littoral Environnement et Sociétés (LIENSs), UMRi 7266 CNRS - La Rochelle Université, La Rochelle 17000, France*

*5 - Department of Biology, École Normale Supérieure de Lyon, Lyon 69007, France*

*6 - Helmholtz Institute for Functional Marine Biodiversity, 26129 Oldenburg, Germany*

*7 - Alfred Wegener Institute Helmholtz Centre for Polar and Marine Research, 27570 Bremerhaven, Germany*

*8 - Institute for Chemistry and Biology of the Marine Environment, University of Oldenburg, 26129 Oldenburg, Germany*

*9 - Marine 'Omics Bridging Group, Max Planck Institute for Marine Microbiology, 28359 Bremen, Germany*

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

*11 - Centre for Ancient Environmental Genomics, Globe Institute, University of Copenhagen, Copenhagen, Denmark*

*12 - Lundbeck Foundation GeoGenetics Centre, Globe Institute, University of Copenhagen, Copenhagen, Denmark*

*13 - MARUM Center for Marine Environmental Sciences, University of Bremen,28359 Bremen, Germany*

*14 - Aix Marseille Univ., Université de Toulon, CNRS, IRD, MIO UMR 110, Marseille 13288, France.*

*15 - Laboratoire d'Océanographie et du Climat (LOCEAN), CNRS-Sorbonne Université, Paris, France*

*16 - Genoscope, Institut François Jacob, CEA, Université Paris-Saclay, Evry, 91057, France.*

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

*18 - Research Federation for the Study of Global Ocean Systems Ecology and Evolution, FR2022/Tara Oceans GO-SEE, CNRS, Paris, France.*

*19 - Swiss Polar Institute, Ecole Polytechnique Fédérale de Lausanne, 1015 Lausanne, Switzerland*

*20 - School of Architecture, Civil, and environmental engineering, Smart Environmental Sensing in Extreme Environments, ALPOLE, Ecole Polytechnique Fédérale de Lausanne, Sion, Switzerland.*

*21 - Institute of Earth Sciences, University of Lausanne, Lausanne, Switzerland.*

*22 - Department of Life Sciences, Texas A&M University – Corpus Christi, 78412 Corpus Christi, Texas, USA*

*23 - Division of Earth and Climate Sciences, Nicholas School of the Environment, Duke University, 27710 Durham, North Carolina, USA*

*24 - CNRS, Université de Brest, IRD, Ifremer, LEMAR, Plouzané, France*

*Corresponding authors:* emile.faure@sb-roscoff.fr; lois.maignien@univ-brest.fr; Nicolas.Cassar@duke.edu
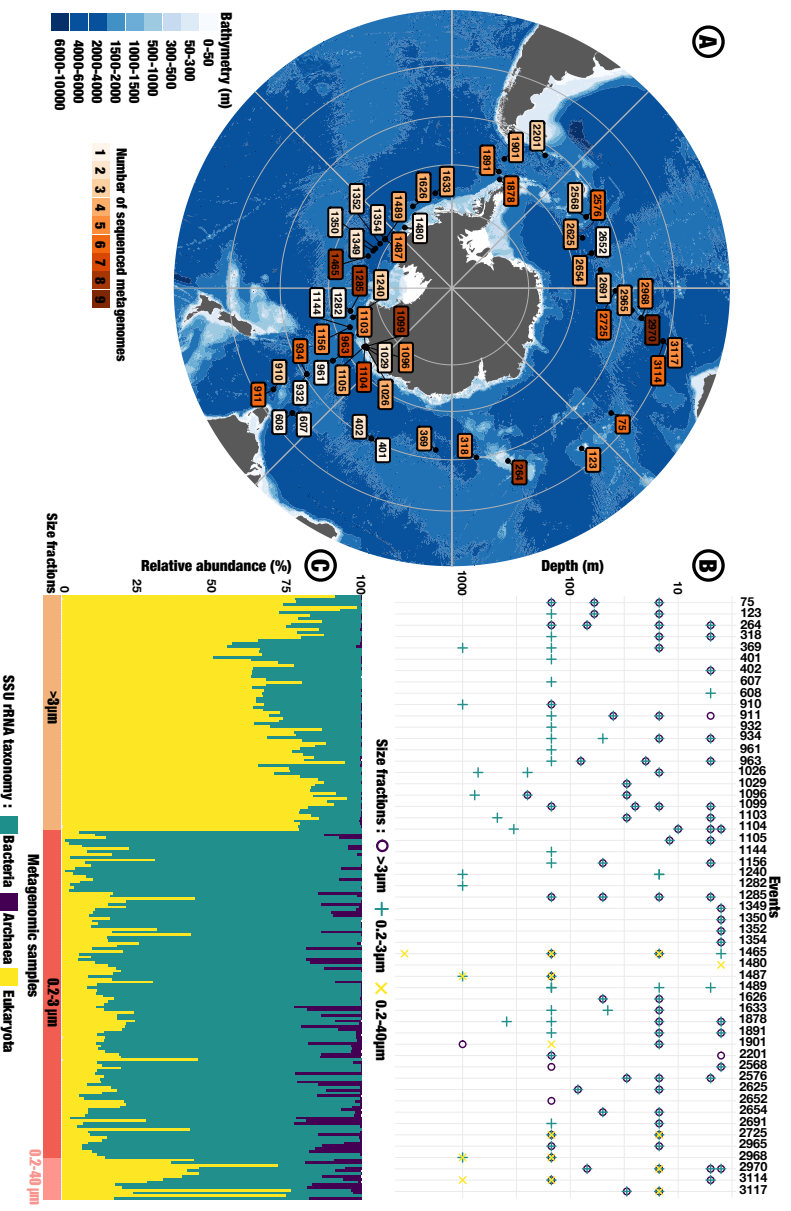
## ABSTRACT

The Southern Ocean (SO) plays a key role in regulating global biogeochemical cycles and climate, yet microbial genes sustaining its biological activity remain poorly characterized. We introduce a comprehensive SO microbial genes collection from 218 metagenomes sampled during the Antarctic Circumnavigation Expedition, the majority of which are missing from functional databases. 38% even lack homologs in current reference marine gene catalogs, defining a singular genetic seascape. We show that SO gene assemblages exhibit a common polar signature with the Arctic Ocean while being structured by water masses at the SO-scale. We analyze genomic markers of diverse SO biomes, focusing on adaptations to organic matter consumption in the blooming Mertz polynya and temperature-dependent trace metal utilization by the ubiquist Bacteria Pelagibacter. Our work takes a step towards a more comprehensive understanding of SO's plankton ecology and evolution, capturing the current state of the unique microbial diversity in this rapidly changing Ocean.

## INTRODUCTION

52  The Southern Ocean (SO) dominates other oceans in heat and carbon uptakes while being
53  particularly exposed to climate change impacts[1]. It is mainly composed of high nutrient low
54  chlorophyll waters (HNLC) where phytoplanktonic growth is limited by trace elements such as
55  iron or manganese[2,3]. In presence of these elements, phytoplankton blooms can reach
56  concentrations of $10^8$ cells per liter[4], playing a key role in carbon sequestration through the
57  biological pump[5]. Beyond phytoplankton, the extent of carbon export is impacted by the
58  consumption and remineralization of organic matter by communities of bacteria and
59  archaea[6,7]. Yet, *in situ* abundance and diversity of microbial communities in the SO remain
60  poorly described.

61  Recent large-scale environmental metagenomics projects highlighted the rich functional and
62  taxonomic diversity of marine plankton and the driving effect of environmental conditions on
63  planktonic communities[8–11]. However, only two sampling locations in the SO were included in
64  recent efforts to compute global genes and genomes catalogs[12,13], underscoring the
65  substantial undersampling of this critical ocean. A study focusing on polar oceans and
66  including 21 metagenomics samples from the SO allowed the construction of a first polar gene
67  catalog in 2020, showing the high prevalence of polar specific genes in the SO[14]. Yet, we still
68  lack a realistic census of SO's microbial diversity and of the environmental factors structuring
69  its planktonic communities. We address this important knowledge gap, identifying drivers of
70  planktonic functional and taxonomic diversity in this area subject to major environmental
71  changes[1].

72  The Antarctic Circumnavigation Expedition (ACE) circumnavigated the Southern Ocean
73  during the 2016-2017 austral summer, producing an unprecedented amount of physical,
74  chemical and biological observations[15]. Analyzing 218 metagenomes, we increase by an order
75  of magnitude the number of SO samples ever considered in a meta-omics study to present
76  the first SO-specific gene catalog (Fig. 1). Building on the seminal work of previous global
77  metagenomics efforts[16–18], we first demonstrate the broadscale uniqueness of the SO
78  compared to other oceans, before diving into its regional variability. To exemplify the
79  uniqueness of biomes in the SO, we focus on the genomic signature of specialist species
80  occurring in the Mertz polynya, before using SAR11 as a case study of genomic adaptations
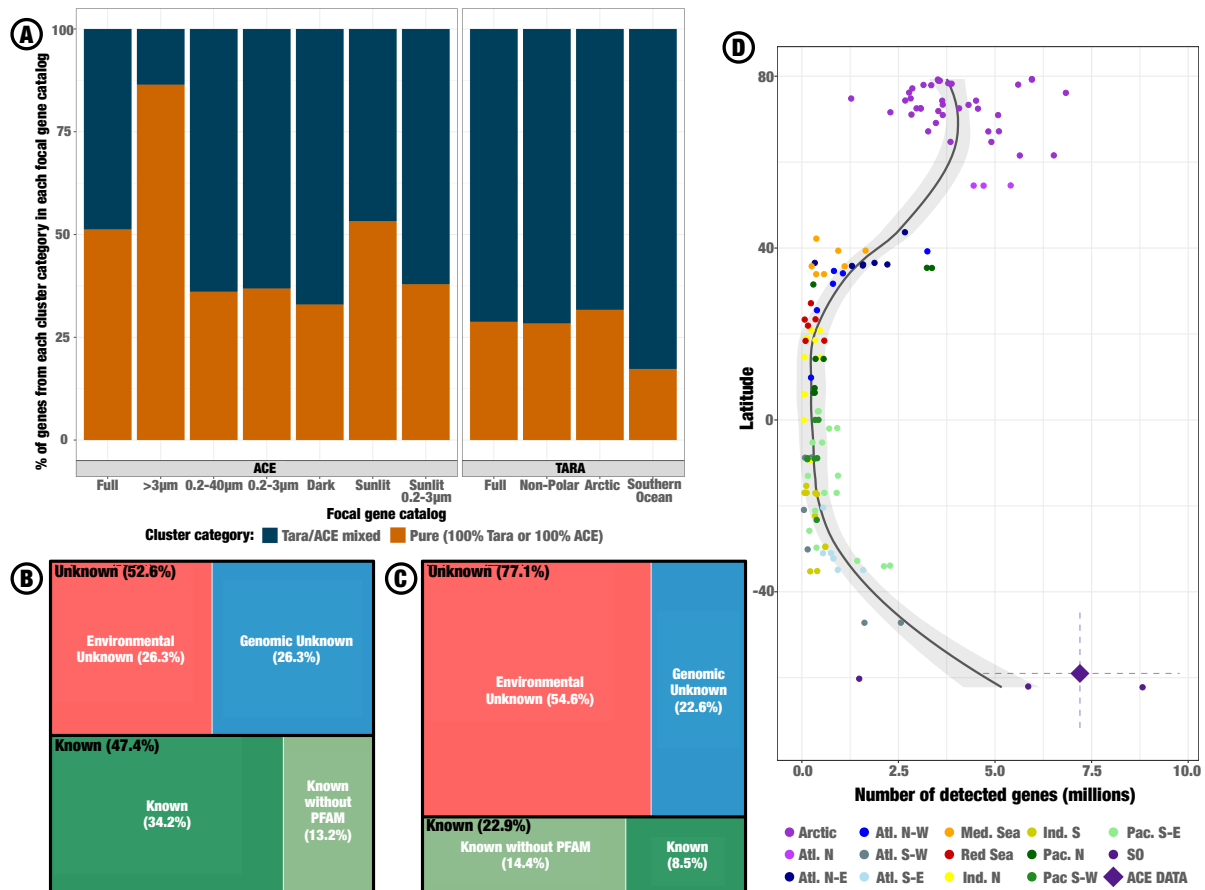81  to polar conditions in a ubiquitous taxon.

*Figure 1: Overview of the ACE metagenomics dataset. (A) Map of CTD downcast events on which metagenomics samples were taken, colored according to the number of samples taken on the cast. (B) Depth and size fraction chart of all ACE metagenomes. Each vertical line corresponds to an event as pictured on the map in A, and each metagenome is represented as a dot at its corresponding depth, with the shape indicating the size fraction. (C) Relative abundance of each domain of life in every sampled metagenome, as estimated through SSU rRNA reconstruction from metagenomics short reads. Samples were separated according to their size fraction on the x axis, as indicated on the bottom layer.*

## RESULTS

### Broadscale novelty of Southern Ocean's microbial genes

*An unsuspected genomic diversity at SO scale*

Individual assemblies of the 218 metagenomes (Fig.1, Fig.S1) produced 68,074,004 contigs in which we identified 175,336,776 Open-Reading Frames (ORFs). We dereplicated these ORFs using 95% similarity and 90% coverage thresholds, producing 89,739,060 dereplicated-genes hereafter called unigenes. 51.3% of ACE unigenes did not cluster with any unigene from the most recent *Tara Ocean and Polar Circle* gene catalog[16] (OM-RGC-v2) at thresholds of 30% similarity and 80% coverage in amino acid sequence (Fig. 2A). This number remained at 37.9% accounting only for ACE unigenes assembled from the 0.2-3μm size fraction in the sunlit layer (Fig. 2A). Conversely, 28.9% of OM-RGC v2 unigenes did not cluster with any ACE unigene using the same thresholds. This strong mutual exclusion between the two catalogs highlights the originality of the SO as compared to other oceans, especially considering that the OM-RGC-v2 includes unigenes assembled from Arctic metagenomes.

*Figure 2: Novelty of Southern Ocean microbial genes. (A) Distribution of genes from either the ACE unigene catalog (left box) or the OM-RGC v2 (right box) into pure or mixed gene clusters. Genes from both catalogs were clustered at 30% similarity and 80% coverage thresholds in amino acid sequences, then classified as either belonging to a pure cluster, i.e. a cluster only containing genes from one catalog, or a mixed one, i.e. a cluster mixing genes from both catalogs. Results are either presented on the full catalogs, or restrained to specific gene subsets involving size fractions, depth, and geography, as described on the x axis. (B) Chart of AGNOSTOS annotations at gene level, i.e. accounting for the number of genes in each cluster annotation category: Environmental Unknowns (EU) lack functional annotations and are absent from any genome recorded in the AGNOSTOS database, Genomic Unknowns (GU) also lack functional annotations yet are recorded in a genomic context in the AGNOSTOS database, Knowns are functionally annotated, either with (K) or without PFAM (KWP) annotations. (C) Chart of AGNOSTOS annotations at AGNOSTOS gene cluster level. (D) Latitudinal gradient of ACE genes' detection in Tara Oceans and Polar Circle samples. Genes were considered as detected if at least 60% of their sequence was covered with a depth of 1X or more. The mean number of detected genes in ACE samples is indicated by the diamond shaped point, with the horizontal dashed line spanning from first to third quartile of detected gene number and the vertical one from minimum to maximum latitudes. A loess curve was fitted to the number of detected genes in Tara samples, not taking into account ACE samples.*

We further explored distant gene homology using AGNOSTOS[19]. We clustered the 175,336,776 ORFs into 30,123,228 AGNOSTOS gene clusters (AGC), of which 64.8% were singletons, 32.6% were good-quality clusters of multiple ORFs as per AGNOSTOS standards, and 2.5% were discarded as low-quality clusters. 52.6% of the ORFs were tagged as unknowns (*i.e.* without functional annotation) and contributed to 77.1% of all AGC, illustrating the high prevalence of singletons among unknown ORFs compared to known ones, which
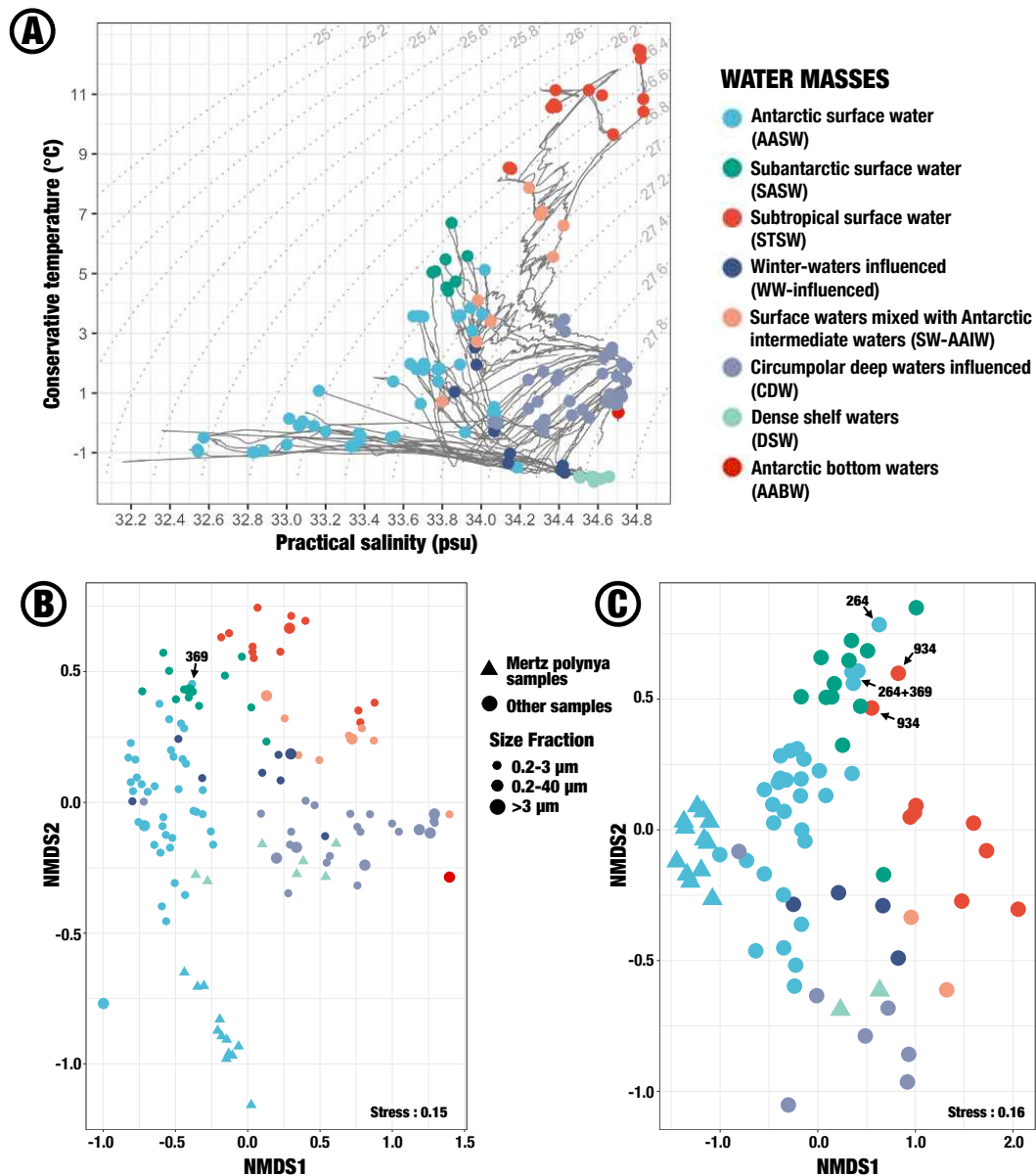
131    clustered better together (Fig. 2B,2C). The asymptotic nature of collector curves drawn at AGC
132    level suggests that the ACE AGC catalog covers most of SO genomic diversity (Figure S2).

133    *Bipolar distribution of Southern Ocean microbial genes*

134    Adaptation to polar conditions is thought to be responsible for a high genomic similarity
135    between Arctic and Antarctic microbiomes despite dispersal isolation[14]. To quantify this bipolar
136    pattern among ACE genes, we mapped 134 *Tara Oceans* (TO) and *Tara Polar Circle* (TPC)
137    metagenomes covering most subtropical and arctic oceanic regions onto ACE contigs. The
138    resulting detection matrix shows a bipolar distribution of SO-genes at global scale (Fig. 2D).
139    The mean number of detected ORF per ACE sample was of 7,198,029, while it was of
140    3,956,716 in TPC samples, illustrating the high level of endemicity of the SO despite
141    similarities in gene content between poles. This mean dropped to 334,896 ORFs for non-polar
142    TO samples. Of the 34,344,531 ACE ORFs detected in at least one sample from the Arctic
143    Ocean, 26,353,298 were absent from all non-polar oceans sampled during TO and therefore
144    identified as polar-specific. Polar specific ORFs were distributed in 14,426,012 unigenes and
145    4,105,973 AGC clusters, of which 61.8% were unknown (39% environmental unknown, 22.8%
146    genomic unknown). We identified 4,314 EggNOG functions significantly enriched in polar
147    specific unigenes compared to the rest of ACE unigenes (over a total of 54,772 functions,
148    unilateral Fisher tests, adjusted p-value < 0.01). The six functions with the highest odds ratio,
149    ranging between 4.0 and 4.3 in favor of polar-specific unigenes, were *Formate dehydrogenase*
150    *(NAD+) activity*, *Excinuclease ABC* (UV-specific endonuclease), *Septum formation initiator*,
151    *cold-shock protein*, *oxidoreductase activity acting on the aldehyde or oxo group of donors,*
152    *iron-sulfur protein as acceptor* and *Iron-binding zinc finger CDGSH type* (See Table S1 to
153    access complete list of enriched functions).

154    **The SO hosts a diversity of unique microbial biomes shaped by oceanographic fronts**
155    **and phytoplankton blooms**

156    We analyzed AGC's biogeography following three steps, (1) an unconstrained analysis of
157    AGC's distribution across samples, (2) a univariate exploration of each AGC to detect those
158    linked to the environment (env-AGC) and (3) a grouping of env-AGCs into co-abundant groups
159    to allow a multivariate exploration of their response to environmental gradients. We worked
160    independently on the free-living (0.2-40 + 0.2-3 µm) and >3µm size fractions considering their
161    different taxonomic profiles (Figure 1C). We focused exclusively on AGCs with non-repeated
162    coverage values in at least 20% of samples (1,906,624 and 2,437,988 clusters in the free-
163    living and >3µm size fractions, resp.), avoiding rare AGCs as well as AGCs with uniform
164    distribution across samples.

*Figure 3: Microbial genes assemblages of the Southern Ocean are water mass specific. (A) Temperature – salinity diagram based on ACE downcast CTD data[20]. Each grey line corresponds to a CTD cast. Dots correspond to depths at which seawater was sampled for metagenomic libraries construction, colored according to their attributed water masses. Dotted lines in the background correspond to isopycnals. (B,C) NMDS computed on AGC abundance matrices of free-living (B) and >3µm (C) size fractions, colored according to their water masses using the same color legend as in (A). Positions of samples in B and C are only determined by their composition in AGC. The event numbers of samples taken above 150m are indicated by black arrows when their positions do not match their water mass classification, as discussed in the results. Events 369 and 264 were taken right on the Polar Front, while event 934 was taken on the Sub-Antarctic Front.*

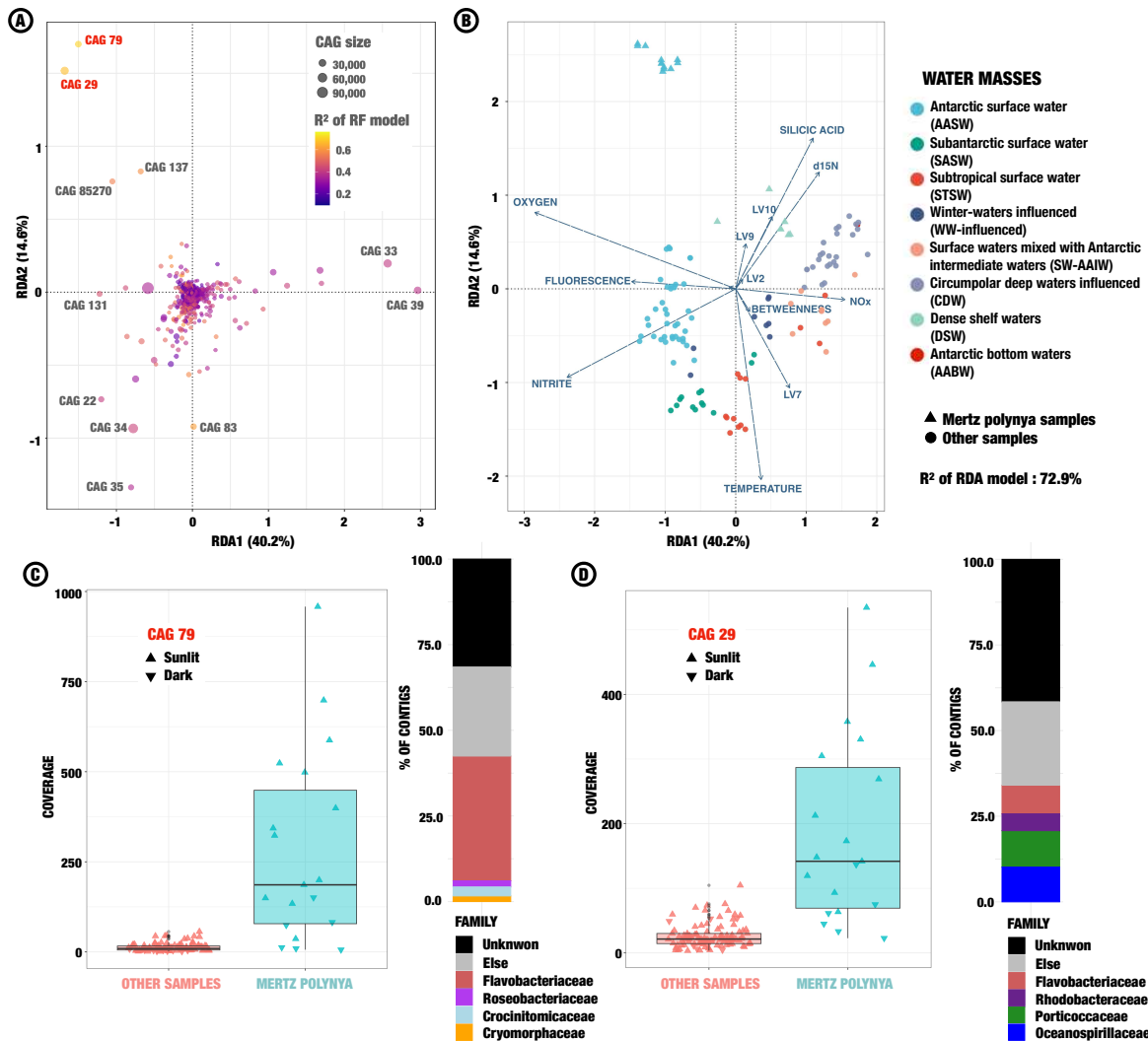### Microbial gene assemblages are distinct across water masses at SO scale

The classification of our samples in water masses based on temperature-salinity-oxygen diagrams was the best grouping variable for predicting AGC abundance (Figure 3, S3). AGC assemblages were significantly distinct across water masses in the both size fractions (Figure 3B, PERMANOVA with 999 permutations, p-value<0.001). Sub-Antarctic surface waters (SASW) samples were separated from both Antarctic surface waters (AASW) samples and

182    sub-tropical surface waters (STSW) samples, suggesting biogeographical barriers at both the
183    Sub-Antarctic Front and the Polar Front. This was already observed for Flavobacteria[21] but
184    lacked confirmation on broader taxonomic range[22]. All surface samples showing potential
185    mismatches between their AGC assemblage and their attributed water masses came from
186    events located on oceanographic fronts (Figure 3B,3C), suggesting a potential mixing of
187    microbial assemblages at these fronts.

188    Samples from Circumpolar Deep Waters (CDW) were well separated from surface water
189    masses in both size fractions, while surface waters influenced by colder (Winter Water-
190    influenced and Antarctic Intermediate-influenced) layers appeared between CDW and surface
191    waters. Samples from Dense Shelf Waters (DSW) were mostly similar to CDW samples,
192    except the two shallowest DSW samples which appeared closer to AASW in genomic
193    composition for the free-living size fraction (Figure 3B). These results suggest a diminution of
194    AGC diversity in deep water masses (see Supplementary Results). Still, the only Antarctic
195    Bottom Water sample (AABW, 3460m depth) had the most extreme coordinate on the NMDS
196    first axis (Figure 3B), suggesting a unique genomic composition. In light of this uniqueness
197    and considering the projected decrease in AABW formation due to increasing influence of
198    meltwater[23] from Antarctica, a better characterization of the functional roles from AABW
199    microbial populations is urgently needed.

200    *Identifying genomic markers following environmental gradients at SO-scale*

201    We built random forest regression models for each AGC, predicting coverage using 50
202    environmental predictors from ACE metadata. We defined $R^2$ thresholds based on permuted
203    repetitions of the analysis to only consider AGCs linked to the environment (env-AGC: $R^2 >$
204    10% in the free-living, 15% in the attached size fraction, Figure S4). 89.0% (resp. 82.1%) of
205    the considered AGC were env-AGC in the free-living (resp. >3µm) size fraction. Over both
206    size fractions, 894,292 models (20.6%) showed $R^2$ values above 50%, indicating predictability
207    of AGC abundance based only on the environmental context and opening the way for
208    genomic-based correlative models at SO scale[9,10]. To analyze env-AGC in a multivariate
209    context, we grouped them into 156,671 and 28,756 co-abundant groups (CAGs)[24] in the free-
210    living and >3µm size fractions, respectively. We then identified CAGs of interest associated
211    with various biomes through a redundancy analysis (Figure 4). We first present CAGs specific
212    to the Mertz polynya, before focusing on 3 CAGs illustrating a gradient of polar adaptation
213    across latitudes. In the supplementary materials, we describe two CAGs linked with specialist
214    species thriving in polar conditions (CAGs 131 and 34), two CAGs associated with deep water
215    masses (CAGs 33 and 39), as well as outliers at the SO scale, including CAGs specific to sub-
216    Antarctic islands (CAGs 136, 73614 and 177401).

217
218 *Figure 4: The response of co-abundant groups of env-AGCs (CAGs) to environmental gradients at SO*
219 *scale highlights Mertz polynya's originality. Redundancy analysis of CAGs abundance in response to*
220 *environmental variables, in the free-living size fraction (A,B). RDA triplot was separated in two parts for*
221 *better readability, (A) showing the distribution of CAGs in the RDA space, colored according to their*
222 *mean random forest R-squared (reflecting the predictability of their abundance using environmental*
223 *data). The size of each dot corresponds to the size of the CAG, in number of env-AGC. The different*
224 *CAGs of interest mentioned in this study are indicated with grey labels, while the two CAGs plotted in*
225 *(C) and (D) are highlighted in red. (B) shows samples and environmental variables distribution in the*
226 *same RDA space. Samples are colored according to their water mass. The first axis of the RDA*
227 *opposed surface samples with high fluorescence and oxygen (RDA1<0) from deep samples showing*
228 *high NOx concentrations (RDA1>0). The second axis was driven by temperature, opposing warm*
229 *STSW samples (RDA2<0) from colder samples, and isolating all AASW samples from Mertz as an*
230 *outlier group (RDA2>0). LV stands here for latent variables, corresponding to the ones described in*
231 *Landwehr et al.[15]. LV2 is linked with cloud condensation, LV7 with seasonal signal, LV9 is linked to*
232 *marginal sea ice zone and snowfall and LV10 to the dial cycle. A similar RDA triplot for the >3µm size*
233 *fraction is presented in Figure S5. Mertz polynya's originality is further illustrated in (C) and (D), showing*
234 *the abundance and taxonomy of the two CAGs most linked to it. Boxplots corresponding to each CAG's*
235 *coverage are plotted at Mertz versus in other samples, with each individual sample plotted as points*
236 *shaped according to categorical depth: sunlit (150m and above) and dark (below 150m). Family-level*
237 *taxonomic profiles are represented next to each boxplot, as estimated through contigs taxonomic*
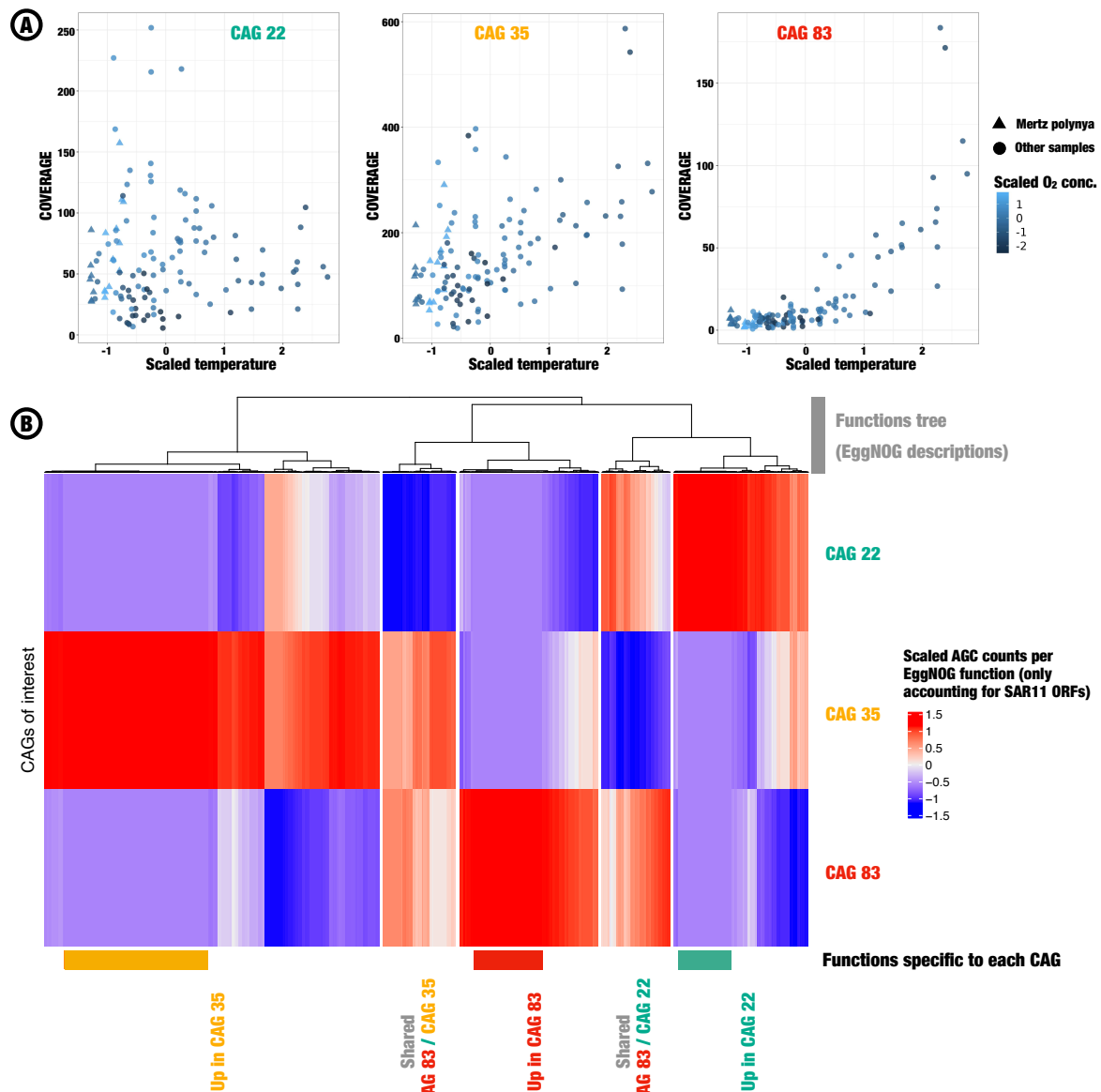
240    *The genomic signature of an active diatom bloom in the Mertz polynya*

241 Four CAGs from the free-living size fraction were enriched in Mertz samples (RDA2 > 0.5,
242 Figure 4A): CAG 79, CAG 29, CAG 137 and CAG 85270 (Figure 4C,4D). They contained
243 below 16% of *environmental unknowns (EU)* and more than 50% of *known (K)* and *known*
244 *without PFAM (KWP)*, when 54.6% of all AGCs were annotated as EU. All four CAGs were
245 significantly enriched (Fisher test, adj. p-value < 0.01) in *TonB* receptors and *TonB*-linked
246 outer membrane proteins specialized in the import of degradation products from proteins or
247 carbohydrates as nutrients (*SusC/RagA*, *SusD/RagB*). They were also enriched in proteins
248 involved in the carbohydrate metabolism, *e.g.* glycosyl transferase, polysaccharide
249 biosynthesis protein or glutamine synthetase, in ABC transporters and in proteins involved in
250 cell motility (*e.g.* gliding motility, morphogenesis and elongation of the flagellar filament). A
251 variety of metallo-protein were enriched in all four CAGs as well, including heme-binding
252 proteins, M6 family metalloprotease or metal-dependent hydrolase. Finally, the four CAGs
253 were enriched in phage integrase, and three out of four were enriched in phage plasmid
254 primase P4, suggesting a strong phage presence in the Mertz polynya. We provide a complete
255 list of enriched functions in each CAG of interest (Table S1).

256 Taxonomic profiles of all CAGs of interest were estimated based on contig-level taxonomic
257 annotation using the UniRef90 database[25] as reference (Figure 4C,4D). To increase
258 taxonomic precision, we also investigated annotations obtained on MAGs binned from the
259 same metagenomic assemblies of ACE data (Pommellec et al., *in prep.*). ORFs from CAGs
260 79 and 85270 were mainly affiliated with the *Polaribacter* genus in both MAG-based and
261 contigs-based annotations. In CAG 137, Rhodobacteraceae dominated both annotation types
262 but Roseobacteraceae were absent from MAGs while abundant in contig-based annotations.
263 In CAG 29, HTCC2207 SAR 92 was the most represented genus using MAGs-based
264 annotation, matching the Porticoccaceae dominance in contigs-based annotations (Figure
265 4D). SAR 92 is a widely distributed oligotrophic clade known for its ability to consume
266 polysaccharides in the epipelagic zone, notably through TonB-dependent receptors[26]. It has
267 been associated with late stages of diatom-induced bacterioplankton blooms in the North Sea,
268 uptaking and degrading specific polysaccharides including chrysolaminarin[27]. *Polaribacter* and
269 SAR 92 have both been associated with Phaeocystis-produced chrysolaminarin degradation
270 in the SO[28]. The second most represented genus in CAG 29 was ASP10-02a, which was
271 identified as the main cobalamin (Vitamin B12) producer in a coastal area of the SO, playing
272 a key role in primary production co-limitation by micronutrients[29].

273 Two CAGs from the >3μm size fraction were enriched in the Mertz polynya (Figure S5,
274 discussed in Supplementary results). Both were linked to *Fragilariopsis cylindrus*, an indicator
275 species of cold water evolutionarily adapted to the polar environment[30]. All CAGs enriched in
276 Mertz polynya samples were thus linked with organisms specifically adapted to polar blooms
277 conditions.

*Figure 5: Functional changes across latitudes in SAR11-related CAGs. (A) Abundance of three CAGs dominated by Candidatus* Pelagibacter *as a function of scaled temperature and oxygen. (B) Heatmap of scaled counts of AGCs per unique EggNOG functional annotations, taking only into account ORFs coming from contigs annotated as Candidatus* Pelagibacter. *The tree on top of the heatmap clusters each unique function according to its AGC count profile across CAGs, using euclidean distance and ward's D2 clustering method. The tree was manually cut to form 5 groups, differentiating functions shared by multiple CAGs from functions being more abundant in one CAG, as described by text labels below the heatmap. Functions found only in one of the three CAGs are highlighted with colored bars on the bottom layer of the heatmap.*

*Global latitudinal shifts across gene groups illustrate adaptation to polar conditions in the ubiquitous SAR11*

We identified three CAGs of interest with a distinct response to environmental data despite having similar taxonomic profiles (Figure 4A, 5). CAG 83 was specific to warm waters, with low abundance in all AASW samples supposing a latitudinal boundary at the polar front (Figure 5A). CAG 35 was positively correlated to temperature yet ubiquitous, even in the coldest waters. Finally, CAG 22 was present in all water masses, with higher abundances in cold AASW waters. CAGs 22, 35 and 83 were all dominated by Pelagibacteraceae, and

296   respectively 88.75%, 83.5% and 59.48% of their AGCs contained at least one ORF from a
297   contig annotated as *Candidatus* Pelagibacter (SAR11 and relatives). To compare the
298   functional annotations of SAR11 genes across the three CAGs, we extracted ORFs coming
299   from SAR11 contigs and compared EggNOG annotations between them at AGC-level (Figure
300   5). CAGs 22 and 83 showed opposite functional profiles, *i.e.* functions highly present in one
301   were rare or absent in the other, while CAG 35 shared functions with the two other CAGs,
302   matching its ubiquitous distribution (Figure 5B). A total of 430 unique functions showed an
303   increased presence in CAG22, among which 170 were only found in CAG22 (Figure 5B). The
304   most observed was the transmembrane NikM subunit, transporting nickel, which was carried
305   by SAR11 ORFs in 10 AGCs of CAG22 and none of CAG83 and CAG35. Among the other
306   functions with increased detection in CAG22, many were related to trace metals (zinc, iron),
307   sulfur (e.g. nucleotide-disulfide oxidoreductase, Sulfotransferase) and phosphorus cycles (e.g.
308   polyphosphate pyrophosphohydrolases, metal-dependent phosphohydrolase). See Table S1
309   to access the complete list.
310

311   **DISCUSSION**
312   Analyzing 218 metagenomes from a circumpolar expedition, we were able to characterize the
313   originality and biogeography of SO's microbial genomic diversity. We identified a set of genes
314   distributed at both poles while being absent from most other latitudes on the planet; some of
315   them are involved in adaptations to polar-specific constraints like UV-exposure and cold
316   temperatures. We illustrated how microbial gene assemblages from the SO are largely
317   endemic and unknown, both at taxonomic and functional level. The high number of singletons
318   in our assemblies suggests the presence of a significant proportion of rare genes in the SO,
319   showing no deep homologies with each other. This could partly be due to ORFs from the >3µm
320   size fraction, which should be treated with caution due to the difficulty of detecting good quality
321   ORFs from eukaryotic contigs[31] (*cf* Methods). Yet, 53.3% of all singletons came from
322   prokaryote-dominated samples of the free-living size fractions, and singletons from both size
323   fractions were largely dominated by unknowns (83.8% in free-living, 93.2% in >3µm),
324   suggesting a limited impact of eukaryotic contigs on our conclusions. Collector curves'
325   asymptotic profiles were stronger when decreasing detection thresholds (Figure S2),
326   suggesting that singletons do recruit reads in multiple samples independently of their size
327   fraction of origin. Otherwise, they would remain undetected in all or most samples whatever
328   the threshold, preventing the asymptotic form of the curve. It suggests they do share distant
329   homologies, *e.g.* at domain level, with unassembled genes across multiple samples. The
330   decrease in taxonomic diversity in the SO compared to subtropical latitudes[32] could then be
331   balanced by an abundance of diverse yet individually rare genomic elements distributed at SO
332   scale. This hypothesis will have to be confirmed by further explorations of ACE singletons, of
333   which the majority was excluded from our biogeographical analysis to focus on widely
334   distributed genes. This could be done through network-based methods allowing the
335   characterization of distant and rare homologs[33].
336   Gene assemblages were structured by water mass at SO scale, supporting the observation
337   that Processes leading to water mass formation and transport exert the strongest control on
338   microbial community composition[22]. Our statistical approach allowed us to identify genes
339   particularly abundant in the Mertz polynya, corroborating previous findings identifying polynya
340   bacterial communities to be mostly heterotrophs exploiting residues from eukaryotic
341   phytoplankton blooms[22], including taxa playing key roles in primary production limitation by
342   iron and other micronutrients like cobalamin[29,34]. The Mertz polynya was iron-limited at the

343 time of ACE sampling[15], and an investigation of metallo-proteins diversity in Mertz samples
344 through specific annotation tools[35] could help to better identify the roles of prokaryotes in trace
345 metal cycling in the context of diatom blooms. A previous study investigating a transect from
346 Tasmania to Mertz identified a significant difference in genomic composition between samples
347 taken at Mertz and samples taken above the Polar Front[36]. They highlighted the polar front as
348 the main biogeographical boundary, acknowledging that the continental shelf could also
349 explain the partitioning considering their lack of samples between the front and Mertz. Our
350 results suggest a greater difference between populations of the polynya *versus* other AASW
351 populations than between populations on both sides of the polar front, highlighting the
352 uniqueness of SO's coastal biomes. Mertz being the only sampling location above the
353 Antarctic shelf in our dataset, it is impossible to state if our observations could be considered
354 representative of shelf conditions at SO scale. We find them more likely to be polynya-specific,
355 as they seem to be driven by the high activity of a *Fragilariopsis cylindrus*-dominated bloom.

356 Phytoplankton dynamics in polynyas usually show dominance of either diatoms or
357 Prymnesiophyceae, mainly *Phaeocystis antarctica*[37,38]. Global warming could be causing a
358 shift from *P. antarctica* to diatom blooms in coastal polynyas[37,39], and the increased sinking
359 rate of diatoms compared to Phaeocystis could impact carbon export[28]. *P. antarctica* did not
360 appear as a significant contributor to any of the 6 CAGs identified as differentially abundant in
361 the Mertz polynya, while it was the main contributor to CAG 34, abundant in waters with low
362 silicic acid and moderately high temperature, far from the Antarctic coastline. Models predict
363 the diatom-*Phaeocystis* competition to mainly depend on iron availability and light sensitivity[40].
364 A eukaryote-focused re-analysis of the key samples identified through our approach, *i.e.* using
365 eukaryote specific gene-callers in combination with genomes from diatoms and *Phaeocystis*
366 isolates, could lead to the detection of functional markers helping to decipher the mechanisms
367 of the diatom-phaeocystis competition at genomic level. Interestingly, our Mertz-associated
368 CAGs were similar to genomic markers of a polynya in the Amundsen Sea dominated *by P.*
369 *antarctica*[28], suggesting bacterial functional redundancy in polynyas independently from the
370 dominant phytoplankton lineage. An analysis of the bacterial transcriptional activity in multiple
371 polynyas combined with measures of estimated carbon export[41] would lead to a better
372 understanding of the impact of planktonic compositional switches on remineralization and
373 sinking rates in polynyas, allowing for better predictions of their potential effect on the SO
374 biological carbon pump in a context of global change.

375 In addition to gene clusters highlighting the functional and taxonomic uniqueness of SO's
376 biomes, we identified gene clusters showing different latitudinal niches and functional profiles
377 despite all being associated with the ubiquitous SAR11. SAR11 subclades adapted to SO's
378 extreme conditions have been observed through amplicon sequencing off the Kerguelen
379 Islands[42], while SAR11 genomes assembled from Arctic metagenomes contained polar-
380 specific genes content, the vast majority of which coded for poorly characterized proteins[43].
381 Our results suggest a genomic adaptation of SAR11 across oceanographic fronts transitioning
382 from subtropical surface waters to Antarctic surface waters, even including the specific
383 conditions of a polynya bloom: strong competition for nutrients, organic matter and trace
384 metals. SAR11 could thus play a role in trace metal cycling in SO polar conditions. A strain-
385 resolved analysis of SAR11 genomes based on ACE metagenomes should provide
386 unprecedented insights into SAR11 Southern Ocean adapted ecotypes.

387 The ACE campaign ran from spring to late summer and some of the variability observed could
388 be temporal, as illustrated by the strong seasonal dynamics of viral communities of Marguerite
389 Bay[44]. The genomic content of microbial populations in the dark winter of the SO remains to

390 be described by future campaigns. Our results will soon be complemented by viral size fraction
391 metagenomics and >3μm size fraction metatranscriptomics samples from the same campaign,
392 which combined with our metagenomic assemblies should allow a better description of viral
393 and eukaryotic functional adaptations in the SO, offering a holistic view of its unique genomic
394 seascape.

395 By compiling catalogs of contigs, unigenes, AGC and CAGs from across the SO, we provide
396 a robust basis for any future polar and/or global-scale meta-omics investigation (Table S1).
397 Doing so, we address a critical gap in the metagenomes currently available in public
398 databases[12]. We notably provide the Southern Ocean Reference Gene Catalog (SO-RGC),
399 focused on the 0.2-3 μm size fraction and complementary to the OM-RGC[16,17], and a catalog
400 of polar-specific ACE ORFs, *i.e.* detected in at least one Arctic sample while being absent
401 from non-polar TOPC samples. Using these catalogs, we quantified the novelty of SO
402 microbial genes, demonstrating their high endemicity. By linking gene-level abundance and
403 environmental metadata, we were able to describe the biogeography of prokaryotes at SO-
404 scale, identifying distinct gene assemblages in different water-masses and defining genomic
405 markers of diverse biomes, from the blooming Mertz polynya to the Southernmost Sub-
406 Tropical waters. Overall, our results advocate for the development of regional-scale
407 descriptions and models of planktonic diversity in the Southern Ocean, distinguishing coastal
408 and offshore systems, and implementing the specific response of prokaryotes to localized
409 eukaryotic blooms. Our statistical results suggest that our gene catalog, combined with
410 extensive environmental and biogeochemical monitoring, could lead to correlative models of
411 gene abundance at SO scale, offering new tools to predict the future of this rapidly changing
412 ecosystem. Existing Antarctic time-series (*e.g.*, the Palmer LTER[45] or the Rothera time-
413 series[46]) should thus be complemented by genomic time series to provide valuable insights
414 into seasonal cycles and enhance our ability to monitor and predict the impact of climate
415 change on Southern Ocean microbial communities.

416

417 **MATERIAL AND METHODS**

418 A list of all publicly available resources is available in Table S1.

419 *Sampling and sequencing protocols*

420 218 samples for metagenomics analyses were collected at 34 stations during the ACE
421 campaign in the Austral summer 2016-2017. 197 of the 218 samples, thereafter called CTD
422 samples, were collected from Niskin bottles during rosette upcast and separated into three
423 size fractions (0.2-3 μm, 3-200 μm, and 0.2-40 μm). The remaining 21 samples, thereafter
424 called UDW samples, correspond to water pumped directly from the surface and separated
425 into the same three size fractions. Samples were sent for DNA extraction and shotgun
426 sequencing to Genoscope, the French National Platform for DNA Sequencing, following
427 protocols used by *Tara* expeditions[47]. Briefly, after filter cryogrinding, DNA was extracted using
428 total RNA/DNA Purification and Nucleospin RNA/DNA Buffer Set (MACHEREY-NAGEL).
429 Metagenomic libraries were prepared using the Illumina kit according to manufacturer
430 instructions. DNA libraries were sequenced on a Novaseq 4000 instrument, with a target of
431 100M paired-end reads per library (2x150bp; 500bp insert size).

432

433 *Environmental metadata compilation*

434 The ACE campaign hosted 22 scientific projects encompassing biology, oceanography,
435 climatology, glaciology, and biochemistry. For each CTD sample, all available metadata from

436   the corresponding cast and depth were retrieved from [SPI-ACE repository](). Similarly, metadata
437   from each pumping event were retrieved for UDW samples, but considering the limited number
438   of sequenced UDW samples and the lack of homogeneity in measured variables across CTD
439   and UDW samples, these metadata and their corresponding samples could not be used in
440   statistical investigations based on environmental variables (*i.e.*, random forest models, RDA).
441   For 10 of the 21 UDW samples, surface CTD metadata from the same sampling event were
442   available, enabling us to incorporate these samples in statistical investigations along with CTD
443   samples, while the remaining 11 UDW samples could not be considered.  Up to 56 variables
444   were retrieved per CTD samples, including basic physico-chemical variables (*e.g.*,
445   temperature, salinity, nutrients, depth), trace metals concentrations (*e.g.*, dissolved Fe, Cu,
446   Ni, Zn), isotopes (*e.g.* $^{13}$C, $^{15}$N) and pigment-based measures (*e.g.* concentrations of
447   cyanobacteria, diatoms or haptophytes). Variables measured in less than half of the samples
448   were dropped for further statistical explorations, leading to the selection of 33 variables. In
449   addition to these data retrieved *in situ*, physical variables were calculated at each sampling
450   using a Lagrangian approach and an integration time of 10 days. These included current
451   velocity, Okubo-Weiss (a proxy of eddy presence) and Lagrangian betweenness (a proxy of
452   bottleneck presence which has been related to biodiversity[48]. 14 latent variables computed
453   through a sparse PCA for each ACE station to summarize the global biogeochemical context[15]
454   were added to the metadata set. Please refer to the original Landwehr et al.[15] paper for a full
455   description of each latent variable. Finally, each sample was associated to a Longhurst
456   biogeographical province based on its coordinates, and to a water mass based on
457   temperature-salinity-oxygen diagrams computed from CTD downcast profiles (Figure 3A).
458   When needed, missing values in the CTD metadata set were imputed using the k-nearest-
459   neighbors approach encoded in the caret R package[49], with the default value of k=5. For a full
460   list of available metadata variables, a precise description of their compilation and of their pre-
461   processing, please refer to this GitHub repository: [ACE_gene_centric_scripts]().

462

463   *Assembly of metagenomic short reads and the profiling of resulting contigs*

464   Short-reads were quality-filtered using the Minoche[50] approach implemented in illumina-utils[51]
465   with default parameters, and sample-by-sample assemblies were obtained from MEGAHIT
466   v1.2.9[52]. The 68,074,004 contigs from the 218 single assemblies were concatenated into a
467   FASTA file from which a single anvi'o contigs database, hereafter called the ACE Contigs-DB,
468   was generated using the program anvi-gen-contigs-database as implemented in anvi'o v8[53].
469   During the generation of the ACE Contigs-DB open-reading frames were detected in all
470   contigs using Prodigal v2.6.3[54] which resulted in 175,336,776 non-dereplicated ORFs that
471   represented the raw ACE gene catalog for downstream analyses. To estimate the fraction of
472   eukaryotic organisms sampled, especially in the size fraction >3 µm, Phyloflash v3.4[55] was
473   used on quality-filtered reads. Considering that some samples were dominated by eukaryotes,
474   it is likely that some contigs in the ACE contigs database are from eukaryotic organisms. To
475   assess this likelihood, Eukrep v0.6.7[56] (West et al., 2018) and Whokaryote[57] were used to try
476   and detect eukaryotic contigs. However, only 2,343,800 contigs (3.4%) were classified as
477   eukaryotic by both tools, clearly underestimating the eukaryotic fraction of contigs. The
478   annotation of these contigs using the UniRef90 database and MMSeqs v14.7e284[58]
479   demonstrated the presence of 229,179 (9.8%) potential false positives annotated as bacteria.
480   We thus decided to keep all contigs in the database for the rest of the pipeline, while tagging
481   the ones identified as eukaryotic by EukRep as potentially eukaryotic.

482

483     *Generation and annotation of Southern Ocean's microbial reference gene catalog*

484     Open-reading frames were detected in all contigs using Prodigal v2.6.3[54]. The 175,336,776
485     non-dereplicated ORFs constitute the raw ACE gene catalog. Nucleotide sequences were
486     then clustered at 95% similarity and 90% coverage using CD-Hit V4.8.1[59], to produce unigenes
487     comparable to those of the OM-RGC computed from Tara Oceans and Tara Polar Circle
488     expeditions. The 89,739,060 unigenes produced constitute the full ACE reference gene
489     catalog (ACE-RGC). The ACE-RGC was annotated with EggNOG-mapper v2.1.8[60] and
490     KOFamSCAN v1.3.0[61]. To allow easier usage in conjunction with the OM-RGC, in which only
491     the 0.2-3 μm size fraction is included, the SO-RGC was defined as the unigenes from the
492     ACE-RGC that contained at least one ORF detected in a contig assembled in the 0.2-3 μm
493     size fraction. Finally, to produce coarser yet functionally homogeneous clusters, the
494     AGNOSTOS clustering pipeline[19] was used on the raw ACE gene catalog to produce
495     30,123,228 AGNOSTOS gene clusters (AGC), of which 765,003 were discarded as low
496     quality. The 29,358,225 good quality AGC were classified in 4 categories based on their PFAM
497     annotation and their similarity with the members of the AGNOSTOS-DB: Known (K), Known
498     without PFAM (KWP), Genomic unknown (GU; genes of unknown function yet found in a
499     genomic context - MAG, SAG, isolate genome...) or Environmental unknown (EU; genes of
500     unknown function never integrated in a genomic context). For a detailed description of these
501     categories and of the methodology for clustering and annotating within the AGNOSTOS
502     pipeline, please refer to Vanni et al.[19]. Please note that the AGC we use in this study are
503     issued from an AGNOSTOS-based clustering and annotation of ACE ORFs, and not to an
504     integration of ACE ORFs within the public AGNOSTOS gene database. AGC-level EggNOG
505     and KEGG annotations were defined as the modal value from the annotations of all cluster's
506     members.

507

508     *Computation of gene- and cluster-level coverage and detection*

509     Quality-filtered short reads were mapped on the ACE contigs DB to produce contigs-level
510     coverage and detection (% of the contigs covered at least at 1X) profiles, using Bowtie2
511     v2.4.5[62], in competitive mode with equivalent mapping scores across different references
512     distributed at random. Gene-level metrics were obtained for all the raw ACE gene catalog
513     through the program anvi-profile-blitz ([https://anvio.org/m/anvi-profile-blitz](https://anvio.org/m/anvi-profile-blitz)) implemented in
514     anvi'o[53,63] for this project. By deriving gene-level metrics from the larger genomic context
515     afforded by contigs, rather than using read recruitment to individual gene sequences, we were
516     able to (1) avoid bell-shaped coverage signal that would dwindle around ORF extremities, (2)
517     avoid mapping errors due to assembled sequences being removed from the reference during
518     pre-mapping de-replication, and (3) use exhaustive contigs-level metrics to build direct links
519     between gene-level results obtained in this study and MAGs-level results obtained in parallel
520     work (Pommellec et al, in preparation). The coverage values reported from anvi-profile-blitz
521     were expressed per base-pair, *i.e.* normalized by gene length. Outputs from all samples were
522     then concatenated into a coverage matrix and a detection matrix of each 218 columns and
523     175,336,776 lines where each line represented an individual ORF.

524     The coverage of each unigene was defined as the sum of the per-base pair coverages from
525     all members of its dereplication cluster. Similarly, per-base pair coverages of all members of
526     each AGNOSTOS cluster were summed to obtain AGC-level coverages. To avoid false-
527     positive coverage values due to mapping mistakes and read dilution across conserved
528     domains, a threshold of detection was applied at cluster-level. Detection at cluster level was
529     defined as $Detection_{Cluster} = \max(Detection_{Cluster\ members})$.

530 Increasing the threshold of detection at cluster level caused both the mean slope of the
531 collector curve and the amount of undetected AGNOSTOS clusters to increase (Figure S2). A
532 flat collector curve is likely to be the result of false positives considering the many singletons
533 that are likely to be rare, but a high number of undetected clusters is likely to be due to false
534 negatives since their sequences should be present at least in the samples in which they were
535 assembled. We then decided to use a threshold of 60% detection, as it was the highest
536 threshold allowing to detect more than 95% of AGNOSTOS clusters in at least one sample.
537 To apply this threshold, all AGC-level coverage values corresponding to AGC-level detection
538 scores below 60% were turned to 0.

539

540 *Pre-processing and normalization of cluster-level abundances*

541 After applying the 60% detection threshold, all remaining coverage values were rounded to
542 the nearest integer. The whole AGNOSTOS cluster-level coverage matrix was then
543 normalized using the *rlog* method from the DESeq2 R package[64]. Relative log expression
544 normalization method was identified as one of the most adapted to metagenomics-based
545 microbiome studies[65].

546

547 *Highlighting novelty in Southern Ocean's microbial genes*

548 Protein sequences from the ACE-RGC were clustered with those of the OM-RGC at 30%,
549 50% and 80% similarity thresholds, with a fixed coverage threshold of 80%. Clusters were
550 separated in three categories: pure ACE when only composed of sequences assembled in
551 ACE samples, pure TARA when only composed of sequences from the OM-RGC, and mixed
552 for the rest. Clusters were further characterized based on their members' origin of assembly,
553 mainly distinguishing sunlit (<150m) and dark (>150m) samples as well as the different size
554 fractions.

555 To better estimate the global presence of genes from the ACE-RGC, short reads from 134
556 samples from Tara Oceans and Tara Polar Circle expeditions corresponding to the 0.2-3 µm
557 size fraction were quality-filtered and mapped on the ACE contigs DB using the protocol
558 described in *Computation of gene- and cluster-level coverage and detection.* Description of
559 the Tara samples is available in Salazar et al.[16].

560

561 *Identifying environmental drivers of gene-clusters distribution at Southern Ocean's scale*

562 For further biogeographical explorations, the global matrix was split into two parts, the free-
563 living part corresponding to 0.2-3 and 0.2-40 µm size fractions, and the >3µm part
564 corresponding to the 3-200 µm size fraction. As stated in the *Environmental metadata*
565 *compilation* section, UDW samples were removed from the biogeographical investigations due
566 to differences in environmental metadata availability. Finally, clusters showing near-zero
567 variance abundance profiles were removed from each matrix using the preProcess function
568 from the Caret R Package[49]. The near zero variance definition was set at a minimal threshold
569 of 20% unique abundance values and a maximal ratio of 95 to 5 between the most abundant
570 and second most abundant values.

571

572 *Identification of AGNOSTOS clusters highly linked with the environment*

573 A random forest regression model was fitted for each cluster of the free-living and >3µm
574 matrices that passed the near zero variance threshold. Normalized coverage values were

575  used as interest variables, and the 50 environmental variables from the CTD metadata as
576  predictors. For each model, the number of predictors tried at each split was optimized between
577  5 and 8 (default being the rounded down square root of the total number of predictors), the
578  number of trees was set at 501 and other parameters were left at default in the ranger function
579  from the ranger R package[66]. Each model went through 3 repetitions of 4-fold cross-validation
580  using the train function from the Caret package. Variable importance, based on permutations,
581  and adjusted cross-validation R-squared values from each selected model were retrieved.
582  Density of R-squared values were drawn for free-living and >3µm results, separately. To
583  estimate a threshold of R-squared at which it is unlikely that the link between coverage and
584  metadata could be observed by chance, the same runs of random-forest models were
585  computed on four matrices with randomly permuted rows, two of the free-living matrices and
586  two of the >3µm ones. Based on the 95th centile value for each set of permutations, R-squared
587  thresholds were set at 10% for the free-living matrix and 15% for the >3µm one. AGNOSTOS
588  clusters meeting these thresholds were defined as highly linked with the environment (env-
589  AGC).
590
591  *Grouping of co-abundant AGNOSTOS clusters*
592  To further reduce the dimensionality of the two matrices of interest without removing any env-
593  AGC, the approach described by Minot and Willis[24] was used to group them into groups of Co-
594  Abundant env-AGC (CAG). This approach, based around the Approximate Nearest Neighbor
595  heuristic, allows to cluster millions of genes/gene clusters into co-abundant groups with limited
596  computer power and in reasonable time. The clustering python scripts available at
597  https://github.com/FredHutch/find-cags were used with default parameters independently on
598  the free-living and >3µm env-AGC matrices. CAG-level coverage matrices were created by
599  summing coverage across all members of a CAG.
600
601  *Constrained ordination and further investigation of CAGs*
602  Redundancy analysis was fitted on the Hellinger-transformed free-living and >3µm CAG-level
603  matrices. Again, coverage values were used as interest variables, and environmental
604  variables as explanatory variables. Both analyses were significant (ANOVA, p-value=0.001 for
605  both free-living and >3µm), allowing us to go further by selecting environmental variables
606  through a two-directional stepwise selection based on the Akaike Information Criterion (AIC).
607  The selected models were again significant (ANOVA, p-value=0.001 for both free-living and
608  >3µm). For each model, CAGs appearing on the extremities of axis 1 to 5 were individually
609  selected to be analyzed in depth. Taxonomic annotations of genes within each CAG were
610  retrieved through the annotation of their contigs of origin through MMSeqs2 taxonomic
611  annotation tool with the UniRef90 database as reference[58]. Since it took 15 to 20 hours to
612  annotate splits of 25,000 contigs using 24 CPUs and 80 Gb of memory, only a selection of
613  CAGs of interest were annotated this way. In addition, genes in CAGs were annotated based
614  on the presence of their contig of origin in MAGs from the ACE MAGs database (Pommellec
615  et al., in prep.), allowing precise taxonomic annotations for all genes found in MAGs. To
616  estimate a potential functional enrichment in a CAG, AGNOSTOS cluster-level EggNOG and
617  KEGG annotations were retrieved for all members of the CAG, and compared to the rest of
618  AGNOSTOS cluster-level annotations through a one-sided Fisher test. Obtained p-values
619  were corrected for multiple testing using the Benjamini-Hochberg method, and p-value
620  threshold was set to 0.01 for enrichment.

621

*Code and data availability*

The 218 metagenomic samples were deposited on the European Nucleotide Archive (ENA), under the submission code ERA30995399. A correspondence table of all samples linking their ACE unique ID, BioSample code, ENA run and experiment codes is available on Zenodo. Links to all raw and processed data are available in Table S1. The scripts used to produce the results of this study are available at https://github.com/EmileFaure/ACE_gene_centric_scripts. The different steps to go from quality-filtered reads to gene-level per-base coverage and detection matrices were integrated into a Nextflow workflow available at https://gitlab.ifremer.fr/bioinfo/workflows/noemie.

*References*

1. Gray, A. R. The Four-Dimensional Carbon Cycle of the Southern Ocean. Annu. Rev. Mar. Sci. **16**, 163–190 (2024).

2. Hassler, C. S., Sinoir, M., Clementson, L. A. & Butler, E. C. V. Exploring the Link between Micronutrients and Phytoplankton in the Southern Ocean during the 2007 Austral Summer. Front. Microbiol. **3**, 202 (2012).

3. Tagliabue, A. et al. The integral role of iron in ocean biogeochemistry. Nature **543**, 51–59 (2017).

4. Deppeler, S. L. & Davidson, A. T. Southern Ocean Phytoplankton in a Changing Climate. Front. Mar. Sci. **4**, (2017).

5. Hauck, J. et al. On the Southern Ocean CO2 uptake and the role of the biological carbon pump in the 21st century. Glob. Biogeochem. Cycles **29**, 1451–1470 (2015).

6. Christaki, U. et al. Seasonal microbial food web dynamics in contrasting Southern Ocean productivity regimes. Limnol. Oceanogr. **66**, 108–122 (2021).

7. Landa, M., Blain, S., Christaki, U., Monchy, S. & Obernosterer, I. Shifts in bacterial community composition associated with increased carbon cycling in a mosaic of phytoplankton blooms. ISME J. **10**, 39–50 (2016).

662    8.  Doré, H. et al. Differential global distribution of marine picocyanobacteria gene clusters reveals
663        distinct niche-related adaptive strategies. ISME J. **17**, 720–732 (2023).

664    9.  Faure, E., Ayata, S.-D. & Bittner, L. Towards omics-based predictions of planktonic functional
665        composition from environmental data. Nat. Commun. **12**, 4361 (2021).

666   10.  Frémont, P. et al. Restructuring of plankton genomic biogeography in the surface ocean under
667        climate change. Nat. Clim. Change **12**, 393–401 (2022).

668   11.  Sunagawa, S. et al. Tara Oceans: towards global ocean ecosystems biology. Nat. Rev. Microbiol.
669        **18**, 428–445 (2020).

670   12.  Laiolo, E. et al. Metagenomic probing toward an atlas of the taxonomic and metabolic foundations
671        of the global ocean genome. Front. Sci. **1**, (2024).

672   13.  Paoli, L. et al. Biosynthetic potential of the global ocean microbiome. Nature **607**, 111–118 (2022).

673   14.  Cao, S. et al. Structure and function of the Arctic and Antarctic marine microbiota as revealed by
674        metagenomics. Microbiome **8**, 47 (2020).

675   15.  Landwehr, S. et al. Exploring the coupled ocean and atmosphere system with a data science
676        approach applied to observations from the Antarctic Circumnavigation Expedition. Earth Syst.
677        Dyn. **12**, 1295–1369 (2021).

678   16.  Salazar, G. et al. Gene Expression Changes and Community Turnover Differentially Shape the
679        Global Ocean Metatranscriptome. Cell **179**, 1068-1083.e21 (2019).

680   17.  Sunagawa, S. et al. Structure and function of the global ocean microbiome. Science **348**, 1261359–
681        1261359 (2015).

682   18.  Acinas, S. G. et al. Deep ocean metagenomes provide insight into the metabolic architecture of
683        bathypelagic microbial communities. Commun. Biol. **4**, 1–15 (2021).

684   19.  Vanni, C. et al. Unifying the known and unknown microbial coding sequence space. eLife **11**,
685        e67667 (2022).

686   20.  Henry, T. et al. Physical and biogeochemical oceanography data from Conductivity, Temperature,
687        Depth (CTD) rosette deployments during the Antarctic Circumnavigation Expedition (ACE). Zenodo
688        https://doi.org/10.5281/zenodo.3813646 (2020).

689   21.  Abell, G. C. J. & Bowman, J. P. Ecological and biogeographic relationships of class Flavobacteria in
690        the Southern Ocean. FEMS Microbiol. Ecol. **51**, 265–277 (2005).

691   22.  Wilkins, D., van Sebille, E., Rintoul, S. R., Lauro, F. M. & Cavicchioli, R. Advection shapes Southern
692        Ocean microbial assemblages independent of distance and environment effects. Nat. Commun.
693        **4**, (2013).

694   23.  Li, Q., England, M. H., Hogg, A. M., Rintoul, S. R. & Morrison, A. K. Abyssal ocean overturning
695        slowdown and warming driven by Antarctic meltwater. Nature **615**, 841–847 (2023).

696   24.  Minot, S. S. & Willis, A. D. Clustering co-abundant genes identifies components of the gut
697        microbiome that are reproducibly associated with colorectal cancer and inflammatory bowel
698        disease. Microbiome **7**, 110 (2019).

699  25. Suzek, B. E., Huang, H., McGarvey, P., Mazumder, R. & Wu, C. H. UniRef: comprehensive and non-
700      redundant UniProt reference clusters. Bioinformatics **23**, 1282–1288 (2007).

701  26. Xue, C. et al. Polysaccharide utilization by a marine heterotrophic bacterium from the SAR92
702      clade. FEMS Microbiol. Ecol. **97**, fiab120 (2021).

703  27. Kappelmann, L. (Meta-)genomic Analysis of the Diversity and the Carbohydrate Degradation
704      Potential of the SAR92 Clade during a Diatom-induced Bacterioplankton Bloom. (University of
705      Bremen Bremen / Germany, 2013).

706  28. Kim, S.-J. et al. Genomic and metatranscriptomic analyses of carbon remineralization in an
707      Antarctic polynya. Microbiome **7**, 29 (2019).

708  29. Bertrand, E. M. et al. Phytoplankton–bacterial interactions mediate micronutrient colimitation at
709      the coastal Antarctic sea ice edge. Proc. Natl. Acad. Sci. **112**, 9938–9943 (2015).

710  30. Mock, T. et al. Evolutionary genomics of the cold-adapted diatom Fragilariopsis cylindrus. Nature
711      **541**, 536–540 (2017).

712  31. Holst, F. et al. Helixer–de novo Prediction of Primary Eukaryotic Gene Models Combining Deep
713      Learning and a Hidden Markov Model. 2023.02.06.527280 Preprint at
714      https://doi.org/10.1101/2023.02.06.527280 (2023).

715  32. Ibarbalz, F. M. et al. Global Trends in Marine Plankton Diversity across Kingdoms of Life. Cell **179**,
716      1084-1097.e21 (2019).

717  33. Sussfeld, D. et al. Network studies unveil new groups of highly divergent proteins in families as
718      old as cellular life with important biological functions in the ocean. 2024.01.08.574615 Preprint
719      at https://doi.org/10.1101/2024.01.08.574615 (2024).

720  34. Debeljak, P., Toulza, E., Beier, S., Blain, S. & Obernosterer, I. Microbial iron metabolism as revealed
721      by gene expression profiles in contrasted Southern Ocean regimes. Environ. Microbiol. **21**, 2360–
722      2374 (2019).

723  35. Garber, A. I. et al. FeGenie: A Comprehensive Tool for the Identification of Iron Genes and Iron
724      Gene Neighborhoods in Genome and Metagenome Assemblies. Front. Microbiol. **11**, 37 (2020).

725  36. Wilkins, D. et al. Biogeographic partitioning of Southern Ocean microorganisms revealed by
726      metagenomics. Environ. Microbiol. **15**, 1318–1333 (2013).

727  37. Arrigo, K. R. & van Dijken, G. L. Phytoplankton dynamics within 37 Antarctic coastal polynya
728      systems. J. Geophys. Res. Oceans **108**, (2003).

729  38. Schofield, O. et al. In situ phytoplankton distributions in the Amundsen Sea Polynya measured by
730      autonomous gliders. Elem. Sci. Anthr. **3**, 000073 (2015).

731  39. Arrigo, K. R., Lowry, K. E. & van Dijken, G. L. Annual changes in sea ice and phytoplankton in
732      polynyas of the Amundsen Sea, Antarctica. Deep Sea Res. Part II Top. Stud. Oceanogr. **71–76**, 5–
733      15 (2012).

734  40. Nissen, C. & Vogt, M. Factors controlling the competition between &lt;i&gt;Phaeocystis&lt;/i&gt;
735      and diatoms in the Southern Ocean and implications for carbon export fluxes. Biogeosciences **18**,
736      251–283 (2021).

737  41. Guidi, L. et al. Plankton networks driving carbon export in the oligotrophic ocean. Nature (2016).

738  42. Dinasquet, J., Landa, M. & Obernosterer, I. SAR11 clade microdiversity and activity during the early
739      spring blooms off Kerguelen Island, Southern Ocean. Environ. Microbiol. Rep. **14**, 907–916 (2022).

740  43. Kraemer, S., Ramachandran, A., Colatriano, D., Lovejoy, C. & Walsh, D. A. Diversity and
741      biogeography of SAR11 bacteria from the Arctic Ocean. ISME J. **14**, 79–90 (2020).

742  44. Piedade, G. J. et al. Seasonal dynamics and diversity of Antarctic marine viruses reveal a novel viral
743      seascape. Nat. Commun. **15**, 9192 (2024).

744  45. Smith, R. C. et al. The Palmer LTER: A Long-Term Ecological Research Program at Palmer Station,
745      Antarctica. Oceanography **8**, 77–86 (1995).

746  46. Venables, H. et al. Sustained year-round oceanographic measurements from Rothera Research
747      Station, Antarctica, 1997–2017. Sci. Data **10**, 265 (2023).

748  47. Alberti, A. et al. Viral to metazoan marine plankton nucleotide sequences from the Tara Oceans
749      expedition. Sci. Data **4**, 170093 (2017).

750  48. Ser-Giacomi, E. et al. Lagrangian betweenness as a measure of bottlenecks in dynamical systems
751      with oceanographic examples. Nat. Commun. **12**, 4935 (2021).

752  49. Kuhn, M. Building Predictive Models in R Using the caret Package. J. Stat. Softw. **28**, 1–26 (2008).

753  50. Minoche, A. E., Dohm, J. C., Himmelbauer, H., & others. Evaluation of genomic high-throughput
754      sequencing data generated on Illumina HiSeq and genome analyzer systems. Genome Biol **12**,
755      R112 (2011).

756  51. Eren, A. M., Vineis, J. H., Morrison, H. G. & Sogin, M. L. A Filtering Method to Generate High Quality
757      Short Reads Using Illumina Paired-End Technology. PLoS ONE **8**, e66643 (2013).

758  52. Li, D., Liu, C.-M., Luo, R., Sadakane, K. & Lam, T.-W. MEGAHIT: An ultra-fast single-node solution
759      for large and complex metagenomics assembly via succinct de Bruijn graph. Bioinformatics btv033
760      (2015).

761  53. Eren, A. M. et al. Anvi'o: an advanced analysis and visualization platform for 'omics data. PeerJ **3**,
762      e1319 (2015).

763  54. Hyatt, D. et al. Prodigal: prokaryotic gene recognition and translation initiation site identification.
764      BMC Bioinformatics **11**, 1 (2010).

765  55. Gruber-Vodicka, H. R., Seah, B. K. B. & Pruesse, E. phyloFlash: Rapid Small-Subunit rRNA Profiling
766      and Targeted Assembly from Metagenomes. mSystems **5**, (2020).

767  56. West, P. T., Probst, A. J., Grigoriev, I. V., Thomas, B. C. & Banfield, J. F. Genome-reconstruction for
768      eukaryotes from complex natural microbial communities. Genome Res. **28**, 569–580 (2018).

769  57. Pronk, L. J. U. & Medema, M. H. Whokaryote: distinguishing eukaryotic and prokaryotic contigs in
770      metagenomes based on gene structure. Microb. Genomics **8**, 000823 (2022).

771  58. Steinegger, M. & Söding, J. MMseqs2 enables sensitive protein sequence searching for the
772      analysis of massive data sets. Nat. Biotechnol. **35**, 1026–1028 (2017).

773    59.   Li, W. & Godzik, A. Cd-hit: a fast program for clustering and comparing large sets of protein or
774         nucleotide sequences. Bioinforma. Oxf. Engl. **22**, 1658–1659 (2006).

775    60.   Huerta-Cepas, J. et al. eggNOG 4.5: a hierarchical orthology framework with improved functional
776         annotations for eukaryotic, prokaryotic and viral sequences. Nucleic Acids Res. **44**, D286–D293
777         (2016).

778    61.   Aramaki, T. et al. KofamKOALA: KEGG Ortholog assignment based on profile HMM and adaptive
779         score threshold. Bioinformatics **36**, 2251–2252 (2020).

780    62.   Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. Nat. Methods **9**, 357–
781         359 (2012).

782    63.   Eren, A. M. et al. Community-led, integrated, reproducible multi-omics with anvi'o. Nat. Microbiol.
783         **6**, 3–6 (2021).

784    64.   Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for RNA-
785         seq data with DESeq2. Genome Biol. **15**, 550 (2014).

786    65.   Lloréns-Rico, V., Vieira-Silva, S., Gonçalves, P. J., Falony, G. & Raes, J. Benchmarking microbiome
787         transformations favors experimental quantitative approaches to address compositionality and
788         sampling depth biases. Nat. Commun. **12**, 3562 (2021).

789    66.   Wright, M. N. & Ziegler, A. ranger: A Fast Implementation of Random Forests for High Dimensional
790         Data in C++ and R. J. Stat. Softw. **77**, (2017).
791

# Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- FaureetalWaterMassSpecificPolarGenesDominateTheSOMicrobiomeSupplementaryMaterials.docx