



Metabolic Roles of Uncultivated Bacterioplankton Lineages in the Northern Gulf of Mexico “Dead Zone”

J. Cameron Thrash,^a Kiley W. Seitz,^b Brett J. Baker,^b Ben Temperton,^c Lauren E. Gillies,^d Nancy N. Rabalais,^{e,f} Bernard Henrissat,^{g,h,i} Olivia U. Mason^d

Department of Biological Sciences, Louisiana State University, Baton Rouge, Louisiana, USA^a; Department of Marine Science, Marine Science Institute, University of Texas at Austin, Port Aransas, Texas, USA^b; School of Biosciences, University of Exeter, Exeter, United Kingdom^c; Department of Earth, Ocean, and Atmospheric Science, Florida State University, Tallahassee, Florida, USA^d; Department of Oceanography and Coastal Sciences, Louisiana State University, Baton Rouge, Louisiana, USA^e; Louisiana Universities Marine Consortium, Chauvin, Louisiana, USA^f; Architecture et Fonction des Macromolécules Biologiques, CNRS, Aix-Marseille Université, 13288 Marseille, France^g; INRA, USC 1408 AFMB, F-13288 Marseille, France^h; Department of Biological Sciences, King Abdulaziz University, Jeddah, Saudi Arabiaⁱ

ABSTRACT Marine regions that have seasonal to long-term low dissolved oxygen (DO) concentrations, sometimes called “dead zones,” are increasing in number and severity around the globe with deleterious effects on ecology and economics. One of the largest of these coastal dead zones occurs on the continental shelf of the northern Gulf of Mexico (nGOM), which results from eutrophication-enhanced bacterioplankton respiration and strong seasonal stratification. Previous research in this dead zone revealed the presence of multiple cosmopolitan bacterioplankton lineages that have eluded cultivation, and thus their metabolic roles in this ecosystem remain unknown. We used a coupled shotgun metagenomic and metatranscriptomic approach to determine the metabolic potential of Marine Group II *Euryarchaeota*, SAR406, and SAR202. We recovered multiple high-quality, nearly complete genomes from all three groups as well as candidate phyla usually associated with anoxic environments—*Parcubacteria* (OD1) and *Peregrinibacteria*. Two additional groups with putative assignments to ACD39 and PAUC34f supplement the metabolic contributions by uncultivated taxa. Our results indicate active metabolism in all groups, including prevalent aerobic respiration, with concurrent expression of genes for nitrate reduction in SAR406 and SAR202, and dissimilatory nitrite reduction to ammonia and sulfur reduction by SAR406. We also report a variety of active heterotrophic carbon processing mechanisms, including degradation of complex carbohydrate compounds by SAR406, SAR202, ACD39, and PAUC34f. Together, these data help constrain the metabolic contributions from uncultivated groups in the nGOM during periods of low DO and suggest roles for these organisms in the breakdown of complex organic matter.

IMPORTANCE Dead zones receive their name primarily from the reduction of eukaryotic macrobiota (demersal fish, shrimp, etc.) that are also key coastal fisheries. Excess nutrients contributed from anthropogenic activity such as fertilizer runoff result in algal blooms and therefore ample new carbon for aerobic microbial metabolism. Combined with strong stratification, microbial respiration reduces oxygen in shelf bottom waters to levels unfit for many animals (termed hypoxia). The nGOM shelf remains one of the largest eutrophication-driven hypoxic zones in the world, yet despite its potential as a model study system, the microbial metabolisms underlying and resulting from this phenomenon—many of which occur in bacterioplankton from poorly understood lineages—have received only preliminary study. Our work details the metabolic potential and gene expression activity for uncultivated

Received 15 June 2017 Accepted 7 August 2017 Published 12 September 2017

Citation Thrash JC, Seitz KW, Baker BJ, Temperton B, Gillies LE, Rabalais NN, Henrissat B, Mason OU. 2017. Metabolic roles of uncultivated bacterioplankton lineages in the northern Gulf of Mexico “dead zone”. *mBio* 8:e01017-17. <https://doi.org/10.1128/mBio.01017-17>.

Editor Mary Ann Moran, University of Georgia

Copyright © 2017 Thrash et al. This is an open-access article distributed under the terms of the [Creative Commons Attribution 4.0 International license](https://creativecommons.org/licenses/by/4.0/).

Address correspondence to J. Cameron Thrash, thrashc@lsu.edu, or Brett J. Baker, acidophile@gmail.com.

lineages across several low DO sites in the nGOM, improving our understanding of the active biogeochemical cycling mediated by these “microbial dark matter” taxa during hypoxia.

KEYWORDS candidate phyla, hypoxia, metagenomics, microbial ecology, microbial metabolism

Hypoxia (dissolved oxygen [DO] below $2 \text{ mg} \cdot \text{liter}^{-1}$ / $\sim 62.5 \text{ } \mu\text{mol} \cdot \text{kg}^{-1}$) is dangerous or lethal to a wide variety of marine life, including organisms of economic importance (1). Hypoxia results from oxygen consumption by aerobic microbes combined with strong stratification that prevents reoxygenation of bottom waters. These taxa are fueled primarily by autochthonous organic matter generated from phytoplankton responding to nitrogen input (1). Hypoxic zones have become more widespread globally through the proliferation of nitrogen-based fertilizers and the resulting increases in transport to coastal oceans via runoff (2). In the northern Gulf of Mexico (nGOM), nitrogen runoff from the Mississippi and Atchafalaya Rivers leads to bottom water hypoxia that can extend over 20,000 km²—one of the world’s largest seasonal “dead zones” (1). Action plans to mitigate nGOM hypoxia have stressed that increasing our “understanding of nutrient cycling and transformations” remains vital for plan implementation (3). These needs motivated our current study of the engines of hypoxic zone nutrient transformation: microorganisms.

Much of our current knowledge regarding microbial contributions to regions of low DO comes from numerous studies investigating naturally occurring, deep-water oxygen minimum zones (OMZs), such as those in the Eastern Tropical North and South Pacific, the Saanich Inlet, and the Arabian, Baltic, and Black Seas (4–11). In many of these systems, continual nutrient supply generates permanent or semipermanent decreases in oxygen, sometimes to the point of complete anoxia (4). During these conditions, anaerobic metabolisms, such as nitrate and sulfate reduction and anaerobic ammonia oxidation, become prevalent (5, 9, 11–13). In contrast, nGOM hypoxia is distinguished by a seasonal pattern of formation, persistence, and dissolution (1); by benthic contributions to bottom water oxygen consumption (14, 15); and by a shallow shelf that places much of the water column within the euphotic zone (16). While parts of the nGOM hypoxic zone can become anoxic (1, 17), many areas maintain low oxygen concentrations even during peak hypoxia while the upper water column remains oxygenated (18–20).

The first studies of bacterioplankton assemblages during nGOM hypoxia showed that nitrifying *Thaumarchaea* dominated (21) and could be highly active (22), suggesting a major role for these taxa in nGOM nitrogen cycling. However, many more poorly understood organisms from cosmopolitan, but still uncultivated “microbial dark matter” (23) lineages, such as Marine Group II *Euryarchaeota* (MGII), SAR406, and SAR202, also occurred in abundance (21, 22). While the likely functions of some of these groups have become clearer recently, all of them contain multiple sublineages that may have distinct metabolic roles. For example, the SAR202 lineage of *Chloroflexi* contains at least five subclades with distinct ecological profiles (24, 25), and the best understood examples have been examined in the context of complex carbon degradation in the deep ocean (25). Likewise, SAR406 represents a distinct phylum with numerous sublineages, and the bulk of metabolic inference comes from taxa in deep-water OMZs (23, 26–28).

None of these groups have been studied in detail in shallow coastal waters, particularly in the context of seasonal hypoxia. Thus, we pursued a combined metagenomic/metatranscriptomic approach to (i) elucidate the specific contributions of these uncultivated lineages to biogeochemical cycling in the nGOM during hypoxia, (ii) evaluate the relative similarity of these organisms to their counterparts elsewhere, and (iii) determine whether other uncultivated lineages had eluded previous microbial characterization in the region due to confounding factors such as primer bias (29), 16S rRNA gene introns (30), or low abundance. Metagenomic binning recovered 20 ge-

nomes from seven uncultivated lineages, including MGII, SAR406, and SAR202, and from candidate phyla previously uncharacterized in the nGOM: *Parcubacteria* (23), *Peregrinibacteria* (31), and possibly PAUC34f (32) and ACD39 (33). Our results provide the first information on the likely potential function and activity of these taxa during hypoxia in the shallow nGOM, and suggest novel roles for some of these groups that possibly reflect sublineage-specific adaptations.

RESULTS

Study area. Our previous work used 16S rRNA gene amplicon data and quantitative PCR (qPCR) to examine correlations between whole microbial communities, nutrients, and DO across the geographic range of the 2013 seasonal hypoxia (21). Here we selected six of those samples from offshore of the region between Atchafalaya Bay and Terrebonne Bay (D', D, and E transects). These sites ranged considerably in DO concentration (~ 2.2 to $132 \mu\text{mol} \cdot \text{kg}^{-1}$), and we chose them to facilitate a detailed investigation of the metabolic repertoire of individual taxa across the span of suboxic (1 to $20 \mu\text{mol} \cdot \text{kg}^{-1}$ DO) to oxic ($>90 \mu\text{mol} \cdot \text{kg}^{-1}$ DO) (5) water. Microbial samples from these sites were collected at the oxygen minimum near the bottom. Site depth ranged from 8 to 30 m, with the hypoxic ($<2 \text{ mg} \cdot \text{liter}^{-1}/62.5 \mu\text{mol} \cdot \text{kg}^{-1}$) layer (at sites D2, D3, E2A, and E4) extending up to ~ 5 m off the bottom (see Table S1 at <http://thethrashlab.com/publications>).

Metagenomic assembly yielded high-quality genomes from multiple uncultivated lineages. Our initial assembly and binning efforts recovered 76 genomes. Using a concatenated ribosomal protein tree that included members of the candidate phylum radiation (CPR) (34) (see Fig. S1 in the supplemental material), CheckM (35) (Fig. S2), 16S rRNA genes and other single-copy markers where available, and analyses of individual gene taxonomy (Fig. S3), we assigned 20 genomes to uncultivated “microbial dark matter” groups. These were six Marine Group II *Euryarchaeota* (MGII), five *Marinimicrobia* (SAR406), three in the SAR202 clade of *Chloroflexi*, and within candidate phyla (CP), one *Parcubacteria* (OD1), two *Peregrinibacteria*, and putatively, one ACD39 and two PAUC34f (Table 1 and Text S1). We further defined the MGII, SAR406, and SAR202 genomes into sublineages based on average amino acid identity (AAI), GC content, clade structure in the ribosomal protein tree, and 16S rRNA genes (Text S1). SAR406 genomes belonged to two groups, groups A and B, corresponding to the previously established Arctic96B-7 and SHBH1141 16S rRNA gene clades (27). The three SAR202 genomes belonged to the previously established subclade I 16S rRNA gene clade (24). All genomes, with the exception of *Parcubacteria* Bin 40, had estimated contamination of less than 6%, and in the majority of cases, less than 2%. Four of the six MGII genomes had estimated completeness (via CheckM) of greater than 61%, four of the five SAR406 genomes had greater than 73%, and all three SAR202 genomes were estimated to be greater than 83% complete. All CP lineages had at least one genome estimated to be greater than 71% complete (Table 1).

Unique roles for the ubiquitous MGII, SAR406, and SAR202 lineages in nGOM hypoxia. MGII comprised more than 10% of the total community in some samples from 2013, and one MGII operational taxonomic unit (OTU) also had a strong negative correlation with DO during 2013 hypoxia (21). Within our metagenomic data set, MGII were more abundant in lower oxygen samples than in fully oxic samples and the most abundant of the lineages reported here (Fig. S7). The majority of genomes encoded for aerobic, chemoheterotrophic metabolism, with no predicted genes for nitrogen or sulfur respiration except for a putative nitrite reductase (*nirK*) in a single genome, Bin 15 (Fig. 1; see Table S1 at <http://thethrashlab.com/publications>). MGII genomic abundance correlated well with transcriptional abundance in most samples (Fig. 2), and we specifically found MGII cytochrome *c* oxidase expression throughout, though the levels and patterns differed depending on the gene and the source genome (Fig. 3; see Table S1 at <http://thethrashlab.com/publications>). Expression of the *nirK* gene occurred in the D2 and E2A samples—both suboxic. All but the most incomplete genome encoded for ammonia assimilation, making this a likely nitrogen source. Aggregate

TABLE 1 Genome characteristics for the 20 bins associated with uncultivated lineages

IMG genome ID ^a	Bin ID	Taxonomy	Compl. (%) ^b	Contam. (%) ^c	Strain het. ^d	No. of scaff. ^e	Longest scaff. (bp)	Size (bp)	No. of genes	GC content (fract.) ^f	Coding density (fract.)	Estim. compl. size (Mbp) ^g
2651870035	43-1	<i>Chloroflexi</i> (SAR202)	90.3	0	0	69	161,242	1,972,793	1,882	0.52	0.88	2.2
2651870036	43-2	<i>Chloroflexi</i> (SAR202)	88.6	4.1	15.4	134	157,345	2,402,386	2,373	0.52	0.91	2.7
2651870034	43	<i>Chloroflexi</i> (SAR202)	83.2	0.1	100	179	90,503	2,475,308	2,392	0.53	0.89	3.0
2693429801	45	<i>Marinimicrobia</i> (SAR406) group B	89.8	5.5	71.4	124	105,140	2,811,623	2,487	0.47	0.93	3.1
2651870052	45-1	<i>Marinimicrobia</i> (SAR406) group B	85.2	0.1	100	144	100,582	2,410,233	2,235	0.49	0.93	2.8
2693429802	45-2	<i>Marinimicrobia</i> (SAR406) group B	79.3	1.8	12.5	287	61,408	2,811,444	2,554	0.46	0.94	3.5
2651870053	51-1	<i>Marinimicrobia</i> (SAR406) group A	73.6	1.7	50	75	191,525	1,901,306	1,835	0.39	0.95	2.6
2651870051	51	<i>Marinimicrobia</i> (SAR406) group A	21.3	0	0	115	16,923	578,802	686	0.41	0.95	2.7
2651870038	15	<i>Euryarchaeota</i> (MGII)	83.2	1.6	0	43	255,599	1,885,130	1,614	0.62	0.95	2.3
2651870039	17	<i>Euryarchaeota</i> (MGII)	81.9	0.1	50	88	137,319	1,803,861	1,564	0.43	0.96	2.2
2651870037	14	<i>Euryarchaeota</i> (MGII)	71.2	0.8	100	80	90,069	1,389,909	1,305	0.54	0.94	2.0
2651870040	18	<i>Euryarchaeota</i> (MGII)	61.1	0.8	100	155	20,633	1,033,226	1,025	0.50	0.96	1.7
2651870042	38	<i>Euryarchaeota</i> (MGII)	27.1	0	0	138	12,893	615,290	610	0.55	0.96	2.3
2651870041	17-1	<i>Euryarchaeota</i> (MGII)	16.6	2.8	33.3	123	16,952	538,052	566	0.41	0.95	3.2
2693429807	13	Unclassified (ACD39)	89.8	5.1	20	401	75,104	4,269,849	3,686	0.47	0.93	4.8
2693429799	40	<i>Parcubacteria</i> (OD1)	71.9	15.1	75	97	65,104	1,086,283	1,208	0.52	0.92	1.5
2693429797	16	<i>Peregrinibacteria</i>	83.2	0.3	100	67	59,308	1,384,712	1,318	0.39	0.93	1.7
2693429798	39	<i>Peregrinibacteria</i>	49.9	0.3	100	131	18,806	747,520	809	0.45	0.95	1.5
2693429804	50	Unclassified (PAUC34f)	84.8	1.4	0	455	76,269	5,346,994	5,484	0.58	0.92	6.3
2693429803	48	Unclassified (PAUC34f)	51.9	2.2	0	470	20,176	2,566,149	2,596	0.55	0.94	4.9

^aID, identifier.^bCompl., estimated completeness.^cContam., estimated contamination.^dStrain het., strain heterogeneity.^escaff., scaffolds.^ffract., fraction.^gEstim. compl. size, estimated complete genome size.

metabolic construction from multiple bins also indicated a complete tricarboxylic acid (TCA) cycle, glycolysis via the pentose phosphate pathway, and gluconeogenesis (Fig. 1). Carbohydrate active enzyme (CAZy) genes can provide critical information on the relationships between microbes and possible carbon sources (36). We found few CAZy genes, and these were largely restricted to glycosyltransferases (GT) in families 2 and 4, with activities related to cellular synthesis. In general, CAZy expression occurred for at least one gene in every genome, and we detected expression of GT cellular synthesis genes in the E2A sample (Fig. 4), likely indicating actively growing cells.

SAR406 represented more than 5% of the population in some locations during hypoxia in 2013, and one abundant OTU was negatively correlated with DO (21). Metagenomic read recruitment to the SAR406 bins confirmed this trend, with greater recruitment in the suboxic samples relative to dysoxic or oxic samples (Fig. S7). Total RNA recruitment was strongest to Bins 45 and 51-1, though most bins showed an RNA-to-DNA recruitment ratio of >1 in at least one sample, indicating that these taxa were likely active (Fig. 2). Despite their affinity for low-oxygen environments, the SAR406 genomes encoded a predicted capacity for aerobic respiration (Fig. 1), and we found expression of cytochrome *c* oxidases in even the lowest oxygen samples (Fig. 3). The group B genomes encoded both high- and low-affinity cytochrome *c* oxidases (37), whereas the high-affinity (*ccb*₃-type) oxidases were not recovered in the group A genomes (see Table S1 at <http://thethrashlab.com/publications>), which may indicate sublineage-specific optimization for different oxygen regimes.

Sublineage variation also appeared in genes for the nitrogen and sulfur cycles. Group B genomes all contained predicted nitrous oxide reductases (*nosZ*) and *nrfAH* genes for dissimilatory nitrite reduction to ammonium (defined here as DNRA, although

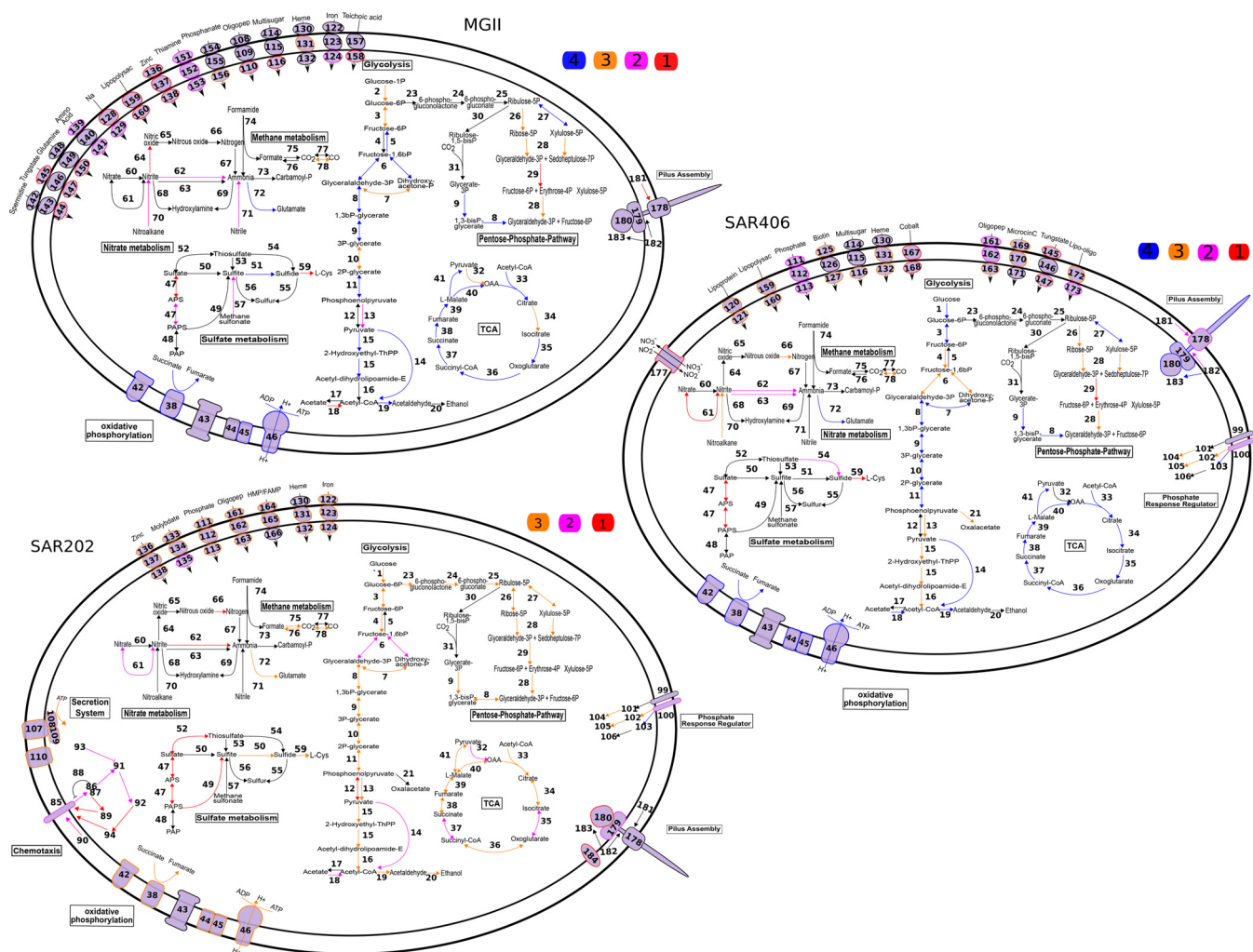


FIG 1 Metabolic reconstruction of Marine Group II *Euryarchaeota*, SAR406, and SAR202, based on the top three or four most complete genomes. Colors indicate pathway elements based on the number of genomes in which they were recovered, according to the key. Black outlines and/or arrows indicate genes that were not observed. Boldface black numbers correspond to annotations supplied in Table S1 at <http://thetrashlab.com/publications>.

this acronym frequently refers to nitrate, even though that is a misnomer [38]). The *nrfA* genes formed a monophyletic group with *Anaeromyxobacter dehalogenans* 2CP-1, an organism with demonstrated DNRA activity (38) (Fig. S8A). The genes also contained conserved motifs diagnostic of the *nrfA* gene (38) (Fig. S8B and C). We observed expression of *nrfAH* and *nosZ* at the sites with the lowest DO concentrations (D2, E2A, and E4), and expression appeared to have a negative relationship with DO concentration (Fig. 3). The Bin 51-1 group A genome contained predicted *narHI* genes for dissimilatory nitrate reduction, which we did not find in the group B genomes. We observed expression of SAR406 *narHI* only in the lowest DO sample from station E2A (Fig. 3). Two group B SAR406 genomes had predicted *phsA* genes for thiosulfate reduction to sulfide (and/or polysulfide reduction [39]), as previously described from fosmid sequences (27). We detected transcripts for these genes only in samples E2A and E4, the two lowest DO samples (Fig. 3). Many of the anaerobic respiratory genes were coexpressed with cytochrome *c* oxidases, indicating a potential for either coreduction of these alternative terminal electron acceptors or poisoning of these organisms for rapid switching between aerobic and anaerobic metabolism (40).

All SAR406 genomes had numerous genes for heterotrophy. We found CAZy genes in all major categories except polysaccharide lyases, and expression for most of these genes in both group A and group B genomes in one or more samples (Fig. 4). Notable

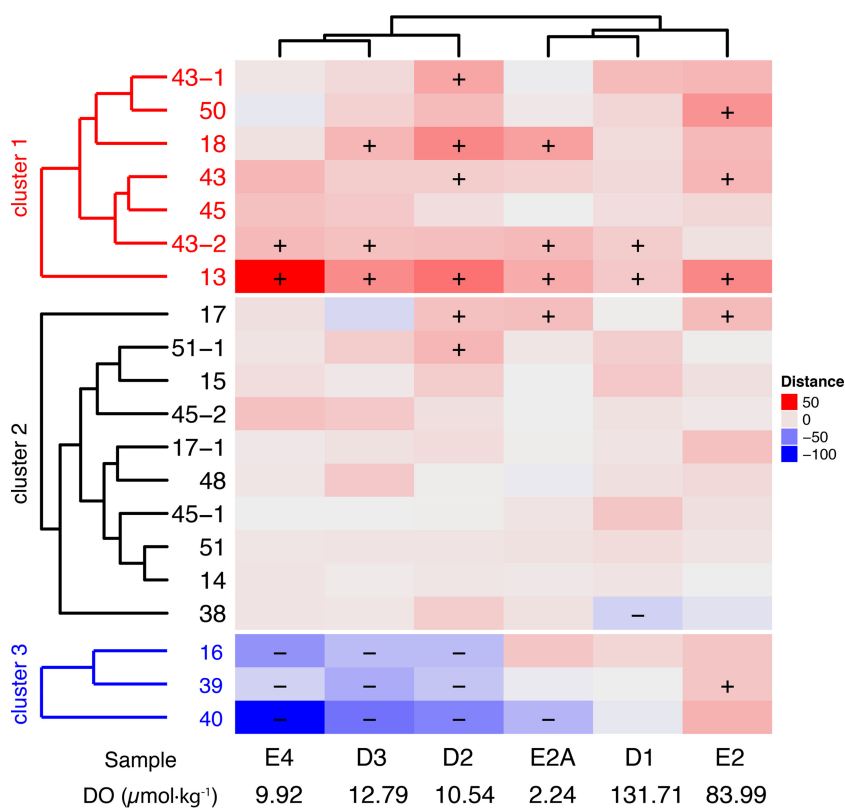


FIG 2 Relative DNA to RNA recruitment rank for each genome, by sample. Colors indicate the relative difference in the ratio of rank based on total RNA and DNA mapping. Red indicates a higher RNA recruitment rank compared to DNA recruitment rank, and vice versa for blue. The plus and minus symbols indicate bins where the rank residual from the identity in RNA versus DNA read mapping was more or less than 1 standard deviation beyond 0, respectively. Dendrograms were calculated using an unweighted-pair group method with arithmetic mean (UPGMA) from Euclidean distances of rank residuals across all samples and bins.

carbohydrate compounds for which degradation capacity was predicted include cellulose (glycoside hydrolase [GH] families GH3 and GH5; carbohydrate binding module [CBM] family CBM6), starch (GH13), agar and other sulfated galactans (GH2 and GH16), chitin (GH18), xylan (GH30 and CBM9), and peptidoglycan (GH23, GH103, and CBM50). The genomes contained putative transporters for a variety of dissolved organic matter (DOM) components including nucleosides, amino and fatty acids, and oligopeptides (see Table S1 at <http://thethrashlab.com/publications>). We also found numerous outer membrane transporters (including cation symporters), outer membrane receptors (OMR) (TonB dependent) (which play important roles in transport of metals, vitamins, colicins, and other compounds), and outer membrane factors (OMF). Most genomes also had large numbers of duplicated genes (24 in Bin 45-2), identified via hidden Markov model searches against the Sifted Families (SFam) database (41), annotated as "Por secretion system C-terminal sorting domain-containing protein," some of which were associated with GH16. These genes likely play a role in sorting C-terminal tags of proteins targeted for secretion via the Por system, which is essential for gliding motility and chitinase secretion in some *Bacteroidetes* (42). The extensive gene duplication may indicate expanded and/or specialized sorting functionality and suggests an emphasis on protein secretion in this group. Expression of a membrane-bound lytic murein transglycosylase D (GH23) involved in membrane remodeling also supports the idea of active and growing cells from group A in all samples (see Table S1 at <http://thethrashlab.com/publications>).

We detected *Chloroflexi* 16S rRNA gene sequences during 2013 hypoxia at up to 5% of the community (21) and recovered three mostly complete SAR202 *Chloroflexi* genomes in this work. Although present at a lower abundance than MGII and SAR406

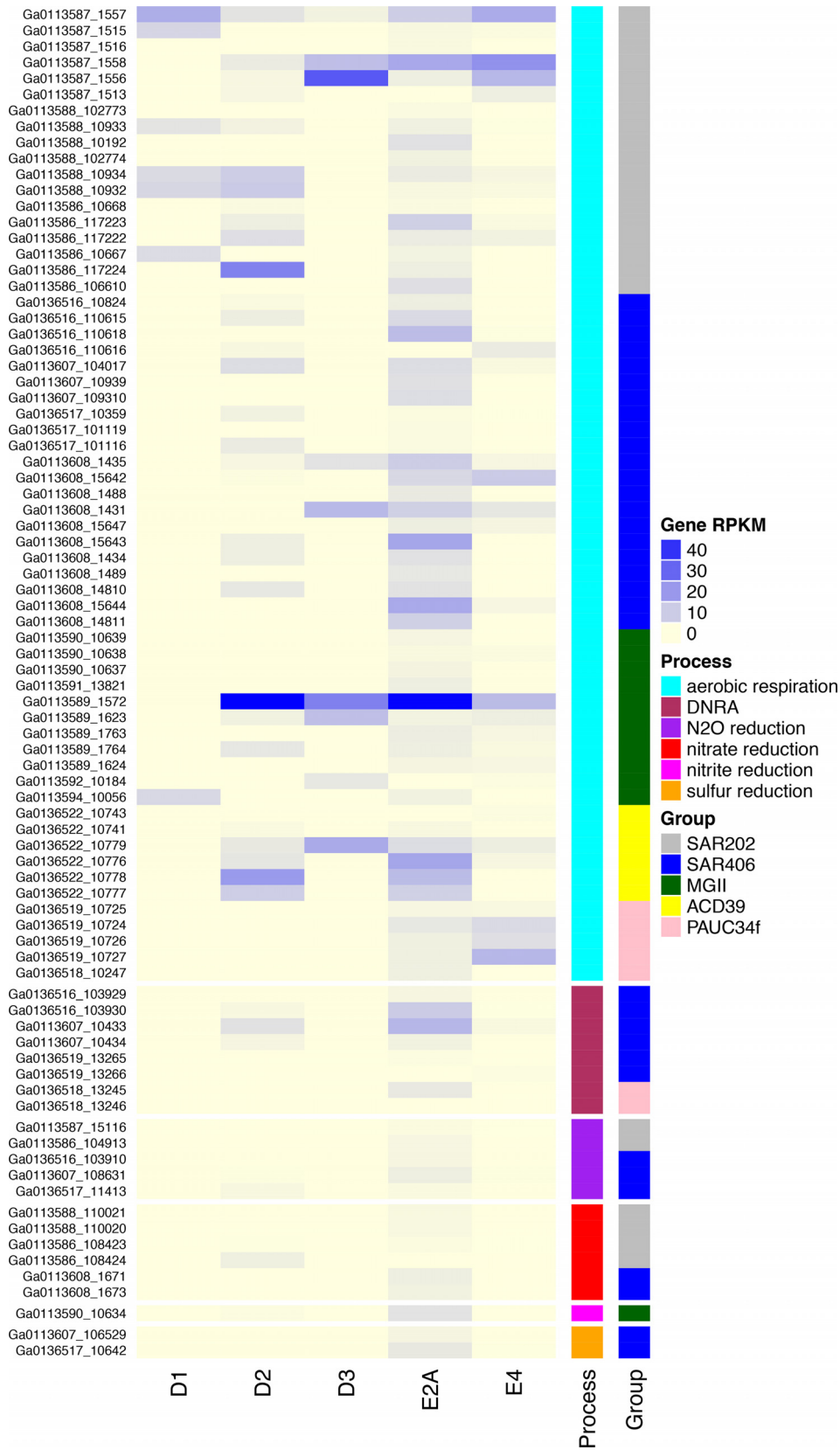


FIG 3 Expression of predicted respiratory genes. RPKM values of RNA recruitment for each gene, by sample, are depicted with colors according to the Gene RPKM key (pale yellow to blue shows increasing intensity). Genes are grouped by bin, taxonomic affiliation, and specific respiratory process. DNRA, dissimilatory nitrite reduction to ammonia.

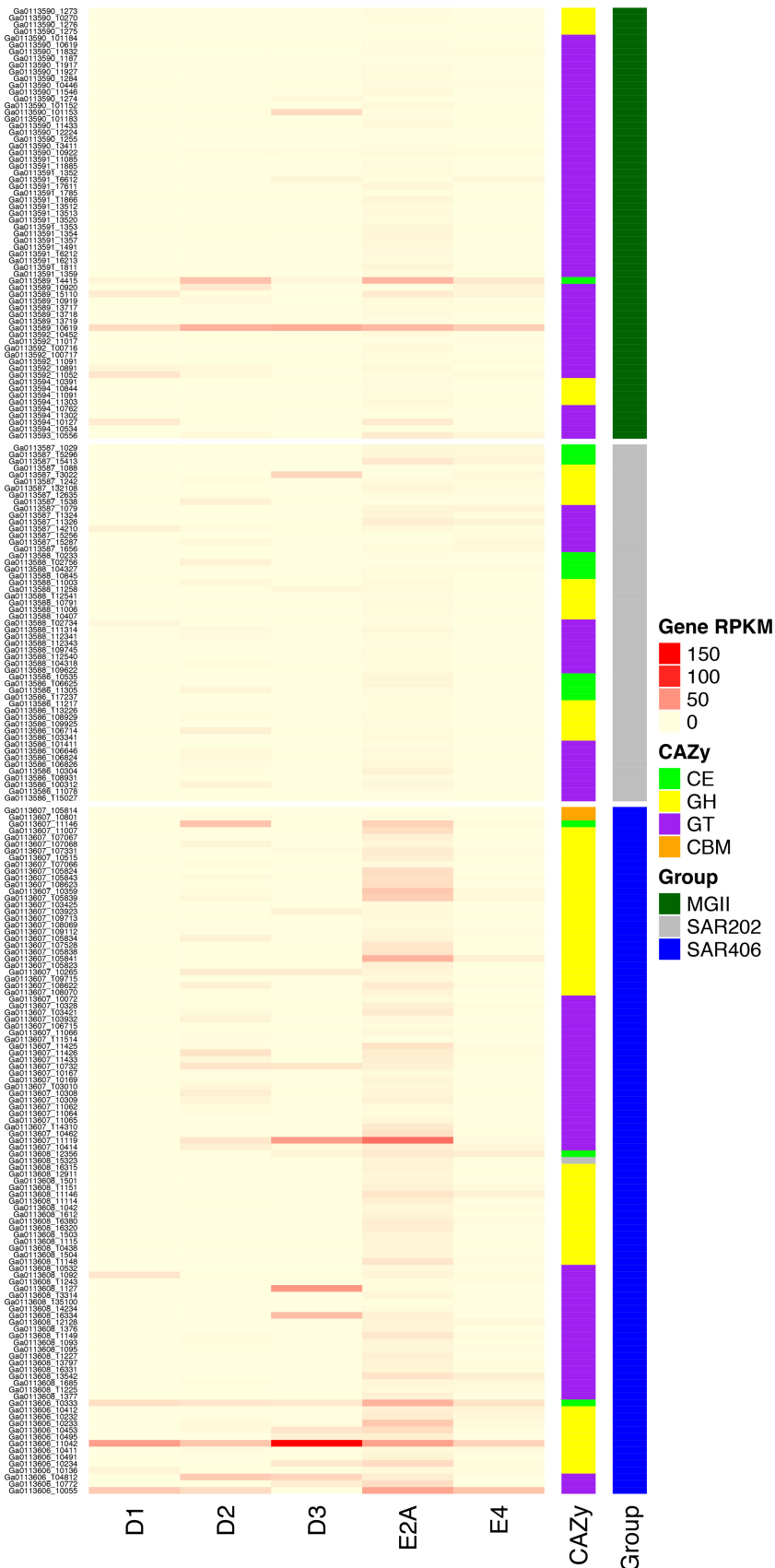


FIG 4 Expression of predicted CAZy genes. RPKM values or RNA recruitment for each gene, by sample, are depicted with colors according to the Gene RPKM key (pale yellow to red shows increasing intensity).

(Continued on next page)

(Fig. S7), these genomes showed relatively high activity in some samples (Fig. 2). Like subclade III and V, subclade I organisms likely respire oxygen. However, we also found *napAB* and *nosZ* genes for nitrate and nitrous oxide reduction, respectively (Fig. 1). As in SAR406, we detected concurrent expression of these genes with cytochrome *c* oxidases in the lowest DO samples (Fig. 3) (see Table S1 at <http://thethrashlab.com/publications>).

The SAR202 genomes have numerous transporters, many with predicted roles in organic matter transport, which supports previous observations of DOM uptake (43). In particular, SAR202 genomes had considerably more major facilitator superfamily (MFS) transporters than the other genomes in this study (see Table S1 at <http://thethrashlab.com/publications>) and those of the subclade III genomes (25), and SFam searches revealed that the majority of these shared annotation as a “Predicted arabinose efflux permease” (SFam 346742). MFS genes transport numerous diverse substrates, such as sugars and amino acids, through coupling with an ion gradient, and can be associated with either uptake or export of compounds (44). SAR202 genomes also had between 53 and 66 predicted ABC transporters.

The SAR202 genomes carried a number of duplicated genes in specific gene families. The largest gene family expansion that we observed was associated with SFam 6706, with between 46 and 48 genes in this family in each genome. Most of these genes (121/142) were annotated as either a “galactonate dehydratase” or a “L-alanine-DL-glutamate epimerase.” Galactonate dehydratase catalyzes the first step of the pathway to utilize D-galactonate in central carbon metabolism via the pentose phosphate pathway. The large number of genes in these categories likely indicates some divergence for alternative roles, as this group belongs broadly to the COG4948 “L-alanine-DL-glutamate epimerase or related enzyme of enolase superfamily.” All genomes also had numerous dehydrogenase genes as reported for the subclade III genomes (25). Specifically, SFams 346640 and 1639 were the third and fourth most abundant, with 16 to 18 and 13 to 15 genes in each family, respectively, in the three genomes. Genes in these families were annotated as “short-chain alcohol dehydrogenase family,” “3-alpha (or 20-beta)-hydroxysteroid dehydrogenase,” “meso-butanediol dehydrogenase,” and others. These match the annotations of the subclade III genomes and suggest a similar role in conversion of alcohols to ketones (25). The SAR202 genomes have comparatively few CAZy genes relative to the other genomes. GH15 and GH63 suggest starch degradation and GH105 pectin degradation, and we detected expression of multiple genes in these categories across samples (Fig. 4; see Table S1 at <http://thethrashlab.com/publications>).

Other candidate phylum organisms in nGOM hypoxia. In contrast to the abundant and cosmopolitan MGII, SAR406, and SAR202 clades, we also recovered genomes from several groups that were either previously undetected in the nGOM or very rare. Although these taxa likely do not contribute the biomass of more populous clades, their genomes provide important insight into their functional potential during hypoxia. The Bin 13 genome (possibly ACD39) also had the highest relative activity compared to all the other genomes in our study (Fig. 2), underlining the point that low abundance does not automatically equate to low metabolic impact. Bin 13 had predicted aerobic respiration with both high- and low-affinity cytochrome *c* oxidases (Fig. 5). The low-affinity oxidases contributed more reads in the samples where we could detect expression (see Table S1 at <http://thethrashlab.com/publications>). The genome contained numerous predicted CAZy genes in the glycosyltransferase and glycoside hydrolase categories, spread across multiple families in each (see Table S1 at <http://thethrashlab.com/publications>). Notable degradation capacity included starch (GH13) and peptidoglycan (GH23, GH103, and GH104).

FIG 4 Legend (Continued)

Genes are grouped by bin, taxonomic affiliation, and general CAZy categories. CE, carbohydrate esterase; GH, glycoside hydrolase; GT, glycosyltransferase; CBM, carbohydrate binding module.

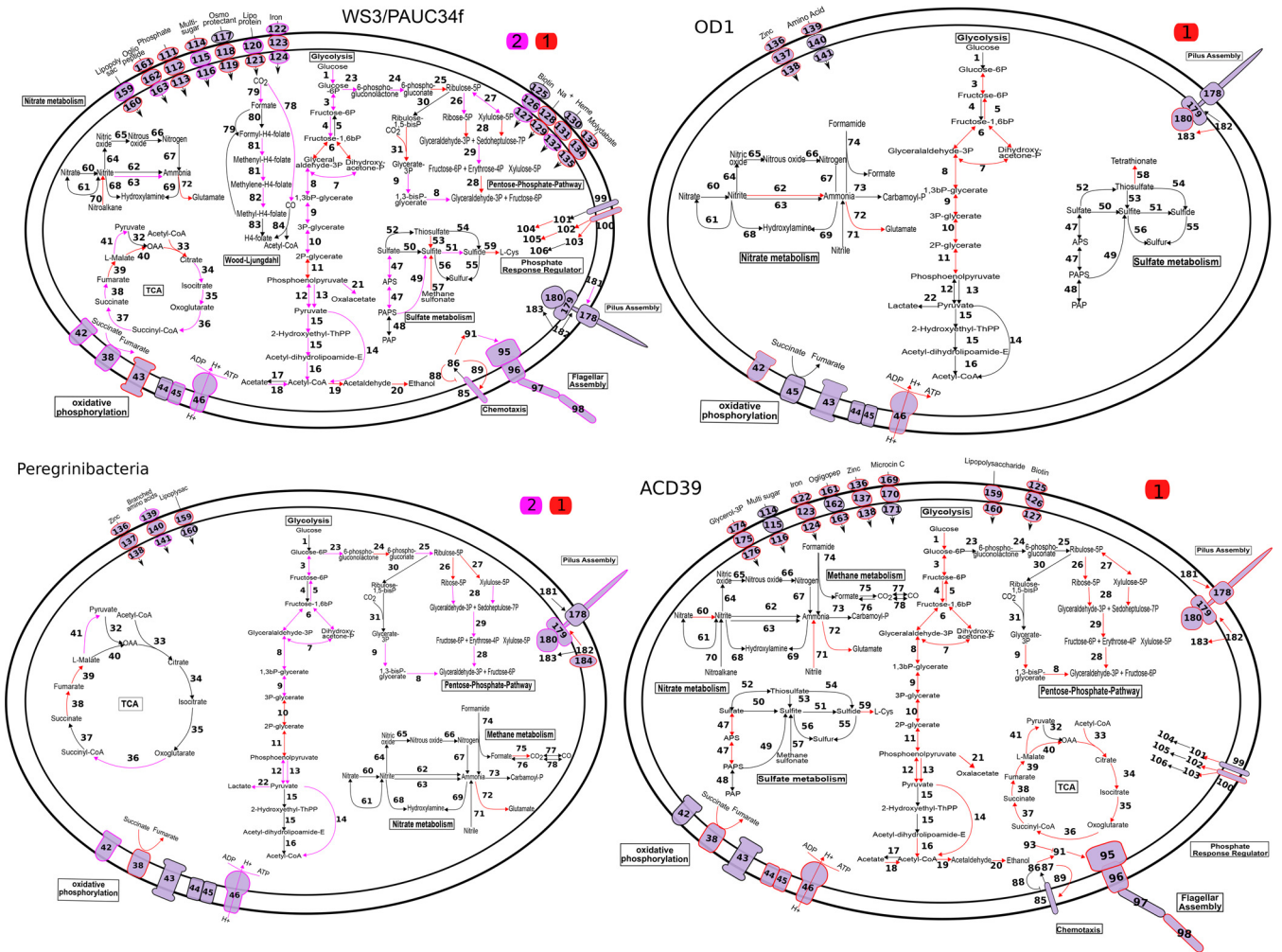


FIG 5 Metabolic reconstruction of the candidate phylum members PAUC34f, *Parcubacteria* (OD1), *Peregrinibacteria*, and ACD39 (Bin 13). Colors indicate pathway elements based on the number of genomes in which they were recovered, according to the key. Black outlines and/or arrows indicate genes that were not observed. Boldface black numbers correspond to annotations supplied in Table S1 at <http://thethrashlab.com/publications>.

Bin 13 had ~80 ABC transporter genes, and similar to the SAR406 genomes, numerous outer membrane transporters, including the OMR and OMF families. We predict complete glycolysis/gluconeogenesis pathways and a TCA cycle. We recovered paralogous pilus subunit genes, chemotaxis genes, and a partial flagellar assembly. Furthermore, we detected relatively high expression of the *flgLN* flagellin genes in samples D2, D3, E2A, and E4 (see Table S1 at <http://thethrashlab.com/publications>), suggesting active motility in these environments. Several other Bin 13 genes were among the most highly expressed in all samples but could be classified only as hypothetical (see Table S1 at <http://thethrashlab.com/publications>). Similarly, the three most populous SFams in Bin 13, according to the numbers of genes ($n = 16, 15,$ and 13) also linked to genes annotated as hypothetical proteins with either tetratricopeptide, HEAT, TPR, or Sel1 repeats. Although currently obscure, these and the highly expressed hypothetical genes represent important targets for future research into the function of this group.

Bins 50 and 48 were lower in abundance than SAR202 genomes (Fig. S7, PAUC34f), with no observable trend associated with oxygen levels (Fig. S7). These genomes encoded flagellar motility, aerobic respiration, glycolysis via the pentose-phosphate pathway, gluconeogenesis, assimilatory sulfate reduction, and DNRA (Fig. 5). The *nrfA* subunit from both genomes grouped in the same monophyletic clade as those from SAR406 (Fig. S8A) and had similar conserved motifs (Fig. S8B and C). However, we note

that the *nrfAH* gene sets for Bins 50 and 48 occurred on relatively short contigs (5,650 and 5,890 bp, respectively), so the metabolic assignment cannot be corroborated as definitively as that for SAR406. The Bin 50 genome was among the more active in our analysis (Fig. 2), and we detected the highest expression of cytochrome *c* oxidase components in samples E2A and E4 (Fig. 3). DNRA gene expression was low but observable in the same samples. We also recovered a partial gene for the ribulose-bisphosphate carboxylase (RuBisCO) large subunit, but this fragment was on a very short contig (3,954 bp), and we did not detect expression in any of our samples, so we cannot rule out that this gene occurred on a contaminating contig.

The Bin 50 and 48 genomes had abundant CAZy genes in all categories, suggesting a highly flexible metabolic repertoire for carbon acquisition. They contain possible capacity for breakdown of starch (GH13 and CBM48), peptidoglycan (GH23 and CBM50), fructose-based oligosaccharides (GH32), and hemicellulose (GH2, GH3, and GH43). Notably, these genomes were the only ones with predicted polysaccharide lyases (PL) among those compared (with the exception of a single predicted PL gene in SAR406 [see Table S1 at <http://thethrashlab.com/publications>]). PL genes cleave uronic-acid containing polysaccharides (45). These organisms seem particularly adapted for pectin (PL1, PL2, PL9, PL10, PL11, PL22, and GH78) and alginate (PL15 and PL17) degradation—both compounds are common cell wall components of green and brown algae, respectively.

In line with the algal cell wall degradation ability, we detected a large expansion (102 genes in Bin 50) of sulfatase genes in SFam 1534, annotated predominantly as either “arylsulfatase A” or “choline-sulfatase.” Arylsulfatases cleave sulfate esters, usually to supply microbes with a source of sulfur, and can be located intracellularly or in membranes (46). Choline sulfatases cleave choline sulfate to choline and sulfate, with downstream use for the former as a carbon source or osmoprotectant and the latter as a sulfur source (47). Given the predicted assimilatory sulfate reduction pathway in Bins 50 and 48, this is a logical means to obtain sulfur for the group. We observed large expansions in galactonate and other dehydratases (as in SAR202, above; SFam 6706, 42 genes in Bin 50), as well as numerous ABC transporter permeases (SFam 4442), which match the transporter predictions via IMG (Integrated Microbial Genomes): 117 predicted genes for ABC transporters in all. These genomes also had numerous OMF and OMR transporter genes (see Table S1 at <http://thethrashlab.com/publications>). The large number of transporters and protein family expansions correspond to the relatively large expected genome sizes (between 5 and 6 Mbp).

We also recovered genomes associated with CPR taxa usually associated with anoxic environments: two *Peregrinibacteria* and one from the *Urbacteria* subclade of the *Parcubacteria* (formerly OD1). All three genomes could be assigned taxonomically with high confidence based on their positions in the ribosomal protein tree (Fig. S1) and via gene annotations (Fig. S3). We note that although the *Peregrinibacteria* bins (16 and 39) had very low predicted contamination, the *Parcubacteria* Bin 40 had 15% predicted contamination (75% of which we attributed to strain heterogeneity in the bin) (Table 1). Recovery of *Parcubacteria* from a coastal marine system is unusual, but not unprecedented. *Parcubacteria* single-cell genomes have been identified in marine and brackish sources (23), and we previously identified 26 rare OTUs assigned to the phylum in nGOM hypoxia (21). That number of OTUs may explain why we observed 20 single-copy marker genes present in two copies in Bin 40 (see Table S1 at <http://thethrashlab.com/publications>).

In contrast to *Parcubacteria*, *Peregrinibacteria* have thus far been found only in terrestrial subsurface aquifers (31, 33, 48, 49) and remained undetected in our amplicon survey (21). Both groups occurred in low relative abundance compared to the other taxa in this study (Fig. S7) and showed the lowest activity (Fig. 2). Consistent with previous reports of obligate fermentative metabolism by *Parcubacteria* and *Peregrinibacteria* (23, 30, 31, 48), we identified no respiratory pathways for these taxa (Fig. 5), and they trended toward greater abundances in the lowest DO samples (Fig. S7). In spite of relatively high predicted genome completion, we found very few CAZy genes, and

those genes were mostly restricted to glycosyltransferases (see Table S1 at <http://thethrashlab.com/publications>) probably involved in capsular polysaccharide synthesis. While these organisms had low relative abundance to the other groups (Fig. S7), we did observe activity in some samples (Fig. 2, samples E2, E2A, and D1).

DISCUSSION

This work provides the first reconstruction of multiple nearly complete genomes from uncultivated bacterioplankton during nGOM hypoxia. Although we define roles for MGII, SAR406, SAR202, Bin 13, and Bins 50/48 as aerobic heterotrophs, we also observed concurrent expression of genes associated with anaerobic metabolism in SAR406 (nitrate reduction, DNRA, nitrous oxide reduction, and sulfur reduction), SAR202 (nitrate and nitrous oxide reduction), MGII (nitrite reduction), and Bins 50/48 (DNRA) in suboxic samples with the lowest measured DO concentrations. Simultaneous utilization of multiple electron acceptors with different redox potentials likely indicates an abundant supply of electron donors (50), may denote niche partitioning within group sublineages at a finer level of taxonomic resolution than we observed, or indicates poisoning of taxa for rapidly changing chemical gradients (40). An organism's set of CAZy genes often gives insights into its biology, in particular into nutrient sensing and acquisition. All taxa examined in this study had predicted chemoorganoheterotrophic metabolism, and the CAZy genes found in these genomes suggest that SAR406, SAR202, Bin 13, and Bins 50/48 participate in the degradation of complex organic matter resulting from the detritus of larger organisms. This matches the general model of hypoxic zone oxygen consumption resulting from sinking organic matter provided by algal blooms in surface waters (1). The observed activity of obligate fermentative groups *Parcubacteria* and *Peregrinibacteria* also suggests that anoxic pockets occur in the water column where these organisms can thrive.

Marine Group II (MGII) is a broadly distributed archaeal clade, with members found in different marine (51, 52) and sedimentary (53) environments. Previous work during 2012 and 2013 hypoxia indicated a proliferation of archaeal taxa in both the *Thaumarchaea* and MGII phyla (21, 22). The prevalence of MGII among lower oxygen samples in the hypoxic zone is somewhat surprising, considering that they are commonly associated with aerobic environments (52). However, oxygen was still present in even the lowest DO samples (Fig. 2), and MGII success likely had more to do with the carbon content than oxygen levels. These nGOM MGII appear to be metabolically similar to those described in previous work: MGII have been shown to be dominant in water column environments associated with blooms in productivity, for example at deep-sea hydrothermal plumes (51). Thus, the increased availability of organic matter (proteins and carbohydrates), thought to be preferred substrates for MGII (54, 55), probably explains their abundance.

Another cosmopolitan group found in our samples was SAR406 or Marine Group A. These organisms were discovered more than 20 years ago (28, 56), and the clade has recently been proposed as the phylum "*Marinimicrobia*" (23). SAR406 occur in numerous marine (5, 23, 26, 28, 57), sedimentary (23), and even oil reservoir (58) environments. They are prevalent in deeper ocean waters (28, 57, 59) and prefer lower oxygen concentrations in OMZs (5, 26, 60). Our genomes had larger estimated genome sizes—2.6 to 2.7 Mbp (group A) and 2.8 to 3.5 Mbp (group B)—compared to 1.1 to 2.4 Mbp from single-cell genomes (23). Overall GC content, however, was in the range of the 30 to 48% reported for fosmids (27) and single-cell genomes (23). The lower-GC group A genomes specifically had GC contents similar to that of the Arctic96B-7 fosmids, matching their predicted phylogenetic affiliation (see below) (27).

Our data now also define roles for SAR406 in the eutrophication-driven hypoxia of the nGOM. Previous metabolic reconstructions of SAR406 predicted aerobic metabolism (23) and sulfur reduction (27), which our data confirm, although the sulfur reduction genes were found only in group B organisms (see Table S1 at <http://thethrashlab.com/publications>). Our genomes also suggest multiple nitrogen cycling roles that appear to be organized by sublineages within the phylum, and sublineage-

specific presence of both high- and low-affinity cytochrome *c* oxidases. The group B organisms group with the early diverging SBH1141 clade (27) for which no previous genome data exist. Group B organisms contained both types of cytochrome *c* oxidase genes, *nosZ* and *nrfAH* genes, whereas group A organisms, sister to the Arctic96B-7 clade, contained only the low-affinity cytochrome *c* oxidases, and additionally *narHI* genes not found in group B. The unique roles predicted for these taxa are not surprising given the diversity of the SAR406 clade and the genetic distances between group A and B (see Fig. S4 in the supplemental material). The fosmids associated with the Arctic96B-7 clade contained genes for oxidative stress and sulfur reduction (27), although we only found sulfur reduction genes in the distantly related group B genomes. The Arctic96B-7 clade may be diverse enough to encompass differing metabolic strategies, but the variable presence of *phsA* genes in this group may simply be due to incomplete genomic data. In addition to sublineage-specific respiratory characteristics, our results also generate specific hypotheses about organic matter metabolism in SAR406: likely degradation capacity for cellulose, starch, agar, xylan, and peptidoglycan; transport of nucleosides, amino and fatty acids, and oligopeptides; and substantial gene duplication associated with protein secretion for possible extracellular metabolism. Together, these data suggest that during nGOM hypoxia, SAR406 members degrade complex carbohydrates fueled by aerobic respiration and supplemented with facultative anaerobic respiration of nitrate, nitrite, or sulfur compounds.

Members of the SAR202 clade of *Chloroflexi* also inhabit a wide variety of marine environments (24), frequently in deeper waters (24, 43, 57, 59, 61) and remain functionally understudied because genome data for SAR202 have been lacking. Landry and colleagues recently described the properties for several single-cell genomes representing SAR202 subclades III and V recovered from the mesopelagic zone (25). Our genomes have generally higher GC content and much lower expected genome sizes than those predicted by Landry et al. (25), although these calculations are likely complicated by the relative incompleteness of their genomes (8 to 56%). The Landry et al. genomes indicated a role for SAR202 in the oxidation of recalcitrant dissolved organic matter, and specifically cyclic alkanes, via flavin mononucleotide monooxygenases (FMNOs) and different dehydrogenases that occurred in paralogous groups (25). We observed many of the same gene expansions, namely that of MFS transporters and short-chain dehydrogenases (and related genes), but we did not recover any FMNOs of SFam 4832 or 4965, suggesting subclade- and/or niche-specific adaptations. Furthermore, we observed *napAB* and *nosZ* genes for nitrate and nitrous oxide reduction (and expression of these genes), which were not reported for subclade III or V. Our nGOM hypoxia SAR202 genomes had CAZy genes implicating them in degradation of complex compounds such as chitin and pectin. The emerging picture of these taxa from both shallow hypoxic waters and the mesopelagic zone is one of recalcitrant carbon degraders, with overlapping suites of paralogous genes, but that may be specialized for specific compounds more commonly available in their respective habitats.

This study has also developed roles for CP taxa in a shallow marine water column during hypoxia. The most active organism in our survey based on the ratio of RNA to DNA reads recruited, Bin 13, putatively belongs to a group with little genomic data—ACD39. The original ACD39 genome was reconstructed from an aquifer community (33). Although this was only a partial genome, it shared some features with our putative ACD39 member, namely pilin and chemotaxis genes, those containing TPR and tetratricopeptide repeats, and CAZy genes for degradation of complex compounds such as starch (33). Our study provides evidence that these taxa have relatively large genomes (~4.8 Mbp), are active aerobes in nGOM hypoxia, and have chemotaxis and motility genes that could facilitate scavenging and surface attachment. However, most of the highly expressed genes in this organism were annotated as hypothetical proteins, so much of the function of these organisms remains to be uncovered.

Bins 50 and 48 provide novel genome data for bacterioplankton in nGOM hypoxia, although the exact taxonomic position of these bins remains in conflict. The ribosomal protein tree provides evidence that these taxa belong to the *Latescibacteria* (WS3)

(Fig. S1), but 16S rRNA genes (Fig. S6) and our amplicon data point toward membership in the more poorly understood PAUC34f clade. Since no previous genome data exist for PAUC34f, we cannot rule out erroneous assignment in the ribosomal protein tree due to insufficient taxon selection. Bins 50 and 48 represented the largest genomes of the study, with estimated complete sizes of ~5 to 6 Mbp, and numerous genes suggesting degradation of a wider suite of complex organic matter than any of the other genomes examined. For example, they were the only genomes with numerous polysaccharide lyase genes, and these likely facilitate breakdown of algal cell wall components like pectin and alginate. The Bin 50 genome was among the most active across all samples (Fig. 2), and we detected expression of cytochrome *c* oxidase genes, and those for DNRA, in both the Bin 50 and 48 genomes. Thus, we expect these organisms to have an aerobic, potentially facultatively anaerobic, multifaceted chemoorganoheterotrophic metabolism with roles in complex carbon compound degradation (like that of algal cell walls) and the nitrogen cycle.

If these bins belong to PAUC34f, they represent the first genomic data for the group. Although originally discovered, and commonly found, in marine sponges (32, 62–64), this putative bacterial phylum (via GreenGenes/SILVA) has been detected as a rare group in other marine invertebrates (65) and stream sediment (66), and we identified 18 distinct but rare PAUC34f OTUs in nGOM hypoxia compared to just 3 from WS3 (21). Although the majority of studies suggest an endosymbiotic lifestyle for PAUC34f, our representative genome data point toward a free-living existence with multiple terminal electron-accepting processes, motility genes for seeking more favorable conditions, and a large metabolic repertoire for degradation of complex compounds. On the other hand, if these genomes represent WS3, the sister clade to PAUC34f (Fig. S6), they have many similarities to the lifestyles inferred from recent metagenomic investigations (67, 68). Specifically, while this group was previously considered anaerobic (67), new data have supported an aerobic lifestyle for some members (48) and revealed complete electron transport chains and both high- and low-affinity cytochrome *c* oxidases (68). The Bin 50 and 48 genomes predict aerobic metabolism as well, although only with low-affinity cytochrome *c* oxidases. Farag et al. also found little evidence of these taxa in host-associated environments, contrary to PAUC34f sequence data (68). The enrichment of PL family genes in Bins 50 and 48, polysaccharide degradation capability in general, and specific genes for degradation of cell wall components, all corroborate previous findings on WS3 as well (68). Bin 50 had 78 annotated peptidases, nearly double that in all other genomes in the study (Bin 48 had 46), which also concurs with metagenomic predictions for WS3 (68). Our genomes differed from WS3 metagenomes principally in the predicted DNRA metabolism and the dramatic expansion of sulfatases. Although sulfatases were observed in WS3 metagenomes (68), they were not present in the numbers associated with Bin 50 ($n = 102$). A large cadre of sulfatases has been previously reported for *Lentisphaera* ($n = 267$) and *Pirellula* ($n = 110$) genomes (69, 70) and suggests specialization for degradation of sulfate esters to satisfy carbon and/or sulfur requirements.

Although *Parcubacteria* and *Peregrinibacteria* occurred in low abundance (Fig. S7) and we detected activity in only a few samples, their recovery in the hypoxic zone is notable because these organisms have generally been associated with anoxic environments. Our predicted genome sizes (~1.5 Mbp) corroborate previous reports of these organisms having small genomes (31, 48). We did not observe any genes associated with nitrogen or sulfur redox transitions, although we cannot rule out these capabilities entirely due to incomplete genomes. Regardless, we can hypothesize that *Parcubacteria* and *Peregrinibacteria* persist as members of the rare biosphere until they can take advantage of microanoxic niches in the water column where they participate in carbon cycling as obligately fermentative organisms.

Excluding *Parcubacteria* and *Peregrinibacteria*, the other uncultivated groups in the nGOM hypoxic zone had one or more genomes that encoded cytochrome *c* oxidases (and other electron transport chain components) for respiring oxygen, making these taxa likely only facultative anaerobes. Pervasive aerobic metabolism in an oxygen-

depleted water column may seem counterintuitive, yet despite DO being as low as $2.2 \mu\text{mol} \cdot \text{kg}^{-1}$ in the E2A sample, oxygen probably remained high enough to sustain aerobic microbes. As little as $\sim 0.3 \mu\text{mol} \cdot \text{kg}^{-1}$ oxygen inhibited denitrification in OMZ populations by 50% (71), and even *Escherichia coli* K-12 could grow aerobically at oxygen concentrations as low as 3 nM (72). Thus, for many organisms, active aerobic respiration likely persists even in suboxic waters during nGOM hypoxia.

Nevertheless, our data also suggest pervasive coreduction of alternative terminal electron acceptors (oxygen, nitrate, nitrite, nitrous oxide, and sulfur), sometimes within the same organism (Fig. 3). Coreduction of electron acceptors with different redox potentials across a community could indicate microniches and/or aggregates in the water column where DO concentrations drop below bulk values (40). Alternatively, this can occur with an abundance of electron donor, and overlapping redox processes have been reported in multiple environments, including aquatic ones (50, 73). Concurrent expression of genes for multiple terminal electron-accepting processes within a single organism has been proposed as a means of improved readiness for dynamic conditions, albeit at the cost of lower productivity (40). Given that many uncultivated taxa likely perform multiple terminal electron-accepting processes (and possibly do so simultaneously) and we found a comparative cornucopia of genes for degradation of chemoorganoheterotrophic energy sources, we hypothesize that niche differentiation within uncultivated hypoxic zone bacterioplankton occurs predominantly via specialization for different oxidizable substrates rather than for distinct roles in the canonical redox cascade (4, 5).

Importantly, many of the active uncultivated taxa also appeared adapted for degradation of complex carbon substrates. Such compounds might comprise the bulk of available organic matter during the later stages of hypoxia after initial oxygen depletion by microorganisms feeding on more labile carbon sources. Selection for chemoorganotrophic microbes adapted to utilize recalcitrant organic matter could also explain why organisms that do not require an exogenous carbon source, such as the chemolithoautotrophic *Nitrosopumilus*, proliferate during hypoxia (21, 22) compared to their levels during spring before DO decreases (74, 75). Temporal data on the relative abundance and activity of these nGOM microbial dark matter organisms, and of organic matter composition in the water column, will be critical to more fully understand the relationship of bacterioplankton to the creation, maintenance, and dissolution of nGOM hypoxia.

MATERIALS AND METHODS

Sample selection and nucleic acid processing. Six samples representing hypoxic ($n = 4$) and oxic ($n = 2$) dissolved oxygen (DO) concentrations were picked from among those previously reported (21) at stations D1, D2, D3, E2, E2A, and E4 (see Table S1 at <http://thethrashlab.com/publications>). DO and nutrient collection information is detailed in the study by Gillies et al. (21). Nucleic acids were collected as follows. At these six stations, 10 liters of seawater was collected and filtered with a peristaltic pump. A $2.7 \mu\text{m}$ Whatman GF/D prefilter was used, and samples were concentrated on $0.22 \mu\text{m}$ Sterivex filters (EMD Millipore). Sterivex filters were immediately sparged, filled with RNAlater, and placed at -20°C , where they were maintained until extraction. DNA and RNA were extracted directly off the filter by placing half of the Sterivex filter in a Lysing matrix E (LME) glass/zirconia/silica beads tube (MP Biomedicals, Santa Ana, CA) using the protocol described by Gillies et al. (21) which combines phenol:chloroform:isoamyl alcohol (25:24:1) and bead beating. Genomic DNA and RNA were stored at -80°C until purified. DNA and RNA were purified using a Qiagen (Valencia, CA) AllPrep DNA/RNA kit. DNA quantity was determined using a Qubit2.0 fluorometer (Life Technologies, Grand Island, NY). RNA with an RNA integrity number (RIN) (16S/23S rRNA ratio determined with the Agilent TapeStation) of ≥ 8 (on a scale of 1 to 10, with 1 being degraded and 10 being undegraded RNA) was selected for metatranscriptomic sequencing. Using a Ribo-Zero kit (Illumina), rRNA was subtracted from total RNA. Subsequently, mRNA was reverse transcribed to cDNA as described by Mason et al. (76).

Sequencing, assembly, and binning. DNA and RNA were sequenced separately, six samples per lane, with Illumina HiSeq 2000 chemistry to generate 100 bp, paired-end reads (180 bp insert size) at the Argonne National Laboratory Next Generation Sequencing Facility. The data are available at the NCBI SRA repository under the BioSample accession numbers [SAMN05791315](https://www.ncbi.nlm.nih.gov/sra/SAMN05791315) to [SAMN05791320](https://www.ncbi.nlm.nih.gov/sra/SAMN05791320) (DNA) and [SAMN05791321](https://www.ncbi.nlm.nih.gov/sra/SAMN05791321) to [SAMN05791326](https://www.ncbi.nlm.nih.gov/sra/SAMN05791326) (RNA). DNA sequencing resulted in a total of 416,924,120 reads that were quality trimmed to 413,094,662 reads after adaptors were removed using Scythe (<https://github.com/vsbuffalo/scythe>), and low-quality reads ($Q < 30$) were trimmed using Sickle (<https://github.com/najoshi/sickle>). Reads with three or more N's or with an average quality score of less than Q20 and a length of < 50 bp were removed. Metagenomic reads from all six samples were pooled, assembled, and binned using previously described methods (77, 78). Briefly, quality-filtered reads were assembled with IDBA-UD (79) on a 1TB RAM, 40-core node at the Louisiana State University high-performance computing

cluster SuperMikell, using the following settings: `-mink 65 -maxk 105 -step 10 -pre_correction -seed_kmer 55`. Initial binning of the assembled fragments was performed using tetranucleotide frequency signatures using 5 kbp fragments of the contigs. Emergent self-organizing maps (ESOM) were manually delineated and curated based on clusters within the map. The primary assembly utilized all reads and produced 28,080 contigs ≥ 3 kbp totaling 217,715,956 bp. Of these, 303 contigs were over 50 kbp, 72 were over 100 kbp, and the largest contig was just under 495 kbp. Binning produced 76 genomes, of which 20 genomes were assigned to lineages with uncultivated representatives using CheckM, ribosomal protein trees, and 16S rRNA gene sequences (below).

DNA and RNA mapping. Metagenomic and metatranscriptomic sequencing reads from each sample were separately mapped to binned contigs using BWA (80) to compare bin abundance across samples and facilitate bin cleanup (below). Contigs within each bin were concatenated into a single fasta sequence, and BWA was used to map the reads from each sample to all bins. All commands used for these steps are available in the supplemental information at <http://thethrashlab.com/publications>.

Bin quality control. Bins were examined for contamination and completeness with CheckM (35), and we attempted to clean bins with $>10\%$ estimated contamination using a combination of methods. First, the CheckM modify command removed contigs determined to be outliers by GC content, coding density, and tetranucleotide frequency. Next, in bins that still showed $>10\%$ contamination, contigs were separated according to the comparative relative abundance of mean DNA read coverage by sample. Final bins were evaluated with CheckM again to generate the statistics in Table S1 at <http://thethrashlab.com/publications> and final bin placements in the CheckM concatenated gene tree (Fig. S2).

Ribosomal protein tree. The concatenated ribosomal protein tree was generated using 16 syntenic genes that have been shown to undergo limited lateral gene transfer (rpl2, 3, 4, 5, 6, 14, 15, 16, 18, 22, and 24 and rp53, 8, 10, 17, and 19) (81). Ribosomal proteins for each bin were identified with PhyloSift (82). Amino acid alignments of the individual ribosomal proteins were generated using MUSCLE (83) and trimmed using BMGE (84) (with the following settings: `-m BLOSUM30 -g 0.5`). The curated alignments were then concatenated for phylogenetic analyses, and phylogeny was inferred via RAxML v 8.2.8 (85) with 100 bootstrap runs (with the following settings: `mpirun -np 4 -npernode 1 raxmlHPC-HYBRID-AVX -f a -m PROTCATLG -T 16 -p 12345 -x 12345 -# 100`). Note that this is similar to the number utilized in a previous publication for this tree with automated bootstrapping (34), and required just over 56 h of wall clock time. The alignment is available in the supplemental information at <http://thethrashlab.com/publications>.

Average amino acid identity. Average amino acid identity (AAI) was calculated with GET_HOMOLOGUES (86) v. 02032017, with the following settings: `-M -t 0 -n 16 -A`.

Taxonomic assignment. Taxonomy for each bin was assigned primarily using the ribosomal protein tree. However, for bins that did not have enough ribosomal proteins to be included in the tree, or for bins for which the placement within the tree was poorly supported, assignments were made based on the concatenated marker gene tree as part of the CheckM analysis (Fig. S2) or via 16S rRNA gene sequences, when available. 16S rRNA genes were identified via CheckM, and these sequences were aligned against the NCBI nr database using BLASTN to corroborate CheckM assignments. In the case of the SAR202 genomes, which did not have representative genomes in either the ribosomal protein tree or the CheckM tree, the 16S rRNA gene sequences for two of the three bins (43-1 and 43-2) were available and aligned with the sequences used to define the SAR202 clade (24) (Fig. S5). Alignment, culling, and inference were completed with MUSCLE (83), Gblocks (87), and FastTree2 (88), respectively, with the FT_pipe script. The script is provided in the supplemental information at <http://thethrashlab.com/publications>. The 16S rRNA gene tree for subclade assignment of SAR406 (Fig. S4) was assembled by subjecting the four 16S sequences predicted by CheckM to a BLAST search against a local GenBank nucleotide database using blastn (v. 2.2.28+) (89), selecting the top 100 nonredundant hits to each sequence, and manually removing all hits to genome sequences. These sequences were combined with previously defined SAR406 subclade reference sequences (26), fosmid 16S sequences (27), single-cell genome sequences (23), and run through alignment, culling, and inference with FT_pipe. Taxa with identical alignments were removed with RAxML v 8.2.8 (85) using default settings, and the final tree was inferred using FastTree2 (88). For putative candidate phylum (CP) genomes, taxonomy was also evaluated by examining the taxonomic identification for each of the predicted protein sequences after a BLASTP search against the NCBI nr database. After the BLAST search, the number of assignments to the dominant one or two taxonomic names, along with the number of assignments to “uncultured bacterium,” was plotted for each genome according to the bit score quartile (Fig. S3). Quartiles were determined in R using the summary function. Bin 56 has two ribosomal protein operons on scaffold_2719/Ga0113622_1153 and scaffold_21777/Ga0113622_1009. In the ribosomal protein tree, the former placed the organism in the *Planctomycetaceae*, while the latter (which was much smaller) placed the organism in CP WS3. The majority of BLASTP annotations to the nr database matched *Planctomycetaceae* taxa, as did the 16S rRNA gene sequences found in the genome, so the Bin 56 organism was designated a *Planctomycetes* and not WS3 and excluded from this study. The 16S rRNA gene from Bin 50 was also used to infer taxonomic identity using an established phylogeny for the WS3 clade (68) and relevant outgroups. The Bin 50 sequence was subjected to a BLAST search against the GreenGenes database (December 2013) with megablast, and since many of the top hits belonged to the PAUC34f clade, these were included with the sequences from Farag et al. (68). Alignment, culling, and inference were completed with FT_pipe. Node labels were constructed with the newick utilities (90) script nw_rename.

Metabolic reconstruction. After binning, genomes were submitted individually to IMG (Integrated Microbial Genomes) (91) for annotation. Genome accession numbers are in Table 1, and all are publically available. Metabolic reconstruction found in Fig. S5 to S7 and in Table S1 and Fig. S11 to S13 at

<http://thethrashlab.com/publications> came from these annotations and inspection with IMG's analysis tools, including Kyoto Encyclopedia of Genes and Genomes (KEGG) pathway assignments and transporter predictions. Transporters highlighted for dissolved organic matter (DOM) uptake were identified based on information at the Transporter Classification Database (92). Carbohydrate-active enzymes (CAZymes) were predicted using the same routines as those in operation for the updates of the carbohydrate-active enzymes database (<http://www.cazy.org>) (93).

RPKM abundance of taxa and genes. Abundance of taxa within the sample was quantified by evaluating mapped reads using RPKM (reads per kilobase per million) normalization (94) according to $A_{ij} = (N_{ij}/L_i) \times (1/T_j)$, where A_{ij} is the abundance of bin i in sample j , N_{ij} is the number of reads that map to bin from sample j , L_i is the length of bin i in kilobases, and T_j is the total number of reads in sample j divided by 10^6 . These values were generated for all bins, with only the data for the 20 uncultivated bins reported here. All contigs within a given bin were artificially concatenated into "supercontigs" prior to mapping. N_{ij} was generated using the samtools (80) `idxstats` function after mapping with BWA. The data in Fig. S7 were created by summing (N_{ij}/L_i) for groups of taxa defined in Table S1 at <http://thethrashlab.com/publications> prior to multiplying by $(1/T_j)$. RNA coverage was used to evaluate both bin and gene activity for all bins. Mean coverage for each supercontig was calculated using bedtools (95), and bins were assigned a rank from lowest mean recruitment to highest mean recruitment. Bins with particularly high or low activity (transcript abundance) relative to their abundance (genome abundance) were identified using rank residuals, which were calculated as follows. On a plot of DNA coverage rank versus RNA coverage rank, residuals for each bin or gene were calculated from the identity. As the rank residuals followed a Gaussian distribution, bins with a residual that was >1 standard deviation (SD) from the rank residual mean were classified as having higher-than-expected transcriptional activity; bins with a residual that was <1 SD from the mean were classified as having lower-than-expected transcriptional activity. RPKM values were also calculated for every gene in every bin analogously to that for bins, using RNA mapping values extracted with the bedtools `multicov` function. Sample E2 was omitted from gene-specific calculations as only 4,588 transcriptomic reads mapped successfully from this sample, compared to $>100,000$ from other samples. Of 140,347 genes, 17,827 had no evidence of expression in any sample and so were removed from further analysis. A total of 3,840 genes recruited reads in all remaining samples. All calculations are available in Table S1 at <http://thethrashlab.com/publications> or the R markdown document `Per.gene.RPKM.Rmd` in supplemental information. Table S1 at <http://thethrashlab.com/publications> includes only analyzed data for the uncultivated bins reported in this study. Note that RPKM values indicate abundance measurements across a small number of samples. While we can evaluate the relative expression of genes for those samples, our data set lacks sufficient power to evaluate estimates of significance in differential expression.

***nrfA* sequence assessment.** Initial annotation of our bins identified putative homologs to the *nrfAH* genes associated with dissimilatory nitrite reduction to ammonia. Since *nrfA*-type nitrite reductases can be misannotated due to homology with other nitrite reductases, annotation for these genes was curated with phylogenetic analysis using known *nrfA* genes (38) obtained via A. Welsh (personal communication). Alignment, culling, and inference were completed with the FT_pipe script. The tree was rooted on the designated outgroup octaheme nitrite reductase sequence from *Thioalkalivibrio nitratireducens* ONR. Node labels were constructed with the newick utilities (90) script `nw_rename`. Visualization of the alignment (Fig. S8B and C) to confirm the presence of the first CXXCK/CXXCH and highly conserved KXQH/KXRH catalytic site was completed with the MSAViewer (96) online using the uncultured *nrfA* alignment as input.

SFam homology searches. To identify group-specific expansions in particular gene families, we performed a homology search of all predicted protein-coding sequences in each bin against the Sifted Families (SFam) database (41) using `hmmsearch` (HMMER 3.1b [97]) with default settings except for the utilization of 16 CPUs per search.

Accession numbers. Sequence read data are available at the NCBI SRA repository under the BioSample accession numbers [SAMN05791315](https://www.ncbi.nlm.nih.gov/sra/SAMN05791315) to [SAMN05791320](https://www.ncbi.nlm.nih.gov/sra/SAMN05791320) (DNA) and [SAMN05791321](https://www.ncbi.nlm.nih.gov/sra/SAMN05791321) to [SAMN05791326](https://www.ncbi.nlm.nih.gov/sra/SAMN05791326) (RNA).

Additional supplemental information. Table S1, scripts, workflows, and key files, including fasta files for each tree, are available at the Thrash Lab website: <http://thethrashlab.com/publications>.

SUPPLEMENTAL MATERIAL

Supplemental material for this article may be found at <https://doi.org/10.1128/mBio.01017-17>.

TEXT S1, DOCX file, 0.1 MB.

FIG S1, PDF file, 0.4 MB.

FIG S2, PDF file, 0.5 MB.

FIG S3, PDF file, 0.01 MB.

FIG S4, PDF file, 0.1 MB.

FIG S5, PDF file, 0.03 MB.

FIG S6, PDF file, 0.04 MB.

FIG S7, PDF file, 0.01 MB.

FIG S8, PDF file, 0.1 MB.

ACKNOWLEDGMENTS

We thank the crew of the R/V *Pelican*, Allana Welsh and Mostafa Elshahed for providing fasta files for *nrfA* and WS3 phylogenetic comparisons, and Mostafa Elshahed for helpful comments regarding WS3 phylogeny. Portions of this research were conducted with high-performance computing resources provided by Louisiana State University (<http://www.hpc.lsu.edu>).

Funding for this work was provided to J.C.T. through the Oak Ridge Associated Universities Ralph E. Powe Junior Faculty Enhancement Award and the Department of Biological Sciences of Louisiana State University. A portion of the funding for this work was provided by a Planning Grant award to O.U.M. from Florida State University. Funding for the research vessel and collection of oceanographic data was provided by the National Oceanic & Atmospheric Administration, Center for Sponsored Coastal Ocean Research award NA09NOS4780204 to N.N.R.

J.C.T. and O.U.M. designed the study. L.E.G., N.N.R., and J.C.T. collected samples. N.N.R. provided processed oceanographic data. L.E.G. and O.U.M. extracted and quantified nucleic acids and determined the quality of nucleic acids. J.C.T., K.W.S., and B.J.B. reconstructed the genomes. K.W.S., B.J.B., B.T., B.H., and J.C.T. conducted downstream analyses. J.C.T. led manuscript writing, and all coauthors evaluated and contributed edits.

REFERENCES

- Rabalais NN, Turner RE, Dortch Q, Justic D, Bierman VJ, Wiseman WJ, Jr. 2002. Nutrient-enhanced productivity in the northern Gulf of Mexico: past, present and future. *Hydrobiologia* 475:39–63.
- Diaz RJ, Rosenberg R. 2008. Spreading dead zones and consequences for marine ecosystems. *Science* 321:926–929. <https://doi.org/10.1126/science.1156401>.
- Mississippi River Gulf of Mexico Watershed Nutrient Task Force. 2008. Action plan 2008 for reducing, mitigating, and controlling hypoxia in the Northern Gulf of Mexico and improving water quality in the Mississippi River basin. Mississippi River Gulf of Mexico Watershed Nutrient Task Force, Office of Wetlands, Oceans, and Watersheds, US Environmental Protection Agency, Washington, DC.
- Ulloa O, Canfield DE, DeLong EF, Letelier RM, Stewart FJ. 2012. Microbial oceanography of anoxic oxygen minimum zones. *Proc Natl Acad Sci U S A* 109:15996–16003. <https://doi.org/10.1073/pnas.1205009109>.
- Wright JJ, Konwar KM, Hallam SJ. 2012. Microbial ecology of expanding oxygen minimum zones. *Nat Rev Microbiol* 10:381–394. <https://doi.org/10.1038/nrmicro2778>.
- Grote J, Jost G, Labrenz M, Herndl GJ, Jürgens K. 2008. Epsilonproteobacteria represent the major portion of chemoautotrophic bacteria in sulfidic waters of pelagic redoxclines of the Baltic and Black Seas. *Appl Environ Microbiol* 74:7546–7551. <https://doi.org/10.1128/AEM.01186-08>.
- Glaubitx S, Kießlich K, Meeske C, Labrenz M, Jürgens K. 2013. SUP05 dominates the gammaproteobacterial sulfur oxidizer assemblages in pelagic redoxclines of the Central Baltic and Black Seas. *Appl Environ Microbiol* 79:2767–2776. <https://doi.org/10.1128/AEM.03777-12>.
- Friedrich J, Janssen F, Aleynik D, Bange HW, Boltacheva N, Çagatay MN, Dale AW, Etiope G, Erdem Z, Geraga M, Gilli A, Gomoiu MT, Hall POJ, Hansson D, He Y, Holtappels M, Kirf MK, Kononets M, Kononov S, Lichtschlag A, Livingstone DM, Marinaro G, Mazlumyan S, Naeher S, North RP, Papatheodorou G, Pfannkuche O, Prien R, Rehder G, Schubert CJ, Soltwedel T, Sommer S, Stahl H, Stanev EV, Teaca A, Tengberg A, Waldmann C, Wehrli B, Wenzhöfer F. 2014. Investigating hypoxia in aquatic environments: diverse approaches to addressing a complex phenomenon. *Biogeosciences* 11: 1215–1259. <https://doi.org/10.5194/bg-11-1215-2014>.
- Lam P, Kuypers MM. 2011. Microbial nitrogen cycling processes in oxygen minimum zones. *Annu Rev Mar Sci* 3:317–345. <https://doi.org/10.1146/annurev-marine-120709-142814>.
- Beman JM, Popp BN, Alford SE. 2012. Quantification of ammonia oxidation rates and ammonia-oxidizing archaea and bacteria at high resolution in the Gulf of California and eastern tropical North Pacific Ocean. *Limnol Oceanogr* 57:711–726. <https://doi.org/10.4319/lo.2012.57.3.0711>.
- Saad EM, Longo AF, Chambers LR, Huang R, Benitez-Nelson C, Dyhrman ST, Diaz JM, Tang Y, Ingall ED. 2016. Understanding marine dissolved organic matter production: compositional insights from axenic cultures of *Thalassiosira pseudonana*. *Limnol Oceanogr* 61:2222–2233. <https://doi.org/10.1002/lno.10367>.
- Canfield DE, Stewart FJ, Thamdrup B, De Brabandere L, Dalsgaard T, DeLong EF, Revsbech NP, Ulloa O. 2010. A cryptic sulfur cycle in oxygen-minimum-zone waters off the Chilean coast. *Science* 330:1375–1378. <https://doi.org/10.1126/science.1196889>.
- Stewart FJ, Ulloa O, DeLong EF. 2012. Microbial metatranscriptomics in a permanent marine oxygen minimum zone. *Environ Microbiol* 14: 23–40. <https://doi.org/10.1111/j.1462-2920.2010.02400.x>.
- Roberts BJ, Doty SM. 2015. Spatial and temporal patterns of benthic respiration and net nutrient fluxes in the Atchafalaya River Delta estuary. *Estuaries Coasts* 38:1918–1936. <https://doi.org/10.1007/s12237-015-9965-z>.
- McCarthy MJ, Carini SA, Liu Z, Ostrom NE, Gardner WS. 2013. Oxygen consumption in the water column and sediments of the northern Gulf of Mexico hypoxic zone. *Estuarine Coast Shelf Sci* 123:46–53. <https://doi.org/10.1016/j.ecss.2013.02.019>.
- Schaeffer BA, Sinclair GA, Lehrter JC, Murrell MC, Kurtz JC, Gould RW, Yates DF. 2011. An analysis of diffuse light attenuation in the northern Gulf of Mexico hypoxic zone using the SeaWiFS satellite data record. *Remote Sens Environ* 115:3748–3757. <https://doi.org/10.1016/j.rse.2011.09.013>.
- Rabalais NN, Turner RE, Sen Gupta BKS, Boesch DF, Chapman P, Murrell MC. 2007. Hypoxia in the northern Gulf of Mexico: does the science support the plan to reduce, mitigate, and control hypoxia? *Estuaries Coasts* 30:753–772. <https://doi.org/10.1007/BF02841332>.
- Rabalais NN, Turner RE, Wiseman WJ, Jr. 2001. Hypoxia in the Gulf of Mexico. *J Environ Qual* 30:320–329. <https://doi.org/10.2134/jeq2001.302320x>.
- Rabalais NN, Turner RE, Wiseman WJ. 2002. Gulf of Mexico hypoxia, A.K.A. “the dead zone”. *Annu Rev Ecol Syst* 33:235–263. <https://doi.org/10.1146/annurev.ecolsys.33.010802.150513>.
- Wiseman WJ, Jr, Rabalais NN, Turner RE, Dinnel SP, MacNaughton A. 1997. Seasonal and interannual variability within the Louisiana coastal current: stratification and hypoxia. *J Mar Syst* 12:237–248. [https://doi.org/10.1016/S0924-7963\(96\)00100-5](https://doi.org/10.1016/S0924-7963(96)00100-5).
- Gillies LE, Thrash JC, deRada S, Rabalais NN, Mason OU. 2015. Archaeal enrichment in the hypoxic zone in the northern Gulf of Mexico. *Environ Microbiol* 17:3847–3856. <https://doi.org/10.1111/1462-2920.12853>.
- Bristow LA, Sarode N, Cartee J, Caro-Quintero A, Thamdrup B, Stewart FJ. 2015. Biogeochemical and metagenomic analysis of nitrite accumulation in the Gulf of Mexico hypoxic zone. *Limnol Oceanogr* 60:1733–1750. <https://doi.org/10.1002/lno.10130>.
- Rinke C, Schwientek P, Sczyrba A, Ivanova NN, Anderson IJ, Cheng J-F, Darling A, Malfatti S, Swan BK, Gies EA, Dodsworth JA, Hedlund BP,

- Tsiamis G, Sievert SM, Liu W-T, Eisen JA, Hallam SJ, Kyrpidis NC, Stepanauskas R, Rubin EM, Hugenholtz P, Woyke T. 2013. Insights into the phylogeny and coding potential of microbial dark matter. *Nature* 499: 431–437. <https://doi.org/10.1038/nature12352>.
24. Morris RM, Rappé MS, Urbach E, Connon SA, Giovannoni SJ. 2004. Prevalence of the Chloroflexi-related SAR202 bacterioplankton cluster throughout the mesopelagic zone and deep ocean. *Appl Environ Microbiol* 70:2836–2842. <https://doi.org/10.1128/AEM.70.5.2836-2842.2004>.
 25. Landry Z, Swan BK, Herndl GJ, Stepanauskas R, Giovannoni SJ. 2017. SAR202 genomes from the dark ocean predict pathways for the oxidation of recalcitrant dissolved organic matter. *mBio* 8:e00413-17. <https://doi.org/10.1128/mBio.00413-17>.
 26. Allers E, Wright JJ, Konwar KM, Howes CG, Beneze E, Hallam SJ, Sullivan MB. 2013. Diversity and population structure of Marine Group A bacteria in the Northeast subarctic Pacific Ocean. *ISME J* 7:256–268. <https://doi.org/10.1038/ismej.2012.108>.
 27. Wright JJ, Mewis K, Hanson NW, Konwar KM, Maas KR, Hallam SJ. 2014. Genomic properties of Marine Group A bacteria indicate a role in the marine sulfur cycle. *ISME J* 8:455–468. <https://doi.org/10.1038/ismej.2013.152>.
 28. Gordon DA, Giovannoni SJ. 1996. Detection of stratified microbial populations related to Chlorobium and Fibrobacter species in the Atlantic and Pacific Oceans. *Appl Environ Microbiol* 62:1171–1177.
 29. Parada AE, Needham DM, Fuhrman JA. 2016. Every base matters: assessing small subunit rRNA primers for marine microbiomes with mock communities, time series and global field samples. *Environ Microbiol* 18:1403–1414. <https://doi.org/10.1111/1462-2920.13023>.
 30. Brown CT, Hug LA, Thomas BC, Sharon I, Castelle CJ, Singh A, Wilkins MJ, Wrighton KC, Williams KH, Banfield JF. 2015. Unusual biology across a group comprising more than 15% of domain Bacteria. *Nature* 523: 208–211. <https://doi.org/10.1038/nature14486>.
 31. Wrighton KC, Thomas BC, Sharon I, Miller CS, Castelle CJ, Verberkmoes NC, Wilkins MJ, Hettich RL, Lipton MS, Williams KH, Long PE, Banfield JF. 2012. Fermentation, hydrogen, and sulfur metabolism in multiple uncultivated bacterial phyla. *Science* 337:1661–1665. <https://doi.org/10.1126/science.1224041>.
 32. Hentschel U, Hopke J, Horn M, Friedrich AB, Wagner M, Hacker J, Moore BS. 2002. Molecular evidence for a uniform microbial community in sponges from different oceans. *Appl Environ Microbiol* 68:4431–4440. <https://doi.org/10.1128/AEM.68.9.4431-4440.2002>.
 33. Wrighton KC, Castelle CJ, Wilkins MJ, Hug LA, Sharon I, Thomas BC, Handley KM, Mullin SW, Nicora CD, Singh A, Lipton MS, Long PE, Williams KH, Banfield JF. 2014. Metabolic interdependencies between phylogenetically novel fermenters and respiratory organisms in an unconfined aquifer. *ISME J* 8:1452–1463. <https://doi.org/10.1038/ismej.2013.249>.
 34. Hug LA, Baker BJ, Anantharaman K, Brown CT, Probst AJ, Castelle CJ, Butterfield CN, Hemsdorf AW, Amano Y, Ise K, Suzuki Y, Dudek N, Relman DA, Finstad KM, Amundson R, Thomas BC, Banfield JF. 2016. A new view of the tree of life. *Nat Microbiol* 1:16048. <https://doi.org/10.1038/nmicrobiol.2016.48>.
 35. Parks DH, Imelfort M, Skennerton CT, Hugenholtz P, Tyson GW. 2015. CheckM: assessing the quality of microbial genomes recovered from isolates, single cells, and metagenomes. *Genome Res* 25:1043–1055. <https://doi.org/10.1101/gr.186072.114>.
 36. Cantarel BL, Coutinho PM, Rancurel C, Bernard T, Lombard V, Henrissat B. 2009. The Carbohydrate-Active EnZymes database (CAZy): an expert resource for glycogenomics. *Nucleic Acids Res* 37:D233–D238. <https://doi.org/10.1093/nar/gkn663>.
 37. Pitcher RS, Watmough NJ. 2004. The bacterial cytochrome cbb3 oxidases. *Biochim Biophys Acta* 1655:388–399. <https://doi.org/10.1016/j.bbabi.2003.09.017>.
 38. Welsh A, Chee-Sanford JC, Connor LM, Löffler FE, Sanford RA. 2014. Refined nrfA phylogeny improves PCR-based nrfA gene detection. *Appl Environ Microbiol* 80:2110–2119. <https://doi.org/10.1128/AEM.03443-13>.
 39. Burns JL, DiChristina TJ. 2009. Anaerobic respiration of elemental sulfur and thiosulfate by *Shewanella oneidensis* MR-1 requires psrA, a homolog of the phsA gene of *Salmonella enterica* serovar typhimurium LT2. *Appl Environ Microbiol* 75:5209–5217. <https://doi.org/10.1128/AEM.00888-09>.
 40. Chen J, Hanke A, Tegetmeyer HE, Kattelman I, Sharma R, Hamann E, Hargesheimer T, Kraft B, Lenk S, Geelhoed JS, Hettich RL, Strous M. 2017. Impacts of chemical gradients on microbial community structure. *ISME J* 11:920–931. <https://doi.org/10.1038/ismej.2016.175>.
 41. Sharpton TJ, Jospin G, Wu D, Langille MGI, Pollard KS, Eisen JA. 2012. Sifting through genomes with iterative-sequence clustering produces a large, phylogenetically diverse protein-family resource. *BMC Bioinformatics* 13:264. <https://doi.org/10.1186/1471-2105-13-264>.
 42. McBride MJ, Zhu Y. 2013. Gliding motility and Por secretion system genes are widespread among members of the phylum Bacteroidetes. *J Bacteriol* 195:270–278. <https://doi.org/10.1128/JB.01962-12>.
 43. Varela MM, van Aken HM, Herndl GJ. 2008. Abundance and activity of Chloroflexi-type SAR202 bacterioplankton in the meso- and bathypelagic waters of the (sub)tropical Atlantic. *Environ Microbiol* 10: 1903–1911. <https://doi.org/10.1111/j.1462-2920.2008.01627.x>.
 44. Law CJ, Maloney PC, Wang D-N. 2008. Ins and outs of major facilitator superfamily antiporters. *Annu Rev Microbiol* 62:289–305. <https://doi.org/10.1146/annurev.micro.61.080706.093329>.
 45. Lombard V, Bernard T, Rancurel C, Brumer H, Coutinho PM, Henrissat B. 2010. A hierarchical classification of polysaccharide lyases for glycogenomics. *Biochem J* 432:437–444. <https://doi.org/10.1042/BJ20101185>.
 46. Cregut M, Piutti S, Slezack-Deschaumes S, Benizri E. 2013. Compartmentalization and regulation of arylsulfatase activities in *Streptomyces* sp., *Microbacterium* sp. and *Rhodococcus* sp. soil isolates in response to inorganic sulfate limitation. *Microbiol Res* 168:12–21. <https://doi.org/10.1016/j.micres.2012.08.001>.
 47. Cregut M, Durand M-J, Thouand G. 2014. The diversity and functions of choline sulphatases in microorganisms. *Microb Ecol* 67:350–357. <https://doi.org/10.1007/s00248-013-0328-7>.
 48. Anantharaman K, Brown CT, Hug LA, Sharon I, Castelle CJ, Probst AJ, Thomas BC, Singh A, Wilkins MJ, Karaoz U, Brodie EL, Williams KH, Hubbard SS, Banfield JF. 2016. Thousands of microbial genomes shed light on interconnected biogeochemical processes in an aquifer system. *Nat Commun* 7:13219. <https://doi.org/10.1038/ncomms13219>.
 49. Luef B, Frischkorn KR, Wrighton KC, Holman H-YN, Birarda G, Thomas BC, Singh A, Williams KH, Siegerist CE, Tringe SG, Downing KH, Comolli LR, Banfield JF. 2015. Diverse uncultivated ultra-small bacterial cells in groundwater. *Nat Commun* 6:6372. <https://doi.org/10.1038/ncomms7372>.
 50. Alewell C, Paul S, Lischeid G, Storck FR. 2008. Co-regulation of redox processes in freshwater wetlands as a function of organic matter availability? *Sci Total Environ* 404:335–342. <https://doi.org/10.1016/j.scitotenv.2007.11.001>.
 51. Li M, Baker BJ, Anantharaman K, Jain S, Breier JA, Dick GJ. 2015. Genomic and transcriptomic evidence for scavenging of diverse organic compounds by widespread deep-sea archaea. *Nat Commun* 6:8933. <https://doi.org/10.1038/ncomms9933>.
 52. Zhang CL, Xie W, Martin-Cuadrado A-B, Rodriguez-Valera F. 2015. Marine Group II Archaea, potentially important players in the global ocean carbon cycle. *Front Microbiol* 6:1108. <https://doi.org/10.3389/fmicb.2015.01108>.
 53. Orsi WD, Smith JM, Liu S, Liu Z, Sakamoto CM, Wilken S, Poirier C, Richards TA, Keeling PJ, Worden AZ, Santoro AE. 2016. Diverse, uncultivated bacteria and archaea underlying the cycling of dissolved protein in the ocean. *ISME J* 10:2158–2173. <https://doi.org/10.1038/ismej.2016.20>.
 54. Ouverney CC, Fuhrman JA. 2000. Marine planktonic archaea take up amino acids. *Appl Environ Microbiol* 66:4829–4833. <https://doi.org/10.1128/AEM.66.11.4829-4833.2000>.
 55. Iverson V, Morris RM, Frazar CD, Berthiaume CT, Morales RL, Armbrust EV. 2012. Untangling genomes from metagenomes: revealing an uncultured class of marine Euryarchaeota. *Science* 335:587–590. <https://doi.org/10.1126/science.1212665>.
 56. Fuhrman JA, McCallum K, Davis AA. 1993. Phylogenetic diversity of subsurface marine microbial communities from the Atlantic and Pacific Oceans. *Appl Environ Microbiol* 59:1294–1302.
 57. Schattenhofer M, Fuchs BM, Amann R, Zubkov MV, Tarran GA, Pernthaler J. 2009. Latitudinal distribution of prokaryotic picoplankton populations in the Atlantic Ocean. *Environ Microbiol* 11:2078–2093. <https://doi.org/10.1111/j.1462-2920.2009.01929.x>.
 58. Hu P, Tom L, Singh A, Thomas BC, Baker BJ, Piceno YM, Andersen GL, Banfield JF. 2016. Genome-resolved metagenomic analysis reveals roles for candidate phyla and other microbial community members in biogeochemical transformations in oil reservoirs. *mBio* 7:e01669-15. <https://doi.org/10.1128/mBio.01669-15>.
 59. DeLong EF, Preston CM, Mincer T, Rich V, Hallam SJ, Frigaard N-UU, Martinez A, Sullivan MB, Edwards R, Brito BR, Chisholm SW, Karl DM. 2006. Community genomics among stratified microbial assemblages in the ocean's interior. *Science* 311:496–503. <https://doi.org/10.1126/science.1120250>.
 60. Fuchs BM, Woebken D, Zubkov MV, Burkil P, Amann R. 2005. Molecular

- identification of picoplankton populations in contrasting waters of the Arabian Sea. *Aquat Microb Ecol* 39:145–157. <https://doi.org/10.3354/ame039145>.
61. Treusch AH, Vergin KL, Finlay LA, Donatz MG, Burton RM, Carlson CA, Giovannoni SJ. 2009. Seasonality and vertical structure of microbial communities in an ocean gyre. *ISME J* 3:1148–1163. <https://doi.org/10.1038/ismej.2009.60>.
 62. Weigel BL, Erwin PM. 2017. Effects of reciprocal transplantation on the microbiome and putative nitrogen cycling functions of the intertidal sponge, *Hymeniacidon heliophila*. *Sci Rep* 7:43247. <https://doi.org/10.1038/srep43247>.
 63. Jeong I-H, Kim K-H, Park J-S. 2013. Analysis of bacterial diversity in sponges collected off Chujado, an island in Korea, using barcoded 454 pyrosequencing: analysis of a distinctive sponge group containing *Chloroflexi*. *J Microbiol* 51:570–577. <https://doi.org/10.1007/s12275-013-3426-9>.
 64. Cuvelier ML, Blake E, Mulheron R, McCarthy PJ, Blackwelder P, Thurber RL, Lopez JV. 2014. Two distinct microbial communities revealed in the sponge *Cinachyrella*. *Front Microbiol* 5:581. <https://doi.org/10.3389/fmicb.2014.00581>.
 65. Erwin PM, Pineda MC, Webster N, Turon X, López-Legentil S. 2014. Down under the tunic: bacterial biodiversity hotspots and widespread ammonia-oxidizing archaea in coral reef ascidians. *ISME J* 8:575–588. <https://doi.org/10.1038/ismej.2013.188>.
 66. Costa PS, Reis MP, Ávila MP, Leite LR, de Araújo FMG, Salim ACM, Oliveira G, Barbosa F, Chartone-Souza E, Nascimento AMA. 2015. Metagenome of a microbial community inhabiting a metal-rich tropical stream sediment. *PLoS One* 10:e0119465. <https://doi.org/10.1371/journal.pone.0119465>.
 67. Youssef NH, Farag IF, Rinke C, Hallam SJ, Woyke T, Elshahed MS. 2015. In silico analysis of the metabolic potential and niche specialization of candidate phylum “Latescibacteria” (WS3). *PLoS One* 10:e0127499. <https://doi.org/10.1371/journal.pone.0127499>.
 68. Farag IF, Youssef NH, Elshahed MS. 2017. Global distribution patterns and pangenomic diversity of the candidate phylum “Latescibacteria” (WS3). *Appl Environ Microbiol* 83:e00521-17. <https://doi.org/10.1128/AEM.00521-17>.
 69. Thrash JC, Cho JC, Vergin KL, Morris RM, Giovannoni SJ. 2010. Genome sequence of *Lentisphaera araneosa* HTCC2155T, the type species of the order *Lentisphaerales* in the phylum *Lentisphaerae*. *J Bacteriol* 192:2938–2939. <https://doi.org/10.1128/JB.00208-10>.
 70. Krüger M, Meyerdierks A, Glöckner FO, Amann R, Widdel F, Kube M, Reinhardt R, Kahnt J, Böcher R, Thauer RK, Shima S. 2003. A conspicuous nickel protein in microbial mats that oxidize methane anaerobically. *Nature* 426:878–881. <https://doi.org/10.1038/nature02207>.
 71. Dalsgaard T, Stewart FJ, Thamdrup B, De Brabandere L, Revsbech NP, Ulloa O, Canfield DE, DeLong EF. 2014. Oxygen at nanomolar levels reversibly suppresses process rates and gene expression in anammox and denitrification in the oxygen minimum zone off Northern Chile. *mBio* 5:e01966-14. <https://doi.org/10.1128/mBio.01966-14>.
 72. Stolper DA, Revsbech NP, Canfield DE. 2010. Aerobic growth at nanomolar oxygen concentrations. *Proc Natl Acad Sci U S A* 107:18755–18760. <https://doi.org/10.1073/pnas.1013435107>.
 73. Eggleston EM, Lee DY, Owens MS, Cornwell JC, Crump BC, Hewson I. 2015. Key respiratory genes elucidate bacterial community respiration in a seasonally anoxic estuary. *Environ Microbiol* 17:2306–2318. <https://doi.org/10.1111/1462-2920.12690>.
 74. King GM, Smith CB, Tolar B, Hollibaugh JT. 2012. Analysis of composition and structure of coastal to mesopelagic bacterioplankton communities in the northern gulf of Mexico. *Front Microbiol* 3:438. <https://doi.org/10.3389/fmicb.2012.00438>.
 75. Tolar BB, King GM, Hollibaugh JT. 2013. An analysis of *Thaumarchaeota* populations from the northern gulf of Mexico. *Front Microbiol* 4:72. <https://doi.org/10.3389/fmicb.2013.00072>.
 76. Mason OU, Hazen TC, Borglin S, Chain PSG, Dubinsky EA, Fortney JL, Han J, Holman H-YN, Hultman J, Lamendella R, Mackelprang R, Malfatti S, Tom LM, Tringe SG, Woyke T, Zhou J, Rubin EM, Jansson JK. 2012. Metagenome, metatranscriptome and single-cell sequencing reveal microbial response to Deepwater Horizon oil spill. *ISME J* 6:1715–1727. <https://doi.org/10.1038/ismej.2012.59>.
 77. Dick GJ, Andersson AF, Baker BJ, Simmons SL, Thomas BC, Yelton AP, Banfield JF. 2009. Community-wide analysis of microbial genome sequence signatures. *Genome Biol* 10:R85. <https://doi.org/10.1186/gb-2009-10-8-r85>.
 78. Sharon I, Morowitz MJ, Thomas BC, Costello EK, Relman DA, Banfield JF. 2013. Time series community genomics analysis reveals rapid shifts in bacterial species, strains, and phage during infant gut colonization. *Genome Res* 23:111–120. <https://doi.org/10.1101/gr.142315.112>.
 79. Peng Y, Leung HCM, Yiu SM, Chin FYL. 2012. IDBA-UD: a de novo assembler for single-cell and metagenomic sequencing data with highly uneven depth. *Bioinformatics* 28:1420–1428. <https://doi.org/10.1093/bioinformatics/bts174>.
 80. Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* 25:1754–1760. <https://doi.org/10.1093/bioinformatics/btp324>.
 81. Sorek R, Zhu Y, Creevey CJ, Francino MP, Bork P, Rubin EM. 2007. Genome-wide experimental determination of barriers to horizontal gene transfer. *Science* 318:1449–1452. <https://doi.org/10.1126/science.1147112>.
 82. Darling AE, Jospin G, Lowe E, Matsen FA, IV, Bik HM, Eisen JA. 2014. PhyloSift: phylogenetic analysis of genomes and metagenomes. *PeerJ* 2:e243. <https://doi.org/10.7717/peerj.243>.
 83. Edgar RC. 2004. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res* 32:1792–1797. <https://doi.org/10.1093/nar/gkh340>.
 84. Criscuolo A, Gribaldo S. 2010. BMGE (Block Mapping and Gathering with Entropy): a new software for selection of phylogenetic informative regions from multiple sequence alignments. *BMC Evol Biol* 10:210. <https://doi.org/10.1186/1471-2148-10-210>.
 85. Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics* 30:1312–1313. <https://doi.org/10.1093/bioinformatics/btu033>.
 86. Contreras-Moreira B, Vinuesa P. 2013. GET_HOMOLOGUES, a versatile software package for scalable and robust microbial pangenome analysis. *Appl Environ Microbiol* 79:7696–7701. <https://doi.org/10.1128/AEM.02411-13>.
 87. Castresana J. 2000. Selection of conserved blocks from multiple alignments for their use in phylogenetic analysis. *Mol Biol Evol* 17:540–552. <https://doi.org/10.1093/oxfordjournals.molbev.a026334>.
 88. Price MN, Dehal PS, Arkin AP. 2010. FastTree 2—approximately maximum-likelihood trees for large alignments. *PLoS One* 5:e9490. <https://doi.org/10.1371/journal.pone.0009490>.
 89. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL. 2009. BLAST+: architecture and applications. *BMC Bioinformatics* 10:421. <https://doi.org/10.1186/1471-2105-10-421>.
 90. Junier T, Zdobnov EM. 2010. The Newick Utilities: high-throughput phylogenetic tree processing in the Unix shell. *Bioinformatics* 26:1669–1670. <https://doi.org/10.1093/bioinformatics/btq243>.
 91. Markowitz VM, Chen IMA, Palaniappan K, Chu K, Szeto E, Pillay M, Ratner A, Huang J, Woyke T, Huntemann M, Anderson I, Billis K, Varghese N, Mavromatis K, Pati A, Ivanova NN, Kyrpides NC. 2014. IMG 4 version of the integrated microbial genomes comparative analysis system. *Nucleic Acids Res* 42:D560–D567. <https://doi.org/10.1093/nar/gkt963>.
 92. Saier MH, Reddy VS, Tsu BV, Ahmed MS, Li C, Moreno-Hagelsieb G. 2016. The Transporter Classification Database (TCDB): recent advances. *Nucleic Acids Res* 44:D372–D379. <https://doi.org/10.1093/nar/gkv1103>.
 93. Lombard V, Golaconda Ramulu H, Drula E, Coutinho PM, Henrissat B. 2014. The carbohydrate-active enzymes database (CAZY) in 2013. *Nucleic Acids Res* 42:D490–D495. <https://doi.org/10.1093/nar/gkt1178>.
 94. Mortazavi A, Williams BA, McCue K, Schaeffer L, Wold B. 2008. Mapping and quantifying mammalian transcriptomes by RNA-Seq. *Nat Methods* 5:621–628. <https://doi.org/10.1038/nmeth.1226>.
 95. Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* 26:841–842. <https://doi.org/10.1093/bioinformatics/btq033>.
 96. Yachdav G, Wilzbach S, Rauscher B, Sheridan R, Sillitoe I, Procter J, Lewis SE, Rost B, Goldberg T. 2016. MSASviewer: interactive JavaScript visualization of multiple sequence alignments. *Bioinformatics* 32:3501–3503. <https://doi.org/10.1093/bioinformatics/btw474>.
 97. Eddy SR. 2011. Accelerated profile HMM searches. *PLoS Comput Biol* 7:e1002195. <https://doi.org/10.1371/journal.pcbi.1002195>.