

Ecosystem-specific selection pressures revealed through comparative population genomics

Maureen L. Coleman^a and Sallie W. Chisholm^{a,b,1}

^aDepartment of Civil and Environmental Engineering, and ^bDepartment of Biology, Massachusetts Institute of Technology, Cambridge, MA 02139

Edited by David M. Karl, University of Hawaii, Honolulu, HI, and approved August 24, 2010 (received for review June 30, 2010)

Bacterial populations harbor vast genetic diversity that is continually shaped by abiotic and biotic selective pressures, as well as by neutral processes. Individuals coexisting in the same geographically defined population often have significantly different gene content, but whether this variation is largely adaptive or neutral remains poorly understood. Here we quantify heterogeneity in gene content for two model marine microbes, *Prochlorococcus* and *Pelagibacter*, within and between populations in the Atlantic and Pacific Oceans, to begin to understand the selective pressures that are shaping these “population genomes.” We discovered a large fraction of genes that are rare in each population, reflecting continual gene transfer and loss. Despite this high variation within each population, only a few genes significantly differ in abundance between the two biogeochemically distinct environments; nearly all of these are related to phosphorus acquisition and are enriched in the Atlantic relative to the Pacific. Moreover, P-related genes from the two sites form phylogenetically distinct clusters, whereas housekeeping genes do not, consistent with a recent spread of adaptive P-related genes in the Atlantic populations. These findings implicate phosphorus availability as the dominant selective force driving divergence between these populations, and demonstrate the promise of this approach for revealing selective agents in more complex microbial systems.

adaptation | biogeography | microbial evolution | natural selection | phosphorus limitation

Perhaps the most surprising lesson we have learned from microbial genome sequences is that closely related isolates of the same species often harbor substantially different gene complements, owing to horizontal gene transfer and gene loss (1, 2). The set of genes unique to a particular strain or subset of strains—the flexible genome—is hypothesized to be responsible for niche-specific adaptations (2, 3). However, in any given isolate, some fraction of the flexible genome might provide no fitness benefit, making it difficult to infer ecology from individual genome sequences. Moreover, snapshots of gene content in some natural populations have revealed high levels of coexisting variation (4–7), even among individuals that appear to experience similar selection pressures, suggesting that much of this variation is neutral. Alternatively, this variation could reflect adaptation by subpopulations specialized for microenvironments (8), or could reflect predation-driven frequency-dependent or diversifying selection (9). Understanding what fraction of gene content variation is in fact acted upon by selection, and what role this variation plays in ecosystem functioning, are central questions in microbial biology (10).

In the absence of selection, genes will be lost from a bacterial genome, owing to a mutational bias that favors deletions (11, 12). Therefore, genes that persist and rise to fixation in a population are inferred to be functional and to enhance organismal fitness. Likewise, genes that are differentially maintained in two populations reveal differential selection pressures acting on the two populations. Based on these principles, quantitative population genomics—in particular, the analysis of gene frequencies and patterns of sequence variation—has the power to highlight salient genetic features amid a background of continual gene transfer and gene erosion (10, 13). Thus, we can begin to quantify the roles of

horizontal gene transfer and the flexible genome in microbial adaptation across environments.

Here we apply this approach to advance our understanding of evolutionary dynamics and the selective pressures facing two model microbes, the marine cyanobacterium *Prochlorococcus* and the heterotroph *Pelagibacter* (a member of the SAR11 clade; ref. 14), in two distinct oceanic regions. These microbial groups numerically dominate the biogeochemically well-characterized North Atlantic and North Pacific subtropical gyres, hereinafter referred to by the representative long-term study stations at which we sampled, the Bermuda Atlantic Time Series (BATS) station and the Hawaii Ocean Time-Series (HOT) Station ALOHA (15). Both locations are oligotrophic with similar rates of primary production and carbon export, but BATS experiences stronger seasonal mixing and nutrient supply compared with HOT (15). BATS has lower phosphate concentrations (16) but higher fluxes of dust inputs, which bring iron and other metals (17). Thus, the two sites could potentially select for a range of distinct microbial traits. Moreover, gene content at these two sites can be framed in the context of reference genomes (12 for *Prochlorococcus* and 3 for *Pelagibacter*).

Specifically, we investigated the degree of heterogeneity in gene content among individuals in marine microbial populations, how much of the flexible genome is being maintained by selection, and what the functions of these adaptive genes tell us about ecosystem-specific selective pressures. Toward this end, we quantitatively compared gene frequencies among *Prochlorococcus* and *Pelagibacter* using pyrosequenced community DNA from HOT and BATS. Because *Prochlorococcus* comprises distinct high- and low-light-adapted clades (18) whose abundance varies with depth (19, 20), we sampled three analogous depths at each site to provide a more representative picture of the entire vertically integrated (meta)population (Table S1). In doing so, we captured similar abundances of high- and low-light-adapted clades at both sites, thereby minimizing the effect of clade structure on our intersite comparison (Table S2).

Results and Discussion

To begin to understand the extent of gene content variability in natural populations, and how much of this variability is likely to be adaptive, we first quantified the occurrence of each *Prochlorococcus* flexible gene relative to the core genome (21). Genes belonging to the core genome—the 1,221 single-copy genes shared by all 12 sequenced isolates of *Prochlorococcus*—appear to be shared by nearly all *Prochlorococcus* cells in these wild pop-

Author contribution: M.L.C. and S.W.C. designed research; M.L.C. performed research; M.L.C. analyzed data; and M.L.C. and S.W.C. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

Database deposition: The sequences reported in this paper have been deposited in the NCBI Sequence Read Archive, <http://www.ncbi.nlm.nih.gov>. For a list of accession numbers, see *SI Text*.

¹To whom correspondence should be addressed. E-mail: chisholm@mit.edu.

This article contains supporting information online at www.pnas.org/lookup/suppl/doi:10.1073/pnas.1009480107/-DCSupplemental.

ulations as well, as evidenced by the single narrow peak in the distribution of relative core gene frequencies (Fig. 1). Note that this distribution for core genes is, by definition, centered around one copy per cell, but the shape of the distribution is free to vary (see *Materials and Methods*). Compared with the core genome, most flexible genes (i.e., genes found in some, but not all, cultured isolates) are rare; 974 genes at HOT and 784 genes at BATS are estimated to occur at < 0.25 copy per cell (Fig. 1). A preponderance of rare flexible genes even at a given depth (Fig. S1) indicates that light adaptation cannot solely account for the distribution shown in Fig. 1. The majority of these rare flexible genes are likely neutral and transient in the population. In addition, some may have been acquired only recently, and some encode cell surface traits that may be driven by frequency-dependent selection (9). It is also possible that some of these rare flexible genes are adaptive in microenvironments within our samples, although such microenvironments are difficult to define for a planktonic organism. In contrast, about 400 flexible genes are estimated to occur at > 0.65 copy per cell at either site, and thus are “nearly core” in these particular populations (Fig. 1 and Table S3). Most of these widespread flexible genes are prevalent at both HOT and BATS, suggesting that they confer a fitness advantage in both open ocean environments, yet nearly three-quarters of these genes have no known function. One notable exception is the urease pathway; a dozen genes encoding urea transport and metabolism are found at about one copy per cell at both sites, suggesting that urea is an important nutrient source common to both these environments.

To identify selective pressures affecting the HOT and BATS populations differentially, we compared relative gene frequencies between the two sites. Only 29 of 2,854 observed *Prochlorococcus* genes have significantly different relative frequencies between the two sites (Fig. 24 and Table S4). Nearly all of the 29 differentially maintained genes are enriched at BATS relative to HOT and have functions related to phosphorus (P) acquisition and metabo-

lism, as inferred from homology to known P-uptake genes, up-regulation during P-starvation conditions in cultured strains MED4 (22) and MIT9301 (Fig. S2), or colocation with such genes in the genome (Fig. 2B). Many of these P-related genes are present in only a small fraction of *Prochlorococcus* cells at HOT but are nearly one copy per cell at BATS (Table S4); for instance, alkaline phosphatase (*phoA*) is estimated to be present in only 3% of *Prochlorococcus* cells at HOT but in 100% of these cells at BATS. A pathway for phosphonate utilization (23) is present in an estimated 13–21% of *Prochlorococcus* cells at BATS, but in virtually none of these cells at HOT (*phnYZ*; Table S4). Thus the “population genome” (i.e., the gene inventory in the *Prochlorococcus* collective) suggests strong selection for maintenance of accessory P- acquisition genes at BATS, but not at HOT. This difference between the two populations is consistent with exceptionally low surface water phosphate concentrations and physiological indicators of microbial phosphate limitation in the western North Atlantic (16, 24–27) relative to the Pacific. What is striking, however, is that this is the sole difference in gene content that emerged, implying that phosphate scarcity has been the most persistent and influential selective force driving diversification between the HOT and BATS populations.

To test whether these same selective pressures have impacted the genomes of coexisting microbes in the community, we repeated the analysis with putative *Pelagibacter* (a heterotrophic bacterium) sequences (Fig. 2C), because this is another abundant microbe at these two sites for which reference genomes exist. Of 1,667 observed gene clusters, only 31 differ in abundance between BATS and HOT, and 29 of these are enriched at BATS (Table S5). Twenty-seven of these genes are located in one contiguous region of the genome in *Pelagibacter* strain HTCC7211, and nearly all are involved in phosphate or phosphonate metabolism (Fig. 2D), demonstrating that adaptation to P scarcity is a broad feature of the BATS ecosystem. Notably, however, the enriched gene set for *Pelagibacter* is distinct from that of *Prochlorococcus*. For example, genes encoding high-affinity phosphate transport (*pstSCAB*) and polyphosphate storage and breakdown (*ppk* and *ppx*) are represented equally among *Prochlorococcus* at HOT and at BATS, because they are core genes in this group (21). In contrast, the frequency of these genes varies dramatically between HOT and BATS in *Pelagibacter* (Table S5), in which these are not core genes in cultured isolates. The acquisition of beneficial nutrient-scavenging genes, along with their large cellular surface-to-volume ratio, helps explain the observation that *Prochlorococcus* and *Pelagibacter*-like bacteria can account for 90% of phosphate uptake in the North Atlantic (28).

Our results reveal unanticipated adaptations to phosphorus scarcity as well. Arsenate reductase and an arsenite efflux pump are enriched at BATS in *Pelagibacter* and *Prochlorococcus*, respectively (Tables S4 and S5). When cells scavenge phosphate, they sometimes take up the structurally similar arsenate ion, which is then reduced to arsenite and exported as a detoxification mechanism (29). Although average concentrations of surface inorganic arsenic (predominantly arsenate) are similar in the North Pacific and Atlantic, the arsenate:phosphate ratio is several-fold higher at BATS, where there is generally more arsenate (average, 16.3 nM) than phosphate in surface waters (30). Our population genomic analysis suggests that arsenic toxicity is an important selective force influencing *Prochlorococcus* and *Pelagibacter* at BATS, but not at HOT.

Several evolutionary scenarios could have led to this differential gene content between the two sites. If an ancestor of the BATS population had acquired genes that enhanced its fitness (e.g., alkaline phosphatase), a selective sweep could have ensued, leading to ecologically and genetically distinct HOT and BATS ecotypes (31). Alternatively, frequent gene transfer could have introduced P-acquisition genes into diverse genome backgrounds relatively recently (31). To distinguish between these scenarios, we built

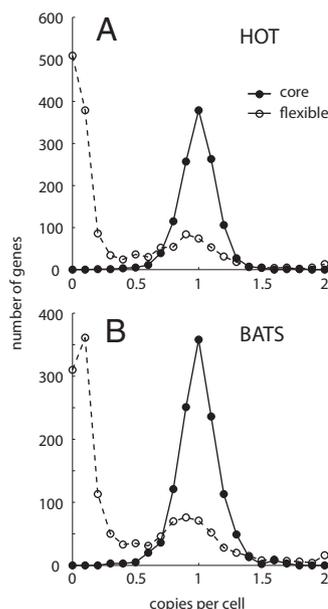


Fig. 1. Distribution of core and flexible genes (21) among *Prochlorococcus* populations at HOT (A) and at BATS (B). Core genes are defined as present in all 12 sequenced cultured isolates of *Prochlorococcus*, and the tight single-peaked distribution around the mean implies that they also are core genes and are present in one-to-one stoichiometry in these wild populations. Most flexible genes (those present in some, but not all, genomes of cultured isolates) are rare in these populations, being present in only a small proportion of the cells. The copy number per cell for each gene was estimated as described in *Materials and Methods*.

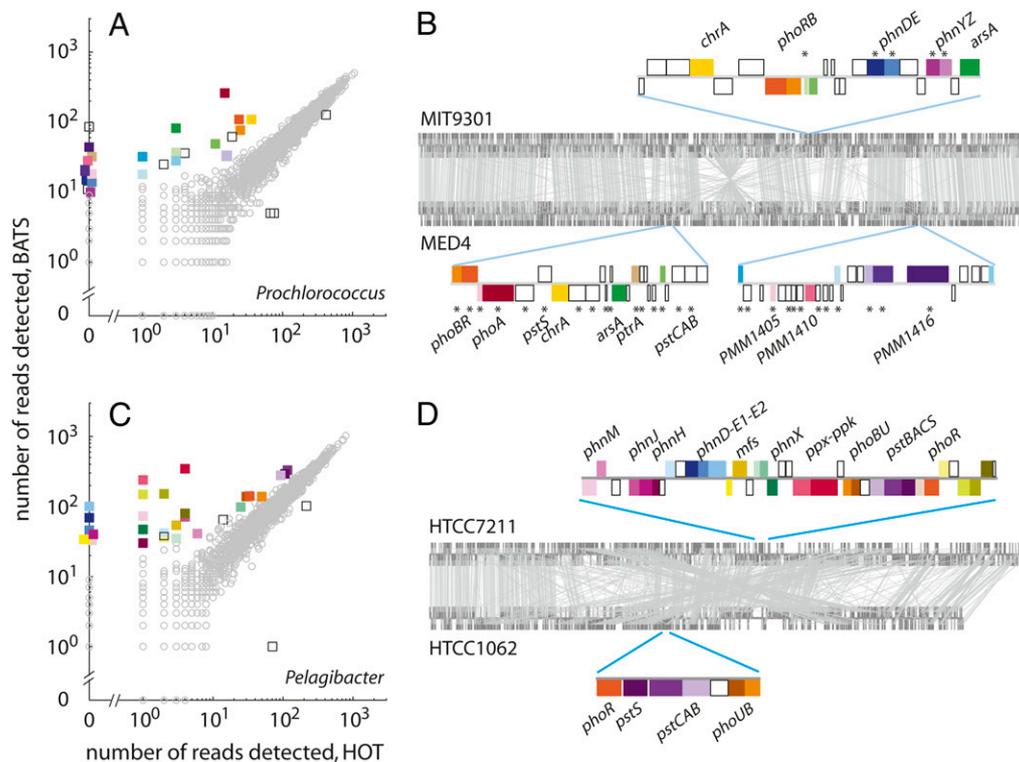


Fig. 2. Relative gene frequency among *Prochlorococcus* and *Pelagibacter* populations in the oligotrophic North Pacific (HOT) and North Atlantic (BATS) subtropical ocean gyres. (A) Detection of each *Prochlorococcus* gene at HOT and BATS, measured as the number of pyrosequencing reads. For each gene, the number of reads detected is proportional to the product of gene length and gene multiplicity per cell. Thus, assuming similar gene lengths at both sites, genes that fall along the diagonal trend have the same relative frequency at both sites, whereas those that fall above or below the diagonal are enriched at one site compared with the other. Squares represent significantly different frequencies between sites (G test; $P < 0.01$); circles represent nonsignificant differences. Colored squares represent genes whose chromosomal positions are depicted in B. (B) Genome comparison of two cultured isolates of *Prochlorococcus* (MED4 and MIT9301), with gray lines connecting homologous genes. Genes represented by corresponding colored squares in A are clustered together in a few distinct regions of the chromosome, including hypervariable genomic islands, in these isolates. Genes marked with an asterisk are up-regulated in response to P starvation (ref. 22 and Fig. S2). MED4 and MIT9301, isolated from the Mediterranean and the North Atlantic, respectively, are depicted because they carry the largest complements of phosphorus uptake genes among *Prochlorococcus* isolates (22). (C) Relative frequency of each *Pelagibacter* gene at HOT and BATS, as in A. (D) Genome comparison of two *Pelagibacter* strains, HTCC7211 and HTCC1062, as in B.

phylogenetic trees from both core genes and P-related flexible genes shared between the two populations and used Unifrac (32) to test whether the HOT- and BATS-derived alleles constituted distinct lineages. As exemplified by the *pntA* (encoding transhydrogenase) gene (Fig. 3A and B), the housekeeping genes tested in both *Prochlorococcus* and *Pelagibacter* do not show divergence between the two populations (Table S6). Several P-uptake genes in both *Prochlorococcus* and *Pelagibacter* do show distinct HOT and BATS lineages, however (e.g., the *pstB* trees shown in Fig. 3C and D). This phylogenetic divergence appears to be specific to P-uptake genes, and is not seen for the iron- and nitrogen-acquisition genes *fur* and *amtB* (Fig. 3E and F and Table S6). This result is consistent with previous studies suggesting that the evolutionary history of the core genome backbone is decoupled from that of the P-associated regions (22, 33, 34). Together these results point to relatively recent gene acquisition, for instance, via phage-mediated gene transfer associated with genomic islands (35), combined with efficient selection in large populations, as the major processes underlying gene content variability.

Over time, acquired genes can become core-like not only in their frequency in the population (i.e., one copy per cell), but also in their arrangement within individual genomes. For example, genes first acquired in hypervariable genomic islands can eventually migrate to more stable regions of the chromosome where insertions and deletions are less likely (36). Phosphate-acquisition genes are clustered in hypervariable regions in the genomes of isolates of both *Prochlorococcus* (22, 35) (Fig. 2B) and *Pelagibacter* (37) (Fig. 2D).

To assess how conserved the gene neighborhoods surrounding accessory P-uptake genes have become within natural populations at BATS, we used small-insert clone libraries, constructed from the same DNA used for pyrosequencing (see Materials and Methods). We identified clones in which at least one read matched a *Prochlorococcus* or *Pelagibacter* BATS-enriched gene (Tables S4 and S5). In 22% of such *Prochlorococcus* clones ($n = 307$), but in 40% of such *Pelagibacter* clones ($n = 554$), the paired-end read matches another BATS-enriched gene, suggesting that these genes are more often clustered on the chromosome in *Pelagibacter* than in *Prochlorococcus*. Furthermore, we counted the number of unique clusters of orthologous groups of proteins (COGs) adjacent to each BATS-enriched gene in *Prochlorococcus* and *Pelagibacter* clones, and found a greater diversity of adjacent COGs in *Prochlorococcus* (Fig. S3). Notably, core genes and BATS-enriched genes in *Pelagibacter* show similar levels of local gene order conservation, whereas in *Prochlorococcus*, local gene order is significantly more variable surrounding BATS-enriched genes compared with core genes (Fig. S3). Together, these results suggest that accessory P-uptake genes have become core-like in BATS populations of *Pelagibacter* in both frequency and arrangement, perhaps because *Pelagibacter* has experienced P limitation for a longer time; indeed, despite its smaller cell and genome size, *Pelagibacter* likely has a higher per-cell quota for P than *Prochlorococcus*, owing to its dependence on phospholipids (24).

Our analysis reveals that P availability is the primary selective force driving genome divergence between these two ocean regions.

Estimating Multiplicity per Cell. The number of DNA fragments, n_i , observed for a gene i of length L_i and average multiplicity per genome, m_i , in a sample of N_{Pro} sequences is expected to be binomially distributed,

$$n_i(L_i, m_i) \sim \text{binomial}\left(N_{Pro}, \frac{L_i m_i}{\sum L_i m_i}\right).$$

Thus, each sequencing read that we sample can be classified as either a “success” or “failure;” either it belongs to gene i or it does not. The probability of success (i.e., that the read came from a given gene i) depends on the length of the gene; longer genes will be detected more frequently. Indeed, using 1,221 single-copy core genes (i.e., $m_i = 1$, defined based on isolate genomes; ref. 21), we found a tight linear relationship between gene length and the number of sequencing reads detected ($r^2 = 0.95$ for HOT and 0.95 for BATS), and also found that the number of reads detected for nearly all core genes fell within the confidence intervals predicted by the binomial distribution assuming that $m_i = 1$. From this relationship, based on all 1,221 known single-copy genes, we then inferred gene multiplicity per cell for the “flexible” genome (*sensu* ref. 21); this approach is more robust than normalizing to a single core gene. Multiplicity per cell for gene i was calculated as

$$m_i = \frac{n_i}{bL_i},$$

where n_i is the number of sequence reads mapped to gene i , L_i is the gene length, and b is the slope of the length-versus-reads detected relationship for core genes. Confidence intervals for m_i were estimated using the binomial distribution. This approach was taken for *Pelagibacter* estimates of multiplicity per genome as well, using its respective core genome.

Validation of Taxon Mapping. An obvious limitation of our approach is that we can never know whether a pyrosequencing read with strong similarity to a particular organism’s genome was truly derived from a cell of that organism, or whether it resides in the genome of another organism due to horizontal gene transfer. To address this issue, we used paired-end reads from small-insert shotgun clone libraries constructed from the same DNA used for pyrosequencing. Clone libraries were constructed and sequenced by the Joint Genome Institute according to standard production protocols. We identified putative *Prochlorococcus* clones, in which at least one of the two paired-end reads matched *Prochlorococcus* as the best hit, and then identified the best-hit taxon for its paired end using BLASTN against a custom database of marine microbial genomes. If both paired-end reads match *Prochlorococcus*, we can be more confident that the clone came from a *Prochlorococcus* cell, although multigene horizontal transfer events are also possible. Only 9.3% of the putative *Prochlorococcus* clones at BATS and 8.6% of those at HOT matched *Prochlorococcus* on one end and a different taxon on the other end. We repeated the analysis for putative *Pelagibacter* clones and found that 18.2% of clones at BATS and 21.6% at HOT matched *Pelagibacter* on one end and a different taxon on the other end. Clones in which only one of two paired reads matched *Prochlorococcus* (or *Pelagibacter*) are still quite likely to have derived from a *Prochlorococcus* (or *Pelagibacter*) cell, given the high abundance of *Prochlorococcus* (and *Pelagibacter*) in these communities. There are more reference genomes available for *Prochlorococcus* than for *Pelagibacter* (12 vs. 3), providing a more comprehensive (although still incomplete) picture of the *Prochlorococcus* pan-genome; this larger set of reference sequences likely explains why paired-end reads were more often coincidentally identified as *Prochlorococcus* than as *Pelagibacter*.

Synteny Analysis. From these same clone libraries, we identified clones in which at least one read matched a BATS-enriched gene from *Prochlorococcus* or *Pelagibacter*, and clones in which at least one read matched a core gene from either taxon. We then compared the diversity of the paired-end reads associated with BATS-enriched genes and with core genes. We used rpsblast (47) to identify the best COG match for each paired-end read, and counted the number of unique COGs associated with each core gene and each BATS-enriched gene as a measure of gene order conservation (Fig. S3). We used the χ^2 goodness-of-fit test to evaluate whether the BATS-enriched sample came from the null distribution specified by the core genes (Poisson).

Phylogenetic Analyses. We constructed phylogenetic trees using small-insert shotgun clone sequences. Environmental gene sequences were aligned to the reference strain MIT9301 for *Prochlorococcus* and strain HTCC7211 for *Pelagibacter* using BlastAlignP (48), allowing a maximum of 70% gaps in the alignment due to the variable length and location of the environmental clone sequences. For *Prochlorococcus*, alignments were pruned to include only sequences most similar to high-light-adapted isolates, for reasons given below. Trees were constructed using PhyML (49) using the HKY85 model, with 10 different random starting trees, 4 gamma rate categories, and 0 invariant sites and parameters estimated from the data, with SH-like aLRT branch supports. Maximum likelihood is the most accurate method when using gene fragments and is the least sensitive to missing data (50). These trees were then used as input for UniFrac (32). The pruning to high-light adapted *Prochlorococcus* was done because long branches in the LL clades, which were detected only occasionally (the high-light-adapted clades eMIT9312 and eMED4 account for about 85% of cells at both sites; Table S2), can strongly influence the UniFrac test (e.g., if a LL clone was detected at HOT but not at BATS). Both the UniFrac significance test and the P test were performed for the pairwise HOT versus BATS comparison with 100 replicates.

Gene Expression. To measure gene expression during P starvation, strain MIT9301 was grown to midexponential phase in Pro99 medium, harvested at 8,000 g, washed twice with $-P$ Pro99 (no added phosphate), and resuspended in the same. Duplicate cultures, grown in continuous light at 20 $\mu\text{mol photon}\cdot\text{m}^{-2}\cdot\text{s}^{-1}$ and 21 $^{\circ}\text{C}$, were analyzed. At each time point, an aliquot of cells was harvested at 10,000 g, resuspended in storage buffer [200 mM sucrose, 10 mM sodium acetate (pH 5.2), and 5 mM EDTA], frozen in liquid N_2 , and stored at -80°C . RNA was extracted using the mirVana kit (Ambion) and Dnase-treated using Turbo DNA-free (Ambion). RNA (2–10 ng per reaction) was reverse-transcribed using SuperScript II (Invitrogen) and gene-specific primers. Transcripts of each gene were quantified using the QuantiTect SYBR Green Kit (Qiagen) and normalized to *rnpB* transcripts. The five genes tested in MIT9301 (*P9301_12441*, *12511*, *12521*, *12551*, and *12561*) were chosen because they are absent from strain MED4, for which whole genome expression data under P starvation are available (22).

ACKNOWLEDGMENTS. We thank Stephan Schuster and Ed DeLong and members of their laboratories for pyrosequencing the BATS216 and HOT186 samples, respectively; the HOT and BATS teams and the captain and crew of the R/V *Kilo Moana* and R/V *Atlantic Explorer*; Matt Sullivan, Jay McCarren, Suzanne Kern, Sarah Bagby, and Yanmei Shi for help with sample collection and processing; Scott Chilton for laboratory assistance; Daniele Veneziano for advice on statistical analyses; Kerrie Barry and the Joint Genome Institute for shotgun sequencing; and Jake Waldbauer, Vanja Klepac-Ceraj, Jesse Shapiro, and Eric Alm for comments and discussion. This work was supported in part by the Gordon and Betty Moore Foundation, the National Science Foundation (NSF), the U.S. Department of Energy, and an NSF Graduate Research Fellowship (to M.L.C.).

1. Welch RA, et al. (2002) Extensive mosaic structure revealed by the complete genome sequence of uropathogenic *Escherichia coli*. *Proc Natl Acad Sci USA* 99:17020–17024.
2. Konstantinidis KT, Tiedje JM (2005) Genomic insights that advance the species definition for prokaryotes. *Proc Natl Acad Sci USA* 102:2567–2572.
3. Hacker J, Carniel E (2001) Ecological fitness, genomic islands and bacterial pathogenicity: A Darwinian view of the evolution of microbes. *EMBO Rep* 2:376–381.
4. Simmons SL, et al. (2008) Population genomic analysis of strain variation in *Leptospirillum* group II bacteria involved in acid mine drainage formation. *PLoS Biol* 6:e177.
5. Rusch DB, et al. (2007) The Sorcerer II Global Ocean Sampling expedition: Northwest Atlantic through eastern tropical Pacific. *PLoS Biol* 5:e77.
6. Martin-Cuadrado AB, et al. (2007) Metagenomics of the deep Mediterranean, a warm bathypelagic habitat. *PLoS ONE* 2:e914.
7. Cuadros-Orellana S, et al. (2007) Genomic plasticity in prokaryotes: The case of the square haloarchaeon. *ISME J* 1:235–245.
8. Hunt DE, et al. (2008) Resource partitioning and sympatric differentiation among closely related bacterioplankton. *Science* 320:1081–1085.
9. Wildschutte H, Wolfe DM, Tamewitz A, Lawrence JG (2004) Protozoan predation, diversifying selection, and the evolution of antigenic diversity in *Salmonella*. *Proc Natl Acad Sci USA* 101:10644–10649.
10. Wilmes P, Simmons SL, Deneff VJ, Banfield JF (2009) The dynamic genetic repertoire of microbial communities. *FEMS Microbiol Rev* 33:109–132.
11. Mira A, Ochman H, Moran NA (2001) Deletional bias and the evolution of bacterial genomes. *Trends Genet* 17:589–596.
12. Kuo C-H, Ochman H (2009) The fate of new bacterial genes. *FEMS Microbiol Rev* 33:38–43.
13. Whitaker RJ, Banfield JF (2006) Population genomics in natural microbial communities. *Trends Ecol Evol* 21:508–516.
14. Morris RM, et al. (2002) SAR11 clade dominates ocean surface bacterioplankton communities. *Nature* 420:806–810.
15. Karl DM, et al. (2001) Building the long-term picture: The U.S. JGOFS time-series programs. *Oceanography* 14:6–17.
16. Wu JF, Sunda W, Boyle EA, Karl DM (2000) Phosphate depletion in the western North Atlantic Ocean. *Science* 289:759–762.

Supporting Information

Coleman and Chisholm 10.1073/pnas.1009480107

SI Text

The pyrosequencing reads have been deposited in the NCBI Sequence Read Archive, <http://www.ncbi.nlm.nih.gov>, under the following accession numbers: SRX007372 (HOT186 25m), SRX007369 (HOT186 75m), SRX007370 (HOT186 110m), SRX008032 (BATS216 20m), SRX008033 (BATS216 50m),

SRX008035 (BATS216 100m). Shotgun sequence libraries have been deposited in the NCBI Trace Archive with the following trace identifiers: 2243992048-2244135183 (BATS216 20m), 2244135184-2244269775 (BATS216 50m), 2281908320-2281966591 (HOT186 25m), 2281966591-2282003263 and 2282006336-2282042335 (HOT186 75m).

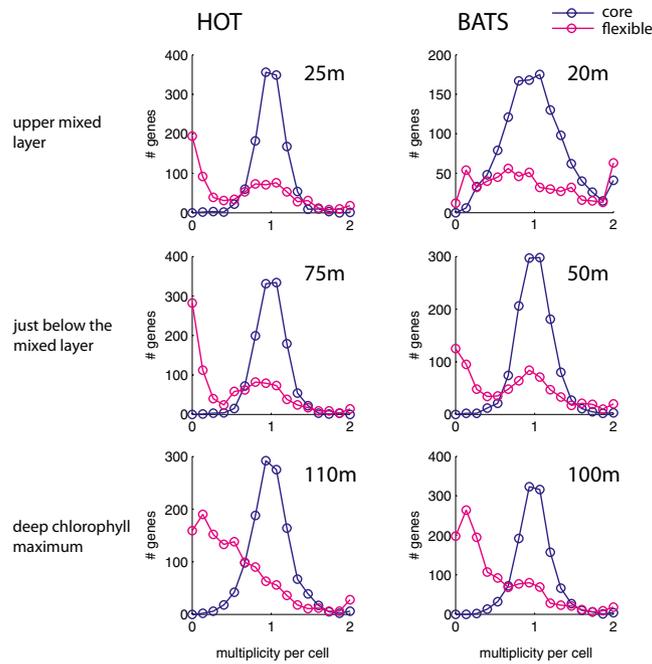


Fig. S1. Copies per *Prochlorococcus* cell of core and flexible genes at each depth. Core genes, defined using whole genome sequences from cultured isolates, also are one copy per cell in these populations, whereas many flexible genes are rare (present in only a subset of cells) at any given depth. Note that the 110-m sample from HOTS and the 20-m sample from BATS each had significantly fewer *Prochlorococcus* reads (Table S1), resulting in the unusual shape of the flexible gene distribution.

