2   Supplementary File 1: Supplementary Text

# 3   Supplementary Materials and Methods

4   ***Sample Collection***

5          Samples were collected in association with Ottesen et al. 2013 off the coast of

6   California from September 16-19, 2010 along the warm side of an upwelling-driven front

7   (1). Briefly, the Environmental Sample Processor (ESP; 2) was suspended at 23 m

8   depth below a semi-Lagrangian surface float, collecting 1L of seawater every ~4 h for

9   61 h (~2.6 days). Retained particulates were size fractionated onto 5 $\mu$m and 0.22 $\mu$m

10  Durapore 25 mm filters (Millipore, Billerica, MA, USA), preserved immediately *in situ* via

11  a 2 min incubation in RNALater (Ambion) and stored at -80°C within 36 hours of ESP

12  recovery. Nutrients, chlorophyll, and other oceanographic metadata was obtained via

13  shipboard CTD/niskin rosette casts (Supplementary Dataset 12) as previously

14  described (3). Drift speed was determined via a surface-float mounted GPS sensor.

15  Water speed relative to the drifter was measured using a surface-float mounted ADCP

16  (Supplementary Dataset 12A) in order to detect deviations from truly Lagrangian

17  sampling (e.g. wind forcing).

18

19  ***Metatranscriptome library preparation and sequencing***

20         Small size class metatranscriptomes were sequenced by Ottesen et al. using a

21  GS Titanium system (Roche) according to their previously published methods (1).

22  Large size class cDNA was prepared as in Ottesen et al. from ribosomal RNA-depleted

23  total RNA (1). 1 ul of cDNA per sample was used to prepare metatranscriptome libraries

24  with the Truseq RNA Sample Prep kit v2 (Illumina™) according to manufacturer's

25  instructions starting from the end repair step. Libraries were paired-end sequenced on

26  the Illumina HiSeq 2000 platform to obtain 2x100bp reads.

27  ***Amplicon library preparation and sequencing***

28      Large size class 16S and 18S ribosomal RNA were sequenced using 454 GS

29  FLX Titanium pyrosequencing. Nearly universal bacterial primers 341F (5'-

30  CCTACGGGNGGCWGCAG-3') (4) and 926R (5'-CCGTCAATTCMTTTRAGT-3')(5)

31  were used to target the v3v5 region of 16S and primers and TAReuk454FWD1 (5'-

32  CCAGCASCYGCGGTAATTCC-3') and TAReukREV3 (5'-ACTTTCGTTCTTGATYRA-3')

33  (6), were used to target the v4 region of 18S, each amplifying an approximately 500 bp

34  region of cDNA. FLX Titanium adapters (A adapter sequence: 5' 127

35  CCATCTCATCCCTGCGTGTCTCCGACTCAG 3'; B adapter sequence: 5' 128

36  CCTATCCCCTGTGTGCCTTGGCAGTCTCAG 3') and 10bp multiplex identifier (MID)

37  barcodes were used for multiplexed 454 sequencing.

38      cDNA was prepared from 50 ng per sample of total RNA using the Life

39  Technologies SuperScript III First Strand Synthesis system with random hexamer

40  primers. cDNA concentration ranged from 312 – 18,440 pg/microliter. 1 μl of cDNA was

41  used as a template amplified using Life Technologies AccuPrime PCR system kit, in a

42  reaction containing 1X AccuPrime Buffer II, .75 units of AccuPrime Taq High Fidelity,

43  and a final primer concentration of 200 nM, alongside a no template negative control for

44  cDNA synthesis.  Amplifications were performed using a Life Technologies ProFlex

45  PCR system, with an initial denaturation at 95°C for 2 minutes, 30 cycles of 95°C for 20

46 seconds, 56°C for 30 seconds, 72°C for 5 minutes. PCR products (2 µl of each sample

47 and 5 µl of negative control) were run on a 1% agarose gel at 105 V for 35 minutes,

48 then cleaned with Ampure XP beads (Beckman Coulter, Brea CA), and resuspended in

49 25 µL of Qiagen elution buffer. 2.5 µL was used for visualization on an agarose gel, 1

50 µL was used in a LifeTechnologies' PicoGreen Quant-IT assay to quantify the final

51 product, and 45 ng of both 16S and 18S amplicons were pooled separately for 454

52 pyrosequencing.

53 　　　The vendor's standard protocols (Roche Diagnostics) were used for library QC,

54 emPCR, enrichment and 454 sequencing with the following modifications: KAPA

55 Biosystems Library Quantification Kit for qPCR was used to accurately estimate the

56 number of molecules needed for emPCR, automation (BioMek FX) was used to "break"

57 the emulsions after emPCR, and butanol was used to for ease of handling during the

58 breaking process. The bead enrichment process was automated by using Roche's REM

59 e (Robotic Enrichment Module).

60

61 **_Bioinformatic analysis of metatranscriptomes_**

62 　　　See Figure S2 for illustration of metatranscriptomic analysis pipeline.

63 _Open reading frame (ORF) calling and annotation_

64 　　　Large fraction Illumina reads and small fraction 454 reads were processed via

65 the RNAseq Annotation Pipeline (rap) v0.4 (7).  Small fraction reads were obtained via

66 DBCLS SRA (http://sra.dbcls.jp/) using accession number SRA062433.  Reads from

67 both fractions were trimmed to remove primers and areas of low sequence quality

68 (reads must be at least 30 base pairs (bp) long and have a quality score of at least 33 to

69    be retained).  Illumina reads were paired. Ribosomal RNA (rRNA) reads were removed

70    using Ribopicker v.0.4.3 (8). Large fraction reads were assembled using CLC Genomics

71    Workbench 9.5.3 (https://www.qiagenbioinformatics.com/) first by library, then overall.

72    Small fraction reads were left unassembled due to the longer read length, lower

73    coverage nature of 454 sequencing.  *Ab initio* ORF prediction was performed with

74    FragGeneScan v1.16 (9) with parameters: complete=0 and train=complete. ORFs were

75    once again screened for contamination in the form of rRNA, ITS, and primers. ITS

76    sequences were downloaded from NCBI, and reduced to 397,062 non-redundant

77    sequences at 0.95 level using cd-hit-est v4.6. Sequences that aligned with an ITS

78    sequence with BLASTN e-value <= 1e-5 were removed. Primer and adapter sequences

79    used in Illumina sequencing were searched using BLASTN and were identified at e-

80    value <= 10. Reads were removed with hits to terminal ends at least 10bp in length, or

81    internal hits at least 15bp in length. Possible organelle genes were classified for query

82    sequences that had closer homology to an organelle gene than a nuclear gene within

83    the organism with the closest known segregated organelle and nuclear genomes based

84    on best BLASTP e-value <= 1e-3.

85          ORFs were annotated via BLASTP (10,11) alignment (e-value threshold $1e^{-3}$) to

86    a comprehensive protein database, *phyloDB*, as well as screened for function *de novo*

87    by assigning Pfams, TIGRfams and transmembrane tmHMMs with hmmer 3.0

88    (http://hmmer.org/; 12) using an e-value threshold of $1.0e^{-4}$. PhyloDB version 1.076

89    consists of 24,509,327 peptides from 19,962 viral, 230 archaeal, 4910 bacterial, and

90    894 eukaryotic taxa (13–15). It includes peptides from the 410 taxa of the Marine

91    Microbial Eukaryotic Transcriptome Sequencing Project

92   ([http://marinemicroeukaryotes.org/](http://marinemicroeukaryotes.org/)), as well as peptides from KEGG, GenBank, JGI,

93   ENSEMBL, CAMERA to KEGG, GenBank, JGI, ENSEMBL, iMicrobe, and the

94   Chloroplast Genome Database (cpbase). Taxonomic annotation of ORFs was also

95   conducted via a BLASTP to phyloDB, and a Lineage Probability Index (LPI) was

96   calculated to avoid biases introduced by classifying ORFs based on best BLAST hit

97   alone (7,16,17). Briefly, LPI was calculated here as a value between 0 and 1 indicating

98   lineage commonality among the top 95-percentile of sequences based on BLAST bit-

99   score.

100      Illumina sequencing of large size fraction total ribosomal-depleted RNA yielded

101   623,461,310 raw reads across 16 time points, of which 265,345,754 were mRNA

102   (~43%). A total of 283,760 contigs were assembled upon which 345,355 ORFs were

103   called.  32,271,421 reads (~12%) mapped to the 111,655 ORFs that remained after

104   strict filtering (~32%). Ottesen et al.'s GS FLX Titanium (Roche) sequencing of small

105   fraction cDNA yielded 9,985,281 raw reads across 13 time points (1). 2,802,084 ORFs

106   were called on the 5,618,280 trimmed reads remaining after quality control (~56%).  Full

107   assembly and annotation statistics in Supplementary Dataset 1. Coverage across taxa

108   groups can be seen in Figure S3A.

109

110   *Mapping to reference transcriptomes*

111      Reference transcriptomes were chosen for read mapping that appeared with high

112   abundance and percent identity among *ab initio* large fraction ORFs. Representative

113   references were chosen from all major taxonomic groups found in *ab initio* ORFs.

114   Large and small fraction reads were aligned to reference ORFs using BWA-MEM

115    version 0.7.12-r1039 (18,19) using default parameters. At least 50% of each read must

116    map to a reference gene at least 80% identity to be considered a hit. References with at

117    least 1000 genes with at least 5 reads mapped were functionally annotated via rap v0.4

118    as above and considered for downstream analysis. Coverage of annotated references

119    can be seen in Figure S3B.

120

121    *Hierarchical clustering*

122          Reference transcriptome ORFs were hierarchically clustered together with all

123    large and small fraction *ab initio* ORFs (including organellar ORFs) to form peptide

124    ortholog groups via the Markov Cluster Algorithm (MCL; https://micans.org/mcl/; 20).

125    Directional edge weights were defined as the ratio of pairwise- to self- BLASTP scores,

126    and default parameters were used to assign ORFs to clusters. Clusters were assigned a

127    consensus annotation if found to be statistically enriched in that annotation with a

128    Fisher's exact test ($p < 0.05$).  Consensus annotations must also represent at least 10%

129    of the reads in the cluster and account for a minimum of 200 reads. Clusters with

130    identical consensus annotations were grouped together into "functional clusters."

131

132    *Identification of significantly periodic ORFs*

133          ORFs with significantly periodic diel expression were identified using harmonic

134    regression analysis (HRA) as previously described (1,21,22). Briefly, for each taxa

135    group of interest, raw ORF counts over time were fit to a generalized linear model (glm)

136    of a sinusoid with a 24-hour period with taxa-specific library sums serving as an offset at

137    each time point. The model was constructed using the "glm" function in R (23) and

138   statistical significance was determined by False Discovery Rate (FDR; (24) adjusted p-

139   values of <= 0.1 (Benjamini-Hochberg), on both a permutation test (500-50,000

140   permutations) and a chi-squared test.

141

142   *Identification of conserved expression modules*

143        The Weighted Gene Correlation Network Analysis (WGCNA) R package (25,26)

144   was used as previously described (22) to identify modules of conserved expression

145   among reference ORFs (Figure S13) and functional clusters (Figures 2, 3E, S5).

146   ORFs/functional clusters with at least 10 raw counts in at least 80% of time points were

147   considered. For Figure 3E, it was further stipulated that clusters must have at least 100

148   raw reads overall to be considered. In all cases, expression was normalized by total

149   pre-filtration counts at a single time point (library) before constructing a Pearson

150   correlation matrix. An adjacency matrix was then constructed from the correlation matrix

151   by applying a power function (AF(s)=s^b). The lowest b value that allowed for a scale-

152   free topology R-squared value above 0.8 was chosen, as recommended in the WGCNA

153   user manual, in order to optimize the mean number of connections of the network while

154   preserving scale independence. A signed Topological Overlap Matrix (TOM) was

155   constructed from the correlation matrix to measure dissimilarity between each pair of

156   nodes based on shared neighbors. Average linkage hierarchical clustering was used to

157   define a dendrogram (cluster tree) of the network via the "blockwiseModules" function. A

158   cut height of .995 as well as a minimum module size of 30 ORFs/functional clusters was

159   used to delineate branches of the hierarchical clustering tree into modules of co-

160   expression. The "moduleEigengenes" function was used with default parameters to

161   calculate "eigengenes" (a measure of "average" expression calculated as the first

162   principal component of the module's expression matrix) for each module.  Modules with

163   correlated eigengenes were merged by setting a "mergeCutHeight" threshold of 0.5.

164   ORFs/functional clusters with a correlation of less than 0.3 to their respective module

165   eigengene were removed and classified as "unassigned" (module 0).  The igraph

166   package (27) was used to visualize expression networks.

167

168   *Differential expression analysis*

169        Differential expression of ORFs, ortholog clusters, taxa groups, and genera

170   across size classes was identified using the R package edgeR (28). Categories (i.e.

171   ORFs, ortholog clusters, taxa groups, genera) with least 1 read per million in at least 3

172   samples were included and used to calculate log fold changes.  Counts were

173   normalized using the "calcNormFactors" function, which accounts for both library size

174   and varied library composition.  An exact test with tagwise dispersion estimation was

175   used to determine ORFs or clusters with significantly different expression across size

176   classes (FDR-corrected p < 0.05).

177

178   *Identification and phylogenetic analyses of LOV domain containing transcripts*

179        A two-step approach was taken to identify LOV domain containing proteins in our

180   reference and *ab initio* transcript sets. LOV domains are a subset of the PAS domain

181   family, and initial survey of a number of known LOV domain proteins using InterproScan

182   (29) suggested the PAS_9 Pfam domain (PF13426) has the highest similarity to the LOV

183   domain. For this reason, we curated a list of transcripts harboring the PAS_9 domain

184    (detail of domain annotation is provided in the 'Bioinformatics analysis of

185    metatranscriptomes' section). LOV domains have a signature motif that has a conserved

186    cysteine at the fourth position, however, some degeneracy can exist at other positions of

187    this domain (30). Given this fact, we further screened the amino acid sequences of the

188    transcripts harboring the PAS_9 domain for the presence of previously-identified LOV

189    specific motifs (30). We constructed a maximum likelihood phylogenetic tree from only

190    the regions of the proteins that aligned to the PAS_9 HMM. Alignment was performed

191    using MUSCLE (31). The tree was constructed in PhyML (32) with aLRT-SH like node

192    support. The tree and the heatmap of the expression profile were visualized in the

193    interactive Tree of Life (33).

194

195    *Analyses of circular "time of day" data*

196        The R package "circular" (34) was used to conduct circular statistics on time-of-

197    day data with a 0-24 hour range.  This includes peak time of day comparisons,

198    calculating the mean peak time of expression of a group of periodic ORFs, and statistics

199    on the photosynthetic cascade (Figure 5): Watson-Wheeler Test of homogeneity of

200    means, Watson-Williams Test of homogeneity of means (35).

201

202    *Figures*

203        Sorting and plotting of data was conducted in R version 3.2.1 (2015-06-18) using

204    the following packages: plyr (36), dplyr (37), reshape2 (38), ggplot2 (39), lubridate (40),

205    ggmap (41), gridExtra (42).

206

207     ***Bioinformatic analysis of amplicon data***

208     *rRNA read processing and annotation*

209          16S and 18S rRNA 454 reads were demultiplexed using Roche/454's sfffile utility

210     and converted from standard flowgram to fasta format using sff2fastq (https://

211     github.com/indraniel/sff2fastq). Primer removal, quality control, trimming, dereplication,

212     and taxonomic annotation were conducted using an in-house rRNA pipeline

213     (https://github.com/allenlab/rRNA_pipeline). Chimeric sequences were removed using

214     USEARCH (43), reads were trimmed to a quality score of 10 over a 2 base window,

215     operational taxonomic units were clustered using SWARM (44) and classified using

216     FASTA36 from the FASTA package

217     (http://faculty.virginia.edu/wrpearson/fasta/fasta36/). Taxonomic annotations were

218     assigned by using GLSEARCH36 (45) with the version 119 of the SILVA reference

219     database (46) for 16S rRNA and a modified PR2 database with updates from Tara

220     Oceans W2 (47) for 18S rRNA.

221          Roche 454 sequencing of large fraction amplicons across all 16 time points

222     yielded 820,700 raw 16S rRNA reads and 970,927 raw 18S rRNA reads, of which

223     36.7% (301,244) and 38.7% (375,904), respectively, remained after filtration.  After de-

224     replication, the large fraction contained 8,522 unique 18S and 5,420 unique 16S reads

225     (1,595 of which were plastid in origin). Full pipeline statistics shown in Supplementary

226     Dataset 2.

227

228     *Phylogenetic placement*

229    rRNA amplicons were processed against rRNA reference covariance models

230    using Infernal (48). A blastn (11) search was performed against SILVA (46) with *e*-value

231    threshold ≤ 1E-100 to identify representative (reference) sequences to be included in

232    the reference phylogenetic trees (eukaryotic 18S, bacterial 16S, and plastidic 16S).

233    Reference sequences were then aligned with MAFFT (49) using the G-INS-i setting for

234    global homology. The generated multiple sequence alignments were visually inspected,

235    manually edited and refined using JalView (50). Maximum likelihood reference trees

236    were inferred under the general time-reversible model with gamma-distributed rate

237    heterogeneity using FastTree (51). Processed rRNA sequences were mapped onto the

238    corresponding reference trees using pplacer (52) with the default settings. The number

239    of the mapped sequences to trees nodes was normalized to the total number of mapped

240    sequences from the corresponding samples. Normalized abundances were visualized

241    as circles mapped onto the reference trees such that the diameters of the circles were

242    proportion to the taxonomic abundances.

243

# 244    Supplementary Results and Discussion

245    ***A molecular window into biogeochemistry***

246    Several lines of evidence indicated that during the drift cells were experiencing

247    iron-limitation and that this factor shaped community composition. WGCNA was used to

248    examine functional clusters related to nutrient cycling (Figure S5). The expression of

249    several low-iron response genes, including iron-starvation-induced proteins ISIP1,

250    ISIP2A (phytotransferrin; (53), ISIP2B, and ISIP3 (54) in diatoms and haptophytes

251    (module 5) indicated cellular iron stress. Phytotransferrin and "silicon transporter"

252    annotations clustered into the same module these iron-response genes (Figure S5;

253    module 5; green) and were dominated by centric diatom expression. Module 5 peaks

254    sharply at the end of the drift track when the measured silica:nitrate ratio, which was

255    initially around 1, dropped most dramatically (Figure S1B). Low silica:nitrate ratios (in

256    the range of 0.8 to 1.1) have been observed in association with iron limitation (55) and

257    are thought to result from silica draw-down by iron-stressed diatoms (56).  In our study,

258    this ratio dropped by an order of magnitude along the drift track. While the main feature

259    of module 5 is its sharp peak at the end of the drift track, there also appears to be some

260    underlying diel periodicity in the signal (upregulated more during night hours). This

261    convolution of expression patterns speaks to the difficulty of teasing apart various

262    physical drivers of transcription in a dynamic natural context.  In the future, sampling for

263    a longer period of time could provide the statistical resolution to address these

264    questions more conclusively.

265         Additional molecular evidence supported iron-limitation, such as high expression

266    of iron complex outer membrane receptor proteins, which are associated with the

267    uptake of siderophores (57), especially in the small fraction. Indeed, "iron complex outer

268    membrane receptor protein" was the 34th most highly expressed annotation across both

269    size classes (Supplementary Data 6). Furthermore, nitrate transporters and reductases

270    were nearly undetectable across phytoplankton lineages whereas ammonium

271    transporters were highly expressed, reflecting a reduced capacity for nitrate

272    assimilation, which requires iron-rich heme cofactors (54). Finally, relative levels of

273    ferredoxin and flavodoxin among photosynthetic organisms are often used as an

274    indicator of iron stress (58,59). Iron-intensive ferredoxin proteins can be substituted by

275 flavodoxin, which performs the same role in photosynthesis but uses flavin

276 mononucleotides in place of iron-sulfur clusters. Strikingly, in the large fraction, 98.12%

277 of expression of these ORFs was attributed to flavodoxin (Pfam PF00258) over

278 ferredoxin (Pfam PF00111) across the five major eukaryotic phytoplankton lineages

279 (diatoms, chlorophytes, dinoflagellates, haptophytes, and pelagophytes; Figure S5).

280 This ratio falls at the extreme edge of the distribution of previously observed cases (60),

281 associated with the lowest iron concentrations. While direct trace-metal clean

282 measurement of iron concentrations was not possible due to the nature of the robotic

283 sampling, the gene sensors, historic oceanographic context (Figure S4), and nutrient

284 proxies we present here establish a high likelihood of iron limitation.

285

286 ***Timing of diel ORF expression across taxonomic groups***

287       Periodic ORFs were only detected in Proteobacteria known to possess

288 proteorhodopsin or carry out anoxygenic photosynthesis. However, previous

289 observations that benefitted from a longer time course were able to detect a greater

290 diversity of periodic ORFs in heterotrophic bacterioplankton (22), indicating that these

291 results only capture the strongest oscillating genes. We observed only a single periodic

292 ORF in the photoheterotrophic bacteria, SAR11 and SAR116, both peaking around

293 12:30p.m. In SAR11, this ORF was in the isocitrate lyase family. This points to a

294 daytime-shifted metabolism resulting from phototrophy; Glyoxylate shunt genes have

295 been found to be expressed 300-fold more in light than darkness in the proteorhodopsin

296 phototroph, *Dokdonia* (61), and were up during the day in photoheterotrophic

297 bacterioplankton in a previous drift track (22). With the exception of a *Pseudomonas*

298 *marR* family transcriptional regulator peaking around 10a.m., all other periodic

299 proteobacterial ORFs occurred in *Rhodobacter* species and were involved in nighttime

300 chelatase, light-independent protochlorophyllide reductase, bacteriochlorophyll

301 synthase, amine oxidases involved in carotenoid biosynthesis, diheme cytochrome c).

302

303 **Viruses in the small size class**

304   In the small fraction, bacteriophages corresponding to several of the abundant

305 bacterial groups were enriched (e.g. *Pelagibacter* phage, $\log_2$FC=11; *Roseobacter*

306 phage, $\log_2$FC=7.4; *Vibrio* phage, $\log_2$FC=5.2; Supplementary Dataset 10), as was a

307 virus best annotated as the uncultivated *Ostreococcus* OsV5 virus for which the host

308 *Ostreococcus* clade is still unclear (15). However, despite the picoprasinophytes

309 *Ostreococcus* ($\log_2$FC=2.25) and *Bathycoccus* ($\log_2$FC=2.36) being enriched in the

310 small size class, the largest signal for viruses most similar to isolates known to infect

311 them (i.e. OlV1, OtV1, and BpV1), came from the large fraction. This could indicate that

312 the close proximity of cells in particle-associated microenvironments promotes infection,

313 that infected cells more easily attach to particles (13), or that infected hosts are larger in

314 size.

315

# References

317 1.   Ottesen EA, Young CR, Eppley JM, Ryan JP, Chavez FP, Scholin CA, et al.

318      Pattern and synchrony of gene expression among sympatric marine microbial

319      populations. Proc Natl Acad Sci. 2013 Feb 5;110(6):E488–97.

320  2.  Scholin CA, Birch J, Jensen S, III RM, Massion E, Pargett D, et al. The quest to

321     develop ecogenomic sensors: A 25-year history of the Environmental Sample

322     Processor (ESP) as a case study. Oceanography. 2017;30(4):100–13.

323  3.  Timothy Pennington J, Chavez FP. Seasonal fluctuations of temperature, salinity,

324     nitrate, chlorophyll and primary production at station H3/M1 over 1989-1996 in

325     Monterey Bay, California. Deep Res Part II Top Stud Oceanogr. 2000;47(5–

326     6):947–73.

327  4.  Lane DJ, Pace B, Olsen GJ, Stahl DA, Sogin ML, Pace NR. Rapid determination

328     of 16S ribosomal RNA sequences for phylogenetic analyses. Proc Natl Acad Sci.

329     1985;82(20):6955–9.

330  5.  Herlemann DPR, Labrenz M, Jürgens K, Bertilsson S, Waniek JJ, Andersson AF.

331     Transitions in bacterial communities along the 2000 km salinity gradient of the

332     Baltic Sea. ISME J. 2011;5(10):1571–9.

333  6.  Stoeck T, Bass D, Nebel M, Christen R, Jones MDM, Breiner HW, et al. Multiple

334     marker parallel tag environmental DNA sequencing reveals a highly complex

335     eukaryotic community in marine anoxic water. Mol Ecol. 2010;19(SUPPL. 1):21–

336     31.

337  7.  Bertrand EM, McCrow JP, Moustafa A, Zheng H, McQuaid JB, Delmont TO, et al.

338     Phytoplankton-bacterial interactions mediate micronutrient colimitation at the

339     coastal Antarctic sea ice edge. Proc Natl Acad Sci U S A. 2015 Aug

340     11;112(32):9938–43.

341  8.  Schmieder R, Lim YW, Edwards R. Identification and removal of ribosomal RNA

342     sequences from metatranscriptomes. Bioinformatics. 2012 Feb 1;28(3):433–5.

343  9.   Rho M, Tang H, Ye Y. FragGeneScan: predicting genes in short and error-prone

344      reads. Nucleic Acids Res [Internet]. 2010 Nov 1 [cited 2016 Dec 1];38(20):e191–

345      e191. Available from: https://academic.oup.com/nar/article-

346      lookup/doi/10.1093/nar/gkq747

347  10.  Protein BLAST: search protein databases using a protein query [Internet].

348      Available from: https://blast.ncbi.nlm.nih.gov/Blast.cgi?PAGE=Proteins

349  11.  Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment

350      search tool. J Mol Biol [Internet]. 1990 Oct [cited 2017 Jul 17];215(3):403–10.

351      Available from: http://linkinghub.elsevier.com/retrieve/pii/S0022283605803602

352  12.  Sonnhammer E, Eddy SR, Birney E, Bateman A, Durbin R. Pfam: multiple

353      sequence alignments and HMM-profiles of protein domains. Nucleic Acids Res.

354      1998 Jan 1;26(1):320–2.

355  13.  Allen LZ, Allen EE, Badger JH, McCrow JP, Paulsen IT, Elbourne LD, et al.

356      Influence of nutrients and currents on the genomic composition of microbes

357      across an upwelling mosaic. ISME J. 2012;6(7):1403–14.

358  14.  Dupont CL, Mccrow JP, Valas R, Moustafa A, Walworth N, Goodenough U, et al.

359      Genomes and gene expression across light and productivity gradients in eastern

360      subtropical Pacific microbial communities. ISME J. 2014;doi(10).

361  15.  Zeigler Allen L, McCrow JP, Ininbergs K, Dupont CL, Badger JH, Hoffman JM, et

362      al. The Baltic Sea virome: diversity and transcriptional activity of DNA and RNA

363      viruses. mSystems. 2017;2(1):e00125-16.

364  16.  Podell S, Gaasterland T. DarkHorse: A method for genome-wide prediction of

365      horizontal gene transfer. Genome Biol. 2007;8(2).

366    17.    Bender SJ, Moran DM, McIlvin MR, Zheng H, McCrow JP, Badger J, et al. Colony

367            formation in *Phaeocystis antarctica*: Connecting molecular mechanisms with iron

368            biogeochemistry. Biogeosciences. 2018;15(16):4923–42.

369    18.    Bayat A, Gaëta B, Ignjatovic A, Parameswaran S. Improved VCF normalization

370            for accurate VCF comparison. Bioinformatics. 2017 Mar 16;33(7):964–70.

371    19.    Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler

372            transform. Bioinformatics. 2009 Jul 15;25(14):1754–60.

373    20.    Enright AJ, Van Dongen S, Ouzounis CA. An efficient algorithm for large-scale

374            detection of protein families. Nucleic Acids Res. 2002;30(7):1575–84.

375    21.    Ottesen EA, Young CR, Gifford SM, Eppley JM, Marin R, Schuster SC, et al.

376            Multispecies diel transcriptional oscillations in open ocean heterotrophic bacterial

377            assemblages. Science (80- ). 2014;345(6193).

378    22.    Aylward FO, Eppley JM, Smith JM, Chavez FP, Scholin CA, DeLong EF.

379            Microbial community transcriptional networks are conserved in three domains at

380            ocean basin scales. Proc Natl Acad Sci. 2015;112(17):5443–8.

381    23.    R Development Core T e. a. m. R: a language and environment for statistical

382            computing [Internet]. Vienna, Austria: R Foundation for Statistical Computing;

383            2011. Available from: http://www.r-project.org/

384    24.    Benjamini Y, Hochberg Y. Controlling the false discovery rate: a practical and

385            powerful approach to multiple testing. J R Stat Soc. 1995;57(1):289–300.

386    25.    Langfelder P, Horvath S. WGCNA: An R package for weighted correlation

387            network analysis. BMC Bioinformatics. 2008;9.

388    26.    Zhang B, Horvath S. A General Framework for Weighted Gene Co-Expression

389         Network Analysis. Stat Appl Genet Mol Biol. 2005;4(1).

390   27.   Csárdi G, Nepusz T. The igraph software package for complex network research.

391         [cited 2017 Jul 26]; Available from:

392         http://www.necsi.edu/events/iccs6/papers/c1602a3c126ba822d0bc4293371c.pdf

393   28.   Robinson MD, McCarthy DJ, Smyth GK. edgeR: A Bioconductor package for

394         differential expression analysis of digital gene expression data. Bioinformatics.

395         2009 Jan 1;26(1):139–40.

396   29.   Jones P, Binns D, Chang HY, Fraser M, Li W, McAnulla C, et al. InterProScan 5:

397         Genome-scale protein function classification. Bioinformatics. 2014;30(9):1236–40.

398   30.   Glantz ST, Carpenter EJ, Melkonian M, Gardner KH, Boyden ES, Wong GK-S, et

399         al. Functional and topological diversity of LOV domain photoreceptors. Proc Natl

400         Acad Sci. 2016;113(11):E1442–51.

401   31.   Edgar RC. MUSCLE: Multiple sequence alignment with high accuracy and high

402         throughput. Nucleic Acids Res. 2004;32(5):1792–7.

403   32.   Guindon S, Lethiec F, Duroux P, Gascuel O. PHYML Online - A web server for

404         fast maximum likelihood-based phylogenetic inference. Nucleic Acids Res.

405         2005;33(SUPPL. 2).

406   33.   Letunic I, Bork P. Interactive tree of life (iTOL) v3: an online tool for the display

407         and annotation of phylogenetic and other trees. Nucleic Acids Res.

408         2016;44(W1):W242–5.

409   34.   Agostinelli C, Lund U. R package "circular": Circular Statistics [Internet]. 2017.

410         Available from: https://r-forge.r-project.org/projects/circular/

411   35.   Tasdan F, Yeniay O. Power study of circular anova test against nonparametric

412          alternatives. Hacettepe J Math Stat. 2014;43(1):97–115.

413    36.    Wickham H. The split-apply-combine strategy for data analysis. J Stat Softw.

414          2011;40(1).

415    37.    Wickham H, Francois R, Henry L, Müller K. dplyr: A grammar of data

416          manipulation. 2017.

417    38.    Wickham H. Reshaping data with the reshape package. 2006 [cited 2017 Jul 25];

418          Available from: http://had.co.nz/reshape

419    39.    Wickham H. Ggplot2 : elegant graphics for data analysis [Internet]. Springer; 2009

420          [cited 2017 Jul 25]. 212 p. Available from:

421          http://www.citeulike.org/group/18896/article/6995399

422    40.    Grolemund G, Wickham H. Dates and times made easy with lubridate. JSS J Stat

423          Softw [Internet]. 2011 [cited 2017 Jul 25];40(3). Available from:

424          http://www.jstatsoft.org/

425    41.    Kahle D, Wickham H. ggmap: Spatial visualization with ggplot2. R J [Internet].

426          2013;5(1):144–61. Available from: http://journal.r-project.org/archive/2013-

427          1/kahle-wickham.pdf

428    42.    Auguie B. Miscellaneous functions for "grid" graphics [Internet]. 2016. p. 10.

429          Available from: https://github.com/baptiste/gridextra

430    43.    Edgar RC. Search and clustering orders of magnitude faster than BLAST.

431          Bioinformatics. 2010;26(19):2460–1.

432    44.    Mahé F, Rognes T, Quince C, de Vargas C, Dunthorn M. Swarm v2: highly-

433          scalable and high-resolution amplicon clustering. PeerJ. 2015;3:e1420.

434    45.    Pearson WR. Finding protein and nucleotide similarities with FASTA. Curr Protoc

435      Bioinforma. 2016;2016(March):3.9.1-3.9.25.

436  46.  Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, et al. The SILVA

437      ribosomal RNA gene database project: Improved data processing and web-based

438      tools. Nucleic Acids Res. 2013;41(D1):590–6.

439  47.  De Vargas C, Audic S, Henry N, Decelle J, Mahé F, Logares R, et al. Eukaryotic

440      plankton diversity in the sunlit ocean. Science (80- ). 2015;348(6237):1261605.

441  48.  Nawrocki EP, Eddy SR. Infernal 1.1: 100-fold faster RNA homology searches.

442      Bioinformatics. 2013 Nov;29(22):2933–5.

443  49.  Katoh K, Standley DM. MAFFT multiple sequence alignment software version 7:

444      improvements in performance and usability. Mol Biol Evol. 2013 Apr;30(4):772–

445      80.

446  50.  Waterhouse AM, Procter JB, Martin DMA, Clamp M, Barton GJ. Jalview Version

447      2-A multiple sequence alignment editor and analysis workbench. Bioinformatics.

448      2009 May;25(9):1189–91.

449  51.  Price MN, Dehal PS, Arkin AP. FastTree 2 - Approximately maximum-likelihood

450      trees for large alignments. PLoS One. 2010 Mar;5(3):e9490.

451  52.  Matsen FA, Kodner RB, Armbrust EV. pplacer: linear time maximum-likelihood

452      and Bayesian phylogenetic placement of sequences onto a fixed reference tree.

453      BMC Bioinformatics. 2010 Oct;11:538.

454  53.  McQuaid JB, Kustka AB, Oborník M, Horák A, McCrow JP, Karas BJ, et al.

455      Carbonate-sensitive phytotransferrin controls high-affinity iron uptake in diatoms.

456      Nature. 2018;555(7697):534–7.

457  54.  Allen AE, LaRoche J, Maheswari U, Lommer M, Schauer N, Lopez PJ, et al.

458       Whole-cell response of the pennate diatom Phaeodactylum tricornutum to iron

459       starvation. Proc Natl Acad Sci. 2008;105(30):10438–43.

460  55.  Hutchins DA, Bruland KW. Iron-limited diatom growth and Si:N uptake ratios in a

461       coastal upwelling regime. Nature. 1998 Jun 11;393(6685):561–4.

462  56.  Brzezinski MA, Krause JW, Bundy RM, Barbeau KA, Franks P, Goericke R, et al.

463       Enhanced silica ballasting from iron stress sustains carbon export in a frontal

464       zone within the California Current. J Geophys Res C Ocean. 2015;120(7):4654–

465       69.

466  57.  Tang K, Jiao N, Liu K, Zhang Y, Li S. Distribution and functions of tonb-dependent

467       transporters in marine bacteria and environments: Implications for dissolved

468       organic matter utilization. PLoS One. 2012;7(7).

469  58.  McKay RML, La Roche J, Yakunin AF, Durnford DG, Geider RJ. Accumulation of

470       ferredoxin and flavodoxin in a marine diatom in response to Fe. J Phycol.

471       1999;35(3):510–9.

472  59.  Erdner DL, Anderson DM. Ferredoxin and flavodoxin as biochemical indicators of

473       iron limitation during open-ocean iron enrichment. Limnol Oceanogr.

474       1999;44(7):1609–15.

475  60.  Carradec Q, Pelletier E, Da Silva C, Alberti A, Seeleuthner Y, Blanc-Mathieu R, et

476       al. A global ocean atlas of eukaryotic genes. Nat Commun. 2018;9(1).

477  61.  Palovaara J, Akram N, Baltar F, Bunse C, Forsberg J, Pedros-Alio C, et al.

478       Stimulation of growth by proteorhodopsin phototrophy involves regulation of

479       central metabolic pathways in marine planktonic bacteria. Proc Natl Acad Sci.

480       2014;111(35):E3650–8.