

Published in final edited form as:

*Nature*. 2016 April 28; 532(7600): 465–470. doi:10.1038/nature16942.

## Plankton networks driving carbon export in the oligotrophic ocean

Lionel Guidi<sup>#1,2</sup>, Samuel Chaffron<sup>#3,4,5</sup>, Lucie Bittner<sup>#6,7,8</sup>, Damien Eveillard<sup>#9</sup>, Abdelhalim Larhlimi<sup>9</sup>, Simon Roux<sup>10,11</sup>, Youssef Darzi<sup>3,4</sup>, Stephane Audic<sup>8</sup>, Léo Berline<sup>1,12</sup>, Jennifer Brum<sup>10,11</sup>, Luis Pedro Coelho<sup>13</sup>, Julio Cesar Ignacio Espinoza<sup>10</sup>, Shruti Malviya<sup>7</sup>, Shinichi Sunagawa<sup>13</sup>, Céline Dimier<sup>8</sup>, Stefanie Kandels-Lewis<sup>13,14</sup>, Marc Picheral<sup>1</sup>, Julie Poulain<sup>15</sup>, Sarah Searson<sup>1,2</sup>, Tara Oceans coordinators, Lars Stemmann<sup>1</sup>, Fabrice Not<sup>8</sup>, Pascal Hingamp<sup>16</sup>, Sabrina Speich<sup>17</sup>, Mick Follows<sup>18</sup>, Lee Karp-Boss<sup>19</sup>, Emmanuel Boss<sup>19</sup>, Hiroyuki Ogata<sup>20</sup>, Stephane Pesant<sup>21,22</sup>, Jean Weissenbach<sup>15,23,24</sup>, Patrick Wincker<sup>15,23,24</sup>, Silvia G. Acinas<sup>25</sup>, Peer Bork<sup>13,26</sup>, Colomban de Vargas<sup>8</sup>, Daniele Iudicone<sup>27</sup>, Matthew B. Sullivan<sup>10,11</sup>, Jeroen Raes<sup>3,4,5</sup>, Eric Karsenti<sup>7,14</sup>, Chris Bowler<sup>7</sup>, and Gabriel Gorsky<sup>1</sup>

<sup>1</sup>Sorbonne Universités, UPMC Université Paris 06, CNRS, Laboratoire d'océanographie de Villefranche (LOV), Observatoire Océanologique, Villefranche-sur-Mer, France

<sup>2</sup>Department of Oceanography, University of Hawaii, Honolulu, Hawaii, USA

<sup>3</sup>Department of Microbiology and Immunology, Rega Institute, KU Leuven, Herestraat 49, 3000 Leuven, Belgium.

<sup>4</sup>Center for the Biology of Disease, VIB, Herestraat 49, 3000 Leuven, Belgium.

<sup>5</sup>Department of Applied Biological Sciences, Vrije Universiteit Brussel, Pleinlaan 2, 1050 Brussels, Belgium.

<sup>6</sup>Sorbonne Universités, UPMC Univ Paris 06, CNRS, Institut de Biologie Paris-Seine (IBPS), Evolution Paris Seine, F-75005, Paris, France.

<sup>7</sup>Ecole Normale Supérieure, PSL Research University, Institut de Biologie de l'Ecole Normale Supérieure (IBENS), CNRS UMR 8197, INSERM U1024, 46 rue d'Ulm, F-75005 Paris, France.

<sup>8</sup>Sorbonne Universités, UPMC Université Paris 06, CNRS, Laboratoire Adaptation et Diversité en Milieu Marin, Station Biologique de Roscoff, Roscoff, France

<sup>9</sup>LINA UMR 6241, Université de Nantes, EMN, CNRS, 44322 Nantes, France.

Reprints and permissions information is available at [www.nature.com/reprints](http://www.nature.com/reprints). Users may view, print, copy, and download text and data-mine the content in such documents, for the purposes of academic research, subject always to the full Conditions of use: [http://www.nature.com/authors/editorial\\_policies/license.html#terms](http://www.nature.com/authors/editorial_policies/license.html#terms)

Correspondence and requests for materials should be addressed to [lguidi@obs-vlfr.fr](mailto:lguidi@obs-vlfr.fr), [samuel.chaffron@vib-kuleuven.be](mailto:samuel.chaffron@vib-kuleuven.be), [lucie.bittner@upmc.fr](mailto:lucie.bittner@upmc.fr), [damien.eveillard@univ-nantes.fr](mailto:damien.eveillard@univ-nantes.fr), [Jeroen.Raes@vib-kuleuven.be](mailto:Jeroen.Raes@vib-kuleuven.be), [karsenti@embl.de](mailto:karsenti@embl.de), [cbowler@biology.ens.fr](mailto:cbowler@biology.ens.fr), [gorsky@obs-vlfr.fr](mailto:gorsky@obs-vlfr.fr).

<sup>11</sup>Current affiliation: Department of Microbiology, The Ohio State University, Columbus OH 43210, USA

<sup>12</sup>Current affiliation: Aix-Marseille Univ., Mediterranean Institute of Oceanography (MIO), 13288, Marseille, Cedex 09, France ; Université du Sud Toulon-Var, MIO, 83957, La Garde cedex, France ; CNRS/INSU, MIO UMR 7294; IRD, MIO UMR235.

Data described herein is available at EBI under the project identifiers PRJEB402, PRJEB6610 and PRJEB7988, PANGAEA<sup>50,51,54</sup>, and a companion website (<http://www.raeslab.org/companion/ocean-carbon-export.html>). The data release policy regarding future public release of Tara Oceans data is described in *Pesant et al.*, [2015]<sup>49</sup>. All authors approved the final manuscript.

The authors declare no competing financial interests.

- <sup>10</sup>Department of Ecology and Evolutionary Biology, University of Arizona, Tucson, AZ, 85721, USA.
- <sup>13</sup>Structural and Computational Biology, European Molecular Biology Laboratory, Meyerhofstr. 1, 69117 Heidelberg, Germany
- <sup>14</sup>Directors' Research European Molecular Biology Laboratory Meyerhofstr. 1 69117 Heidelberg Germany
- <sup>15</sup>CEA - Institut de Génomique, GENOSCOPE, 2 rue Gaston Crémieux, 91057 Evry France.
- <sup>16</sup>Aix Marseille Université CNRS IGS UMR 7256 13288 Marseille France
- <sup>17</sup>Department of Geosciences, Laboratoire de Météorologie Dynamique (LMD), Ecole Normale Supérieure, 24 rue Lhomond 75231 Paris Cedex 05 France.
- <sup>18</sup>Dept of Earth, Atmospheric and Planetary Sciences, Massachusetts Institute of Technology, Cambridge, USA.
- <sup>19</sup>School of Marine Sciences, University of Maine, Orono, USA.
- <sup>20</sup>Institute for Chemical Research, Kyoto University, Gokasho, Uji, Kyoto, 611-0011, Japan.
- <sup>21</sup>PANGAEA, Data Publisher for Earth and Environmental Science, University of Bremen, Bremen, Germany.
- <sup>22</sup>MARUM, Center for Marine Environmental Sciences, University of Bremen, Bremen, Germany.
- <sup>23</sup>CNRS, UMR 8030, CP5706, Evry France.
- <sup>24</sup>Université d'Evry, UMR 8030, CP5706, Evry France.
- <sup>25</sup>Department of Marine Biology and Oceanography, Institute of Marine Sciences (ICM)-CSIC Pg. Marítim de la Barceloneta 37-49 Barcelona E08003 Spain.
- <sup>26</sup>Max-Delbrück-Centre for Molecular Medicine, 13092 Berlin, Germany
- <sup>27</sup>Stazione Zoologica Anton Dohrn, Villa Comunale, 80121, Naples, Italy.
- # These authors contributed equally to this work.

## Abstract

The biological carbon pump is the process by which CO<sub>2</sub> is transformed to organic carbon *via* photosynthesis, exported through sinking particles, and finally sequestered in the deep ocean. While the intensity of the pump correlates with plankton community composition, the underlying ecosystem structure driving the process remains largely uncharacterised. Here we use environmental and metagenomic data gathered during the *Tara* Oceans expedition to improve our understanding of carbon export in the oligotrophic ocean. We show that specific plankton communities, from the surface and deep chlorophyll maximum, correlate with carbon export at 150 m and highlight unexpected taxa such as Radiolaria, alveolate parasites, as well as *Synechococcus* and their phages, as lineages most strongly associated with carbon export in the subtropical, nutrient-depleted, oligotrophic ocean. Additionally, we show that the relative abundance of just a few bacterial and viral genes can predict most of the variability in carbon export in these regions.

Marine planktonic photosynthetic organisms are responsible for approximately fifty percent of Earth's primary production and fuel the global ocean biological carbon pump<sup>1</sup>. The intensity of the pump is correlated to plankton community composition<sup>2,3</sup>, and controlled by the relative rates of primary production and carbon remineralisation<sup>4</sup>. About 10% of this newly produced organic carbon in the surface ocean is exported through gravitational sinking of particles. Finally, after multiple transformations, a fraction of the exported material reaches the deep ocean where it is sequestered over thousand-year timescales<sup>5</sup>.

Like most biological systems, marine ecosystems in the sunlit upper layer of the ocean (denoted the euphotic zone) are complex<sup>6,7</sup>, characterised by a wide range of biotic and abiotic interactions<sup>8-10</sup> and in constant balance between carbon production, transfer to higher trophic levels, remineralisation, and export to the deep layers<sup>11</sup>. The marine ecosystem structure and its taxonomic and functional composition likely evolved to comply with this loss of energy by modifying organism turnover times and by the establishment of complex feedbacks between them<sup>6</sup> and the substrates they can exploit for metabolism<sup>12</sup>. Decades of groundbreaking research have focused on identifying independently the key players involved in the biological carbon pump. Among autotrophs, diatoms are commonly attributed to being important in carbon flux because of their large size and fast sinking rates<sup>13-15</sup> while small autotrophic picoplankton may contribute directly through subduction of surface water<sup>16</sup> or indirectly by aggregating with larger settling particles or consumption by organisms at higher trophic levels<sup>17</sup>. Among heterotrophs, zooplankton such as crustaceans impact carbon flux *via* production of fast-sinking fecal pellets while migrating hundreds of meters in the water-column<sup>18,19</sup>. These observations, focusing on just a few components of the marine ecosystem, highlight that carbon export results from multiple biotic interactions and that a better understanding of the mechanisms involved in its regulation will require an analysis of the entire planktonic ecosystem.

Advanced sequencing technologies offer the opportunity to simultaneously survey whole planktonic communities and associated molecular functions in unprecedented detail. Such a holistic approach may allow the identification of community- or gene-based biomarkers that could be used to monitor and predict ecosystem functions, *e.g.*, related to the biogeochemistry of the ocean<sup>20-22</sup>. Here, we leverage global-scale ocean genomics datasets from the euphotic zone<sup>10,23-25</sup> and associated environmental data to assess the coupling between ecosystem structure, functional repertoire, and carbon export at 150 m.

## Carbon export and plankton community composition

The *Tara* Oceans global circumnavigation crossed diverse ocean ecosystems and sampled plankton at an unprecedented scale<sup>20,26</sup> (see Methods). Hydrographic data were measured *in situ* or in seawater samples at all stations, as well as nutrients, oxygen and photosynthetic pigments (see Methods). Net Primary Production (NPP) was derived from satellite measurements (see Methods). In addition, particle size distributions (100  $\mu\text{m}$  to a few mm) and concentrations were measured using an Underwater Vision Profiler (UVP) from which carbon export, corresponding to the carbon flux (Fig. 1a) at 150 m, was calculated to range from 0.014 to 18.3  $\text{mg}\cdot\text{m}^{-2}\cdot\text{d}^{-1}$  using methods previously described (see Methods). One

should keep in mind that fluxes are calculated from images of particles. These estimates are derived from an approximation of Stokes' law relating the equivalent spherical diameter of particles to carbon flux (see Methods). This exponential approximation is reasonable assuming similar particle composition across all sizes, as highlighted by the standard deviations of parameters in Eq. 5 (see Methods). Furthermore, because of instrument and method limitations, particles  $<250 \mu\text{m}$  were not used, which may underestimate total carbon fluxes. Finally, these fluxes are instantaneous because they do not integrate space and time as sediment traps would. However, the approach allowed us to assemble the largest homogeneous carbon export dataset during a single expedition, corresponding to more than 600 profiles over 150 stations. This dataset is of similar magnitude to the body of historical data available in the literature that includes the 134 deep sediment trap-based carbon flux time-series<sup>27</sup> from the JGOFS program and the 419 thorium-derived particulate organic carbon (POC) export measurements<sup>28</sup>.

From 68 globally distributed sites, a total of 7.2 Tb of metagenomics data, representing ~40 million non-redundant genes, around 35,000 Operational Taxonomic Units (OTUs) of prokaryotes (Bacteria and Archaea) and numerous mainly uncharacterized viruses and picoeukaryotes, have been described recently<sup>23,25</sup>. In addition, a set of 2.3 million eukaryotic 18S rDNA ribotypes was generated from a subset of 47 sampling sites corresponding to approximately 130,000 OTUs<sup>24</sup> (Fig. 1a). Finally, 5,476 viral "populations" were identified at 43 sites from viral metagenomic contigs, only 39 ( $<0.1\%$ ) of which had been previously observed<sup>25</sup> (see Methods). These genomics data combined across all domains of life and viruses together with carbon export estimates and other environmental parameters were used to explore the relationships between marine biogeochemistry and euphotic plankton communities (see Methods) in the top 150 m of the oligotrophic open ocean. Our study did not include high latitude areas due to the current lack of available molecular data and results should not be extrapolated to deeper depths.

Using a method for regression-based modeling of high multidimensional data in biology (specifically a sparse Partial Least Square analysis - sPLS<sup>29</sup>, Extended data Fig. 1), we detected several plankton lineages for which relative sequence abundance correlated with carbon export and other environmental parameters, most notably with NPP, as expected (Fig. 1b and see Supplementary Table 1). These included diatoms, dinoflagellates and metazoa (zooplankton), lineages classically identified as key contributors to carbon export.

## Plankton community networks associated with carbon export

While the analysis presented in Fig. 1b supports previous findings about key organisms involved in carbon export from the euphotic zone<sup>14,15,17-19</sup>, it is not able to capture how the intrinsic structure of the planktonic community relates to this biogeochemical process. Conversely, although other recent holistic approaches<sup>10,30,31</sup> used species co-occurrence networks to reveal potential biotic interactions, they do not provide a robust description of sub-communities driven by abiotic interactions. To overcome these issues, we applied a systems biology approach known as Weighted Gene Correlation Network Analysis (WGCNA<sup>32,33</sup>) to detect significant associations between the *Tara* Oceans genomics data and carbon export. This method delineates communities in the euphotic zone that are the

most associated with carbon export rather than predicting organisms associated with sinking particles.

In brief, the WGCNA approach builds a network in which nodes are features (in this case plankton lineages or gene functions) and links are evaluated by the robustness of co-occurrence scores. WGCNA then clusters the network into modules (hereafter denoted subnetworks) that can be examined to find significant subnetwork-trait relationships. We then filtered each subnetwork using a Partial Least Square (PLS) analysis that emphasizes key nodes (based on the Variable Importance in Projection (VIP) scores; see Methods and Extended data Fig. 1). These particular nodes are mandatory to summarize a subnetwork (or community) related to carbon export. In particular, they are of interest for evaluating (i) subnetwork robustness and (ii) predictive power for a given trait (see Methods and Extended data Fig. 1).

We applied WGCNA to the relative abundance tables of eukaryotic, prokaryotic and viral lineages<sup>23-25</sup> and identified unique subnetworks significantly associated with carbon export within each dataset (see Methods and Supplementary Tables 2, 3, 4). The eukaryotic subnetwork (subnetwork-trait relationship to carbon export, Pearson cor. = 0.81,  $p = 5e^{-15}$ ) contained 49 lineages (Extended data Fig. 2a and Supplementary Table 2) among which 20% represented photosynthetic organisms (Fig. 2a and Supplementary Table 2). Surprisingly, this small subnetwork's structure correlates very strongly to carbon export (Pearson cor. = 0.87,  $p = 5e^{-16}$ , Extended data Fig. 2d) and it predicts as much as 69% (Leave-One-Out Cross-Validated (LOOCV),  $R^2 = 0.69$ ) of the variability in carbon export (Extended data Fig. 2g). Only ~6% of the subnetwork nodes correspond to diatoms and they show lower VIP scores than dinoflagellates (Supplementary Table 2). This is likely because our samples are not from silicate replete conditions where diatoms were blooming. Furthermore, our analysis did not incorporate data from high latitudes, where diatoms are known to be particularly important for carbon export, so this result suggests that dinoflagellates have a heretofore unrecognized role in carbon export processes in subtropical oligotrophic 'type' ecosystems. More precisely four of the five highest VIP scoring eukaryotic lineages that correlated with carbon export at 150 m were heterotrophs such as Metazoa (copepods), non-photosynthetic Dinophyceae, and Rhizaria (Fig. 2a and Supplementary Table 2). These results corroborate recent metagenomics analysis of microbial communities from sediment traps in the oligotrophic North Pacific subtropical gyre<sup>34</sup>. Consistently, *in situ* imaging surveys have revealed Rhizarian lineages, made up of large fragile organisms such as the Collodaria, to represent an until now under-appreciated component of global plankton biomass<sup>35</sup>, which here also appear to be of relevance for carbon export. Another 14% of lineages from the subnetwork correspond to parasitic organisms, a largely under-explored component of planktonic ecosystems when studying carbon export.

The prokaryotic subnetwork that associated most significantly with carbon export at 150 m (subnetwork-trait relationship to carbon export, Pearson cor. = 0.32,  $p = 9e^{-03}$ ) contained 109 OTUs (Extended data Fig. 2b and Supplementary Table 3), its structure correlated well to carbon export (Pearson cor. = 0.47,  $p = 5e^{-06}$ , Extended data Fig. 2e) and it could predict as much as 60% of the carbon export (LOOCV,  $R^2 = 0.60$ ) (Extended data Fig. 2h). By far

the highest VIP score within this community was assigned to *Synechococcus*, followed by *Cobetia*, *Pseudoalteromonas* and *Idiomarina*, as well as *Vibrio* and *Arcobacter* (Fig. 2b and Supplementary Table 3). Noteworthy, *Prochlorococcus* genera and SAR11 clade fall out of this community, while the significance of *Synechococcus* for carbon export could be validated using absolute cell counts estimated by flow cytometry (Pearson cor. = 0.64,  $p = 4e^{-10}$ , Extended data Fig. 2k). Moreover, *Prochlorococcus* cell counts did not correlate with carbon export (Pearson cor. = -0.13,  $p = 0.27$ , Extended data Fig. 2j) whereas the *Synechococcus* to *Prochlorococcus* cell count ratio correlated positively and significantly (Pearson cor. = 0.54,  $p = 4e^{-07}$ , Extended data Fig. 2l), suggesting the relevance of *Synechococcus*, rather than *Prochlorococcus*, to carbon export. Interestingly, *Pseudoalteromonas*, *Idiomarina*, *Vibrio* and *Arcobacter* (of which several species are known to be associated with eukaryotes<sup>36</sup>) have also been observed in live and poisoned sediment traps<sup>34</sup> and display very high VIP scores in the subnetwork associated with carbon export. Additional genera reported as being enriched in poisoned traps (also known as being associated with eukaryotes) include *Enterovibrio* and *Campylobacter*, and are present as well in the carbon export associated subnetwork.

Interestingly, the viral subnetwork (involving 277 populations) most related to carbon export at 150 m (Pearson cor. = 0.93,  $p = 2e^{-15}$ , Extended data Fig. 2c) contained particularly high VIP scores for two *Synechococcus* phages (Fig. 2c and Supplementary Table 4), which represented a 16-fold enrichment (Fisher's exact test  $p = 6.4e^{-09}$ ). Its structure also correlated with carbon export (Pearson cor. = 0.88,  $p = 6e^{-93}$ , Extended data Fig. 2f) and could predict up to 89% of the variability of carbon export (LOOCV,  $R^2 = 0.89$ ) (Extended data Fig. 2i). The significance of these convergent results is reinforced by the fact that sequences from these datasets are derived from organisms collected on independent size filters (see Methods), and further implicates the importance of top-down processes in carbon export.

With the aim of integrating eukaryotic, prokaryotic, and viral communities in the euphotic zone with carbon export at 150 m, we synthesized their respective subnetworks using a single global co-occurrence network established previously<sup>10</sup>. The resulting network focused on key lineages and their predicted co-occurrences (Fig. 3). Lineages with high VIP values (such as *Synechococcus*) are revealed as hubs of the co-occurrence network<sup>10</sup>, illustrating the potentially strategic key roles within the integrated network of lineages underappreciated by conventional methods to study carbon export. Associations between the hub lineages are mostly mutually exclusive which may explain the relatively weak correlation of some of these lineages with carbon export when using standard correlation analyses as shown in Fig. 1b.

## Gene functions associated with carbon export

Given the potential importance of prokaryotic processes influencing the biological carbon pump<sup>22</sup>, we used the same analytical approaches to examine the prokaryotic genomic functions associated with carbon export at 150 m in the annotated Ocean Microbial Reference Gene Catalogue from *Tara Oceans*<sup>23</sup>. We built a global co-occurrence network for functions (*i.e.*, Orthologous Groups of genes or OGs) from the euphotic zone and identified

two subnetworks of functions that are significantly associated with carbon export (light and dark green subnetworks; FNET1 and FNET2, respectively, see Extended data Fig. 3a, 3b and 3c).

The majority of functions in FNET1 and FNET2 correlate well with carbon export (FNET1: mean Pearson cor. = 0.45, s.d. 0.09 and FNET2: mean Pearson cor. = 0.34, s.d. 0.10). Interestingly, FNET2 functions ( $n=220$ ) encode mostly (83%) core functions (i.e., functions observed in all euphotic samples, see Methods) while the majority of FNET1 functions ( $n=441$ ) are non-core (85%) (see Supplementary Tables 5, 6), highlighting both essential and adaptive ecological functions associated with carbon export. Top VIP scoring functions in the FNET1 subnetwork are membrane proteins such as ABC-type sugar transporters (Extended data Fig. 3c). This subnetwork also contains many functions specific to the *Synechococcus* accessory photosynthetic apparatus (e.g., relating to phycobilisomes, phycocyanin and phycoerythrin; see Supplementary Table 5), which is consistent with the major role of this genus for carbon export inferred from the prokaryotic subnetwork (Fig. 2b). In addition, functions related to carbohydrates, inorganic ion transport and metabolism, as well as transcription, are also well represented (Fig. 4), suggesting overall a subnetwork of functions dedicated to photosynthesis and growth.

The FNET2 subnetwork contains several functions encoded by genes taxonomically assigned to *Candidatus pelagibacter* and *Prochlorococcus*, known as occupying similar oceanic regions as *Synechococcus*, but overall most of its relative abundance (74%) is taxonomically unclassified (Extended data Fig. 3e). Top VIP scoring functions in FNET2 are also membrane proteins and ABC-type sugar transporters, as well as functions involved in carbohydrate breakdown such as a chitinase (Extended data Fig. 3c). These features highlight the potential roles of bacteria in the formation and degradation of marine aggregates<sup>37</sup>. Strikingly, 77% and 58%, of OGs with a VIP score > 1 in FNET1 and FNET2, respectively, are functionally uncharacterized<sup>38,39</sup> (Fig. 4), pointing to the strong need for future molecular work to explore these functions (see Supplementary Tables 5, 6).

The relevance of the identified bacterial functions to predict carbon export was also confirmed by PLS regression (Extended data Fig. 3d). As proposed for plankton communities, the functional subnetworks predict 41% and 48% of carbon export variability (LOOCV,  $R^2 = 0.41$  and  $0.48$  for FNET1 and FNET2, respectively) with a minimal number of functions (Fig. 4, 123 and 54 functions with a VIP score > 1 for FNET1 and FNET2, respectively). Finally, higher predictive power was obtained using subnetworks of viral protein clusters (Extended data Fig. 4a, 4b and 4c), predicting 55% and 89% of carbon export variability (LOOCV  $R^2 = 0.55$  and  $0.89$  for VNET1 and VNET2, respectively; Extended data Fig. 4d, Supplementary Tables 7, 8), suggesting the key role, of not only bacteria, but also their phages in processes sustaining carbon export at a global level.

## Discussion

In this report we reveal the potential contribution of unexpected components of plankton communities, and confirm the importance of prokaryotes and viruses in the correlating with carbon export in the nutrient-depleted oligotrophic ocean. Carbon export at 150 m has been

estimated from particle size distribution in a global dataset, but should be taken with caution, as the estimates do not account for particle composition. In addition, these export estimates evaluate how much carbon leaves the euphotic zone, but they are not related and should not be extrapolated to sequestration, which occurs after remineralisation, deeper in the water column, and over longer timescales. Nonetheless, the use of the UVP was the only realistic method to evaluate carbon flux over the 3 years expedition because deployment of sediment traps at all stations would have been impossible. While our findings are consistent with the numerous previous studies that have highlighted the central role of copepods and diatoms in carbon export<sup>14,15,17-19</sup>, they place them in an ecosystem context and reveal hypothetical processes correlating with the intensity of export, such as parasitism and predation. For example, while viruses are commonly assumed to lyse cells and maintain fixed organic carbon in surface waters, thereby reducing the intensity of the biological carbon pump<sup>40</sup>, there are hints that viral lysis may increase carbon export through the production of colloidal particles and aggregate formation<sup>41</sup>. Our current study suggests that these latter roles may be more ubiquitous than currently appreciated. The importance of aggregation and cell stickiness as inferred from gene network analysis should be further explored mechanistically to investigate the biological significance of these findings.

The future evolution of the oceanic carbon sink remains uncertain because of poorly constrained processes, particularly those associated with the biological pump. With current trends in climate change, the size and biodiversity of phytoplankton are predicted to decrease globally<sup>42,43</sup>. Furthermore, in spite of the potential importance of viruses revealed in this study, they have largely been ignored because of limitations in sampling technologies. Consequently, as oligotrophic gyres expand and global mean NPP decreases<sup>44</sup>, the field is currently unable to predict the consequences for carbon export from the ocean's euphotic zone. By pinpointing key lineages and key microbial functions that correlate with carbon export at 150 m in these areas, this study provides a framework to address this critical bottleneck. However, the associations presented do not necessarily suggest a causal effect on carbon export, which will require further investigation.

One of the grand challenges in the life sciences is to link genes to ecosystems<sup>45</sup>, based on the posit that genes can have predictable ecological footprints at community and ecosystem levels<sup>46-48</sup>. The *Tara* Oceans data sets have allowed us to predict as much as 89% of the variability in carbon export from the oligotrophic surface ocean with just a small number of genes, largely with unknown functions, encoded by prokaryotes and viruses. These findings can be used as a basis to include biological complexity and guide experimental work designed to inform climate modeling of the global carbon cycle. Such statistical analyses, scaling from gene-to-ecosystems, may open the way to the development of a new conceptual and methodological framework to better understand the mechanisms underpinning key ecological processes.

## Methods

### Environmental data collection

From 2009-2013, environmental data (Supplementary Table 9) were collected across all major oceanic provinces in the context of the *Tara* Oceans expeditions<sup>20</sup>. Sampling stations



were selected to represent distinct marine ecosystems at a global scale<sup>49</sup>. Note that Southern Ocean stations were not examined herein because they were ranked as outliers due to their exceptional environmental characteristics and biota<sup>23,24</sup>. Environmental data were obtained from vertical profiles of a sampling package<sup>50,51</sup>. It consisted of conductivity and temperature sensors, chlorophyll and CDOM fluorometers, light transmissometer (Wetlabs C-star 25cm), a backscatter sensor (WetLabs ECO BB), a nitrate sensor (SATLANTIC ISUS) and a Hydroptic Underwater Vision Profiler (UVP; Hydroptics<sup>52</sup>. Nitrate and fluorescence to chlorophyll concentrations as well as salinity were calibrated from water samples collected with Niskin bottle<sup>50</sup>. Net Primary Production (NPP) data were extracted from 8 day composites of the Vertically Generalized Production Model (VGPM<sup>53</sup>) at the week of sampling<sup>54</sup>. Carbon fluxes and carbon export, corresponding to the carbon flux at 150 m, were estimated based on particle concentration and size distributions obtained from the UVP<sup>51</sup> and details are presented below.

### From particle size distribution to carbon export estimation

Previous research has shown that the distribution of particle size follows a power law over the  $\mu\text{m}$  to the mm size range<sup>3,55,56</sup>. This *Junge*-type distribution translates into the following mathematical equation, whose parameters can be retrieved from UVP images:

$$n(d) = ad^k \quad (\text{eq. 1})$$

where  $d$  is the particle diameter, and exponent  $k$  is defined as the slope of the number spectrum when equation (2) is log transformed. This slope is commonly used as a descriptor of the shape of the aggregate size distribution.

The carbon-based particle size approach relies on the assumption that the total carbon flux of particles ( $F$ ) corresponds to the flux spectrum integrated over all particle sizes:

$$F = \int_0^{\infty} n(d) \cdot m(d) \cdot w(d) dd \quad (\text{eq. 2})$$

where  $n(d)$  is the particle size spectrum, i.e., equation (1), and  $m(d)$  is the mass (here carbon content) of a spherical particle described as:

$$m(d) = \alpha d^3 \quad (\text{eq. 3})$$

where  $\alpha = \pi\rho/6$ ,  $\rho$  is the average density of the particle, and  $w(d)$  is the settling rate calculated using Stokes Law:

$$w(d) = \beta d^2 \quad (\text{eq. 4})$$

where  $\beta = g(\rho - \rho_0)(18\nu\rho_0)^{-1}$ ,  $g$  is the gravitational acceleration,  $\rho_0$  the fluid density, and  $\nu$  the kinematic viscosity.

In addition, mass and settling rates of particles,  $m(d)$  and  $w(d)$ , respectively, are often described as power law functions of their diameter obtained by fitting observed data,  $m(d)$ .  $w(d) = Ad^B$ . The particles carbon flux can then be estimated using an approximation of Eq. 2 over a finite number ( $x$ ) of small logarithmic intervals for diameter  $d$  spanning from 250  $\mu\text{m}$  to 1.5 mm (particles  $<250 \mu\text{m}$  and  $>1.5 \text{ mm}$  are not considered, consistent with the method presented by *Guidi et al., [2008]*<sup>7</sup>) such as

$$F = \sum_{i=1}^x n_i A d_i^B \Delta d_i \quad (\text{eq. 5})$$

where  $A=12.5 \pm 3.40$  and  $B=3.81 \pm 0.70$  have been estimated using a global dataset that compared particle fluxes in sediment traps and particle size distributions from the UVP images.

### Genomic data collection

For the sake of consistency between all available datasets from the *Tara* Oceans expeditions, we considered subsets of the data recently published in *Science*<sup>23-25</sup>. In brief, one sample corresponds to data collected at one depth (surface (SRF) or Deep Chlorophyll Maximum (DCM) determined from the profile of chlorophyll fluorometer) and at one station. To study the eukaryotic community in our current manuscript, we selected stations at which we had environmental data and carbon export estimated at 150 m with the UVP and all size fractions. Consequently a subset of 33 stations (corresponding to 56 samples) has been created compared to the 47 stations analyzed in *de Vargas et al. [2015]*. A similar procedure has been applied to the prokaryotic and viral datasets, reducing the *Sunagawa et al. [2015]* prokaryotic dataset to a subset of 104 samples from 62 stations and the *Brum et al. [2015]* viral dataset into a subset of 37 samples from 22 stations (See Supplementary Table 10). In addition a detailed table is provided summarizing which samples (depth and station) are available for each domain (Supplementary Table 11).

### Eukaryotic taxa profiling

Photic-zone eukaryotic plankton diversity has been investigated through millions of environmental Illumina reads. Sequences of the 18S ribosomal RNA gene V9 region were obtained by PCR amplification and a stringent quality-check pipeline has been applied to remove potential chimera or rare sequences (details on data cleaning in *de Vargas et al. [2015]*<sup>24</sup>). For 47 stations, and if possible at two depths (SRF and DCM), eukaryotic communities were sampled in the *piconano-* (0.8-5  $\mu\text{m}$ ), *micro-* (20-180  $\mu\text{m}$ ) and *meso-* plankton (180-2000  $\mu\text{m}$ ) fractions (a detailed list of these samples is given in Supplementary Table 12). In the framework of the carbon export study, sequences from all size fractions were pooled in order to get the most accurate and statistically reliable dataset of the eukaryotic community. The 2.3 million eukaryotic ribotypes were assigned to known eukaryotic taxonomic entities by global alignment to a curated database<sup>24</sup>. To get the most accurate vision of the eukaryotic community, sequences showing less than 97% identity with reference sequences were excluded. The final eukaryotic relative abundance matrix used in our analyses included 1,750 lineages (taxonomic assignation has been performed using a last

common ancestor methodology, and had thus been performed down to species level when possible) in 56 samples from 33 stations. Pooled abundance (number of V9 sequences) of each lineage has been normalized by the total sum of sequences in each sample.

### Prokaryotic taxa profiling

To investigate the prokaryotic lineages, communities were sampled in the pico-plankton. Both filter sizes have been used along the *Tara* Oceans transect: up to station #52, prokaryotic fractions correspond to a 0.22-1.6  $\mu\text{m}$  size fraction, and from station #56, prokaryotic fractions correspond to a 0.22-3  $\mu\text{m}$  size fraction. Prokaryotic taxonomic profiling was performed using 16S rRNA gene tags directly identified in Illumina-sequenced metagenomes ( $_{\text{mi}}$ tags) as described in *Logares et al.*, [2014]<sup>58</sup>. 16S  $_{\text{mi}}$ tags were mapped to cluster centroids of taxonomically annotated 16S reference sequences from the SILVA database<sup>59</sup> (release 115: SSU Ref NR 99) that had been clustered at 97% sequence identity using USEARCH v6.0.307<sup>60</sup>. 16S  $_{\text{mi}}$ tag counts were normalized by the total reads count in each sample (further details in *Sunagawa et al.* [2015]<sup>23</sup>). The photic-zone prokaryotic relative abundance matrix used in our analyses included 3,253,962  $_{\text{mi}}$ tags corresponding to 1,328 genera in 104 samples from 62 stations.

### Prokaryotic functional profiling

For each prokaryotic sample, gene relative abundance profiles were generated by mapping reads to the OM-RGC using the MOCAT pipeline<sup>61</sup>. The relative abundance of each reference gene was calculated as gene length-normalized base counts. And functional abundances were calculated as the sum of the relative abundances of these reference genes, annotated to OG functional groups. In our analyses, we used the subset of the OM-RGC that was annotated to Bacteria or Archaea (24.4 M genes). Using a rarefied (to 33 M inserts) gene count table, an OG was considered to be part of the ocean microbial core if at least one insert from each sample was mapped to a gene annotated to that OG. For further details on the prokaryotic profiling please refer to *Sunagawa et al.* [2015]<sup>23</sup>. The final prokaryotic functional relative abundance matrix used in our analyses included 37,832 OGs or functions in 104 samples from 62 stations. Genes from functions of FNET1 and FNET2 subnetworks were taxonomically annotated using a modified dual BLAST-based last common ancestor (2bLCA) approach<sup>62</sup>. We used RAPsearch2<sup>63</sup> rather than BLAST to efficiently process the large data volume and a database of non-redundant protein sequences from UniProt (version: UniRef\_2013\_07) and eukaryotic transcriptome data not represented in UniRef (see Supplementary Table 5, 6, for full annotations).

### Enumeration of prokaryotes by flow cytometry

For prokaryote enumeration by flow cytometry, three aliquots of 1 ml of seawater (pre-filtered by 200- $\mu\text{m}$  mesh) were collected from both SRF and DCM. The samples were fixed immediately using cold 25% glutaraldehyde (final concentration 0.125%), left in the dark for 10 min at room temperature, flash-frozen and kept in liquid nitrogen on board and then stored at  $-80^{\circ}\text{C}$  on land. Two subsamples were taken to separate counts of heterotrophic prokaryotes (not shown herein) and phototrophic picoplankton. For heterotrophic prokaryote determination, 400  $\mu\text{l}$  of sample was added to a diluted SYTO-13 (Molecular Probes Inc., Eugene, OR, USA) stock (10:1) at 2.5  $\mu\text{mol l}^{-1}$  final concentration, left for about 10 min in

the dark to complete the staining and run in the flow cytometer. We used a FACS Calibur (Becton & Dickinson) flow cytometer equipped with a 15 mW Argon-ion laser (488 nm emission). At least 30,000 events were acquired for each subsample (usually 100,000 events). Fluorescent beads (1  $\mu$ m, Fluoresbrite carboxylate microspheres, Polysciences Inc., Warrington, PA) were added at a known density as internal standards. The bead standard concentration was determined by epifluorescence microscopy. For phototrophic picoplankton, we used the same procedure as for heterotrophic prokaryote, but without addition of SYTO-13. Data analysis was performed with FlowJo software (Tree Star, Inc.).

### Profiling of viral populations

In order to associate viruses to carbon export we used viral populations as defined in *Brum et al. [2015]<sup>25</sup>* using a set of 43 *Tara* Oceans viromes. Briefly, viral populations were defined as large contigs (>10 predicted genes and >10 kb) identified as most likely originating from bacterial or archaeal viruses. These 6,322 contigs remained and were then clustered into populations if they shared more than 80% of their genes at >95% nucleotide identity. This resulted in 5,477 ‘populations’ from the 6,322 contigs, where as many as 12 contigs were included per population. For each population, the longest contig was chosen as the ‘seed’ representative sequence. The relative abundance of each population was computed by mapping all quality-controlled reads to the set of 5,477 non-redundant populations (considering only mapping quality scores greater than 1) with Bowtie2<sup>64</sup> and if more than 75% of the reference sequence was covered by virome reads. The relative abundance of a population in a sample was computed as the number of base pairs recruited to the contig normalized to the total number of base pairs available in the virome and the contig length if more than 75% of the reference sequence was covered by virome reads, and set to 0 otherwise (see *Brum et al. [2015]<sup>25</sup>* for further details). The final viral population abundance matrix used in our analyses included 5,291 viral population contigs in 37 samples from 22 stations.

### Viral host predictions

The longest contig in a population was defined as the seed sequence and considered the best estimate of that population’s origin. These seed sequences were used to assess taxonomic affiliation of each viral population. Cases where >50% of the genes were affiliated to a specific reference genome from RefSeq Virus (based on a BLASTp comparison with thresholds of 50 for bit score and  $10^{-5}$  for e-value) with an identity percentage of at least 75% (at the protein sequence level) were considered as confident affiliations to the corresponding reference virus. The viral population host group was then estimated based on these confident affiliations (see Supplementary Table 13 for host affiliation of viral population contigs associated to carbon export).

### Viral protein clusters

Viral protein clusters (PCs) correspond to ORFs initially mapped to existing clusters (POV, GOS and phage genomes). The remaining, unmapped ORFs were self-clustered, using cd-hit as described in *Brum et al. [2015]<sup>25</sup>*. Only PCs with more than two ORFs were considered bona fide and were used for subsequent analyses. To compute PC relative abundance for statistical analyses, reads were mapped back to predicted ORFs in the contigs dataset using

Mosaik as described in *Brum et al. [2015]*<sup>25</sup>. Read counts to PCs were normalized by sequencing depth of each virome. Importantly, we restricted our analyses to 4,294 PCs associated to the 277 viral population contigs significantly associated to carbon export in 37 samples from 22 stations.

### Sparse Partial Least Squares analysis

In order to directly associate eukaryotic lineages to carbon export and other environmental traits (Fig. 1b), we used sparse Partial Least Square (sPLS<sup>65</sup> as implemented in the R package *mixOmics*<sup>29</sup>. We applied the sPLS in regression mode, which will model a causal relationship between the lineages and the environmental traits, *i.e.* PLS will predict environmental traits (*e.g.* carbon export) from lineage abundances. This approach enabled us to identify high correlations (see Supplementary Table 1) between certain lineages and carbon export but without taking into account the global structure of the planktonic community.

### Co-occurrence network model analysis

Weighted correlation network analysis (WGCNA) was performed to delineate feature (lineages, viral populations, PCs or functions) subnetworks based on their relative abundance<sup>66,67</sup>. A signed adjacency measure for each pair of features was calculated by raising the absolute value of their Pearson correlation coefficient to the power of a parameter  $p$ . The default value  $p=6$  was used for each global network, except for the Prokaryotic functional network where  $p$  had to be lowered to 4 in order to optimize the scale-free topology network fit. Indeed, this power allows the weighted correlation network to show a scale free topology where key nodes are highly connected with others. The obtained adjacency matrix was then used to calculate the topological overlap measure (TOM), which for each pair of features, taking into account their weighted pairwise correlation (direct relationships) and their weighted correlations with other features in the network (indirect relationships). For identifying subnetworks a hierarchical clustering was performed using a distance based on the TOM measure. This resulted in the definition of several subnetworks, each represented by its first principal component.

These characteristic components play a key role in weighted correlation network analysis. On the one hand, the closeness of each feature to its cluster, referred to as the subnetwork membership, is measured by correlating its relative abundance with the first principal component of the subnetwork. On the other hand, association between the subnetworks and a given trait is measured by the pairwise Pearson correlation coefficients between the considered environmental trait and their respective principal components. A similar protocol has been performed on the eukaryotic relative abundance matrix, the prokaryotic relative abundance matrix, the prokaryotic functions relative abundance matrix and the viral population and PC relative abundance matrices. All procedures were applied on Hellinger-transformed log-scaled abundances. Noteworthy, the protocol is not sensitive to copy number variation as observed across different eukaryotic species, because the association between two species relies on a correlation score between relative abundance measurements. Computations were carried out using the R package *WGCNA*<sup>33</sup>.

Given the nature of the eukaryotic dataset (three distinct size fractions), the sampling process may lead to the loss of size fractions. In particular, samples #1, #3, #17, #37, #39, #43, #48, #53, #54, #55, #66 are eventually biased by such a loss (Supplementary Table 12). A complementary WGCNA analysis was performed with addition of these samples to evaluate the robustness of our protocol to missing size fractions. The composition of the eukaryotic subnetwork built with an extended dataset (*i.e.*, 67 samples from 37 stations for which size fractions were missing in 11 samples) was compared to the subnetwork as presented above (*i.e.*, 56 samples from 33 stations). Both subnetworks shown an overlap of 75% of lineage, whereas four of the top five VIP lineages with the extended dataset (see Extended data Fig. 5 for details) can be found in the top six VIP lineages of the above subnetwork (Supplementary Table 2), emphasizing highly similar results and a small sensitivity to size fraction loss.

### Extraction of subnetworks related to carbon export

For each subnetwork (called modules within WGCNA) extracted from each global network, pairwise Pearson correlation coefficients between the subnetwork principal components and the carbon export estimation was computed, as well as corresponding p-values corrected for multiple testing using the Benjamini & Hochberg FDR procedure. The subnetworks showing the highest correlation scores are of interest and were investigated. One subnetwork (49 nodes) was significant within the eukaryotic network; one subnetwork (109 nodes) was significant for the prokaryotic network; one subnetwork (277 nodes) was significant within the virus network; two subnetworks (441 and 220 nodes) were significant within the prokaryotic functional network, and two subnetworks (1,879 and 2,147 nodes) were significant within the viral PCs network.

### Partial Least Squares regression

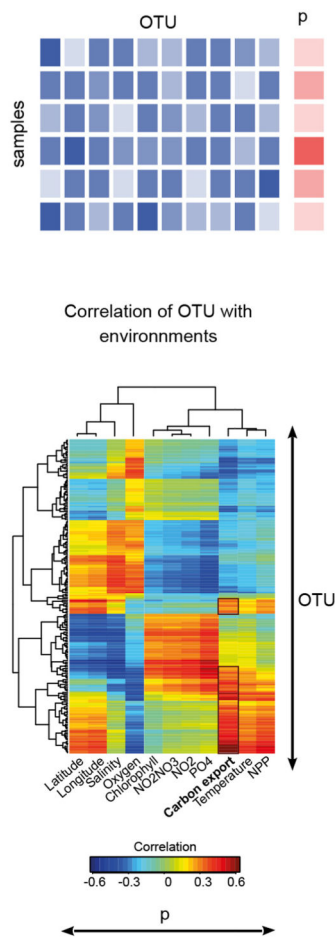
In addition to the network analyses, we asked whether the identified subnetworks can be used as predictors for the carbon export estimations. To answer this question, we used Partial least squares (PLS) regression, which is a dimensionality-reduction method that aims at determining predictor combinations with maximum covariance with the response variable. The identified combinations, called latent variables, are used to predict the response variable. The predictive power of the model is assessed by correlating the predicted vector with the measured values. The significance of the prediction power was evaluated by permuting the data 10,000 times. For each permutation, a PLS model was built to predict the randomized response variable and a Pearson correlation was calculated between the permuted response variable and in Leave-One-Out Cross-Validation (LOOCV) predicted values. The 10,000 random correlations are compared to the performance of the PLS model that were used to predict the true response variable. In addition, the predictors were ranked according to their value importance in projection (VIP)<sup>68</sup>. The VIP measure of a predictor estimates its contribution in the PLS regression. The predictors having high VIP values are assumed important for the PLS prediction of the response variable. The VIP values of the prokaryotic functional subnetworks are provided in Supplementary Tables 5, 6. For the sake of illustration, only lineages or functions with  $VIP > 1$ <sup>68</sup> are discussed and pictured in Figure 2 and 4. Our computations were carried out using the R package *pls*<sup>69</sup>. All programs are available under GPL Licence.

## Subnetwork representations

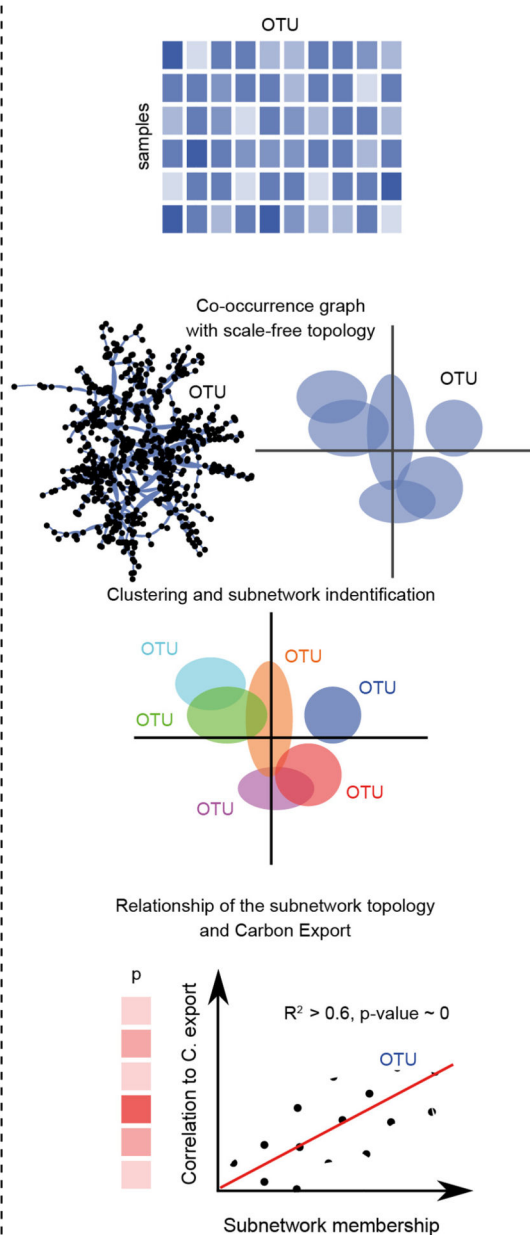
Nodes of the subnetworks represent either lineages (eukaryotic, prokaryotic or viral) or functions (prokaryotic or viral). Subnetworks related to the carbon export have been represented in two distinct formats. Scatter plots represent each nodes based on their Pearson correlation to the carbon export and their respective node centrality within the subnetwork. The latter has been recomputed using significant Spearman correlations above 0.3 (>0.9 for viral PCs) as edges, this is done for visualization purposes since WGCNA subnetworks (based on the Topology Overlap Measure (TOM) between nodes) are hyper-connected. Size representation of nodes are proportional to the VIP score after PLS. The hiveplots depict the same subnetworks by focusing on two main features: x-axis and y-axis depict nodes of subnetworks ranked by their VIP scores and Pearson correlation to the carbon export, respectively.

## Extended Data

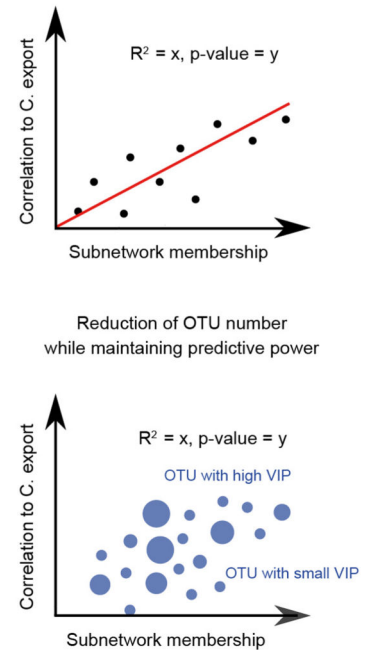
## a Pairwise approach



## b Graph-based approach (WGCNA)



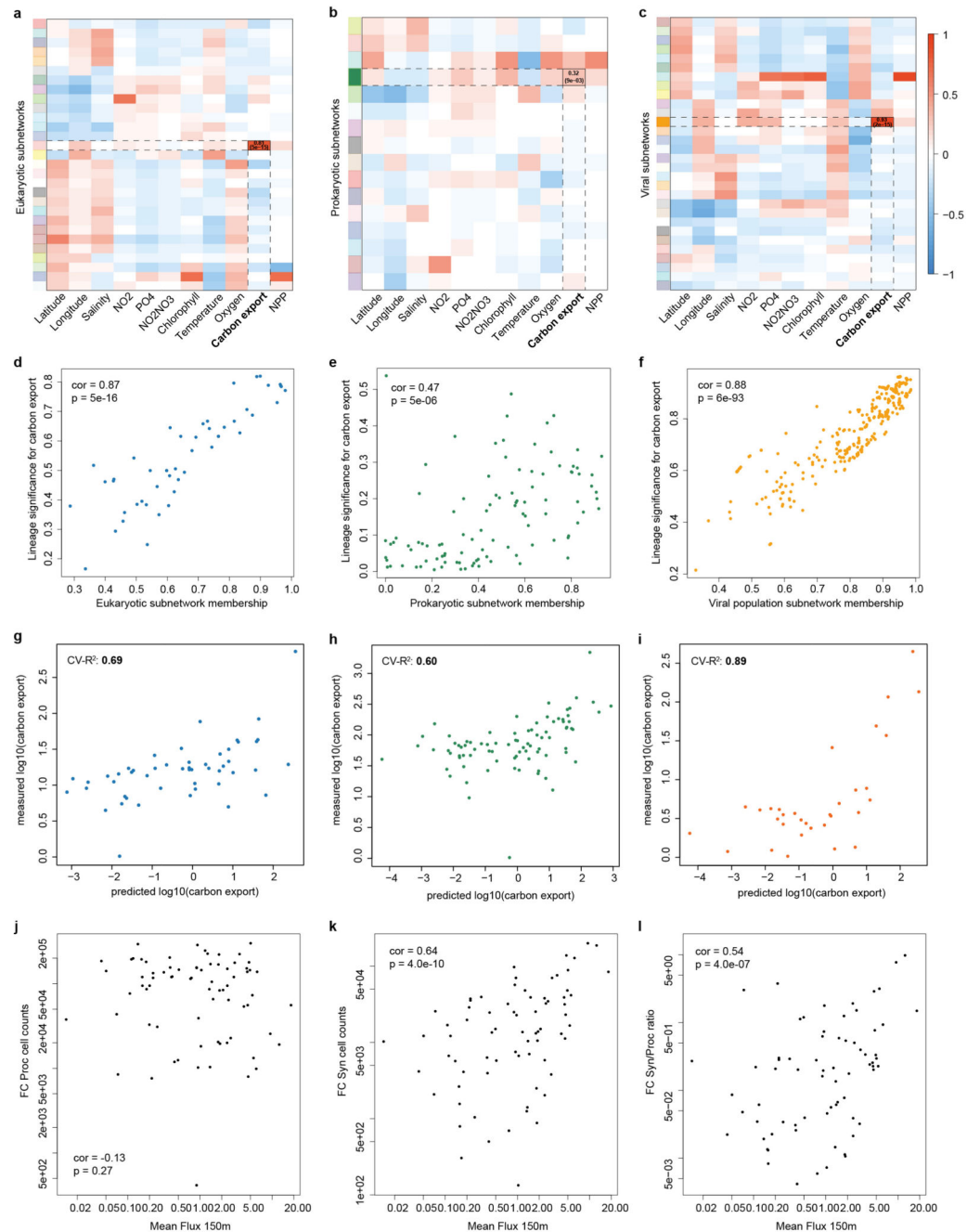
## c Machine learning technique (PLS)

**Extended Data Figure 1:**

Overview of analytical methods used in the manuscript. **a**, Depiction of a standard pairwise analysis that considers a sequence relative abundance matrix for  $s$  samples ( $s \times \text{OTUs}$  (Operational Taxonomic Units)) and its corresponding environmental matrix ( $s \times p$  (parameters)). sPLS results emphasize OTU(s) that are the most correlated to environmental parameters. **b**, Depiction of a graph-based approach. Using only a relative abundance matrix ( $s \times \text{OTUs}$ ), WGCNA builds a graph where nodes are OTUs and edges represent significant

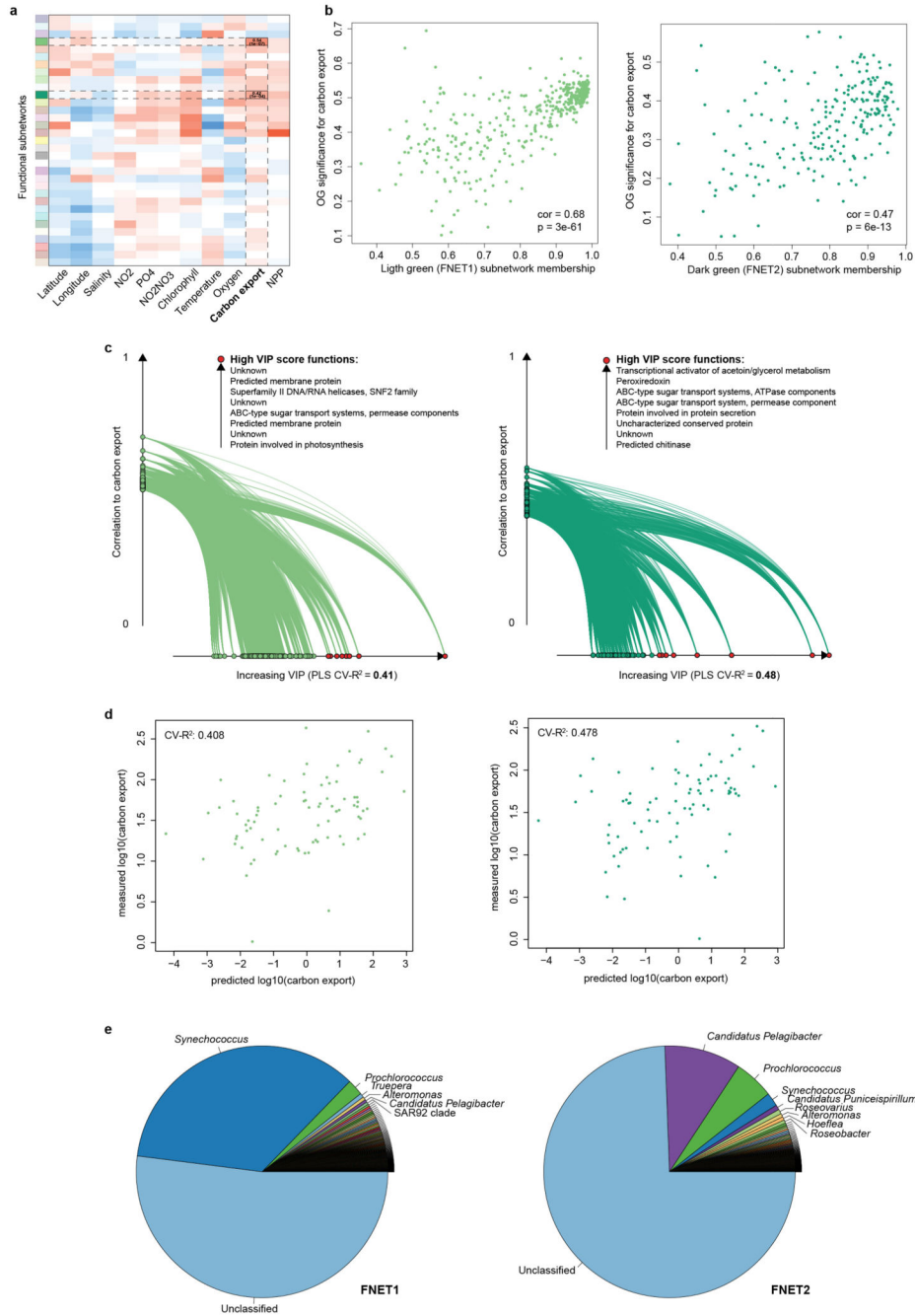


co-occurrence. Co-occurrence scores between nodes are weights allocated to corresponding edges. These weights are magnified by a power-law function until the graph becomes scale-free. The graph is then decomposed within subnetworks (groups of OTUs) that are analyzed separately. One subnetwork (group of OTUs) is considered of interest when its topology is related to the trait of interest; in the current case carbon export. For each subnetwork (for instance the subnetwork related to carbon export), each OTU is spread within a feature space that plots each OTU based on its membership to the subnetwork (x-axis) and its correlation to the environmental trait of interest (i.e., carbon export). A good regression of all OTUs emphasizes the putative relation of the subnetwork topology and the carbon export trait (*i.e.* the more a given OTU defines the subnetwork topology, the more it is correlated to carbon export). **c**, Depiction of the machine learning (PLS) approach that was applied following subnetwork identification and selection. Greater VIP scores (*i.e.* larger circles) emphasized most important OTUs. VIP refers to Variable Importance in Projection and reflects the relative predictive power of a given OTU. OTUs with VIP score greater than 1 are considered as important in the predictive model and their selection do not alter the overall predictive power.

**Extended Data Figure 2:**

Lineage ecological subnetworks associated to environmental parameters and their structures correlating to carbon export. **a,b,c**, Global ecological networks were built using the WGCNA methodology (see methods) and correlated to classical oceanographic parameters as well as carbon export (estimated at 150 m from particles size distribution and abundance). Each domain-specific global network is decomposed into smaller coherent subnetworks (depicted by distinct colours on the y-axis) and their eigen vector is correlated to all environmental parameters. Similar to a correlation at the network scale, this approach

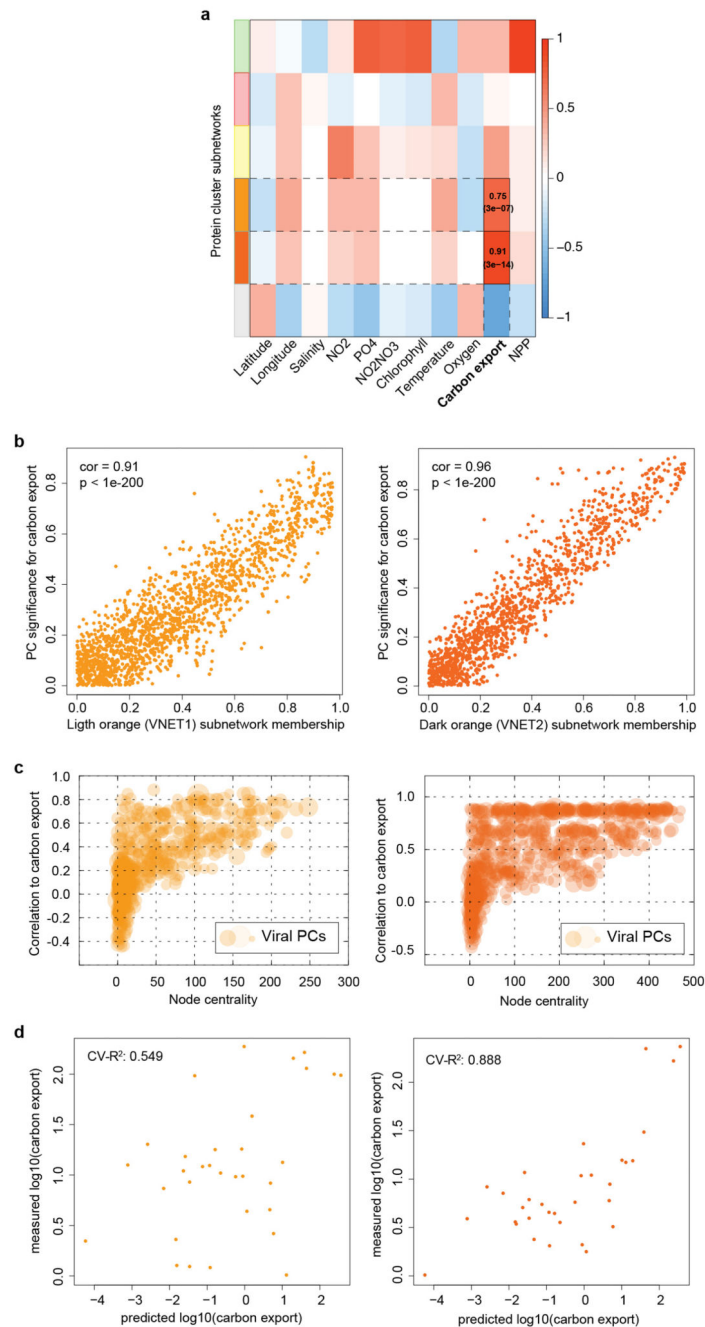
directly links subnetworks to environmental parameters (*i.e.* the more the taxa contribute to the subnetwork structure, the more their abundance are correlated to the parameter). **a**, A single eukaryotic subnetwork ( $n=58$ ,  $N=1'870$ ) is strongly associated to carbon export (Pearson cor. 0.81,  $p = 5e^{-15}$ ). **b**, A single prokaryotic subnetwork ( $n=109$ ,  $N=1'527$ ) is moderately associated to carbon export (Pearson cor. 0.32,  $p = 9e^{-03}$ ). **c**, A single viral subnetwork ( $n=277$ ,  $N=5'476$ ) is strongly associated to carbon export (Pearson cor. 0.93,  $p = 2e^{-15}$ ). **d,e,f**, The WGCNA approach directly links subnetworks to environmental parameters, *i.e.* the more the features contribute to the subnetwork structure (topology), the more their abundance are correlated to the parameter. This measure allows to identify subnetworks for which the overall structure, summarized as the eigen vector of the subnetwork, is related to the carbon export. **d**, The eukaryotic subnetwork structure correlates to carbon export (Pearson cor. = 0.87,  $p = 5e^{-16}$ ). **e**, The prokaryotic subnetwork structure correlates to carbon export (Pearson cor. = 0.47,  $p = 5e^{-06}$ ). **f**, The viral population subnetwork structure correlates to carbon export (Pearson cor. = 0.88,  $p = 6e^{-93}$ ). **g,h,i**, Lineage subnetworks predict carbon export. PLS regression was used to predict carbon export using lineage abundances in selected subnetworks. LOOCV was performed and VIP scores computed for each lineage. **g**, The eukaryotic subnetwork predicts carbon export with a  $R^2$  of 0.69. **h**, The prokaryotic subnetwork predicts carbon export with a  $R^2$  of 0.60. **i**, The viral population subnetwork predicts carbon export with a  $R^2$  of 0.89. **j, k, l**, *Synechococcus* (rather than *Prochlorococcus*) absolute cell counts correlate well to carbon export. **j**, *Prochlorococcus* cell counts estimated by flow cytometry do not correlate to carbon export (mean carbon flux at 150m, Pearson cor. =  $-0.13$ ,  $p = 0.27$ ). **k**, *Synechococcus* cell counts estimated by flow cytometry correlate significantly to carbon export (Pearson cor. = 0.64,  $p = 4.0e^{-10}$ ). **l**, *Synechococcus / Prochlorococcus* cell counts ratio correlates significantly to carbon export (Pearson cor. = 0.54,  $p = 4.0e^{-07}$ ).



**Extended Data Figure 3:**

Prokaryotic function subnetworks associated to environmental parameters and their structure correlate to carbon export. **a,b,c** Global ecological networks were built for the prokaryotic functions using the WGCNA methodology (see methods) and correlated to classical oceanographic parameters as well as carbon export. **a**, Two bacterial functional subnetworks ( $n=441$  and  $n=220$ ,  $N=37'832$ ) are associated to carbon export (Pearson cor. 0.54,  $p = 1e^{-07}$  and 0.42,  $p = 1e^{-04}$ ). **b**, The WGCNA approach directly links subnetworks to environmental parameters, *i.e.* the more the features contribute to the subnetwork structure (topology), the

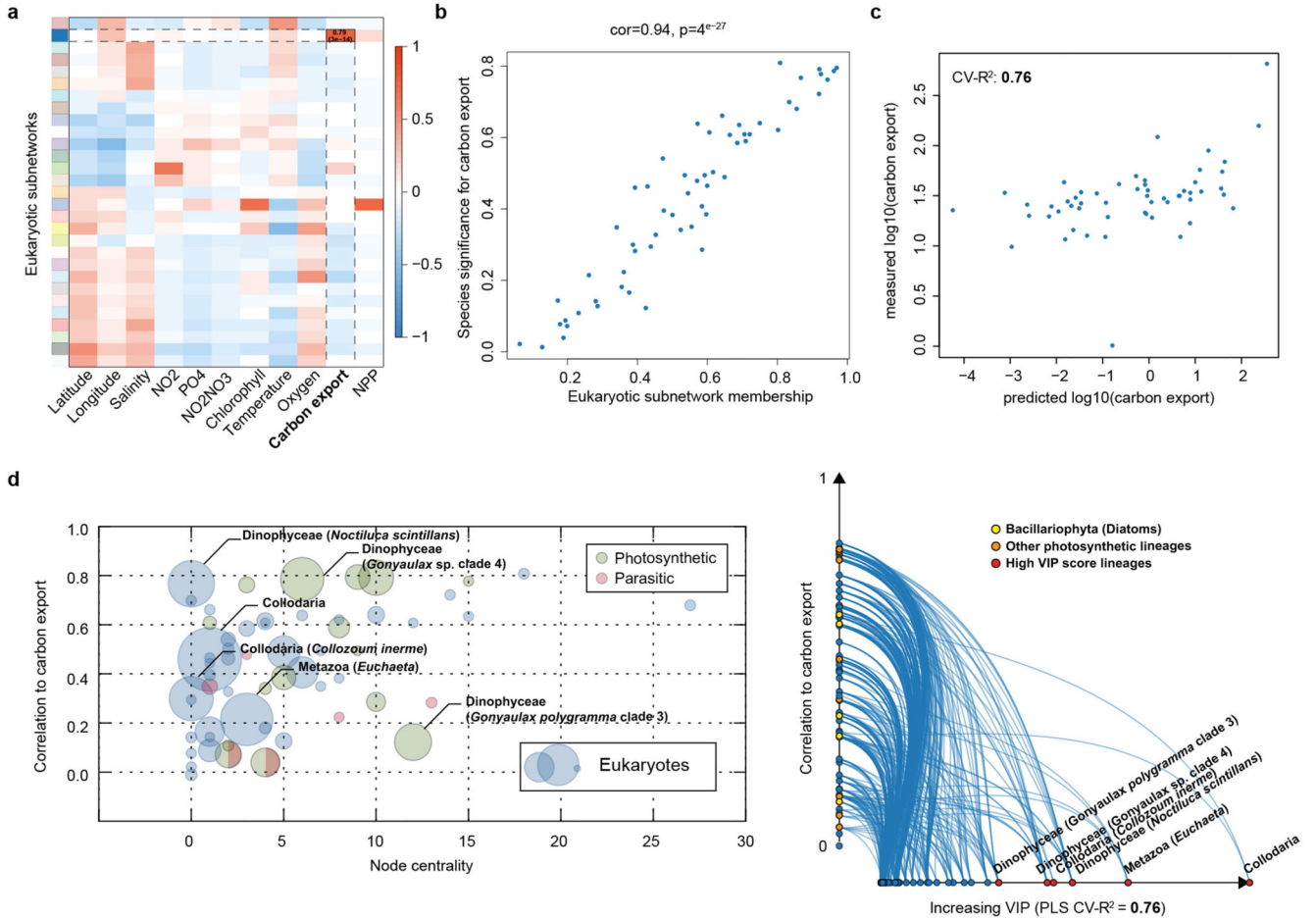
more their abundance are correlated to the parameter. This measure allows to identify subnetworks for which the overall structure, summarized as the eigen vector of the subnetwork, is related to the carbon export. The bacterial function subnetwork structures correlate to carbon export (FNET1 Pearson cor. = 0.68,  $p = 3e^{-61}$ , and FNET2 Pearson cor. = 0.47,  $p = 6e^{-13}$ ). **c**, Two functional subnetworks (light and dark green, FNET1 ( $n=220$ ) and FNET2 ( $n=441$ ), respectively) are significantly associated with carbon export (FNET1: Pearson cor. 0.42,  $p = 4e^{-09}$  and FNET2: 0.54,  $p = 7e^{-06}$ ). The highest VIP score functions from top to bottom correspond to red dots from right to left. **d**, PLS regression was used to predict carbon export using abundances of functions (OGs) in selected subnetworks. LOOCV was performed and VIP scores computed for each function. Light green subnetwork (FNET1) functions predict carbon export with a  $R^2$  of 0.41. Dark green subnetwork (FNET2) functions predict carbon export with a  $R^2$  of 0.48. **e**, Cumulative abundance of genus-level taxonomic annotations of genes encoding functions from FNET1 and FNET2 subnetworks and Bacterial function subnetworks predict carbon export. Genes contributing to the relative abundance of FNET1 and FNET2 subnetwork functions were taxonomically annotated by homology searches against a non-redundant gene reference database using a last common ancestor (LCA) approach (see methods).



#### Extended Data Figure 4:

Viral protein cluster networks reveal potential marker genes for carbon export prediction at global scale. **a**, A viral protein cluster (PC) network was built using abundances of PCs predicted from viral population contigs associated to carbon export (Fig. 2c) using the WGCNA methodology (see methods) and correlated to classical oceanographic parameters. Two viral PC subnetworks ( $n=1'879$  and  $n=2'147$ ,  $N=4'678$ , light and dark orange, VNET1 and VNET2, left and right panel respectively) are strongly associated to carbon export (VNET1: Pearson cor. 0.75,  $p = 3e^{-07}$  and VNET2: 0.91,  $p = 3e^{-14}$ ). **b**, The viral PC

subnetwork structures correlate to carbon export (VNET1 Pearson cor. = 0.91,  $p < 1e^{-200}$ , and VNET2 Pearson cor. = 0.96,  $p < 1e^{-200}$ ). **c**, Size of dots is proportional to the VIP score computed for the PLS regression. **d**, Viral PC subnetworks predict carbon export. PLS regression was used to predict carbon export using abundances of viral protein clusters (PCs) in selected subnetworks. LOOCV was performed and VIP scores computed for each PC. Light orange subnetwork (VNET1, left panel) PCs predict carbon export with a  $R^2$  of 0.55. Dark orange subnetwork (VNET2, right panel) PCs predict carbon export with a  $R^2$  of 0.89.



#### Extended Data Figure 5:

WGCNA and PLS regression analyses for the full Eukaryotic dataset. **a**, A single eukaryotic subnetwork ( $n=58$ , is strongly associated to carbon export (Pearson cor. 0.79,  $p = 3e^{-14}$ ). **b**, The eukaryotic subnetwork structure correlates to carbon export (Pearson cor. = 0.94,  $p = 4e^{-27}$ ). **c**, The eukaryotic subnetwork predicts carbon export with a  $R^2$  of 0.76. **d**, Lineages with the highest VIP score (dots size is proportional to the VIP score in the scatter plot) in the PLS are depicted as red dots corresponding to two rhizaria (Collodaria), one copepod (*Euchaeta*), and three dinophyceae (*Noctiluca scintillans*, *Gonyaulax polygramma* and *Gonyaulax* sp. (clade 4)).

## Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

## Acknowledgements

We thank the commitment of the following people and sponsors: CNRS (in particular Groupement de Recherche GDR3280), European Molecular Biology Laboratory (EMBL), Genoscope/CEA, VIB, Stazione Zoologica Anton Dohrn, UNIMIB, Fund for Scientific Research – Flanders, Rega Institute, KU Leuven, The French Ministry of Research, the French Government 'Investissements d'Avenir' programmes OCEANOMICS (ANR-11-BTBR-0008), FRANCE GENOMIQUE (ANR-10-INBS-09-08), MEMO LIFE (ANR-10-LABX-54), PSL\* Research University (ANR-11-IDEX-0001-02), ANR (projects POSEIDON/ANR-09-BLAN-0348, PHYTBACK/ANR-2010-1709-01, PROMETHEUS/ANR-09-PCS-GENM-217, TARA-GIRUS/ANR-09-PCS-GENM-218, SAMOSA, ANR-13-ADAP-0010), European Union FP7 (MicroB3/No.287589, IHMS/HEALTH-F4-2010-261376), ERC Advanced Grant Award to CB (Diatomite: 294823), Gordon and Betty Moore Foundation grant (#3790 and #2631) and the UA Technology and Research Initiative Fund and the Water, Environmental, and Energy Solutions Initiative to MBS, Spanish Ministry of Science and Innovation grant CGL2011-26848/BOS MicroOcean PANGENOMICS to SGA, TANIT (CONES 2010-0036) from the Agència de Gestió d'Ajuts Universitaris i Reserca to SGA, JSPS KAKENHI Grant Number 26430184 to HO, and FWO, BIO5, Biosphere 2 to MBS. We also thank the support and commitment of Agnès b. and Etienne Bourgois, the Veolia Environment Foundation, Region Bretagne, Lorient Agglomeration, World Courier, Illumina, the EDF Foundation, FRB, the Prince Albert II de Monaco Foundation, the Tara schooner and its captains and crew. We thank MERCATOR-CORIOLIS and ACRI-ST for providing daily satellite data during the expedition. We are also grateful to the French Ministry of Foreign Affairs for supporting the expedition and to the countries who graciously granted sampling permissions. Tara Oceans would not exist without continuous support from 23 institutes (<http://oceans.taraexpeditions.org>). The authors further declare that all data reported herein are fully and freely available from the date of publication, with no restrictions, and that all of the samples, analyses, publications, and ownership of data are free from legal entanglement or restriction of any sort by the various nations whose waters the *Tara* Oceans expedition sampled in. This article is contribution number ZZZ of *Tara* Oceans.

## Author Contributions

L.G., S.C., Lu.B. and D.E. designed the study and wrote the paper. C.D., M.P., J.P. and Sa.S. collected *Tara* Oceans samples. S.K-L managed the logistics of the *Tara* Oceans project. L.G. and M.P. analysed oceanographic data. S.C. and Lu.B. analysed taxonomic data. S.C., Lu.B., D.E. and S.R. performed the genomic and statistical analyses. A.L., Y.D., L.G., S.C., Lu.B. and D.E. produced and analysed the networks. E.K., C.B. and G.G. supervised the study. M.S., J.R., E.K., C.B. and G.G. provided constructive comments, revised and edited the manuscript. *Tara* Oceans coordinators provided a creative environment and constructive criticism throughout the study. All authors discussed the results and commented on the manuscript.

## References and Notes

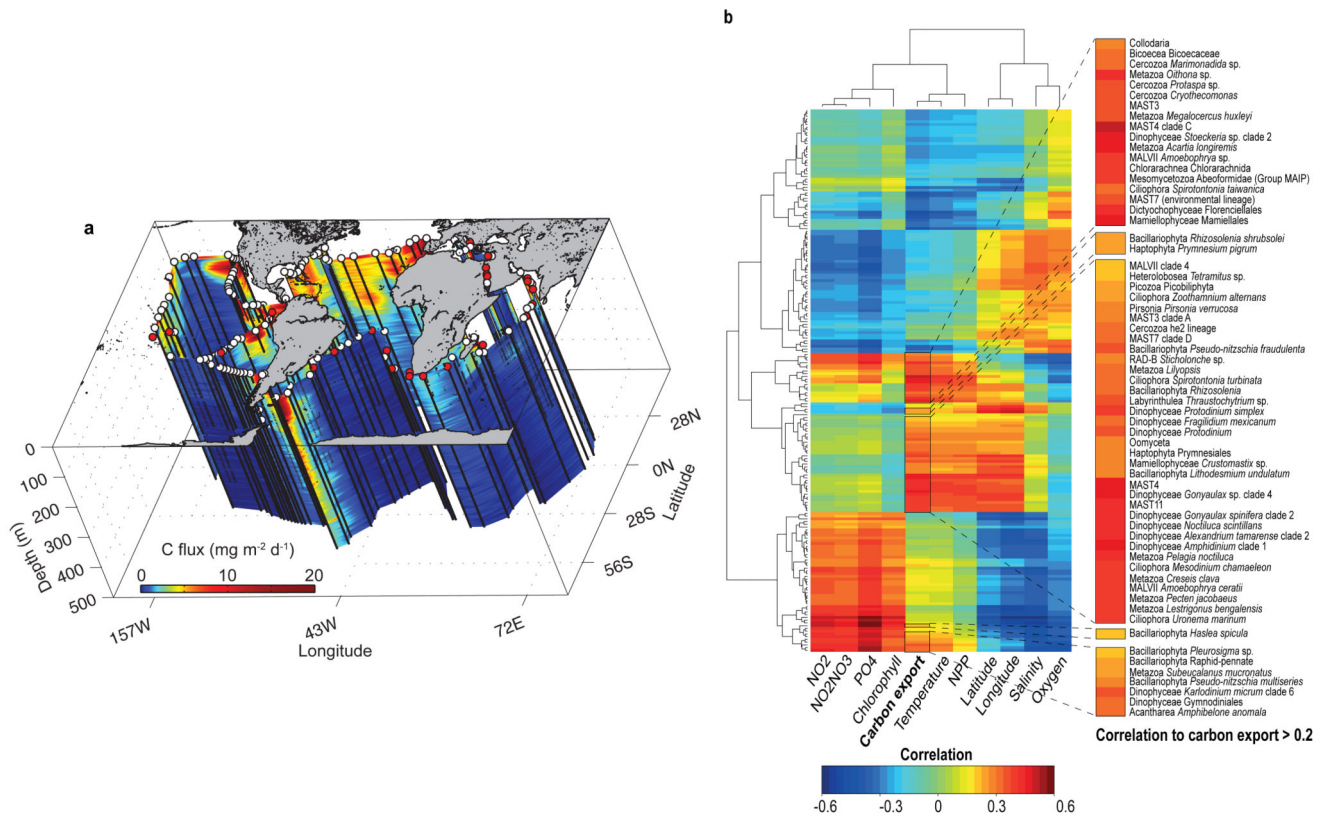
1. Field CB, Behrenfeld MJ, Randerson JT, Falkowski P. Primary production of the biosphere: Integrating terrestrial and oceanic components. *Science*. 1998; 281:237–240. doi:10.1126/Science.281.5374.237. [PubMed: 9657713]
2. Boyd PW, Newton P. Evidence of the potential Influence of planktonic community structure on the interannual variability of particulate organic-carbon flux. *Deep-Sea Res. I*. 1995; 42:619–639.
3. Guidi L, et al. Effects of phytoplankton community on production, size, and export of large aggregates: A world-ocean analysis. *Limnol. Oceanogr.* 2009; 54:1951–1963.
4. Kwon EY, Primeau F, Sarmiento JL. The impact of remineralization depth on the air-sea carbon balance. *Nat Geosci.* 2009; 2:630–635.
5. IPCC. Climate Change 2013: The Physical Science Basis. Contribution of Working Group I to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change. Cambridge University Press; 2013.



6. Kitano H. Biological robustness. *Nat Rev Genet.* 2004; 5:826–837. doi:10.1038/Nrg1471. [PubMed: 15520792]
7. Suweis S, Simini F, Banavar JR, Maritan A. Emergence of structural and dynamical properties of ecological mutualistic networks. *Nature.* 2013; 500:449–452. doi:10.1038/Nature12438. [PubMed: 23969462]
8. Chow CET, Kim DY, Sachdeva R, Caron DA, Fuhrman JA. Top-down controls on bacterial community structure: microbial network analysis of bacteria, T4-like viruses and protists. *ISME J.* 2014; 8:816–829. doi:10.1038/Ismej.2013.199. [PubMed: 24196323]
9. Fuhrman JA. Microbial community structure and its functional implications. *Nature.* 2009; 459:193–199. doi:10.1038/Nature08058. [PubMed: 19444205]
10. Lima-Mendez G, et al. Determinants of community structure in the global plankton interactome. *Science.* 2015; 348 doi:10.1126/science.1262073.
11. Giering SLC, et al. Reconciliation of the carbon budget in the ocean's twilight zone. *Nature.* 2014; 507:480–483. [PubMed: 24670767]
12. Azam F. Microbial control of oceanic carbon flux: The plot thickens. *Science.* 1998; 280:694–696.
13. Agusti S, et al. Ubiquitous healthy diatoms in the deep sea confirm deep carbon injection by the biological pump. *Nat Commun.* 2015; 6 doi:10.1038/Ncomms8608.
14. Sancetta C, Villareal T, Falkowski P. Massive Fluxes of Rhizosolenid Diatoms - a Common Occurrence. *Limnol. Oceanogr.* 1991; 36:1452–1457.
15. Scharek R, Tupas LM, Karl DM. Diatom fluxes to the deep sea in the oligotrophic north Pacific gyre at station ALOHA. *Mar. Ecol. Prog. Ser.* 1999; 182:55–67. doi:10.3354/meps182055.
16. Omand MM, et al. Eddy-driven subduction exports particulate organic carbon from the spring bloom. *Science.* 2015; 348:222–225. doi:10.1126/science.1260062. [PubMed: 25814062]
17. Richardson TL, Jackson GA. Small phytoplankton and carbon export from the surface ocean. *Science.* 2007; 315:838–840. [PubMed: 17289995]
18. Steinberg DK, et al. Bacterial vs. Limnol. *Oceanogr.* 2008; 53:1327–1338.
19. Turner JT. Zooplankton fecal pellets, marine snow, phytodetritus and the ocean's biological pump. *Prog. Oceanogr.* 2015; 130:205–248. doi:10.1016/j.pocean.2014.08.005.
20. Karsenti E, et al. A Holistic Approach to Marine Eco-Systems Biology. *Plos Biol.* 2011; 9 doi: 10.1371/journal.pbio.1001177.
21. Strom SL. Microbial ecology of ocean biogeochemistry: A community perspective. *Science.* 2008; 320:1043–1045. doi:10.1126/Science.1153527. [PubMed: 18497289]
22. Worden AZ, et al. Rethinking the marine carbon cycle: Factoring in the multifarious lifestyles of microbes. *Science.* 2015; 347:1257594. doi:10.1126/Science.1257594. [PubMed: 25678667]
23. Sunagawa S, et al. Structure and function of the global ocean microbiome. *Science.* 2015; 348 doi: 10.1126/science.1261359.
24. de Vargas C, et al. Eukaryotic plankton diversity in the sunlit ocean. *Science.* 2015; 348 doi: 10.1126/science.1261605.
25. Brum JR, et al. Patterns and ecological drivers of ocean viral communities. *Science.* 2015; 348 doi: 10.1126/science.1261498.
26. Bork P, et al. Tara Oceans studies plankton at PLANETARY SCALE. *Science.* 2015; 348:873–873. doi:10.1126/science.aac5605. [PubMed: 25999501]
27. Honjo S, Manganini SJ, Krishfield RA, Francois R. Particulate organic carbon fluxes to the ocean interior and factors controlling the biological pump: A synthesis of global sediment trap programs since 1983. *Prog. Oceanogr.* 2008; 76:217–285. doi:10.1016/j.pocean.2007.11.003.
28. Henson SA, Sanders R, Madsen E. Global patterns in efficiency of particulate organic carbon export and transfer to the deep ocean. *Global. Biogeochem. Cy.* 2012; 26 doi: 10.1029/2011GB004099.
29. Lê Cao KA, Rossouw D, Robert-Granié C, Besse P. A Sparse PLS for Variable Selection when Integrating Omics Data. *Stat Appl Genet Mol.* 2008; 7 doi:10.2202/1544-6115.1390.
30. Chaffron S, Rehrauer H, Pernthaler J, von Mering C. A global network of coexisting microbes from environmental and whole-genome sequence data. *Genome Res.* 2010; 20:947–959. doi: 10.1101/Gr.104521.109. [PubMed: 20458099]

31. Faust K, Raes J. Microbial interactions: from networks to models. *Nat. Rev. Microbiol.* 2012; 10:538–550. doi:10.1038/Nrmicro2832. [PubMed: 22796884]
32. Aylward FO, et al. Microbial community transcriptional networks are conserved in three domains at ocean basin scales. *Proceedings of the National Academy of Sciences.* 2015 doi:10.1073/pnas.1502883112.
33. Langfelder P, Horvath S. WGCNA: an R package for weighted correlation network analysis. *Bmc Bioinformatics.* 2008; 9
34. Fontanez KM, Eppley JM, Samo TJ, Karl DM, DeLong EF. Microbial community structure and function on sinking particles in the North Pacific Subtropical Gyre. *Front Microbiol.* 2015; 6 ArtN 469, doi:10.3389/Fmicb.2015.00/169.
35. Biard T, et al. *In situ* imaging reveals the biomass of large protists in the global ocean. *Nature.* (submitted).
36. Thomas T, et al. Analysis of the *Pseudoalteromonas tunicata* Genome Reveals Properties of a Surface-Associated Life Style in the Marine Environment. *PLoS ONE.* 2008; 3 doi:10.1371/journal.pone.0003252.
37. Azam F, Malfatti F. Microbial structuring of marine ecosystems. *Nat. Rev. Microbiol.* 2007; 5:782–791. doi:10.1038/nrmicro1747. [PubMed: 17853906]
38. Shi YM, Tyson GW, DeLong EF. Metatranscriptomics reveals unique microbial small RNAs in the ocean's water column. *Nature.* 2009; 459:266–U154. doi:10.1038/nature08055. [PubMed: 19444216]
39. Yooseph S, et al. The Sorcerer II Global Ocean Sampling expedition: Expanding the universe of protein families. *Plos Biol.* 2007; 5:432–466. doi:10.1371/journal.pbio.0050016.
40. Suttle CA. Marine viruses - major players in the global ecosystem. *Nat. Rev. Microbiol.* 2007; 5:801–812. doi:10.1038/Nrmicro1750. [PubMed: 17853907]
41. Weinbauer MG. Ecology of prokaryotic viruses. *Fems Microbiol Rev.* 2004; 28:127–181. doi: 10.1016/j.femsre.2003.08.001. [PubMed: 15109783]
42. Finkel ZV, et al. Phytoplankton in a changing world: cell size and elemental stoichiometry. *J. Plankton Res.* 2010; 32:119–137.
43. Sommer U, Lewandowska A. Climate change and the phytoplankton spring bloom: warming and overwintering zooplankton have similar effects on phytoplankton. *Glob. Change Biol.* 2011; 17:154–162. doi:10.1111/J.1365-2486.2010.02182.X.
44. Behrenfeld MJ, et al. Climate-driven trends in contemporary ocean productivity. *Nature.* 2006; 444:752–755. [PubMed: 17151666]
45. DeLong EF, et al. Community genomics among stratified microbial assemblages in the ocean's interior. *Science.* 2006; 311:496–503. doi:10.1126/Science.1120250. [PubMed: 16439655]
46. Gianoulis TA, et al. Quantifying environmental adaptation of metabolic pathways in metagenomics. *P. Natl. Acad. Sci. USA.* 2009; 106:1374–1379. doi:10.1073/Pnas.0808022106.
47. Tilman D, et al. The influence of functional diversity and composition on ecosystem processes. *Science.* 1997; 277:1300–1302. doi:10.1126/Science.277.5330.1300.
48. Wymore AS, et al. Genes to ecosystems: exploring the frontiers of ecology with one of the smallest biological units. *New Phytol.* 2011; 191:19–36. doi:10.1111/J.1469-8137.2011.03730.X. [PubMed: 21631507]
49. Pesant S, et al. Open science resources for the discovery and analysis of Tara Oceans data. *Scientific Data.* 2015; 2:150023. doi:10.1038/sdata.2015.23. [PubMed: 26029378]
50. Picheral, M., et al. Vertical profiles of environmental parameters measured on discrete water samples collected with Niskin bottles during the Tara Oceans expedition 2009-2013. 2014. doi: 10.1594/PANGAEA.836319
51. Picheral, M., et al. Vertical profiles of environmental parameters measured from physical, optical and imaging sensors during Tara Oceans expedition 2009-2013. 2014. doi:10.1594/PANGAEA.836321
52. Picheral M, et al. The Underwater Vision Profiler 5: An advanced instrument for high spatial resolution studies of particle size spectra and zooplankton. *Limnol. Oceanogr. Meth.* 2010; 8:462–473. doi:10.4319/lom.2010.8.462.

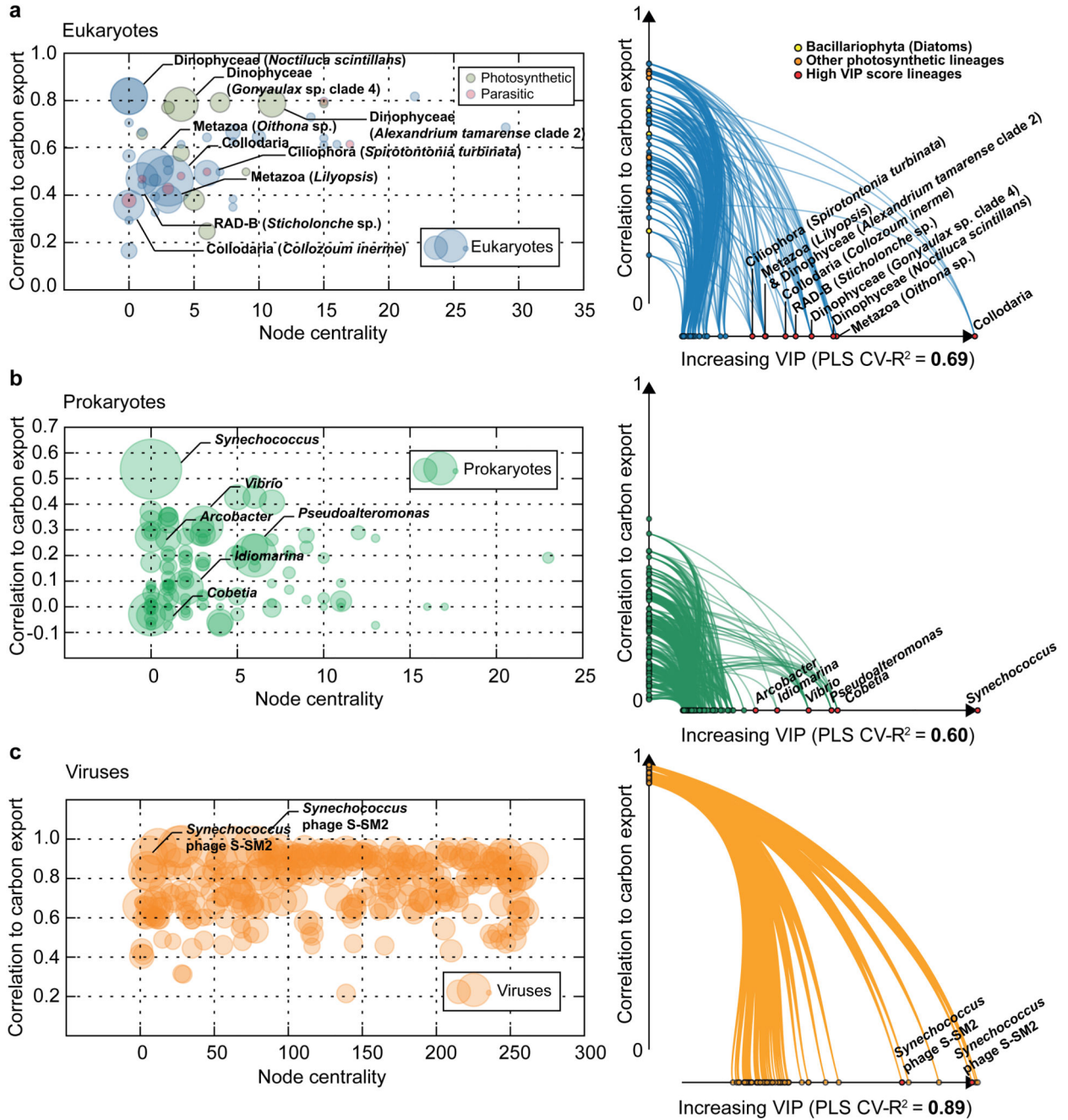
53. Behrenfeld MJ, Falkowski PG. Photosynthetic rates derived from satellite-based chlorophyll concentration. *Limnol. Oceanogr.* 1997; 42:1–20.
54. Chaffron, S., et al. Contextual environmental data of selected samples from the Tara Oceans Expedition (2009–2013). 2014. doi:10.1594/PANGAEA.840718
55. McCave IN. Size spectra and aggregation of suspended particles in the deep ocean. *Deep-Sea Res. I.* 1984; 31:329–352.
56. Sheldon RW, Prakash A, Sutcliffe WH. Size distribution of particles in ocean. *Limnol. Oceanogr.* 1972; 17:327–340.
57. Guidi L, et al. Relationship between particle size distribution and flux in the mesopelagic zone. *Deep-Sea Res. I.* 2008; 55:1364–1374. doi:10.1016/j.dsr.2008.05.014.
58. Logares R, et al. Metagenomic 16S rDNA Illumina tags are a powerful alternative to amplicon sequencing to explore diversity and structure of microbial communities. *Environ Microbiol.* 2014; 16:2659–2671. doi:Doi 10.1111/1462-2920.12250. [PubMed: 24102695]
59. Quast C, et al. The SILVA ribosomal RNA gene database project: improved data processing and web-based tools. *Nucleic Acids Res.* 2013; 41:D590–D596. doi:10.1093/Nar/Gks1219. [PubMed: 23193283]
60. Edgar RC. Search and clustering orders of magnitude faster than BLAST. *Bioinformatics.* 2010; 26:2460–2461. doi:10.1093/Bioinformatics/Btq461. [PubMed: 20709691]
61. Kultima JR, et al. MOCAT: A Metagenomics Assembly and Gene Prediction Toolkit. *PLoS ONE.* 2012; 7 ARTN e47656, doi:10.1371/journal.pone.0047656.
62. Hingamp P, et al. Exploring nucleo-cytoplasmic large DNA viruses in Tara Oceans microbial metagenomes. *ISME J.* 2013; 7:1678–1695. doi:10.1038/Ismej.2013.59. [PubMed: 23575371]
63. Zhao YA, Tang HX, Ye YZ. RAPSearch2: a fast and memory-efficient protein similarity search tool for next-generation sequencing data. *Bioinformatics.* 2012; 28:125–126. doi:10.1093/Bioinformatics/Btr595. [PubMed: 22039206]
64. Langmead B, Salzberg S. L. *Nat Methods.* 2012; 9:357–U354. doi:10.1038/Nmeth.1923. [PubMed: 22388286]
65. Shen HP, Huang JHZ. Sparse principal component analysis via regularized low rank matrix approximation. *J Multivariate Anal.* 2008; 99:1015–1034. doi:10.1016/J.Jmva.2007.06.007.
66. Langfelder P, Horvath S. Eigengene networks for studying the relationships between co-expression modules. *Bmc Syst Biol.* 2007; 1 Artn 54, doi:10.1186/1752-0509-1-54.
67. Li A, Horvath S. Network neighborhood analysis with the multi-node topological overlap measure. *Bioinformatics.* 2007; 23:222–231. doi:10.1093/Bioinformatics/Btl581. [PubMed: 17110366]
68. Chong IG, Jun CH. Performance of some variable selection methods when multicollinearity is present. *Chemometr. Intell. Lab.* 2005; 78:103–112. doi:10.1016/J.Chemolab.2004.12.011.
69. Mevik BH, Wehrens R. The pls package: Principal component and partial least squares regression in R. *J Stat Softw.* 2007; 18:1–23.



**Figure 1. Global view of carbon fluxes along the Tara Oceans circumnavigation route and associated eukaryotic lineages**

**a**, Carbon flux in  $\text{mg m}^{-2} \text{d}^{-1}$  and carbon export at 150 m estimated from particles size distribution and abundance measured with the Underwater Vision Profiler 5 (UVP5).

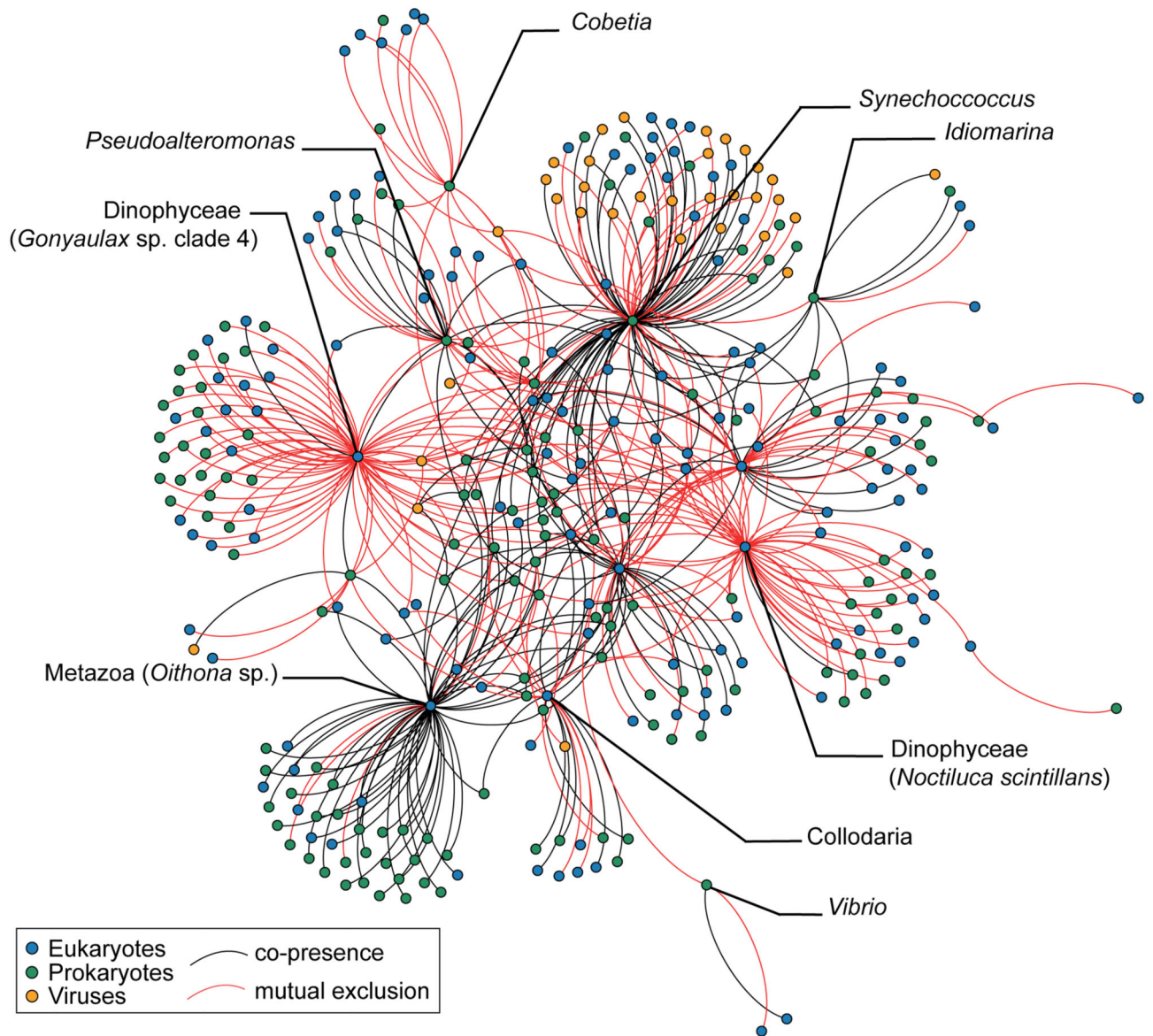
Stations at which environmental data are available (Supplementary Table 9) are depicted by white dots. Stations at which eukaryotic samples are available are colored in red (Supplementary Tables 10 and 12). **b**, Eukaryotic lineages associated to carbon export as revealed by standard methods for regression-based modeling (sPLS analysis). Correlations between lineages and environmental parameters are depicted as a clustered heatmap and lineages with a correlation to carbon export higher than 0.2 are highlighted (detailed results in Supplementary Table 1).



**Figure 2. Ecological networks reveal key lineages associated with carbon export at 150 m at global scale**

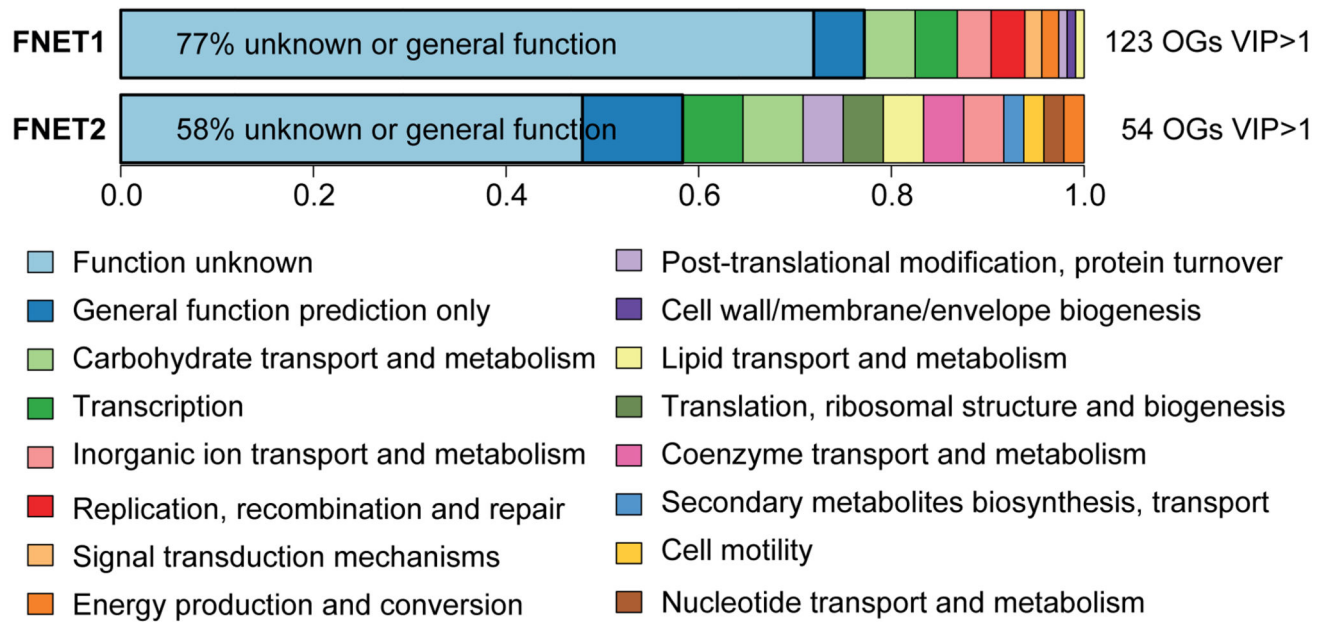
The relative abundances of taxa in selected subnetworks were used to estimate carbon export and to identify key lineages associated with the process. **a**, The selected eukaryotic subnetwork ( $n=49$ , see Supplementary Table 2) can predict carbon export with high accuracy (PLS regression, LOOCV,  $R^2=0.69$ , see Extended data Fig. 2g). Lineages with the highest VIP score (dots size is proportional to the VIP score in the scatter plot) in the PLS are depicted as red dots corresponding to three Rhizaria (Colloclaria, *Collozoum inerme* and

*Sticholonche* sp.), one copepod (*Oithona* sp.), one siphonophore (*Lilyopsis*), three Dinophyceae and one ciliate (*Spirotonionia turbinata*). **b**, The selected prokaryotic subnetwork ( $n=109$ , see Supplementary Table 3) can predict carbon export with good accuracy (PLS regression, LOOCV,  $R^2=0.60$ , see Extended data Fig. 2h). **c**, The selected viral population subnetwork ( $n=277$ , see Supplementary Table 4) can predict carbon export with high accuracy (PLS regression, LOOCV,  $R^2=0.89$ , see Extended data Fig. 2i). Two viral populations with a high VIP score (red dots) are predicted as *Synechococcus* phages (see Supplementary Table 4).



**Figure 3. Integrated plankton community network built from eukaryotic, prokaryotic and viral subnetworks related to carbon export at 150 m**

Major lineages were selected within the three subnetworks ( $VIP > 1$ ) (Supplementary Tables 2, 3 and 4). Co-occurrences between all lineages of interest were extracted, if present, from a previously established global co-occurrence network (see methods). Only lineages discussed within the study are pinpointed. The resulting graph is composed of 329 nodes, 467 edges, with a diameter of 7, and average weighted degree of 4.6.



**Figure 4. Key bacterial functional categories associated with carbon export at 150 m at global scale**

A bacterial functional network was built based on Orthologous Group/Gene (OG) relative abundances using the WGCNA methodology (see Methods) and correlated to classical oceanographic parameters. Two functional subnetworks (FNET1 ( $n=220$ ) and FNET2 ( $n=441$ ), respectively, Extended data Fig. 3a) are significantly associated with carbon export (FNET1: Pearson cor. 0.42,  $p = 4e^{-09}$  and FNET2: 0.54,  $p = 7e^{-06}$ , see Extended data Fig. 3b). Higher functional categories are depicted for functions with a VIP score  $>1$  (PLS regression, LOOCV, FNET1  $R^2=0.41$  and FNET2  $R^2=0.48$ , see Extended data Fig. 3d) in both subnetworks.