

Bioinformatics approaches for assessing microbial communities in the surface ocean

Kara Martin

Supervisors:

Prof.Vincent Moulton

Prof.Thomas Mock

Dr.Richard Leggett

A thesis submitted for the degree of Doctor of Philosophy

University of East Anglia,

School of Computing Sciences,

Norwich, United Kingdom.

September, 2018

©This copy of the thesis has been supplied on condition that anyone who consults it is understood to recognise that its copyright rests with the author and that use of any information derived therefrom must be in accordance with current UK Copyright Law. In addition, any quotation or extract must include full attribution.

Abstract

Microbes are vital for life on Earth. Within the oceans, they are the major primary producers of oxygen and contribute greatly to the other biogeochemical cycling of the elements which in turn influence the global climate. These microbes can be found inhabiting the oceans throughout the world and they cover over $\sim 70\%$ of the surface of the Earth.

Microbes have evolved in different environments in the oceans and in different ways. To gain an understanding of the microbial communities in the surface oceans in the Arctic and Atlantic oceans environmental scientists based at the University of East Anglia, the University of Groningen and Royal Netherlands Institute for Sea Research collected ocean samples from 68 stations along a transect of the Arctic Ocean, North Atlantic Ocean and South Atlantic Ocean. In addition, they recorded environmental data at the time of sampling, such as temperature and salinity. Genomic DNA from filtered samples was sequenced using high-throughput sequencing.

This thesis contains a comprehensive analysis of this sequencing data with the aim of understanding the composition and distribution of microbial communities in the surface of the ocean. To this end, we designed bioinformatic pipelines in order to analyse metatranscriptome, 18S and 16S rDNA datasets from the set of stations. In addition, we developed a novel methodology for normalising 18S and 16S rDNA copy numbers. This enabled us to perform additional analyses such as biodiversity, co-occurrence and breakpoint analyses. The breakpoint analysis is the first of this type performed for microbes in the ocean across a temperature gradient.

In our results, we observed a greater diversity of 18S and 16S rDNA taxa in the tropical regions of the South Atlantic Ocean, versus the polar regions of the Arctic Ocean. Moreover, in the co-occurrence analysis of the 18S and 16S rDNA datasets, we found two community networks, one positively correlated to temperature and the other

negatively. We also performed a breakpoint analysis on our metatranscriptome, 18S and 16S rDNA datasets and found a shift in diversity occurring in the North Atlantic Ocean. In particular, the shift occurs in the temperate region of the North Atlantic Ocean, between the polar Arctic Ocean and tropical South Atlantic Ocean.

These results are important because the co-occurrence analysis enables us to hypothesise that different microbial communities have different preferences for temperature. Moreover, as global warming is predicted to raise the temperatures in the ocean, our results could potentially enable forecasts of how climate change will affect these microbial communities using climate models underpinned by genetic information.

Acknowledgements

I would like to thank my supervisors Prof.Vincent Moulton, Prof.Thomas Mock and Dr.Richard Leggett for their continuous support and guidance. I would also like to thank Dr.Andrew Toseland for all of his support and advice. I would like to thank Dr.Katrin Schmidt, Dr.Willem van de Poll and Dr.Klaas Timmermans for collecting the samples during the three expeditions, and I would like to thank the Joint Genome Institute for sequencing the samples. Finally, I would like to thank the University of East Anglia and the Earlham Institute for funding my PhD.

Contents

Abstract	1
Acknowledgements	3
Contents	4
List of figures	16
List of tables	17
1 Introduction	18
1.1 Marine microbes	18
1.2 Thesis scope	22
1.3 Summary of thesis	24
2 Sequencing and sequence analysis	26
2.1 Summary	26
2.2 Next-generation sequencing technologies	26
2.3 Datasets	30
2.3.1 18S and 16S rDNA	30
2.3.2 Metagenomes and Metatranscriptomes	32
2.3.3 Dataset applications	33
2.4 Sequence preprocessing methods	34
2.4.1 Adapter searching and trimming	35
2.4.2 Read merging	37
2.4.3 Read clustering	38
2.5 Database searching	39

2.6	Phylogenetic analysis	41
2.6.1	Alignments	42
2.6.2	Tree building and placing	45
2.6.3	Visualisation	47
2.7	Discussion	49
3	18S rDNA and 16S rDNA analysis	50
3.1	Summary	50
3.2	Sampling and sequencing	51
3.2.1	Sampling	51
3.2.2	Sequencing and preprocessing of the 18S rDNA and 16S rDNA .	52
3.3	Methods	53
3.3.1	Computational pipeline for 18S rDNA analysis	53
3.3.2	Computational pipeline for 16S rDNA analysis	62
3.3.3	Further analysis methods	66
3.4	Results	69
3.4.1	Rarefaction curves	69
3.4.2	Principal Coordinates Analysis	71
3.4.3	Heatmaps	74
3.4.4	Evenness and occupancy	76
3.4.5	Environmental plots	78
3.4.6	Breakpoint analysis	88
3.4.7	Co-occurrence analysis	90
3.5	Discussion	94
4	Metatranscriptomics analysis	97
4.1	Summary	97
4.2	Sequencing and preprocessing	98
4.3	Methods	100
4.3.1	Computational pipeline for taxonomic classification analysis . .	100
4.3.2	Computational pipeline for functional analysis	101
4.3.3	Further analysis	102
4.4	Results	103

4.4.1	Taxonomic classification heatmap	103
4.4.2	Rarefaction curves	106
4.4.3	Principle Components Analysis (PCA)	107
4.4.4	Canonical Correspondence Analysis (CCA)	108
4.4.5	Breakpoint analysis	110
4.4.6	Co-occurrence analysis	111
4.5	Discussion	112
5	Discussion and future work	115
5.1	Summary	115
5.2	Future work	116
5.3	Conclusions	118
	Bibliography	121
	Appendices	146
A		147
A.A	Metadata	147
A.B	18S rDNA reference databases taxa IDs	148
A.C	Evenness and occupancy taxonomy names to numbers	152
A.D	18S and 16S rDNA correlation heatmap coefficient values and p-values	155
A.E	Co-occurrence node numbers to taxon and class membership	161
A.F	Co-occurrence module correlation heatmaps	166
B		168
B.A	PhymmBL four genomes locations online	168

List of Figures

1.1	Diatoms species, starting from the left is <i>Thalassiosira pseudonana</i> , <i>Emiliana huxleyi</i> , <i>Fragilariopsis cylindrus</i> , <i>Thalassiosira pseudonana</i> and <i>Fragilariopsis cylindrus</i> . (Photo from MOCK RESEARCH LAB, http://mocklab.com/)	18
1.2	The Argentine Sea northeast off the coast of the Falkland Islands. The spiralling pattern of greens and blues are phytoplankton growing on the surface of the Argentine Sea. The image was captured with the Moderate Resolution Imaging Spectroradiometer (MODIS) on NASA's Terra satellite on December 2, 2015. (Photo taken by Jeff Schmaltz, LANCE/EOSDIS Rapid Response, NASA's Earth Observatory)	21
1.3	Photos taken on the boat during an expedition to collect the Arctic Ocean samples. (a) photo of the conductivity, temperature, and depth (CTD) rosette sampler that was used to collect water. (b) two crew who are helping by positioning the CTD in the water. (c) Dr.Schmidt working in the lab preparing the samples	23
2.1	Workflow diagram of a typical computational analysis. We begin by sequencing the samples (pink) on next-generation sequencing (blue) platform which results in raw reads (turquoise). The raw reads (turquoise) undergo sequencing preprocessing (red) to prepare them for analysis. We conduct phylogenetic analysis (green) which includes database searching (yellow). Finally, we perform statistical analyses (purple)	30
2.2	18S rDNA dataset displayed on MEGAN's taxonomy tree at the taxonomic rank of class	49

3.1	Arctic and Atlantic Ocean sampling sites and measured metadata, (a) Sites of three expeditions, April to May 2011 in red, June to July 2012 in black and November to December 2012 in yellow. At each station, microbial communities were sampled at the deep chlorophyll maximum (DCM), corresponding to 68 samples altogether. (b) Isosurface plot of temperature (°C) measured at sampling depth. (c) Salinity (practical salinity unit(PSU)) measured at sampling depth for all stations. (d) Dissolved silicate (mol/L) concentrations measured at sampling depth for each station. (e) Concentration of dissolved phosphate (mol/L) measured at sampling depth for each station. (f) Nitrate and Nitrite (mol/L) concentrations measured at sampling depth for each station. (Figure was generated with Ocean Data view, R. Schlitzer, www.odv.awi.de, 2016)(Plot generated by Dr.Katrin Schmidt)	51
3.2	Pipeline diagram of Pplacer 18S rDNA classification analysis. The pipeline at various stages incorporates databases (blue), software tools (green), processed files (grey) and runtimes (yellow). Boxes a , b , c and d refer to sections in the text	54
3.3	The reference tree consists of 1636 species from the groups: Opisthokonta, Cryptophyta, Glaucocystophyceae, Rhizaria, Stramenopiles, Haptophyceae, Viridiplantae, Alveolata, Amoebozoa and Rhodophyta	55
3.4	The graph of 18S rDNA copy number and their related genome size (Mb) for 185 species across the eukaryote tree. We investigated 18S rDNA gene copy number and their related genome sizes. We observed a significant correlation of R^2 0.5480435268 with a p-value $< 2.2e-16$ between genome size and 18S rDNA copy number. Based on the log10 transformed data a regression equation was determined, $f(x)=0.66X+0.75$	59

3.5	Part of the 18S rDNA dataset displayed on MEGAN’s taxonomy tree. All the nodes highlighted in yellow at the leaves of the tree are class nodes. Nodes at the leaves of the tree that are not highlighted do not have a taxonomic classification of class in their lineage. The nodes between the leaves of the tree and the eukaryotic node are the internal nodes that are representing higher taxonomic levels such as phylum. The colour key represent each colour on the nodes and corresponds to a sample in the 18S rDNA dataset	61
3.6	JGI’s pipeline diagram of 16S rDNA classification analysis. The pipeline at various stages incorporates databases (blue), software tools (green) and processed files (grey). Box a and b refer to sections in the text . .	63
3.7	(a) rarefaction curves for 18S rDNA species level ($n=54$) and (b) the rarefaction curves for 16S rDNA genus level ($n=57$). In each panel the colours correspond to the sample sites of the three expeditions, April to May 2011 in red (North Atlantic Ocean), June to July 2012 in black (Arctic Ocean) and November to December 2012 in yellow (South Atlantic Ocean). The rarefaction curves were generated using MEGAN . .	70
3.8	Principal coordinates analysis (PCoA) of a represents the eukaryotic communities ($n=54$) and b , represents the prokaryotic communities ($n=57$) at the taxonomic rank of class. In MEGAN, communities are clustered according to their similarity based on Bray-Curtis distances. The samples are numbered and coloured by region; these numbers and colours correspond to the map of the region sampling sites in figure 3.1a, where the Arctic Ocean samples are coloured black, the North Atlantic Ocean samples are coloured red and the South Atlantic Ocean samples are coloured yellow	72

- 3.9 Panel **a** and **b** represents heatmaps of abundances arranged by latitude to the taxonomic rank of class, **a** represents the 18S rDNA taxonomy and **b** represents the 16S rDNA taxonomy. The taxonomy names in the dataset are displayed along the right side. The numbers at the bottom correspond to sample locations as shown in figure 3.1 **a**. The three regions of the Arctic Ocean, North Atlantic Ocean and South Atlantic Ocean are displayed underneath their corresponding sample numbers. The colours correspond to log₁₀-scaled abundances of the 18S and 16S rDNA, where red colours are high values and blue colours are low values. The heatmaps were generated using the heatmap.2 function, which is part of the gplots package, in R 75
- 3.10 Panels **a** and **b** represent abundance taxonomy evenness and occupancy plots for the 18S and 16S rDNA datasets, respectively. The numbers in the plots correspond to taxon names which can be found in the Appendix A.C. The *x*-axis represents the number of times that class taxonomy occurs across the stations. The *y*-axis represents the evenness of that class taxonomy across stations it occurs in. Each circle represents a class taxonomy abundance. The size of each circle corresponds to the total abundance for that class, calculated by taking the square root of the abundance divided by π 77
- 3.11 Correlation heatmaps of identified taxonomic classes and environmental data, **a** represents the Arctic 18S rDNA classes, **b** represents the North Atlantic 18S rDNA classes and **c** represents the South Atlantic 18S rDNA classes. The class names in the dataset are displayed along the right side and hierarchical clustering dendrogram on the opposite side. The environmental parameters are displayed at the bottom and hierarchical clustering dendrogram on the opposite side. The colours correspond to the Pearson correlation coefficient, where blue indicates a positive and red a negative correlation. The grey colour corresponds to no results, due to insufficient abundance data across the samples. The actual coefficients and p-values are given in Appendix A.D 80

- 3.12 Correlation heatmaps of identified taxonomic classes and environmental data, **a** represents the Arctic 16S rDNA classes, **b** represents the North Atlantic 16S rDNA classes and **c** represents the South Atlantic 16S rDNA classes. The class names in the dataset are displayed along the right side and hierarchical clustering dendrogram on the opposite side. The environmental parameters are displayed at the bottom and hierarchical clustering dendrogram on the opposite side. The colours correspond to the Pearson correlation coefficient, where blue indicates a positive and red a negative correlation. The grey colour corresponds to no results, due to insufficient abundance data across the samples. The actual coefficients and p-values are given in Appendix A.D 83
- 3.13 Panel **a**, NMDS of sampled stations with each number representing one sample of the 18S rDNA community. Significant environmental vectors for temperature (p=0.001), salinity (p=0.001), phosphate (p=0.001) and silicate (p=0.001) were fixed. Panel **b**, NMDS of sampled stations with each number representing one sample of 16S rDNA community. Significant environmental vectors for temperature (p=0.001), salinity (p=0.001), phosphate (p=0.001) and silicate (p=0.001) were fixed. NMDS was performed with the metaMDS function and the environmental variables were fitted with the envfit function (permutation test, 999 permutations). Both functions are part of the vegan package, in R. The numbers correspond to sample locations and the colours of the numbers correspond to ocean regions, black corresponds to the Arctic Ocean, red corresponds to the North Atlantic Ocean and yellow corresponds to the South Atlantic Ocean as shown in figure 3.1a 86

- 3.14 Panel **a**, a positive correlation of 18S rDNA diversity based on Shannon index with temperature. Based on backward model selection temperature was the only significant environmental covariate determined. Panel **b**, a positive correlation of 16S rDNA diversity, based on Shannon index with temperature. Based on backward model selection where temperature was the only significant environmental covariate determined. In panels **a** and **b**, the numbers correspond to sample locations and the colours of the numbers correspond to ocean regions, black corresponds to the Arctic Ocean, red corresponds to the North Atlantic Ocean and yellow corresponds to the South Atlantic Ocean as shown in figure 3.1a 88
- 3.15 Panels **a** and **b** represent breakpoint analysis to the taxonomic rank of class, **a** represents the 18S rDNA dataset and **b** represents the 16S rDNA dataset. The numbers correspond to sample locations as shown in figure 3.1a. The breakpoint analysis was generated using piecewise regression in R as detailed in [Castro-Insua et al., 2016]. The *y*-axis represents the beta diversity across the stations. The *x*-axis represents the temperature. In each plot, the horizontal line marks the breakpoint. For the 18S rDNA dataset in panel **a** the breakpoint is 13.96°C with a p-value of 8.407e-11. For the 16S rDNA dataset in panel **b** the breakpoint is 9.49°C with a p-value of 1.413e-4 89
- 3.16 In the co-occurrence analysis with WGCNA on the log10-scaled abundances of 18S rDNA species level and 16S rDNA genus level, two modules were found. Depicted is the correlation heatmap between the modules' eigengene and environmental variables. Along the left side, the two modules are displayed in turquoise (*n*=70) and blue (*n*=51). The environmental variables are displayed at the bottom. The colours correspond to the correlation values; red is positively correlated and blue is negatively correlated. The values in each of the squares correspond to the assigned Pearson correlation coefficient value on top and p-value in brackets below 91

- 3.17 In the co-occurrence analysis with WGCNA on the log₁₀-scaled abundances of 18S rDNA species level and 16S rDNA genus level, two modules were found. Depicted is the correlation heatmap for each of the two modules **a** (turquoise ($n=70$)) and **b** (blue ($n=51$)) of their species log₁₀-scaled abundances and environmental variables. Along the left side on each of the two modules **a** and **b** are displayed the species name and environmental variables are displayed at the bottom. The colours correspond to the correlation values; red is positively correlated and blue is negatively correlated. The Pearson correlation coefficient values and p-values can be found in Appendix A.F 92
- 3.18 Co-occurrence analysis with WGCNA on the log₁₀-scaled abundances of 18S rDNA species level and 16S rDNA genus level. In panels **a1** (turquoise ($n=70$)) and **b1** (blue ($n=51$)) are the two modules that we found depicted as network diagrams. These were generated with Cytoscape [Shannon et al., 2003]. The edge distance indicates the correlation strength between the nodes and the top five most highly connected nodes are coloured in orange. The numbers on the nodes correspond to the taxa names, the list of names to numbers of the taxa can be found in the Appendix A.E. In panels **a2** (turquoise ($n=70$)) and **b2** (blue ($n=51$)) the modules are depicted as word clouds. These consist of the member species and genus names, generated in WordCloud.com. The larger the name, the more connections that taxa has. In panel **a2**, taxa *Crocinitomix*, *Salinirepens*, *Bacillus*, *Bradyrhizobium*, *Rubritalea*, *Arcobacter*, *Delftia*, *Marinicella*, *Nonlabens* and *Psychroflexus* were removed as they are unreadable due to their frequency being too low to display with the others. In panels **a3** (turquoise ($n=70$)) and **b3** (blue ($n=51$)) are pie charts of the classes of the modules taxa. The list of names of the taxa to their class can be found in the Appendix A.E . . . 93
- 4.1 Diagram of JGI's computational pipeline for preprocessing the metatranscript reads. The pipeline at various stages incorporates databases (blue), BBTools tools (green) and processed files (grey) 98

4.2	A heatmap of the metatranscriptomic dataset taxonomically classified and arranged by latitude versus the taxonomic rank of phylum. The taxonomy names in the dataset are displayed along the right side. The numbers at the bottom correspond to sample locations as shown in figure 3.1 a . The three regions of the Arctic Ocean, North Atlantic Ocean and South Atlantic Ocean are displayed underneath their corresponding sample numbers. The colours correspond to log ₁₀ -scaled read counts, where purple colours are high values and blue colours are low values. The heatmap was generated using the heatmap.2 function, which is part of the gplots package, in R	104
4.3	Panel a and b represents Venn diagrams of the number of taxa entities in the 18S, 16S rDNA and metatranscriptome datasets at the taxonomic rank of phylum. Panel a is comparing the 18S rDNA dataset and metatranscriptome dataset. Panel b is comparing the 16S rDNA dataset and metatranscriptome dataset. (Figure was generated with http://bioinformatics.psb.ugent.be/webtools/Venn/)	105
4.4	Rarefaction curves for Pfam protein families. The numbers displayed in the plot correspond to sample location. The x-axis represents the random subsample size taken from the dataset and the y-axis indicates the number of unique Pfam protein families found. The R package VEGAN using the rarecurve function was employed to perform the rarefaction curves analysis	106
4.5	A Principle Components Analysis (PCA) for the Pfam protein families (log ₁₀ transformed) gene counts. The samples are numbered and coloured by region; these numbers and colours correspond to the map of the region sampling sites in figure 3.1 a , where the Arctic Ocean samples are coloured black, the North Atlantic Ocean samples are coloured red and the South Atlantic Ocean samples are coloured yellow. The R package VEGAN using the princomp function was employed to perform the PCA	107

- 4.6 A CCA for the Pfam protein families. The numbers in the plots correspond to sample locations as given in figure 3.1a. The numbers are coloured by region, red is for the North Atlantic Ocean, black is for the Arctic Ocean and yellow is for South Atlantic Ocean. We used the R package VEGAN to perform a Canonical Correspondence Analysis (CCA) between the Pfam dataset and the environmental data. The y -axis represents the CCA2 and the x -axis represents the CCA1. The arrows represent the direction and the length of the vector. Each vector represents an environmental factor variable 109
- 4.7 A breakpoint analysis for the Pfam protein families. The numbers in the plots correspond to sample locations as given in figure 3.1a. The breakpoint analysis was generated using piecewise regression in R as outlined in section 3.3.3. The y -axis represents the beta diversity across the stations. The x -axis represents the temperature. In the plot, the horizontal line marks the breakpoint. The Pfam protein families breakpoint is 18.06°C with a p-value of 1.24e-07 110
- 4.8 In the WGCNA analysis of the log10-scaled gene counts of Pfam protein families, thirteen modules (and a grey module) were found. In the figure, we present a correlation heatmap for the modules. The fourteen modules are displayed as coloured blocks labelled along the left hand side of the plot. The environmental parameters are displayed at the bottom. The colours correspond to the correlation values, red is positively correlated and blue is negatively correlated. The values in each of the squares correspond to the assigned Pearson correlation coefficient value on top and p-value in brackets below 112

- A.1 In the co-occurrence analysis with WGCNA on the log10-scaled abundances of 18S rDNA species level and 16S rDNA genus level, two modules were found. Depicted is the correlation heatmap for the turquoise module in figure 3.17 **a** ($n=70$). Along the left hand side is the species/genus name and environmental variables are displayed at the bottom. The colours correspond to the correlation values, red is positively correlated and blue is negatively correlated. The Pearson correlation coefficient values and p-values in brackets are displayed in each square 166
- A.2 In the co-occurrence analysis with WGCNA on the log10-scaled abundances of 18S rDNA species level and 16S rDNA genus level, two modules were found. Depicted is the correlation heatmap for the blue module in figure 3.17 **b** ($n=51$). Along the left hand side is the species/genus name and environmental variables are displayed at the bottom. The colours correspond to the correlation values, red is positively correlated and blue is negatively correlated. The Pearson correlation coefficient values and p-values in brackets are displayed in each square 167

List of Tables

A.1	Evenness and occupancy taxonomy names to numbers	152
A.2	Arctic 18S rDNA correlation heatmap coefficients values and p-values .	155
A.3	North Atlantic 18S rDNA correlation heatmap coefficients values and p-values	156
A.4	South Atlantic 18S rDNA correlation heatmap coefficients values and p-values	157
A.5	Arctic 16S rDNA correlation heatmap coefficients values and p-values .	158
A.6	North Atlantic 16S rDNA correlation heatmap coefficients values and p-values	159
A.7	South Atlantic 16S rDNA correlation heatmap coefficients values and p-values	160
A.8	Co-occurrence module ($n=70$) node numbers to taxon and class mem- bership	161
A.9	Co-occurrence module ($n=51$) node numbers to taxon and class mem- bership	163

Chapter 1

Introduction

1.1 Marine microbes

The surface of the planet is made up of about 70% water, therefore providing the largest habitat for life on the surface of the planet, in particular, microbes [Das et al., 2006]. Among these marine microbes inhabiting this vast marine ecosystem are prokaryote species such as alphaproteobacteria and eukaryote species such as diatoms [Aryal et al., 2015], [Brierley, 2017]. Phytoplankton is a class of unicellular photosynthetic organisms, that is composed of a wide range of organisms of both prokaryotes and eukaryotes, and include thousands of different species [Collins et al., 2014], [Brierley, 2017]. Diatoms are eukaryotic phytoplankton and are the most diverse of the eukaryotic phytoplankton as they have about 200,000 different species [Armbrust, 2009]. (Figure 1.1).

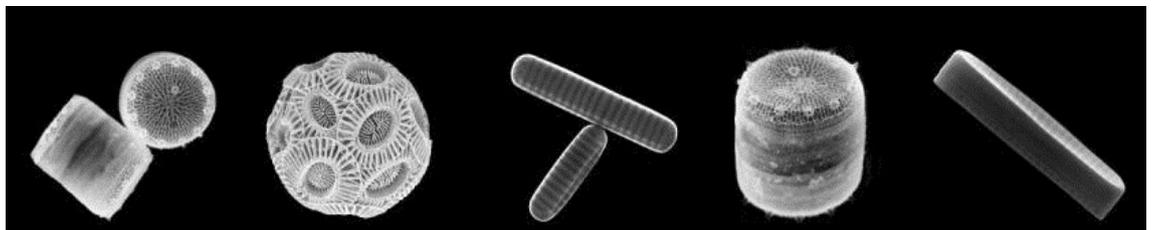


Figure 1.1: Diatoms species, starting from the left is *Thalassiosira pseudonana*, *Emiliana huxleyi*, *Fragilariopsis cylindrus*, *Thalassiosira pseudonana* and *Fragilariopsis cylindrus*. (Photo from MOCK RESEARCH LAB, <http://mocklab.com/>)

Communities can be defined as a large number of species that interact with each other in a countless number of ways [Godfray and May, 2014]. Communities of marine prokaryote and eukaryote species have co-existed throughout their evolution in

common habitats [Cordero and Datta, 2016], [Brussaard et al., 2016]. In the environment, two principal forces can be seen to act on these microbes - biotic, which refers to the interactions between the microbes, and abiotic, which refers to environmental influences such as temperature and salinity [Bijlsma and Loeschke, 2005]. The interaction within marine microbial communities occur under these biotic and abiotic influences and therefore shape their evolution and adaptation [Brussaard et al., 2016]. The microbes have a range of interactions such as endosymbiosis (which refers to the merging of distinct cells so one is inside the other), mutualism (which refers to the interaction between different species that results in a mutually beneficial outcome, be it for reproduction and/or survival), parasitism (which refers to a species that obtains sustenance from its host) and commensalism (which refers to when one species gains from its association with another species who is not affected in a positive or negative manner) [Baron, 1996], [Holland and Bronstein, 2008], [Boon et al., 2014], [Haque and Haque, 2017].

A notable example, as outlined next, is how microbial interactions resulted in phytoplankton acquiring the ability to convert light into energy in order to thrive and the subsequent diversification of phytoplankton species [Wernegreen, 2012]. Phytoplankton obtained their photosynthesis capabilities through the process of endosymbiosis from a class of bacteria called cyanobacteria [Simon et al., 2009]. About 1.5 billion years ago an endosymbiosis event is hypothesized to have produced the first photosynthetic eukaryote, a process by which a heterotrophic eukaryote engulfed a cyanobacterium [Shemi et al., 2015]. An intracellular gene transfer occurred between the primitive host of heterotrophic eukaryote and its symbiont cyanobacterium, after which the heterotrophic eukaryote retained the cyanobacterium by converting it into a plastid [Simon et al., 2009]. This primary endosymbiosis event resulted in 3 distinct clades of unicellular algae. One of these clades is the viridiplantae, also known as the green plastid lineage, which is thought to be the source from which all land plants and green algae evolved. Another is the rhodophyta, also known as the red plastid lineage. This clade includes an ancient group of marine red microalgae and seaweed. The other clade is the glaucophyta, which is made up of a small group of freshwater algae. About 1 billion years ago, a second endosymbiosis event is hypothesized to have occurred, of a second heterotroph engulfing and retaining a member of one of these clades. During

this secondary endosymbiosis event, the clades rhodophyta and chlorophytes gave rise to the dominant lineages of algae, such as stramenopiles, dinoflagellates, cryptophytes, and haptophytes. In present-day oceans, the most abundant, diverse, and ecologically important microbes are coccolithophores, diatoms, and dinoflagellates [Shemi et al., 2015].

These microbial community interactions that support their evolution and adaptation, are also the engines that drive the biogeochemical cycle of elements that occur in the oceans, such as the nitrogen and carbon cycles [Brussaard et al., 2016], [Falkowski et al., 2008]. Nitrogen is vital for life on earth and while nitrogen (N_2) is abundant in the atmosphere, this form is inaccessible for use, therefore a conversion into ammonia (NH_3) is required so it may be utilised and this is achieved through a process called nitrogen fixation [Wernegreen, 2012], [Voss et al., 2013]. There are certain bacteria and archaeal groups able to perform nitrogen fixation, and also some eukaryotic lineages have this ability as they required it through endosymbiosis in order to live in nitrogen-poor habitats. For example, microbial communities of cyanobacteria that live within the eukaryote phytoplankton species called diatoms perform this nitrogen fixation by their endosymbiosis mutualism interactions [Foster et al., 2011], [Wernegreen, 2012].

The world's surface ocean is divided into latitudinal temperature zones, ranging from about $30^\circ C$ in the tropics to about $-1.8^\circ C$ in the ocean's polar sea and ice interface, and in the sea ice, the temperature can even be well below $-1.8^\circ C$ [Toseland et al., 2013]. Phytoplankton species can be found inhabiting all these temperature zones and due to their photosynthetic capabilities can be found inhabiting the surface layers of the ocean [Simon et al., 2009], [Aryal et al., 2015]. In the ocean, mixing occurs mainly during storms by wind and waves, and stratification of the water column occurs as the surface waters warm and storm-driven mixing decreases. Stratification layers consist of warm, low nutrient and illuminated water over deeper, darker, and cooler high nutrient water. These layers are divided by rapid changes over depth of water density, called the pycnocline, and temperature called the thermocline [Brierley, 2017]. In the ocean, the supply of nutrients is dependent on temperature-driven stratification and mixing of the ocean [Toseland et al., 2013]. The growth and diversity of phytoplankton is dependent on their optimal growth temperature and on the supply of nutrients [Toseland et al., 2013]. Rapid phytoplankton growth episodes known as blooms are so big that they

can even be seen from space as shown in figure 1.2 [Brierley, 2017].

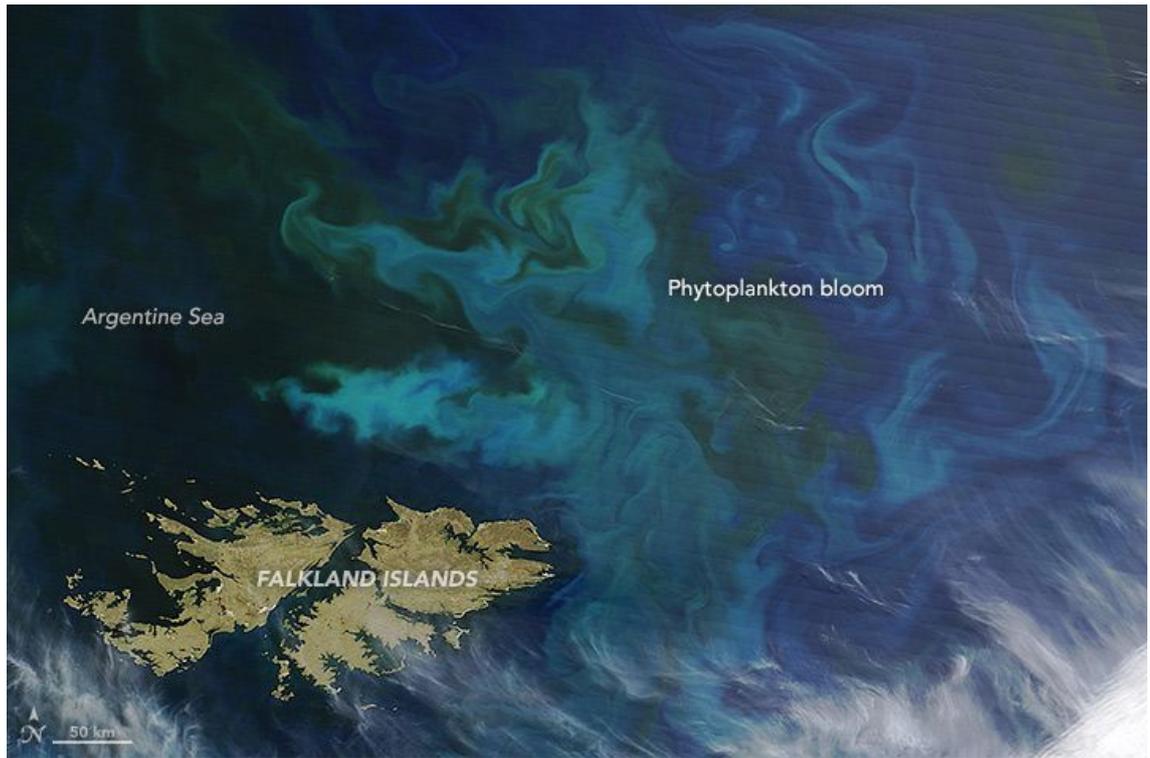


Figure 1.2: The Argentine Sea northeast off the coast of the Falkland Islands. The spiralling pattern of greens and blues are phytoplankton growing on the surface of the Argentine Sea. The image was captured with the Moderate Resolution Imaging Spectroradiometer (MODIS) on NASA's Terra satellite on December 2, 2015. (Photo taken by Jeff Schmaltz, LANCE/EOSDIS Rapid Response, NASA's Earth Observatory)

In recent work [Collins et al., 2014], the effect of recent human activity has been observed. In particular, the amount of carbon dioxide and bicarbonate in the world's oceans is increasing because of anthropogenic (pollution from human activity) carbon dioxide and the marine ecosystems will likely be affected by this increase in the oceans and atmosphere. This is causing ocean acidification, which is a decrease in the oceans pH. The mean pH of the surface oceans has decreased by about 0.1 pH units since the industrial revolution and is likely to decrease by an extra 0.3 pH units by the end of the 21st century, therefore this will result in an increase in acidity of about 150%. It is likely that the biogeochemical cycle of elements will be affected by ocean acidification, as it can change the community composition and can push for physiological and evolutionary change. Additionally, since the industrial revolution, the average ocean surface water temperature has already increased by 0.7°C and is likely to increase by an extra 3°C by the end of the 21st century. An increase of stratification of the oceans surface water is caused by an increase in temperature, which affects the light regime and also decreases

the amounts of nutrient supplied from below [Collins et al., 2014].

Phytoplankton divide their cells asexually at a rapid rate in an order of hours to days and have enormous population sizes. It is these features that enable phytoplankton to evolve in response to changing environmental conditions on time scales of weeks, months or years. For many years phytoplankton such as coccolithophores, diatoms, and dinoflagellate have been used as environmental indicators, as they are abundant in the oceans and have a substantial fossil record. Each year marine phytoplankton are responsible for about 50% of the carbon dioxide that is fixed in the atmosphere [Toseland et al., 2013]. Also, phytoplankton contribute to the base of the marine food web [Sarmiento et al., 2010]. Therefore it is crucial to know how these marine microbes will respond to changing environmental conditions as this will affect the marine food web and biogeochemical cycles [Ribeiro et al., 2013], [Collins et al., 2014]. Marine microbes are extremely difficult to isolate from the environment and one possible reason is that during laboratory culturing, community interactions which are important for growth may be destroyed [Joint et al., 2010], [Kazamia et al., 2016]. We have known for some time that standard laboratory culturing techniques can only isolate a very small proportion of microbes [Joint et al., 2010]. Genomic data analysis of microbes has resulted in many insights being discovered such as how they evolved to be the biogeochemical engineers of life [Falkowski et al., 2008].

1.2 Thesis scope

We use multiple sequencing data types and develop bioinformatic approaches to assess microbial communities in the surface ocean. This thesis is a collaboration between the University of East Anglia (UEA), Earlham Institute (EI) and Joint Genome Institute (JGI). Dr.Katrin Schmidt, an environmental scientist who was based at UEA, sampled stations from a transect of the Arctic Ocean and the South Atlantic Ocean. In figure 1.3 are photos taken while Dr.Schmidt was on the expedition to collect the samples. Dr.Willem van de Poll of the University of Groningen, Netherlands and Dr.Klaas Timmermans of the Royal Netherlands Institute for sea research sampled stations from the North Atlantic Ocean.

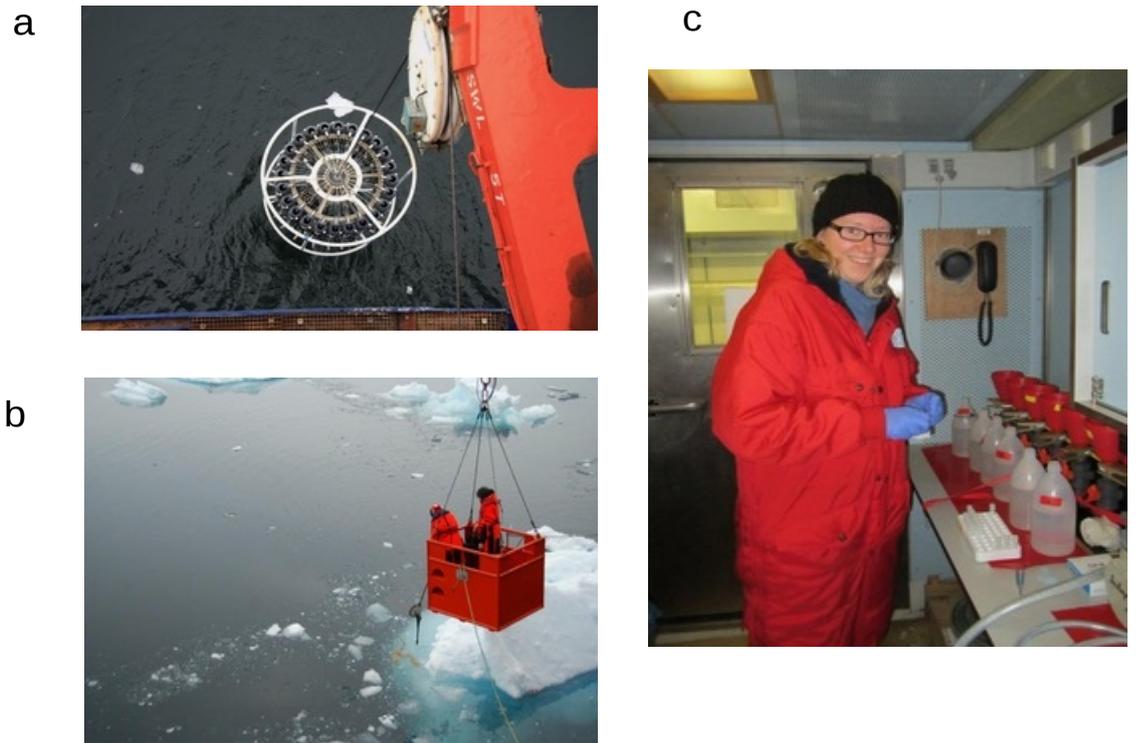


Figure 1.3: Photos taken on the boat during an expedition to collect the Arctic Ocean samples. (a) photo of the conductivity, temperature, and depth (CTD) rosette sampler that was used to collect water. (b) two crew who are helping by positioning the CTD in the water. (c) Dr.Schmidt working in the lab preparing the samples

The samples were sequenced with Illumina technology for metatranscriptome, metagenome, 18S and 16S rDNA analysis. This provided us with large amounts of sequence data, which enabled analysis across a wide range of environmental conditions. This could potentially enable us to forecast how climate change will affect these phytoplankton communities. The main challenge of this project was to develop a strategy to analyse the huge amounts of sequence data and to develop computational tools when necessary to extract pertinent information from the data. To address this challenge, we developed pipelines to analyse 18S rDNA data, 16S rDNA data and metatranscriptome data.

There have been other studies in the areas of phylogenetics and metatranscriptomics of marine samples, for example, [Alemzadeh et al., 2014] and [Alexander et al., 2015]. However, these focus on single locations such as the Persian Gulf and Narragansett Bay, United States of America, respectively. Most recent and most notable is the Tara Oceans project [Bork et al., 2015]. This paper describes data taken from sample locations from India around South Africa and over to South America into the Pacific

Ocean [Bork et al., 2015]. These are tropical and temperate regions in the oceans and sampled mainly in the open ocean. The novelty of our work is that we sampled close to the coast for a transect of the Arctic Ocean down through the North Atlantic Ocean and to Cape Town in the South Atlantic Ocean as shown in figure 3.1a. No other study has done this before. Additionally, we developed novel analysis pipelines such as the 18S rDNA taxonomic classification pipeline as well as our development of new approaches to normalising 18S and 16S rDNA copy numbers in our 18S and 16S rDNA datasets. In addition, environmental data at the time of sample collection, such as temperature and salinity, was recorded, which we included in our analysis, including co-occurrence analysis to identify community networks and breakpoint analysis which have not been used in this context before.

1.3 Summary of thesis

In chapter 2 we will summarise bioinformatics tools and methods that we use for the analysis of our 18S and 16S rDNA datasets. We will discuss these tools under their respective headings of sequencing processing, database searching and phylogenetic analysis, and in terms of what these are and their importance. Also, we will give a brief background on next-generation sequencing technology.

In chapter 3 we will discuss the computational pipeline and the analysis for the 18S and 16S rDNA datasets. We describe the sampling and preparation that was performed by Dr.Katrin Schmidt, Dr.Klaas Timmermans and Dr.Willem van de Poll. We also describe the sequencing and preprocessing of the 18S rDNA and 16S rDNA datasets, and the computational pipeline of the 16S rDNA dataset that was performed by JGI. Our contribution was the implementation of the computational pipeline for the 18S rDNA dataset, and the analysis of 18S and 16S rDNA datasets, as well as additional methods for the normalisation of the 18S and 16S rDNA copy number.

In chapter 4 we will discuss the computational pipeline and the analysis for the metatranscriptomic datasets. We also describe the sequencing, preprocessing and the computational pipeline of the metatranscriptomic dataset that was performed by JGI. Our contribution was the implementation of the computational pipeline for the analysis of the metatranscriptomic dataset.

In chapter 5 we will discuss our conclusions on our analyses of the metatranscriptomic, 18S and 16S rDNA datasets. We also discuss future work, such as additional analyses to be performed on our metatranscriptomic dataset, and the analysis of associated metagenomic datasets.

Chapter 2

Sequencing and sequence analysis

2.1 Summary

In this chapter, we shall summarise a number of the bioinformatics tools and methods that have been used for the analysis of our 18S and 16S rDNA datasets. We also give some background on 18S and 16S rDNA, metatranscriptomics and next-generation sequencing technology. We discuss sequence processing, database searching and phylogenetic analysis, describing what they are, why they are important and the tools that are involved in these approaches that are applied in this thesis.

2.2 Next-generation sequencing technologies

All organisms on Earth can be defined by their genome. The genome contains the biological blueprints for the construction and maintenance of that organism. For cellular organisms, the genome is made up of deoxyribonucleic acid (DNA) but the genome of a few viruses consist of ribonucleic acid (RNA). DNA and RNA are polymeric molecules that consist of a chain of monomeric subunits called nucleotides. There are five chemically distinct nucleotides, for DNA, these are Adenine, Thymine, Cytosine and Guanine and in the case of RNA, these are Adenine, Cytosine, Guanine and Uracil instead of Thymine [Brown, 2002]. Segments of DNA that encode for proteins with function or phenotype are called genes [Wain et al., 2002]. These protein coding genes are expressed by the process of transcription into messenger RNA (mRNA) transcripts and these, in turn, are translated into proteins. These proteins specify the type of the

biochemical processes that the cell is able to carry out [Brown, 2002].

Sequencing is a method for determining the DNA or RNA sequences within organisms [Sanger et al., 1977], [Wang et al., 2009]. Sequencing technology has revolutionised biological research, making possible the sequencing of entire genomes or a particular area of interest within a genome and have facilitated the further development of fields such as phylogenetics [Behjati and Tarpey, 2013], [van Dijk et al., 2014]. Additionally, sequencing technology has enabled the study of the transcriptome, which is the full set of mRNA transcripts in a cell. This consists of quantifying the types and amounts of transcripts in a cell under varying conditions [Wang et al., 2009].

In the 1970's, Sanger et al. developed the first generation DNA sequencing technology known as chain termination [Sanger et al., 1977]; and Maxam and Gilbert developed an alternative method called the chemical degradation method [Maxam and Gilbert, 1977] [Brown, 2002], [van Dijk et al., 2014]. Due to the radioisotopes and level of toxic chemicals involved in Maxam and Gilbert's chemical degradation method, the chain termination method became the predominant DNA sequencing method for the next 30 years [van Dijk et al., 2014]. Since the 1990's, Sanger sequencing biochemistry has mostly been carried out by a capillary based method and is semi-automated [Shendure and Ji, 2008].

Within Sanger sequencing high throughput pipelines, DNA can be prepared by a method called targeted resequencing, which consist of Polymerase chain reaction (PCR) amplification with primers that flank the target DNA. This results in many PCR amplicons within a single reaction volume. Cycles of template denaturation, primer annealing and primer extension are performed during the cycle sequencing reaction. The primer sequence is complementary and known to the region flanking the sequence of interest. Each round of primer extension is randomly terminated when fluorescently labelled dideoxynucleotides (ddNTPs) are incorporated. This results in a mixture of end labeled extension products. The product's label on the terminating ddNTP each corresponds to the nucleotide identity of its terminal position. The sequences of the single stranded products in a capillary based polymer gel are determined using high resolution electrophoresis. As they exit the capillary the fluorescent labels are excited by a laser. Also attached is a colour detection of emission spectrum in order to provide an output in the form of a plot called a Sanger sequencing trace. The Sanger

sequencing trace is then converted into DNA sequence with error probabilities for each base call [Shendure and Ji, 2008].

Though accurate, Sanger sequencing is low yielding and expensive [Reuter et al., 2015]. By comparison, the second generation of DNA sequencing technologies often referred to as next-generation sequencing (or NGS for short) can perform reactions in parallel producing much more data at a lower cost [van Dijk et al., 2014]. Our samples were sequenced with Illumina NGS [Bennett, 2004] by Joint Genome Institute (JGI), for metatranscriptome, metagenome, 18S and 16S rDNA sequences. Illumina HiSeq2000, for example, can produce 150-200 Gb (gigabase) of data per run and at a cost of \$0.02 per million bases [Pillai et al., 2017].

The Illumina methodology differs from Sanger sequencing in that it uses the technology of sequencing by synthesis (SBS) [Liu et al., 2012]. For Illumina, sequencing occurs within a flow cell which contains one, two or eight separate lanes [Buermans and den Dunnen, 2014]. Adapters which are essentially sequencing primer annealing sequences are ligated to each of the DNA fragments [Kozarewa et al., 2009]. A library of DNA molecules with attached sequencing adapters is denatured to produce single stands. These are passed through a flowcell and attached to the complementary oligonucleotides that are spread over the flowcell. This is followed by a procedure called bridge amplification, which is a solid phase PCR which forms clusters of clonal DNA fragments [Liu et al., 2012], [Heather and Chain, 2016]. SBS involves additions of fluorescent reversible-terminator dNTPs, this results in no more nucleotides being able to bind to the DNA molecule. Before polymerisation can advance, these must be cleaved off thus allowing the sequencing to occur at the same time throughout. In cycles, the altered dNTPs and DNA polymerase are passed through the flowcell containing the primed single stranded clusters. For each cycle, the incorporating nucleotide is identified with a charge coupled device (CCD) by exciting the fluorophores with suitable lasers, before enzymatic removal of the blocking fluorescent component and then continues to the next position [Heather and Chain, 2016].

While Illumina's second generation technology still dominates the market, two third generation technologies from Pacific Biosciences and Oxford Nanopore Technologies have begun to gain traction. These technologies are both single molecule sequencers, requiring no artificial amplification and are characterised by long reads (in the thou-

sands or tens of thousands of bases) and a higher error rate than second generation technologies [Liu et al., 2012], [van Dijk et al., 2014]. Nanopore sequencing technology was first developed by David Deamer (University of California Santa Cruz), George Church and Daniel Branton (Harvard University) [Jain et al., 2016]. A number of companies have developed nanopore based sequencing technologies, but Oxford Nanopore Technologies (ONT) is the only one thus far to bring a product to market (in 2014) with the release of MinION.

The MinION is the smallest portable sequencing device to date, weighing 90g and measuring 10 x 3 x 2 cm and powered from a standard USB3 port [Jain et al., 2016], [Lu et al., 2016]. There is no charge for the device itself, with labs paying only for consumables. It can output read lengths of tens of kilobases limited only by the length of DNA molecules inserted into it and with a single read accuracy of 95% [Laver et al., 2015], [Jain et al., 2016], [Carter and Hussain, 2017]. The MinION contains a flowcell with 2048 nanopores arranged in groups of 4 under 512 current sensors. Before sequencing, adapters are ligated to both ends of the DNA or cDNA fragments which enables the strands to be captured and loading of the processive enzyme (motor protein) at the 5' -end of each strand, thus ensuring unidirectional single nucleotide displacement along the strand. The adapters increase the DNA capture rate by several thousand fold by concentrating the DNA substrates at the membrane surface near to the nanopore. When the DNA molecule is captured in the nanopore, the motor protein advances the template strand through the nanopore. Once the enzyme passes through the hairpin, this is repeated for the complementary strand. Each sensor monitors the changes in ionic current as the DNA moves through the pore. The changes in the ionic current are divided into distinct events to which a duration, mean amplitude, and variance can be associated. Using probabilistic models this is interpreted computationally as a sequence of 3 to 6 nucleotide kmers [Jain et al., 2016].

In figure 2.1 we give a high-level overview of how a typical analysis of data proceeds in this thesis. After samples have been sequenced as outlined here in section 2.2, sequence preprocessing takes place, which we outline in section 2.3. Once the reads have been prepared the next step is to identify the reads by phylogenetic analysis, which includes searching databases. We outline database searching in section 2.4 and phylogenetic analysis in section 2.5.

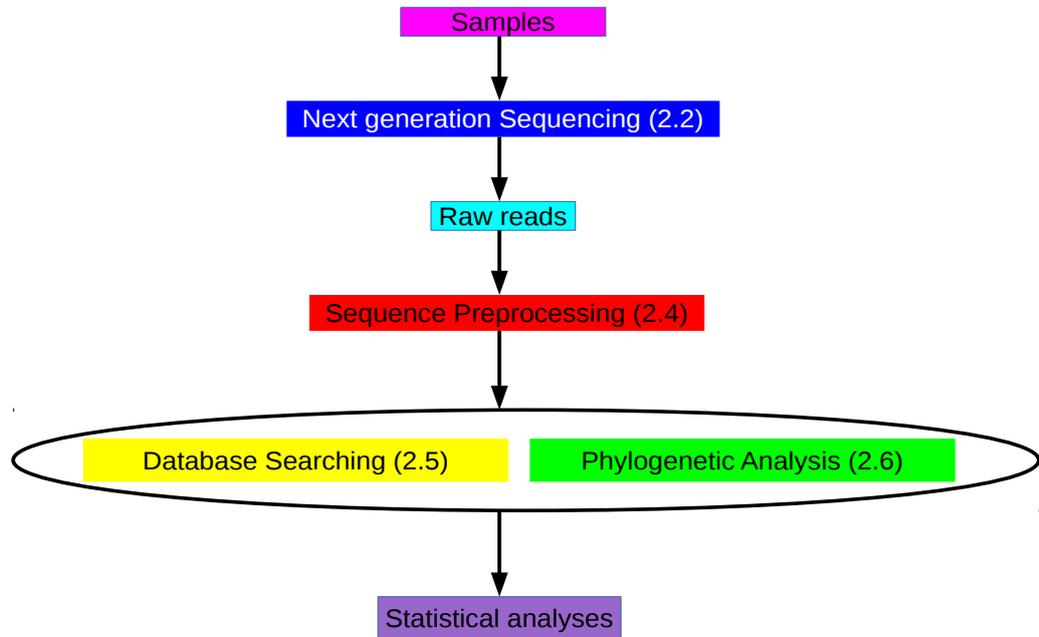


Figure 2.1: Workflow diagram of a typical computational analysis. We begin by sequencing the samples (pink) on next-generation sequencing (blue) platform which results in raw reads (turquoise). The raw reads (turquoise) undergo sequencing preprocessing (red) to prepare them for analysis. We conduct phylogenetic analysis (green) which includes database searching (yellow). Finally, we perform statistical analyses (purple)

2.3 Datasets

2.3.1 18S rDNA and 16S rDNA

Within eukaryotic and prokaryotic cells, the sites of protein synthesis are the ribosomes. In prokaryote cells, the ribosomes are called 70 Svedberg units (S) and in eukaryote cells, the ribosomes which are slightly larger than the prokaryote ribosomes are called 80S. Overall the structure of the ribosomes for eukaryotes and prokaryotes cells are similar. In eukaryotes and prokaryotes, the ribosome is composed of a large subunit and a small subunit. For prokaryote ribosome 70S, the smaller subunit referred to as 30S consists of 16S ribosomal RNA (rRNA) and 21 proteins. The larger subunit referred to as 50S consists of 23S and 5S rRNAs and 34 proteins. For eukaryotic ribosome 80S, the smaller subunit referred to as 40S consists of 18S rRNA and 30 proteins. The larger subunit referred to as 60S consists of 28S, 5.8S, and 5S rRNAs and about 45 proteins [Cooper, 2000]. Eukaryotic and prokaryotic cells typically contain many ribosomes and within each small ribosomal subunit a single RNA species exists of 18S

rRNA and 16S rRNA, respectively [Cooper, 2000], [Amit Roy et al., 2014].

The ribosomal DNA (rDNA) sequences encode for rRNAs, and these are tandemly repeated [Gibbons et al., 2014]. The rDNA copy number can vary greatly among prokaryotes and eukaryotes. For prokaryotes, the 16S rDNA copy number can vary as much as from one to fifteen and for eukaryotes, the 18S rDNA copy number can vary even greater from one to thousands [de Vargas et al., 2015], [Perisin et al., 2016]. The sequence that is composed of repeated rRNA gene clusters separated by intergenic spacers (IGS) is referred to as the ribosomal gene array. For eukaryotes, each cluster is made up of sequences of rRNA encoding highly conserved genes for 18S rRNA, 5.8S rRNA and 28S rRNA and these are separated by internal transcribed spacers which are referred to as ITS1 and ITS2, and also external transcribed spacers called 5ETS and 3ETS, which are located downstream of the 18S rRNA gene and upstream of the 28S rRNA gene respectively [Dyomin et al., 2016], [Fernández-Pérez et al., 2018]. The 5S rDNA sequence which codes for 5S rRNA is clustered in tandem arrays and is separated by flanking DNA sequences called non-transcribed spacers (NTSs) [Fernández-Pérez et al., 2018]. All these components are transcribed into a single RNA precursor, pre-rRNA [Dyomin et al., 2016]. The ETS, ITS and NTSs sequences are lost during maturation [Mandal, 1984], [Fernández-Pérez et al., 2018]. For prokaryotes the rDNA encoding the rRNA are traditionally arranged in a single operon in the order of 16S rDNA, 23S rDNA and 5S rDNA and these are separated by ITSs. All the components are transcribed into a single RNA precursor, which is then separated and processed by RNases [Brewer et al., 2019].

For taxonomic classification, the small subunit rRNA gene is a standard reference sequence, the 18S rDNA for eukaryotes and 16S rDNA for prokaryotes [Wang et al., 2014]. The 16S and 18S rDNA molecular markers are the most popular and ideal for phylogenetic studies [Fu and Gong, 2017]. Features that make rDNA so ideal for phylogenetic studies, for example, are that rDNA contains highly conserved and variable domains, it evolves slower than protein coding genes and is universal among organisms [Johnston, 2006], [Amit Roy et al., 2014]. The 18S and 16S rDNA are highly conserved within eukaryotes and prokaryotes, respectively [Amit Roy et al., 2014], [Wu et al., 2015]. The 18S rDNA is about 1800bp long and 16S rDNA is about 1550 bp long. Each of the 18S and 16S rDNA sequences contains both variable

and conserved regions with characteristic oligonucleotide signature sequences that are distinctive to a specific phylogenetic group [Iwen and Hinrichs, 2002], [Amit Roy et al., 2014]. SILVA [Quast et al., 2013a] is an online web resource for databases of ribosomal RNA (rRNA) gene sequences. SILVA contain sequences from the Bacteria, Archaea and Eukaryota domains [Quast et al., 2013a]. The SILVA database taxonomic rank assignments are manually curated [Balvočit and Huson, 2017]. SILVA is a popular choice for the thousands of researchers around the world who included SILVA in their work, such as [Lambrechts et al., 2019], in which samples were taken from the Antarctic to study the composition of prokaryotic communities [Quast et al., 2013a].

2.3.2 Metagenomes and Metatranscriptomes

Metagenomics, the direct genetic analysis of genomes within an environmental DNA sample, enables the novel exploration of functional gene composition and diversity of these microbial communities [Thomas et al., 2012b]. The same holds for metatranscriptomics which targets the mRNA within an environmental sample, thus enabling the investigation of taxonomic composition and the activity of biochemical functions of microbial community [Gilbert and Hughes, 2011], [Jiang et al., 2016a]. There are three main aims for the analysis of metatranscriptomic datasets: who is there, what are they doing and how do different samples compare? [Huson et al., 2011]. Metatranscriptomic analysis gives a new insight into microbial communities, and in particular how these microbial communities react to changes in the environment [Gilbert and Hughes, 2011], [Narayanasamy et al., 2016]. There have been a number of advances in microbial ecology, evolution, and diversity over the past decade, therefore a substantial number of research laboratories are actively engaged in the field [Thomas et al., 2012b]. Even so, there are a number of limitations, such as the small number of reference genomes, pipelines and computational tools, thus making it challenging to analyse and interpret metatranscriptomics datasets [Jiang et al., 2016b].

Metatranscriptomic approaches are being applied to investigate microbial communities in numerous different habitats including, for example, marine, soil and the human gut [Yu and Zhang, 2013]. Of most relevance to this thesis, the Tara Oceans project has sampled a large number of stations around the world's open oceans, upon which they

have published a number of papers including their metatranscriptomic analysis [Caradec et al., 2018], and in which they have characterised the changing genetic activity in the different organisms' size fractions across the open oceans. Another example is [Toseland et al., 2013], which included samples from the Arctic Ocean, North Atlantic Ocean, the North Pacific Ocean, the Equatorial Pacific and the Southern Ocean. This analysis included environmental data, which enabled the authors/researchers to investigate and identify that cellular resource allocation is significantly impacted by temperature [Toseland et al., 2013].

2.3.3 Dataset applications

It has been estimated that approximately 99% of microbes in the environment cannot be cultured under laboratory conditions. This greatly limits our understanding of microbial diversity, genetics and community ecology [Singh et al., 2009]. Advances in high-throughput sequencing technologies have vastly altered our ability to investigate these unculturable microbes [Bik, 2014]. Taxonomic classification involving the 18S and 16S rDNA provides identification of the reads in a sample for eukaryote and prokaryote species, respectively. These samples can also be compared for differences and similarity in composition, distribution and abundance [Wang et al., 2014]. While with metatranscriptomics, this provides insight into microbial communities, with a particular interest in how these microbial communities react to changes in the environment [Gilbert and Hughes, 2011], [Narayanasamy et al., 2016].

The 18S and 16S rDNA genes have been employed for decades for taxonomic classification studies [Wang et al., 2014]. In that time there have been numerous studies such as [Hunt et al., 2013] which employed 16S rDNA to investigate Bacterioplankton in the surface ocean, specifically looking at the relationship between abundance and their activity. Another example is [Stecher et al., 2016] which employed 18S rDNA to examine protist diversity in sea ice from the central Arctic Ocean. Default parameters were chosen for these and many other studies. When constructing our 18S rDNA pipeline we were confident that the default parameters would provide the best output as these parameters worked well in other studies. For instance in our 18S rDNA pipeline we employed a multiple sequence aligner called ClustalW [Thompson et al., 1994] and

used the default parameters. Likewise in the 18S rDNA study of [Shull et al., 2001] ClustalW was employed with default parameters to investigate the phylogenetic relationship in the suborder group of Adephaga [Shull et al., 2001]. Aligning protein coding genes such as 18S and 16S rDNA genes which are highly conserved within eukaryotes and prokaryotes, respectively, is not as difficult or problematic as aligning noncoding DNA sequences [Wang et al., 2006], [Amit Roy et al., 2014], [Wu et al., 2015]. There is a minimum amount of difficulty in producing convincing alignments of protein coding DNA sequences when sequence divergence is low, this is because indels mostly occur in multiples of three base pairs, and rarely within codon regions [Keightley and Johnson, 2004].

2.4 Sequence preprocessing methods

Sequence preprocessing is a common first step that usually involves quality control (QC), as well as identifying and filtering unwanted data. One source of the unwanted data is low quality sequences, which may occur as a result of sequencing instrument limitations or sample preparation problems [Zhou et al., 2014]. A major step in Illumina sequencing involves the ligation of the target DNA fragments to specific adapters for clonal amplification [Aird et al., 2011]. Artifacts such as adapter sequences which have not been fully removed are another source of unwanted data [Schmieder and Edwards, 2011]. Raw sequencing files can contain millions of reads and downstream analysis is computationally intensive, the demand on computational resources is one of the major bottlenecks for analysis [Schmieder and Edwards, 2011], [Yu and Zhang, 2013]. Unwanted data, if present, can drain resources and result in misassembly or erroneous conclusions (see for example [Schmieder and Edwards, 2011]).

Another aspect of sequence preprocessing is merging paired end reads. Illumina performs paired-end sequencing which can produce reads from both ends of target DNA fragments. If the insert size is sufficiently small, the two paired-end reads can be merged and this results in an increased overall read length that can have a significant and favourable affect on the quality of the analysis [Magoč and Salzberg, 2011], [Zhang et al., 2014]. A large number of short reads are generated from NGS technologies such as Illumina. Short read lengths are problematic for analysis such as *de novo* assemblies,

even for very deep genome coverage [Magoč and Salzberg, 2011]. Illumina sequences can produce single-end reads that range from 75 to 300 bp, and there is an increase in error rates as the reads get longer. When paired-end reads are merged sequencing errors can be corrected when the paired-end reads overlap and therefore potentially give higher quality reads [Zhang et al., 2014].

There are numerous tools available for sequencing preprocessing of raw sequencing files, such as those reviewed in [Zhou et al., 2015]. We do not describe all of the tools that are available, but only those used in our own pipelines.

2.4.1 Adapter searching and trimming

Duk

Duk [Li et al., 2011] (which stands for Decontamination Using kmers) is a matching tool for DNA sequences [Li et al., 2011], [DOE Joint Genome Institute, 2017]. This tool was implemented in our 18S rDNA pipeline. Duk (otherwise known as BBDuk) is part of a suite of bioinformatics tools called BBtools developed by the Joint Genome Institute [DOE Joint Genome Institute, 2017]. Duk is used to screen for contamination such as adapter sequences in raw reads after sequencing. Besides contamination removal, Duk has a range of other applications, such as organelle genome separation and assembly refinement [Li et al., 2011]. The matching technique used by Duk is similar to aligning sequences, but instead of creating an alignment between sequences, Duk takes a query sequence to search against a reference sequence for partial or total matches. Many traditional contamination searching tools are usually performed with an alignment technique. But since it is only necessary to know if a match is present or not, and not which bases of a query sequence match to which position of a reference sequence, Duk's performance is much faster than tools that employ alignment techniques [Li et al., 2011].

Duk uses a kmer (substring of size k) hashing method to index reference sequences to identify matching DNA in the query sequences [Li et al., 2011], [Mapleson et al., 2017]. Duk estimates the occurrence of the match by calculating p-values from a Poisson distribution. In order to calculate this Duk assumes that the kmers in the DNA sequence are randomly distributed. Assuming that there are M kmers in reference

sequences, N kmers in the query sequence and that the kmer size is k , the probability of having c or more common kmers between the reference sequence and the query sequence is given by the formula $1 - \text{pois}(c, u)$, where pois is the Poisson distribution function and $u = M * N / 4^k$ [Li et al., 2011]. An output with a high p-value implies that the match between the query sequence and the reference sequence is probably the result of random sequence variation, as opposed to the match event being the result of sequence homology [Li et al., 2011].

Cutadapt

Cutadapt [Martin, 2011] is a software tool that searches and removes adapter sequences from reads that were left on after sequencing. This tool was implemented in our 18S rDNA pipeline. Adapter sequences are considered a form of contamination and must be removed so only the relevant part of the read is used for further analysis. A read which contains an adapter sequence can be trimmed or completely discarded. After trimming if a read falls outside a specified length range, the read can be discarded [Martin, 2011].

For each read, Cutadapt begins by calculating the optimal alignments between the read and all given adapter sequences. Cutadapt's algorithm is called regular semi-global alignment and it does not penalise initial or trailing gaps, thus allowing the read and the given adapter sequence to shift freely relative to one another. The user can tell Cutadapt to look for the adapter at the 3' end of the molecule by giving the "-a" parameter to provide the adapter sequence. This results in Cutadapt removing the adapter sequence and all the nucleotides after it. The adapter sequence must begin at the start or within the read. This is done by penalising initial gaps in the read sequence. If the location of the adapter sequence is unknown, the user can provide the "-b" parameter. Therefore, if the adapter sequence is overlapping the beginning of the read, all nucleotides before the first non-adapter nucleotide are deleted. After Cutadapt has aligned all adapters to the read, the alignment with the greatest number of matching nucleotides between the read and adapter is taken to be the best one. An error rate of e/l is calculated, where e is the number of errors and l is the length of the matching piece between read and adapter. The read is trimmed when the error rate is below the allowed maximum [Martin, 2011].

2.4.2 Read merging

USEARCH merging of paired reads

USEARCH merging of paired reads [Edgar, 2010b] is a tool that merges overlapping paired end reads into single sequences [Edgar, 2010e]. This tool was implemented in the Joint Genome Institute (JGI) 16S rDNA pipeline. The output gives increased length and higher quality of reads, thus improving the quality of the analysis [Zhang et al., 2014].

USEARCH merging of paired reads is accomplished by aligning the forward and reverse reads and also rectifying any mismatches that may have occurred during aligning. The reverse read is the reverse complement to the forward read. Therefore the reverse read is positioned on the same strand as the forward read. Generally, the alignment between the forward and reverse reads is covered between them partially. This, therefore, results in unaligned segments at the start of both reads. The alignment is staggered in the event that the sequencing construct is shorter than the read length, which results in the unaligned segments at the end of both reads rather than the start. By default, USEARCH with the command “fastq mergepairs” will trim the unaligned segments in the event that the alignment is staggered. When mismatches occur in the alignment, the base call with the biggest Q (Phred) score is selected, but if both bases have the same Q score then the forward read base is chosen [Edgar, 2010e].

FLASH

FLASH [Magoč and Salzberg, 2011] is a software tool that locates the correct overlap between paired-end reads and produces a merged read of greater length than a single read [Magoč and Salzberg, 2011]. This tool was implemented in our 18S rDNA pipeline. The output of increased read length results in a significant and beneficial affect on the quality of the analysis such as genome assembly [Magoč and Salzberg, 2011].

Paired-end reads are produced from both ends of the target DNA fragments. FLASH takes each read pair separately and looks for the right overlap between the paired-end reads. The two reads are merged when the right overlap is found, thus producing a longer read. The new longer read is of the same length as the original target DNA fragment from which the paired-end reads were generated. FLASH examines and scores

every possible ungapped alignment overlap between paired-end reads, in order to find the right overlap. FLASH's scoring system accounts for the number of bases overlapping between the paired-end reads, this parameter is called min-olap and the default is set to 10 base pairs (bp) [Magoč and Salzberg, 2011].

2.4.3 Read clustering

CD-HIT

CD-HIT [Li and Godzik, 2006] is a tool for clustering raw sequencing data. This tool was implemented in our 18S rDNA pipeline. Clustering is a data reduction strategy that reduces sequence redundancy and therefore improves the performance of downstream analysis, for example reducing storage space and computational time [Fu et al., 2012]. At first, CD-HIT was designed to cluster protein sequences to build reduced redundancy reference databases such as UniProt, but it was later expanded to support clustering nucleotide sequences [Li and Godzik, 2006], [Fu et al., 2012]. Since CD-HIT's release, it has become very popular for a large range of applications including, for example, protein family classification, metagenomics annotation and identifying artifacts in datasets [Fu et al., 2012].

CD-HIT implements a greedy incremental algorithm. Firstly, CD-HIT orders the sequences by length. The longest sequence becomes the seed for the first cluster and each remaining sequence is compared with established seeds. If there is a similarity with any seed which meets a pre-defined cutoff, it is grouped into that cluster; else, it begins a new cluster [Li et al., 2012]. The similarity between seeds is determined by common word counting, by using word indexing and counting tables. This results in filtering out the undesirable sequence alignment, which can be used to calculate exact similarity. Therefore, a big redundant dataset can be represented by a smaller non-redundant dataset, then each cluster can be represented by a single entry [Fu et al., 2012].

USEARCH cluster otus

USEARCH cluster otus [Edgar, 2010b] is a tool that is employed to cluster sequences, specifically reads from a marker gene amplicon sequencing experiment such as 16S

rDNA [Edgar, 2010a]. This tool was implemented in the JGI 16S rDNA pipeline. As stated above clustering is a data reduction strategy [Fu et al., 2012].

USEARCH cluster otus performs operational taxonomic unit (OTU) clustering, with identity threshold fixed at 97% using the UPARSE-OTU algorithm [Edgar, 2010a]. The aim of the UPARSE-OTU algorithm is to find a set of OTU representative sequences and these must satisfy four criteria points. The criteria points are as follows, one, all pairs of OTU sequences must possess 97% or greater pairwise sequence identity. Two, the OTU sequence must be the most abundant within the 97% group [Edgar, 2010c]. Three, a chimeric amplicon occurs when a partially complete DNA strand anneals to a different template. A new template is synthesised by the primers based on the two different biological sequences and therefore any chimeric sequences identified are discarded [Edgar, 2010c], [Edgar, 2016]. Four, all non-chimeric input sequences must match at least one OTU with 97% or greater pairwise sequence identity. UPARSE-OTU is a greedy algorithm. The input to the UPARSE-OTU algorithm is a set of sequences and each sequence is labelled with a number representing its abundance of reads having a given unique sequence. The sequences are ordered in decreasing abundance, as OTU centroids (representative sequences) are chosen from the more abundant reads [Edgar, 2010c].

Every input sequence is compared to the current OTU database and using the UPARSE-REF [Edgar, 2010d] algorithm a maximum parsimony model of the sequence is determined [Edgar, 2010c]. There are three cases as follows, one, the UPARSE-REF model has 97% or greater sequence identity to an existing OTU, and therefore that input sequence becomes a member of the OTU [Edgar, 2010c]. Two, the model is chimeric, and then the input sequence is discarded [Edgar, 2010c]. Three, the model is greater than 97% and possesses the highest sequence identity to any current OTU. Therefore the input sequence is added to the database and also becomes the centroid of a new OTU [Edgar, 2010c].

2.5 Database searching

In a typical analysis, the datasets of sequences first undergo a similarity search in order to classify these sequences. This is achieved by comparing the sequences to public

databases of annotated sequences such as GenBank [NCBI Resource Coordinators, 2016] or SILVA [Quast et al., 2013b] [Bazinet and Cummings, 2012], [Yu and Zhang, 2013]. Sequence alignments are used to identify similarity between sequences [Eric et al., 2014]. For instance, metagenomic environmental samples contain DNA sequences from many different species and the typical aims of a metagenomic analysis are to try to identify what species and genes are present [Thomas et al., 2012a], [Suzuki et al., 2015]. There are a number of other fields where sequence homology searches are common, including phylogenetic analysis [Pearson, 2013], [Suzuki et al., 2015]. A well founded evolutionary hypothesis based on molecular sequences is dependent largely on the ability to align sequences of nucleotides or amino acids in the same position [Barta, 1997]. There are a number of tools and databases available and we outline the tool we used in the next section.

HMMER

HMMER [Eddy, 1996] is a tool for sequence searching, which uses probabilistic methods known as Hidden Markov Models (HMMs) in the analysis of homologous amino acid and nucleotide sequences with a high degree of sensitivity [Finn et al., 2011], [Ferreira et al., 2014], [Jiang et al., 2016b]. This tool was implemented in our 18S rDNA pipeline. HMMER, which employs Profile Hidden Markov Models, provides a sensitive approach for the detection of distant homologs since sequence profiles can give an improved representation of a set of homologous sequences compared with a single sequence [Sinha and Lynn, 2014].

There are two main stages for profile HMM homology detection; first model building and then database searching. Building the model entails the conversion of a multiple sequence alignment into a probabilistic model, and database searching is the scoring of a sequence to the profile HMM [Wistrand and Sonnhammer, 2005]. As compared to searching with a single query sequence, a previously built sequence consensus is used. This is called a consensus profile, which gives a more flexible method for the identification of homologs of a particular family, by highlighting the features they have in common and by lessening the value of the divergences between the sequences [Ferreira et al., 2014]. The profile HMM is mainly comprised of three types of states, one for each of the labels we could assign to a nucleotide therefore corresponding to matches

or mismatches, insertions and deletions, all with precise transitions between the three states [Eddy, 2004], [Ferreira et al., 2014]. The most common algorithms to process HMMs are the Forward algorithm and the Viterbi algorithm. The Forward algorithm computes a full probability for all possible model state paths and the Viterbi algorithm provides the best possible sequence of model states for the generation of the query sequence. The Viterbi algorithm obtains the entire path of states, which corresponds to an optimal alignment of the query sequence to the profile model [Ferreira et al., 2014]. Profile HMMs have a reputation for generating good results, and therefore are employed by a number of databases, such as Pfam [El-Gebali et al., 2019] and Superfamily [Gough et al., 2001]. Within such databases there are a large collection of protein families where each family is represented by a profile HMM and this is what is commonly used to represent the family in database searches [Wistrand and Sonnhammer, 2005].

USEARCH oligodb

USEARCH oligodb [Edgar, 2010b] is a tool that is employed to search for matches to nucleotide sequences in a database of short nucleotide sequences. This tool was implemented in the JGI 16S rDNA pipeline. USEARCH oligodb is typically employed in the search of matches of primers or probes to genome sequences or to gene databases [Edgar, 2010f]. The success of any polymerase chain reaction (PCR) based method is greatly dependant on the correct nucleic acid sequence. The sensitivity and specificity of primers or probes are predicted by searching a database to find sequences that contain the ideal number of mismatches and similarity [Kalendar et al., 2017]. The USEARCH oligodb algorithm is not heuristic, it is therefore exact. The alignments are global, with no gaps allowed except in the case of terminal gaps in the query sequence [Edgar, 2010f].

2.6 Phylogenetic analysis

Phylogenetic analysis is the study of the evolution of species, by examining the relationships between molecules, phenotypes and organisms [Singh et al., 2009], [Rokas, 2011]. The goal of a phylogenetic study is to establish which tree out of all possible

trees gives the best estimate for the true evolutionary relationships of the dataset analysed according to some optimisation criterion. Unfortunately, due to computational expense, it is not typically possible to determine the best tree amidst all possible trees for the sequence data. This is even true for a small number of sequences, as the number of alternative trees are extraordinarily large since the number of all possible trees grows exponentially with the number of sequences. For instance, the number of different phylogenetic trees that portray the evolutionary relationships of 50 sequences is roughly the number of atoms in the known universe [Rokas, 2011].

Phylogenetic analysis is a standard and very important tool for any bioinformatician, since it helps in understanding big evolutionary questions, such as the origins and history of macromolecules, developmental systems, phenotypes, and of course life [Rokas, 2011]. In our approach, phylogenetic analysis is achieved by cloning and sequencing the ribosomal RNA (rRNA) genes, in particular the 16S/18S small subunit rRNA [von Mering et al., 2007]. This can closely estimate the level of species diversity and unusual organisms can also be identified by this approach [von Mering et al., 2007], [Medlar et al., 2014]. Also, determination of the taxonomic composition of environmental samples can provide important indicators into the underlying communities' ecology and function [von Mering et al., 2007]. Phylogenetic analysis is also essential to gene discovery and annotation, to prediction of gene function, and the identification and construction of gene families [Rokas, 2011].

There are a wide numbers of tools available for phylogenetic analysis [Pavlopoulos et al., 2010]. In the following sections we describe the tools that we used in our bioinformatics pipelines. These sections are ordered as in a typical phylogenetic analysis. First sequences are aligned, then a tree is built, so that new sequences can be placed onto this tree, and finally this tree is visualised.

2.6.1 Alignments

ClustalW

ClustalW [Thompson et al., 1994] is a heuristic multiple sequence alignment (MSA) program. This tool was implemented in our 18S rDNA pipeline. MSAs are vital to many areas in bioinformatics, for example in homology modelling and phylogenetic

analysis [Capella-Gutierrez et al., 2009]. The aim of MSA is to output an arrangement of a set of sequences, with the aim that similar sequence features are aligned together, so that patterns can be identified that may be common among many sequences, or changes revealed that may clarify functional and phenotypic variability. A feature can be defined as any relevant biological information, that is, structure, function or homology to the common ancestor [Kemena and Notredame, 2009]. The quality of MSAs for these applications is critical for the reliability and accuracy of the analyses. A large number of algorithms for MSA are presently available, which apply different heuristic algorithms to find the optimal solutions to the alignment problem. 80-90% accuracies have been reported for the best MSA algorithms, but even these algorithms can fail at specific regions in the alignment. For large scale analysis the problem gets worse, due to the implementation of faster algorithms that are less reliable [Capella-Gutierrez et al., 2009].

The basic ClustalW algorithm consists of three key parts. First all pairs of sequences are aligned separately and from this a distance matrix is calculated, thus giving the divergence of each pair of sequences. Second, from the distance matrix and using the Neighbor-Joining [Saitou and Nei, 1987] method a guide tree is calculated. Initially an unrooted tree is constructed with branch lengths proportional to the approximate divergence for each branch. By employing the mid-point method a root is placed at a point on the tree, in which the branch lengths on either side of the root are equal. Third, according to the branching order in the guide tree, from the leaves of the rooted tree towards the root, the sequences are progressively aligned. A dynamic programming algorithm at each stage of the alignment is performed with a residue weight matrix and also penalties for opening and extending gaps. Each part is made up of aligning two existing alignments or sequences and gaps that are introduced in previous alignments remain unaltered. New gaps that are added at each stage get full gap opening and extension penalties, regardless of whether or not they are added inside an original gap location. The score at a position from one sequence or alignment and another sequence or alignment is calculated based on the average of all the pairwise weight matrix scores from the sets of sequences used. If any set of sequences has one or more gaps in one of the locations being considered, this gets scored a zero if it is a gap versus a residue. The default amino acid weight matrices used are adjusted to

be assigned positive values. Consequently, this treatment of gaps results in the score of a residue versus a gap ending up with the worst possible score. Therefore when the sequences are weighted, each weight matrix value is multiplied by the weights from the two sequences [Thompson et al., 1994].

trimAL

trimAl [Capella-Gutierrez et al., 2009] is an automated trimming tool for multiple sequence alignment. This tool was implemented in our 18S rDNA pipeline. It has been reported that the removal of poorly aligned regions from an alignment increases the quality of further analyses [Capella-Gutierrez et al., 2009]. trimAl firstly reads all the columns in the alignment and calculates a score, a gap score, a similarity score or a consistency score for each of the columns. The score for each column is calculated on information from that column or, if a window size is given, it relates to the average value of the given window size columns around the position being considered. The gap score for a given column is the fraction of sequences with no gap in that specified position. The residue similarity score uses the mean distance score between pairs of residues, as defined by a given scoring matrix. The consistency score is only calculated when more than one alignment for the same set of sequences is given. The consistency score is the level of consistency for all residue pairs located in a given column as compared with other alignments. The alignment with the highest consistency score is trimmed to remove the columns that are less conserved.

trimAL can proceed in two ways after all the column scores have been calculated. A conservation threshold relates to the minimum percentage of columns, from the initial alignment, that the user would like to have in the trimmed multiple sequence alignment. If a score and a minimum conservation threshold are provided, trimAL will output a trimmed alignment. This alignment will consist only of the columns with scores greater than the score threshold. If the number falls below the conservation threshold, in a decreasing order of scores trimAl will add more columns to the trimmed alignment until the conservation threshold is hit. Alternatively, trimAl has three modes for the automated selection of parameters- gappyout, strict and strictplus- these are based on the use of gap and similarity scores. trimAl will calculate the specific score thresholds based on the characteristics of each alignment. trimAL also has an option

which implements a heuristic in order to decide the appropriate mode depending on the alignment characteristics [Capella-Gutierrez et al., 2009].

2.6.2 Tree building and placing

RAxML

RAxML [Stamatakis, 2014] (Randomised Axelerated Maximum Likelihood), is a tool for phylogenetic analysis of large datasets under maximum likelihood [Stamatakis, 2014]. This tool was implemented in our 18S rDNA pipeline. There has been a staggering accumulation of genetic information from various organisms in recent years [Stamatakis et al., 2005]. There are a number of approaches such as maximum likelihood and Bayesian methods which can be used to compute these relationships. Maximum likelihood is incorporated in RAxML and is considered to represent one of the more accurate approaches available for phylogenetic reconstruction [Stamatakis et al., 2005].

RAxML has a fast maximum likelihood tree search algorithm, and has been shown to return trees with “good” likelihood scores [Stamatakis, 2014]. RAxML generates a starting tree, which is built by adding sequences one at a time at random, using the parsimony optimality criterion to identify their optimal location on the tree. Due to the random order in which the sequences are added, this is likely to generate several different starting trees each time a new analysis is run. This can result in a improved exploration of the tree space. Moreover, if multiple analyses are run which use different starting trees and all result in the same tree, this gives more confidence that this is close to the true tree. The next step in the search strategy is the implementation of the method called “lazy subtree rearrangement”. This means that all possible subtrees are cut and reinserted at all possible locations of a tree. The number of branches between the cut and insertion points must be smaller than N branches. RAxML automatically estimates the N value for a data set or the user can give the value of N. The lazy subtree rearrangement method is applied on the starting tree, and there after multiple times on the currently best tree as the program continues. When RAxML reaches the point where a better tree cannot be found, RAxML ends the search and the tree is returned [Rokas, 2011].

Pplacer

Pplacer [Matsen et al., 2010] is a tool for performing phylogenetic placement, thereby assigning query sequences to taxa and providing an evolutionary understanding of the sequence data. This tool was implemented in our 18S rDNA pipeline. Phylogenetic placement is an alternative to the classic phylogenetic analysis for dealing with datasets with a very large number of sequences. Pplacer assigns the unknown query sequences to a fixed reference tree via a reference alignment according to the maximum likelihood criterion. Since pplacer uses a fixed reference tree, there are just two tree searches needed to precalculate the information required from the reference tree. Therefore, all likelihood calculations are at once performed on the set of three taxon trees. Hence the query sequence placement part of pplacer has linear time and space complexity for the number of taxa in the reference tree.

A problem with likelihood based phylogenetic analysis is that it cannot be applied to the very large number of short reads that are produced by next-generation sequencers [Matsen et al., 2010]. This is due to computational complexity because the maximum likelihood phylogenetic problem is NP-hard (nondeterministic polynomial time) and therefore it is probably not possible to find maximum likelihood trees in a reasonable amount of time [Matsen et al., 2010], [Papadimitriou, 2014]. Also the lack of phylogenetic signal is a problem, because employing the classic method of maximum likelihood phylogeny to a single alignment of shotgun reads with the full-length reference sequences can result in the inaccurate grouping of a short read because of its position in the alignment.

Pplacer can apply the inferential strength of likelihood based methodologies, that enables the fast placement of the large amount of short query sequences and avoid some of the problems associated in applying phylogenetics to a very large number of taxa. This is accomplished as the computing complexity is greatly reduced, resulting in a program that can assign large amounts of query sequences per processor each hour on to a fixed reference phylogeny. Pplacer performs the calculation in parallel, because each query sequence can be processed independently and also the relationships between the query sequences are not investigated, thus reducing exponential time to a linear time. Short length query sequences are less of an issue for pplacer, as they are

compared to the full length of the reference sequences [Matsen et al., 2010].

Output from pplacer is a series of assignments of every query sequence to branches on the reference tree and a confidence score [Matsen et al., 2010], [Matsen et al., 2012]. A query sequence can be assigned to more than one branch to show placement uncertainty for that sequence [Matsen et al., 2012]. Pplacer computes edge uncertainty using a likelihood weight ratio, which measures the uncertainty edge by edge. It does this by comparing the best placement locations for each of the edges. Pplacer applies expected distance between placement locations (EDPL) to overcome the difficult situation in distinguishing between uncertainty of local (where placements are all located in a small area of the reference tree) and global (where possible placements are spread throughout the tree) placements. This is a problem because it depends on the confidence scores computed on an edge by edge basis.

A naive algorithm would place the query sequence onto each edge of the tree and execute a complete branch length optimization by carrying out the cached likelihood vectors. However, Pplacer achieves linear time and space scaling for the size of the reference tree because it performs an initial calculation of likelihood vectors at both ends of each edge of the reference tree. Pplacer executes a two part search algorithm for the query sequences to speed up placements. In particular, an initial quick evaluation of the tree is performed, then a more complete search is made on high scoring regions of the tree. The first part is carried out by calculating likelihood vectors for the middle of each edge. The calculated likelihood vectors are used to rapidly sort the edges in a rough order of fit for each query sequence [Matsen et al., 2010].

2.6.3 Visualisation

MEGAN

Earlier metagenomic projects were focused on the identification of species and their function from individual data sets, but over the past few years there has been a growing emphasis on comparative analysis. MEGAN [Huson et al., 2007] (MEtaGenomic ANalyzer) facilitates interactively examining, analysing and comparing multiple environmental datasets [Huson et al., 2009]. MEGAN was implemented in our 18S rDNA pipeline and the JGI 16S rDNA pipeline. MEGAN establishes its taxonomic classifi-

cation based on the NCBI taxonomy. This is a hierarchically structured classification of all species that are currently represented in NCBI databases. MEGAN performs the taxonomic analysis by assigning each read onto different taxa in the NCBI taxonomy. For determining gene content a sequence comparison of the query reads to one or more databases of determined reads is performed using a comparison tool such as BLAST. Since different projects need to use different alignment tools and databases, MEGAN gives the users the freedom to choose which ever suits their project needs. The results of the sequence comparison are processed by MEGAN. This entails the collection of all matching reads to the sequences of the NCBI database and assigning a taxon ID to each sequence based on the NCBI taxonomy.

MEGAN employs the Lowest Common Ancestor (LCA) algorithm to assign reads, which involves assigning reads to the node representing the lowest ancestor of all high-quality matches for the sequence. As a result, the species-specific sequences are assigned to taxa closer to the leaves and generally conserved sequences end up being assigned to higher order taxa that are nearer the root of the NCBI tree. A MEGAN file is generated that consists of all the information for analysing the output and to produce graphical and statistical outputs [Huson et al., 2011]. The end result is represented on a rooted tree, so that each node displays different taxa, and nodes are scaled and labelled by the number of reads that are assigned to that taxon as shown in figure 2.2 [Mitra et al., 2011].

MEGAN's functional analysis of a microbial community can help to understand the biochemical processes and to estimate the impact of environmental changes in the various ecosystems [Mitra et al., 2011]. MEGAN accomplishes this by using the SEED classification of subsystems and the Kyoto Encyclopedia for Genes and Genomes (KEGG) classification of pathways and enzymes [Huson et al., 2011]. Based on a BLAST file imported into MEGAN, a SEED classification is determined by assigning each read of the highest scoring gene in a BLAST comparison against a protein database to the functional role. Different functional roles are then grouped into assigned subsystems. A rooted tree is used to display the SEED classification, the internal nodes are the different subsystems and the leaves are the functional roles. The tree can be multi-labelled, which means that different leaves on the tree can have the same functional role, as they may occur in different subsystems. The nodes of the tree display the numbers of reads

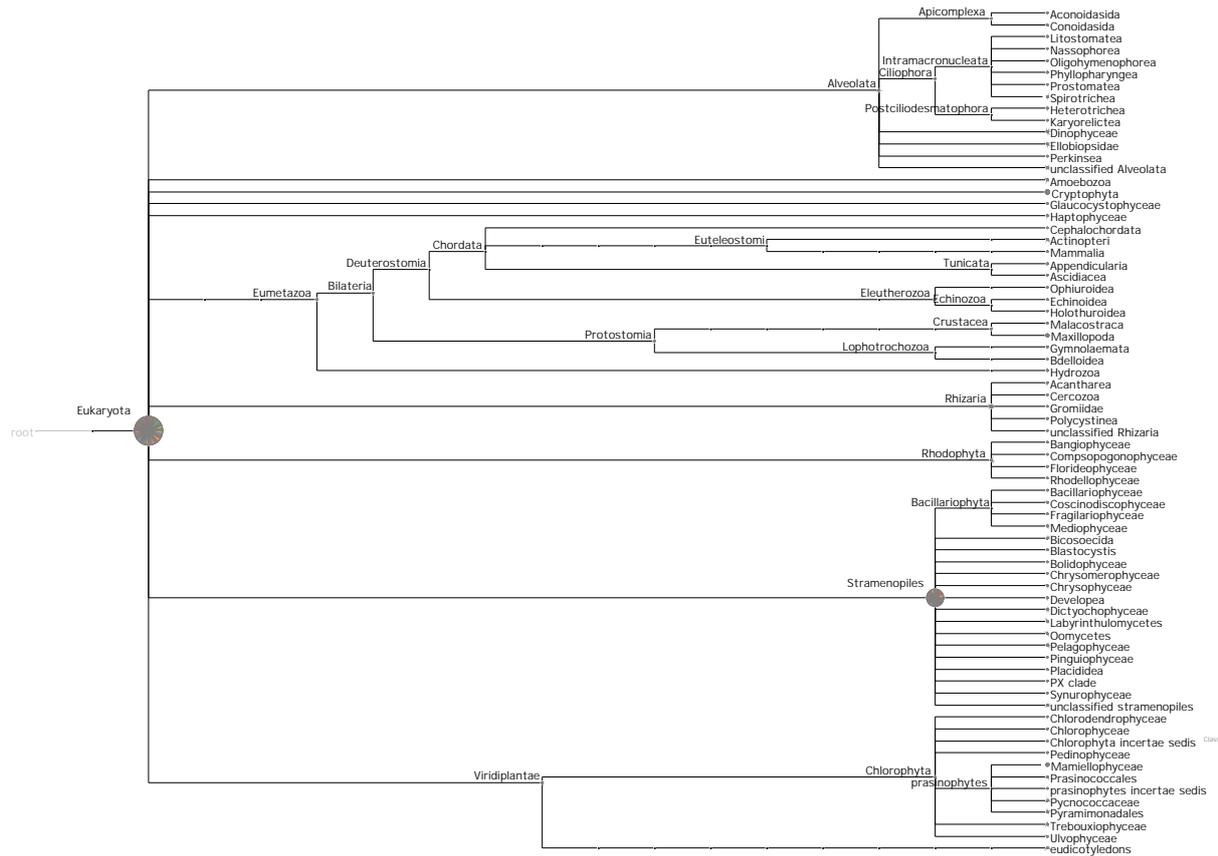


Figure 2.2: 18S rDNA dataset displayed on MEGAN’s taxonomy tree at the taxonomic rank of class

assigned to each functional role. For a comparative analysis, it is possible to map a number of datasets onto the SEED hierarchy and, also based on their SEED content, calculate the distance matrices on the datasets. For the KEGG analysis, MEGAN aims to match each read to a KEGG orthology (KO) accession number, taking the highest scoring match to a reference sequence that a KO accession number is known. A KEGG analysis window in MEGAN presents which KEGG pathways are in the dataset. The user can examine these pathways; for example, visualising reads that are mapped to a pathway of interest [Mitra et al., 2011].

2.7 Discussion

In this chapter, we have given a summary of 18S and 16S rDNA, metatranscriptomics, next-generation sequencing technology and bioinformatic tools. In the next chapter, we will describe in more detail the methodology that we developed, specifically to analyse our 18S and 16S rDNA datasets.

Chapter 3

18S rDNA and 16S rDNA analysis

3.1 Summary

This chapter outlines the computational pipeline and analysis of the 18S rDNA and 16S rDNA data collected in the expedition mentioned in chapter 1. As explained in the next section, samples were collected from a range of latitudes, from the South Atlantic Ocean, spanning the West African and European coasts up to the Arctic Ocean as shown in figure 3.1a. Also at the time of sampling environmental data was recorded, such as temperature and salinity as shown in figure 3.1b-f. This enabled us to study the composition and distribution of the marine microbial communities in the upper ocean.

Two different pipelines were implemented to taxonomically classify the 18S and 16S rDNA data. Our 18S rDNA pipeline which is outlined in section 3.3.1 is based on pplacer. We choose pplacer for our 18S rDNA pipeline because for unknown sequences likelihood-based phylogenetic inference is commonly regarded to be the most reliable classification method [Matsen et al., 2010]. While JGI's 16S rDNA pipeline which is outlined in section 3.3.2 is based on USEARCH. USEARCH is a popular software package for the analysis of operational taxonomic units (OTUs). USEARCH performs a blast like mapping against a reference database such as SILVA. This works well for microbes that are represented well in ribosomal RNA databases [Hugert and Andersson, 2017].

In the next section, we describe the sampling and sequencing of the 18S rDNA and 16S rDNA datasets. Then, in Section 3.3 we describe the pipelines for the 18S and

16S analysis, as well as additional methods, after which we present the results of our analysis in Section 3.4. In Section 3.5, we end with a discussion of the results that we obtained.

3.2 Sampling and sequencing

3.2.1 Sampling

Samples were collected during three field expeditions across latitude ranges as shown in figure 3.1a. The first set of samples was collected from April to May 2011 in the

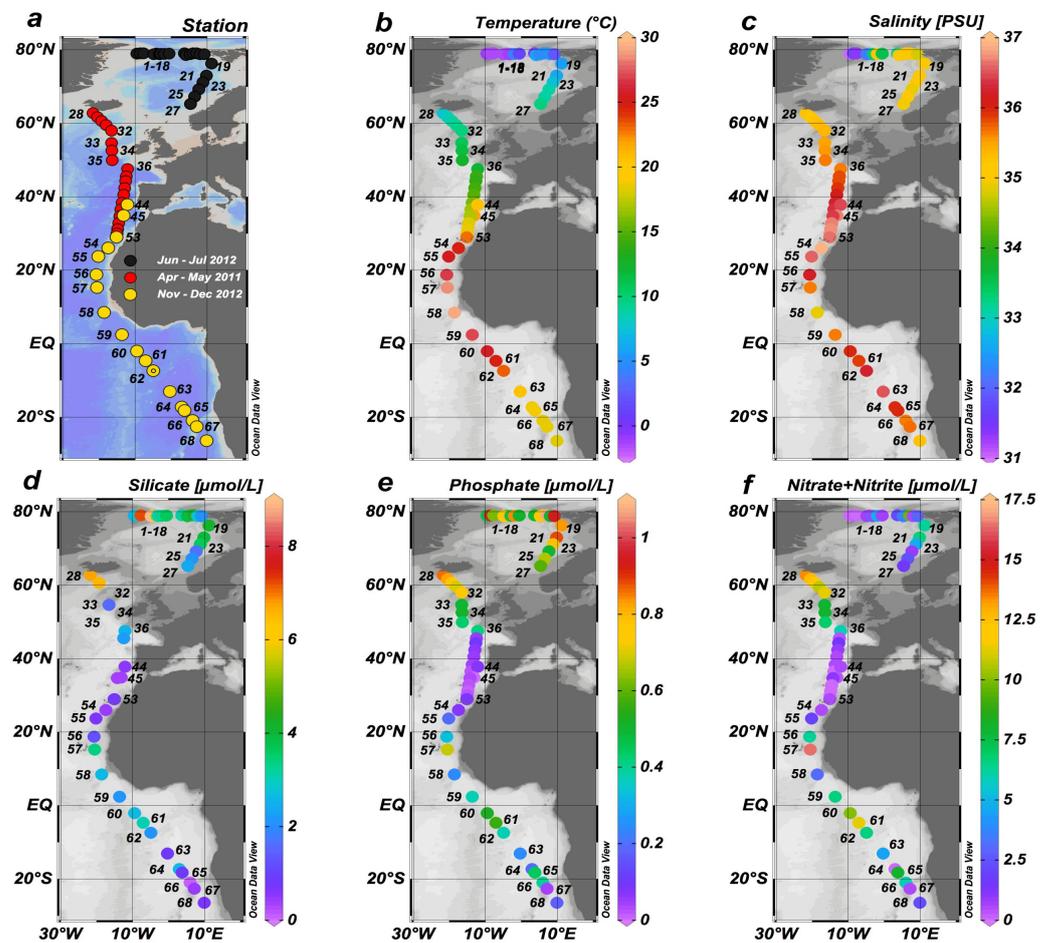


Figure 3.1: Arctic and Atlantic Ocean sampling sites and measured metadata, (a) Sites of three expeditions, April to May 2011 in red, June to July 2012 in black and November to December 2012 in yellow. At each station, microbial communities were sampled at the deep chlorophyll maximum (DCM), corresponding to 68 samples altogether. (b) Isosurface plot of temperature ($^{\circ}\text{C}$) measured at sampling depth. (c) Salinity (practical salinity unit(PSU)) measured at sampling depth for all stations. (d) Dissolved silicate (mol/L) concentrations measured at sampling depth for each station. (e) Concentration of dissolved phosphate (mol/L) measured at sampling depth for each station. (f) Nitrate and Nitrite (mol/L) concentrations measured at sampling depth for each station. (Figure was generated with Ocean Data view, R. Schlitzer, www.odv.awi.de, 2016)(Plot generated by Dr.Katrin Schmidt)

North Atlantic Ocean spanning the Canary Islands to Iceland, by Dr. Willem van de Poll of the University of Groningen, Netherlands and Dr. Klaas Timmermans of the Royal Netherlands Institute for Sea Research. The second collection of samples was carried out from June to July 2012 by Dr. Katrin Schmidt in the Arctic Ocean, spanning the West Spitsbergen current, east Greenland current and Norwegian Atlantic current. The third set of samples was also collected by Dr. Schmidt from November to December 2012, spanning the Canary Islands down to Cape Town in the South Atlantic Ocean.

The samples were taken either at the chlorophyll maximum (10-110m) and/or surface of the ocean (0-10m). The samples were filtered and frozen in liquid nitrogen and stored at -80°C until further analysis. A full description of the materials and methods used for sampling can be found in [Schmidt, 2016]. Also at the time of sampling environmental data was collected (see Appendix A.A) as shown in figure 3.1**b-f**.

3.2.2 Sequencing and preprocessing of the 18S rDNA and 16S rDNA

All extracted DNA samples were sequenced and preprocessed by the Joint Genome Institute (JGI) (Department of Energy, Walnut Creek, CA, USA). Amplicon sequencing was performed with primers for the V4 region of the 16S and 18S rRNA gene on an Illumina MiSeq instrument with a 2x300 bp read configuration [Tremblay et al., 2015].

18S sequences were preprocessed, this consisted of scanning for contamination with the tool Duk [DOE Joint Genome Institute, 2017] and quality trimming of reads with cutadapt [Martin, 2011]. Paired end reads were merged using FLASH [Magoc and Salzberg, 2011] with a max mismatch set to 0.3 and min overlap set to 20. A total of 54 18S samples out of 68 passed quality control after sequencing. After read trimming, there was an average of 142,693 read pairs per 18S sample with an average length of 367bp and 2.8 Gb of data over all samples.

16S sequences were preprocessed, this consisted of merging the overlapping read pairs using USEARCH's merge pairs [Edgar, 2017] with the parameter minimum number of differences (merge max diff pct) set to 15.0 into unpaired consensus sequences. Any reads that could not be merged are discarded. JGI then applied the tool USEARCH's search oligodb with the parameters length mean (len mean) set to 292, length

standard deviation (len stdev) set to 20, primer trimmed max difference (primer trim max diffs) set to 3, a list of primers and length filter max difference (len filter max diffs) set to 2.5 to ensure the Polymerase Chain Reaction (PCR) primers were located with the correct direction and inside the expected spacing. Reads that did not pass this quality control step were discarded. With a max expected error rate (max exp err rate) set to 0.02, JGI evaluated the quality score of the reads and those with too many expected errors were discarded. Any identical sequence was de-replicated. These are then counted and sorted alphabetically for merging with other such files later. A total of 57 16S samples passed quality control after sequencing. There was an average 393,247 read pairs per sample and an average base length of 253bp for each sequence with a total of 5.6 Gb.

3.3 Methods

In this section, we describe the pipelines that we developed for the 18S and 16S rDNA analyses, as well as further analysis methods that were employed.

3.3.1 Computational pipeline for 18S rDNA analysis

We first describe the computational pipelines that we developed for taxonomically classifying the 18S rDNA data. Also, we describe how we account for 18S rDNA copy number variation in order to give an estimate of abundance for the species in our dataset. An overview of the 18S rDNA classification pipeline is shown in figure 3.2. This consists of the construction of the 18S reference database highlighted in box (a), phylogenetic analysis highlighted in box (b), phylogenetic placement highlighted in box (c) and normalising the 18S rDNA copy number in box (d) of the figure. Each of these is explained in the following subsections.

Constructing the 18S reference database (box a)

We began by compiling a reference dataset of 18S rRNA gene sequences that represent algae taxa for the construction of a phylogenetic tree. We retrieved sequences of algae and outgroup taxa from the SILVA database [Quast et al., 2013b] and Marine Microbial Eukaryote Transcriptome Sequencing Project (MMETSP) database [Keeling et al.,

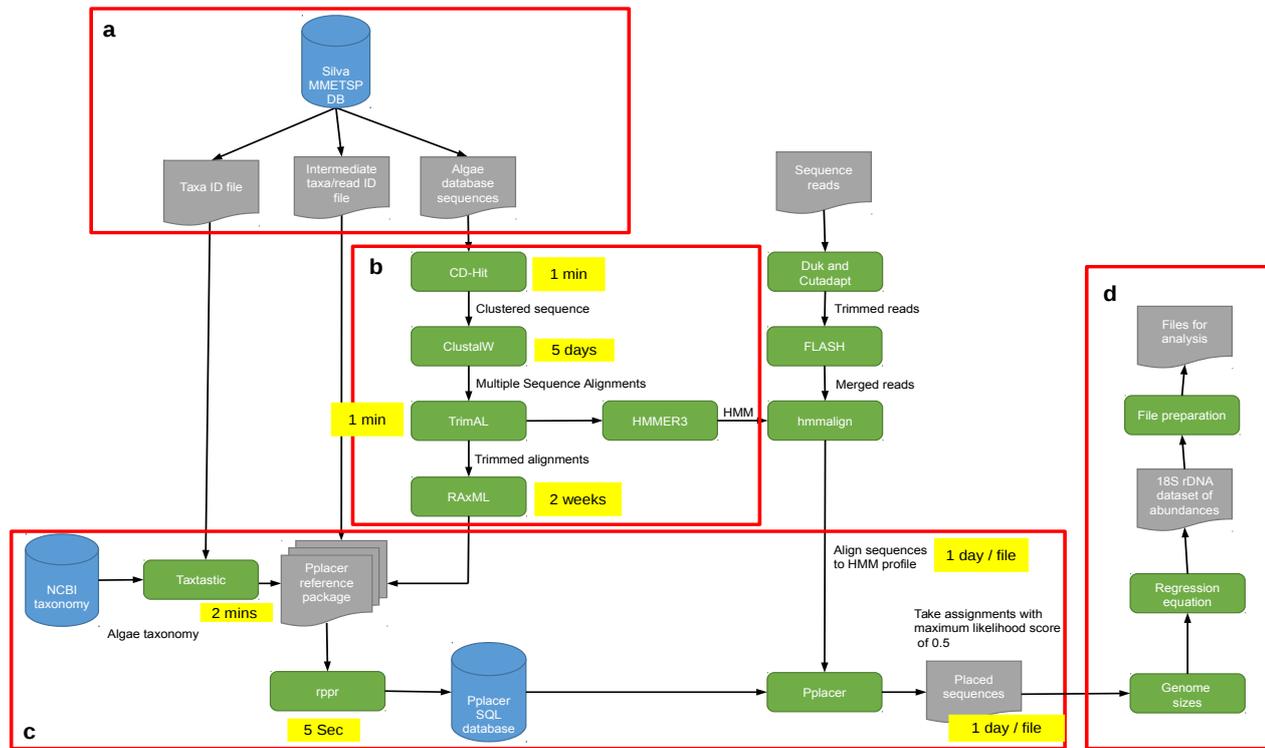


Figure 3.2: Pipeline diagram of Pplacer 18S rDNA classification analysis. The pipeline at various stages incorporates databases (blue), software tools (green), processed files (grey) and runtimes (yellow). Boxes **a**, **b**, **c** and **d** refer to sections in the text

2014]. For our outgroup species, these species came from marine invertebrates, plants, fish, zooplankton and fungi, for example, *Aurelia aurita*, *Zea mays*, *Salmo salar*, *Acartia tonsa* and *Saccharomyces cerevisiae*, respectively and also *Homo sapiens*. The inclusion of outgroups is important for the identification of potential sources of contamination that may have occurred during the earlier stages. The algae reference database consists of 1636 species from the following groups: Opisthokonta, Cryptophyta, Glaucocystophyceae, Rhizaria, Stramenopiles, Haptophyceae, Viridiplantae, Alveolata, Amoebozoa and Rhodophyta as shown in figure 3.3. A full list of the algae reference database taxa IDs can be found in Appendix A.B.

An accurate estimate of phylogeny is essential, not just for ecology but for other areas of research such as genomics. Under certain conditions, phylogenetic methods could return incorrect estimates [Wiens and Tiu, 2012]. A possible common scenario is when too few taxa are included in the reference database, thus leading to such situations as conflicting phylogenetic signals which can cause a lack of resolution and

incongruent phylogenies [Nabhan and Sarkar, 2011], [Wiens and Tiu, 2012]. Adding more taxa even with highly incomplete character data to the reference database can improve phylogenetic accuracy in scenarios where the analysis was misled by a limited number of taxa in the reference database [Wiens and Tiu, 2012].

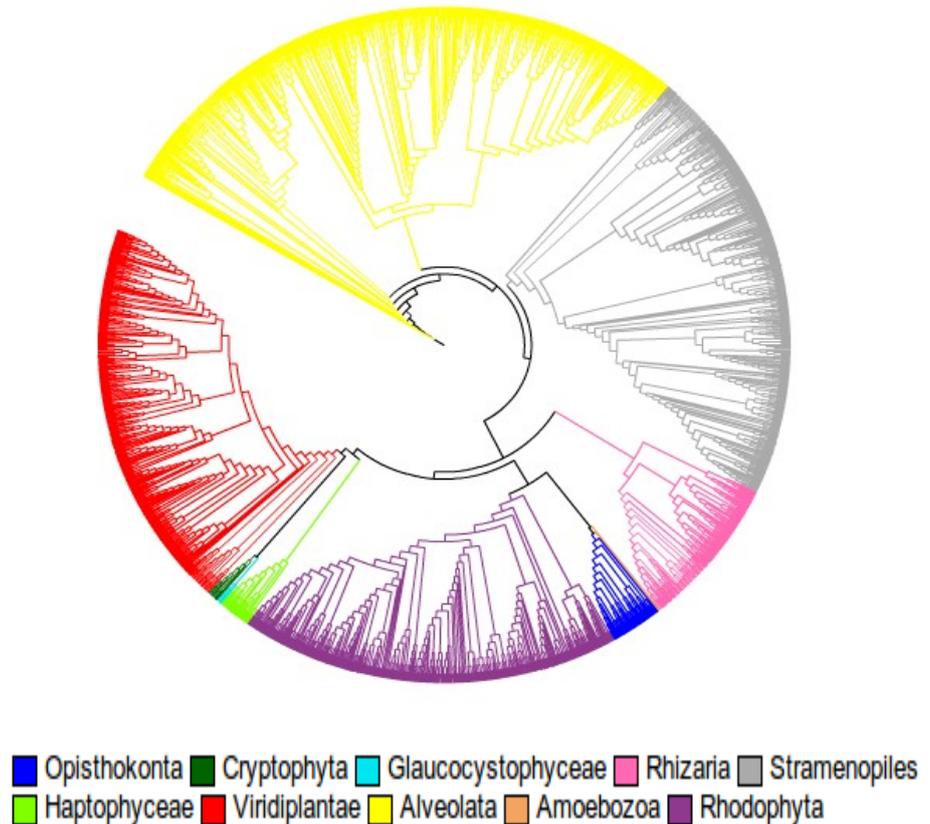


Figure 3.3: The reference tree consists of 1636 species from the groups: Opisthokonta, Cryptophyta, Glaucocystophyceae, Rhizaria, Stramenopiles, Haptophyceae, Viridiplantae, Alveolata, Amoebozoa and Rhodophyta

In order to construct the algae 18S reference database, we first retrieved all eukaryote species from the SILVA database with a sequence length of ≥ 1500 base pairs (bp) and converted all base letters of U to T. Under each genus, we took the first encountered species to represent that genus. Using a custom script by Dr. Andrew Toseland, the species of interest (as stated above) were selected from the SILVA databases, classified with NCBI taxa IDs and a sequence information file was produced that describes each of the algae sequences by their sequence ID and NCBI species ID. The taxonomy database from the NCBI, eukaryote sequences from the SILVA database and a list of algal taxa including outgroups were used as input for the script. This information was

combined with the MMETSP database, excluding duplications.

Phylogenetic analysis (box b)

As depicted in figure 3.2 box (b), we clustered the algae reference database to remove closely related sequences with CD-HIT [Li and Godzik, 2006] using a similarity threshold of 97%. Then using ClustalW [Thompson et al., 1994], we aligned the algae reference sequences of the database with the addition of the parameter iteration numbers set to 5. The alignment was examined by colour coding each species to their groups and visualising in iTOL [Letunic and Bork, 2007]. We observed that a few species were misaligning to other groups and these were then deleted using Jalview [Waterhouse et al., 2009]. The resulting alignment was tidied up with TrimAL [Capella-Gutierrez et al., 2009] by applying parameters to delete any positions in the alignment that contain gaps in 10% or more of the sequence, except if this results in less than 60% of the sequence remaining [Capella-Gutierrez et al., 2009]. Our algae reference phylogenetic tree is displayed in figure 3.3. We constructed a maximum likelihood phylogenetic reference tree and statistics file based on our algae reference alignment by employing RAxML [Stamatakis, 2014] with a general time reversible model of nucleotide substitution along with the GAMMA model of rate heterogeneity. Based on the algae reference multiple sequence alignment, with HMMER3 [Eddy, 2009] for the 18S rDNA gene we created a Profile HMM (pHMM). Our pHMM differs from other pHMMs constructed by other phylogenetic groups in that our pHMM is based on a reference database composed of an update SILVA database at the time of our analysis and we included the MMETSP database. A full description of our reference database is outlined above in section “Constructing the 18S reference database (box a)”

Phylogenetic placement (box c)

As depicted in figure 3.2 box (c), we used the NCBI taxtastic tool [NCBI Resource Coordinators, 2016] to create a taxtastic file, which is a description of the lineages of all species back to the root in our algae reference database. We created this taxtastic file by submitting the taxa IDs for each species to extract a subset of the NCBI taxonomy. A placer reference package using the NCBI taxtastic tool was generated, which produced an organized collection of all the files and taxonomic information into one directory.

With the reference package, an SQLite database was created using pplacer's Reference Package PReparer (rppr). With hmmlalign, we aligned the query sequences to the reference set and created a combined Stockholm format alignment. Pplacer [Matsen et al., 2010] was used to place the query sequences on the phylogenetic reference tree by means of the reference alignment according to a maximum likelihood model. The placefiles were converted to CSV with pplacer's guppy tool. With the use of our custom made script, we took reads that had a taxonomic assignment with a maximum likelihood score of ≥ 0.5 and counted the number of reads assigned to each classification. This resulted in 6,053,291 reads that were taxonomically assigned for further analysis.

Pplacer's reference tree is fixed, in regards to the topology and branch length. Also, only two tree searches are necessary to precalculate the information that is required from the reference tree. From here all the likelihood calculations are performed on a set of three taxon trees, the number of this is linear to the number of reference taxa in the reference database. Therefore the placement part of the pplacer algorithm has linear time and space complexity for the number of taxa n in the reference tree [Matsen et al., 2010]. In figure 3.2 are the runtimes highlighted in the yellow boxes for the tools as described above that were implemented in our 18S rDNA pipeline. It took a day for an 18S rDNA file containing an average of 142,693 read pairs with an average length of 367bp to run with our pplacer (18S rDNA) pipeline. We were able to perform our pplacer pipeline even faster due to parallel computing on the Earlham Institute cluster thus enabling us to run all 54 files within 3 days. If another similar dataset was to be submitted to our pplacer pipeline, this new dataset would also finish in this time. If the number of taxa in our reference database was to be altered then this would affect the time and space complexity.

Ribosomal RNA (rRNA) gene copy number

The rRNA gene is a marker for taxonomic diversity but the relationship between amplicon and species abundance is indeterminate due to rDNA copy number variation within the genomes of different species. For bacteria, the 16S rDNA copy number can vary as much as from one to fifteen [Perisin et al., 2016]. For eukaryotes, the 18S rDNA copy number can vary even more greatly from one to thousands [de Vargas et al., 2015]. This is a hindrance to effectively analyse our rDNA copy number datasets. In order to

get an estimate of abundances of species in the samples, we had to normalise the data, that is we had to adjust for the rDNA copy number.

Even though there is a 16S rDNA copy number database called the Ribosomal RNA Operon Copy Number Database (rrnDB) [Klappenbach et al., 2001] it only contained 2,876 species in 2014. While this is a considerable amount, it compares little to the actual number of bacteria species in the world. This was evident to us when we attempted to apply the rrnDB database to our 16S rDNA dataset. In order to fill in the missing copy numbers in our 16S rDNA dataset, we used averages between closely related known taxa. We detail this in the section “Normalising 16S rDNA copy number”.

There is no database of 18S rDNA copy number available for eukaryote species. Previous work has explored the link between rDNA copy number and genome size [Prokopowich et al., 2003]. We decided to build on this approach in order to get a rough estimate of read count to 18S rDNA copy number. The regression line is an imperfect approach, due to rDNA copy number variation within species that can exist thus being a source of error. In the following sections, we detail our work in determining a regression equation based on 18S rDNA copy and genome size. We also explain how we retrieved where possible the genome sizes for our 18S rDNA dataset, and determined the genome sizes when they were unknown, in order to apply the regression equation to our 18S rDNA dataset.

18S rDNA copy number and genome size regression equation (box d)

In order to address the varying copy number among eukaryotes, we investigated the gene copy number and related genome sizes of 185 species across the eukaryote microbial tree [Godhe et al., 2008], [Carlton et al., 2013], [Torres-Machorro et al., 2010], [Oliver et al., 2007], [Moreau et al., 2012], [Boucher et al., 1991], [Hauser et al., 2010], [Prokopowich et al., 2003], [Rödström, 2017], [NCBI Resource Coordinators, 2016] and [Nordberg et al., 2014]. Species with single genomes were investigated for their 18S rDNA copy number. Outliers were found but all 185 species were included, as we researched as many published 18S rDNA copy numbers as we could find with known genome size, but published 18S rDNA copy numbers are very few especially in comparison to 16S rDNA copy numbers and with known genome size.

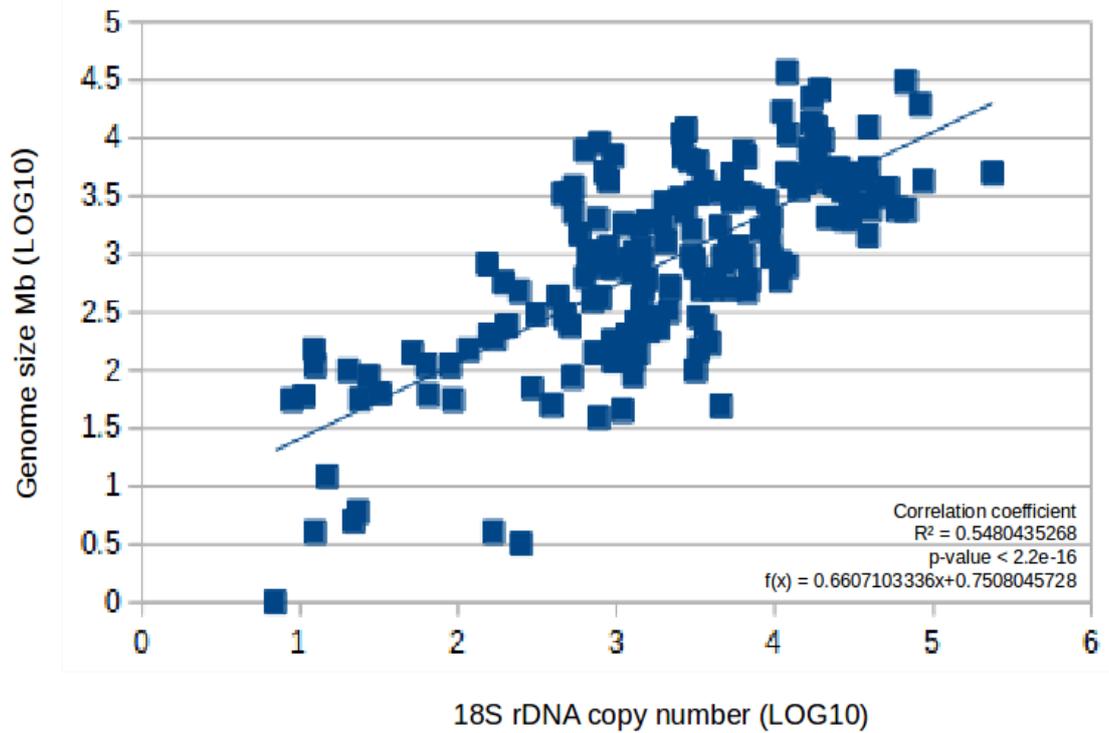


Figure 3.4: The graph of 18S rDNA copy number and their related genome size (Mb) for 185 species across the eukaryote tree. We investigated 18S rDNA gene copy number and their related genome sizes. We observed a significant correlation of R^2 0.5480435268 with a p-value < 2.2e-16 between genome size and 18S rDNA copy number. Based on the log10 transformed data a regression equation was determined, $f(x)=0.66X+0.75$

We performed a log10 transform of the data of 185 18S rDNA copy number and related genome sizes so that the data's scale and distribution complied better to the assumptions of a parametric statistical test. Log transformation ($x'_i = \log_b(x_i + c)$ where \log_b can be 2, e, or 10 and c is a small number when $x_i=0$) is commonly used [Paliy and Shankar, 2016]. Based on the log10 transformed data, a significant correlation with an R^2 of 0.55 and a p-value < 2.2e-16 between genome size and 18S copy number was observed. A regression equation was determined ($f(x)=0.66X+0.75$) as shown in figure 3.4.

Genome sizes for the 18S rDNA dataset to normalise the 18S rDNA copy number (box d)

In order to apply the equation of the line to our 18S rDNA dataset, we retrieved the genome sizes for the species in the dataset from the NCBI genome database. The NCBI genome database consisted at the time of 2,477 eukaryote entries. Firstly, since multiple entries of a species are in the NCBI genomes database due to different strains,

we calculated an average genome size for each species in the database, which resulted in 2,059 species entries.

The higher taxonomic levels for the NCBI genomes species needed to be established so that we could calculate the average of genome sizes. For a description of the lineages of all species back to the root in NCBI genomes database, we submitted the species names for each entry to extract a subset of the NCBI taxonomy with the NCBI taxtastic tool, thus producing a taxtastic file. The taxtastic file based on species from the NCBI genomes database was used to calculate the average genome sizes for higher taxonomic levels from the known genome size species level, with the assistance of the parent id and taxa id layout in the taxtastic file.

Using the taxtastic file based on our algae reference database, we assigned our algae entries a genome size from species to root from the prepared average genome sizes NCBI genomes taxtastic file. Not all genome sizes in the algae reference database were known. We, therefore, took the average of closely related species from the above taxonomic level of those we could get and took that as the genome size for those that were missing in our dataset.

As depicted in figure 3.2 box (d) the 18S rDNA dataset was normalised by assigning their genome sizes and applying the equation of the line. A normalisation procedure called the hits per million reads method was applied to the files, which entails scaling the files to a common value [Robinson and Oshlack, 2010].

File preparation for 18S rDNA analysis (box d)

In our 18S rDNA dataset, we had a total abundance of 53,750,176 taxa from the eukaryote node down to the species nodes. We employed MEGAN to cut out specific taxonomic levels. In MEGAN, we extracted the classifications at the taxonomic rank of species. This consisted of a file being generated for each station that contained the species names and their assigned abundances. The files were further normalised to hits per million.

In MEGAN, we extracted the leaves of the taxonomy tree at the rank of class and above but excluded assignments to the eukaryote node. Firstly, this consisted of a file being generated for each station that contained all assignments to the class nodes as well as any assignments under their respective lineages down to species being summed

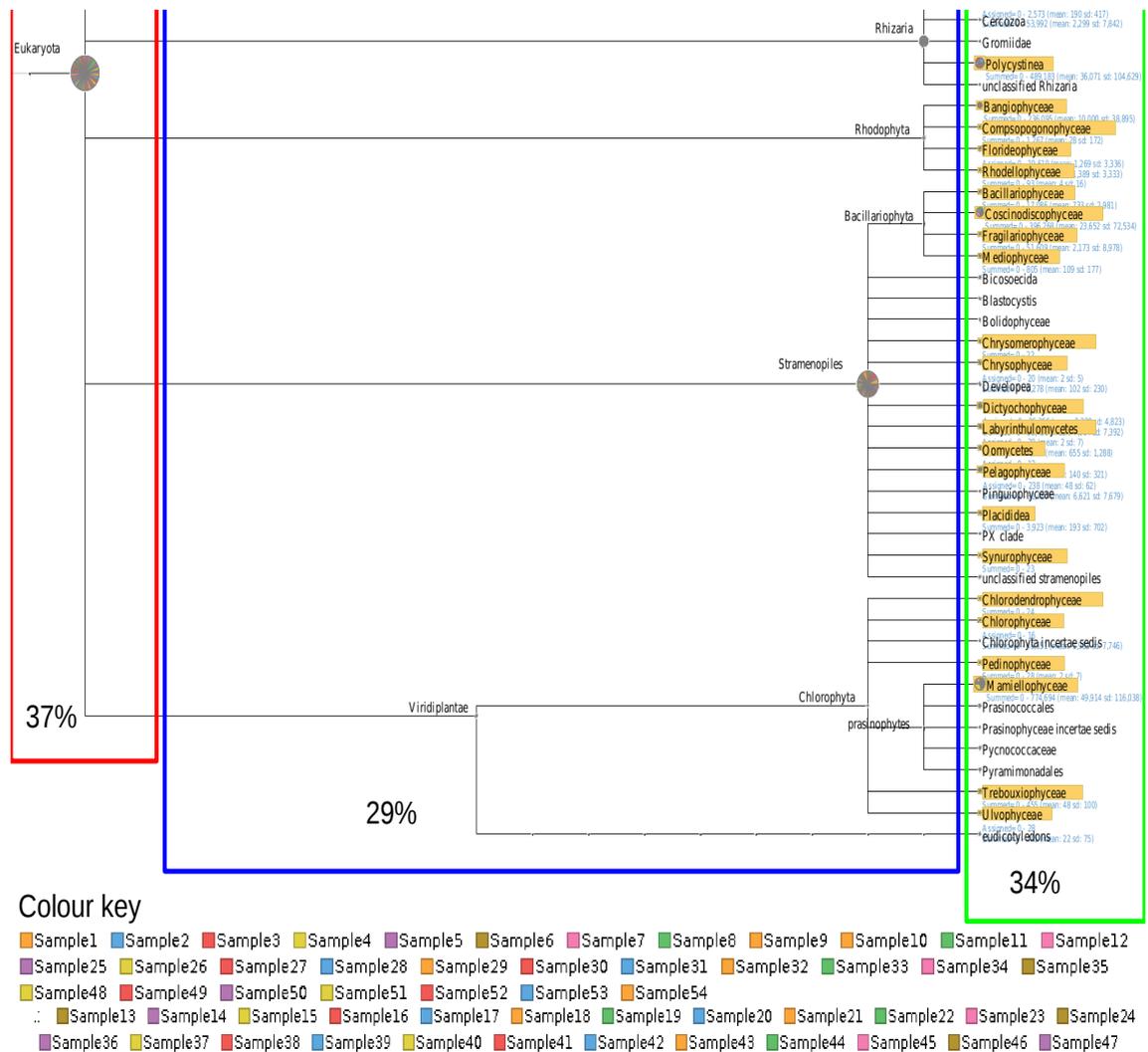


Figure 3.5: Part of the 18S rDNA dataset displayed on MEGAN’s taxonomy tree. All the nodes highlighted in yellow at the leaves of the tree are class nodes. Nodes at the leaves of the tree that are not highlighted do not have a taxonomic classification of class in their lineage. The nodes between the leaves of the tree and the eukaryotic node are the internal nodes that are representing higher taxonomic levels such as phylum. The colour key represent each colour on the nodes and corresponds to a sample in the 18S rDNA dataset

up under the individual class node. These are displayed in figure 3.5 as the highlighted nodes on the leaves of the tree.

Secondly, we included nodes that were not highlighted on the leaves of the tree, as displayed in figure 3.5. In NCBI taxonomy there are species that do not have a taxonomy designation at every taxonomy level. In our 18S rDNA dataset, we had species that do not have a taxonomic rank of class and these are displayed in figure 3.5 as the leaves of the tree that are not highlighted. We took the nodes that were not highlighted on leaves of the tree and summed them together within their respective lineages and placed them under a new name. For example, under the phylum Rhizaria, on the leaves of the tree, there is Cercozoa, Gromiidae and unclassified Rhizaria which

are not highlighted. Their abundance was summed together and renamed Nc.Rhizaria, “Nc.” standing for “No class”. The abundances assigned to Rhizaria were not included in this calculation. The leaves of the tree as displayed in figure 3.5, made up 34% of the total 18S rDNA dataset, as it resulted in an abundance of 18,332,601.

The internal nodes between the leaves of the tree and the eukaryote node as displayed in figure 3.5 was given a “U.” in front of their name, “U.” standing for “Unknown”. This was done to highlight that while they are of course associated with the lower lineages they are in fact considered separate, as those assignments to those nodes could not be determined any lower. The internal nodes made up 29% of the total 18S rDNA dataset, as it resulted in an abundance of 15,678,138.

The abundance assigned to the eukaryote node was excluded from our analysis as these sequences could not be classified lower. This comprised of a total abundance of 19,739,437, which is 37% of the 18S rDNA dataset. The assignments to the eukaryote node were excluded as these reads could not be classified to lower taxonomic ranks. Our algae reference database contains an up to date SILVA database at the time of our analysis, but still 37% of the 18S rDNA dataset is essentially unclassifiable. We have lost valuable information that could have potentially provided more insight and understanding into our analysis.

A file was generated for each station that contained the class nodes, “Nc.” nodes and “U.” nodes with their respective abundances. The files were further normalised to hits per million.

3.3.2 Computational pipeline for 16S rDNA analysis

In this section, we describe JGI’s computational pipeline for taxonomically classifying the 16S rDNA dataset. Also, we describe how we account for 16S rDNA copy number variation in order to give an estimate of abundance for the species in our dataset. An overview of JGI’s 16S rDNA classification pipeline is shown in figure 3.6 and a description of the pipeline is below.

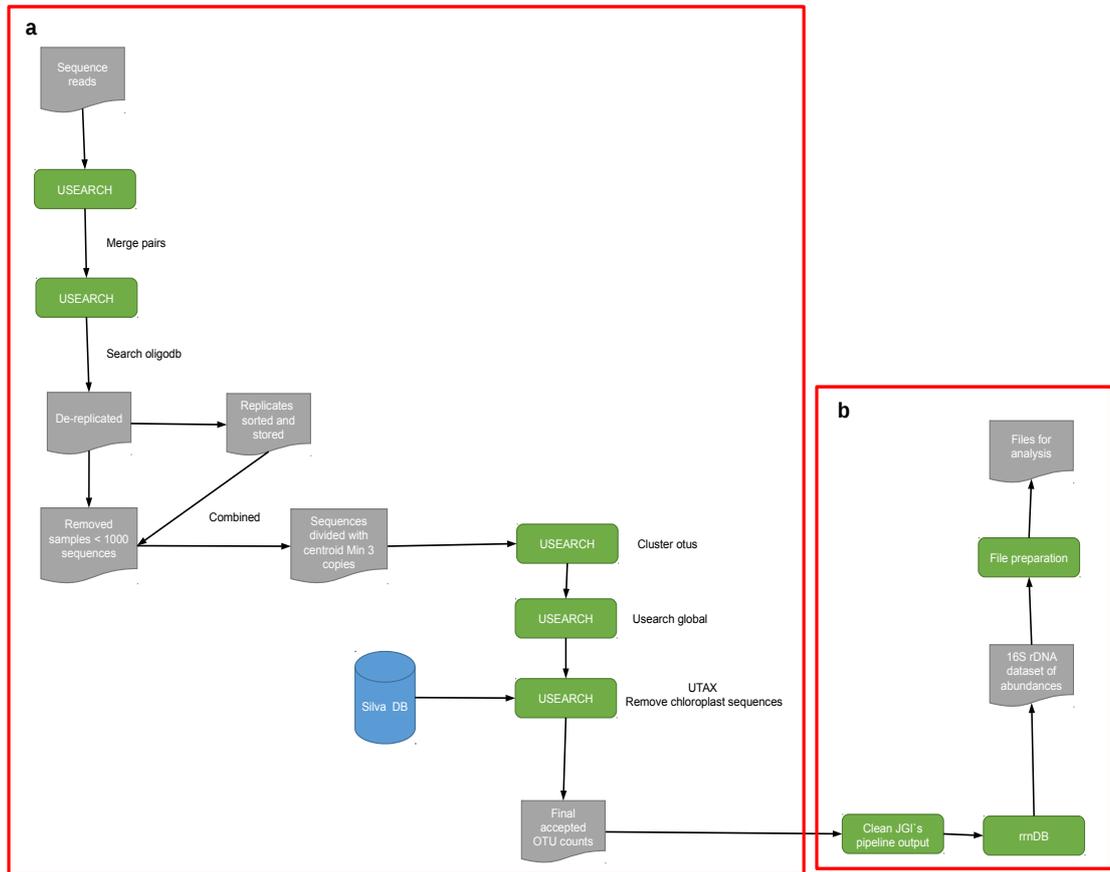


Figure 3.6: JGI's pipeline diagram of 16S rDNA classification analysis. The pipeline at various stages incorporates databases (blue), software tools (green) and processed files (grey). Box **a** and **b** refer to sections in the text

JGI's computational pipeline for taxonomical classifying the 16S rDNA dataset (box a)

The JGI's 16S rDNA classification pipeline consists of firstly removing samples with less than 1000 sequences. The remaining samples and the de-replicated identical sequences from the preprocessing step (as outlined in Section 3.2.2) are then combined and their sequences organized by decreasing abundance. The sequences are divided out based on the criterion as to whether they contained a cluster centroid with a minimum size of at least 3 copies. The low abundance sequences are put aside and not used for clustering. USEARCH's cluster otus command is employed to incrementally cluster the clusterable sequences. This begins at 99% identity and the radius is increased by 1% for each iteration until a OTU clustering identity of 97% is reached.

At each step, the sequences are sorted by decreasing abundance. Once clustering is complete, USEARCH's usearch global is used to map the low-abundance sequences to the cluster centroids. These are added to OTU counts if they were in the prescribed

percent identity threshold. If they do not fall within this prescribed percent identity threshold they are discarded. USEARCH's USTAX along with the SILVA database is used to evaluate the clustered centroid sequences. The predicted taxonomic classifications are then filtered with a cutoff of 0.5. Any chloroplast sequences identified are removed. The final accepted OTUs and read counts for each sample are finally placed in a taxonomic classification file.

Cleaning the JGI's pipeline output (box b)

The JGI output taxonomic classification file contained a matrix of the 57 station's names to their counts of read assignments. Also given were the taxonomic assignments in the form of the SILVA taxonomic path from domain to the taxonomic level it was assigned. Our SILVA taxonomic assignments had a number of issues that needed resolving before we could account for the 16S rDNA copy number variation in our datasets. SILVA contains entries without a formal taxonomic name and these are entered in SILVA's database as for example "unknown". When we had such an assignment we moved the assignment up a taxonomic level until a taxonomic assignment with a "real" taxonomic name was given.

In addition, SILVA taxonomic names are formatted slightly differently from NCBI entries, for example, some entries have numbers attached. We identified such entries in our dataset and edited the names. SILVA taxonomic names are also not updated with NCBI entries. We identified such entries in our dataset and edited the names to ensure they were up to date with the NCBI taxonomy. We created a custom made script to accomplish this, to tidy up SILVA taxonomic names of any format issues, to have "real" up to date taxonomic names for each assignment and to have a single taxonomic name rather than a taxonomic path for each assignment. We combined any duplicated taxonomic assignments for each station.

Normalising 16S rDNA copy number (box b)

In order to normalise the 16S copy number, the 16S copy numbers for the species in the dataset were retrieved from the Ribosomal RNA Operon Copy Number Database (rrnDB) [Klappenbach et al., 2001]. The rrnDB database consisted at the time of 3,021 bacterial entries. Firstly, since multiple entries of a species are in the rrnDB

database due to the presence of different strains, we obtained an average copy number for each species in the rrnDB database, which resulted in 2,876 species entries. The higher taxonomic levels for the rrnDB species needed to be established so that we could calculate their average copy number. For a description of the lineages of all species back to the root in the rrnDB database, we submitted the species names for each entry to extract a subset of the NCBI taxonomy with the NCBI taxtastic tool [NCBI Resource Coordinators, 2016], thus producing a Taxtastic file. The Taxtastic file based on species from the rrnDB database was used to calculate the average copy number for higher taxonomic levels from the known copy number species level, with the assistance of the parent id and taxa id layout in the Taxtastic file. A Taxtastic file based on 16S rDNA species from our dataset was generated and we assigned our 16S species entries a copy number from species to root from the prepared average copy number rrnDB Taxtastic file. Not all copy numbers in the 16S rDNA dataset were known. We therefore took the average of closely related species from the above taxonomic level of those we could get and took that as the copy number for those that were missing in our dataset. The 16S dataset was normalised by dividing by the assigned copy number. The files were normalised to hits per million.

File preparation for 16S rDNA analysis (box b)

In our 16S rDNA dataset, we had a total abundance of 56,999,957 taxonomic assignments to nodes from the bacteria node down to the genus nodes. We prepared the 16S rDNA taxonomic levels in the same manner as the 18S rDNA taxonomic levels as outlined in Section 3.3.1.

We extracted the leaves of the tree that include class nodes and “Nc.” nodes with their respective abundances. This step resulted in an abundance of 53,723,979 (94%). Also, we extracted the internal nodes and placed “U.” in front of their names. This resulted in an abundance of 1,627,260 (3%). The abundance assigned to the bacteria node was excluded from our analysis and this comprised of a total of 1,648,718 (3%). We generated a file for each station that contained the class nodes, “Nc.” nodes and “U.” nodes with their respective abundances. The files were further normalised by applying the hits per million reads method.

We extracted the classifications at the taxonomic rank of genus. This consisted of a

file being generated for each station that contained the genus names and their assigned abundances. The files were further normalised by applying the hits per million reads method.

3.3.3 Further analysis methods

In this section, we describe the methodology that we used for our analysis of the 18S and 16S rDNA datasets.

Evenness and occupancy

Alpha diversity determines the diversity of individuals in local communities [Marcon et al., 2014]. Indices are used to describe the general properties of a community and then this enables us to compare different regions or taxa [Morris et al., 2014]. There are a number of alpha diversity indices available, such as the Shannon diversity ($H' = - \sum P_i \ln(P_i)$), where P_i is the proportion of individuals belonging to species i) [MacArthur and MacArthur, 1961] which is sensitive to the different number of species present in a community [Morris et al., 2014], [Johnson and Burnet, 2016]. Taxonomy evenness is the similarity in relative abundance across the sample locations [Zhang et al., 2012]. An evenness value ranges between 1 and 0, with an evenness value of 1 corresponding to complete evenness and 0 no evenness. The occupancy refers to how many sample sites that species occurs in.

We produced an abundance, species evenness and occupancy plot for each 18S rDNA class level ($n=54$) and 16S rDNA class level ($n=57$). The x -axis represents the number of times that taxon occurs across the stations. The y -axis represents the evenness of that taxon across stations it occurs in. Each circle represents a taxon abundance. The size of each circle is resized by replacing the area of the circle which represented the total abundance for that taxon with the square root of the abundance divided by π . The evenness and occupancy plot was calculated using a Dispersion index, which is a variant of Pielou's evenness [Pielou, 1966], and based on Shannon diversity index [Payne et al., 2005].

There are multiple indices available to quantify biodiversity, but there is no consensus regarding which is more appropriate and informative [Morris et al., 2014]. For our

analysis, we choose the Shannon diversity index which is equally sensitive to rare and abundant species [Morris et al., 2014]. We also choose the Dispersion index which is also sensitive to samples containing rare species [Payne et al., 2005]. It is important for us to consider both abundant and rare species in our analysis. Rare species are defined as those with very low abundance and are not present in every sample [Chapman et al., 2018]. Our 18S and 16S rDNA datasets contain a considerable portion of rare species, $\sim 42\%$ and $\sim 40\%$ respectively.

The presence and abundances of some species in a particular location may not be independent of some other species [Schluter, 1984]. In this circumstance, there would be a potential need to adjust the Shannon index to account for non-independent species. We would then combine the non-independent species into modules instead of considering them separately. We would then apply the combined presence and abundance in the Shannon diversity index.

Breakpoint analysis

We performed a breakpoint analysis based on the methodology from [Castro-Insua et al., 2016]. This approach plots beta diversity against temperature. Beta diversity is a measure of the amount of variation in species composition for a community among samples [Ricotta, 2017]. Beta diversity enables us to compare and contrast communities. There are a number of different indices for beta diversity. The beta diversity indices that we use in our breakpoint analysis is the Sørensen indices ($\beta_{sor} = \frac{b+c}{2a+b+c} = \frac{b}{b+a} + (\frac{c-b}{2a+b+c}) (\frac{a}{b+a})$), where “a” is the number of species two sites have in common, “b” is the number of unique species in the poorest site and “c” is the number of unique species in the richest site [Baselga and Orme, 2012].

The breakpoint analysis enabled us to search for temperature breakpoints in the 18S and 16S rDNA datasets. This analysis provided insight about changes in the biodiversity of the identified prokaryotic and eukaryotic species across the different temperatures of the polar region in the Arctic Ocean, through the temperate region in the North Atlantic Ocean and into the tropical region in the South Atlantic Ocean.

A breakpoint was determined and plotted for each of the 18S rDNA class level ($n=54$) and 16S rRNA class level ($n=57$) datasets. This was calculated by firstly producing a presence absence matrix for each dataset. A multiple-site dissimilarity was

performed on the presence absence matrix with `beta.pair`, a function from the `betapart` R package and a dissimilarity index set to `sorensen`. These values were then plotted against temperature, to enable us to get a range of values in which the breakpoint might be located. We then searched through these possible breakpoints for the one with the lowest mean squared error. In each of the 18S rDNA and 16S rDNA datasets plots the y -axis represents the beta diversity and the x -axis represents temperature with piecewise regression lines and breakpoints shown.

Co-occurrence analysis

We also undertook a co-occurrence analysis with weighted Gene Co-Expression Network Analysis (WGCNA) [Langfelder et al., 2008], a method for finding modules (a WGCNA term used to describe networks) of highly correlated individuals. The prokaryotes at the taxonomic rank of genus and eukaryotes at the taxonomic rank of species normalised files were combined for each station ($n=50$). Using the R package WGCNA on samples of combined prokaryotes and eukaryotes we obtained modules derived from their log10-scaled abundances.

A network can be described by its adjacency matrix. The adjacency matrix (a_{ij}) is calculated by first defining a co-expression similarity ($s_{ij} = |\text{cor}(x_i, x_j)|$, where `cor` is correlation, and x_i and x_j are gene (species) expression profiles, consisting of expression of genes (species) i and j across a number of samples). We performed a signed co-expression measure to keep track of the sign of positively correlated genes (species) of the co-expression measure. Also, we transformed the co-expression similarity into a weighted adjacency (an adjacency that keeps track of the correlation values determined between genes (species)), as we wanted the adjacency to keep track of the connection strength (correlation coefficient value) between the genes (species) [Langfelder et al., 2008].

To construct a signed weighted network we first determined a soft threshold which is a power called beta ($\beta \geq 1$) [Langfelder et al., 2008]. The soft threshold was determined by using the `pickSoftThreshold` function, which is part of the WGCNA package, in R [Langfelder et al., 2008]. We use the power to raise the correlation calculation ($a_{ij} = (|1 + \text{cor}(x_i, x_j)|/2)^\beta$) so the network is constructed with the emphasis on high correlations at the expense of low correlations, therefore the network is more

robust [Zhang and Horvath, 2005], [Langfelder et al., 2008]. We choose the lowest power that results in approximate scale free topology as measured by the scale free topology fitting index [Zhang and Horvath, 2005]. For our data, we got a power of 11. Therefore a signed weighted adjacency measure for each pair of species was determined by raising the absolute value of the correlation coefficient to the power of 11. A topological overlap measure (TOM) was calculated from the resulting adjacency matrix. Hierarchical clustering was carried out on the TOM measure. This resulted in two modules being found ((a) $n=70$ and (b) $n=51$).

Highly correlated modules are not distinct, therefore using the modules' eigengene, which is the first principal component of a module, to be a representative of that module, in order to further examine the modules, to merge highly correlated (≥ 0.75) modules based on their eigengene [Langfelder and Horvath, 2007]. This did not result in the two modules being merged and therefore the two modules were taken for analysis. When incorporating environmental data, latitude values were redefined, so that the North pole is 0° , the Equator is 90° and the South pole is 180° .

3.4 Results

3.4.1 Rarefaction curves

Samples taken from a population should be a representative of that population, but we are unable to visually confirm the organisms we are attempting to capture. Rarefaction curves are a standard tool for analysis, in order to generate a rarefaction curve, the number of species is plotted as a function of the number of individuals sampled. Rarefaction curves are a means of determining if sequencing depth is sufficient [Wooley et al., 2010]. Also with rarefaction curves, we can compare the diversity among the samples [Hughes et al., 2001]. Initially, if the curve begins with a steep slope and then at some point, the curve begins to level off this indicates that fewer species are being discovered in the sample. If the slope increases more gently, this indicates less contribution of the sampling to the total number of species. Interpreting a comparison of the diversity among the samples in the rarefaction curves plot can be achieved by taking the highest shared sample size (x-axis) between the curves and examining the

curves position to each others in relation to the number of leaves in taxonomy (y -axis), the curves that are highest along the y -axis are more diverse [Wooley et al., 2010].

The rarefaction curves in figure 3.7a and 3.7b were generated using MEGAN and are based on the taxonomic classification of 18S and 16S rDNA species and genus level, respectively. In each rarefaction curve, the individual curves represent a single sample, as we ran rarefaction curves on each sample and then the curves are plotted together per dataset. The samples are coloured by region, which corresponds to the map of the region sampling sites in figure 3.1a, where the polar Arctic Ocean samples are coloured black, the temperate North Atlantic Ocean samples are coloured red and the tropical South Atlantic Ocean samples are coloured yellow.

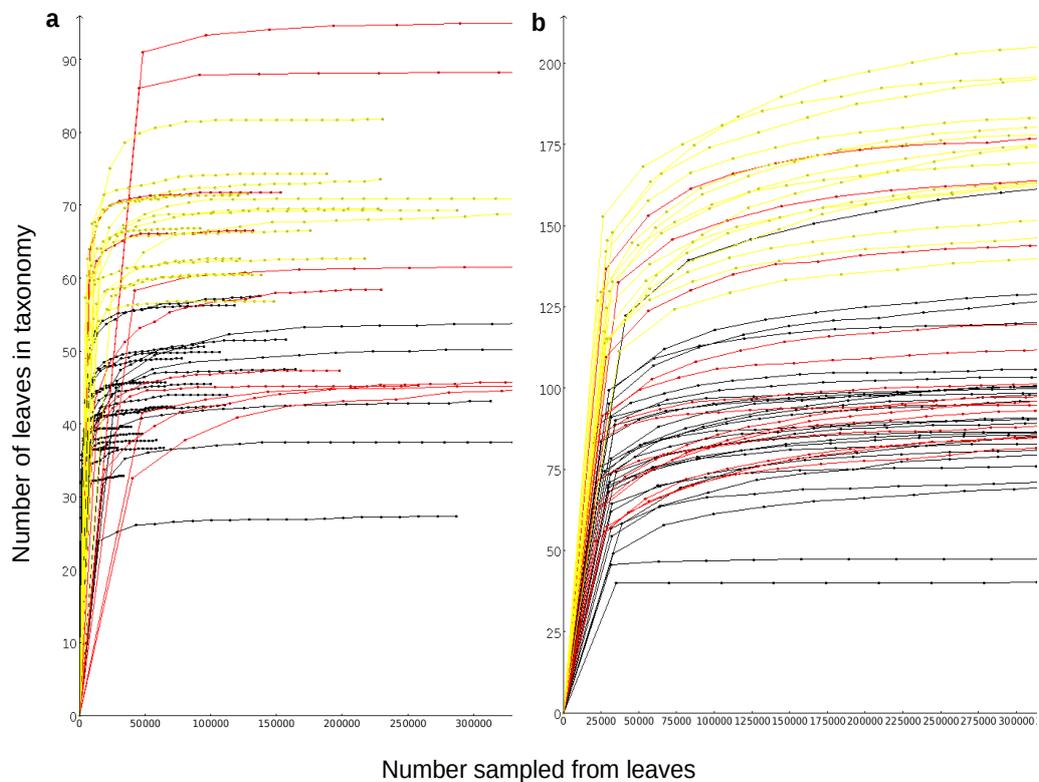


Figure 3.7: (a) rarefaction curves for 18S rDNA species level ($n=54$) and (b) the rarefaction curves for 16S rDNA genus level ($n=57$). In each panel the colours correspond to the sample sites of the three expeditions, April to May 2011 in red (North Atlantic Ocean), June to July 2012 in black (Arctic Ocean) and November to December 2012 in yellow (South Atlantic Ocean). The rarefaction curves were generated using MEGAN

As seen in figure 3.7a for the 18S rDNA species, there is a sharp rise at first in the curves of the samples from the Arctic Ocean coloured in black and the curves of the samples from the South Atlantic Ocean coloured in yellow. The curves of the samples from the North Atlantic Ocean coloured in red rise slightly less sharply at

first compared to the other samples from the Arctic Ocean and South Atlantic Ocean. Then all the curves begin to rise more slowly as rarer species are added and then the curves level off. The levelling off of all of the curves happens quite quickly, therefore we conclude that sufficient sampling was achieved during the three expeditions.

In figure 3.7a for the 18S rDNA species, the curves for the samples from the South Atlantic Ocean coloured in yellow contain the greatest amount of diversity. The curves for the samples from the Arctic Ocean and North Atlantic Ocean coloured in black and red, respectively, contain in general about an equal amount of diversity, with the exception of several red samples which have a higher amount of diversity.

For the 16S rDNA species as shown in figure 3.7b there is an even steeper rise at first in the curves of the samples from all three expeditions. Then all curves begin to rise more slowly as rarer species are added before the curves level off. The levelling off of the curves happens quite quickly. Therefore we conclude that we adequately sampled during the three expeditions.

In figure 3.7b for the 16S rDNA species, we see a similar pattern of diversity among the samples as we did in 3.7a for the 18S rDNA species. The curves for the samples from the South Atlantic Ocean coloured in yellow contain the greatest amount of diversity. The curves for the samples from the Arctic Ocean and North Atlantic Ocean coloured in black and red, respectively, contain in general about an equal amount of diversity, with the exception of a few red samples which have a higher amount of diversity.

3.4.2 Principal Coordinates Analysis

A clustering analysis was performed in MEGAN for each of our 18S and 16S rDNA datasets. Principal Coordinates Analysis (PCoA) is a multidimensional scaling technique which enables us to compare a large number of samples at once. PCoA attempts to order the objects across the axes of principal coordinates. PCoA accomplishes this by employing a linear transformation of the distance or dissimilarities between the objects onto the plot, while also seeking to explain as much of the variance in the initial dataset [Ramette, 2007], [Paliy and Shankar, 2016]. PCoA summarises the community compositional differences between the samples, the principal coordinates from a PCoA are plotted against each other and each point in the plot represents a sample. The

relative positioning of the point represents the relationships among the variables measured in the samples [Goodrich et al., 2014], [Paliy and Shankar, 2016]. The distance of the groups from one another cannot be quantified with a high degree of reliability as a reflection of the dataset. This is due to the fact that PCoA summarises the dataset in a two dimensional scatterplot, the two principal coordinates used in the plot only explain a fraction of the variability in the dataset [Goodrich et al., 2014]. The groups within the PCoA can be determined by using betadisper, a function from the vegan R package, which calculates how compact the objects are by calculating the average distance of group members to the centroid of that group [Simpson, 2006].

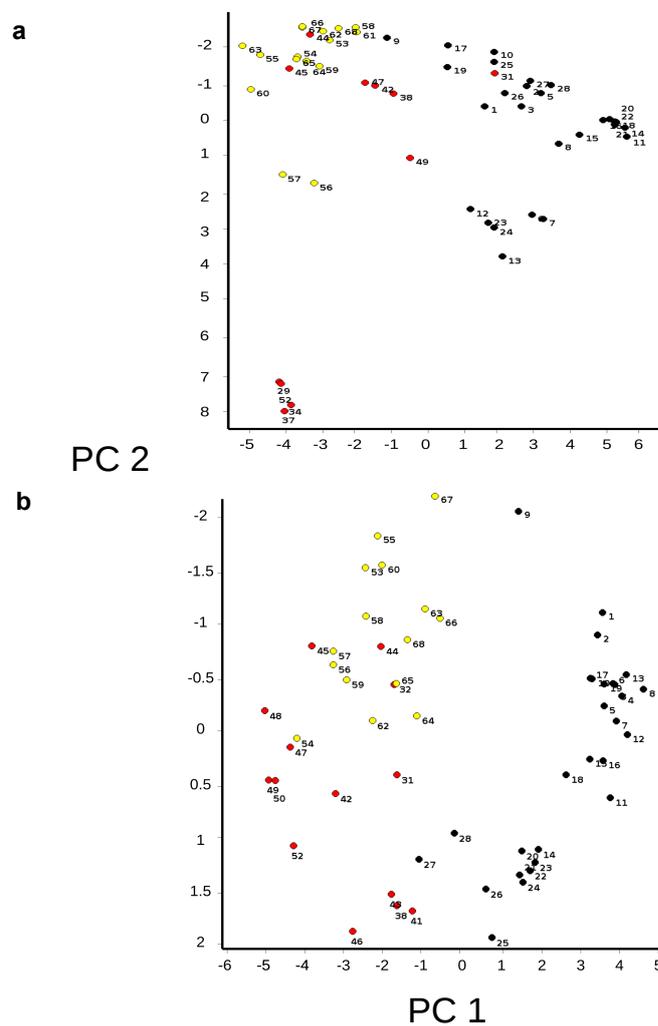


Figure 3.8: Principal coordinates analysis (PCoA) of **a** represents the eukaryotic communities ($n=54$) and **b**, represents the prokaryotic communities ($n=57$) at the taxonomic rank of class. In MEGAN, communities are clustered according to their similarity based on Bray-Curtis distances. The samples are numbered and coloured by region; these numbers and colours correspond to the map of the region sampling sites in figure 3.1a, where the Arctic Ocean samples are coloured black, the North Atlantic Ocean samples are coloured red and the South Atlantic Ocean samples are coloured yellow

Any distance or dissimilarity matrix can be used in a PCoA [Paliy and Shankar, 2016]. Multivariate analyses of ecological data are often based on a dissimilarity matrix, such as Bray Curtis [Anderson and Santana-Garcon, 2015]. Bray Curtis accounts for both species presence and abundances at each site [Pos et al., 2014]. A distance matrix based on Bray Curtis ($\sum |x_i - x_j| / \sum (x_i + x_j)$, where x are the counts of species in samples i and j) was used to calculate the distance between samples to be plotted [Ricotta and Podani, 2017].

Displayed in figure 3.8a is the 18S rDNA ($n=54$) dataset at the taxonomic rank of class along with higher level taxonomic assignments but excluding those assigned to the eukaryote node. Displayed in figure 3.8b is the 16S rDNA ($n=57$) dataset at the taxonomic rank of class along with higher level taxonomic assignments but excluding those assigned to the prokaryotes node. In each PCoA plot the samples are numbered and coloured by region; these numbers and colours correspond to the map of the region sampling sites in figure 3.1a, where the Arctic Ocean samples are coloured black, the North Atlantic Ocean samples are coloured red and the South Atlantic Ocean samples are coloured yellow.

The 18S rDNA are presented in figure 3.8a with PC1 accounting for 24% of sample variation, while PC2 accounts for 19.9% of sample variation. Overall the 18S rDNA samples are to a reasonable extent clustering well by region, as the matching colours are grouped together. Also for the 18S rDNA samples, there is a transition of the samples from black to red to yellow, which coincides with how the samples are positioned by latitude as can be seen in figure 3.1a. There are 4 samples that are outliers, these are samples 29, 52, 34 and 37 which are shown in figure 3.8a at the bottom left side of the plot. These samples came from the North Atlantic Ocean expedition as explained in section 3.2.1. Dr. Willem van de Poll and Dr. Klaas Timmermans noted that at the time of sampling a bloom may have been occurring, which could explain why these samples are clustering differently since their composition and abundance is remarkably different from the other samples in that region of the North Atlantic Ocean.

The 16S rDNA are shown in figure 3.8b with PC1 accounting for 51.6% of sample variation, while PC2 accounts for 26.3% of sample variation. The 16S rDNA samples are also to a reasonable extent clustering well by region, as the matching colours are grouping together. The 16S rDNA samples transition from black to red to yellow in

the shape of a horseshoe, this shape is indicative of an underlying linear gradient. Also this transition of colour coincides with how the samples are positioned by latitude as can be seen in figure 3.1a.

3.4.3 Heatmaps

Figure 3.9a and 3.9b display heatmaps arranged by latitude for the 18S and 16S rDNA datasets, respectively. The numbers at the bottom of the heatmap correspond to sample site numbers as shown in figure 3.1a. The heatmaps are arranged by placing the most abundant taxa at the top of the plot down to the least abundant at the bottom. Heatmaps enabled us to overview the distribution, composition and abundance of our datasets. The heatmaps were generated on log₁₀-scaled abundances of the 18S and 16S rDNA datasets, using the heatmap.2 function, which is part of the gplots package, in R.

In figure 3.9a is the 18S rDNA at the taxonomic rank of class along with higher level taxonomic assignments but excluding those assigned to the eukaryote node. The most abundant and constant eukaryotic taxon across the samples is the U.Stramenopiles. The U.Stramenopiles represent Stramenopiles species that could not be identified to the taxonomic level of class or below. Stramenopiles are a very diverse group of marine microbes and have been found inhabiting the oceans of the world, both in open ocean and coastal areas [Lin et al., 2012]. Therefore this is not a surprising result, that Stramenopiles are the most abundant and constant group throughout our samples. Also for example, in figure 3.9a, Dinophyceae abundance increases from the polar regions of the Arctic Ocean as we move down into the tropical region of the South Atlantic Ocean. This result is as expected, as Dinophyceae are found in the polar and tropical regions, with a higher abundance found in the tropical regions [Okolodkov and Dodge, 1996], [Taylor et al., 2007].

The taxonomy entries at the bottom of the heatmap from around the entry U.Rhodophyta and below are the ones that are driving diversity. There are a greater number of taxonomy entries in samples from around 29 to 68 which are located in the tropical regions compared to those samples around 1 to 28 located in the colder regions, which corresponds to what one would expect as the tropical regions are more diverse than the

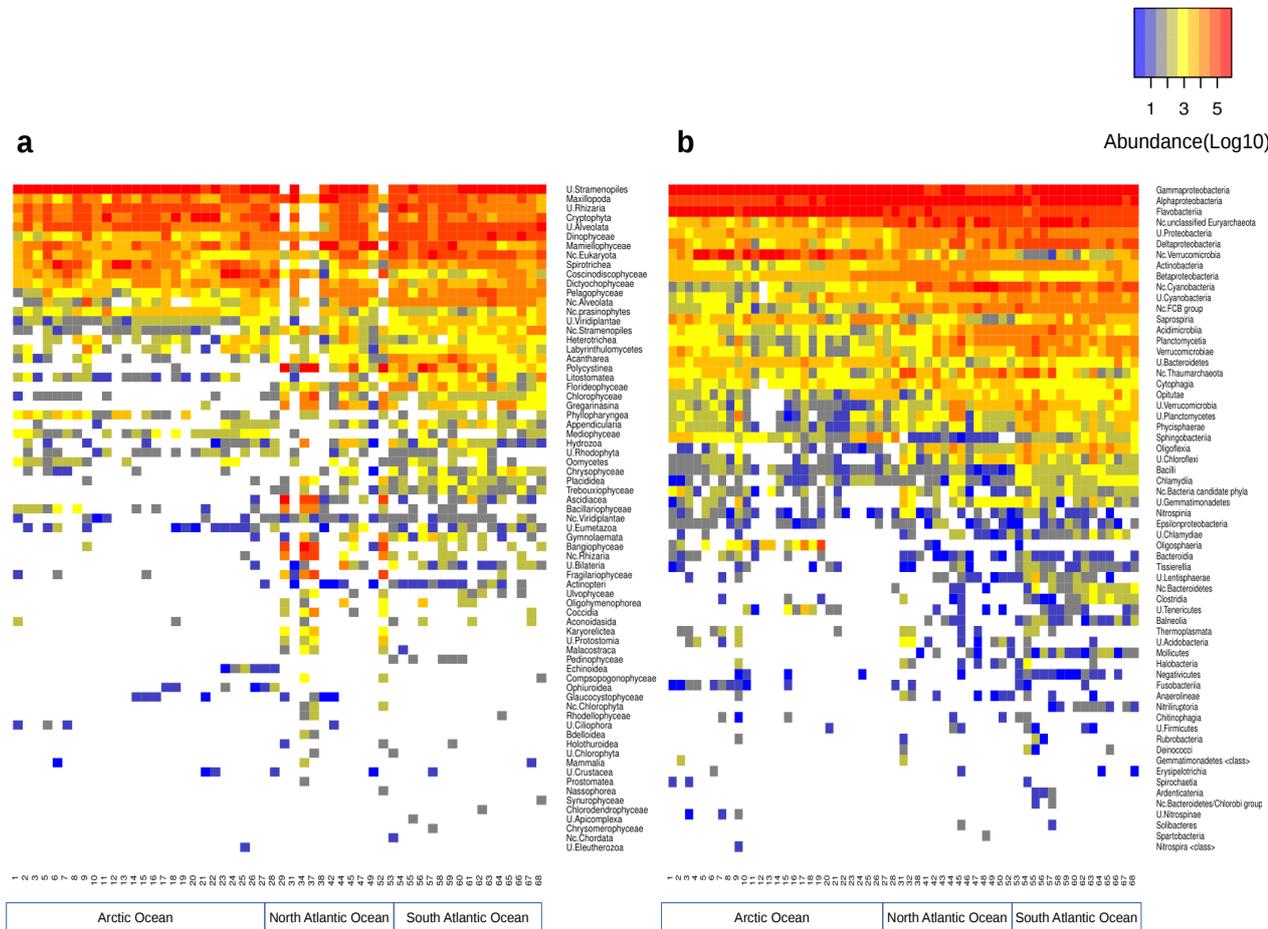


Figure 3.9: Panel **a** and **b** represents heatmaps of abundances arranged by latitude to the taxonomic rank of class, **a** represents the 18S rDNA taxonomy and **b** represents the 16S rDNA taxonomy. The taxonomy names in the dataset are displayed along the right side. The numbers at the bottom correspond to sample locations as shown in figure 3.1 **a**. The three regions of the Arctic Ocean, North Atlantic Ocean and South Atlantic Ocean are displayed underneath their corresponding sample numbers. The colours correspond to log₁₀-scaled abundances of the 18S and 16S rDNA, where red colours are high values and blue colours are low values. The heatmaps were generated using the heatmap.2 function, which is part of the gplots package, in R

colder regions.

In our PCoA plots in figure 3.8a at the bottom left side of the plot we identified 4 samples that are outliers, these are samples 29, 52, 34 and 37. In the heatmap in figure 3.9a it can be seen clearly how these samples' composition and abundance are different from those of the surrounding samples. These samples came from the North Atlantic Ocean expedition as explained in section 3.2.1. As noted before, Dr.Willem van de Poll and Dr.Klaas Timmermans noted that at the time of sampling a bloom was occurring, which could explain why these samples are clustering differently, since their composition and abundance are remarkably different from the other samples in

that region of the North Atlantic Ocean.

Displayed in figure 3.9**b** is the 16S rDNA at the taxonomic rank of class along with higher level taxonomic assignments but excluding those assigned to the bacteria node. Likewise, for 16S rDNA, the top entries are the most abundant and constant throughout our samples going down to our least abundant and constant. The most abundant and constant taxon across our 16S rDNA samples is the Gammaproteobacteria. The Gammaproteobacteria are a large class of bacteria and can be found throughout the world's oceans [Williams et al., 2010], [Franco et al., 2017]. Therefore it is somewhat reassuring that Gammaproteobacteria are the most abundant and constant group throughout our 16S rDNA samples.

Also, the taxonomy entries at the bottom of the heatmap from around the entry U.Chlamydiae and below are the ones that are driving diversity. There are a greater number of taxonomy entries in samples from around 1 to 13 and 48 to 68 which are located in the polar and tropical regions, respectively. The pattern appears “n” shaped because diversity is decreasing in taxonomy entries from the polar regions as we pass into the temperate regions and then increasing as we pass into the tropical regions. This pattern is different to what is shown in the 18S rDNA heatmap in figure 3.9**a**.

3.4.4 Evenness and occupancy

For each of the species in our 18S and 16S rDNA datasets an evenness score (J) was calculated ($J = H' / \log(\text{number of species})$, where H' is the Shannon diversity index) [Morris et al., 2014]. The occupancy refers to how many sample sites that species occurs in. A description of the methodology for the evenness and occupancy plots can be found in Section 3.3.3.

In figure 3.10**a** and 3.10**b**, we present an evenness and occupancy plot of the 18S and 16S rDNA datasets, respectively. In figure 3.10**a** is the 18S rDNA at the taxonomic rank of class along with higher level taxonomic assignments but excluding those assigned to the eukaryote node. Displayed in figure 3.10**b** is the 16S rDNA at the taxonomic rank of class along with higher level taxonomic assignments but excluding those assigned to the prokaryotes node. The numbers in the figure 3.10**a** and 3.10**b** correspond to taxonomy names which can be found in the Appendix A.C.

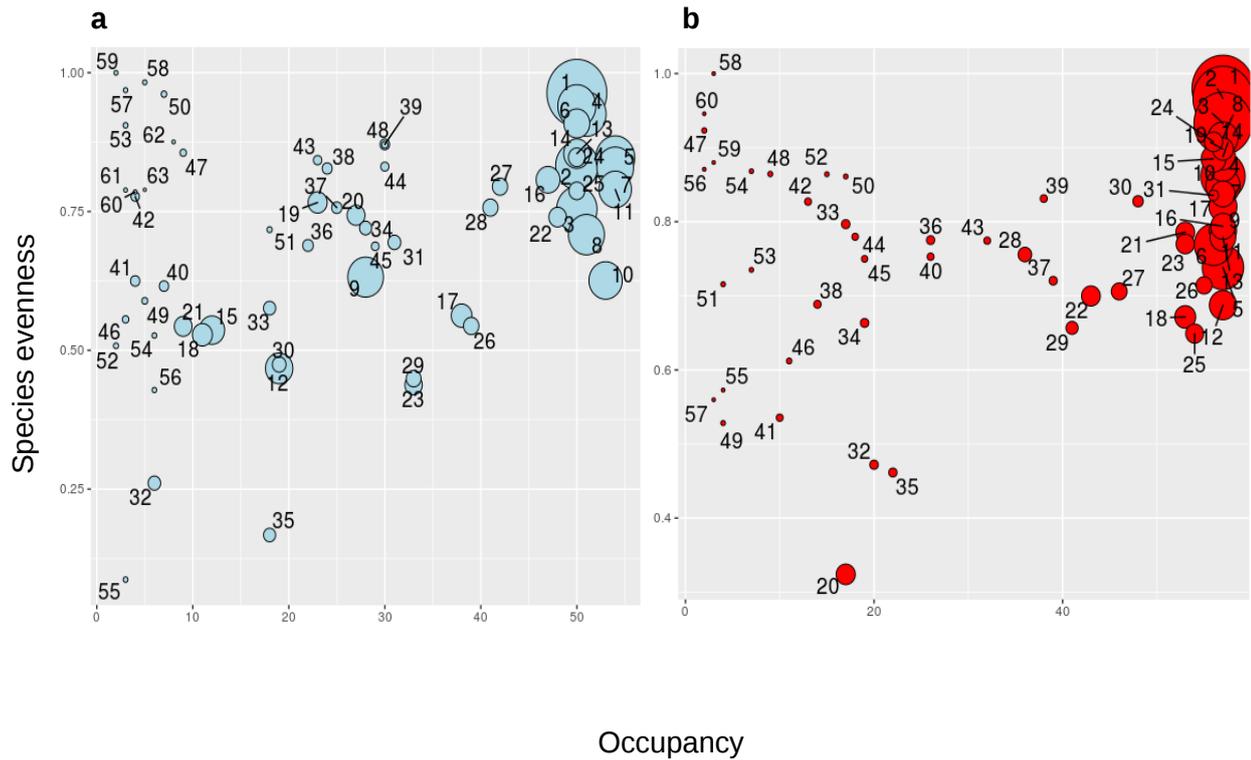


Figure 3.10: Panels **a** and **b** represent abundance taxonomy evenness and occupancy plots for the 18S and 16S rDNA datasets, respectively. The numbers in the plots correspond to taxon names which can be found in the Appendix A.C. The x -axis represents the number of times that class taxonomy occurs across the stations. The y -axis represents the evenness of that class taxonomy across stations it occurs in. Each circle represents a class taxonomy abundance. The size of each circle corresponds to the total abundance for that class, calculated by taking the square root of the abundance divided by π

The class evenness occupancy plot displayed in figure 3.10a for 18S rDNA shows that the more than 50% of the 18S rDNA classes have an evenness value between 0.5 to 0.85. This result indicates that the majority of the 18S rDNA classes are moderately to constantly present throughout our samples. A gradient pattern of high to low abundance and occupancy is observed of the 18S rDNA classes. The largest abundance is indicated by the greatest size of the circles and is located on the upper right hand corner of the plot. As the occupancy goes from high on the right hand side of the plot to low on the left hand side, the sizes of the circles decrease, indicating abundance also decreases. The Stramenopiles (1), Cryptophyta (2) and Mamiellophyceae (3) are the most abundant and widespread classes. From the 18S rDNA heatmap in figure 3.10a,

Stramenopiles are shown to dominate communities in the polar and temperate regions and decline only slightly in the tropical region. To a lesser degree, Cryptophyta also dominate communities in the polar and temperate regions and decline in the tropical region. Whereas Mamiellophyceae are more abundant in communities from the tropical region than in the polar and temperate regions.

The class evenness occupancy plot displayed in figure 3.10**b** of 16S rDNA shows the vast majority of the 16S rDNA taxonomy have an evenness value around 0.6 to 1. Therefore as we observed with the 18S rDNA species, the majority of the 16S rDNA species are also moderately to constantly present throughout our samples. Also, a gradient pattern of high to low abundance and occupancy is observed for the 16S rDNA species. As the occupancy goes from high on the right hand side of the plot to low on the left hand side, the sizes of the circles decrease, indicating abundance also decreases. The Gammaproteobacteria (1), Alphaproteobacteria (2) and Flavobacteriia (3) are the most abundant and widespread. From the 16S rDNA heatmap in figure 3.10**b**, Gammaproteobacteria are shown to dominate communities in the polar and temperate regions and decline very slightly in the tropical region. To a lesser degree, Flavobacteriia also dominate communities in the polar and temperate regions and decline in the tropical region. Whereas Alphaproteobacteria are more abundant in communities from the tropical and temperate regions than in the polar region.

Note that in panel **a** of figure 3.10 the 18S rDNA ($n=54$) for Prostomatea, Nasophorea, Synurophyceae, Chlorodendrophyceae, Apicomplexa, Chrysomerophyceae, Cephalochordata and Eleutherozoa are excluded from the analysis due to insufficient data to calculate the Shannon index. Also for 16S rDNA ($n=57$) in panel **b** Spartobacteria and Nitrospira were excluded from the analysis due to insufficient data to calculate the Shannon index.

3.4.5 Environmental plots

In this section, we present environmental plots. These plots were generated in collaboration with Dr.Schmidt.

Correlation heatmaps

In figure 3.11 panel **a** is the correlation heatmap for the Arctic Ocean samples, in panel **b** is the correlation heatmap for the North Atlantic Ocean samples and in panel **c** is the correlation heatmap for the South Atlantic Ocean samples. The correlation heatmaps enabled us to understand the statistical relationship of each taxon to the environmental variables. We produced the correlation heatmap with hierarchical clustering dendrograms on log₁₀-scaled abundances of the 18S and 16S rDNA using the `cor` function, which is part of the WGCNA package and `heatmap.2` function, which is part of the `gplots` package, in R. Displayed in figure 3.11 is the 18S rDNA at the taxonomic rank of class along with higher-level taxonomic assignments but excluding those assigned to the eukaryote node. Displayed in figure 3.12 is the 16S rDNA at the taxonomic rank of class along with higher-level taxonomic assignments but excluding those assigned to the prokaryotes node.

In figure 3.11**a**, which represents the correlation heatmap for the 18S rDNA of the Arctic Ocean samples, we observe that under temperature, salinity, longitude and latitude the heatmap is divided into two parts. We also can see according to the dendrogram on the left-hand side of the heatmap that the taxa form two large clusters. The cluster of taxa represented on the top part of the correlation heatmap is predominantly taxa that are positively correlating with temperature, salinity, longitude and latitude, while the cluster of taxa represented in the bottom half of the correlation heatmap is predominantly taxa that are negatively correlating with temperature, salinity, longitude and latitude. Under phosphate and silicate, there is no discernible pattern that can be determined for either of the clusters of taxa.

In figure 3.11**b**, which represents the correlation heatmap for the 18S rDNA of the North Atlantic Ocean samples, we can see according to the dendrogram on the left-hand side of the heatmap that the taxa form three clusters. The correlation heatmap is roughly divided into eight blocks of positively correlating and negatively correlating variables. The cluster of taxa represented on the top part of the correlation heatmap is predominantly taxa that are positively correlating with temperature, salinity, longitude and latitude, and negatively correlating with phosphate and silicate. The largest cluster of taxa represented in the middle part of the correlation heatmap is formed by two

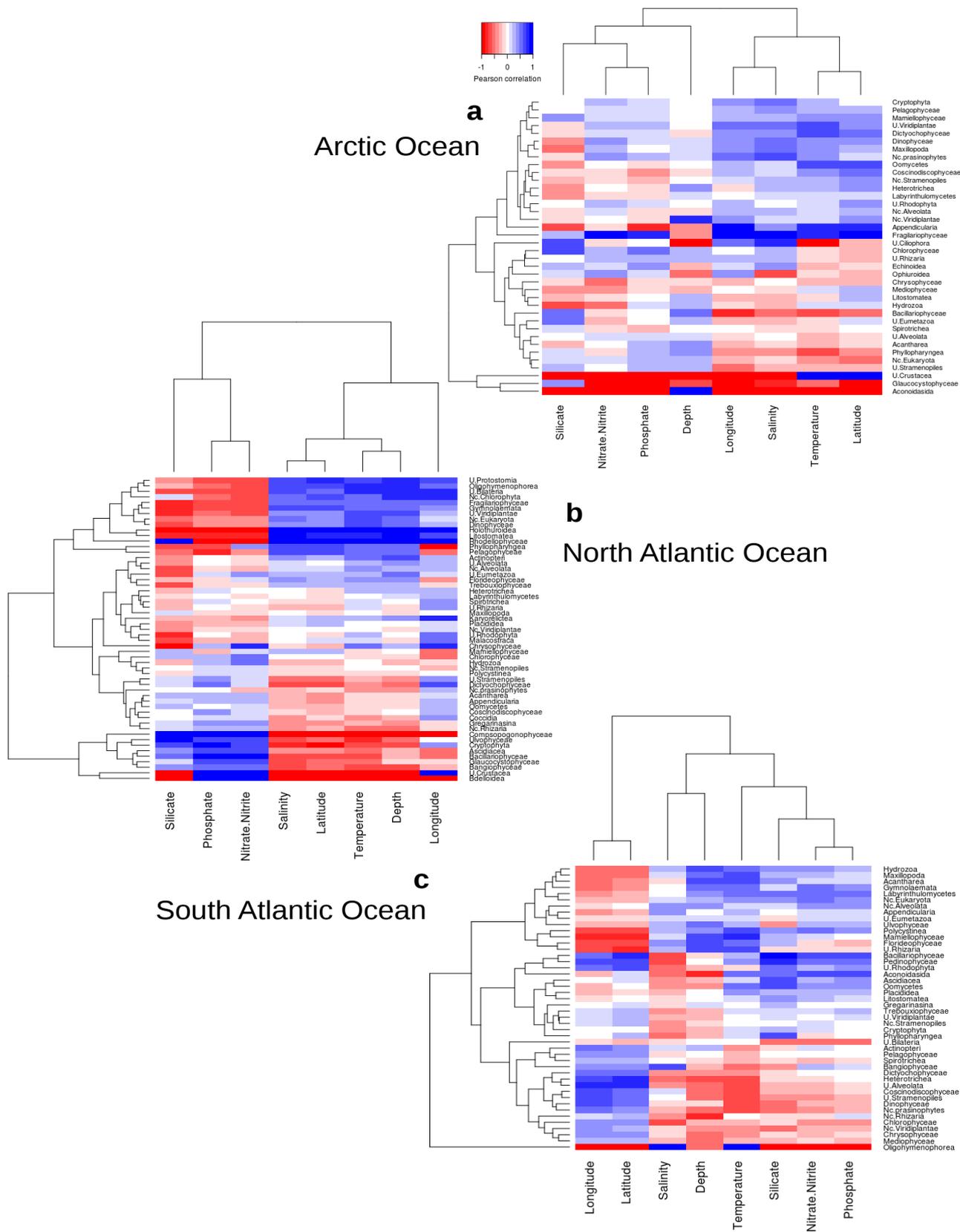


Figure 3.11: Correlation heatmaps of identified taxonomic classes and environmental data, **a** represents the Arctic 18S rDNA classes, **b** represents the North Atlantic 18S rDNA classes and **c** represents the South Atlantic 18S rDNA classes. The class names in the dataset are displayed along the right side and hierarchical clustering dendrogram on the opposite side. The environmental parameters are displayed at the bottom and hierarchical clustering dendrogram on the opposite side. The colours correspond to the Pearson correlation coefficient, where blue indicates a positive and red a negative correlation. The grey colour corresponds to no results, due to insufficient abundance data across the samples. The actual coefficients and p-values are given in Appendix A.D

smaller clusters, each with their own preference for the environmental variables. For the cluster of taxa at the top of the middle cluster of the correlation heatmap, the majority of the taxa are correlating positively with temperature, salinity, longitude and latitude, and negatively correlating with phosphate and silicate. For the cluster of taxa at the bottom of the middle cluster of the correlation heatmap, this consists predominantly of taxa that are negatively correlating with temperature, salinity, longitude and latitude, and positively correlating with phosphate and silicate. The cluster of taxa represented on the bottom part of the correlation heatmap is predominantly taxa that are negatively correlating with temperature, salinity, longitude and latitude, and positively correlating with phosphate and silicate.

In figure 3.11c, which represents the correlation heatmap for the 18S rDNA of the South Atlantic Ocean samples, we can see according to the dendrogram on the left-hand side of the heatmap that the taxa form three large clusters plus a singleton. The cluster of taxa represented on the top part of the correlation heatmap is predominantly taxa that are positively correlating with temperature, salinity, phosphate and silicate, while negatively correlating with longitude and latitude. The cluster of taxa represented in the middle of the correlation heatmap is predominantly taxa that are positively correlating with temperature, phosphate and silicate, longitude and latitude, while negatively correlating with salinity. The cluster of taxa represented at the bottom of the correlation heatmap is predominantly taxa that are negatively correlating with temperature, salinity, phosphate and silicate, while positively correlating with longitude and latitude. The singleton at the very bottom of the correlation heatmap is correlating negatively with phosphate, silicate, longitude and latitude, while correlating positively with salinity and temperature.

Specifically in the example in figure 3.11a, Chlorophyceae have a slightly negative correlation relationship with temperature and a slightly positive correlation relationship with salinity. In figure 3.11b, Chlorophyceae have a slightly negative correlation relationship with temperature and no correlation relationship with salinity. In figure 3.11c, Chlorophyceae have a negative correlation relationship with temperature and salinity. These results are not what is observed elsewhere in our analysis, for example, Chlorophyceae are more abundant in a tropical region, according to the 18S rDNA heatmap in figure 3.9a. It has been observed that species belonging to the group Chlorophyta,

of which Chlorophyceae is a member, can live even in extremely hot environments such as hot springs where temperatures can reach as high as 60°C [Mezhoud et al., 2014]. Therefore other relationships, potentially even positive correlations with temperature and salinity, would be possible in regions such as the South Atlantic Ocean.

In figure 3.12a, which represents the correlation heatmap for the 16S rDNA of the Arctic Ocean samples, we observe that under temperature, salinity, longitude and latitude the heatmap is divided into two parts. We also can see according to the dendrogram on the left-hand side of the heatmap that the taxa form two large clusters. The cluster of taxa represented on the top part of the correlation heatmap is predominantly taxa that are positively correlating with temperature, salinity, longitude and latitude, while the cluster of taxa represented in the bottom half of the correlation heatmap is predominantly taxa that are negatively correlating with temperature, salinity, longitude and latitude. Under phosphate and silicate, there is no discernible pattern that can be determined for either of the clusters of taxa.

In figure 3.12b, which represents the correlation heatmap for the 16S rDNA of the north Atlantic Ocean samples, we observe that the correlation heatmap is divided into four parts. We also can see according to the dendrogram on the left-hand side of the heatmap that the taxa form two large clusters. The cluster of taxa represented on the top part of the correlation heatmap is predominantly taxa that are positively correlating with temperature, salinity and latitude, while negatively correlating with phosphate, silicate and longitude. The cluster of taxa represented in the bottom half of the correlation heatmap consists predominantly of taxa that are negatively correlating with temperature, salinity and latitude, while positively correlating with phosphate, longitude and silicate.

In figure 3.12c, which represents the correlation heatmap for the 16S rDNA of the South Atlantic Ocean samples, we can see according to the dendrogram on the left-hand side of the heatmap that the taxa form five clusters, three small clusters and two larger clusters. The two small clusters of taxa represented on the top part of the correlation heatmap, are predominantly positively correlating with latitude, salinity and temperature, while negatively correlating with longitude, silicate and phosphate. The larger cluster of taxa represented in the middle part of the correlation heatmap, are predominantly negatively correlating with latitude, salinity and temperature, while

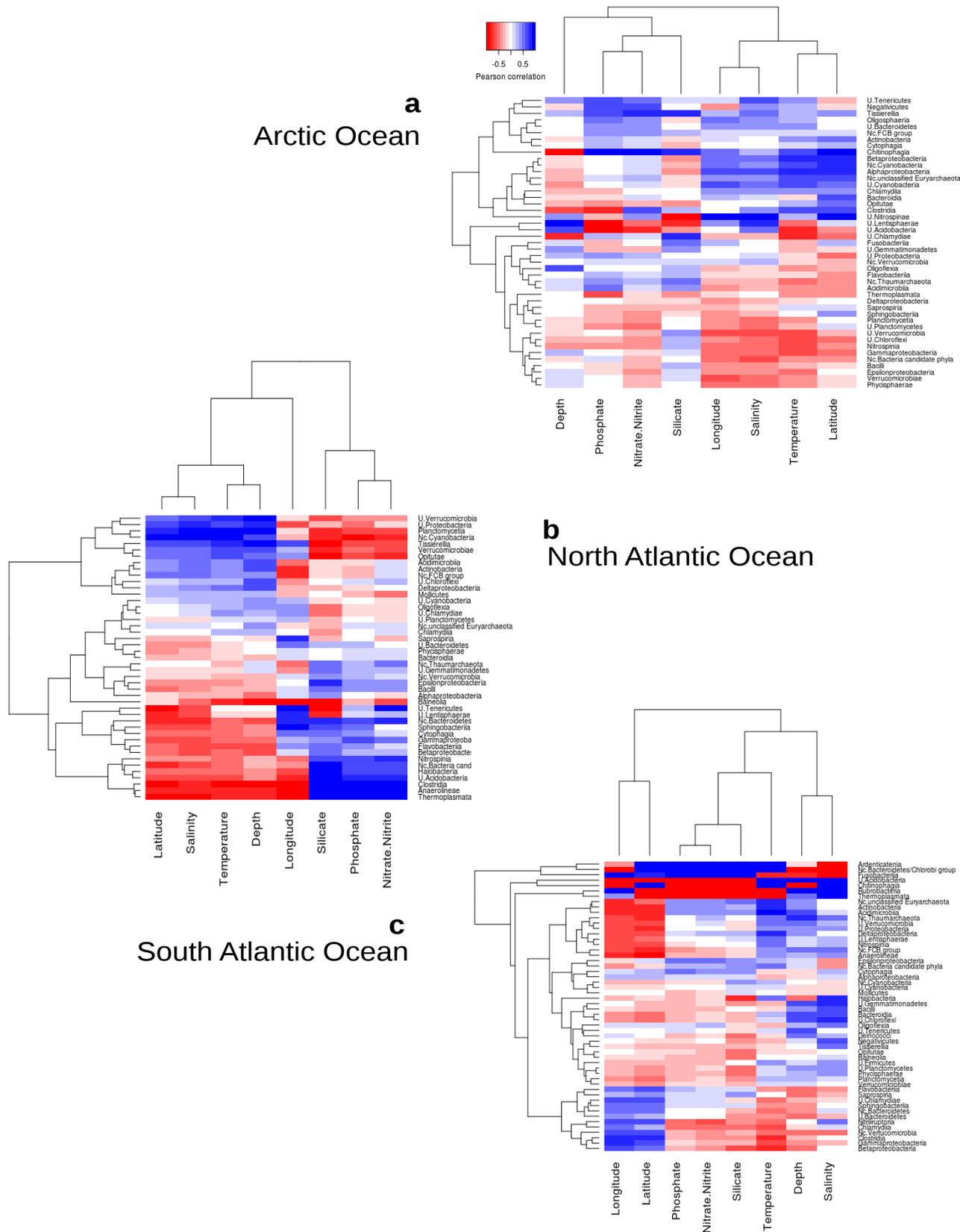


Figure 3.12: Correlation heatmaps of identified taxonomic classes and environmental data, **a** represents the Arctic 16S rDNA classes, **b** represents the North Atlantic 16S rDNA classes and **c** represents the South Atlantic 16S rDNA classes. The class names in the dataset are displayed along the right side and hierarchical clustering dendrogram on the opposite side. The environmental parameters are displayed at the bottom and hierarchical clustering dendrogram on the opposite side. The colours correspond to the Pearson correlation coefficient, where blue indicates a positive and red a negative correlation. The grey colour corresponds to no results, due to insufficient abundance data across the samples. The actual coefficients and p-values are given in Appendix A.D

positively correlating with longitude, silicate and phosphate. The two small clusters of taxa represented on the bottom part of the correlation heatmap, are predominantly negatively correlating with latitude, salinity and temperature, while negatively correlating with longitude, silicate and phosphate.

Specifically for the example in figure 3.12**a**, **b** and **c**, Nc.Cyanobacteria have a positive correlation relationship with temperature and salinity. Cyanobacteria are more abundant in the tropical regions, according to the 16S rDNA heatmap in figure 3.9**a**. Cyanobacteria can be found throughout the oceans of the world with a higher abundance for some of the species under Cyanobacteria in the tropical regions of the world [Flombaum et al., 2013].

Note that in figure 3.11 and figure 3.12 a number of the taxa are excluded from the analysis due to insufficient data to calculate the correlation coefficient values. In figure 3.11 panel **a** which represents the 18S rDNA from the Arctic Ocean the taxa excluded are Actinopteri, Ascidiacea, Bangiophyceae, Gregarinasina, Gymnolae-mata, Mammalia, Nc.Rhizaria, U.Bilateria and U.Eleutherozoa. In figure 3.11 panel **b** which represents the 18S rDNA from the North Atlantic Ocean the taxa excluded are Aconoidasida, Echinoidea, Mammalia, Mediophyceae, Nassophorea, Ophiuroidea, Prostomatea, U.Chlorophyta and U.Ciliophora. In figure 3.11 panel **c** which represents the 18S rDNA from the South Atlantic Ocean the taxa excluded are Chloroden-drophyceae, Chrysomerophyceae, Coccidia, Compsopogonophyceae, Fragilariophyceae, Holothuroidea, Malacostraca, Mammalia, Nc.Chordata, Rhodellophyceae, Synurophyceae, U.Apicomplexa and U.Chlorophyta. In figure 3.12 panel **a** which represents the 16S rDNA from the Arctic Ocean the taxa excluded are Anaerolineae, Erysipelotrichia, Gemmatimonadetes, Halobacteria, Nitrospira, Rubrobacteria and U.Firmicutes. In figure 3.12 panel **b** which represents the 16S rDNA from the North Atlantic Ocean the taxa excluded are Deinococci, Erysipelotrichia, Fusobacteriia, Gemmatimonadetes, Nega-tivicutes, Oligosphaeria, Rubrobacteria, Solibacteres, Spartobacteria and U.Firmicutes. In figure 3.12 panel **c** which represents the 16S rDNA from the South Atlantic Ocean the taxa excluded are Erysipelotrichia, Oligosphaeria, Solibacteres and Spirochaetia.

Non-metric multidimensional scaling

In figure 3.13a and 3.13b we present the non-metric multidimensional scaling (NMDS) plots for the 18S and 16S rDNA datasets, respectively, with environmental factors fitted. The numbers correspond to sample locations as displayed in figure 3.1a. NMDS plots allow us to visualise the distance matrix that is based on Bray-Curtis. Also, environmental factors that significantly correlate to 18S and 16S rDNA communities were fitted to the plot. Only environmental variables with a p-value of < 0.05 were selected. We performed NMDS with the metaMDS function and fitted the environmental variables with the envfit function, both functions are part of the vegan package, in R. This allowed us to investigate which environmental variables are driving the diversity of 18S and 16S rDNA communities. Figure 3.13a shows the 18S rDNA dataset at the taxonomic rank of class along with higher level taxonomic assignments but excluding those assigned to the eukaryote node. Similarly, figure 3.13b shows the 16S rDNA dataset at the taxonomic rank of class along with higher level taxonomic assignments but excluding those assigned to the bacteria node. We log₁₀ transformed the 18S and 16S rDNA datasets for our NMDS analysis.

In figure 3.13a which displays the 18S rDNA dataset, temperature, salinity, phosphate and silicate are the environmental factors fitted, each with a p-value of 0.001. The majority of the samples cluster together, distinct from a sub-population on the right hand side of the plot. This sub-population is composed of samples 29, 34, 37 and 52. These are the same samples that cluster separately in the ordination analysis of the PCoA plots in figure 3.8a. In figure 3.13a, the short distance between the vectors representing the environmental variables temperature and salinity suggests that they are strongly positively correlated with one another. Likewise, for phosphate and silicate, the represented vectors are also positioned very close together, indicating they have a very strong positive correlation. Also given the opposing orientation of the pairs of environmental variables of temperature and salinity to phosphate and silicate this suggests a strong negative correlation between the two pairs. In the larger cluster, the samples are grouped by their regions of Arctic Ocean, North Atlantic Ocean and South Atlantic Ocean. Samples from the polar region of the Arctic Ocean correlate strongly to phosphate and silicate, and samples from the tropical region of the South Atlantic

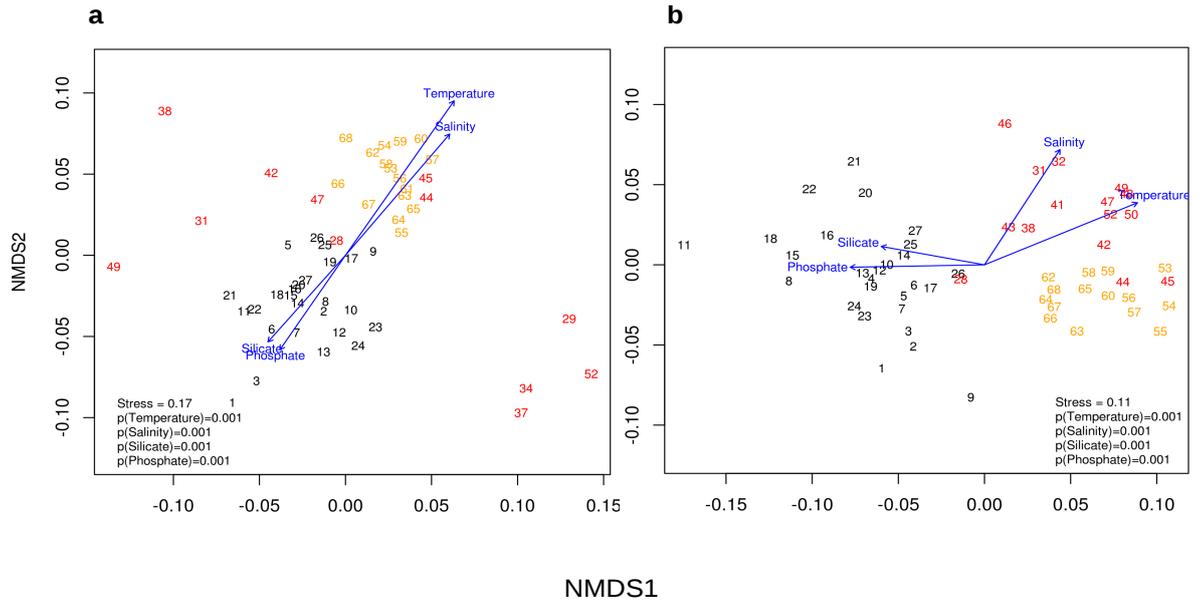


Figure 3.13: Panel **a**, NMDS of sampled stations with each number representing one sample of the 18S rDNA community. Significant environmental vectors for temperature ($p=0.001$), salinity ($p=0.001$), phosphate ($p=0.001$) and silicate ($p=0.001$) were fixed. Panel **b**, NMDS of sampled stations with each number representing one sample of 16S rDNA community. Significant environmental vectors for temperature ($p=0.001$), salinity ($p=0.001$), phosphate ($p=0.001$) and silicate ($p=0.001$) were fixed. NMDS was performed with the metaMDS function and the environmental variables were fitted with the envfit function (permutation test, 999 permutations). Both functions are part of the vegan package, in R. The numbers correspond to sample locations and the colours of the numbers correspond to ocean regions, black corresponds to the Arctic Ocean, red corresponds to the North Atlantic Ocean and yellow corresponds to the South Atlantic Ocean as shown in figure 3.1a

Ocean correlate strongly to temperature and salinity.

In figure 3.13b, which displays the 16S rDNA dataset, temperature, salinity, phosphate and silicate are again the environmental factors fitted, each with a p-value of 0.001. In figure 3.13b, the distance between the vectors representing the environmental variables temperature and salinity suggests that they are positively correlated, though not so strongly as in the 18S rDNA dataset. Likewise, for phosphate and silicate, the vectors are positioned even closer together, indicating they are strongly positively correlated. Also given the position of the vector representing temperature to those of phosphate and silicate this suggests a negative correlation, while salinity has a negative correlation to phosphate but a very weak correlation to silicate. The samples are

grouped by their regions of Arctic Ocean, North Atlantic Ocean and South Atlantic Ocean. Samples from the polar region of the Arctic Ocean correlate with phosphate and silicate, while samples from the tropical region of the South Atlantic Ocean correlate with temperature and salinity.

Alpha diversity versus temperature

We investigated alpha-diversity (Shannon index) (as explained in section 3.3.3) in relation to the environmental covariates. To determine which environmental covariates were significant in the 18S and 16S rDNA datasets, the environmental covariates were related to the datasets' Shannon index by fitting generalized linear models. We then used a step-by-step backwards selection of the environmental covariates for model building and removed non-significant environmental covariates until the remaining environmental covariates were significant.

For the 18S and 16S rDNA datasets, temperature was the only significant environmental covariate with a p-value of $5.73e-2$ and a p-value of $2.6e-05$, respectively, that explained significant amounts of variation in the diversity for each of the 18S and 16S rDNA datasets. In figure 3.14a and 3.14b we present alpha diversity plotted against temperature for 18S and 16S rDNA respectively. The y -axis represents the alpha diversity for the stations, the x -axis represents the temperature. The numbers in the plots correspond to sample locations, according to the map in figure 3.1a. For each of the datasets, a significant positive correlation was observed, for the 18S rDNA dataset a R^2 of 0.4 with a p-value of $1.99e-07$ and for the 16S rDNA dataset a R^2 of 0.7 with a p-value $< 2.2e-16$. The alpha diversity was lower in the polar and temperate communities and highest in the tropical communities.

This relationship is also observed in the heatmaps in figure 3.9a and 3.9b for the 18S and 16S rDNA datasets, respectively. In the heatmaps, diversity increases as we move from the cold regions in the Arctic Ocean to the tropical regions of the Atlantic Ocean.

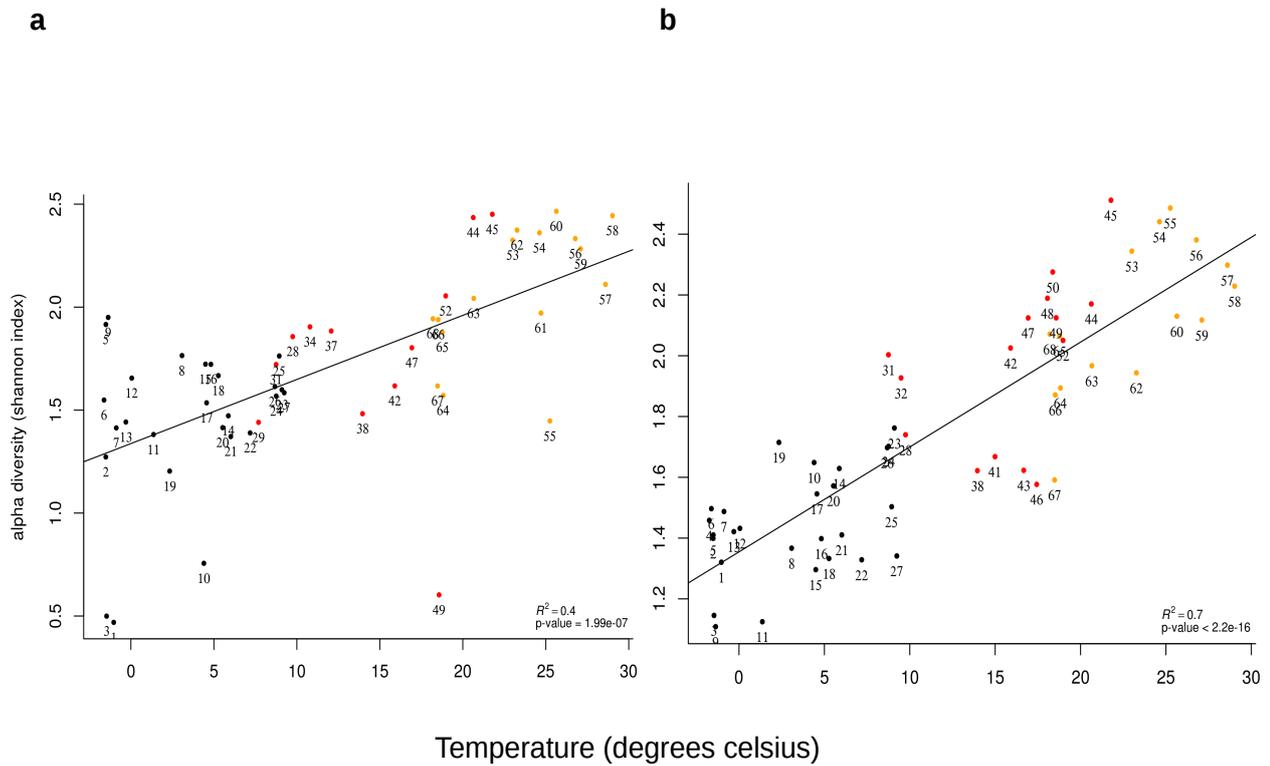


Figure 3.14: Panel **a**, a positive correlation of 18S rDNA diversity based on Shannon index with temperature. Based on backward model selection temperature was the only significant environmental covariate determined. Panel **b**, a positive correlation of 16S rDNA diversity, based on Shannon index with temperature. Based on backward model selection where temperature was the only significant environmental covariate determined. In panels **a** and **b**, the numbers correspond to sample locations and the colours of the numbers correspond to ocean regions, black corresponds to the Arctic Ocean, red corresponds to the North Atlantic Ocean and yellow corresponds to the South Atlantic Ocean as shown in figure 3.1a

3.4.6 Breakpoint analysis

In our previous analysis we investigated alpha-diversity in relation to the environmental covariates and determined temperature to be the only significant environmental covariate with a p-value of $5.73e-2$ for our 18S rDNA dataset and a p-value of $2.6e-05$ for our 16S rDNA dataset. We, therefore, chose to perform a breakpoint analysis in relation to temperature alone. The breakpoint analysis enabled us to investigate how increasing temperature affects the changing diversity of 18S and 16S rDNA across the Arctic Ocean down to the South Atlantic Ocean.

Displayed in figure 3.15a is the 18S rDNA dataset at the taxonomic rank of class

along with higher level taxonomic assignments but excluding those assigned to the eukaryote node. In figure 3.15**b** is the 16S rDNA dataset at the taxonomic rank of class along with higher level taxonomic assignments but excluding those assigned to the prokaryote node. In figure 3.15**a** and 3.15**b** is the breakpoint analysis for the 18S and 16S rDNA datasets, respectively. The numbers in the plots correspond to sample location as shown in figure 3.1**a**. The y -axis represents the beta diversity across the stations, the x -axis represents the temperature and the horizontal line marks the breakpoint.

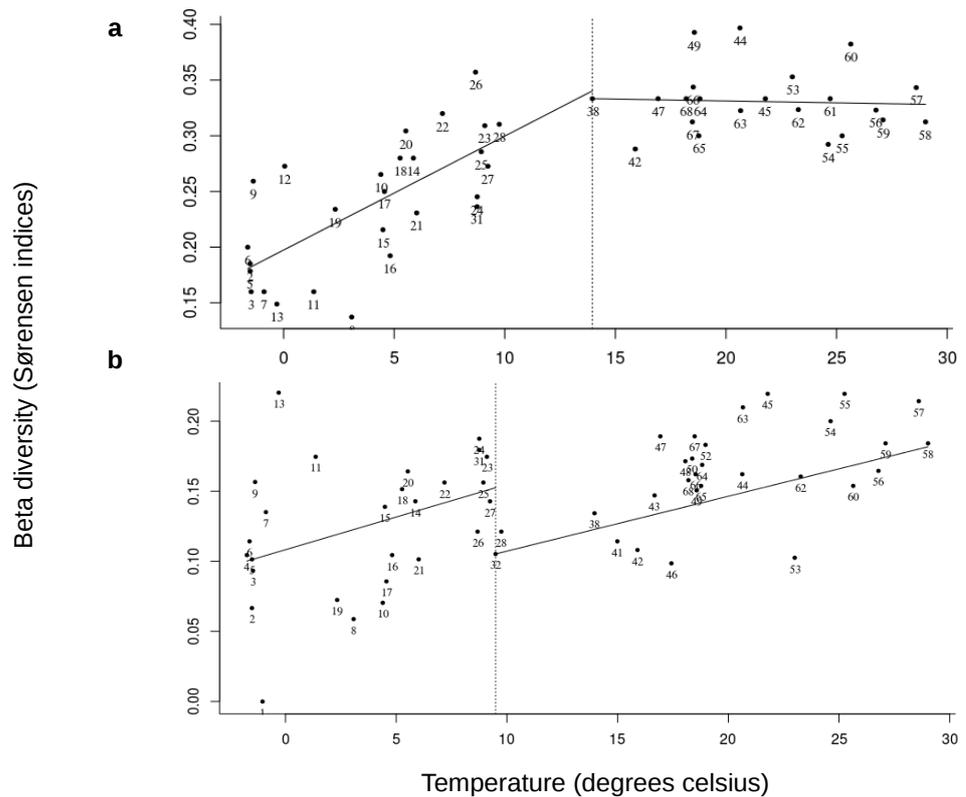


Figure 3.15: Panels **a** and **b** represent breakpoint analysis to the taxonomic rank of class, **a** represents the 18S rDNA dataset and **b** represents the 16S rDNA dataset. The numbers correspond to sample locations as shown in figure 3.1**a**. The breakpoint analysis was generated using piecewise regression in R as detailed in [Castro-Insua et al., 2016]. The y -axis represents the beta diversity across the stations. The x -axis represents the temperature. In each plot, the horizontal line marks the breakpoint. For the 18S rDNA dataset in panel **a** the breakpoint is 13.96°C with a p-value of 8.407e-11. For the 16S rDNA dataset in panel **b** the breakpoint is 9.49°C with a p-value of 1.413e-4

The breakpoint for the 18S rDNA dataset was determined to be 13.96°C with a p-value of 8.407e-11 and for the 16S rDNA dataset, the breakpoint was determined to be 9.49°C with a p-value of 1.413e-4. In figure 3.15**a** and 3.15**b**, the left hand side of the plots represent the lower temperatures, containing sample sites from the Arctic

Ocean and North Atlantic Ocean. The right hand side of the plots represent the higher temperatures and therefore contain sample sites from the North Atlantic Ocean and South Atlantic Ocean. There is a clear shift in the diversity for both the 18S and 16S rDNA dataset at their respective breakpoints. For both datasets of 18S and 16S rDNA, this shift in diversity occurs around the sample sites numbered in the 20's moving into the 30's, which positions the breakpoints in the North Atlantic Ocean as shown in figure 3.1a.

But the results are inconclusive due to the sparsity of the data points within the temperature range 9°C to 14°C, therefore we cannot conclusively determine breakpoints to exist within the temperate region of the North Atlantic Ocean. A change in diversity for both the 18S and 16S rDNA can be seen in the heatmaps in figure 3.9a and 3.9b. At the sample sites moving from the 20's to the 30's in the heatmaps which correspond to the temperature range 9°C to 14°C, the diversity can be seen to begin to change as it moves through the North Atlantic Ocean. Therefore potentially breakpoints within this region of the North Atlantic could exist but without greater sampling across this region, and at sufficient depth, we cannot determine this from our breakpoint analysis.

3.4.7 Co-occurrence analysis

In our co-occurrence analysis using the WGCNA package in R on the log10-scaled abundances of 18S rDNA species level and 16S rDNA genus level, two modules (a WGCNA term used to describe networks of highly correlating individuals) were found. In figure 3.16 we present the correlation heatmap between the two modules' eigengenes and the environmental variables. The module represented as blue ($n=51$) has a strong positive relationship with phosphate and a moderate positive relationship with silicate. Also, the blue module has a strong negative relationship with temperature and moderate negative relationship with salinity. The module represented as turquoise ($n=70$) has a strong negative relationship with phosphate and a moderate negative relationship with silicate. This turquoise module also has a strong positive relationship with temperature and a moderate positive relationship with salinity. Both the turquoise and the blue modules exhibit a very weak relationship with the environmental variable nitrate/nitrite. Above all the environmental variables, both the turquoise and blue

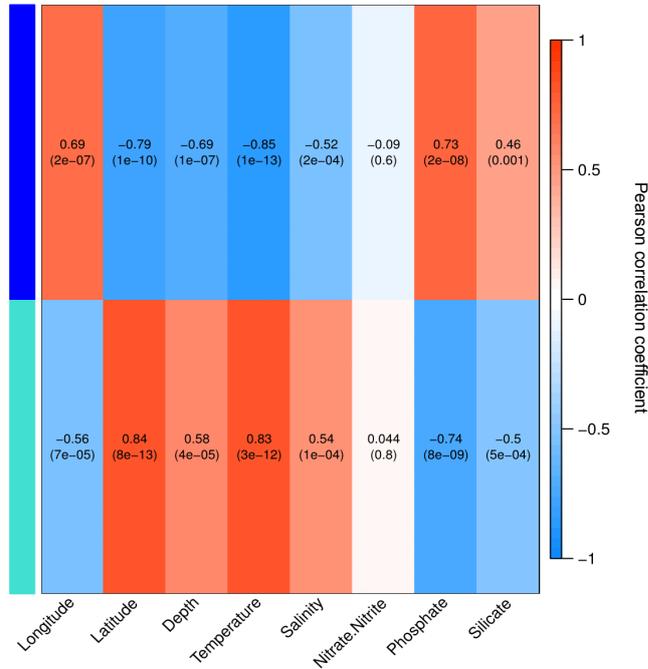


Figure 3.16: In the co-occurrence analysis with WGCNA on the log10-scaled abundances of 18S rDNA species level and 16S rDNA genus level, two modules were found. Depicted is the correlation heatmap between the modules' eigengene and environmental variables. Along the left side, the two modules are displayed in turquoise ($n=70$) and blue ($n=51$). The environmental variables are displayed at the bottom. The colours correspond to the correlation values; red is positively correlated and blue is negatively correlated. The values in each of the squares correspond to the assigned Pearson correlation coefficient value on top and p-value in brackets below

modules have the highest correlation to temperature. Also, the modules exhibit the opposite preference to temperature, the turquoise module has a strong positive correlation to temperature, while the blue module has a strong negative correlation to temperature.

We further examined each module's relationship to the environmental variables as depicted in figure 3.17. The correlation heatmap shows how each individual species of the turquoise module **a** ($n=70$) and the blue module **b** ($n=51$) correlate to the environmental variables. The taxa of the turquoise module represented in figure 3.17 **a** ($n=70$) have an overall negative relationship with phosphate and silicate and an overall positive relationship with temperature and salinity in varying degrees. The taxa of the blue module represented in figure 3.17 **b** ($n=51$) have an overall positive relationship with phosphate and silicate and an overall negative relationship with temperature and salinity in varying degrees. Both modules exhibit a very weak relationship with the environmental variable nitrate/nitrite.

In figure 3.18**a1** (turquoise ($n=70$)) and 3.18**b1** (blue ($n=51$)) the modules are

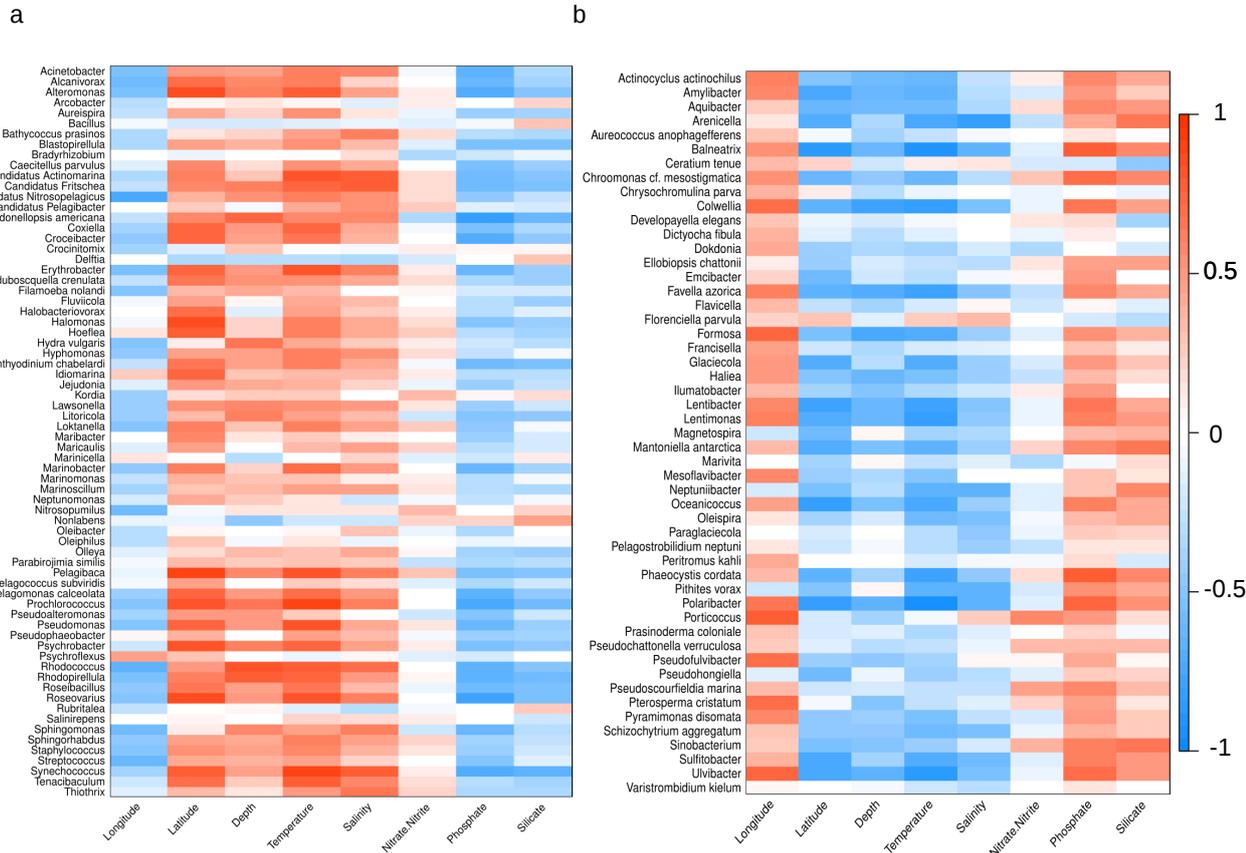


Figure 3.17: In the co-occurrence analysis with WGCNA on the log₁₀-scaled abundances of 18S rDNA species level and 16S rDNA genus level, two modules were found. Depicted is the correlation heatmap for each of the two modules **a** (turquoise ($n=70$)) and **b** (blue ($n=51$)) of their species log₁₀-scaled abundances and environmental variables. Along the left side on each of the two modules **a** and **b** are displayed the species name and environmental variables are displayed at the bottom. The colours correspond to the correlation values; red is positively correlated and blue is negatively correlated. The Pearson correlation coefficient values and p-values can be found in Appendix A.F

depicted as network diagrams. The edge distance between the nodes is an indication of the correlation strength. The numbers on the nodes correspond to the taxa names; the list of numbers to species can be found in the Appendix A.E. Degrees of connectivity refers the nodes (species) with the highest number of connections to other nodes (species). This does not automatically indicate the most abundant species or species that appear in all the stations. We could have connectivity between species that are lowly abundant and appear in few stations. The top five species/nodes with the greatest degrees of connectivity are coloured in orange. For the module in figure 3.18a1 the five most highly connected taxa are Erythrobracter(1), Alteromonas(2), Roseovarius(3), Marinobacter(4) and Pelagomonas calceolata(5). For the module in

figure 3.18**b1** the five most highly connected taxa are Colwellia(1), Polaribacter(2), Balneatrix(3), Ulvibacter(4) and Amylibacter(5).

We also depict the two modules as word clouds in figure 3.18**a2** (turquoise ($n=70$)) and 3.18**b2** (blue ($n=51$)). These consist of the member taxa names with the size of a word reflecting the number of connections that taxa possess. Displayed in figure 3.18**a3** (turquoise ($n=70$)) and **b3** (blue ($n=51$)) are pie charts of the modules' taxa

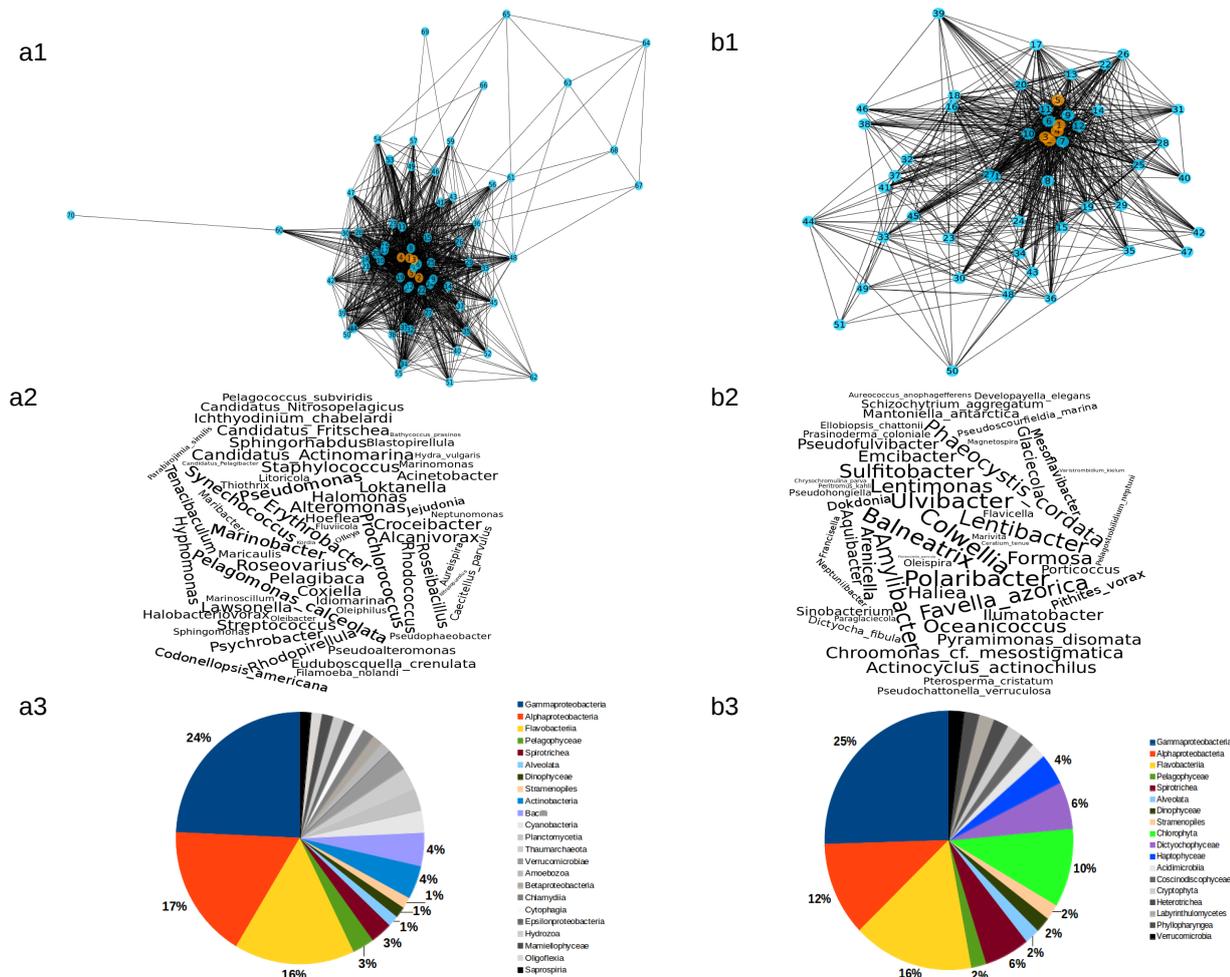


Figure 3.18: Co-occurrence analysis with WGCNA on the log₁₀-scaled abundances of 18S rDNA species level and 16S rDNA genus level. In panels **a1** (turquoise ($n=70$)) and **b1** (blue ($n=51$)) are the two modules that we found depicted as network diagrams. These were generated with Cytoscape [Shannon et al., 2003]. The edge distance indicates the correlation strength between the nodes and the top five most highly connected nodes are coloured in orange. The numbers on the nodes correspond to the taxa names, the list of names to numbers of the taxa can be found in the Appendix A.E. In panels **a2** (turquoise ($n=70$)) and **b2** (blue ($n=51$)) the modules are depicted as word clouds. These consist of the member species and genus names, generated in WordCloud.com. The larger the name, the more connections that taxa has. In panel **a2**, taxa Crocinitomix, Salinirepens, Bacillus, Bradyrhizobium, Rubritalea, Arcobacter, Delftia, Marinicella, Nonlabens and Psychroflexus were removed as they are unreadable due to their frequency being too low to display with the others. In panels **a3** (turquoise ($n=70$)) and **b3** (blue ($n=51$)) are pie charts of the classes of the modules' taxa. The list of names of the taxa to their class can be found in the Appendix A.E

class membership; the list of names of the taxa to their class can be found in the Appendix A.E. Gammaproteobacteria, Alphaproteobacteria, Flavobacteriia, Pelagophyceae, Spirotrichea, Alveolata, Dinophyceae and Stramenopiles are the greatest contributors in each module and are similarly proportionate. In figure 3.18a3 (turquoise ($n=70$)), Actinobacteria and Bacilli contribute significantly but are unique to that module. Also in figure 3.18b3 (blue ($n=51$)) Chlorophyta, Dicictyochophyceae and Haptophyceae contribute significantly and are unique to that module. These results are consistent with what we see in the heatmaps in figure 3.9a and 3.9b. The heatmaps show the distribution and abundance across the samples from the polar Arctic Ocean down to the tropical South Atlantic Ocean.

3.5 Discussion

In this chapter, we have presented a phylogenetic analysis of phytoplankton 18S and 16S rDNA. We described the methodology of the phylogenetic analysis and normalised copy number for both the 18S and 16S rDNA dataset, and their analysis with rarefaction curves, heatmaps, PCoA, evenness and occupancy plots, environmental plots, breakpoint analysis and co-occurrence analysis. This was a unique large scale examination of the 18S and 16S rDNA species communities taken from a transect of the Arctic and Atlantic oceans that was sampled in close proximity to the coast. There are significant differences between coastal waters and the open ocean. Coastal waters have a lower temperature and higher nutrient content in comparison to the open ocean [Toseland et al., 2013]. Therefore coastal waters are more productive than the open ocean for phytoplankton, in terms of, for example, Chlorophyll *a* concentrations and primary productivity [Trimborn et al., 2015]. This has therefore given us a new understanding of how the environmental conditions affect phytoplankton species communities.

While to the best of our abilities we have extensively sampled and precisely designed our experimental approaches there are limitations in our analysis. One such limitation is that the North Atlantic Ocean samples were acquired from two other collaborators; Dr.Willem van de Poll of the University of Groningen, Netherlands and Dr.Klaas Timmermans of the Royal Netherlands Institute for Sea Research. Their sampling procedures were slightly different to those performed by Dr.Katrin Schmidt

who performed an additional pre-filtering step with a 100m mesh to remove larger organisms. Also another limitation is that a number of the samples from the North Atlantic Ocean failed to pass quality control steps during sequencing. From our 18S rDNA samples a total of 13 out of 24 samples failed from the North Atlantic Ocean and also station 4 from the Arctic Ocean. From our 16S rDNA samples, a total of 10 samples failed out of 24 from the North Atlantic Ocean and also station 61 from the South Atlantic Ocean. In addition, only a single sample was obtained at each station; we did not obtain replicate samples. However, we regard these as minimal limitations given the number of high quality samples that we obtained and the relatively close proximity of the samples to one another.

For our analysis of the 18S and 16S rDNA datasets, we identified a gradient of increasing diversity, as we moved from the cold temperatures of the Arctic Ocean down to the tropical temperatures of the South Atlantic Ocean. These findings are what we expected to observe, as it has been known for years that diversity is higher in the warm tropical regions than in the cold polar regions of the world [Brown, 2014]. We also identified a number of 18S and 16S rDNA species that were occurring continuously throughout the samples, which included, for example, Stramoenopiles and Gammaproteobacteria, and these findings are consistent with previewed published works have found [Lin et al., 2012], [Franco et al., 2017].

We further identified in both 18S and 16S rDNA datasets, that the samples from the polar region of the Arctic Ocean are correlated strongly to phosphate and silicate, and samples from the tropical region of the South Atlantic Ocean are correlated strongly to temperature and salinity. In addition, for both the 18S and 16S rDNA datasets we found that temperature was the only significant environmental covariate for alpha diversity. In our co-occurrence analysis with WGCNA on the 18S rDNA species level and 16S rDNA genus level, we found two modules. The larger of the modules (turquoise ($n=70$)) was found to have a strong positive correlation relationship to temperature, indicating that this module is likely to be found in a warm climate. In contrast, the smaller module (blue ($n=51$)) was found to have an overall strong negative correlation relationship to temperature, indicating that this module is likely to be found in a cold climate. These co-occurrence analysis results are important as we can hypothesise that different microbial communities have different preferences for temperature. Moreover,

as global warming is predicted to raise the temperatures in the ocean, our results could potentially enable us to forecast how climate change will affect these microbial communities using climate models underpinned by genetic information.

In our breakpoint analysis of the 18S and 16S rDNA datasets, we identified a putative breakpoint of 13.96°C for the 18S rDNA dataset and 9.49°C for the 16S rDNA dataset. This positioned the breakpoint in the North Atlantic Ocean off the coast of France. It is interesting to note that the breakpoints were all located in the temperate region of the North Atlantic Ocean. Therefore as you move from the cold Arctic Ocean to the warm tropical regions in the South Atlantic Ocean there is a radical shift in the diversity and interactions of the 18S and 16S rDNA species communities. There have been various types of studies that have included a breakpoint analysis, such as [Campra and Morales, 2016] which looks at surface air temperature records in southeastern Spain over a number of years. Also, another breakpoint study is [Castro-Insua et al., 2016], which is an analysis of various terrestrial animals such as bats and birds across latitudes in America. They determined a number of breakpoints for each of their subject animals, and these were found in the range of 25° to 58° latitude [Castro-Insua et al., 2016]. Our breakpoints are also located in this latitude range, as we can determine from our metadata that the 18S rDNA breakpoint of 13.96°C is located at 59.5° latitude and the 16S rDNA breakpoint is located at 45.5° latitude. But our results are currently inconclusive due to the sparsity of the data points within the temperature range 9 to 14°C, therefore we cannot conclusively determine breakpoints to exist within the temperate region of the North Atlantic Ocean. To explicitly determine breakpoints in our 18S and 16S rDNA datasets we require further sampling to be uniformly distributed across the entire sampling range and of sufficient sampling depth throughout.

Chapter 4

Metatranscriptomics analysis

4.1 Summary

In the last chapter, we analysed 18S and 16S rDNA datasets from the Arctic Ocean, North Atlantic Ocean and the South Atlantic Ocean and found a greater diversity of microbes in the tropical regions of the South Atlantic Ocean, versus the polar regions of the Arctic Ocean. Additionally, in our co-occurrence analysis on the 18S and 16S rDNA datasets, we found two community networks, one positively correlated to temperature and the other negatively correlated to temperature. Based on these results, we can hypothesise that different microbial communities have different preferences for temperature. We also performed a breakpoint analysis on our 18S and 16S rDNA datasets and found a shift in diversity occurring in the North Atlantic Ocean. In particular, the shift occurs in the temperate region of the Ocean, between the polar Arctic Ocean and the tropical South Atlantic Ocean.

Now we focus on the question of what are the microbes doing, and in particular gaining an understanding of their genetic activity. Moreover, it will be of interest to see if the different data types (16S/18S rDNA data and metatranscriptome data) are in agreement or not. To address this question, we shall perform a metatranscriptomic analysis. In the next section, we describe the sequencing of the metatranscriptomic dataset. In Section 4.3 we describe our pipelines for the metatranscriptomic analysis, as well as additional methods, then we will present the results of our analysis in Section 4.4. In Section 4.5, we end with a discussion of the results presented in this chapter.

4.2 Sequencing and preprocessing

The samples were collected as described in chapter 3. All samples were sequenced and preprocessed by the Joint Genome Institute (JGI) (Department of Energy, Walnut Creek, CA, USA). Metatranscript sequencing was performed on an Illumina HiSeq-2000 instrument [Huntemann et al., 2016]. A total of 65 samples passed quality control after sequencing with 5.7 Gb of sequence read data over all samples for analysis. Here we describe the JGI’s computational pipeline for preprocessing the metatranscript reads.

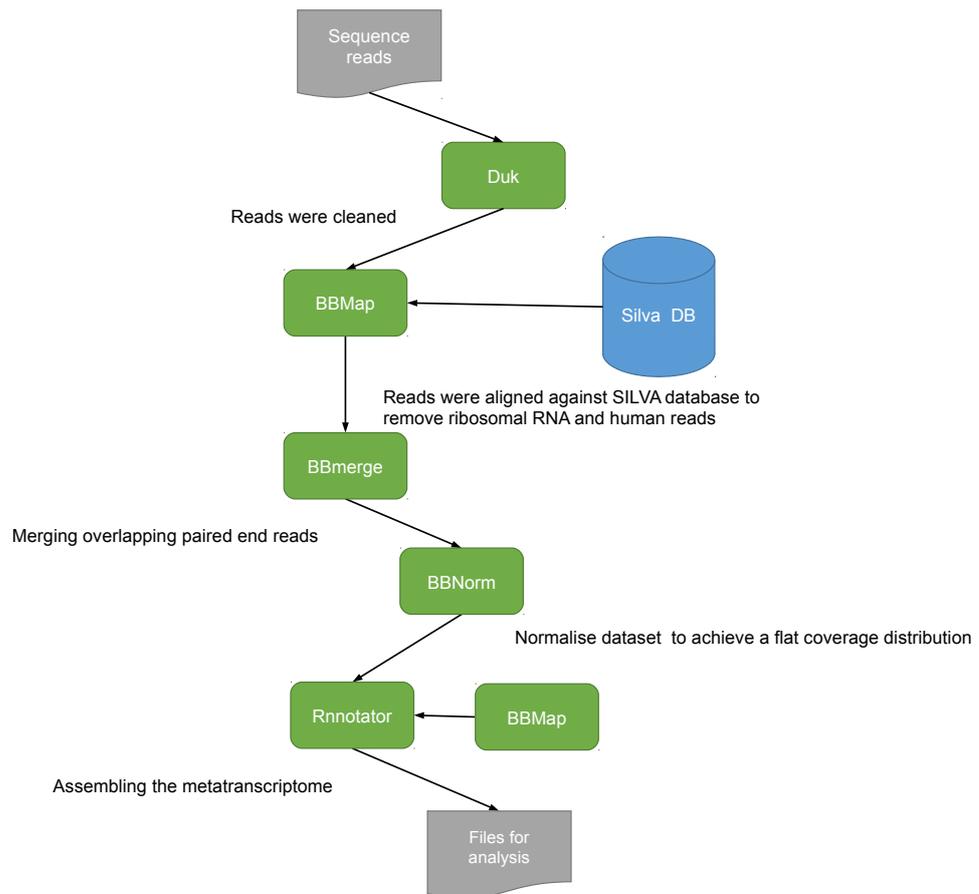


Figure 4.1: Diagram of JGI’s computational pipeline for preprocessing the metatranscript reads. The pipeline at various stages incorporates databases (blue), BBTools tools (green) and processed files (grey)

In figure 4.1 we provide an overview of the JGI pipeline. JGI employed their suite of tools called BBTools [DOE Joint Genome Institute, 2017] for preprocessing the sequences. As shown in figure 4.1, first the sequences were cleaned using a tool called Duk [DOE Joint Genome Institute, 2017]. Duk (see Section 2.4.1) is a tool in the BBTools suite, that performs various data quality procedures such as quality trimming and filtering by kmer matching [DOE Joint Genome Institute, 2017]. In our dataset, Duk identified and removed adapter sequences, and also quality trimmed the raw reads

to a phred score of Q10. In Duk the parameters were; kmer-trim (ktrim) was set to r, kmer (k) was set to 25, shorter kmers (mink) set to 12, quality trimming (qtrim) was set to r, trimming phred (trimq) set to 10, average quality below (maq) set to 10, maximum Ns (maxns) set to 3, minimum read length (minlen) set to 50, the flag “tpe” was set to t, so both reads are trimmed to the same length and the “tbo” flag was set to t, so to trim adapters based on pair overlap detection. The reads were further filtered to remove process artefacts also using Duk with the kmer (k) parameter set to 16.

BBMap [DOE Joint Genome Institute, 2017] is also a tool in the BBTools suite, that performs other operations such as making sequence alignments of DNA and RNA reads to a database. BBMap aligns the reads by using a multi-kmer-seed-and-extend approach [DOE Joint Genome Institute, 2017]. To remove ribosomal RNA reads, the reads were aligned against a trimmed version of the SILVA database using BBMap with parameters set to; minratio (minid) set to 0.90, local alignment converter flag (local) set to t and fast flag (fast) set to t. Also any human reads identified were removed using BBMap [DOE Joint Genome Institute, 2017].

BBmerge [DOE Joint Genome Institute, 2017] is a tool in the BBTools suite that performs the merging of overlapping paired end reads [DOE Joint Genome Institute, 2017]. For assembling the metatranscriptome, the reads were first merged with the tool BBmerge, and then BBNorm was used to normalise the coverage so as to generate a flat coverage distribution. This type of operation can speed up assembly and can even result in an improved assembly quality [DOE Joint Genome Institute, 2017].

Finally as shown in figure 4.1, Rnnotator [Martin et al., 2010] was employed for assembling the metatranscriptome. Rnnotator assembles the transcripts by using a de novo assembly approach of RNA-Seq data and it accomplishes this without a reference genome [Martin et al., 2010]. The tool BBMap was used for reference mapping, the cleaned reads were mapped to metagenome/isolate reference(s) and the metatranscriptome assembly [DOE Joint Genome Institute, 2017].

4.3 Methods

4.3.1 Computational pipeline for taxonomic classification analysis

PhymmBL is a hybrid taxonomic classifier, that combines Phymm composition-based taxonomic predictions and BLAST based homology results to label each of the input sequences [Mande et al., 2012]. Phymm employs interpolated Markov models (IMMs) which is a form of the Markov chain that uses a variable number of states to calculate the probability of the next state. Phymm builds IMMs to characterize the variable length oligonucleotides that are distinct for a particular phylogenetic clade, whether it be a species, genus or a higher taxonomic level. During construction of the IMMs, the IMM algorithm builds a probability distribution based on the observed patterns of nucleotides that describe each species in the reference database. During classification for each of the input sequences, each of the IMMs is used as a scoring method by inspecting the nucleotides in the input sequence and then outputs a score that corresponds to the probability that the input sequence was generated from the same distribution as that used to construct the IMM. The input sequence is classified with the clade labels that belong to the organism whose IMM produced the best score for that input sequence. During BLAST, each input sequence is submitted as a BLASTN query, searching against the same reference database used to generate the IMMs, and clade labels are assigned for the best BLAST hit. The combined score of Phymm and BLAST for PhymmBL is determined by using the function $\text{score} = \text{IMM} + 1.2(4 - \log(E))$, where IMM is the score from the best matching IMM and E is the smallest E-value given by BLAST. It has been demonstrated by [Brady and Salzberg, 2009a] that PhymmBL's hybrid method outperforms Phymm and BLAST separately and also that BLAST outperforms Phymm [Brady and Salzberg, 2009a].

The metatranscriptomic dataset was taxonomically classified using Dr. Andrew Toseland's taxonomic classification pipeline as described in [Toseland et al., 2013]. This pipeline employ's PhymmBL [Brady and Salzberg, 2009b] which contains phytoplankton in its database, and the contents of PhymmBL's default database were reduced by taking the first occurrence of a species under each genus. We ran the pipeline with

Dr. Toseland's assistance.

The PhymmBL taxonomic classification pipeline contains a representative set of 44 eukaryote organisms. This set consists of genomes and expressed sequence tags (EST) of the major eukaryote groups with a focus on algal species. From NCBI-dbEST the EST sequences were downloaded. Then with CD-HIT-est these EST sequences were clustered with a 95% similarity in order to ensure non-redundancy of the sequences. The genome sequences were downloaded from NCBI GenBank, JGI and separately four genomes; *Cyanidioschyzon merolae*, *Danio rerio*, *Homo sapiens* and *Strongylocentrotus purpuratus* were downloaded. These four genomes were obtained from specific locations as outlined in Appendix B.A. Species such as *Danio rerio* and *Homo sapiens* were included to check for contamination. Taxonomic classifications were taken from the NCBI taxonomy [NCBI Resource Coordinators, 2016] and AlgaeBase [Guiry, M.D. & Guiry, 2008] to produce a PhymmBL configuration file. In batch mode, the sequence files and taxonomic details were added to PhymmBL and interpolated Markov models (IMMs) were created for each new organism [Toseland et al., 2013].

For each sequence, PhymmBL outputs a taxonomic label of genus, family, order, class and phylum with a confidence score between 0 and 1 [Brady and Salzberg, 2011]. For our PhymmBL results, we took a confidence score cutoff of ≥ 0.9 at the phylum level to order to ensure our results were as accurate as possible. For each sample, the number of sequences under each taxon was summed up. The files were further normalised by applying hits per million.

4.3.2 Computational pipeline for functional analysis

JGI performed the functional analysis on the metatranscriptomic dataset. A functional analysis for a metatranscriptomic dataset is a standard pipeline but we liaised with JGI, and they updated their pipeline based on our feedback. For example, functional analysis results obtained using JGI's Integrated Microbial Genomes (IMG) standard pipeline were different to those obtained using the standard JGI pipeline. It was found that this was due to different versions of databases being employed. This feedback induced JGI to update their databases in their standard pipeline.

Our datasets were submitted to JGI's IMG. JGI's annotation system is called the

Metagenome Annotation Pipeline (MAP) (v4.15.2). JGI used HMMER 3.1b2 [Eddy, 1996] and the Pfam v30 [Finn et al., 2016] database for the functional analysis of our metatranscriptomic dataset. There are other databases that perform similar roles but we only discuss Pfam because that is the one we use for this thesis.

4.3.3 Further analysis

In this section, we describe the methodology that we developed for our analysis of the metatranscriptomic dataset. From JGI's IMG, we downloaded a Pfam file for each sample. This resulted in 7,453 Pfam functional assignments and their gene counts across the 65 samples. The files were further normalised by applying hits per million.

Canonical Correspondence Analysis (CCA)

We used the R package VEGAN to perform a Canonical Correspondence Analysis (CCA) between the Pfam dataset and the environmental data. CCA uses a dataset of measured variables such as Pfam gene counts in our case, and a dataset of additional explanatory variables such as temperature and salinity during the analysis, in order to find the relationship between the two datasets, and therefore enable us to determine how might the environmental variables determine the response variable values [Paliy and Shankar, 2016].

Environmental variables such as temperature and salinity can influence microbial communities [Hou et al., 2017]. Therefore in chapter 3 for the analysis of our 18S and 16S rDNA datasets, we used NMDS plots to visualise the distance matrix based on Bray-Curtis [Paliy and Shankar, 2016] with environmental factors fitted that significantly correlate to 18S and 16S rDNA communities. To analyse the metatranscriptome data we instead performed a CCA on our Pfam protein families datasets, because not all genes are affected by environmental variables such as housekeeping genes which are constantly expressed [Eisenberg and Levanon, 2013]. CCA enables us to determine what proportion of our dataset is affected by environmental variables and identify the environmental variables that significantly explain the variation of our dataset [Paliy and Shankar, 2016].

We log10 transformed the Pfam normalised gene counts, as explained in chapter

3 so that the data complies better to the assumptions of a parametric statistical test [Paliy and Shankar, 2016]. The environmental data consisted of temperature, salinity, nitrate/nitrite, phosphate and silicate.

Co-occurrence analysis

The methodology of how we performed the co-occurrence analysis is outlined in section 3.3.3. Based on the Pfam log₁₀-scaled gene counts dataset the power beta ($\beta \geq 1$) was determined to be 15.

The co-occurrence analysis resulted in fifteen modules being found, that included a grey module which represents those taxa that could not be assigned to a module. The modules were further examined to determine if any modules were highly correlated (≥ 0.75) to one another based on their eigengene. This resulted in two modules being merged and therefore thirteen modules were taken for analysis.

4.4 Results

4.4.1 Taxonomic classification heatmap

In figure 4.2 we present a heatmap that we generated, arranged by latitude for the taxonomic classified metatranscriptomic dataset at the taxonomic rank of phylum. The numbers at the bottom correspond to sample site numbers as shown in figure 3.1a. The heatmap is arranged by placing the most abundant read counts at the top of the plot down to the least abundant read counts at the bottom. Heatmaps enabled us to overview the distribution, composition and abundant read counts of our dataset. In figure 4.2 we basically observe no gradient of increasing diversity of taxa across the samples as we move from the polar Arctic Ocean through the temperate North Atlantic Ocean and into the tropical South Atlantic Ocean. This is in contrast to what we observed in the heatmaps of the 18S and 16S rDNA datasets in figure 3.9a and 3.9b, respectively. In figure 4.2 the majority of the taxa display a consistent abundance across the samples. There is variation between taxa abundance, as the heatmap is roughly divided into four blocks of colour.

It is not meaningful to compare the read counts of the metatranscriptomic dataset

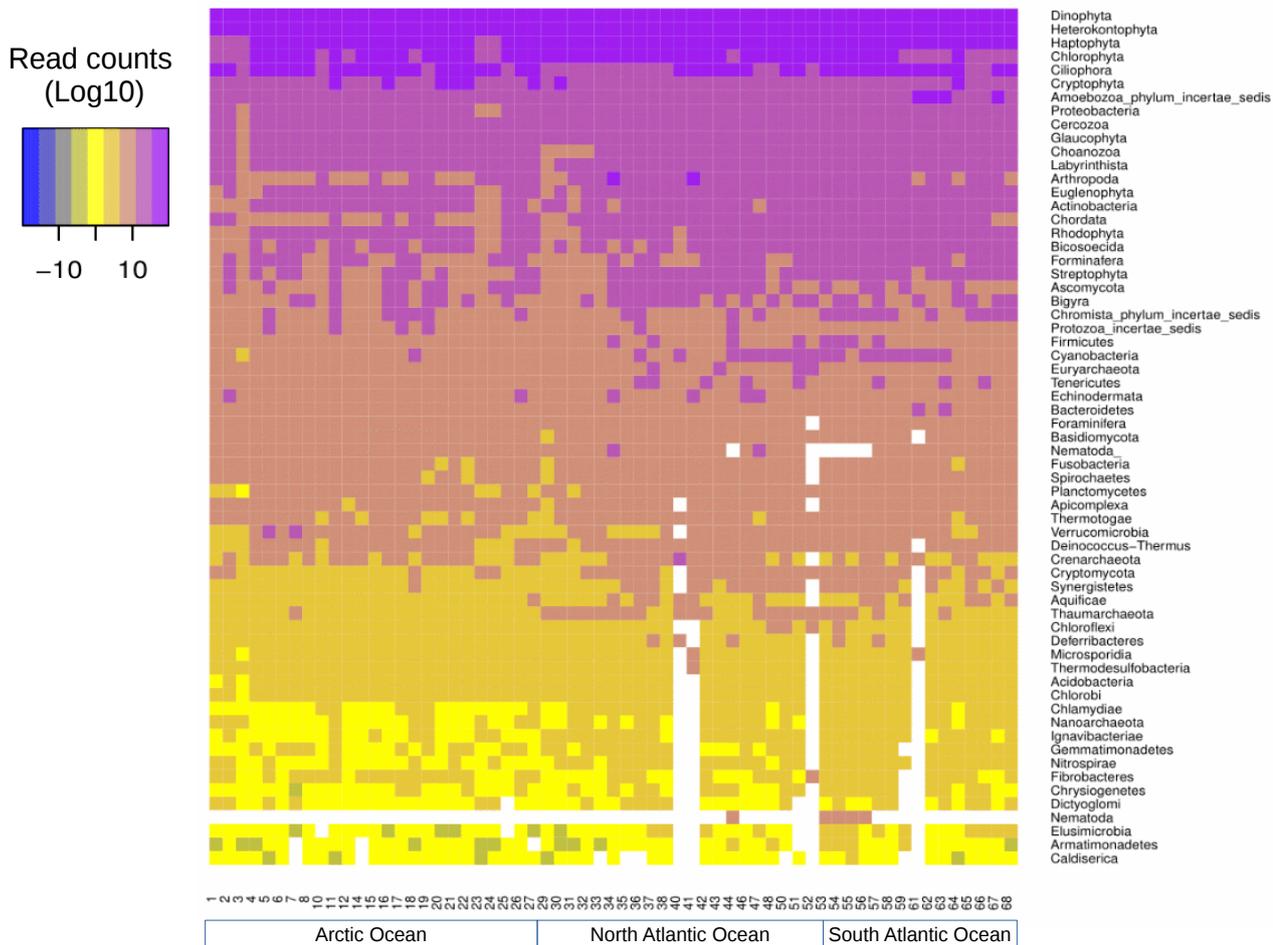


Figure 4.2: A heatmap of the metatranscriptomic dataset taxonomically classified and arranged by latitude versus the taxonomic rank of phylum. The taxonomy names in the dataset are displayed along the right side. The numbers at the bottom correspond to sample locations as shown in figure 3.1 a. The three regions of the Arctic Ocean, North Atlantic Ocean and South Atlantic Ocean are displayed underneath their corresponding sample numbers. The colours correspond to log₁₀-scaled read counts, where purple colours are high values and blue colours are low values. The heatmap was generated using the heatmap.2 function, which is part of the gplots package, in R

in figure 4.2 with the abundances of the 18S and 16S rDNA datasets in figure 3.9. This is because abundance gives an indication on the density of a species, whereas read counts only give an indication of presence or absence. However with metagenomic data that we will have in the near future for these sample sites, we can use this data to compare with the abundances of the 18S and 16S rDNA datasets. We do not see a gradient on species in the metatranscriptomic dataset as we did in the 18S and 16S rDNA datasets because the 18S and 16S rDNA analysis targets specific genes and then taxonomically classifies the samples against a reference database based on the SILVA database as described in Section 3.3.1 and Section 3.3.2, respectively [Reller et al., 2007]. In contrast the metatranscriptomic analysis employs direct cDNA sequencing of

the sample and then taxonomically classifies it against a reference databases of genomes and NCBI expressed sequence tags (EST) as described in Section 4.4.1 [Leimena et al., 2013].

However, we can compare identified taxa entities between the metatranscriptomic dataset and the 18S and 16S rDNA datasets, as a means of confirmation between the two analyses. In figure 4.3 **a** and **b** we represent Venn diagrams, comparing the number of taxa entities at the taxonomic rank of phylum between the 18S rDNA dataset to the metatranscriptomic dataset and the 16S rDNA dataset to the metatranscriptomic dataset, respectively. In figure 4.3 **a**, in the comparison of the 18S rDNA taxa entities to the metatranscriptomic taxa entities we see they share 10 taxa entities such as Dinophyta and Chlorophyta, while the 18S rDNA dataset contains 27 unique taxa entities such as Rotifera and the metatranscriptomic dataset contains 64 unique taxa entities such as Crenarchaeota. In figure 4.3 **b**, the comparison of the 16S rDNA taxa entities to the metatranscriptomic taxa entities we see they share 18 taxa entities such as Cyanobacteria and Proteobacteria, while the 16S rDNA dataset contains 13 unique taxa entities such as Balneolaeota and the metatranscriptomic dataset contain 56 unique taxa entities such as Ignavibacteriae.

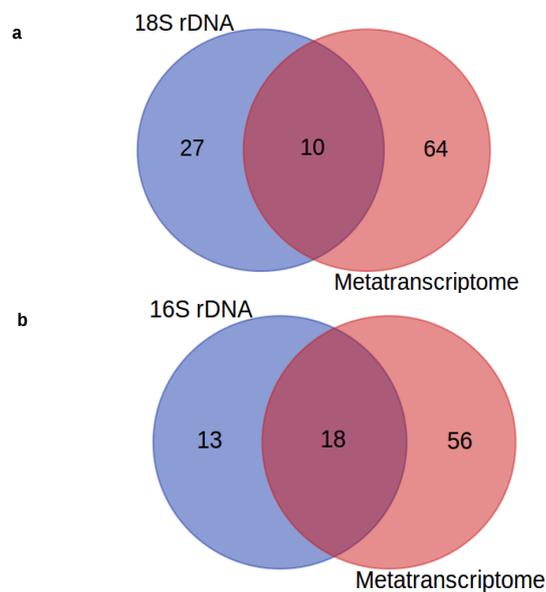


Figure 4.3: Panel **a** and **b** represents Venn diagrams of the number of taxa entities in the 18S, 16S rDNA and metatranscriptome datasets at the taxonomic rank of phylum. Panel **a** is comparing the 18S rDNA dataset and metatranscriptome dataset. Panel **b** is comparing the 16S rDNA dataset and metatranscriptome dataset. (Figure was generated with <http://bioinformatics.psb.ugent.be/webtools/Venn/>)

4.4.2 Rarefaction curves

The rarefaction curves in figure 4.4 are based on the Pfam protein families. These curves were generated using the rarecurve function, which is part of the VEGAN package, in R. The numbers displayed in the plot correspond to sample location, as shown in figure 3.1a. As outlined in section 3.4, rarefaction curves enable us to investigate if we have sufficient data to justify the results of our analyses.

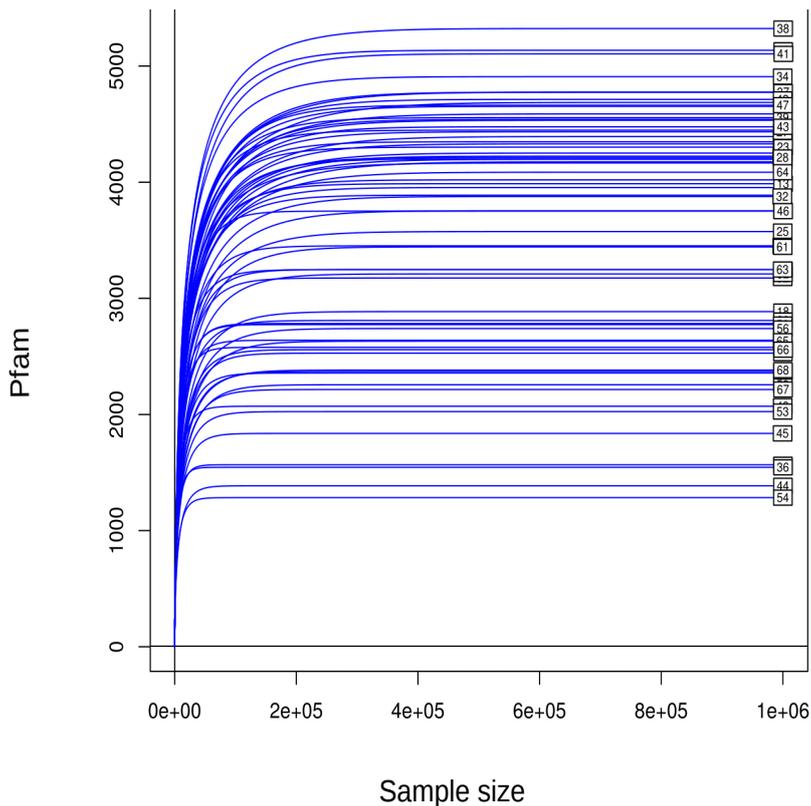


Figure 4.4: Rarefaction curves for Pfam protein families. The numbers displayed in the plot correspond to sample location. The x-axis represents the random subsample size taken from the dataset and the y-axis indicates the number of unique Pfam protein families found. The R package VEGAN using the rarecurve function was employed to perform the rarefaction curves analysis

In figure 4.4, there is a sharp rise at first in all the curves for the 65 samples. The majority of the curves, for example samples 47, 64 and 25, rise more slowly as more rare species are added, and then the curves level off. There are a number of curves, for example samples 53, 66 and 45, which level off immediately. The levelling off of all of the curves happens quite quickly, therefore we conclude that sufficient sampling was achieved.

4.4.3 Principle Components Analysis (PCA)

The Principle Components Analysis (PCA) in figure 4.5 is based on the Pfam protein families (log₁₀ transformed) gene counts. This PCA was generated using the princomp function, which is part of the VEGAN package, in R. The numbers and colours displayed in the plot correspond to sample location, as shown in figure 3.1a.

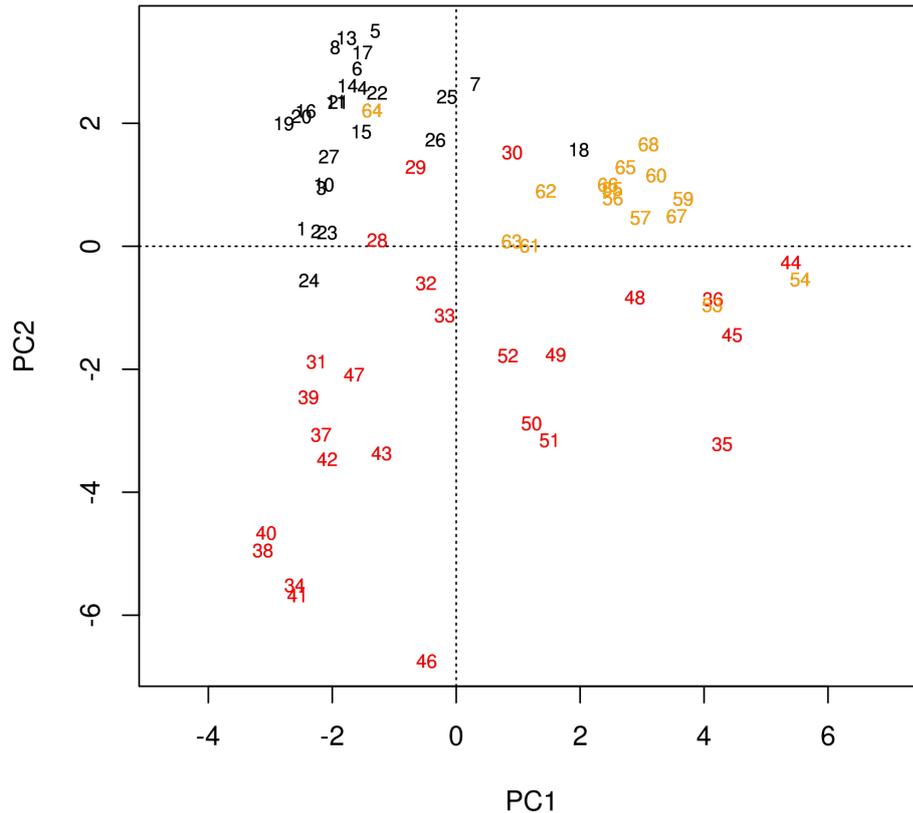


Figure 4.5: A Principle Components Analysis (PCA) for the Pfam protein families (log₁₀ transformed) gene counts. The samples are numbered and coloured by region; these numbers and colours correspond to the map of the region sampling sites in figure 3.1a, where the Arctic Ocean samples are coloured black, the North Atlantic Ocean samples are coloured red and the South Atlantic Ocean samples are coloured yellow. The R package VEGAN using the princomp function was employed to perform the PCA

PCA is one of the most popular methods for exploratory analyses, as it is a simple visualization tool to summarize dataset variance [Ramette, 2007], [Paliy and Shankar, 2016]. The general principle of PCA is to calculate new synthetic variables called principal components. Principal components are linear combinations of the initial variables and these attempts to account for as much of the variance of the initial dataset

as possible. The objective is to represent the objects and variables of the dataset in a new system of coordinates. This is generally on two axes where the maximum amount of variation from the initial dataset can be displayed [Ramette, 2007]. The largest gradient of variability in the dataset is represented by the first principal component axis of the PCA. The second principal component represents the second largest and so on, till all the dataset variability has been accounted for [Paliy and Shankar, 2016].

Displayed in figure 4.5 is the Pfm protein families (log10 transformed) gene counts for the first two principal components. In the PCA plot the samples are numbered and coloured by region; these numbers and colours correspond to the map of the region sampling sites in figure 3.1a, where the Arctic Ocean samples are coloured black, the North Atlantic Ocean samples are coloured red and the South Atlantic Ocean samples are coloured yellow.

The Pfm protein families dataset is derived from the metatranscriptomic dataset. The functional analysis which generated the Pfm protein families dataset is outlined in Section 4.4.2. The activity of the organisms within each sample is reflected in the functional composition of transcripts, any changes may indicate a metabolic response to conditions [Klingenberg and Meinicke, 2017]. Therefore the Pfm protein families dataset represented in figure 4.5 is based on similarity of activity level.

The Pfm protein families dataset displayed in figure 4.5 has PC1 accounting for 26.97% of sample variation, while PC2 accounts for 8.4% of sample variation. Overall the Pfm protein families samples are to a reasonable extent clustering well by region, as the matching colours are grouped together. Also for the Pfm protein families samples, there is a general transition of the samples from black to red to yellow, which coincides with how the samples are positioned by latitude as can be seen in figure 3.1a.

4.4.4 Canonical Correspondence Analysis (CCA)

In figure 4.6 we present a CCA for the Pfm protein families dataset. The numbers in the CCA plot correspond to sample site numbers as shown in figure 3.1a. The y -axis represents the CCA2, the x -axis represents the CCA1, and the arrows represent the direction and the length of the vectors for the environmental variables. A CCA enables us to find the relationship between the Pfm protein families and the environmental

variables, and therefore enable us to determine how the environmental variables might determine the response variable values of the Pfm protein families [Paliy and Shankar, 2016].

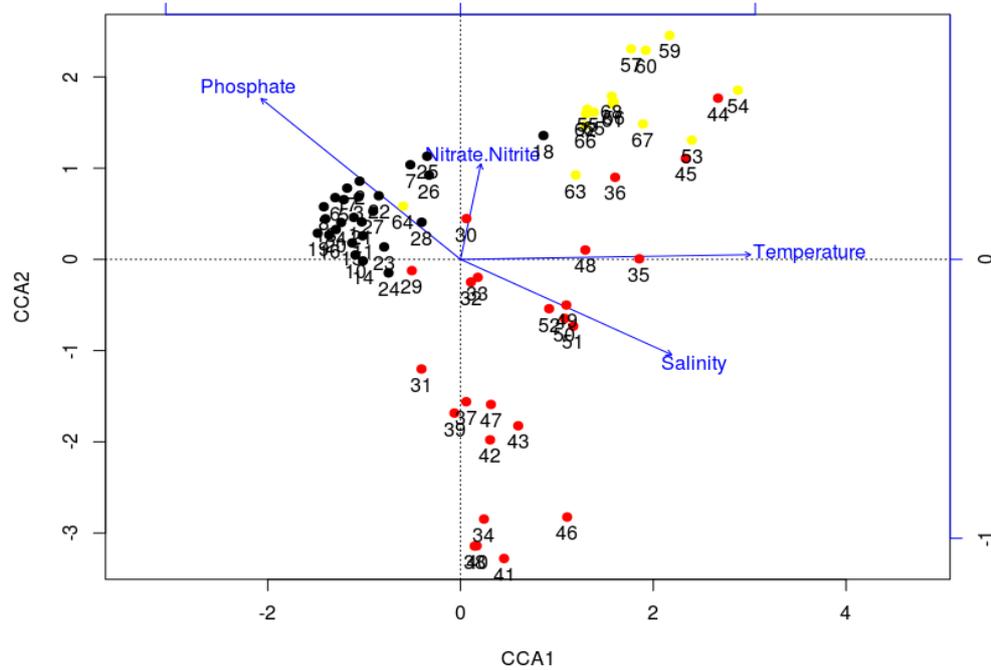


Figure 4.6: A CCA for the Pfm protein families. The numbers in the plots correspond to sample locations as given in figure 3.1a. The numbers are coloured by region, red is for the North Atlantic Ocean, black is for the Arctic Ocean and yellow is for South Atlantic Ocean. We used the R package VEGAN to perform a Canonical Correspondence Analysis (CCA) between the Pfm dataset and the environmental data. The y -axis represents the CCA2 and the x -axis represents the CCA1. The arrows represent the direction and the length of the vector. Each vector represents an environmental factor variable

CCA captured 13.3% of the total variability within the dataset. CCA1 accounts for approximately 45.7% of the constrained variability, with CCA2 accounting for 28.9%, CCA3 accounting for 12.8% and CCA4 accounting for 12.6%. In figure 4.6, the first axis CCA1 is associated with increasing temperature, while the second axis CCA2 is associated with decreasing salinity, increasing nitrate/nitrite and increasing phosphate. The samples are plotted in relation to the arrows, indicating how they are influenced by these environmental variables. The samples from the tropical region of the South Atlantic Ocean, plotted on the right hand side of the figure are strongly influenced by temperature. The samples from the polar region of the Arctic Ocean, plotted on the left hand side of the figure, are poorly influenced by temperature.

4.4.5 Breakpoint analysis

In figure 4.7 we present a breakpoint analysis for the Pfam protein families dataset. The breakpoint analysis was generated using piecewise regression in R as outlined in section 3.3.3. The numbers in the plot correspond to sample location as shown in figure 3.1a. The y -axis represents the beta diversity across the stations, the x -axis represents the temperature and the horizontal line marks the breakpoint. The breakpoint analysis enabled us to investigate how increasing temperature affects the changing diversity of the Pfam protein families from the polar Arctic Ocean through the temperate North Atlantic Ocean and down to the tropical South Atlantic Ocean.

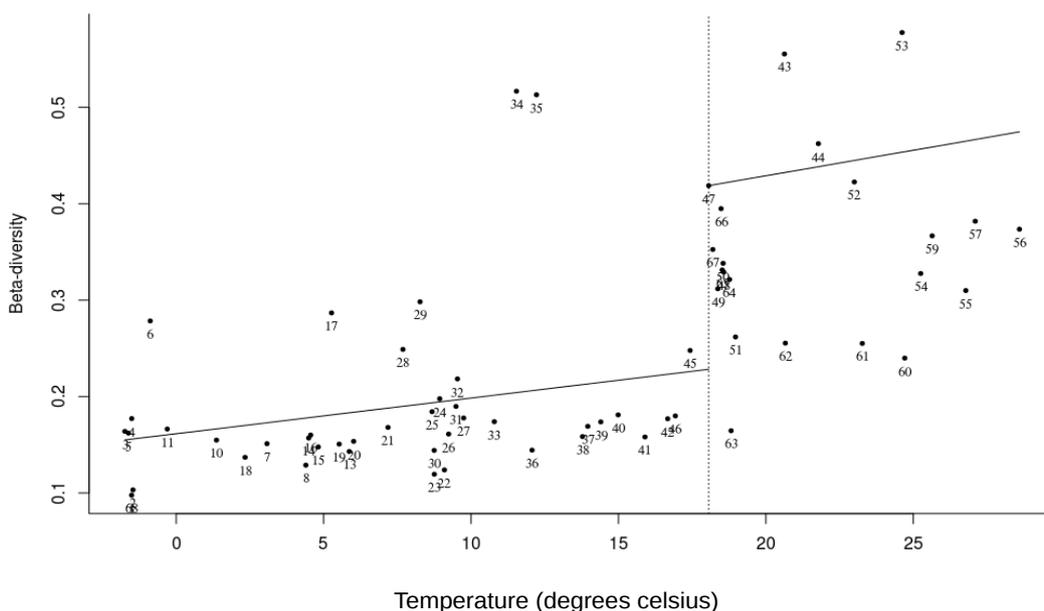


Figure 4.7: A breakpoint analysis for the Pfam protein families. The numbers in the plots correspond to sample locations as given in figure 3.1a. The breakpoint analysis was generated using piecewise regression in R as outlined in section 3.3.3. The y -axis represents the beta diversity across the stations. The x -axis represents the temperature. In the plot, the horizontal line marks the breakpoint. The Pfam protein families breakpoint is 18.06°C with a p -value of $1.24\text{e-}07$

In figure 4.7, the left hand side of the plot represents the lower temperatures containing sample sites from the Arctic Ocean and North Atlantic Ocean. As we move across to the right hand side of the plot, which represents the higher temperatures, we move to sample sites from the North Atlantic Ocean and South Atlantic Ocean. The breakpoint for the Pfam protein families was determined to be 18.06°C with a p -value of $1.24\text{e-}07$. There is a clear shift in the diversity for the Pfam protein families at the

breakpoint. This shift in diversity occurs around the samples located in the temperate region of the North Atlantic Ocean just before we move into the tropical region of the South Atlantic Ocean.

The location of the Pfam protein families breakpoint in the North Atlantic Ocean is consistent with the breakpoint results for the 18S and 16S rDNA datasets. These were also located in the North Atlantic Ocean, with breakpoints determined at 13.96°C (p-value of 3.121e-06) and 9.49°C (p-value of 8.114e-03), respectively (see chapter 3).

In figure 4.2 the heatmap of taxonomically classified sequences belonging to the metatranscriptomic dataset, we observed no gradient of increasing diversity of taxa across the samples as we move from the polar Arctic Ocean through the temperate North Atlantic Ocean and into the tropical South Atlantic Ocean. In figure 4.7 the breakpoint analysis which is based on the Pfam protein families dataset, we observed changes in the diversity of the Pfam as we move across the polar Arctic Ocean through the temperate North Atlantic Ocean and down to the tropical South Atlantic Ocean. The Pfam protein families dataset is derived from the metatranscriptomic dataset. The functional analysis which generated the Pfam protein families dataset is outlined in Section 4.3.2. The activity of the organisms within each sample is reflected in the functional composition of transcripts, any changes may indicate a metabolic response to conditions [Klingenberg and Meinicke, 2017].

4.4.6 Co-occurrence analysis

In our co-occurrence analysis using WGCNA on the Pfam protein family (log10 transformed) gene counts, thirteen modules (networks) were found and a grey module which represents those protein families that could not be assigned to a module. We call the modules black ($n=174$), blue ($n=547$), brown ($n=515$), cyan ($n=83$), green ($n=403$), greenyellow ($n=116$), pink ($n=162$), purple ($n=132$), red ($n=205$), salmon ($n=85$), tan ($n=100$), turquoise ($n=768$), yellow ($n=264$) and grey ($n=7$).

In figure 4.8 we present a correlation heatmap generated between each module's eigengene and environmental parameters. A number of modules were highly correlated, either positively or negatively with the environmental variables. For example, for temperature, the highly correlated modules are tan, blue, turquoise, yellow and pink.

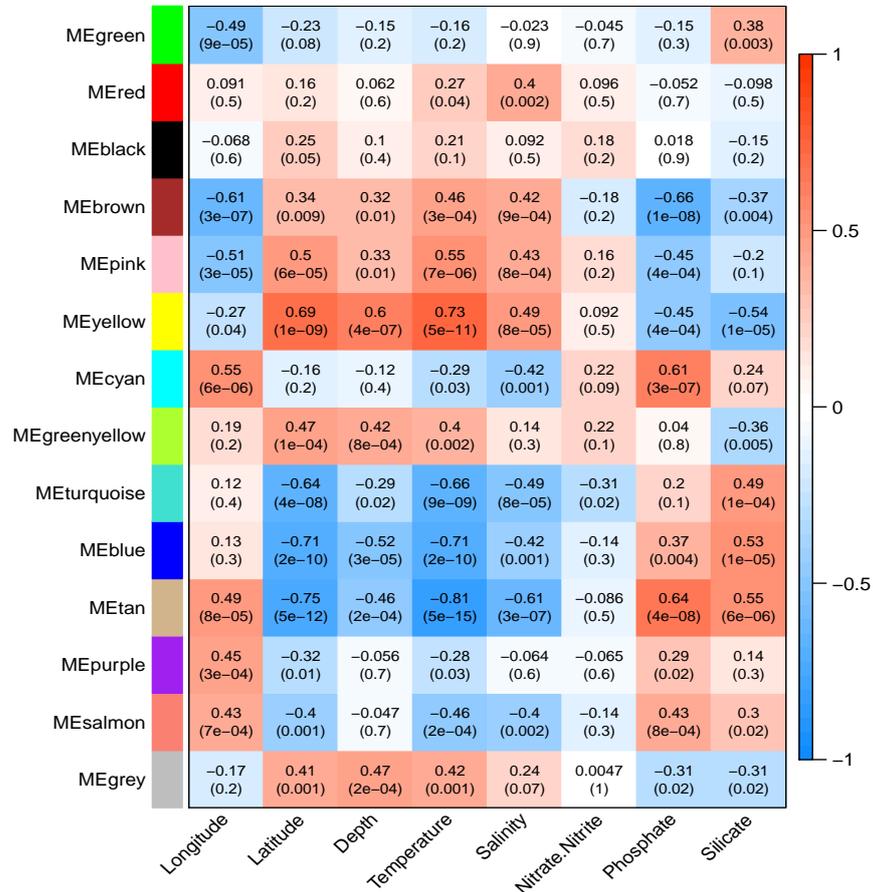


Figure 4.8: In the WGCNA analysis of the log₁₀-scaled gene counts of Pfam protein families, thirteen modules (and a grey module) were found. In the figure, we present a correlation heatmap for the modules. The fourteen modules are displayed as coloured blocks labelled along the left hand side of the plot. The environmental parameters are displayed at the bottom. The colours correspond to the correlation values, red is positively correlated and blue is negatively correlated. The values in each of the squares correspond to the assigned Pearson correlation coefficient value on top and p-value in brackets below

For salinity only the tan module was highly correlated. For phosphate, the highly correlated modules are tan, cyan and brown. For silicate, the highly correlated modules are tan, blue and yellow. No modules correlated significantly with nitrate/nitrite.

The module tan ($n=100$) will be an ideal module for further analysis as it had the highest correlation to the environmental variables of temperature, phosphate and salinity as shown in figure 4.8.

4.5 Discussion

In this chapter, we have presented a metatranscriptomic analysis. We described the dataset and the methodology of the metatranscriptomic analysis, including heatmaps, rarefaction curves, canonical correspondence analysis, breakpoint analysis and co-

occurrence analysis. This was a unique large-scale examination of the protein composition of the dataset, taken from a transect of the Arctic and Atlantic oceans, the sampling of which is described in section 3.2.1. This has given us new insights into what these marine microbial communities are probably doing in response to environmental conditions.

From each sample metatranscriptomic, metagenomic, 18S and 16S rDNA sequencing was performed. As explained in section 3.5, only a single sample was obtained at each station; we did not obtain replicate samples. In addition due to time constraints, a full and extensive analysis of the metatranscriptomic dataset was not performed. Further analysis is still required, as well as a more in-depth examination of the Pfam proteins families co-occurrence analysis.

For our metatranscriptomic analysis, we performed a CCA on the Pfam protein families dataset. With the CCA we captured 13.3% of the total variability in the Pfam protein families dataset, and of this CCA1 accounts for approximately 45.74% of the constrained variability. From the plot in figure 4.6 we identified CCA1 to have an association with increasing temperature for about half of the samples.

We performed a breakpoint analysis on our Pfam protein families dataset and determined the breakpoint to be 18.06°C with a p-value of 1.24e-07. This positions the breakpoint in the temperate region of the North Atlantic Ocean off the coast of France and is in agreement with our 18S and 16S rDNA datasets, as we identified breakpoint of 13.96°C for the 18S rDNA dataset and 9.49°C for the 16S rDNA dataset. These results indicate that as you move from the cold Arctic Ocean to the warm tropical regions in the South Atlantic Ocean there is a radical shift in the diversity of the 18S and 16S rDNA species communities and a radial shift in activity according to the Pfam protein families.

In our co-occurrence analysis with WGCNA on the Pfam protein families dataset, we found thirteen modules for further analysis. A number of modules were highly correlated either positively or negatively with the environmental variables but the tan module ($n=100$) had the highest correlations to the environmental variables (temperature, phosphate and salinity). For future work we will continue the examination of the thirteen modules, beginning with the tan module. For each environmental variable temperature, phosphate and salinity, we will plot gene significance against module

membership, to identify Pfam protein families in the tan module that have a high significance to that environmental variable and analyse them to understand their connection to that environmental variable.

Chapter 5

Discussion and future work

5.1 Summary

In chapter 3 we outlined the computational pipelines and analysis of the 18S rDNA and 16S rDNA datasets that were derived from samples collected from a transect of the Arctic Ocean, North Atlantic Ocean and South Atlantic Ocean. Firstly, this involved constructing a computational pipeline for taxonomically classifying the 18S rDNA dataset. Then we devised a methodology to normalise the 18S and 16S rDNA copy number, in order to interpret the data in terms of species abundance rather than read counts and therefore conduct various analyses that also included the environmental data that was recorded during the expeditions.

From our analysis of the 18S rDNA and 16S rDNA datasets, we observed a greater diversity of microbes in the tropical regions of the South Atlantic Ocean, in comparison to the polar regions of the Arctic Ocean. From analyses that included environmental data, we identified temperature to be the driving force of diversity. Furthermore, a breakpoint analysis was performed on our 18S and 16S rDNA datasets in which we found a shift in diversity occurring in the temperate region of the North Atlantic Ocean, between the polar Arctic Ocean and tropical South Atlantic Ocean. In addition, from our co-occurrence analysis on the 18S and 16S rDNA datasets, we identified two community networks. Each of these networks was found to have a temperature preference, as one positively correlated to temperature and the other negatively correlated to temperature.

In chapter 4 we outlined the computational pipeline of the metatranscriptomic

dataset that was also derived from the samples that were collected from a transect of the Arctic Ocean, North Atlantic Ocean and South Atlantic Ocean. We also outlined the analyses of the Pfam protein families dataset. For our analysis of this dataset, we performed a canonical correspondence analysis (CCA) and we observed which environmental variables and by how much they explained the variation in our dataset. Furthermore, a breakpoint analysis was performed in which we found a shift in diversity occurring in the temperate region of the North Atlantic Ocean, between the polar Arctic Ocean and tropical South Atlantic Ocean. The Pfam protein families breakpoint is in the same region as that of the 18S and 16S rDNA datasets as described in chapter 3. In addition, from our co-occurrence analysis on the Pfam protein families dataset we identified thirteen networks for further analysis.

5.2 Future work

For future work on our 68 samples, we will continue the examination of our metatranscriptomic dataset. As mentioned in chapter 4, we performed a co-occurrence analysis on Pfam protein families and this resulted in several modules being found. These modules will be taken for further analysis and research. We aim to determine Pfam protein families' relationships with one another and how environmental factors may be affecting these relationships. Also, we want to perform a GO enrichment analysis which provides defined GO terms to genes. GO terms cover for the highest level cellular components, molecular functions, and biological processes, and at the lowest level, these can be assigned to genes when relevant. Then we will perform various analyses such as heatmaps in order to examine their composition and distribution.

Our 18S, 16S rDNA and metatranscriptomic analyses have given us new insights into microbial communities, but studies of metatranscriptomes have found that functional genes predicted from metagenomic studies are not necessarily expressed [Heintz-Buschart et al., 2016], [Narayanasamy et al., 2016]. Ideally, we therefore need to integrate metatranscriptomics with metagenomics analysis for a more conclusive link between the genetic potential and the actual phenotype *in situ* [Narayanasamy et al., 2016].

Bearing this in mind, from the 68 stations sampled, we have chosen 11 samples for

metagenomic analysis. These are samples 1, 7, 18, 21, 22 and 23 which come from the Arctic Ocean, sample 45 from the North Atlantic Ocean and samples 54, 57, 59 and 64 from the South Atlantic Ocean, as shown in figure 3.1a. These samples were selected based on our 18S rDNA analysis as we were attempting to avoid samples that contain large amounts of dinoflagellates, as they possess some of the largest nuclear genomes. These samples have been sequenced by JGI. From the metagenomic analysis, this should enable us to understand their genetic potential. These metagenomic samples will be analysed by taxonomic classification and functional analysis. They will then be compared to their corresponding metatranscriptomic samples.

In related work, Emma Langan, a PhD student in the Environmental Sciences department at UEA will be making an expedition to the Antarctic Ocean at the end of December 2018, just before the ocean begins to freeze. The ship will float with the ice during which time a number of samples will be collected periodically, thus enabling the *in situ* real-time sequencing on Oxford nanopore sequencing instrument of polar microbes. The aim of this study is to monitor polar microbes in relation to their composition, distribution and abundance, and therefore monitor how climate change and different oceanographic features may affect these polar microbes. This will directly build on the results presented in this thesis since we examined the composition, distribution and abundance of microbes in the polar regions of the Arctic Ocean, the temperate regions of the North Atlantic Ocean and the tropical regions of South Atlantic Ocean. Sequencing *in situ* overcomes a number of common problems for researchers such as samples degrading. Also, the sampling can be directed based on real-time results rather than best guess of where to sample. Furthermore, if high quality genome assemblies can be achieved from the nanopore long reads, they can be used to produce reference genomes, which will be a great benefit to the research community as there are currently only a few phytoplankton reference genomes available. This will also enable them to perform comparative genomic analysis between temperate and polar species, in order to find the evolutionary mechanisms for polar adaptations.

Recently, more samples have been collected and sequenced by JGI. These samples cover the gap in our data between Cape Town in the South Atlantic Ocean and the Antarctic Ocean. Metatranscriptomic and metagenomic analysis will be performed on these samples. The Antarctic Ocean has a seasonal temperature that rises above

0°C and therefore these polar microbes experience reduced kinetic energy that imposes constraints on cellular activity. But despite this, the Antarctic Ocean contributes high levels of microbial primary production. For example, it has been estimated that while the Antarctic Ocean represents only about 10% of the total surface area of the world's ocean, it contributes about 30% of the global ocean uptake of carbon dioxide. Also, the Antarctic Ocean forms a significant amount of the oceanic food web [Wilkins et al., 2013]. This study will, therefore, provide an important and interesting insight into how these polar microbes live and function and how they compare with the microbes we have analysed in this thesis.

An exciting possibility is that our data from our metatranscriptomic, 18S and 16S rDNA datasets may be used to generate models by T. M. Lenton from the University of Exeter. These models could potentially enable us to forecast how climate change will affect these microbial communities [Toseland et al., 2013].

5.3 Conclusions

High-throughput sequencing technology has enabled us to study species that cannot be grown in the laboratory. In the past, culturing a species was necessary in order to study that species, but due to the advances of high-throughput sequencing we can examine these species, study their expression and genetics and compare and contrast multiple samples in metatranscriptomic and metagenomic analyses. However, there are challenges with implementing these analyses.

The first difficulty we had was preparing the 18S rDNA datasets for analysis. In the NCBI taxonomy not every species has a taxonomy designation at every rank. This caused issues when trying to describe and analyse our data at the taxonomic rank of class. We overcame this problem by placing those that do not have a taxonomic rank of class into a new name to represent them at the class level. This situation occurs because it is possible that there is currently no assignment at a particular taxonomic rank that the species fits into or there is no agreement among scientists into which higher taxonomy a species should be placed. The creation of a temporary name until a permanent rank can be assigned would be one possible solution to this problem.

An important additional issue was caused by the 18S and 16S rDNA copy number.

The relationship between amplicon and species abundance is indeterminate due to rDNA copy number variation within the genomes of different species [Perisin et al., 2016]. While there is a 16S rDNA database available, it is not very large. The 18S rDNA gene has no database available at all. NCBI could provide a solution to this, by creating an 18S rDNA database. NCBI actively collects all kinds of information and would be well positioned to accomplish this.

We were not able to taxonomically classify about a third of our 18S rDNA dataset. We lost a considerable amount of information due to the fact that those species were not in the NCBI database. The 16S rDNA species are much better represented in databases such as SILVA. Only 3% of our 16S rDNA dataset could not be taxonomically classified. This is changing, as there is an increase in the number of research studies that are focusing on eukaryote species.

Another difficulty we had was due to the size of our datasets. We had difficulties with the time it took to analyse our datasets. This was very apparent when we tried to perform a GO enrichment analysis with the tools online such as WEGO 2.0 [Ye et al., 2018]. While there are a number of websites available, very few take multiple files such as in this study. It is therefore important to develop fast new bioinformatics tools to handle large datasets.

Despite these difficulties, we thoroughly analysed the metatranscriptomic, 18S and 16S rDNA datasets. In the surface ocean, we identified temperature to be the important environmental factor that is driving marine microbial diversity and affecting how these microbial communities interact with each other. While diversity increases with increasing temperature, the composition of the microbial communities appear to change once a threshold is reached. The breakpoint analyses of the 16S, 18S rDNA and metatranscriptomic datasets identified breakpoints to be 9.49°C, 13.96°C and 18.06°C, respectively. These breakpoints are located in the North Atlantic Ocean and this might be the location of the threshold for the microbial communities to change. Since microbial community composition on either side of the temperature threshold is different, this may have consequences in terms of how global warming will affect these interactions and therefore have implications for the biogeochemical cycle of elements. As more data is collected and analysed and our understanding of ocean microbes improves, it will be interesting to see how these predictions of our breakpoint analysis will work

out.

Bibliography

- [Aird et al., 2011] Aird, D., Ross, M. G., Chen, W.-S., Danielsson, M., Fennell, T., Russ, C., Jaffe, D. B., Nusbaum, C., and Gnirke, A. (2011). Analyzing and minimizing PCR amplification bias in Illumina sequencing libraries. *Genome Biology*, 12(2):R18.
- [Alemzadeh et al., 2014] Alemzadeh, E., Haddad, R., and Ahmadi, A.-R. (2014). Phytoplanktons and DNA barcoding: Characterization and molecular analysis of phytoplanktons on the Persian Gulf. *Iranian Journal of Microbiology*, 6(4):296–302.
- [Alexander et al., 2015] Alexander, H., Jenkins, B. D., Rynearson, T. A., and Dyhrman, S. T. (2015). Metatranscriptome analyses indicate resource partitioning between diatoms in the field. *Proceedings of the National Academy of Sciences*, 112(17):E2182–E2190.
- [Amit Roy et al., 2014] Amit Roy, S. R., Ray, S., and Roy, A. (2014). Molecular Markers in Phylogenetic Studies-A Review. *Journal of Phylogenetics & Evolutionary Biology*, 02(02):1–9.
- [Anderson and Santana-Garcon, 2015] Anderson, M. J. and Santana-Garcon, J. (2015). Measures of precision for dissimilarity-based multivariate analysis of ecological communities. *Ecology Letters*, 18(1):66–73.
- [Armbrust, 2009] Armbrust, E. V. (2009). The life of diatoms in the world’s oceans. *Nature*, 459(7244):185–192.
- [Aryal et al., 2015] Aryal, S., Karki, G., and Pandey, S. (2015). Microbial Diversity in Freshwater and Marine Environment. *Nepal Journal of Biotechnology*, 3(1):68.

- [Balvočit and Huson, 2017] Balvočit, M. and Huson, D. H. (2017). SILVA, RDP, Greengenes, NCBI and OTT how do these taxonomies compare? *BMC Genomics*, 18(S2):114.
- [Baron, 1996] Baron, S. (1996). *Introduction to Parasitology*. University of Texas Medical Branch at Galveston.
- [Barta, 1997] Barta, J. R. (1997). Investigating Phylogenetic Relationships within the Apicomplexa Using Sequence Data: The Search for Homology. *Methods*, 13(2):81–88.
- [Baselga and Orme, 2012] Baselga, A. and Orme, C. D. L. (2012). betapart : an R package for the study of beta diversity. *Methods in Ecology and Evolution*, 3(5):808–812.
- [Bazinet and Cummings, 2012] Bazinet, A. L. and Cummings, M. P. (2012). A comparative evaluation of sequence classification programs. *BMC Bioinformatics*, 13(1):92.
- [Behjati and Tarpey, 2013] Behjati, S. and Tarpey, P. S. (2013). What is next generation sequencing? *Archives of disease in childhood. Education and practice edition*, 98(6):236–8.
- [Bennett, 2004] Bennett, S. (2004). Solexa Ltd. *Pharmacogenomics*, 5(4):433–438.
- [Bijlsma and Loeschcke, 2005] Bijlsma, R. and Loeschcke, V. (2005). Environmental stress, adaptation and evolution: an overview. *Journal of Evolutionary Biology*, 18(4):744–749.
- [Bik, 2014] Bik, H. M. (2014). Deciphering diversity and ecological function from marine metagenomes. *The Biological Bulletin*, 227(2):107–16.
- [Boon et al., 2014] Boon, E., Meehan, C. J., Whidden, C., Wong, D. H.-J., Langille, M. G. I., and Beiko, R. G. (2014). Interactions in the microbiome: communities of organisms and communities of genes. *FEMS Microbiology Reviews*, 38(1):90–118.

- [Bork et al., 2015] Bork, P., Bowler, C., de Vargas, C., Gorsky, G., Karsenti, E., and Wincker, P. (2015). Tara Oceans studies plankton at planetary scale. *Science*, 348(6237):873–873.
- [Boucher et al., 1991] Boucher, N., Vaulot, D., and Partensky, F. (1991). Flow cytometric determination of phytoplankton DNA in cultures and oceanic populations. *Marine Ecology Progress Series*, 71(1):75–84.
- [Brady and Salzberg, 2011] Brady, A. and Salzberg, S. (2011). PhymmBL expanded: confidence scores, custom databases, parallelization and more. *Nature Methods*, 8(5):367.
- [Brady and Salzberg, 2009a] Brady, A. and Salzberg, S. L. (2009a). Phymm and PhymmBL: metagenomic phylogenetic classification with interpolated Markov models. *Nature Methods*, 6(9):673–6.
- [Brady and Salzberg, 2009b] Brady, A. and Salzberg, S. L. (2009b). Phymm and PhymmBL: metagenomic phylogenetic classification with interpolated Markov models. *Nature Methods*, 6.
- [Brewer et al., 2019] Brewer, T. E., Albertsen, M., Edwards, A., Kirkegaard, R. H., Rocha, E. P. C., and Fierer, N. (2019). Unlinked rRNA genes are widespread among Bacteria and Archaea. *bioRxiv*, page 705046.
- [Brierley, 2017] Brierley, A. S. (2017). Plankton. *Current Biology*, 27(11):R478–R483.
- [Brown, 2014] Brown, J. H. (2014). Why are there so many species in the tropics? *Journal of Biogeography*, 41(1):8–22.
- [Brown, 2002] Brown, T. A. (2002). *Genomes. 2nd edition*. Wiley-Liss.
- [Brussaard et al., 2016] Brussaard, C. P. D., Bidle, K. D., Pedrós-Alió, C., and Legrand, C. (2016). The interactive microbial ocean. *Nature Microbiology*, 2(1):16255.
- [Buermans and den Dunnen, 2014] Buermans, H. and den Dunnen, J. (2014). Next generation sequencing technology: Advances and applications. *Biochimica et Biophysica Acta (BBA) - Molecular Basis of Disease*, 1842(10):1932–1941.

- [Campra and Morales, 2016] Campra, P. and Morales, M. (2016). Trend analysis by a piecewise linear regression model applied to surface air temperatures in Southeastern Spain. (May):1–25.
- [Capella-Gutierrez et al., 2009] Capella-Gutierrez, S., Silla-Martinez, J. M., and Gabaldon, T. (2009). trimAl: a tool for automated alignment trimming in large-scale phylogenetic analyses. *Bioinformatics*, 25(15):1972–1973.
- [Carlton et al., 2013] Carlton, J. M., Perkins, S. L., and Deitsch, K. W. (2013). *Malaria parasites : comparative genomics, evolution and molecular biology*.
- [Carradec et al., 2018] Carradec, Q., Pelletier, E., Da Silva, C., Alberti, A., Seeleuthner, Y., Blanc-Mathieu, R., Lima-Mendez, G., Rocha, F., Tirichine, L., Labadie, K., Kirilovsky, A., Bertrand, A., Engelen, S., Madoui, M.-A., Méheust, R., Poulain, J., Romac, S., Richter, D. J., Yoshikawa, G., Dimier, C., Kandels-Lewis, S., Picheral, M., Searson, S., Jaillon, O., Aury, J.-M., Karsenti, E., Sullivan, M. B., Sunagawa, S., Bork, P., Not, F., Hingamp, P., Raes, J., Guidi, L., Ogata, H., de Vargas, C., Iudicone, D., Bowler, C., and Wincker, P. (2018). A global ocean atlas of eukaryotic genes. *Nature Communications*, 9(1):373.
- [Carter and Hussain, 2017] Carter, J.-M. and Hussain, S. (2017). Robust long-read native DNA sequencing using the ONT CsgG Nanopore system. *Wellcome Open Research*, 2:23.
- [Castro-Insua et al., 2016] Castro-Insua, A., Gómez-Rodríguez, C., and Baselga, A. (2016). Break the pattern: breakpoints in beta diversity of vertebrates are general across clades and suggest common historical causes. *Global Ecology and Biogeography*, 25(11):1279–1283.
- [Chapman et al., 2018] Chapman, A. S. A., Tunnicliffe, V., and Bates, A. E. (2018). Both rare and common species make unique contributions to functional diversity in an ecosystem unaffected by human activities. *Diversity and Distributions*, 24(5):568–578.

- [Collins et al., 2014] Collins, S., Rost, B., and Ryneerson, T. A. (2014). Evolutionary potential of marine phytoplankton under ocean acidification. *Evolutionary Applications*, 7(1):140–155.
- [Cooper, 2000] Cooper, G. M. (2000). *The Cell: A Molecular Approach*. Sinauer Associates.
- [Cordero and Datta, 2016] Cordero, O. X. and Datta, M. S. (2016). Microbial interactions and community assembly at microscales. *Current Opinion in Microbiology*, 31:227–234.
- [Das et al., 2006] Das, S., Lyla, P. S., and Khan, S. A. (2006). Marine microbial diversity and ecology: importance and future perspectives.
- [de Vargas et al., 2015] de Vargas, C., Audic, S., Henry, N., Decelle, J., Mahe, F., Logares, R., Lara, E., Berney, C., Le Bescot, N., Probert, I., Carmichael, M., Poulain, J., Romac, S., Colin, S., Aury, J.-M., Bittner, L., Chaffron, S., Dunthorn, M., Engelen, S., Flegontova, O., Guidi, L., Horak, A., Jaillon, O., Lima-Mendez, G., Lukeš, J., Malviya, S., Morard, R., Mulo, M., Scalco, E., Siano, R., Vincent, F., Zingone, A., Dimier, C., Picheral, M., Searson, S., Kandels-Lewis, S., Acinas, S. G., Bork, P., Bowler, C., Gorsky, G., Grimsley, N., Hingamp, P., Iudicone, D., Not, F., Ogata, H., Pesant, S., Raes, J., Sieracki, M. E., Speich, S., Stemmann, L., Sunagawa, S., Weisenbach, J., Wincker, P., Karsenti, E., Boss, E., Follows, M., Karp-Boss, L., Krzic, U., Reynaud, E. G., Sardet, C., Sullivan, M. B., and Velayoudon, D. (2015). Eukaryotic plankton diversity in the sunlit ocean. *Science*, 348(6237):1261605–1261605.
- [DOE Joint Genome Institute, 2017] DOE Joint Genome Institute (2017). BBDuk Guide - DOE Joint Genome Institute.
- [Dyomin et al., 2016] Dyomin, A. G., Koshel, E. I., Kiselev, A. M., Saifitdinova, A. F., Galkina, S. A., Fukagawa, T., Kostareva, A. A., and Gaginskaya, E. R. (2016). Chicken rRNA Gene Cluster Structure. *PLOS ONE*, 11(6):e0157464.
- [Eddy, 1996] Eddy, S. R. (1996). Hidden Markov models. *Current Opinion in Structural Biology*, 6(3):361–365.

- [Eddy, 2004] Eddy, S. R. (2004). What is a hidden Markov model? *Nature Biotechnology*, 22(10):1315–1316.
- [Eddy, 2009] Eddy, S. R. (2009). A new generation of homology search tools based on probabilistic inference. *Genome informatics. International Conference on Genome Informatics*, 23(1):205–11.
- [Edgar, 2010a] Edgar, R. (2010a). *USEARCH cluster otus*.
https://www.drive5.com/usearch/manual/cmd_cluster_otus.html.
- [Edgar, 2010b] Edgar, R. C. (2010b). Search and clustering orders of magnitude faster than BLAST. *Bioinformatics*, 26(19):2460–2461.
- [Edgar, 2010c] Edgar, R. C. (2010c). *UPARSE-OTU algorithm*.
- [Edgar, 2010d] Edgar, R. C. (2010d). *UPARSE-REF algorithm*.
- [Edgar, 2010e] Edgar, R. C. (2010e). *USEARCH Merging paired reads*.
https://www.drive5.com/usearch/manual/merge_pair.html.
- [Edgar, 2010f] Edgar, R. C. (2010f). *USEARCH oligodb*.
https://www.drive5.com/usearch/manual/cmd_search_oligodb.html.
- [Edgar, 2016] Edgar, R. C. (2016). UCHIME2: improved chimera prediction for amplicon sequencing. *doi: 10.1101/074252*.
- [Edgar, 2017] Edgar, R. C. (2017). Updating the 97% identity threshold for 16S ribosomal RNA OTUs. *bioRxiv*, page 192211.
- [Eisenberg and Levanon, 2013] Eisenberg, E. and Levanon, E. Y. (2013). Human housekeeping genes, revisited. *Trends in Genetics : TIG*, 29(10):569–74.
- [El-Gebali et al., 2019] El-Gebali, S., Mistry, J., Bateman, A., Eddy, S. R., Luciani, A., Potter, S. C., Qureshi, M., Richardson, L. J., Salazar, G. A., Smart, A., Sonnhammer, E. L., Hirsh, L., Paladin, L., Piovesan, D., Tosatto, S. C., and Finn, R. D. (2019). The Pfam protein families database in 2019. *Nucleic Acids Research*, 47(D1):D427–D432.

- [Eric et al., 2014] Eric, S. D., Nicholas, T. K. D. D., and Theophilus, K. A. (2014). Bioinformatics with basic local alignment search tool (BLAST) and fast alignment (FASTA). *Journal of Bioinformatics and Sequence Analysis*, 6(1):1–6.
- [Falkowski et al., 2008] Falkowski, P. G., Fenchel, T., and Delong, E. F. (2008). The Microbial Engines That Drive Earth’s Biogeochemical Cycles. *Science*, 320(5879):1034–1039.
- [Fernández-Pérez et al., 2018] Fernández-Pérez, J., Nantón, A., and Méndez, J. (2018). Sequence characterization of the 5S ribosomal DNA and the internal transcribed spacer (ITS) region in four European Donax species (Bivalvia: Donacidae). *BMC Genetics*, 19(1):97.
- [Ferreira et al., 2014] Ferreira, M., Roma, N., and Russo, L. M. S. (2014). Cache-Oblivious parallel SIMD Viterbi decoding for sequence search in HMMER. *BMC Bioinformatics*, 15(1):165.
- [Finn et al., 2011] Finn, R. D., Clements, J., and Eddy, S. R. (2011). HMMER web server: interactive sequence similarity searching. *Nucleic Acids Research*, 39(Web Server issue):W29–37.
- [Finn et al., 2016] Finn, R. D., Coggill, P., Eberhardt, R. Y., Eddy, S. R., Mistry, J., Mitchell, A. L., Potter, S. C., Punta, M., Qureshi, M., Sangrador-Vegas, A., Salazar, G. A., Tate, J., and Bateman, A. (2016). The Pfam protein families database: towards a more sustainable future. *Nucleic Acids Research*, 44(D1):D279–D285.
- [Flombaum et al., 2013] Flombaum, P., Gallegos, J. L., Gordillo, R. A., Rincón, J., Zabala, L. L., Jiao, N., Karl, D. M., Li, W. K. W., Lomas, M. W., Veneziano, D., Vera, C. S., Vrugt, J. A., and Martiny, A. C. (2013). Present and future global distributions of the marine Cyanobacteria Prochlorococcus and Synechococcus. *Proceedings of the National Academy of Sciences of the United States of America*, 110(24):9824–9.
- [Foster et al., 2011] Foster, R. A., Kuypers, M. M. M., Vagner, T., Paerl, R. W., Musat, N., and Zehr, J. P. (2011). Nitrogen fixation and transfer in open ocean diatom-cyanobacterial symbioses. *The ISME Journal*, 5(9):1484–93.

- [Franco et al., 2017] Franco, D. C., Signori, C. N., Duarte, R. T. D., Nakayama, C. R., Campos, L. S., and Pellizari, V. H. (2017). High Prevalence of Gammaproteobacteria in the Sediments of Admiralty Bay and North Bransfield Basin, Northwestern Antarctic Peninsula. *Frontiers in Microbiology*, 8:153.
- [Fu et al., 2012] Fu, L., Niu, B., Zhu, Z., Wu, S., and Li, W. (2012). CD-HIT: accelerated for clustering the next-generation sequencing data. *Bioinformatics*, 28(23):3150–3152.
- [Fu and Gong, 2017] Fu, R. and Gong, J. (2017). Single Cell Analysis Linking Ribosomal (r)DNA and rRNA Copy Numbers to Cell Size and Growth Rate Provides Insights into Molecular Protistan Ecology. *Journal of Eukaryotic Microbiology*, 64(6):885–896.
- [Gibbons et al., 2014] Gibbons, J. G., Branco, A. T., Yu, S., and Lemos, B. (2014). Ribosomal DNA copy number is coupled with gene expression variation and mitochondrial abundance in humans. *Nature Communications*, 5:1–12.
- [Gilbert and Hughes, 2011] Gilbert, J. A. and Hughes, M. (2011). Gene Expression Profiling: Metatranscriptomics. In *Methods in Molecular Biology (Clifton, N.J.)*, volume 733, pages 195–205.
- [Godfray and May, 2014] Godfray, H. C. J. and May, R. M. (2014). Open questions: are the dynamics of ecological communities predictable? *BMC Biology*, 12:22.
- [Godhe et al., 2008] Godhe, A., Asplund, M. E., Härnström, K., Saravanan, V., Tyagi, A., and Karunasagar, I. (2008). Quantification of diatom and dinoflagellate biomasses in coastal marine seawater samples by real-time PCR. *Applied and Environmental Microbiology*, 74(23):7174–82.
- [Goodrich et al., 2014] Goodrich, J. K., Di Rienzi, S. C., Poole, A. C., Koren, O., Walters, W. A., Caporaso, J. G., Knight, R., and Ley, R. E. (2014). Conducting a microbiome study. *Cell*, 158(2):250–262.
- [Gough et al., 2001] Gough, J., Karplus, K., Hughey, R., and Chothia, C. (2001). Assignment of homology to genome sequences using a library of hidden Markov mod-

els that represent all proteins of known structure. *Journal of Molecular Biology*, 313(4):903–919.

[Guiry, M.D. & Guiry, 2008] Guiry, M.D. & Guiry, G. (2008). AlgaeBase.

[Haque and Haque, 2017] Haque, S. Z. and Haque, M. (2017). The ecological community of commensal, symbiotic, and pathogenic gastrointestinal microorganisms - an appraisal. *Clinical and Experimental Gastroenterology*, 10:91–103.

[Hauser et al., 2010] Hauser, P. M., Burdet, F. X., Cissé, O. H., Keller, L., Taffé, P., Sanglard, D., Pagni, M., Jr, C. T., Limper, A., Jr, C. T., Limper, A., Davis, J., Fei, M., Huang, L., Wakefield, A., Demanche, C., Berthelemy, M., Petit, T., Polack, B., Wakefield, A., Wakefield, A., Stringer, J., Tamburrin, E., Dei-Cas, E., Gigliotti, F., Harmsen, A., Haidaris, C., Haidaris, P., Aliouat-Denis, C., Chabé, M., Demanche, C., El, M. A., Viscogliosi, E., Keely, S., Renauld, H., Wakefield, A., Cushion, M., Smulian, A., Kutty, G., Maldarelli, F., Achaz, G., Kovacs, J., Joffrion, T., Cushion, M., Kaneshiro, E., Rodrigues, M., Fonseca, A., Keeling, P., Fast, N., Law, J., Williams, B., Slamovits, C., Payne, S., Loomis, W., Gardner, M., Shallom, S., Carlton, J., Salzberg, S., Nene, V., Omsland, A., Cockrell, D., Howe, D., Fischer, E., Virtaneva, K., Ewann, F., Hoffman, P., Stanke, M., Schöffmann, O., Morgenstern, B., Waack, S., Sugiyama, J., Cushion, M., Smulian, A., Slaven, B., Sesterhenn, T., Arnold, J., Andersson, J., Andersson, S., Sakharkar, K., Dhar, P., Chow, V., Cushion, M., Corradi, N., Pombert, J., Farinelli, L., Didier, E., Keeling, P., Nahimana, A., Francioli, P., Blanc, D., Bille, J., Wakefield, A., Choi, M., Chung, B., Chung, Y., Yu, J., Cho, S., Ambrose, H., Keely, S., Aliouat, E., Dei-Cas, E., Wakefield, A., Basselin, M., Qiu, Y., Lipscomb, K., Kaneshiro, E., Reichard, U., Jousson, O., Monod, M., Atzori, C., Angeli, E., Mainini, A., Agoston, F., Micheli, V., Stringer, J., Cushion, M., Nowrousian, M., Nowrousian, M., Stajich, J., Chu, M., Engh, I., Espagne, E., Korf, I., Birney, E., Clamp, M., Durbin, R., Ter-Hovhannisyan, V., Lomsadze, A., Chernoff, Y., Borodovsky, M., Cantarel, B., Korf, I., Robb, S., Parra, G., Ross, E., Altschul, S., Madden, T., Schäffer, A., Zhang, J., Zhang, Z., Hau, J., Muller, M., Pagni, M., Okuda, S., Yamada, T., Hamajima, M., Itoh, M., Katayama, T., Schnoes, A., Brown, S., Dodevski, I., and Babbitt, P. (2010). Comparative Ge-

- nomics Suggests that the Fungal Pathogen *Pneumocystis* Is an Obligate Parasite Scavenging Amino Acids from Its Host's Lungs. *PLoS ONE*, 5(12):e15152.
- [Heather and Chain, 2016] Heather, J. M. and Chain, B. (2016). The sequence of sequencers: The history of sequencing DNA. *Genomics*, 107(1):1–8.
- [Heintz-Buschart et al., 2016] Heintz-Buschart, A., May, P., Laczny, C. C., Lebrun, L. A., Bellora, C., Krishna, A., Wampach, L., Schneider, J. G., Hogan, A., de Beaufort, C., and Wilmes, P. (2016). Integrated multi-omics of the human gut microbiome in a case study of familial type 1 diabetes. *Nature Microbiology*, 2:16180.
- [Holland and Bronstein, 2008] Holland, J. and Bronstein, J. (2008). Mutualism. *Encyclopedia of Ecology*, pages 2485–2491.
- [Hou et al., 2017] Hou, D., Huang, Z., Zeng, S., Liu, J., Wei, D., Deng, X., Weng, S., He, Z., and He, J. (2017). Environmental Factors Shape Water Microbial Community Structure and Function in Shrimp Cultural Enclosure Ecosystems. *Frontiers in Microbiology*, 8:2359.
- [Hugerth and Andersson, 2017] Hugerth, L. W. and Andersson, A. F. (2017). Analysing Microbial Community Composition through Amplicon Sequencing: From Sampling to Hypothesis Testing. *Frontiers in Microbiology*, 8:1561.
- [Hughes et al., 2001] Hughes, J. B., Hellmann, J. J., Ricketts, T. H., and Bohannan, B. J. (2001). Counting the uncountable: statistical approaches to estimating microbial diversity. *Applied and Environmental Microbiology*, 67(10):4399–406.
- [Hunt et al., 2013] Hunt, D. E., Lin, Y., Church, M. J., Karl, D. M., Tringe, S. G., Izzo, L. K., and Johnson, Z. I. (2013). Relationship between abundance and specific activity of bacterioplankton in open ocean surface waters. *Applied and Environmental Microbiology*, 79(1):177–84.
- [Huntemann et al., 2016] Huntemann, M., Ivanova, N. N., Mavromatis, K., Tripp, H. J., Paez-Espino, D., Tennessen, K., Palaniappan, K., Szeto, E., Pillay, M., Chen, I.-M. A., Pati, A., Nielsen, T., Markowitz, V. M., and Kyrpides, N. C. (2016). The standard operating procedure of the DOE-JGI Metagenome Annotation Pipeline (MAP v.4). *Standards in Genomic Sciences*, 11(1):17.

- [Huson et al., 2007] Huson, D. H., Auch, A. F., Qi, J., and Schuster, S. C. (2007). MEGAN analysis of metagenomic data. *Genome Research*, 17(3):377–86.
- [Huson et al., 2011] Huson, D. H., Mitra, S., Ruscheweyh, H.-J., Weber, N., and Schuster, S. C. (2011). Integrative analysis of environmental sequences using MEGAN4. *Genome Research*, 21(9):1552–60.
- [Huson et al., 2009] Huson, D. H., Richter, D. C., Mitra, S., Auch, A. F., and Schuster, S. C. (2009). Methods for comparative metagenomics. *BMC Bioinformatics*, 10(Suppl 1):S12.
- [Iwen and Hinrichs, 2002] Iwen, P. C. and Hinrichs, S. H. (2002). Utilization of the internal transcribed spacer regions as molecular targets to detect and identify human fungal. (July 2001):87–109.
- [Jain et al., 2016] Jain, M., Olsen, H. E., Paten, B., and Akeson, M. (2016). The Oxford Nanopore MinION: delivery of nanopore sequencing to the genomics community. *Genome Biology*, 17(1):239.
- [Jiang et al., 2016a] Jiang, Y., Xiong, X., Danska, J., and Parkinson, J. (2016a). Metatranscriptomic analysis of diverse microbial communities reveals core metabolic pathways and microbiome-specific functionality. *Microbiome*, 4(1):2.
- [Jiang et al., 2016b] Jiang, Y., Xiong, X., Danska, J., and Parkinson, J. (2016b). Metatranscriptomic analysis of diverse microbial communities reveals core metabolic pathways and microbiome-specific functionality. *Microbiome*, 4(1):2.
- [Johnson and Burnet, 2016] Johnson, K. V.-A. and Burnet, P. W. J. (2016). Microbiome: Should we diversify from diversity? *Gut Microbes*, 7(6):455–458.
- [Johnston, 2006] Johnston, J. S. (2006). Characteristics of the nuclear (18S, 5.8S, 28S and 5S) and mitochondrial (12S and 16S) rRNA genes of. *Insect Molecular Biology*, 15:657–686.
- [Joint et al., 2010] Joint, I., Mühling, M., and Querellou, J. (2010). Culturing marine bacteria - an essential prerequisite for biodiscovery. *Microbial Biotechnology*, 3(5):564–75.

- [Kalendar et al., 2017] Kalendar, R., Khassenov, B., Ramankulov, Y., Samuilova, O., and Ivanov, K. I. (2017). FastPCR: An in silico tool for fast primer and probe design and advanced sequence analysis. *Genomics*, 109(3-4):312–319.
- [Kazamia et al., 2016] Kazamia, E., Helliwell, K. E., Purton, S., and Smith, A. G. (2016). How mutualisms arise in phytoplankton communities: building eco-evolutionary principles for aquatic microbes. *Ecology Letters*, 19(7):810–822.
- [Keeling et al., 2014] Keeling, P. J., Burki, F., Wilcox, H. M., Allam, B., Allen, E. E., Amaral-Zettler, L. A., Armbrust, E. V., Archibald, J. M., Bharti, A. K., Bell, C. J., Beszteri, B., Bidle, K. D., Cameron, C. T., Campbell, L., Caron, D. A., Cattolico, R. A., Collier, J. L., Coyne, K., Davy, S. K., Deschamps, P., Dyhrman, S. T., Edvardsen, B., Gates, R. D., Gobler, C. J., Greenwood, S. J., Guida, S. M., Jacobi, J. L., Jakobsen, K. S., James, E. R., Jenkins, B., John, U., Johnson, M. D., Juhl, A. R., Kamp, A., Katz, L. A., Kiene, R., Kudryavtsev, A., Leander, B. S., Lin, S., Lovejoy, C., Lynn, D., Marchetti, A., McManus, G., Nedelcu, A. M., Menden-Deuer, S., Miceli, C., Mock, T., Montresor, M., Moran, M. A., Murray, S., Nadathur, G., Nagai, S., Ngam, P. B., Palenik, B., Pawlowski, J., Petroni, G., Piganeau, G., Posewitz, M. C., Rengefors, K., Romano, G., Rumpho, M. E., Ryneerson, T., Schilling, K. B., Schroeder, D. C., Simpson, A. G. B., Slamovits, C. H., Smith, D. R., Smith, G. J., Smith, S. R., Sosik, H. M., Stief, P., Theriot, E., Twary, S. N., Umale, P. E., Vaulot, D., Wawrik, B., Wheeler, G. L., Wilson, W. H., Xu, Y., Zingone, A., and Worden, A. Z. (2014). The Marine Microbial Eukaryote Transcriptome Sequencing Project (MMETSP): Illuminating the Functional Diversity of Eukaryotic Life in the Oceans through Transcriptome Sequencing. *PLoS Biology*, 12(6):e1001889.
- [Keightley and Johnson, 2004] Keightley, P. D. and Johnson, T. (2004). MCALIGN: stochastic alignment of noncoding DNA sequences based on an evolutionary model of sequence evolution. *Genome research*, 14(3):442–50.
- [Kemena and Notredame, 2009] Kemena, C. and Notredame, C. (2009). Upcoming challenges for multiple sequence alignment methods in the high-throughput era. *Bioinformatics (Oxford, England)*, 25(19):2455–65.

- [Klappenbach et al., 2001] Klappenbach, J. A., Saxman, P. R., Cole, J. R., and Schmidt, T. M. (2001). rrndb: the Ribosomal RNA Operon Copy Number Database. *Nucleic Acids Research*, 29(1):181–4.
- [Klingenberg and Meinicke, 2017] Klingenberg, H. and Meinicke, P. (2017). How to normalize metatranscriptomic count data for differential expression analysis. *PeerJ*, 5:e3859.
- [Kozarewa et al., 2009] Kozarewa, I., Ning, Z., Quail, M. A., Sanders, M. J., Berri-man, M., and Turner, D. J. (2009). Amplification-free Illumina sequencing-library preparation facilitates improved mapping and assembly of (G+C)-biased genomes. *Nature Methods*, 6(4):291–5.
- [Lambrechts et al., 2019] Lambrechts, S., Willems, A., and Tahon, G. (2019). Uncov-ering the Uncultivated Majority in Antarctic Soils: Toward a Synergistic Approach. *Frontiers in Microbiology*, 10:242.
- [Langfelder and Horvath, 2007] Langfelder, P. and Horvath, S. (2007). Eigengene net-works for studying the relationships between co-expression modules. *BMC Systems Biology*, 1:54.
- [Langfelder et al., 2008] Langfelder, P., Horvath, S., Fisher, R., Zhou, X., Kao, M., Wong, W., Steffen, M., Petti, A., Aach, J., D’haeseleer, P., Church, G., Stuart, J., Segal, E., Koller, D., Kim, S., Zhang, B., Horvath, S., Carey, V., Gentry, J., Whalen, E., Gentleman, R., Schaefer, J., Strimmer, K., Chuang, C., Jen, C., Chen, C., Shieh, G., Cokus, S., Rose, S., Haynor, D., Gronbech-Jensen, N., Pellegrini, M., Horvath, S., Zhang, B., Carlson, M., Lu, K., Zhu, S., Felciano, R., Laurance, M., Zhao, W., Shu, Q., Lee, Y., Scheck, A., Liau, L., Wu, H., Geschwind, D., Febbo, P., Kornblum, H., Cloughesy, T., Nelson, S., Mischel, P., Horvath, S., Dong, J., Langfelder, P., Horvath, S., Carlson, M., Zhang, B., Fang, Z., Horvath, S., Mishel, P., Nelson, S., Ghazalpour, A., Doss, S., Zhang, B., Plaisier, C., Wang, S., Schadt, E., Thomas, A., Drake, T., Lusi, A., Horvath, S., Fuller, T., Ghazalpour, A., Aten, J., Drake, T., Lusi, A., Horvath, S., Emilsson, V., Thorleifsson, G., Zhang, B., Leonardson, A., Zink, F., Zhu, J., Carlson, S., Helgason, A., Walters, G., Gunnarsdottir, S., Mouy, M., Steinthorsdottir, V., Eiriksdottir, G., Bjornsdottir, G., Reynisdottir, I.,

Gudbjartsson, D., Helgadóttir, A., Jonasdóttir, A., Jonasdóttir, A., Styrkarsdóttir, U., Gretarsdóttir, S., Magnusson, K., Stefansson, H., Fossdal, R., Kristjansson, K., Gislason, H., Stefansson, T., Leifsson, B., Thorsteinsdóttir, U., Lamb, J., Gulcher, M., null Reitman, Kong, A., Schadt, E., Stefansson, K., van Nas, A., Guhathakurta, D., Wang, S., Yehya, S., Horvath, S., Zhang, B., Drake, L. I., Chaudhuri, G., Schadt, E., Drake, T., Arnold, A., Lusic, A., Oldham, M., Horvath, S., Geschwind, D., Miller, J., Oldham, M., Geschwind, D., Oldham, M., Konopka, G., Iwamoto, K., Langfelder, P., Kato, T., Horvath, S., Geschwind, D., Keller, M., Choi, Y., Wang, P., Davis, D. B., Rabaglia, M., Oler, A., Stapleton, D., Argmann, C., Schueler, K., Edwards, S., Steinberg, H., Neto, E. C., Kleinhanz, R., Turner, S., Hellerstein, M., Schadt, E., Yandell, B., Kendzioriski, C., Attie, A., Presson, A., Sobel, E., Papp, J., Suarez, C., Whistler, T., Rajeevan, M., Vernon, S., Horvath, S., Weston, D., Gunter, L., Rogers, A., Wullschleger, S., Wilcox, R., Yip, A., Horvath, S., Ravasz, E., Somera, A., Mongru, D., Oltvai, Z., Barabási, A., Li, A., Horvath, S., Kaufman, L., Rousseeuw, P., Langfelder, P., Zhang, B., Horvath, S., Dudoit, S., Fridlyand, J., Hastie, T., Tibshirani, R., Sherlock, G., Eisen, M., Brown, P., Botstein, D., Troyanskaya, O., Cantor, M., Sherlock, G., Brown, P., Hastie, T., Tibshirani, R., Botstein, D., Altman, R., Dong, J., Horvath, S., Watts, D., Strogatz, S., Dudoit, S., Yang, Y., Callow, M., Speed, T., Hu, Z., Snitkin, E., DeLisi, C., Shannon, P., Markiel, A., Ozier, O., Baliga, N., Wang, J., Ramage, D., Amin, N., Schwikowski, B., Ideker, T., Frohlich, H., Speer, N., Poustka, A., BeiSZbarth, T., Dennis, G., Sherman, B., Hosack, D., Yang, J., Gao, W., Lane, H., Lempicki, R., Ashburner, M., Ball, C., Blake, J., Botstein, D., Butler, H., Cherry, J., Davis, A., Dolinski, K., Dwight, S., Eppig, J., Harris, M., Hill, D., Issel-Tarver, L., Kasarskis, A., Lewis, S., Matese, J., Richardson, J., Ringwald, M., Rubin, G., Sherlock, G., Zhang, B., Kirov, S., Snoddy, J., Liu, M., Liberzon, A., Kong, S., Lai, W., Park, P., Kohane, I., Kasif, S., Henegar, C., Clement, K., Zucker, J., Gentleman, R., Huber, W., Carey, V., Irizarry, R., Dudoit, S., Opgen-Rhein, R., Strimmer, K., Aten, J., Fuller, T., Lusic, A., Horvath, S., Neto, E. C., Ferrara, C., Attie, A., and Yandell, B. (2008). WGCNA: an R package for weighted correlation network analysis. *BMC Bioinformatics*, 9(1):559.

- [Laver et al., 2015] Laver, T., Harrison, J., O’Neill, P., Moore, K., Farbos, A., Paszkiewicz, K., and Studholme, D. (2015). Assessing the performance of the Oxford Nanopore Technologies MinION. *Biomolecular Detection and Quantification*, 3:1–8.
- [Leimena et al., 2013] Leimena, M. M., Ramiro-Garcia, J., Davids, M., van den Bogert, B., Smidt, H., Smid, E. J., Boekhorst, J., Zoetendal, E. G., Schaap, P. J., and Kleerebezem, M. (2013). A comprehensive metatranscriptome analysis pipeline and its validation using human small intestine microbiota datasets. *BMC Genomics*, 14(1):530.
- [Letunic and Bork, 2007] Letunic, I. and Bork, P. (2007). Interactive Tree Of Life (iTOL): an online tool for phylogenetic tree display and annotation. *Bioinformatics*, 23(1):127–128.
- [Li et al., 2011] Li, M., Copeland, A., and Han, J. (2011). DUK A Fast and Efficient Kmer Matching Tool DUK A Fast and Efficient Kmer Based Sequence Matching Tool.
- [Li et al., 2012] Li, W., Fu, L., Niu, B., Wu, S., and Wooley, J. (2012). Ultrafast clustering algorithms for metagenomic sequence analysis. *Briefings in Bioinformatics*, 13(6):656–68.
- [Li and Godzik, 2006] Li, W. and Godzik, A. (2006). Cd-hit: a fast program for clustering and comparing large sets of protein or nucleotide sequences. *Bioinformatics*, 22(13):1658–1659.
- [Lin et al., 2012] Lin, Y.-C., Campbell, T., Chung, C.-C., Gong, G.-C., Chiang, K.-P., and Worden, A. Z. (2012). Distribution patterns and phylogeny of marine stramenopiles in the north pacific ocean. *Applied and Environmental Microbiology*, 78(9):3387–99.
- [Liu et al., 2012] Liu, L., Li, Y., Li, S., Hu, N., He, Y., Pong, R., Lin, D., Lu, L., and Law, M. (2012). Comparison of Next-Generation Sequencing Systems. *Journal of Biomedicine and Biotechnology*, 2012:1–11.

- [Lu et al., 2016] Lu, H., Giordano, F., and Ning, Z. (2016). Oxford Nanopore MinION Sequencing and Genome Assembly. *Genomics, Proteomics & Bioinformatics*, 14(5):265–279.
- [MacArthur and MacArthur, 1961] MacArthur, R. H. and MacArthur, J. W. (1961). On Bird Species Diversity. *Ecology*, 42(3):594–598.
- [Magoč and Salzberg, 2011] Magoč, T. and Salzberg, S. L. (2011). FLASH: fast length adjustment of short reads to improve genome assemblies. *Bioinformatics (Oxford, England)*, 27(21):2957–63.
- [Magoč and Salzberg, 2011] Magoč, T. and Salzberg, S. L. (2011). FLASH: fast length adjustment of short reads to improve genome assemblies. *Bioinformatics*, 27(21):2957–2963.
- [Mandal, 1984] Mandal, R. K. (1984). The Organization and Transcription of Eukaryotic Ribosomal RNA Genes. *Progress in Nucleic Acid Research and Molecular Biology*, 31:115–160.
- [Mande et al., 2012] Mande, S. S., Mohammed, M. H., and Ghosh, T. S. (2012). Classification of metagenomic sequences: methods and challenges. *Briefings in Bioinformatics*, 13(6):669–681.
- [Mapleson et al., 2017] Mapleson, D., Garcia Accinelli, G., Kettleborough, G., Wright, J., and Clavijo, B. J. (2017). KAT: a K-mer analysis toolkit to quality control NGS datasets and genome assemblies. *Bioinformatics (Oxford, England)*, 33(4):574–576.
- [Marcon et al., 2014] Marcon, E., Scotti, I., Héroult, B., Rossi, V., and Lang, G. (2014). Generalization of the partitioning of shannon diversity. *PloS one*, 9(3):e90289.
- [Martin et al., 2010] Martin, J., Bruno, V. M., Fang, Z., Meng, X., Blow, M., Zhang, T., Sherlock, G., Snyder, M., and Wang, Z. (2010). Rnnotator: an automated de novo transcriptome assembly pipeline from stranded RNA-Seq reads. *BMC Genomics*, 11(1):663.
- [Martin, 2011] Martin, M. (2011). Cutadapt removes adapter sequences from high-throughput sequencing reads. *journal.embnet.org*, 17.

- [Matsen et al., 2012] Matsen, F. A., Hoffman, N. G., Gallagher, A., and Stamatakis, A. (2012). A Format for Phylogenetic Placements. *PLoS ONE*, 7(2):e31009.
- [Matsen et al., 2010] Matsen, F. A., Kodner, R. B., and Armbrust, E. V. (2010). pplacer: linear time maximum-likelihood and Bayesian phylogenetic placement of sequences onto a fixed reference tree. *BMC Bioinformatics*, 11(1):538.
- [Maxam and Gilbert, 1977] Maxam, A. M. and Gilbert, W. (1977). A new method for sequencing DNA. *Proceedings of the National Academy of Sciences of the United States of America*, 74(2):560–4.
- [Medlar et al., 2014] Medlar, A., Aivelo, T., and Löytynoja, A. (2014). S{é}ance: reference-based phylogenetic analysis for 18S rRNA studies. *BMC Evolutionary Biology*, 14(1):235.
- [Mezhoud et al., 2014] Mezhoud, N., Zili, F., Bouzidi, N., Helaoui, F., Ammar, J., and Ouada, H. B. (2014). The effects of temperature and light intensity on growth, reproduction and EPS synthesis of a thermophilic strain related to the genus *Graesiella*. *Bioprocess and Biosystems Engineering*, 37(11):2271–2280.
- [Mitra et al., 2011] Mitra, S., Rupek, P., Richter, D. C., Urich, T., Gilbert, J. A., Meyer, F., Wilke, A., and Huson, D. H. (2011). Functional analysis of metagenomes and metatranscriptomes using SEED and KEGG. *BMC Bioinformatics*, 12 Suppl 1(Suppl 1):S21.
- [Moreau et al., 2012] Moreau, H., Verhelst, B., Couloux, A., Derelle, E., Rombauts, S., Grimsley, N., Van Bel, M., Poulain, J., Katinka, M., Hohmann-Marriott, M. F., Piganeau, G., Rouzé, P., Da Silva, C., Wincker, P., de Peer, Y., and Vandepoele, K. (2012). Gene functionalities and genome structure in *Bathycoccus prasinos* reflect cellular specializations at the base of the green lineage. *Genome Biology*, 13(8):R74.
- [Morris et al., 2014] Morris, E. K., Caruso, T., Buscot, F., Fischer, M., Hancock, C., Maier, T. S., Meiners, T., Müller, C., Obermaier, E., Prati, D., Socher, S. A., Sonnemann, I., Wäschke, N., Wubet, T., Wurst, S., and Rillig, M. C. (2014). Choosing and using diversity indices: insights for ecological applications from the German Biodiversity Exploratories. *Ecology and Evolution*, 4(18):3514–24.

- [Nabhan and Sarkar, 2011] Nabhan, A. R. and Sarkar, I. N. (2011). The impact of taxon sampling on phylogenetic inference: a review of two decades of controversy. *Briefings in Bioinformatics*, 13(1):122–34.
- [Narayanasamy et al., 2016] Narayanasamy, S., Jarosz, Y., Muller, E. E. L., Heintz-Buschart, A., Herold, M., Kaysen, A., Laczny, C. C., Pinel, N., May, P., and Wilmes, P. (2016). IMP: a pipeline for reproducible reference-independent integrated metagenomic and metatranscriptomic analyses. *Genome Biology*, 17(1):260.
- [NCBI Resource Coordinators, 2016] NCBI Resource Coordinators, N. R. (2016). Database resources of the National Center for Biotechnology Information. *Nucleic Acids Research*, 44(D1):D7–19.
- [Nordberg et al., 2014] Nordberg, H., Cantor, M., Dusheyko, S., Hua, S., Poliakov, A., Shabalov, I., Smirnova, T., Grigoriev, I. V., and Dubchak, I. (2014). The genome portal of the Department of Energy Joint Genome Institute: 2014 updates. *Nucleic Acids Research*, 42(D1):D26–D31.
- [Okolodkov and Dodge, 1996] Okolodkov, Y. B. and Dodge, J. D. (1996). Biodiversity and biogeography of planktonic dinoflagellates in the Arctic Ocean. *Journal of Experimental Marine Biology and Ecology*, 202(1):19–27.
- [Oliver et al., 2007] Oliver, M. J., Petrov, D., Ackerly, D., Falkowski, P., and Schofield, O. M. (2007). The mode and tempo of genome size evolution in eukaryotes. *Genome Research*, 17(5):594–601.
- [Paliy and Shankar, 2016] Paliy, O. and Shankar, V. (2016). Application of multivariate statistical techniques in microbial ecology. *Molecular Ecology*, 25(5):1032–57.
- [Papadimitriou, 2014] Papadimitriou, C. (2014). Algorithms, complexity, and the sciences. *Proceedings of the National Academy of Sciences*, 111(45):15881–15887.
- [Pavlopoulos et al., 2010] Pavlopoulos, G. A., Soldatos, T. G., Barbosa-Silva, A., and Schneider, R. (2010). A reference guide for tree analysis and visualization. *BioData Mining*, 3(1):1.

- [Payne et al., 2005] Payne, L. X., Schindler, D. E., Parrish, J. K., and Temple, S. A. (2005). QUANTIFYING SPATIAL PATTERN WITH EVENNESS INDICES. *Ecological Applications*, 15(2):507–520.
- [Pearson, 2013] Pearson, W. R. (2013). An introduction to sequence similarity (“homology”) searching. *Current Protocols in Bioinformatics*, Chapter 3:Unit3.1.
- [Perisin et al., 2016] Perisin, M., Vetter, M., Gilbert, J. A., and Bergelson, J. (2016). 16Stimator: statistical estimation of ribosomal gene copy numbers from draft genome assemblies. *The ISME Journal*, 10(4):1020–4.
- [Pielou, 1966] Pielou, E. (1966). The measurement of diversity in different types of biological collections. *Journal of Theoretical Biology*, 13:131–144.
- [Pillai et al., 2017] Pillai, S., Gopalan, V., and Lam, A. K.-Y. (2017). Review of sequencing platforms and their applications in pheochromocytoma and paragangliomas. *Critical Reviews in Oncology/Hematology*, 116:58–67.
- [Pos et al., 2014] Pos, E., Guevara Andino, J. E., Sabatier, D., Molino, J.-F., Pitman, N., Mogollón, H., Neill, D., Cerón, C., Rivas, G., Di Fiore, A., Thomas, R., Tirado, M., Young, K. R., Wang, O., Sierra, R., García-Villacorta, R., Zagt, R., Palacios, W., Aulestia, M., and Ter Steege, H. (2014). Are all species necessary to reveal ecologically important patterns? *Ecology and Evolution*, 4(24):4626–36.
- [Prokopowich et al., 2003] Prokopowich, C. D., Gregory, T. R., and Crease, T. J. (2003). The correlation between rDNA copy number and genome size in eukaryotes. *Genome*, 46(1):48–50.
- [Quast et al., 2013a] Quast, C., Pruesse, E., Yilmaz, P., Gerken, J., Schweer, T., Yarza, P., Peplies, J., and Glöckner, F. O. (2013a). The SILVA ribosomal RNA gene database project: improved data processing and web-based tools. *Nucleic Acids Research*, 41(Database issue):D590–6.
- [Quast et al., 2013b] Quast, C., Pruesse, E., Yilmaz, P., Gerken, J., Schweer, T., Yarza, P., Peplies, J., and Glöckner, F. O. (2013b). The SILVA ribosomal RNA gene database project: improved data processing and web-based tools. *Nucleic Acids Research*, 41(Database issue):D590—6.

- [Ramette, 2007] Ramette, A. (2007). Multivariate analyses in microbial ecology. *FEMS Microbiology Ecology*, 62(2):142–60.
- [Reller et al., 2007] Reller, L. B., Weinstein, M. P., and Petti, C. A. (2007). Detection and Identification of Microorganisms by Gene Amplification and Sequencing. *Clinical Infectious Diseases*, 44(8):1108–1114.
- [Reuter et al., 2015] Reuter, J. A., Spacek, D. V., and Snyder, M. P. (2015). High-throughput sequencing technologies. *Molecular cell*, 58(4):586–97.
- [Ribeiro et al., 2013] Ribeiro, S., Berge, T., Lundholm, N., Ellegaard, M., and Richardson, B. (2013). Hundred Years of Environmental Change and Phytoplankton Ecophysiological Variability Archived in Coastal Sediments. *PLoS ONE*, 8(4):e61184.
- [Ricotta, 2017] Ricotta, C. (2017). Of beta diversity, variance, evenness, and dissimilarity. *Ecology and Evolution*, 7(13):4835–4843.
- [Ricotta and Podani, 2017] Ricotta, C. and Podani, J. (2017). On some properties of the Bray-Curtis dissimilarity and their ecological meaning. *Ecological Complexity*, 31:201–205.
- [Robinson and Oshlack, 2010] Robinson, M. D. and Oshlack, A. (2010). A scaling normalization method for differential expression analysis of RNA-seq data. *Genome Biology*, 11(3):R25.
- [Rödström, 2017] Rödström, E. M. (2017). *Skeletonema marinoi*. <http://cemeb.science.gu.se/research/target-species-imago+/skeletonema-marinoi>.
- [Rokas, 2011] Rokas, A. (2011). Phylogenetic Analysis of Protein Sequence Data Using the Randomized Accelerated Maximum Likelihood (RAXML) Program. In *Current Protocols in Molecular Biology*, volume Chapter 19, page Unit19.11. John Wiley & Sons, Inc., Hoboken, NJ, USA.
- [Saitou and Nei, 1987] Saitou, N. and Nei, M. (1987). The neighbor-joining method: a new method for reconstructing phylogenetic trees. *Molecular Biology and Evolution*, 4(4):406–25.

- [Sanger et al., 1977] Sanger, F., Nicklen, S., and Coulson, A. R. (1977). DNA sequencing with chain-terminating inhibitors. *Proceedings of the National Academy of Sciences of the United States of America*, 74(12):5463–7.
- [Sarmiento et al., 2010] Sarmiento, H., Montoya, J. M., Vázquez-Domínguez, E., Vaqué, D., and Gasol, J. M. (2010). Warming effects on marine microbial food web processes: how far can we go when it comes to predictions? *Philosophical transactions of the Royal Society of London. Series B, Biological Sciences*, 365(1549):2137–49.
- [Schluter, 1984] Schluter, D. (1984). A Variance Test for Detecting Species Associations, with Some Example Applications. *Ecology*, 65(3):998–1005.
- [Schmidt, 2016] Schmidt, K. (2016). Thermal adaptation of *Thalassiosira pseudonana* using experimental evolution approaches. (June).
- [Schmieder and Edwards, 2011] Schmieder, R. and Edwards, R. (2011). Quality control and preprocessing of metagenomic datasets. *Bioinformatics*, 27(6):863–864.
- [Shannon et al., 2003] Shannon, P., Markiel, A., Ozier, O., Baliga, N. S., Wang, J. T., Ramage, D., Amin, N., Schwikowski, B., and Ideker, T. (2003). Cytoscape: a software environment for integrated models of biomolecular interaction networks. *Genome Research*, 13(11):2498–504.
- [Shemi et al., 2015] Shemi, A., Ben-Dor, S., and Vardi, A. (2015). Elucidating the composition and conservation of the autophagy pathway in photosynthetic eukaryotes. *Autophagy*, 11(4):701–715.
- [Shendure and Ji, 2008] Shendure, J. and Ji, H. (2008). Next-generation DNA sequencing. *Nature Biotechnology*, 26(10):1135–1145.
- [Shull et al., 2001] Shull, V. L., Vogler, A. P., Baker, M. D., Maddison, D. R., and Hammond, P. M. (2001). Sequence alignment of 18S ribosomal RNA and the basal relationships of adephagan beetles: Evidence for monophyly of aquatic families and the placement of trachypachidae. *Systematic Biology*, 50(6):945–969.
- [Simon et al., 2009] Simon, N., Cras, A.-L., Foulon, E., and Lemée, R. (2009). Diversity and evolution of marine phytoplankton. *Comptes Rendus Biologies*, 332(2-3):159–170.

- [Simpson, 2006] Simpson, G. L. (2006). *betadisper function* — *R Documentation*. <https://www.rdocumentation.org/packages/vegan/versions/2.4-2/topics/betadisper>.
- [Singh et al., 2009] Singh, J., Behal, A., Singla, N., Joshi, A., Birbian, N., Singh, S., Bali, V., and Batra, N. (2009). Metagenomics: Concept, methodology, ecological inference and recent advances. *Biotechnology Journal*, 4(4):480–494.
- [Sinha and Lynn, 2014] Sinha, S. and Lynn, A. M. (2014). HMM-ModE: implementation, benchmarking and validation with HMMER3. *BMC Research Notes*, 7(1):483.
- [Stamatakis, 2014] Stamatakis, A. (2014). RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics (Oxford, England)*, 30(9):1312–3.
- [Stamatakis et al., 2005] Stamatakis, A., Ludwig, T., and Meier, H. (2005). RAxML-III: a fast program for maximum likelihood-based inference of large phylogenetic trees. *Bioinformatics*, 21(4):456–463.
- [Stecher et al., 2016] Stecher, A., Neuhaus, S., Lange, B., Frickenhaus, S., Beszteri, B., Kroth, P. G., and Valentin, K. (2016). rRNA and rDNA based assessment of sea ice protist biodiversity from the central Arctic Ocean. *European Journal of Phycology*, 51(1):31–46.
- [Suzuki et al., 2015] Suzuki, S., Kakuta, M., Ishida, T., and Akiyama, Y. (2015). Faster sequence homology searches by clustering subsequences. *Bioinformatics*, 31(8):1183–1190.
- [Taylor et al., 2007] Taylor, F. J. R., Hoppenrath, M., and Saldarriaga, J. F. (2007). Dinoflagellate diversity and distribution. pages 173–184. Springer, Dordrecht.
- [Thomas et al., 2012a] Thomas, M. K., Kremer, C. T., Klausmeier, C. A., and Litchman, E. (2012a). A Global Pattern of Thermal Adaptation in Marine Phytoplankton. *Science*, 338(6110):1085–1088.
- [Thomas et al., 2012b] Thomas, T., Gilbert, J., and Meyer, F. (2012b). Metagenomics - a guide from sampling to data analysis. *Microbial Informatics and Experimentation*, 2(1):3.

- [Thompson et al., 1994] Thompson, J. D., Higgins, D. G., and Gibson, T. J. (1994). CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice. *Nucleic Acids Research*, 22(22):4673–80.
- [Torres-Machorro et al., 2010] Torres-Machorro, A. L., Hernández, R., Cevallos, A. M., and López-Villaseñor, I. (2010). Ribosomal RNA genes in eukaryotic microorganisms: witnesses of phylogeny? *FEMS Microbiology Reviews*, 34(1):59–86.
- [Toseland et al., 2013] Toseland, A., Daines, S. J., Clark, J. R., Kirkham, A., Strauss, J., Uhlig, C., Lenton, T. M., Valentin, K., Pearson, G. A., Moulton, V., and Mock, T. (2013). The impact of temperature on marine phytoplankton resource allocation and metabolism. *Nature Climate Change*, 3(11):979–984.
- [Tremblay et al., 2015] Tremblay, J., Singh, K., Fern, A., Kirton, E. S., He, S., Woyke, T., Lee, J., Chen, F., Dangl, J. L., and Tringe, S. G. (2015). Primer and platform effects on 16S rRNA tag sequencing. *Frontiers in Microbiology*, 6:771.
- [Trimborn et al., 2015] Trimborn, S., Hoppe, C. J., Taylor, B. B., Bracher, A., and Hassler, C. (2015). Physiological characteristics of open ocean and coastal phytoplankton communities of Western Antarctic Peninsula and Drake Passage waters. *Deep Sea Research Part I: Oceanographic Research Papers*, 98:115–124.
- [van Dijk et al., 2014] van Dijk, E. L., Auger, H., Jaszczyszyn, Y., and Thermes, C. (2014). Ten years of next-generation sequencing technology. *Trends in Genetics*, 30(9):418–426.
- [von Mering et al., 2007] von Mering, C., Hugenholtz, P., Raes, J., Tringe, S. G., Doerks, T., Jensen, L. J., Ward, N., and Bork, P. (2007). Quantitative Phylogenetic Assessment of Microbial Communities in Diverse Environments. *Science*, 315(5815):1126–1130.
- [Voss et al., 2013] Voss, M., Bange, H. W., Dippner, J. W., Middelburg, J. J., Montoya, J. P., and Ward, B. (2013). The marine nitrogen cycle: recent discoveries, uncertainties and the potential relevance of climate change. *Philosophical transactions of the Royal Society of London. Series B, Biological Sciences*, 368(1621):20130121.

- [Wain et al., 2002] Wain, H. M., Bruford, E. A., Lovering, R. C., Lush, M. J., Wright, M. W., and Povey, S. (2002). Guidelines for Human Gene Nomenclature. *Genomics*, 79(4):464–470.
- [Wang et al., 2006] Wang, J., Keightley, P. D., and Johnson, T. (2006). MCALIGN2: Faster, accurate global pairwise alignment of non-coding DNA sequences based on explicit models of indel evolution. *BMC Bioinformatics*, 7(1):292.
- [Wang et al., 2014] Wang, Y., Tian, R. M., Gao, Z. M., Bougouffa, S., and Qian, P.-Y. (2014). Optimal eukaryotic 18S and universal 16S/18S ribosomal RNA primers and their application in a study of symbiosis. *PloS one*, 9(3):e90053.
- [Wang et al., 2009] Wang, Z., Gerstein, M., and Snyder, M. (2009). RNA-Seq: a revolutionary tool for transcriptomics. *Nature Reviews. Genetics*, 10(1):57–63.
- [Waterhouse et al., 2009] Waterhouse, A. M., Procter, J. B., Martin, D. M. A., Clamp, M., and Barton, G. J. (2009). Jalview Version 2—a multiple sequence alignment editor and analysis workbench. *Bioinformatics*, 25(9):1189–1191.
- [Wernegreen, 2012] Wernegreen, J. J. (2012). Endosymbiosis. *Current Biology*, 22(14):R555–R561.
- [Wiens and Tiu, 2012] Wiens, J. J. and Tiu, J. (2012). Highly Incomplete Taxa Can Rescue Phylogenetic Analyses from the Negative Impacts of Limited Taxon Sampling. *PLoS ONE*, 7(8):e42925.
- [Wilkins et al., 2013] Wilkins, D., Yau, S., Williams, T. J., Allen, M. A., Brown, M. V., DeMaere, M. Z., Lauro, F. M., and Cavicchioli, R. (2013). Key microbial drivers in Antarctic aquatic environments. *FEMS Microbiology Reviews*, 37(3):303–335.
- [Williams et al., 2010] Williams, K. P., Gillespie, J. J., Sobral, B. W. S., Nordberg, E. K., Snyder, E. E., Shallom, J. M., and Dickerman, A. W. (2010). Phylogeny of gammaproteobacteria. *Journal of Bacteriology*, 192(9):2305–14.
- [Wistrand and Sonnhammer, 2005] Wistrand, M. and Sonnhammer, E. L. L. (2005). Improved profile HMM performance by assessment of critical algorithmic features in SAM and HMMER. *BMC Bioinformatics*, 6(1):99.

- [Wooley et al., 2010] Wooley, J. C., Godzik, A., and Friedberg, I. (2010). A Primer on Metagenomics. *PLoS Computational Biology*, 6(2):e1000667.
- [Wu et al., 2015] Wu, S., Xiong, J., and Yu, Y. (2015). Taxonomic Resolutions Based on 18S rRNA Genes: A Case Study of Subclass Copepoda. *PLOS ONE*, 10(6):e0131498.
- [Ye et al., 2018] Ye, J., Zhang, Y., Cui, H., Liu, J., Wu, Y., Cheng, Y., Xu, H., Huang, X., Li, S., Zhou, A., Zhang, X., Bolund, L., Chen, Q., Wang, J., Yang, H., Fang, L., and Shi, C. (2018). WEGO 2.0: a web tool for analyzing and plotting GO annotations, 2018 update. *Nucleic Acids Research*, 46(W1):W71–W75.
- [Yu and Zhang, 2013] Yu, K. and Zhang, T. (2013). Construction of Customized Sub-Databases from NCBI-nr Database for Rapid Annotation of Huge Metagenomic Datasets Using a Combined BLAST and MEGAN Approach. *PLoS ONE*, 8(4):e59831.
- [Zhang and Horvath, 2005] Zhang, B. and Horvath, S. (2005). A General Framework for Weighted Gene Co-Expression Network Analysis. *Statistical Applications in Genetics and Molecular Biology*, 4(1):Article17.
- [Zhang et al., 2012] Zhang, H., John, R., Peng, Z., Yuan, J., Chu, C., Du, G., and Zhou, S. (2012). The relationship between species richness and evenness in plant communities along a successional gradient: a study from sub-alpine meadows of the Eastern Qinghai-Tibetan Plateau, China. *PloS one*, 7(11):e49024.
- [Zhang et al., 2014] Zhang, J., Kobert, K., Flouri, T., and Stamatakis, A. (2014). PEAR: a fast and accurate Illumina Paired-End reAd mergeR. *Bioinformatics (Oxford, England)*, 30(5):614–20.
- [Zhou et al., 2014] Zhou, Q., Su, X., and Ning, K. (2014). Assessment of quality control approaches for metagenomic data analysis. *Scientific Reports*, 4(1):6957.
- [Zhou et al., 2015] Zhou, Q., Su, X., and Ning, K. (2015). Assessment of quality control approaches for metagenomic data analysis. *Scientific Reports*, 4(1):6957.

Appendices

Appendix A

A.A Metadata

<u>Station</u>	<u>Longitude</u>	<u>Latitude</u>	<u>Depth meter</u>	<u>Temperature degrees celsius</u>	<u>Salinity</u>	<u>Nitrate.Nitrite</u>	<u>Phosphate</u>	<u>Silicate</u>
1	-9.52472	79.0225	17	-1.0337	31.0274	0	0.47	2.48
2	-8.52472	79.07611	26	-1.5122	31.708	0.14	0.95	2.86
3	-7.67278	79.04278	20	-1.4645	31.3282	0	0.62	7.17
4	-4.08556	79.00056	35	-1.7398	33.9123	0.81	0.58	5.49
5	-4.08556	79.00056	10	-1.5083	32.9517	0.81	0.58	5.49
6	-4.78556	78.85611	20	-1.6191	32.2039	0.6	0.74	8.97
7	-3.22861	78.86694	15	-0.8805	32.3454	3.3	0.79	5.47
8	-2.83722	78.89667	110	3.0802	34.9887	10.44	1.05	4.46
9	-2.83722	78.89667	25	-1.3686	33.3389	0.7	0.46	3.03
10	-1.84	78.93667	16	4.3973	35.062	4.77	0.84	3.26
11	-0.55917	78.98167	10	1.3665	33.5306	0.69	0.5	3.64
12	3.73583	79.07278	10	0.0522	33.586	2.8	0.61	2.49
13	3.73583	79.07278	5	-0.3027	33.1465	2.21	0.51	3.12
14	4.1527	78.62153	15	5.87	35.05	NA	NA	NA
15	5.32861	78.83889	15	4.4939	35.0994	3.32	0.73	3.95
16	6.10556	79.08944	7	4.8199	35.1105	3.9	0.76	3.88
17	7.07639	79.06667	18	4.5614	35.0886	9.69	1.07	4.31
18	8.11222	78.86972	10	5.269	35.0693	0.51	0.49	2.75
19	9.23306	78.85139	25	2.3374	34.7579	2.93	0.94	2.16
20	11.30917	76.25389	15	5.5282	35.1514	6.26	0.83	4.08
21	9.85667	73.01889	20	6.0186	35.1528	6.55	0.88	3.9
22	8.86667	71.20083	10	7.1834	35.1344	4.83	0.77	3.62
23	7.73028	69.23028	10	9.0976	34.7321	2.46	0.46	1.52
24	7.73028	69.23028	5	8.7545	34.8661	0.86	0.49	1.68
25	6.53028	67.23028	15	8.9384	35.091	4.16	0.74	2.48
26	6.53028	67.23028	20	8.6781	35.1	2.38	0.66	2.12
27	5.41917	65.24611	20	9.24	34.93	2.85	0.6	2.44
28	5.41917	65.24611	5	9.75	34.94	1.58	0.59	2.35
29	-21.7407	62.8	10	7.69	35.2	13.35	0.862	6.69
30	-20.4903	61.7103	11	8.27	35.18	12.55	0.82	NA
31	-19.3398	60.6801	10	8.75	35.25	11.75	0.755	6.57
32	-18.0699	59.5	10	9.49	35.33	10.61	0.69	NA
33	-16.5201	58.0002	11	9.54	35.37	11.23	0.729	NA
34	-16.509	54.6333	10	10.79	35.45	7.79	0.488	1.59
35	-16.4972	52.6218	26	11.54	35.55	7.95	0.512	NA
36	-16.3567	49.9151	50	12.22	35.64	6.86	0.428	NA
37	-12.1104	47.5701	40	12.07	35.64	5.92	0.412	2.54
38	-12.4301	45.5298	16	13.96	35.81	0.33	0.049	2.18
39	-12.6098	44.2798	20	13.78	35.8	1.47	0.087	NA
40	-12.8799	42.3401	21	14.4	35.9	0.76	0.062	NA
41	-13.1901	40.5296	25	14.99	36.06	0.45	0.047	NA
42	-13.576	38.4205	67	15.9	36.26	0.52	0.045	NA
43	-13.94	36.5297	25	16.67	36.32	0.02	0.019	NA
44	-12.0872	37.833	80	20.63	36.44	0.426330313	0.065720386	0.478187565
45	-13.1352	34.876	80	21.78	36.65	0.269427308	0	0.292683967
46	-14.2597	34.7198	21	17.43	36.46	0.04	0.015	NA
47	-14.2601	34.7197	48	16.93	36.42	1.02	0.031	0.2926
48	-14.5898	32.8199	73	18.06	36.66	0.05	0.022	NA
49	-14.8699	31.2202	75	18.57	36.75	0.02	0.015	NA
50	-14.8704	31.2193	72	18.37	36.71	0.03	0.014	NA
51	-15.0705	30.0186	30	18.55	36.71	0.03	0.021	NA
52	-15	29	81	18.97	36.83	0.02	0.023	0.75105
53	-15.1548	28.937	90	23	36.6	0.543261784	0.076334058	0.751057865
54	-17.4585	26.049	80	24.62	36.9	0.516096977	0.066033049	0.405116155
55	-20.1823	23.69	60	25.25	36.57	1.954317104	0.194632089	0.689491487
56	-20.7015	18.755	45	26.78	36.17	6.127505044	0.331454155	1.187283885
57	-20.515	15.249	55	28.6	35.64	16.42369279	0.682339671	3.367016894
58	-18.6202	8.472	46	29.02	34.88	3.106489549	0.241043741	2.590527277
59	-13.602	2.405	80	27.1	35.61	6.458482733	0.3717252	2.064853304
60	-9.4218	-2.045	63	25.64	35.94	10.02985752	0.51653754	2.617416148
61	-7.0612	-4.668	47	24.71	35.83	11.32366081	0.55562155	3.066000484
62	-4.9043	-7.394	45	23.27	36.23	6.480985475	0.366034527	2.020429743
63	-0.3178	-13.104	75	20.66	36.5	4.395764083	0.236727944	0.826883927
64	2.9768	-17.283	30	18.82	35.97	0	0.19	2.63
65	3.7993	-18.25	43	18.77	35.99	8.005879218	0.444372737	0.85904911
66	5.9978	-20.987	30	18.52	35.62	5.951945303	0.392739675	0.065775733
67	7.1922	-22.644	40	18.48	35.7	0.502052705	0.044543045	0.43401924
68	9.9828	-26.441	30	18.2	35.2	2.283057742	0.173644994	0.688728588

A.B 18S rDNA reference databases taxa IDs

1032745	127567	498767	650286	1072577	670782	858356	1171960	434030
639210	127563	67961	110372	1072585	651459	210735	35113	181586
156996	119481	101922	238096	1072581	110316	2923	408072	197858
157007	342563	221933	66473	75742	348448	425803	5932	223571
296670	402028	48970	632150	203142	650289	425804	37092	197862
44058	1003327	179290	160611	1072566	650283	1003175	104774	274053
1104430	397339	707317	1434936	237895	693284	387435	863766	524897
96791	512344	710654	2957	5875	650331	225107	5927	524900
89044	396045	71746	43686	27996	650333	189332	289478	524895
1117030	49980	138175	217026	5865	651455	326572	289481	693931
135473	219168	332216	8030	476202	113556	510609	210825	197857
413849	6009	44656	7739	300409	240359	1162085	5928	394799
416808	219837	889457	201616	86281	133427	376358	311390	1004046
262226	219841	170400	248472	84966	66800	107758	1081457	240170
31296	653507	1082342	261852	88456	259926	258031	230077	1166918
45884	1071533	56002	31408	84969	2925	39447	122942	455295
1054987	652930	51328	436072	700707	58156	673113	1002335	155143
1054391	223996	56009	143010	482538	73915	271680	1002334	168254
278983	708628	464287	42008	162052	66791	143672	278832	155140
1213618	515478	243130	173452	99158	497816	414396	1002222	168251
178368	497703	1085968	31461	29176	160621	459342	1002226	346245
432302	33640	173495	197877	5811	2916	400756	692909	703568
507873	913975	63592	38533	94643	66465	79894	693871	168244
1247838	303408	3190	309167	689363	325318	410738	880990	168257
985835	97221	926286	522302	62968	326149	71001	35100	446414
614059	265540	204415	536595	333133	685761	858351	211662	1126575
502092	285082	204408	536601	294749	244960	424504	881006	172019
46462	211638	515480	221821	31383	505693	66468	223573	346241
247150	35117	515477	38836	282351	358189	210733	274051	168242
1170558	279580	44430	67781	110459	1118043	39450	435890	168248
693439	216975	643632	36881	110473	158379	160623	669200	151077
696972	425790	652932	94289	541018	5992	1165720	693053	1170556
358025	363325	1050086	552938	46026	696187	327389	35097	37474
601995	142111	197538	578126	244810	651135	339968	1071701	5965
693140	382377	329754	81603	5817	47934	566470	385028	211642
153251	216766	97102	332487	169156	622439	5960	57509	159163
47889	197897	395882	47887	346247	182090	467015	39461	415605
47894	860613	435893	324857	346240	181653	467023	394801	273907
358022	1173301	641229	338344	497731	641228	358021	497729	264369
693421	375585	513281	686996	1173416	513283	197901	990178	1176393
550469	375583	211016	692885	990236	1049793	1170553	375587	1292461
550467	5997	211025	754198	990234	513204	482403	990189	1176361
692521	993385	1169006	268522	614050	1071509	754193	1229648	5974

89957	9606	1134686	142497	257561	35206	332489	400980	132687
584794	60559	101920	173468	164319	81598	332491	400982	864286
122588	7604	182762	123987	257567	479265	332493	356898	592675
160615	82378	1094580	257592	105601	504345	332497	152445	5936
261834	7668	257374	40391	257545	397053	332495	652658	152452
302480	206668	257375	40385	257551	81617	189623	330974	5940
160619	34765	326073	680465	88417	81584	284005	693924	138849
261837	569450	2786	31369	152008	81611	332501	99922	331608
79898	425018	1094582	291393	239144	1074216	332498	382958	59999
388225	97517	35688	367044	257547	1074204	332500	385034	311320
402582	230730	101926	257542	257549	39628	189644	860611	1004043
388228	491138	2771	173450	257554	322167	332483	423615	125641
388232	6669	45157	173447	28020	81592	536594	170499	94673
107036	202087	82540	159503	31421	197880	332484	346229	70075
373098	6661	31354	326253	305493	42024	31367	1174532	1001742
66792	194544	35151	122415	38544	327986	31472	211028	1287487
72554	439682	468936	122413	88422	239160	42480	1292964	5950
261842	178832	82561	35157	88412	931254	257814	1229645	114681
2866	136180	204479	536613	31447	31455	191046	181651	558275
393032	299778	204480	339588	41686	257607	225046	1229734	114678
1132513	560977	753684	339604	196371	164055	31474	1292960	99924
160617	231624	110510	257588	257575	81613	282360	513210	268526
2966	1234259	139984	142502	29232	81600	110475	1176324	278822
763934	104782	362230	31417	29222	81590	332481	1170544	686991
340370	1435206	139907	142499	29241	91060	189642	247154	71585
412152	322853	139980	142505	95361	81594	536593	1071508	40806
418113	136452	282340	35173	29227	81587	536603	526559	1144346
88552	1197702	101924	142512	1269974	81607	536604	717095	163346
95749	314080	37198	128537	257571	38331	536610	1170546	163350
103983	6087	82843	142492	197835	31500	536616	1170548	163357
51511	209422	436068	142489	197840	81615	536591	238901	71594
223366	201645	82839	128534	197837	406714	189638	1297756	5987
350068	201679	70838	217487	38372	81619	536586	692525	993383
1176391	99915	488251	139859	467021	513215	414913	398671	375580
180940	5947	1071522	39464	758568	211022	57513	414910	857030
385029	552936	488247	864284	124797	402891	469766	694303	692523
157072	33674	278986	227086	188973	197904	40330	459245	1170439
101203	87102	366611	552664	188950	527218	502090	55726	1030620
135477	463366	366598	552666	149084	394493	39564	587064	485330
92981	126728	1054399	67809	188943	693161	352448	197894	40328
92972	42467	987159	238782	686740	92969	463362	366596	332299
92976	463364	1054387	238789	45107	687159	985833	54108	88563
92966	4773	366601	238793	563755	391043	350847	74790	1170431
986738	1127205	123356	160259	987156	332305	188946	984054	459518
1054397	278125	983661	37095	1408141	332301	999463	340085	1003332

536599	159727	643655	31496	164601	327057	42696	238774	60001
536597	94299	320783	31481	77548	50045	162320	238761	696133
536609	282348	433419	79259	52965	50044	162317	70182	1005899
536614	282353	1196376	498007	710656	470007	3092	167538	1005901
173553	103713	177373	2801	87090	3097	183309	37360	88570
1261578	110477	109050	700918	647327	52964	1035567	167535	459270
335259	159725	244126	38265	327050	104531	75803	933485	5982
173540	189624	31363	131155	52035	104533	113520	188971	459237
1261580	189634	122417	1070855	269637	104536	113536	5855	459276
1261571	94295	389187	35153	158507	3175	113522	110365	941344
173548	189646	89943	2762	132247	51320	91193	2969	459265
173550	1206573	231748	38271	52685	55999	3088	672928	941345
1261582	158647	871653	38269	142647	33095	55409	340708	459235
1261569	131068	389190	77922	52679	329040	91197	230744	459268
89212	131064	48608	55529	69401	337949	91196	251331	63136
155556	131094	48948	3032	436124	163323	202681	663229	1170433
110469	282363	1034349	478117	361666	993091	795116	340200	485360
121067	189632	464653	464988	579148	204991	183317	340204	909417
228269	35161	1034350	77928	55997	1008953	326141	283649	459257
110466	76903	48942	46947	361668	44654	651586	283647	485357
155558	189640	48951	437768	56006	156110	132186	326279	164621
228262	189636	48615	2898	332218	1148060	875623	1127215	1044902
110471	189648	464689	195067	47790	132188	875621	1127417	523209
1261574	189628	35170	57475	47788	271407	875620	1127213	127229
35163	189652	464649	40526	51718	165818	876697	1127100	188964
945030	268567	48962	193549	56011	55410	145388	1127415	151026
159598	31486	871656	167772	1034561	271398	191687	1127103	361427
282366	67958	231757	551846	361670	271396	307507	1127203	981202
282357	111861	348030	3046	179866	302391	117505	85466	361426
282355	35202	48959	3156	3192	160063	82291	310810	188939
110462	31490	1034343	52028	1158268	132190	889455	135480	87111
860634	209631	348081	71744	34115	299577	31300	114742	33676
99899	209632	48975	107616	3095	170393	356782	100874	12968
98045	877583	279125	1127141	100933	98068	159317	492103	1127216
98064	78393	220110	1127137	687161	64707	377434	35217	1128110
420600	136833	299204	278118	1044904	438413	1127223	39714	178372
52559	280858	534455	64708	42384	536091	1127220	178364	432305
52556	281460	366605	98060	114254	862249	1127211	707168	420607
98055	33653	413940	98066	230409	626141	1127108	707170	3000
67593	127148	413938	2996	42750	1108494	1127133	178375	308878
143451	357350	629711	52241	104657	941460	1069743	39717	44056
126844	272144	629713	590968	159342	210589	44432	27967	88167
707994	167961	1191181	137466	382380	216777	643630	2876	1117027
65357	1486930	366592	216741	443638	178366	515472	29207	210620
1054395	98058	515487	278979	278121	72520	2885	55585	55587

179863	41300	3159	262509	221848	195969	70452	505996	2850
164533	63684	185965	173376	111460	195967	561169	517775	431345
52677	3072	138176	374114	221826	156128	119497	232859	431593
164532	155715	230557	173371	204410	88271	216820	505992	431588
3081	188049	44573	231078	204422	631452	1158023	915350	210449
558848	154508	138169	231080	204419	81844	515468	196682	37319
76111	173497	73033	31301	34119	687948	432268	1115528	1003042
202684	1034627	398706	445995	34121	70448	216617	487668	431324
133488	37433	189347	507620	34140	41875	1003029	1115533	431370
532145	152768	230559	3075	204412	156133	1003039	216745	432551
247495	1293078	230561	38881	162054	156131	33645	265532	1003108
797674	187458	241068	120749	162063	41886	216618	186025	1003036
797671	1034482	63412	163307	35845	539821	426669	941921	1003053
55402	413988	163314	120751	219609	221441	515488	431322	1003046
798524	312843	160070	1065486	219614	35139	426676	216747	515482
3080	306440	348768	202679	3140	127547	426664	216734	1003145
247497	312849	43941	577621	35855	418912	186020	221720	2853
3099	200484	160076	577622	35853	424526	515470	515474	216749
247499	82176	53263	120745	35857	53265	515464	431333	431368
221620	1065496	173490	1212498	35859	38817	515479	216798	197754
75799	312850	173499	1034818	2788	2903	1158036	303458	265561
3074	577485	216028	183676	168183	424539	515471	1003069	197753
797666	63683	153906	51329	3702	127562	1003094	431358	216896
797637	480381	216033	926342	4081	97492	196633	1003020	3003
91190	160068	63410	926344	3659	156173	426662	303409	502569
163309	34148	262505	162066	4558	373042	210441	265537	210618
163317	3093	262511	162062	33129	97495	196638	303462	232512
160069	1232723	439779	162076	4577	1136786	515462	1051708	451786
63675	415184	1081497	926348	4530	284051	210602	718187	29205
634749	34116	1108641	162059	3760	259385	528331	1003112	83371
190057	60131	220632	333302	3871	13221	210616	49249	94617
191674	3171	1081495	162067	284941	37099	455043	449131	423563
249350	284972	220631	185010	1486889	118079	515484	449130	425076
86898	92126	96789	74375	1003075	44451	1003142	216908	441170
142656	92128	37475	2835	59812	658122	1003199	2829	246117
3012	82369	96788	73013	210444	162275	210609	90936	246119
185805	308770	1104315	195974	327391	1278069	374047	375454	243268
2880	240564	98049	73019	134680	475233	233771	310674	1034831
27963	240566	45116	236787	45653	1003208	246123	308883	214821
49246	216801	112064	479719	157126	1003064	1003062	66222	116065
265543	216773	243166	77034	38822	1003144	1003137	62316	157001
1003027	216790	88149	240383	215586	1158022	91992	105404	258581
1003085	265556	64927	82161	35684	265572	3005	105414	196139
1090589	265552	62313	82166	35682	265563	188541	38824	1003033
655750	1003084	3022	420584	172671	172669	124430	157128	

A.C Evenness and occupancy taxonomy names to numbers

Table A.1: Evenness and occupancy taxonomy names to numbers

<u>18S rDNA taxonomy</u>	<u>Number</u>	<u>16S rDNA taxonomy</u>	<u>Number</u>
Stramenopiles	1	Gammaproteobacteria	1
Cryptophyta	2	Alphaproteobacteria	2
Mamiellophyceae	3	Flavobacteriia	3
Rhizaria	4	Marine Group II	4
Maxillopoda	5	unclassified Verrucomicrobia	5
Alveolata	6	Synechococcales	6
Dinophyceae	7	Deltaproteobacteria	7
Spirotrichea	8	Proteobacteria	8
Polycystinea	9	Candidatus Marinimicrobia	9
Coscinodiscophyceae	10	Actinobacteria	10
Amoebozoa	11	Planctomycetia	11
Ascidacea	12	Acidimicrobiia	12
Pelagophyceae	13	Betaproteobacteria	13
Dictyochophyceae	14	Cyanobacteria	14
Bangiophyceae	15	Verrucomicrobiae	15
unclassified Alveolata	16	Saprospira	16
Acantharea	17	Marine Group III	17
Haptophyceae	18	unclassified Thaumarchaeota	18
Cercozoa	19	Nitrosopumilales	19
Gregarinasina	20	Verrucomicrobia	20
Florideophyceae	21	Bacteroidetes	21
Fragilariophyceae	22	Oligosphaeria	22
Chlorophyceae	23	Opitutae	23
Bicosoecida	24	Oligoflexia	24
Viridiplantae	25	Planctomycetes	25

Litostomatea	26	Cytophagia	26
Heterotrichea	27	Sphingobacteriia	27
Labyrinthulomycetes	28	Phycisphaerae	28
Hydrozoa	29	Chloroflexi	29
Pyramimonadales	30	Gemmatimonadetes	30
Bacillariophyceae	31	Oscillatoriophyceidae	31
Phyllopharyngea	32	Patescibacteria group	32
Prasinococcales	33	Chlamydiia	33
Coccidia	34	Bacilli	34
Ellobiopsidae	35	Tenericutes	35
Gymnolaemata	36	Bacteroidetes Order II.Incertae sedis	36
Appendicularia	37	Clostridia	37
Bilateria	38	Lentisphaerae	38
PXclade	39	Chlamydiae	39
Placididea	40	Nitrospina	40
Blastocystis	41	Acidobacteria	41
unclassified Rhizaria	42	Archaea	42
Chrysophyceae	43	Epsilonproteobacteria	43
Pycnococcaceae	44	Tissierellia	44
Oomycetes	45	Halobacteria	45
Mediophyceae	46	Thermoplasmata	46
Perkinsea	47	Bacteroidia	47
Oligohymenophorea	48	Balneolia	48
Karyorelictea	49	Mollicutes	49
Developea	50	Candidatus Saccharibacteria	50
Protostomia	51	Anaerolineae	51
Trebouxiophyceae	52	Candidatus Hydrogenedentes	52
Rhodophyta	53	Gemmatimonadetes <class>	53
Eumetazoa	54	Nitriliruptoria	54
Gromiidae	55	Deinococci	55
unclassified stramenopiles	56	Negativicutes	56

Compsopogonophyceae	57	Rubrobacteria	57
Ulvophyceae	58	Fusobacteriia	58
eudicotyledons	59	Firmicutes	59
Malacostraca	60	Chitinophagia	60
Aconoidasida	61	Erysipelotrichia	61
Actinopteri	62	Chlorobi	62
Bdelloidea	63	Nitrospinae	63
Chlorophyta incertae sedis	64	Spirochaetia	64
Echinoidea	65	Nostocales	65
Mammalia	66	Ardeenticatenia	66
Ophiuroidea	67	Solibacteres	67
Prasinophyceae incertae sedis	68		
Rhodellophyceae	69		
Pedinophyceae	70		
Chlorophyta	71		
Ciliophora	72		
Holothuroidea	73		
Pinguiphyceae	74		
Bolidophyceae	75		
Glaucozystophyceae	76		
Crustacea	77		

Table A.6: North Atlantic 16S rDNA correlation heatmap coefficients values and p-values

Taxonomy	Longitude correlation value	Latitude p-value	Depth correlation value	Temperature p-value	Salinity correlation value	Salinity p-value	Nitrate, Nitrite correlation value	Nitrate, Nitrite p-value	Phosphate correlation value	Phosphate p-value	Silicate correlation value	Silicate p-value
Actinobacteria	-0.586598701	0.432743271	0.624426837	0.152586988	0.212460084	0.014818862	0.131171153	0.648508873	-0.117187115	0.6258783039	0.6112180537	
Actinobacteria	-0.777589116	0.057153055	0.012501695	0.388317775	0.194096885	0.039657187	0.131171153	0.648508873	-0.117187115	0.6258783039	0.6112180537	
Actinobacteria	-0.134118056	0.731049879	0.023535519	0.623351972	0.251869248	0.078927976	0.131171153	0.648508873	-0.117187115	0.6258783039	0.6112180537	
Actinobacteria	0.692922521	-0.095678865	-0.5791222855	0.259292009	-0.210362327	0.0175251791	0.131171153	0.648508873	-0.117187115	0.6258783039	0.6112180537	
Bacilli	0.695711884	-0.095678865	-0.5791222855	0.259292009	-0.210362327	0.0175251791	0.131171153	0.648508873	-0.117187115	0.6258783039	0.6112180537	
Bacteroidia	0.111839145	0.624292169	-0.321410132	0.291192641	0.19320315	0.084836514	0.342901125	0.21032268	0.391349688	0.145590051	0.554692984	
Bacteroidia	0.587823281	-0.167099984	-0.321410132	0.291192641	0.19320315	0.084836514	0.342901125	0.21032268	0.391349688	0.145590051	0.554692984	
Bacteroidia	-0.118332691	0.514101356	0.624292169	-0.321410132	0.291192641	0.19320315	0.084836514	0.342901125	0.21032268	0.391349688	0.145590051	
Chlamydia	0.874462396	-0.028852306	0.918776674	0.217840036	0.076782913	0.021570335	0.105846831	0.013997361	0.707313371	0.885353282	0.026011789	
Cytophaga	0.40064099	-0.878989084	0.153E+005	0.321410132	0.291192641	0.19320315	0.084836514	0.342901125	0.21032268	0.391349688	0.145590051	
Cytophaga	0.324809341	-0.397255506	0.037541091	0.321410132	0.291192641	0.19320315	0.084836514	0.342901125	0.21032268	0.391349688	0.145590051	
Deltaproteobacteria	0.320729313	0.287105789	0.28936675	0.02839155	0.170348015	0.117308374	-0.007014436	0.077675078	0.250240928	0.285626928	0.290970713	
Epsilonproteobacteria	-0.273487987	0.112486721	-0.281736895	0.300611615	-0.108832416	0.141797669	0.437297769	0.218932292	0.350148292	0.206711541	0.42E+005	
Epsilonproteobacteria	0.00885271	NA	NA	NA	NA	NA	NA	0.437297769	0.218932292	0.350148292	0.206711541	
Epsilonproteobacteria	0.408652551	-1.396275148	0.00158852	0.40E+001	0.681206389	-0.677271189	0.035454285	0.394882929	0.144258665	0.014258665	0.014258665	
Epsilonproteobacteria	0.207193635	-0.731071863	0.019376392	0.019376392	0.019376392	-0.731071863	0.035454285	0.394882929	0.144258665	0.014258665	0.014258665	
Gammaimonobacteria	0.002181036	-0.397249831	-0.449616126	0.02969284	-0.320681746	0.02969284	0.02969284	0.02969284	0.02969284	0.02969284	0.02969284	
Gammaimonobacteria	0.128639314	-0.784037884	-0.242610123	0.884212905	-0.542823419	0.030316913	0.030316913	0.030316913	0.030316913	0.030316913	0.030316913	
Nc. Bacteroides	0.007889992	-0.769407094	-0.658240903	0.007462381	-0.481770188	0.017519445	0.816870821	0.000201386	0.709801817	0.000765448	0.862949574	
Nc. Bacteroides	-0.204575516	0.951729333	0.762828706	0.021703412	0.888849346	0.396979788	0.764877424	0.000903192	-0.918287532	1.37E+006	0.004132321	
Nc. Bacteroides	0.001629231	0.001629231	0.001629231	0.001629231	0.001629231	0.001629231	0.001629231	0.001629231	0.001629231	0.001629231	0.001629231	
Nc. Bacteroides	-0.501829576	0.666558249	0.666558249	0.666558249	0.666558249	0.666558249	0.666558249	0.666558249	0.666558249	0.666558249	0.666558249	
Nc. unclassified Euryarchaeota	0.65731576	0.084183287	0.463127785	0.095349031	0.141703982	0.614103125	0.149683734	0.6348309594	0.072610157	0.778872832	0.328791226	
Nc. Verduccium	0.098911803	-0.166947024	0.466260884	0.302466207	-0.306590056	0.263394694	0.176279884	0.852439254	0.2998879243	0.278831555	0.457691154	
Negativetes	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	
Negativetes	0.6384088	0.0785397	0.0785397	0.0785397	0.0785397	0.0785397	0.0785397	0.0785397	0.0785397	0.0785397	0.0785397	
Oligosphaeria	1.88E+008	0.971534203	0.141938032	0.22778989	0.090673183	0.179227062	-0.09023798	0.740100928	0.753141152	-0.158862539	0.047835855	
Oligosphaeria	0.519166877	0.039847709	0.264381041	0.000290813	0.564904876	0.564904876	0.564904876	0.000614071	-0.061934832	0.018087593	1.66E+006	
Oligosphaeria	0.021880919	0.220790689	0.037014189	0.683621532	-0.326781451	0.235112935	0.190126119	0.855656986	0.158292025	0.514287493	0.860103806	
Oligosphaeria	0.14809548	0.76E+005	0.14809548	0.14809548	0.14809548	0.14809548	0.14809548	0.14809548	0.14809548	0.14809548	0.14809548	
Rhodobacteria	0.761611764	-0.268212985	-0.137711394	0.343778486	0.004425294	-0.286304878	0.3008940217	NA	NA	0.962491972	0.174096884	
Solibacteria	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	
Solibacteria	0.93389912	-0.708450222	-0.464711408	0.029790071	-0.10390916	0.002957324	0.038984689	0.892791213	0.585209553	0.021001766	0.004788275	
Thermomicrobium	-0.909779329	3.81E+008	-0.953639147	0.00797005	-0.116595158	2.32E+006	0.942928848	0.947643198	0.947643198	8.16E+008	0.955666589	
Thermomicrobium	0.704294574	0.619088858	0.9171289	0.000319779	0.679107045	0.000448172	-0.618451196	0.01477223	-0.615690413	0.01569613	1.92E+006	
Thermomicrobium	-0.550585793	0.038018104	0.038018104	0.038018104	0.038018104	0.038018104	0.038018104	0.038018104	0.038018104	0.038018104	0.038018104	
Thermomicrobium	0.105614279	0.001710019	0.001710019	0.001710019	0.001710019	0.001710019	0.001710019	0.001710019	0.001710019	0.001710019	0.001710019	
Thermomicrobium	0.354433898	0.011700110	0.277400935	0.629874847	0.298874847	0.770852684	0.0801841584	0.737290982	-0.0801841584	0.477384549	0.477384549	
Thermomicrobium	-0.40680225	0.1066544813	0.621273984	0.238529005	0.298828558	0.341957245	0.222230276	0.4290130387	0.1994124925	0.1994124925	0.3595303108	
Thermomicrobium	0.249665509	0.110516943	0.065716106	0.358893799	0.299520847	0.1535848827	-0.129021715	0.6467195488	0.052326179	0.931260751	0.638193088	
Thermomicrobium	-0.197596919	-0.724099618	-0.197596919	-0.197596919	-0.197596919	-0.197596919	-0.197596919	-0.197596919	-0.197596919	-0.197596919	-0.197596919	
Thermomicrobium	0.830413935	-0.795753271	0.830413935	0.830413935	0.830413935	0.830413935	0.830413935	0.830413935	0.830413935	0.830413935	0.830413935	
Thermomicrobium	0.252062958	-0.1630179118	0.363707657	0.690157574	-0.492173145	0.555801426	0.350980929	-0.074508979	0.96487648	0.06487648	0.23003281	
Thermomicrobium	0.002770036	0.178839759	0.002770036	0.002770036	0.002770036	0.002770036	0.002770036	0.002770036	0.002770036	0.002770036	0.002770036	
Thermomicrobium	0.466793416	0.07731946	0.07731946	0.07731946	0.07731946	0.07731946	0.07731946	0.07731946	0.07731946	0.07731946	0.07731946	
Thermomicrobium	0.68552423	0.660794694	0.68552423	0.68552423	0.68552423	0.68552423	0.68552423	0.68552423	0.68552423	0.68552423	0.68552423	
Thermomicrobium	0.302793812	0.480789694	0.302793812	0.302793812	0.302793812	0.302793812	0.302793812	0.302793812	0.302793812	0.302793812	0.302793812	

Table A.7: South Atlantic 16S rDNA correlation heatmap coefficients values and p-values

Taxonomy	Longitude correlation value	Latitude p-value	Longitude p-value	Latitude correlation value	Depth correlation value	Depth p-value	Temperature correlation value	Temperature p-value	Salinity correlation value	Salinity p-value	Nitrate:Nitrite correlation value	Nitrate:Nitrite p-value	Phosphate correlation value	Phosphate p-value	Silicate correlation value	Silicate p-value
Actinobacteria	-0.877649733	0.000271001	-0.744061829	0.001465794	0.398173045	0.141826179	0.827254736	0.003127792	-0.0503116494	0.173275782	0.832655969	0.027335590	0.344003013	0.203553509	0.330318126	0.229158975
Alphaproteobacteria	0.286481674	0.294162298	0.289828246	0.129623592	0.129623592	0.144227218	-0.252524913	0.358711332	0.759910963	0.1066911684	0.705606415	0.565052844	0.184591416	0.336184565	0.150581733	0.581088275
Anaerolineae	0.001871186	-0.950118553	-0.235693725	0.165717903	0.3570709884	0.397119222	0.165717903	0.086985131	0.193195197	-0.280725442	0.119936607	0.311999607	-0.3139149521	0.2244831224	-0.178842327	0.323944819
Bacteroidia	0.189445854	0.000118818	0.189445854	0.000118818	0.189445854	0.189445854	0.189445854	0.189445854	0.189445854	0.189445854	0.189445854	0.189445854	0.189445854	0.189445854	0.189445854	0.189445854
Betaproteobacteria	0.372119345	0.000118818	0.372119345	0.000118818	0.372119345	0.115851264	-0.02512089	0.790412489	0.037230462	-0.26716321	0.485781822	0.435822065	-0.2623272375	0.345582065	-0.028822208	0.018811694
Bacteroidia	-0.115632999	0.037882984	-0.115632999	0.037882984	0.115632999	0.628644487	0.345186795	0.245186795	0.079417961	-0.301274054	0.461803198	-0.3032127643	0.2719431855	0.461803198	0.3032127643	0.2719431855
Bacteroidia	0.122840871	0.338884279	0.122840871	0.338884279	0.122840871	0.172829714	-0.023721724	0.933109324	0.6666451991	-0.0812302491	0.273254651	0.2939146547	-0.2926297019	0.2939146547	-0.338817966	0.039224143
Beijerinkeellum	0.054949594	-0.3045059178	0.054949594	-0.3045059178	0.054949594	0.054949594	0.000310622	0.000310622	0.050202908	-0.3603254473	0.1868202014	-0.2869274591	0.3901383358	0.050202908	-0.672828797	0.009818155
Bifidobacteriales	0.077817446	0.0054918038	0.077817446	0.0054918038	0.077817446	0.077817446	0.000310622	0.000310622	0.050202908	-0.3603254473	0.1868202014	-0.2869274591	0.3901383358	0.050202908	-0.672828797	0.009818155
Chloroflexi	0.057814313	0.271164555	0.057814313	0.271164555	0.057814313	0.782091191	-0.612387246	0.093857648	0.551011537	-0.531199048	0.042353937	0.065106239	0.065106239	0.065106239	0.532929709	0.049552596
Cyanobacteria	0.822811442	0.000167744	0.822811442	0.000167744	0.822811442	0.434475624	-0.815867307	0.000283474	0.846629245	-0.344378048	0.209621745	-0.2666294182	0.336184565	0.209621745	-0.344378048	0.207148718
Cytophaga	0.181904155	0.358902278	0.181904155	0.358902278	0.181904155	0.72748721	-0.160925754	0.566110642	0.296422585	0.131372289	0.1919198574	0.0490909629	0.336184565	0.1919198574	0.336184565	0.229643496
Deinococcota	0.002851335	-0.623416988	0.002851335	-0.623416988	0.002851335	0.846629245	0.000118818	0.000118818	0.846629245	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818
Deinococcales	0.002851335	-0.623416988	0.002851335	-0.623416988	0.002851335	0.846629245	0.000118818	0.000118818	0.846629245	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818
Deinococcales	0.002851335	-0.623416988	0.002851335	-0.623416988	0.002851335	0.846629245	0.000118818	0.000118818	0.846629245	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818
Deinococcales	0.002851335	-0.623416988	0.002851335	-0.623416988	0.002851335	0.846629245	0.000118818	0.000118818	0.846629245	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818
Deinococcales	0.002851335	-0.623416988	0.002851335	-0.623416988	0.002851335	0.846629245	0.000118818	0.000118818	0.846629245	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818
Deinococcales	0.002851335	-0.623416988	0.002851335	-0.623416988	0.002851335	0.846629245	0.000118818	0.000118818	0.846629245	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818
Deinococcales	0.002851335	-0.623416988	0.002851335	-0.623416988	0.002851335	0.846629245	0.000118818	0.000118818	0.846629245	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818
Deinococcales	0.002851335	-0.623416988	0.002851335	-0.623416988	0.002851335	0.846629245	0.000118818	0.000118818	0.846629245	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818
Deinococcales	0.002851335	-0.623416988	0.002851335	-0.623416988	0.002851335	0.846629245	0.000118818	0.000118818	0.846629245	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818
Deinococcales	0.002851335	-0.623416988	0.002851335	-0.623416988	0.002851335	0.846629245	0.000118818	0.000118818	0.846629245	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818
Deinococcales	0.002851335	-0.623416988	0.002851335	-0.623416988	0.002851335	0.846629245	0.000118818	0.000118818	0.846629245	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818
Deinococcales	0.002851335	-0.623416988	0.002851335	-0.623416988	0.002851335	0.846629245	0.000118818	0.000118818	0.846629245	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818
Deinococcales	0.002851335	-0.623416988	0.002851335	-0.623416988	0.002851335	0.846629245	0.000118818	0.000118818	0.846629245	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818
Deinococcales	0.002851335	-0.623416988	0.002851335	-0.623416988	0.002851335	0.846629245	0.000118818	0.000118818	0.846629245	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818
Deinococcales	0.002851335	-0.623416988	0.002851335	-0.623416988	0.002851335	0.846629245	0.000118818	0.000118818	0.846629245	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818
Deinococcales	0.002851335	-0.623416988	0.002851335	-0.623416988	0.002851335	0.846629245	0.000118818	0.000118818	0.846629245	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818
Deinococcales	0.002851335	-0.623416988	0.002851335	-0.623416988	0.002851335	0.846629245	0.000118818	0.000118818	0.846629245	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818
Deinococcales	0.002851335	-0.623416988	0.002851335	-0.623416988	0.002851335	0.846629245	0.000118818	0.000118818	0.846629245	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818
Deinococcales	0.002851335	-0.623416988	0.002851335	-0.623416988	0.002851335	0.846629245	0.000118818	0.000118818	0.846629245	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818
Deinococcales	0.002851335	-0.623416988	0.002851335	-0.623416988	0.002851335	0.846629245	0.000118818	0.000118818	0.846629245	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818
Deinococcales	0.002851335	-0.623416988	0.002851335	-0.623416988	0.002851335	0.846629245	0.000118818	0.000118818	0.846629245	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818
Deinococcales	0.002851335	-0.623416988	0.002851335	-0.623416988	0.002851335	0.846629245	0.000118818	0.000118818	0.846629245	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818
Deinococcales	0.002851335	-0.623416988	0.002851335	-0.623416988	0.002851335	0.846629245	0.000118818	0.000118818	0.846629245	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818
Deinococcales	0.002851335	-0.623416988	0.002851335	-0.623416988	0.002851335	0.846629245	0.000118818	0.000118818	0.846629245	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818
Deinococcales	0.002851335	-0.623416988	0.002851335	-0.623416988	0.002851335	0.846629245	0.000118818	0.000118818	0.846629245	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818
Deinococcales	0.002851335	-0.623416988	0.002851335	-0.623416988	0.002851335	0.846629245	0.000118818	0.000118818	0.846629245	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818
Deinococcales	0.002851335	-0.623416988	0.002851335	-0.623416988	0.002851335	0.846629245	0.000118818	0.000118818	0.846629245	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818
Deinococcales	0.002851335	-0.623416988	0.002851335	-0.623416988	0.002851335	0.846629245	0.000118818	0.000118818	0.846629245	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818
Deinococcales	0.002851335	-0.623416988	0.002851335	-0.623416988	0.002851335	0.846629245	0.000118818	0.000118818	0.846629245	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818	0.000118818
Deinococcales	0.002851335	-0.623416988	0.002851335	-0.623416988	0.002851335	0.846629245	0.000118818	0.000118818	0.846629245	0.000118818	0.000118818					

A.E Co-occurrence node numbers to taxon and class membership

Table A.8: Co-occurrence module ($n=70$) node numbers to taxon and class membership

<u>Node number</u>	<u>Species/Genus</u>	<u>Class</u>
1	Erythrobacter	Alphaproteobacteria
2	Alteromonas	Gammaproteobacteria
3	Roseovarius	Alphaproteobacteria
4	Marinobacter	Gammaproteobacteria
5	Pelagomonas calceolata	Pelagophyceae
6	Prochlorococcus	Cyanobacteria
7	Pseudomonas	Gammaproteobacteria
8	Synechococcus	Cyanobacteria
9	Pelagibaca	Alphaproteobacteria
10	Staphylococcus	Bacilli
11	Candidatus actinomarina	Actinobacteria
12	Alcanivorax	Gammaproteobacteria
13	Croceibacter	Flavobacteriia
14	Halomonas	Gammaproteobacteria
15	Loktanella	Alphaproteobacteria
16	Rhodococcus	Actinobacteria
17	Sphingorhabdus	Alphaproteobacteria
18	Streptococcus	Bacilli
19	Coxiella	Gammaproteobacteria
20	Hyphomonas	Alphaproteobacteria
21	Lawsonella	Actinobacteria
22	Roseibacillus	Verrucomicrobiae
23	Candidatus fritschea	Chlamydiia
24	Psychrobacter	Gammaproteobacteria
25	Rhodopirellula	Planctomycetia

26	Tenacibaculum	Flavobacteriia
27	Ichthyodinium chabelardi	Alveolata
28	Acinetobacter	Gammaproteobacteria
29	Hoeflea	Alphaproteobacteria
30	Candidatus nitrosopelagicus	Thaumarchaeota
31	Codonellopsis americana	Spirotrichea
32	Euduboscquella crenulata	Dinophyceae
33	Halobacteriovorax	Oligoflexia
34	Pseudoalteromonas	Gammaproteobacteria
35	Idiomarina	Gammaproteobacteria
36	Jejudonia	Flavobacteriia
37	Blastopirellula	Planctomycetia
38	Maricaulis	Alphaproteobacteria
39	Pelagococcus subviridis	Pelagophyceae
40	Maribacter	Flavobacteriia
41	Caecitellus parvulus	Stramenopiles
42	Litoricola	Gammaproteobacteria
43	Aureispira	Saprospira
44	Filamoeba nolandi	Amoebozoa
45	Marinomonas	Gammaproteobacteria
46	Thiothrix	Gammaproteobacteria
47	Sphingomonas	Alphaproteobacteria
48	Oleiphilus	Gammaproteobacteria
49	Pseudophaeobacter	Alphaproteobacteria
50	Marinoscillum	Cytophagia
51	Neptunomonas	Gammaproteobacteria
52	Fluviicola	Flavobacteriia
53	Hydra vulgaris	Hydrozoa
54	Oleibacter	Gammaproteobacteria
55	Parabirojimia similis	Spirotrichea
56	Olleya	Flavobacteriia

57	Bathycoccus prasinus	Mamiellophyceae
58	Candidatus pelagibacter	Alphaproteobacteria
59	Kordia	Flavobacteriia
60	Nitrosopumilus	Thaumarchaeota
61	Nonlabens	Flavobacteriia
62	Psychroflexus	Flavobacteriia
63	Arcobacter	Epsilonproteobacteria
64	Delftia	Betaproteobacteria
65	Marinicella	Gammaproteobacteria
66	Bradyrhizobium	Alphaproteobacteria
67	Rubritalea	Verrucomicrobiae
68	Bacillus	Bacilli
69	Salinirepens	Flavobacteriia
70	Crocinitomix	Flavobacteriia

Table A.9: Co-occurrence module ($n=51$) node numbers to taxon and class membership

<u>Node number</u>	<u>Species/Genus</u>	<u>Class</u>
1	Colwellia	Gammaproteobacteria
2	Polaribacter	Flavobacteriia
3	Balneatrix	Gammaproteobacteria
4	Ulvibacter	Flavobacteriia
5	Amylibacter	Alphaproteobacteria
6	Lentibacter	Alphaproteobacteria
7	Favella azorica	Spirotrichea
8	Phaeocystis cordata	Haptophyceae
9	Lentimonas	Verrucomicrobia
10	Sulfitobacter	Alphaproteobacteria
11	Oceanicoccus	Gammaproteobacteria
12	Formosa	Flavobacteriia
13	Haliea	Gammaproteobacteria

14	<i>Actinocyclus actinochilus</i>	Coscinodiscophyceae
15	<i>Chroomonas cf. mesostigmatica</i>	Cryptophyta
16	<i>Arenicella</i>	Gammaproteobacteria
17	<i>Emcibacter</i>	Alphaproteobacteria
18	<i>Illumatobacter</i>	Acidimicrobiia
19	<i>Pyramimonas disomata</i>	Chlorophyta
20	<i>Aquibacter</i>	Flavobacteriia
21	<i>Glaciecocola</i>	Gammaproteobacteria
22	<i>Pseudofulvibacter</i>	Flavobacteriia
23	<i>Pithites vorax</i>	Phyllopharyngea
24	<i>Dokdonia</i>	Flavobacteriia
25	<i>Mantoniella antarctica</i>	Chlorophyta
26	<i>Mesoflavibacter</i>	Flavobacteriia
27	<i>Porticoccus</i>	Gammaproteobacteria
28	<i>Sinobacterium</i>	Gammaproteobacteria
29	<i>Schizochytrium aggregatum</i>	Labyrinthulomycetes
30	<i>Dictyocha fibula</i>	Dictyochophyceae
31	<i>Flavicella</i>	Flavobacteriia
32	<i>Oleispira</i>	Gammaproteobacteria
33	<i>Pseudochattonella verruculosa</i>	Dictyochophyceae
34	<i>Pterosperma cristatum</i>	Chlorophyta
35	<i>Prasinoderma coloniale</i>	Chlorophyta
36	<i>Pseudoscourfieldia marina</i>	Chlorophyta
37	<i>Developayella elegans</i>	Stramenopiles
38	<i>Ellobiopsis chattonii</i>	Alveolata
39	<i>Pseudohongiella</i>	Gammaproteobacteria
40	<i>Francisella</i>	Gammaproteobacteria
41	<i>Neptuniibacter</i>	Gammaproteobacteria
42	<i>Marivita</i>	Alphaproteobacteria
43	<i>Paraglaciecola</i>	Gammaproteobacteria
44	<i>Aureococcus anophagefferens</i>	Pelagophyceae

45	<i>Pelagostrobilidium neptuni</i>	Spirotrichea
46	<i>Magnetospira</i>	Alphaproteobacteria
47	<i>Peritromus kahli</i>	Heterotrichea
48	<i>Ceratium tenue</i>	Dinophyceae
49	<i>Chrysochromulina parva</i>	Haptophyceae
50	<i>Varistrombidium kielum</i>	Spirotrichea
51	<i>Florenciella parvula</i>	Dictyochophyceae

A.F Co-occurrence module correlation heatmaps

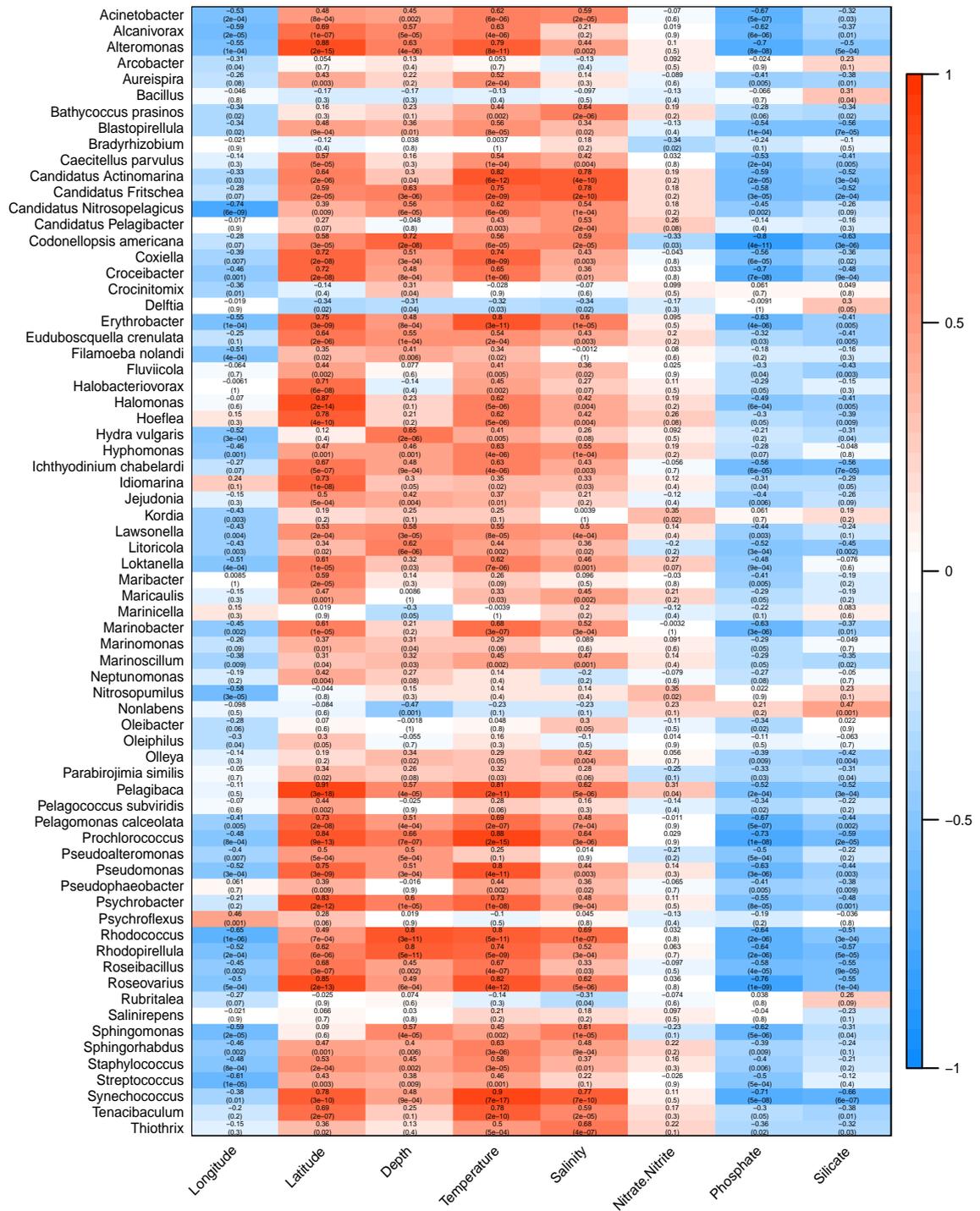


Figure A.1: In the co-occurrence analysis with WGCNA on the log₁₀-scaled abundances of 18S rDNA species level and 16S rDNA genus level, two modules were found. Depicted is the correlation heatmap for the turquoise module in figure 3.17 a ($n=70$). Along the left hand side is the species/genus name and environmental variables are displayed at the bottom. The colours correspond to the correlation values, red is positively correlated and blue is negatively correlated. The Pearson correlation coefficient values and p-values in brackets are displayed in each square

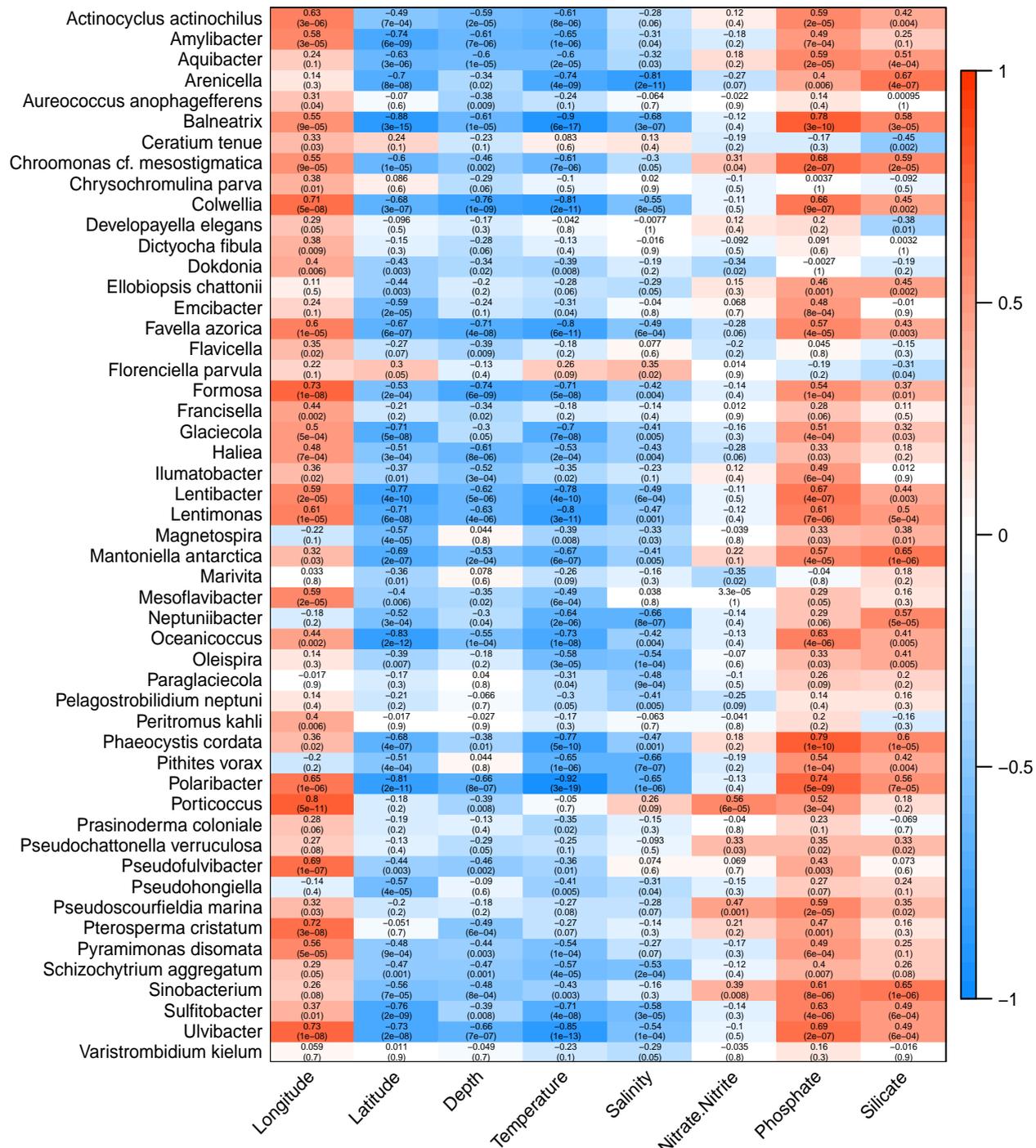


Figure A.2: In the co-occurrence analysis with WGCNA on the log10-scaled abundances of 18S rDNA species level and 16S rDNA genus level, two modules were found. Depicted is the correlation heatmap for the blue module in figure 3.17 b ($n=51$). Along the left hand side is the species/genus name and environmental variables are displayed at the bottom. The colours correspond to the correlation values, red is positively correlated and blue is negatively correlated. The Pearson correlation coefficient values and p-values in brackets are displayed in each square

Appendix B

B.A PhymmBL four genomes locations online

Cyanidioschyzon merolae from Cyanidioschyzon merolae Genome Project <http://merolae.biol.s.u-tokyo.ac.jp/download>;

Danio rerio from UCSC <http://genome.ucsc.edu/cgi-bin/hgGateway?db=danRer5>;

Homo sapiens from Genome Reference Consortium <http://hgdownload.cse.ucsc.edu/goldenPath/hg19/chromosomes/>

Strongylocentrotus purpuratus from Sea Urchin Genome Project <http://www.hgsc.bcm.tmc.edu/projectspecies-o-Strongylocentrotus>.