



Prochlorococcus in Context:

Metagenomic Analysis of Sargasso Sea Microbiomes

Reid Pryzant '16 and Claire Ting - Department of Biology, Williams College



Introduction

Prochlorococcus is a numerically dominant photosynthetic prokaryote residing in subtropical and tropical oceans. It plays a key role in global carbon and energy cycles. Although members of the *Prochlorococcus* lineage are closely related, they have evolved significant genetic differences. The question our laboratory is addressing involves how this genetic diversity translates into variation in photosynthetic capacity and the ability to survive environmental stress.

In order to extend our understanding from the laboratory to the open ocean, we have conducted a metagenomic analyses of the Sargasso Sea, where *Prochlorococcus* is often predominant among bacterioplankton. We aim to characterize the Sargasso Sea microbiome and understand how environmental selection shapes *Prochlorococcus* populations at different depths in the water column. We hypothesize that the partitioning of bacterioplankton (and specifically *Prochlorococcus*) populations and functions will exhibit significant differences between near-surface (40 m) and deeper (100 m) waters. These differences should illuminate micro- and macro-scale heterogeneity in key physico-chemical properties and biological interactions in the open ocean. We will attempt to answer these questions through comparative taxonomic and functional analysis along with contrasting cultured and uncultured genomes.

Metagenomics

Metagenomics is the application of genomics to uncultured assemblies of organisms in the natural environment. It is a fairly recent field; the first paper on marine metagenomics was published in 2004 (1).

During August of 2009, the Ting Lab conducted field work in the Sargasso Sea (31° 40.00'N, 64° 10.00'W) on board the R/V Atlantic Explorer. For each environmental sample, we collected over 240 liters of sea water. We subsequently filtered the water (between 0.2 um and 1.6 um) to collect bacterioplankton with Steripak filters (Millipore). All samples were immediately stored in liquid nitrogen and chemical buffers. We then shipped the samples back to our laboratory at Williams College. The Ting lab spent close to year optimizing environmental DNA isolation before extracting the genetic material from our ocean water. We then sent the extracted DNA to the Department of Energy's Joint Genome Institute (Walnut Creek, CA) for Illumina sequencing. In February 2012, the Joint Genome Institute began the process of uploading the metagenomes into the IMG/M annotation pipeline.

Further analysis was carried out with the MG-RAST pipeline (2), MEGAN (3), STAMP (4), and in-house programs written in Python, Perl, and R. Computation was outsourced to Luis, a server in Jesup Hall running CentOS 6 with 64 AMD 6380 cores at 2500 MHz, 64 gigabytes of memory, and 2.7 terabytes of disk space. Some of the analyses were extremely computationally intensive. Taxonomic profiling of the 100 m dataset, for example, took 10¹³ comparisons and three weeks on Luis.

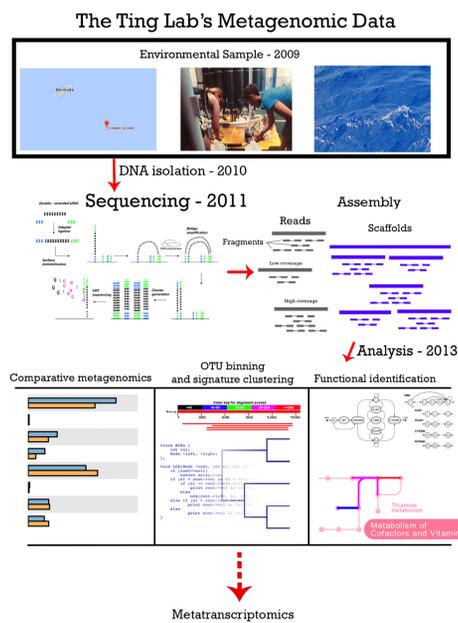


Figure 1: A schematic for the procedure used in metagenomics. It shows key steps from initial field work to the functional identification and taxonomic profiling of over 375,000 proteins. Future work involves metatranscriptomic analysis: the study of gene expression in these samples.

Table 1: Summary of the 40 m and 100 m Sargasso Sea data.

Table 2: The number significant alignments is reported for each strain of Prochlorococcus, along with the Prochlorococcus core genome.

Taxonomic Profiling

We hypothesized that the composition of bacterial phytoplankton at 40 m and 100 m will differ significantly. In order to describe the taxonomic distribution of our environmental samples, we identified genes in each scaffold and then assigned them to operational taxonomic units (OTU's) in a process called binning. Gene calling was carried out by comparing each scaffold to the NCBI-nr reference database. Genes were subsequently binned into OTU's with the least common ancestor labeled as "Other".

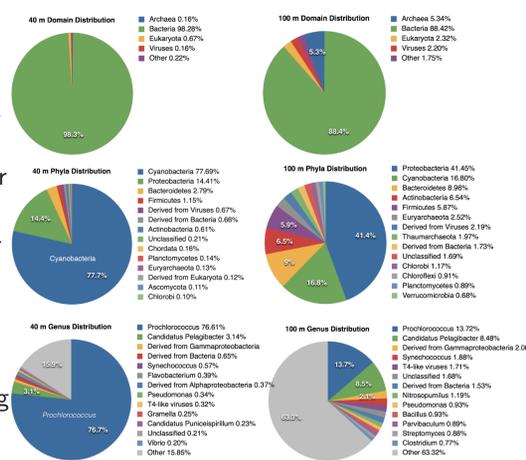


Figure 2: Taxonomic distribution of the 14 most populous groups at the Domain, Phyla, and Genus level. Groups outside the top 14 have been summed with the least common ancestor labeled as "Other".

12-strain Comparison

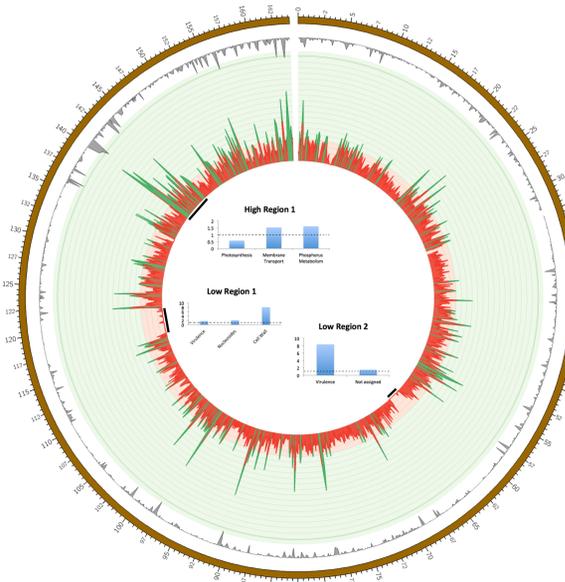


Figure 3: Map of the distribution of shared sequences between MIT 9301 and our environmental samples. This distribution is modeled by inner (40 m) and outer (100 m) lineplots with 0.5% scale. Position along the genome (bp x10,000) is shown around the edge. Red lines connect points <0.01%. Barplots of odds ratios between expected (genomic) and regional proportions show functional over- and underrepresentation in islands of interest. Hits were filtered by e-score (1e-10), ensuring that highly similar sequences are registered along with identical matches. The number of hits per gene is indicative of abundance and primary sequence conservation (8). "Low Regions" mark the locations of novel, unconserved genomic islands: hotspots of recombination and critical contributors to genetic diversity. "High regions" mark areas of conservation.

The number of matches between the environmental sequences and cultured strains appropriately reflects this distribution (6) (Table 2). While members of the large clade of recently differentiated lineages have more matches to subsequences of the 40 m sample, members of deeply branched lineages have more matches with the 100 m sample.

Functional Characterization

We expect the functional profiles of the water column to reflect micro- and macro-scale physico-chemical and biological characteristics. To describe the functional landscape of our metagenomes, we categorized proteins into hierarchical subsystems. Subsystems are a generalization of the term "pathway." Protein-coding genes were queried against the SEED database for subsystem classification (6) (Fig 4).

- The functional landscapes of the 40 m and 100 m bacterial communities are strongly correlated (r = .962) indicating conservation of broad community functions.
Prochlorococcus at 40 m has a similar functional landscape as its community (d = .89) but Prochlorococcus at 100 m is dissimilar to its community (d = 3.07). This is consistent with Prochlorococcus's dominance of the 40 m sample.
In both communities, Prochlorococcus accounts for the majority of photosynthetic potential.
Interestingly, the number of stress response genes associated with Prochlorococcus is higher at 100 m than at 40 m.
Communities lose their photosynthetic potential with depth.

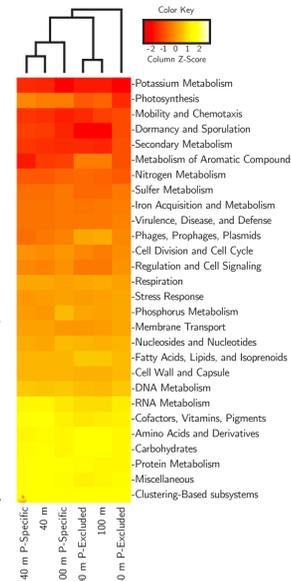


Figure 4: SEED subsystem distribution of all genes, Prochlorococcus-specific genes, and Prochlorococcus-excluded genes at 40 m and 100 m. Values are standardized log scores and hierarchical clustering is based on euclidian distance (d).

Signature Analysis

To learn more about uncultured, environmental Prochlorococcus at 40 m and 100 m, we used sequence signatures to examine inter-genus differentiation (7). We computed tetranucleotide profiles for Prochlorococcus-specific scaffolds with a 1-bp sliding window and summed pairs of reverse complementary tetranucleotides. Principal components analysis revealed distinct partitioning of Prochlorococcus gene signatures at 100 m. The smaller group (cluster 1) is closely related to MIT 9303, while the larger

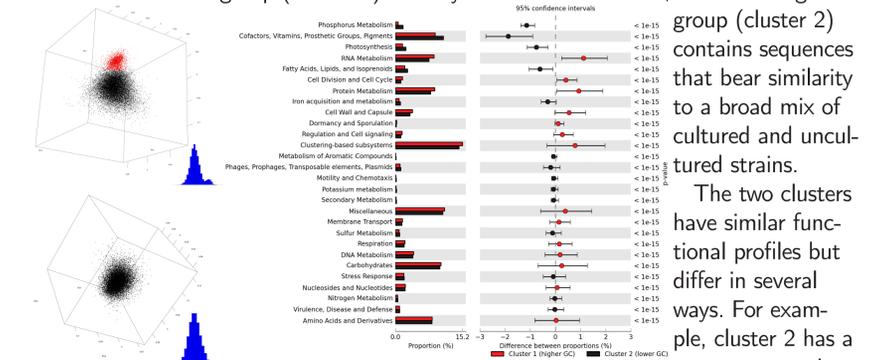


Figure 5: Plots of the first three principal components of Prochlorococcus's tetranucleotide profile at 100 m and 40 m (%[G+C] histograms, inset). The 100 m plot has been colored by its k-means clusters. On the right is a comparative functional analysis between the two 100 m clusters. Wilson confidence intervals and Bootstrap t-tests (100,000 replicates) were computed for each subsystem.

References and Acknowledgements

This work was funded by the National Science Foundation, award number MCB-0850900 to C.S. Ting. I would like to thank professor Ting for her constant help, knowledge and patience, professor Klingenberg, Adam Wang from OIT, and Jonathon Morgan-Leaman from OIT for access to Luis, and Kathleen for putting up with my antics in lab.
1. Venter, J. C. et al. (2004). Environmental genome shotgun sequencing of the Sargasso Sea. Science 304:66-74.
2. F. Meyer et al. (2008) The Metagenomics RAST server - A public resource for the automatic phylogenetic and functional analysis of metagenomes. BMC Bioinformatics, 9:386.
3. Mitra S, Klar B, Hoon DH (2011) Visual and statistical comparison of metagenomic communities. Bioinformatics 26:143-144.
4. Parks, D.H. and Beiko, R.G. (2010). Identifying biologically relevant differences between metagenomic communities. Bioinformatics, 26, 715-721.
5. Johnson ZI, Zinser ER, et al. (2006) Niche partitioning among Prochlorococcus ecotypes along ocean-scale environmental gradients. Science 311: 1737-1740.
6. Overbeek R, et al. (2005). The subsystems approach to genome annotation and its use in the project to annotate 1000 genomes. Nucleic Acids Res. 33(17):5691-5702.
7. Dick G, Andersson A, Baker B, Simmons S, Thomas B, Yelton AP et al. (2009). Community-wide analysis of microbial genome sequence signatures. Genome Biol 10: R85.
8. Coleman ML, Sullivan MB, Martiny AC, Steglich C, Barry K, et al. (2006) Genomic islands and the ecology and evolution of Prochlorococcus. Science 311: 1768-1770.
Photos in Fig. 1: Map: Google Earth. Illumina sequencing: http://www.genome.gov/aboutgov/content/figures/1297-9656-44-21-21.jpg. COG Markov Model: http://www.aquas.com/stefan/plc/LocApplication.png. BLAST output: NCBI web-BLAST. Assembly: adapted from http://genomebiology.com/2009/10/8/R85/figure/F1.