



Genome Streamlining, Proteorhodopsin, and Organic Nitrogen Metabolism in Freshwater Nitrifiers

Justin C. Podowski,^a Sara F. Paver,^a Ryan J. Newton,^b  Maureen L. Coleman^a

^aDepartment of the Geophysical Sciences, University of Chicago, Chicago, Illinois, USA

^bSchool of Freshwater Sciences, University of Wisconsin Milwaukee, Milwaukee, Wisconsin, USA

ABSTRACT Microbial nitrification is a critical process governing nitrogen availability in aquatic systems. Freshwater nitrifiers have received little attention, leaving many unanswered questions about their taxonomic distribution, functional potential, and ecological interactions. Here, we reconstructed genomes to infer the metabolism and ecology of free-living picoplanktonic nitrifiers across the Laurentian Great Lakes, a connected series of five of Earth's largest lakes. Surprisingly, ammonia-oxidizing bacteria (AOB) related to *Nitrosospira* dominated over ammonia-oxidizing archaea (AOA) at nearly all stations, with distinct ecotypes prevailing in the transparent, oligotrophic upper lakes compared to Lakes Erie and Ontario. Unexpectedly, one ecotype of *Nitrosospira* encodes proteorhodopsin, which could enhance survival under conditions where ammonia oxidation is inhibited or substrate limited. Nitrite-oxidizing bacteria (NOB) "*Candidatus Nitrotoga*" and *Nitrosospira* fluctuated in dominance, with the latter prevailing in deeper, less-productive basins. Genome reconstructions reveal highly reduced genomes and features consistent with genome streamlining, along with diverse adaptations to sunlight and oxidative stress and widespread capacity for organic nitrogen use. Our findings expand the known functional diversity of nitrifiers and establish their ecological genomics in large lake ecosystems. By elucidating links between microbial biodiversity and biogeochemical cycling, our work also informs ecosystem models of the Laurentian Great Lakes, a critical freshwater resource experiencing rapid environmental change.

IMPORTANCE Microorganisms play critical roles in Earth's nitrogen cycle. In lakes, microorganisms called nitrifiers derive energy from reduced nitrogen compounds. In doing so, they transform nitrogen into a form that can ultimately be lost to the atmosphere by a process called denitrification, which helps mitigate nitrogen pollution from fertilizer runoff and sewage. Despite their importance, freshwater nitrifiers are virtually unexplored. To understand their diversity and function, we reconstructed genomes of freshwater nitrifiers across some of Earth's largest freshwater lakes, the Laurentian Great Lakes. We discovered several new species of nitrifiers specialized for clear low-nutrient waters and distinct species in comparatively turbid Lake Erie. Surprisingly, one species may be able to harness light energy by using a protein called proteorhodopsin, despite the fact that nitrifiers typically live in deep dark water. Our work reveals the unique biodiversity of the Great Lakes and fills key gaps in our knowledge of an important microbial group, the nitrifiers.

KEYWORDS biogeochemistry, ecological genomics, freshwater, metagenomics, nitrification

The oxidation of ammonia to nitrate powers the growth of nitrifying microorganisms and represents a critical flux in the global nitrogen cycle. Microbial nitrification of ammonia released from organic matter degradation produces nitrate, which can then be removed from the system by denitrification (1). As chemolithoautotrophs,

Editor Stephen J. Giovannoni, Oregon State University

Copyright © 2022 Podowski 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 Maureen L. Coleman, mlcoleman@uchicago.edu.

The authors declare no conflict of interest.

Received 7 August 2021

Accepted 10 March 2022

Published 18 April 2022

nitrifiers are also a major source of dark carbon fixation (2), which may contribute significant organic carbon to the microbial food web of the ocean's interior (3–5) and of deep freshwater lakes (6).

Microbial nitrifiers are found in *Archaea* and several phyla of *Bacteria*, spanning diverse physiology and ecology. Ammonia-oxidizing archaea (AOA) in the phylum *Thaumarchaeota* dominate the mesopelagic oceans (7), likely due to their high affinity for ammonia (8) and streamlined genomes (9). In freshwater systems, AOA are abundant in some oligotrophic lakes, while ammonia-oxidizing bacteria (AOB) affiliated with the *Nitrosomonadaceae* (*Betaproteobacteria*) tend to dominate more eutrophic systems (10–16). Complicating this picture, however, there is considerable physiological variation within both the AOA and AOB, such as low-nutrient-adapted clades of AOB (17, 18) and the ability of some strains to use alternative substrates like urea (18, 19). Within the AOA, there are also distinct ecotypes that appear to segregate with depth in the water column, in both marine (7) and freshwater systems (10). In freshwaters especially—which are poorly characterized compared to the oceans—it remains difficult to predict which AOA and AOB taxa are likely to dominate in a given system (16).

For aquatic nitrite oxidizing bacteria (NOB), which span the phyla *Nitrospira*, *Nitrospinae*, and *Proteobacteria*, niche differentiation is even less clear. The oceans are dominated by exclusively marine lineages (2, 20), consistent with ancient salinity-associated divergence. Among non-marine NOB, cultivated strains show variation in substrate affinity and other physiological traits (20–22), but connecting these culture-based studies to natural ecosystems remains a challenge. Moreover, recent studies have discovered that NOB are capable of alternative energy metabolisms (23, 24) and can access nitrogen from cyanate and urea (25, 26), expanding their ecological potential. In freshwater systems, the NOB "*Candidatus Nitrotoga*" (*Betaproteobacteria*) was only recently discovered to be widespread (27), and the diversity of this genus and factors favoring its success are unknown.

Here, we use the Laurentian Great Lakes (GL) as a model system to examine niche partitioning among planktonic freshwater nitrifiers. The Great Lakes hold 20% of Earth's surface freshwater, and more than half of this volume receives little to no light (<1% surface irradiance). This system, while hydrologically connected, spans strong trophic and chemical gradients: ultraoligotrophic Lake Superior supports low rates of primary production and nitrification comparable to the ocean gyres (28, 29), while Lake Erie supports greater production (30) and more than 70-fold-higher nitrification rates (31). Between these extremes, Lake Ontario has low ambient ammonium concentrations like Lake Superior (32) but nitrification rates up to four times higher (33). While previous studies reported that AOA and AOB dominate Lakes Superior and Erie, respectively (14, 29), recent community profiling has revealed broader diversity in both ammonia-oxidizing and nitrite-oxidizing lineages (34–36). We sought to link taxonomic, genomic, and metabolic diversity of nitrifiers with the varied biogeochemistry of the Great Lakes, using genome reconstructions and abundance profiling. Our results uncover novel lineages and metabolic capabilities and provide the first large-scale assessment of freshwater nitrifier genomics.

RESULTS AND DISCUSSION

Niche partitioning of nitrifiers across the Great Lakes. To map free-living picoplanktonic (here defined as cells that pass through a 1.6- μ m filter) nitrifiers across the Great Lakes, we searched our recent 16S rRNA data sets for known nitrifying taxa (34). We detected putative AOB in the genus *Nitrosospira* (*Betaproteobacteria*, family *Nitrosomonadaceae*) and AOA in the genus *Nitrosarchaeum* (family *Nitrosopumilaceae*), along with putative NOB in the genera "*Ca. Nitrotoga*" (*Betaproteobacteria*, family *Gallionellaceae*) and *Nitrospira* (family *Nitrospiraceae*). We did not detect 16S rRNA amplicons from *Nitrosococcus*, *Nitrococcus*, *Nitrospina*, or *Nitrobacter*. The highest relative abundances of picoplanktonic nitrifiers were observed in deep samples from eastern Lake Erie and Lake Ontario (9 to 24% of total amplicons), compared to 2 to 14% in

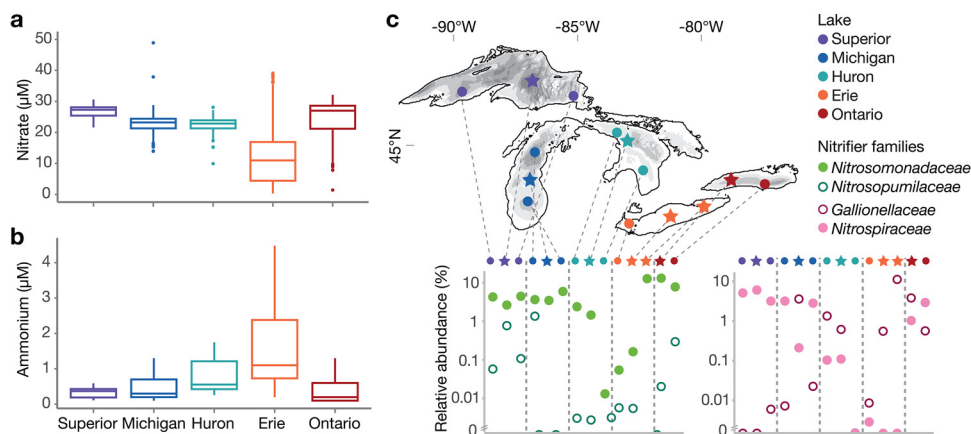


FIG 1 Dissolved inorganic nitrogen availability and distribution of nitrifiers across the Great Lakes. (a) Oxidized nitrogen concentrations. Values include NO_x concentrations from published studies ($n = 128$) (14, 33, 35, 137, 138), U.S. EPA Water Quality Surveys in 2012 and 2013 ($n = 1,626$ from GLENDa database), and this study ($n = 20$). (b) Ammonium concentrations. Values are derived from the literature as described for panel a ($n = 118$) and from this study ($n = 20$). (c) Distribution of nitrifiers across the Great Lakes. Top panel, map of sampling stations; stars indicate stations chosen for metagenome analysis. Bottom panel, relative abundance of ammonia-oxidizing (green) and nitrite-oxidizing (pink) families based on 16S rRNA V4-V5 amplicon sequencing, sampled in the mid-hypolimnion (except western Lake Erie, sampled 1 m from bottom). Data are plotted roughly West to East as indicated on the map.

Lakes Michigan, Huron, and Superior. Lakes Erie and Ontario also have higher cell concentrations and higher surface chlorophyll (see Data Set S1 at <https://doi.org/10.6084/m9.figshare.15130350.v4>). The relative abundance of nitrifiers was negatively correlated with photosynthetically active radiation (PAR; Spearman's $\rho = -0.89$, $P < 2.2e-16$) and reached a maximum below the depth of 1% PAR in each lake, up to 20% of amplicon sequences (see Fig. S1a in the supplemental material). The relative abundances of ammonia- and nitrite-oxidizing taxa were strongly correlated (Spearman's $\rho = 0.918$, $P < 2.2e-16$) (Fig. S1b). Picoplanktonic nitrifiers were rare ($<0.1\%$ relative abundance) in bottom water samples from the southern basin of Lake Huron (HU15M) and the western basin of Lake Erie (ER91M); these two stations are the shallowest in our data set and have relatively high light penetration to the bottom ($\sim 1\%$ PAR). Chlorophyll *a* concentration was also negatively correlated with the relative abundance of nitrifiers (Spearman's $\rho = -0.677$, $P < 1.7e-7$) (Fig. S1c). These findings are consistent with previous work demonstrating photoinhibition of nitrification (37–40), as well as potential competition with phototrophs for ammonium (41).

The taxonomic assemblage of nitrifiers differed across lakes and even among stations within a lake (Fig. 1; see Data Set S1 at <https://doi.org/10.6084/m9.figshare.15130350.v4>), in association with variable productivity and nitrogen availability. Surface ammonium is typically below 300 nM except in Lake Erie, where it is several-fold higher and spatially variable; nitrate, on the other hand, is very high across the lakes but lowest in Erie due to biological uptake (42, 43). Few measurements of urea exist, but it can exceed ammonium (44) (see Data Set S2 at <https://doi.org/10.6084/m9.figshare.15130350.v4>). AOB (*Nitrosomonadaceae*) were observed across all lakes. In contrast, AOA (*Nitrosopumilaceae*) sequences exceeded 0.5% relative abundance only at the three deepest stations (SU08M, MI41M, ON55M), where the ratio of AOB to AOA ranged from 10:1 to 1:3. We found pronounced shifts in the dominant NOB across stations (Fig. 1), and all stations except those in Lake Ontario showed strong dominance (greater than 10-fold) of either “*Ca. Nitrotoga*” (family *Gallionellaceae*) or *Nitrospira*. *Nitrospira* was the only nitrite oxidizer detected in Lake Superior and the dominant nitrite oxidizer in parts of Lake Michigan (MI41M, MI18M). In contrast, “*Ca. Nitrotoga*” was the only nitrite oxidizer observed in Lake Erie and the dominant nitrite oxidizer in Lake Huron and at the shallowest station in Lake Michigan (MI27M). Within each taxon, a single 16S rRNA oligotype dominated the AOA, “*Ca. Nitrotoga*,”

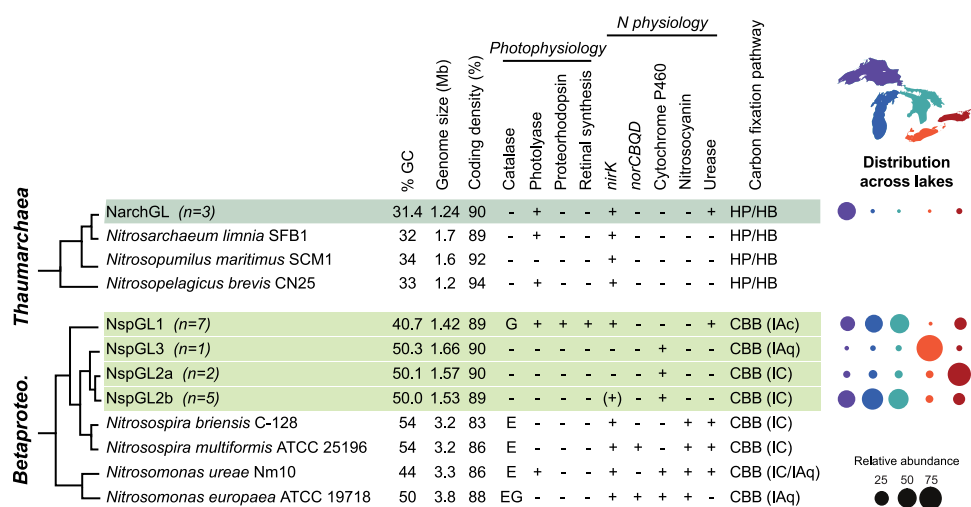


FIG 2 Genome properties and cross-lake distribution of ammonia-oxidizing organisms, showing both *Archaea* (top) and *Betaproteobacteria* (bottom). Rows highlighted in green represent clusters of genomes reconstructed from the Great Lakes, and median values are shown for genome size, GC content, and coding density. For catalase, “E” indicates monofunctional catalase *katE* and “G” indicates bifunctional catalase-peroxidase *katG*. For carbon fixation, the RuBisCO type is shown in parentheses (61). HP/HB, 3-hydroxypropionate/4-hydroxybutyrate cycle; CBB, Calvin-Benson-Bassham cycle. A bubble plot shows the composition of ammonia oxidizers in hypolimnion samples, using MAGs as probes to recruit metagenomic reads (values sum to 100% for each lake column). Genes identified in only a subset of genomes are indicated by (+).

and *Nitrosospira*, while several oligotypes of *Nitrosomonadaceae* shifted abundance across samples (Fig. S2), consistent with ecotypic diversity as discussed below.

Ecotypic variation in abundant streamlined *Nitrosospira*. We reconstructed 15 genomes of the AOB *Nitrosospira*, substantially expanding genome descriptions for this genus (45–47). Based on a phylogenomic tree, free-living Great Lakes *Nitrosospira* falls into two major clades, both of which are distinct from published species; each of these clades also includes metagenome-assembled genomes (MAGs) recovered from Lakes Biwa and Baikal, suggesting novel globally distributed freshwater lineages (Fig. S3). One clade, which we call NspGL1, has a highly reduced genome (median, 1.42 Mb) and low G+C content (40.7%). The second clade was resolved into three subclades (denoted NspGL2a, -2b, and -3) (Fig. S3) based on phylogeny and average nucleotide identity (ANI), all with small genome sizes of 1.45 to 1.68 Mb and 50% G+C content (Fig. 2; and see Data Set S3 at <https://doi.org/10.6084/m9.figshare.15130350.v4>). Compared to 86 reference *Nitrosomonadaceae* genomes, Great Lakes *Nitrosospira* genomes are not only smaller (median estimated complete genome size for the reference = 3.21 Mb, for GL = 1.45 Mb) (Table 1) but also have shorter intergenic spacers, fewer paralogs, fewer pseudogenes, and fewer sigma factors (Table 1 and Fig. S4; see Data Set S4 at <https://doi.org/10.6084/m9.figshare.15130350.v4>), consistent with genome streamlining to reduce resource demands (48). Based on short read mapping, these subclades are ecologically distinct: NspGL1 and NspGL2b—with the smallest genomes—are the dominant AOB in the upper oligotrophic lakes, while NspGL2a is abundant only in Lake Ontario and NspGL3 is abundant only in Lake Erie (Fig. 2). Hereafter, we refer to these subclades as ecotypes due to their phylogenetic and ecological divergence.

We next compared gene content between our Great Lakes *Nitrosospira* and 86 *Nitrosomonadaceae* reference genomes (see Data Set S5 at <https://doi.org/10.6084/m9.figshare.15130350.v4>). On average, Great Lakes *Nitrosospira* genomes encode far fewer two-component signal transduction systems (NspGL = 5 to 8, mean reference = 19), transposases (NspGL = 0 to 7, mean reference = 39), motility genes (NspGL = 0 to 4, mean reference = 52), pilus and secretion genes (NspGL = 2 to 9, mean reference = 27), and defense-related genes (NspGL = 4 to 11, mean reference = 39). They also lack functions related to biofilm formation such as polysaccharide matrix production (e.g., *pel*

TABLE 1 Evidence for genome streamlining in nitrifiers from the Laurentian Great Lakes^a

| Taxonomic group | Genome feature | Median value | | | | No. of genomes | |
|--------------------------------|------------------------------|--------------|-----------|-------------|---------|----------------|----|
| | | Reference | GL | W-statistic | P value | Reference | GL |
| <i>Nitrosomonadaceae</i> (AOB) | Coding fraction | 0.856 | 0.892 | 34 | 5.5E-09 | 86 | 15 |
| | Estimated complete size (bp) | 3,210,560 | 1,450,843 | 1,285 | 1.0E-09 | 86 | 15 |
| | Median i.g. (bp) | 114 | 80 | 1,254 | 6.3E-09 | 86 | 15 |
| | Paralogs | 123 | 29 | 1,189 | 2.1E-07 | 86 | 15 |
| | Pseudogenes | 101 | 11 | 1,215 | 9.0E-10 | 81 | 15 |
| | Sigma factors | 8 | 4 | 1,275 | 1.2E-09 | 86 | 15 |
| <i>Nitrospira</i> (NOB) | Coding fraction | 0.876 | 0.894 | 64 | 0.0074 | 64 | 6 |
| | Estimated complete size | 3,790,956 | 1,828,031 | 373 | 1.5E-04 | 64 | 6 |
| | Median i.g. | 90 | 78 | 299 | 0.026 | 64 | 6 |
| | Paralogs | 212 | 49 | 376 | 1.2E-04 | 64 | 6 |
| | Pseudogenes | 69 | 9 | 182 | 2.9E-04 | 31 | 6 |
| | Sigma factors | 13 | 5 | 379 | 8.3E-05 | 64 | 6 |
| "Ca. Nitrotoga" (NOB) | Coding fraction | 0.857 | 0.910 | 0 | 0.0080 | 5 | 6 |
| | Estimated complete size | 2,858,108 | 1,441,179 | 30 | 0.0081 | 5 | 6 |
| | Median i.g. | 122 | 72 | 30 | 0.0080 | 5 | 6 |
| | Paralogs | 93 | 23 | 30 | 0.0081 | 5 | 6 |
| | Pseudogenes | 18 | 8 | 6 | NS | 1 | 6 |
| | Sigma factors | 8 | 4 | 30 | 0.0054 | 5 | 6 |
| <i>Nitrosopumilaceae</i> (AOA) | Coding fraction | 0.900 | 0.898 | 102 | NS | 62 | 3 |
| | Estimated complete size | 1,398,741 | 1,242,579 | 153 | NS | 62 | 3 |
| | Median i.g. | 61 | 66 | 60 | NS | 62 | 3 |
| | Paralogs | 85 | 38 | 175 | 0.011 | 62 | 3 |
| | Pseudogenes | 22 | 13 | 82 | NS | 34 | 3 |

^aGenome features were compared between Great Lakes MAGs and reference genomes by using a two-sided Wilcoxon/Mann-Whitney test. NS, not significant at 0.05 level. Only genomes with >70% completion and <10% contamination are included. W, Wilcoxon test statistic i.g., intergenic spacers.

genes) and extracellular protein targeting (exosortase and PEP-CTERM motifs). Our 15 new *Nitrosospira* MAGs have high estimated completion (median, 98.6%), and therefore it is unlikely that these gene absences can be entirely attributed to incomplete assemblies. This overall picture of gene content in Great Lakes AOB contrasts with that of *Nitrosospira* isolates from soil (45, 47) and even of oligotrophic *Nitrosomonas* isolates (49) and is consistent with a passive planktonic lifestyle in extremely low-nutrient systems.

We next compared metabolic potential among Great Lakes AOB ecotypes to understand their ecological preferences for upper lakes (NspGL1, NspGL2b), Lake Ontario (NspGL2a), and Lake Erie (NspGL3). Surprisingly, all seven NspGL1 MAGs encode proteorhodopsin, a light-driven proton pump that supports bacterial energy production (50, 51). They also carry the genes necessary to synthesize its chromophore retinal, including those encoding 15,15'-beta-carotene dioxygenase (*blh*), lycopene cyclase (*crtY*), phytoene dehydrogenase (*crtI*), phytoene synthase (*crtB*), and GGPP synthase (*crtE*) (52, 53) (Fig. 3a). We also identified proteorhodopsin in a single-cell amplified genome representing NspGL1 from Lake Michigan (Fig. 3a) (99.8% ANI with NspGL1 MAGs), demonstrating that it is not an artifact of metagenome assembly. To our knowledge, this is the first example of a nitrifier with proteorhodopsin. All NspGL1 proteorhodopsins share residues H95, D127, and E138 along with a short beta-turn (G111-P116) between helices B and C, which are characteristic features of proteorhodopsin as distinct from sensory and other rhodopsins (54), and the presence of leucine at position 135 suggests green light tuning (55) (Fig. 3b). All of the genes in this module have highest similarity to homologs from *Polynucleobacter* but are flanked by *Nitrosomonadaceae*-like genes, suggesting recent horizontal gene transfer (Fig. 3a). The predicted NspGL1 proteorhodopsins cluster with *Polynucleobacter*, *Methylopumilus*, and other freshwater *Betaproteobacteria* in supercluster III as defined by MicRhoDE (56) (Fig. 3c). We compared the homologous genome regions in two highly similar MAGs from Lakes Biwa and Baikal (Fig. S5); these contigs lack the proteorhodopsin module but appear

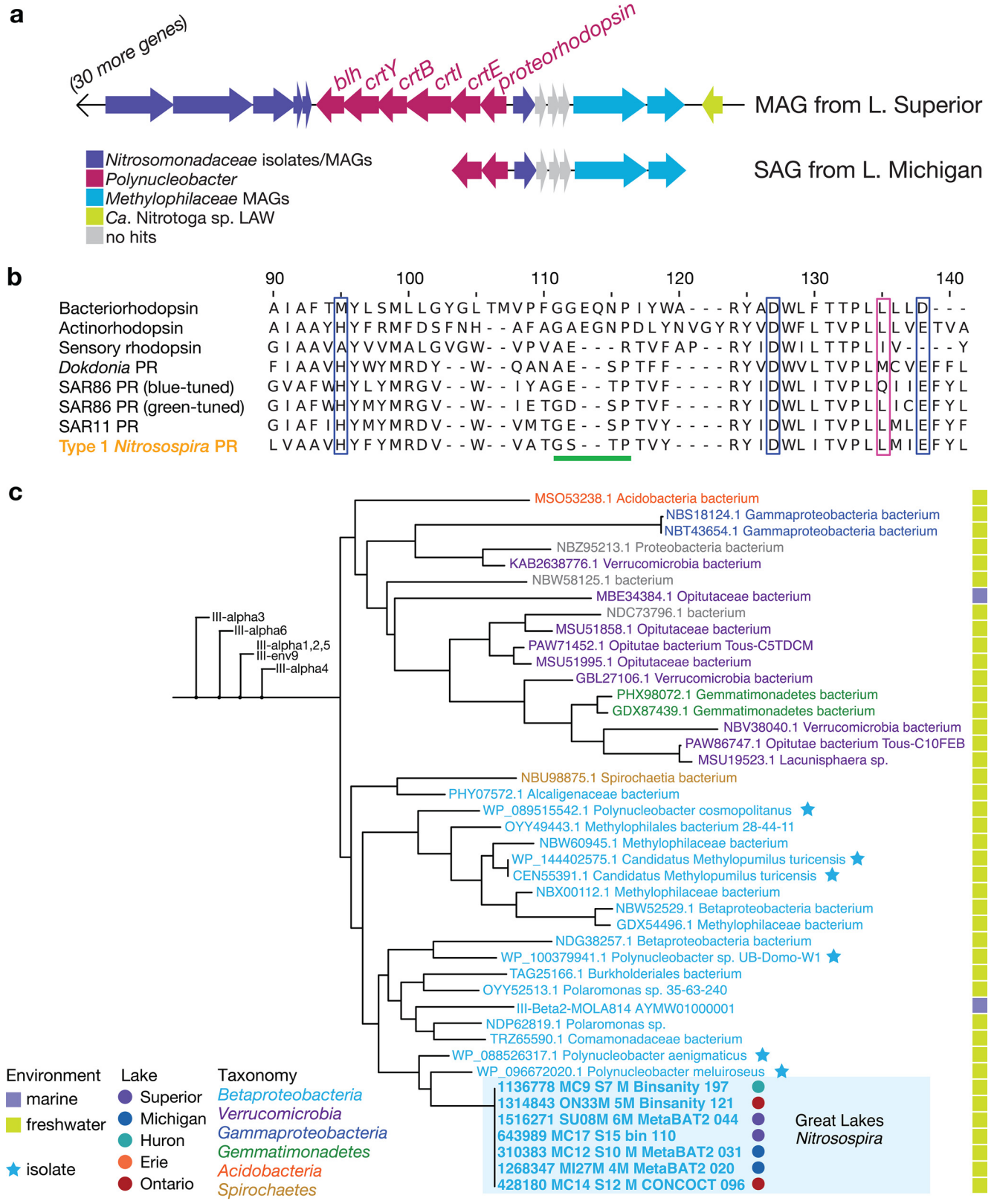


FIG 3 Evidence for proteorhodopsin (PR) in *Nitrosospira* from the Great Lakes. (a) Gene neighborhood surrounding PR in *Nitrosospira* MAG MC17_S15_bin_110 and SAG 207399. Genes are colored according to the best BLAST hit taxonomy in the NCBI nr database. (b) Alignment of predicted *Nitrosospira* PR with reference sequences. Diagnostic features are highlighted (54, 55): blue boxes, diagnostic residues for PR; pink box, residue indicative of blue or green tuning; green underlining, shorter beta-sheet region in PR. Sequence accession numbers: bacteriorhodopsin, P02945; actinorhodopsin,

(Continued on next page)

to flank a variable region where the contig assembly ends. A proteorhodopsin photo-system could support survival of NspGL1 in the presence of sunlight, which has been shown to inhibit ammonia oxidation (37, 57). In the upper lakes where NspGL1 is abundant, light penetration is high well below the thermocline in stratified periods (58), and deep-water taxa are seasonally advected to the surface by water column mixing (34). In addition to proteorhodopsin, NspGL1—but not the other three ecotypes of Great Lakes *Nitrosospira*—encodes a class I cyclopyrimidine dimer photolyase, which uses light energy to repair UV-induced DNA damage, and carries the catalase-peroxidase *katG*, suggesting that the NspGL1 ecotype is adapted to relatively shallow depths in the water column (Fig. 2).

Great Lakes *Nitrosospira* genomes carry a reduced, ecotype-specific complement of nitrogen metabolism genes compared to reference AOB (Fig. 2) (gene absences were verified as described in Materials and Methods). All are presumed to have the core ammonia oxidation enzymes ammonia monooxygenase and hydroxylamine dehydrogenase; these genes were assembled and binned as expected in some MAGs and were manually identified on short unbinned contigs in other cases (see Data Set S6 at <https://doi.org/10.6084/m9.figshare.15130350.v4>) (see Materials and Methods). Surprisingly, all Great Lakes *Nitrosospira* MAGs lack the copper protein nitrosocyanin, whose precise function is unknown but so far has been found in all described AOB except one member of the *Nitrosomonas oligotropha* clade (49). Based on the expanded set of genomes analyzed here, the lack of nitrosocyanin likely extends beyond the Great Lakes MAGs to closely related freshwater and marine strains, along with additional members of the *N. oligotropha* clade (Fig. S3); its absence may be related to the divergence of these clades. Only NspGL1 and NspGL2b encode NO-forming nitrite reductase (NirK), which confers nitrite tolerance (59); this result is surprising, given that these two clades dominate the upper lakes, where productivity and reduced N are lowest. None of the Great Lakes ecotypes encode NO reductase (NorCBQD), and NspGL1 lacks cytochrome P460 family proteins, both of which are common in AOB and implicated in nitrogen oxide metabolism (18, 49). Nitrogen acquisition is also distinct among Great Lakes AOB: NspGL1 lacks an apparent ammonium transporter but has urease structural and accessory genes (*ureABCEFG*) and a high-affinity urea transporter (*urtABCDE*). Further, all Great Lakes ecotypes encode a high-affinity amino acid transporter (*livFGHM*); these genes are rare (<5%) in reference genomes and could supply reduced nitrogen and/or organic carbon. Finally, NspGL1 and NspGL3 have genes for producing cyanophycin, an intracellular storage compound for nitrogen (47, 60). Together, the distinctive gene complements present in Great Lakes *Nitrosospira* illustrate the variability and adaptability of AOB gene content, even across a connected freshwater habitat.

As with nitrogen metabolism, carbon metabolism is also distinct between Great Lakes and reference AOB and among Great Lakes ecotypes (Fig. 2). Unlike most reference AOB, Great Lakes *Nitrosospira* AOB lack two key enzymes of the oxidative pentose phosphate pathway, glucose-6-phosphate dehydrogenase and 6-phosphogluconate dehydrogenase. All ecotypes except Erie-specific NspGL3 also lack genes for glycogen synthesis and degradation, suggesting that they are unable to store and access this carbon reserve. The key enzyme for carbon fixation, RuBisCO, has evolved several kinetically distinct forms whose distribution likely reflects ecological pressures (61). NspGL1 and NspGL3 both contain form IA RuBisCO, while NspGL2a and NspGL2b contain form IC RuBisCO (Fig. 2) (61, 62). NspGL1 genomes also possess an alpha carboxysome-like *cso* operon, similar to *Nitrosomonas eutropha* C91 (62), though our draft assembly lacks the expected carbonic anhydrase (*csoS3/csoSCA*). Carboxysome-associated RuBisCO may allow NspGL1, the ecotype most strongly adapted to energy and nutrient limitation, to more efficiently fix CO₂ by minimizing the waste-

FIG 3 Legend (Continued)

AOA1D9E0H1; sensory rhodopsin, P42196; *Dokdonia* PR, EAQ40507.1; SAR86 blue-tuned PR, Q4PP54; SAR86 green-tuned PR, Q9F7P4; SAR11 PR, A6YQL7. (c) Phylogenetic tree showing close relatives of *Nitrosospira* PR within supercluster III, as defined by MicRhoDE database (56). Neighboring clusters have been collapsed for clarity.

ful oxygenation reaction and reducing the cellular nitrogen allocation to RuBisCO (61). The ranges of kinetic properties observed in other autotrophs for form IAq (found in NspGL3) and form IC (found in NspGL2a and NspGL2b) overlap, and therefore more work is needed to understand the fitness advantages, if any, that this RuBisCO diversity confers on Great Lakes nitrifiers.

Streamlined freshwater *Thaumarchaeota*. We reconstructed three similar genomes (>99% ANI) of *Nitrosarchaeum* (NarchGL) (Fig. S6) from three separate samples (two from Lake Superior, one from Lake Ontario), consistent with our low observed 16S rRNA diversity for *Thaumarchaeota*. These NarchGL genomes are very similar (~99% ANI) to two genomes reconstructed from Lake Baikal, located thousands of kilometers away (63). Their next closest relatives are also from freshwater environments, and phylogenetic clustering suggests that salinity is an important driver of divergence throughout the *Nitrosopumilaceae* (Fig. S6). As a group, the *Thaumarchaeota* tend to have smaller genomes, lower G+C content, higher coding density, and fewer paralogs and pseudogenes than nitrifying bacterial taxa; even within this group, NarchGL genomes fall below the 30th percentile in size and have significantly fewer paralogs than average (Table 1 and Fig. S4) (see Data Set S4 at <https://doi.org/10.6084/m9.figshare.15130350.v4>). By using our reconstructed genomes as probes for metagenomic read recruitment, NarchGL were detected in Lakes Superior, Michigan, and Ontario; they represented roughly one-third of ammonia oxidizers in the mid-hypolimnion of station SU08M (Fig. 2).

NarchGL genomes share nearly 90% of their predicted proteins with close relatives, including “*Ca. Nitrosarchaeum limnia*”; at the same time, they show distinctive patterns in gene content that pinpoint the key selective pressures of deep lakes (Fig. 2) (see Data Set S5 at <https://doi.org/10.6084/m9.figshare.15130350.v4>). All three NarchGL genomes encode urease and a urea transporter, implicating urea as a vital source of nitrogen for energy and/or biosynthesis. Consistent with phosphorus scarcity in much of the Great Lakes (64), NarchGL genomes encode high-affinity transport systems for phosphate and potentially phosphonates, though we did not identify a phosphonate lyase. In addition to both CRISPR/Cas enzymes cas1 and cas4, NarchGL genomes contain several phage proteins, suggesting that viral infection and integration events may be common. DNA photolyases, which have been found in epipelagic clades of marine *Thaumarchaeota* (7), are present in all genomes representing low-salinity *Nitrosarchaeum* including NarchGL, consistent with sunlight exposure due to high water clarity (58) and/or annual mixing in the Great Lakes. NarchGL also lack the common tRNA modification 4-thiouridylation (indicated by KEGG Orthology K04487 and Pfam PF02568-PF18297 [65]); we propose that the absence of this modification, which is susceptible to near-UV radiation (65), is also related to sunlight exposure.

Genomes of NarchGL reveal striking reductions in environmental sensing, response, and regulatory functions, relative to most other *Nitrosarchaeum* and *Nitrosopumilaceae*. NarchGL genomes encode 9 to 12 domains representing the general archaeal transcription factors TATA binding protein (TBP; Pfam PF00352) and transcription factor B (TFB; Pfam PF00382 and PF08271), compared to 21 in “*Ca. Nitrosarchaeum limnia*.” NarchGL genomes lack common domains found in two-component systems that transmit environmental signals to control gene expression or protein activity (Pfam domains PF02743, PF00672, PF00512, and PF00072; NarchGL = 0, “*Ca. N. limnia*” = 53 to 54 copies per genome). Further, they are depleted in ArsR family transcription factors (PF01022; 0 copies in NarchGL versus 2 to 3 in “*Ca. N. limnia*”), P-II proteins for regulation of nitrogen metabolism (PF00543; 1 copy per NarchGL genome versus 5 in “*Ca. N. limnia*”), and other potential regulatory domains (CBS PF00571, 5 in NarchGL versus 18 to 19 in “*Ca. N. limnia*”; USP PF00582, 1 in NarchGL versus 15 in “*Ca. N. limnia*”). This extremely limited regulatory capacity in NarchGL stands in sharp contrast to closely related “*Ca. N. limnia*” and instead parallels the oceanic minimalist “*Ca. Nitrosopelagicus brevis*” (9).

Expanded diversity of “*Ca. Nitrotoga*” with reduced genomes. Despite the broad distribution of “*Ca. Nitrotoga*” in freshwater systems and beyond, only six genomes are available, derived from rivers heavily impacted by urban and agricultural influence, a

wastewater treatment plant, and coastal sediment (27, 66, 67). Hence, the metabolic and phylogenetic diversity of this group is virtually unexplored. We reconstructed six new MAGs of “*Ca. Nitrotoga*,” which form two clusters with >99% ANI within each cluster and ~97% ANI between clusters (NtogaGL1a and NtogaGL1b) (Fig. S3). These new “*Ca. Nitrotoga*” MAGs are far smaller than published genomes (median, GL = 1.44 Mb, reference = 2.61 to 2.98 Mb) and have shorter intergenic regions, fewer sigma factors, and fewer paralogs (Table 1 and Fig. S4) (see Data Set S4 at <https://doi.org/10.6084/m9.figshare.15130350.v4>), consistent with genome streamlining (48). They also have distinctive gene content (see Data Set S5 at <https://doi.org/10.6084/m9.figshare.15130350.v4>). They lack functions such as motility and chemotaxis, pilus biogenesis, and DNA repair (*mutLS*). Great Lakes “*Ca. Nitrotoga*” genomes also encode markedly fewer two-component systems for sensing and responding to environmental cues than four river-derived genomes (NtogaGL = 2 to 6 per genome versus 30 to 35 in four reference genomes; reference strain KNB has 7). Compared to reference genomes, NtogaGL genomes have fewer defense-related genes (restriction-modification, toxin-antitoxin, and CRISPR-Cas systems; mean, NtogaGL = 11 versus 39 for references), and transposases (mean, NtogaGL = 3 versus 19 for references). While incomplete assembly of hypervariable genome regions may explain some of these absences, the overall genome properties are consistent with a relatively stable low-nutrient environment and planktonic lifestyle.

The reduced genomes of NtogaGL1a/b help clarify core features of the genus “*Ca. Nitrotoga*,” along with accessory functions that may enable local adaptation in specific populations. To date, sequenced “*Ca. Nitrotoga*” genomes including NtogaGL1a/b encode similar electron transport pathways, including NADH dehydrogenase-complex I, succinate dehydrogenase-complex II, and alternative complex III, along with high-affinity *cbb3*-type cytochrome oxidases, suggesting adaptation to low-oxygen conditions. They also share the Calvin cycle for carbon fixation, a complete tricarboxylic acid (TCA) cycle, and an evolutionarily distinct nitrite oxidoreductase (NXR) from other NOB (27, 66, 67). All “*Ca. Nitrotoga*” genomes to date also share transporters for amino acids and peptides, potential sources of C and/or N. “*Ca. Nitrotoga*” can also potentially access reduced sulfur compounds for energy via sulfite dehydrogenase, suggesting metabolic flexibility beyond nitrite oxidation.

Beyond these similarities, the small genomes of NtogaGL1a/b are distinct from previously described “*Ca. Nitrotoga*” in many ways. NtogaGL1a/b lack NiFe hydrogenase to use hydrogen as an energy source. They also lack nitrogen metabolism functions, including assimilatory nitrite reductase (*nirBD*) and nitrite reductase to NO (*nirK*). Based on gene content, NtogaGL1a/b appear unable to use hexoses like glucose, since they lack the glycolytic enzyme phosphofructokinase and the Entner-Doudoroff pathway, similar to *Nitrobacter winogradskyi* (68). Consistent with this, they also lack genes for storage and breakdown of glycogen. All but one of the NtogaGL1a/1b genomes encode cyanate lyase (*cynS*), which is found in other NOB but not in “*Ca. Nitrotoga*” to date (25, 69, 70). The *cynS* gene, adjacent to *glnK-amtB* for ammonium sensing and transport, likely functions in N assimilation, as recently described for *Nitrospinae* (71). While cyanase has been shown to mediate reciprocal feeding between some NOB and ammonia oxidizers (25), it remains to be seen whether such an interaction occurs in the free-living (<1.6- μ m) size fraction and dilute environment sampled here. Notably, cyanase from NtogaGL1a/1b, along with predicted *Nitrospira*e proteins from Lake Baikal and soil, form a distinct phylogenetic cluster from most nitrifier cyanase proteins observed to date (Fig. S7).

The two ANI-based clusters we detected, NtogaGL1a and NtogaGL1b, appear to be phylogenetically and ecologically distinct ecotypes. Based on short read mapping, NtogaGL1b dominates Lake Erie, while NtogaGL1a dominates all other “*Ca. Nitrotoga*”-containing samples (Fig. 4). We found several metabolic genes that differentiate the two ecotypes. Lake Erie-specific NtogaGL1b genomes share a region encoding thiosulfate dehydrogenase (*tsdA*), cytochromes, transport of sulfur-containing compounds, lactate

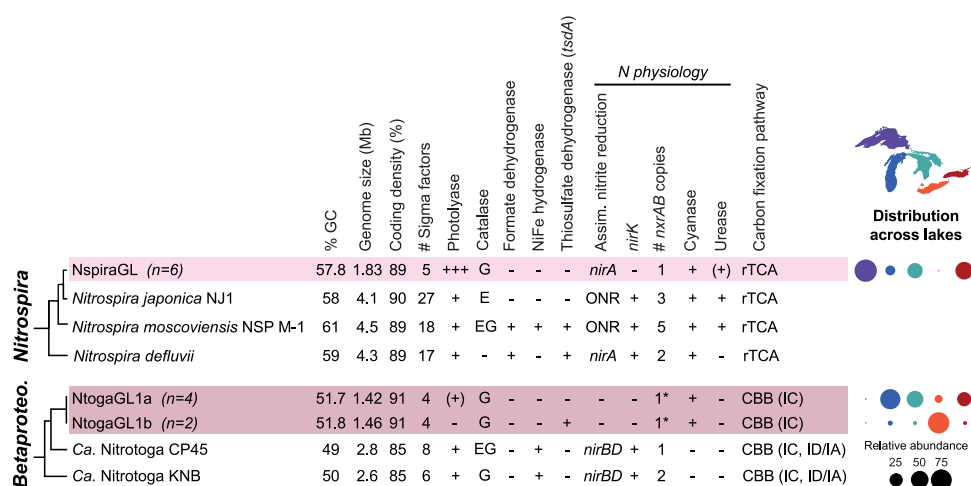


FIG 4 Genome properties and cross-lake distribution of nitrite-oxidizing taxa *Nitrospira* (top) and “*Ca. Nitrotoga*” (*Betaproteobacteria*; bottom). Rows highlighted in pink represent clusters of genomes reconstructed from the Great Lakes, and median values are shown for genome size, GC content, and coding density. rTCA, reductive tricarboxylic acid cycle; CBB, Calvin-Benson-Bassham cycle; ONR, octaheme nitrite reductase. Values in parentheses indicate RuBisCO type (61). A bubble plot shows the composition of NOB per lake based on metagenomic read mapping. Genes identified in only a subset of genomes are indicated by (+). An asterisk indicates that for “*Ca. Nitrotoga*,” one *nrxAB* copy was recovered in genome assemblies, but short read analysis suggests two copies per genome (see Supplemental Text at <https://doi.org/10.6084/m9.figshare.15130350.v4>).

dehydrogenase (*ldh*), a two-component system, and a Crp-family transcription factor (Fig. S8). This region may be involved in oxidizing thiosulfate as an energy source and sensing and responding to redox changes that accompany seasonal hypoxia in Lake Erie. The corresponding region in NtogaGL1a encodes an integrase and photolyase, consistent with greater DNA photodamage in the more transparent waters of Lakes Michigan, Huron, and Ontario, where NtogaGL1a is abundant.

Great Lakes *Nitrospira* genomes reveal adaptations to sunlit oxic environment.

We reconstructed six closely related genomes of *Nitrospira* (~99% ANI) (Fig. S9), representing the predominant NOB throughout Lake Superior and in parts of Lakes Michigan and Ontario (Fig. 4) (see Data Set S1 at <https://doi.org/10.6084/m9.figshare.15130350.v4>). These genomes, which we refer to as NspiraGL, fall within lineage II (Fig. S9), which is broadly distributed across soil, freshwater, and engineered habitats (20); however, genome analyses to date have focused on strains from wastewater and engineered systems, leaving major blind spots. NspiraGL share core features of *Nitrospira* metabolism, including a periplasm-facing NXR that is advantageous under substrate-limiting conditions, multiple cytochrome *bd*-like oxidases, and the reverse TCA cycle for carbon fixation (69). However, as with “*Ca. Nitrotoga*,” the *Nitrospira* genomes we reconstructed in the Great Lakes are markedly smaller than published reference genomes (median for NspiraGL = 1.83 Mb, median for reference = 3.72 Mb), with higher coding density and fewer paralogs, sigma factors, and pseudogenes (Fig. S4) (see Data Set S4 at <https://doi.org/10.6084/m9.figshare.15130350.v4>), consistent with genome streamlining theory (48). Compared to 75 lineage II *Nitrospira* reference genomes, NspiraGL have reduced capacity for environmental sensing (two-component systems: NspiraGL = 7, mean reference = 26), transport (NspiraGL = 76 to 83, mean reference = 140), defense (NspiraGL = 7 to 8, mean reference = 26), and transposition (NspiraGL = 0 to 2, mean reference = 15) and lack pilus or flagellar motility (see Data Set S5 at <https://doi.org/10.6084/m9.figshare.15130350.v4>). NspiraGL encode just five sigma factors, compared to 18 in *Nitrospira moscoviensis*. Further, NspiraGL genomes encode a single NXR, while *N. moscoviensis* carries five copies that are differentially regulated (26, 72). NspiraGL also lack the *glnE* gene for glutamine synthetase (GS) adenylyltransferase, suggesting that GS activity is not repressed by this

mechanism. Together, these features suggest limited regulatory and ecological flexibility, consistent with a relatively constant, oligotrophic environment.

Compared to other *Nitrospira*, NspiraGL exhibit limited energetic flexibility but can access diverse nitrogen sources (Fig. 4). We predict that NspiraGL are unable to grow on hydrogen or formate as alternative energy sources (23, 26), as they lack NiFe-hydrogenase and formate dehydrogenase. The glycolysis and oxidative TCA cycles appear to be incomplete, lacking phosphofructokinase and citrate synthase, respectively; this suggests a limited capacity for organic carbon utilization. NspiraGL lack *nirK*, encoding NO-forming nitrite reductase, which is found in a majority of reference genomes. To obtain nitrogen for biosynthesis, NspiraGL encode a high-affinity nitrate/nitrite/cyanate transporter (*nrtABC*), assimilatory nitrite reductase (*nirA*), and cyanase (*cynS*), along with *amt* family ammonium transporter. Although none of the NspiraGL MAGs include urease (*ureCBA*), one does contain urease accessory proteins (*ureEFGD*) and two contain a urea transporter (*urtABCD*), suggesting incomplete assembly of the urea utilization pathway. As with “*Ca. Nitrotoga*,” we suggest that cyanase, along with urease where present, functions in nitrogen assimilation rather than cross-feeding, given the dilute environment and free-living planktonic cells.

Beyond energy, carbon, and nitrogen metabolism, we discovered striking differences between NspiraGL and reference *Nitrospira* related to DNA repair. NspiraGL encode two additional photolyase-related proteins, along with a class I cyclopyrimidine dimer (CPD) photolyase found in most reference *Nitrospira* taxa (Fig. 5). Photolyases use blue light energy to repair DNA lesions caused by UV radiation (73). The two additional genes in NspiraGL are adjacent and share best hits with *Betaproteobacteria*, suggesting recent horizontal transfer (Fig. S10). One likely encodes an FeS-BCP photolyase, which repairs (6-4) dipyrimidine lesions (74, 75). The other shares an FAD-binding domain with photolyases, but the C-terminal region has no recognizable domains (Fig. 5). This protein is widespread in aquatic bacteria and has not been functionally characterized, though an actinobacterial homolog was suggested to be involved in light sensing and regulation (76). Beyond photolyases, NspiraGL also encode uracil-DNA glycosylase (UNG), which removes misincorporated uracil from DNA. Uracil results from deamination of cytosine, which can occur spontaneously or be induced by NO (77). In addition to the photolyases and UNG that repair DNA lesions, NspiraGL encode translesion DNA polymerase V (*umuCD*), which enables replication to proceed past lesions. Together, these genes indicate that members of *Nitrospira* in the Great Lakes experience significant DNA damage, including UV-induced damage that also requires light for the repair process, in hypolimnion waters with high transparency (58) and/or during seasonal mixing.

Other major differences between NspiraGL and reference *Nitrospira* genomes are related to reactive oxygen species (ROS). Surprisingly, despite their oxic habitat, NspiraGL lack superoxide dismutase (SOD), monofunctional catalase (*katE*), and bacterioferritin, which limits the Fenton reaction by sequestering free iron. However, all six NspiraGL MAGs, but few reference genomes (7% of 75), have recently acquired bifunctional catalase-peroxidase *katG*; interestingly, we also observed *katG* in Great Lakes “*Ca. Nitrotoga*” and *Nitrosospira* (Fig. 2 and 4). The absence of SOD suggests that NspiraGL does not produce damaging levels of endogenous superoxide, perhaps because NspiraGL lack the major respiratory and nonrespiratory flavoproteins that produce ROS in other SOD-containing *Nitrospira* taxa (78). Unlike superoxide, H₂O₂ can cross membranes and is known to be produced by both photooxidation of dissolved organic matter and dark heterotrophic activity (79). The lakes where NspiraGL dominate have high water clarity (58) and low productivity and are fully oxic, consistent with abiotic photochemistry as the primary source of exogenous ROS; this stress may have selected for *katG* as a defense. NspiraGL also lack cytochrome *c* peroxidase, which is found in 70 of 75 reference genomes; this protein is proposed to function in anaerobic respiration of H₂O₂ (80), and therefore its absence in NspiraGL is consistent with a constant oxic environment. Together, these results indicate that members of *Nitrospira* in the Great Lakes face distinct ROS pressures that have shaped their gene content.

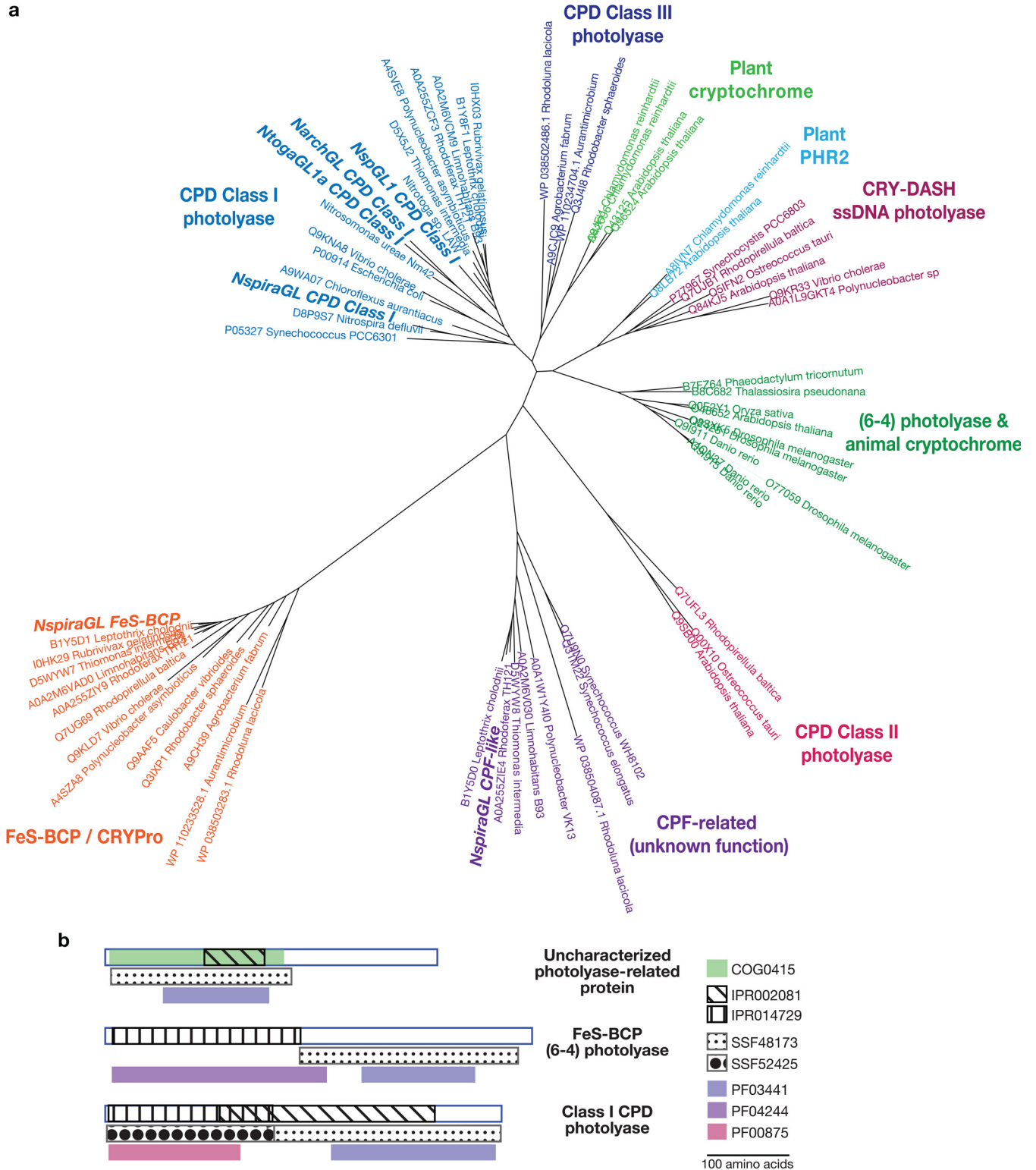


FIG 5 Distinct photolyase proteins in *NspiraGL*. (a) Phylogenetic tree showing families of photolyases. Three families are found in *NspiraGL*: CPD class I photolyase, FeS-BCP/CRYPro family, and an uncharacterized CPF-related family found in diverse *Bacteria*. CPD class I photolyases are also found in other nitrifiers, including “*Ca. Nitrotoga*” *NtogaGL1a*, *Nitrosospira* *NspGL1*, and *Nitrosarchaeum* *NarchGL*. (b) Domain structure of the three photolyase families present in *NspiraGL*.

Conclusions. The Laurentian Great Lakes harbor nitrifiers that are phylogenetically related, but markedly different in genome size and functional capacity, from their well-studied relatives inhabiting wastewater systems, soils, and even other freshwater systems. By examining the entire nitrifier assemblage at once, we detected common features across taxa that illuminate the selective pressures faced by microbes in deep lakes. All the lineages we describe show small genome sizes (1.3 to 1.7 Mb), reduced capacity for environmental sensing and response, and adaptation to a passive (i.e., nonmotile) planktonic lifestyle, features which have not been previously associated with AOB, *Nitrospira*, and “*Ca. Nitrotoga*.” Within the AOB *Nitrospira*, we found ecotypes with a gradient of genome reduction that maps onto their habitats’ trophic gradient: from NspGL1 (1.4 Mb, low GC, upper lakes) to NspGL2b (1.5 Mb, upper lakes) to NspGL2a (1.6 Mb, Lake Ontario) to NspGL3 (1.7 Mb, Lake Erie) (Fig. 2). The thaumarchaeal NarchGL have a markedly reduced regulatory capacity like the open ocean strain *Nitrosopelagicus brevis* (9). The NOB NspiraGL have genomes 50 to 60% smaller than the genomes of described *Nitrospira* taxa and dominate the deeper more oligotrophic basins, while “*Ca. Nitrotoga*” favors shallower, more productive basins. The emergence of Lake Erie-specific ecotypes of both *Nitrospira* (NspGL3) and “*Ca. Nitrotoga*” (NtogaGL1b) demonstrates how distinct this habitat is compared to the other lakes. Importantly, our findings here represent planktonic cells in the smallest size fraction (<1.6 μm); it is likely, especially in Lake Erie, that particle-associated nitrifiers may be abundant and genetically distinct.

Nitrifiers inhabiting the transparent waters of the upper Great Lakes show distinctive adaptations to light, including diverse photolyases, ROS detoxification, and even proteorhodopsin. This discovery is surprising, given that nitrifiers are rare in the surface mixed layer of the Great Lakes (Fig. S1) and that photoinhibition of ammonia oxidation and nitrifier growth is well documented (37, 40, 57). We propose that proteorhodopsin could be used to augment energy metabolism when ammonia oxidation is photoinhibited and/or ammonia oxidation is substrate limited. Water clarity has increased over the past several decades in Lakes Michigan and Huron, now surpassing that of Lake Superior (58). High light penetration along with seasonal mixing likely exposes deep-water cells to damaging levels of light and oxidative stress. Future cultivation and physiological studies should examine photoinhibition and potential phototrophy in Great Lakes nitrifiers.

Our work unveils new clues about the ecological and evolutionary potential of nitrifiers in their natural freshwater habitat. This collective nitrifier diversity undoubtedly influences the cycling of carbon and nitrogen across this ecosystem, and future work will explore the differential contributions to nitrification by the distinct lineages we described here. Understanding what controls the diversity of nitrifiers and other key functional groups, and the consequences of this diversity for biogeochemistry, are essential for forecasting the effects of rapid environmental change across the large lakes of the world (e.g., see reference 81) and predicting impacts on the critical ecosystem services they provide (82).

MATERIALS AND METHODS

Sample collection. Water samples were collected from the Laurentian Great Lakes aboard the R/V *Lake Guardian*, during the biannual Water Quality Surveys conducted by the U.S. EPA Great Lakes National Program Office (83). Station information is provided in Data Set S8 at <https://doi.org/10.6084/m9.figshare.15130350.v4>. Data presented here were collected in April and August 2012. Samples were collected using a conductivity-temperature-depth (CTD) rosette sampler (Sea-Bird Scientific) at the surface (2 m), deep chlorophyll maximum (if present), the mid-hypolimnion (depths ranging from 19 m in Lake Erie to 200 m in Lake Superior) (see Data Set S1 at <https://doi.org/10.6084/m9.figshare.15130350.v4>), and near the bottom of the water column (10 m above the lake bottom at most stations, 1 m above bottom at shallow stations). For each sample, 5 to 8 L of water was prefiltered through a GF/A glass fiber filter (Whatman 1820-047; nominal pore size, 1.6 μm) to exclude eukaryotic phytoplankton and particle-associated microbes, and cells were collected on 0.22- μm Sterivex filters (Millipore SVGP01050). Filters were stored at -80°C . For dissolved nutrient analysis, 0.22- μm filtrate was collected in 125-mL acid-clean high-density polyethylene (HDPE) bottles (Nalgene) and stored at -20°C . Samples for single-cell amplified genomes (SAG) were collected in August 2014. For each sample, 1 mL of raw water was incubated with

100 μ L of glycerol-TE buffer (20 mL 100 \times Tris-EDTA [TE], pH 8, plus 100 mL glycerol plus 60 mL water; final concentrations after sample addition are 10 mM Tris, 1 mM EDTA, 5% glycerol) for 10 min in the dark and then flash frozen in liquid nitrogen and stored at -80°C until processing.

Physicochemical data. CTD profiles, water chemistry, and chlorophyll *a* data were collected by the U.S. EPA according to standard protocols (84) and retrieved from the Great Lakes Environmental Database (<https://cdx.epa.gov/>) for 2012 and 2013. In addition, we measured dissolved nitrogen species from August 2013 samples. Ammonium concentrations were measured using the OPA method in a 96-well plate (85). Nitrate and nitrite concentrations were measured using the Griess reaction method in a 96-well plate (86). Urea concentrations were measured in a 24-well plate using a colorimetric reaction (87).

16S rRNA analysis. The full 16S rRNA amplicon data set was described by Paver and colleagues (34). Here, we focus on data from the V4-V5 region (primers 515F-Y and 926R [88]), collected in 2012 in tandem with metagenome samples from select stations. We classified sequences using the Silva v.132 database (89) and the method of Wang et al. (90) as implemented by mothur (91). Sequences classified to each detected family of nitrifiers (ammonia oxidizer families *Nitrosomonadaceae* and *Nitrosopumilaceae*; nitrite oxidizer families *Gallionellaceae* and *Nitrospiraceae*) with a mothur-assigned confidence score above 90 were delineated into taxonomic units using minimum entropy decomposition with a minimum substantive abundance of 10 (92).

Metagenome and single-cell genome sequencing. One station per lake in Lakes Superior, Michigan, Huron, and Ontario and two stations in Lake Erie were selected for metagenome sequencing. Spring 2012 metagenome samples were collected from the surface, and summer 2012 metagenome samples were collected from the mid-hypolimnion (depths listed in Data Set S1 at <https://doi.org/10.6084/m9.figshare.15130350.v4>). DNA was extracted using a modified phenol-chloroform extraction protocol (34), and libraries were prepared according to the Illumina TruSeq protocol. Samples from spring 2012 were sequenced at the Joint Genome Institute using Illumina HiSeq (2 \times 150 bp). Samples from summer 2012 were sequenced at the University of Chicago Functional Genomics Core Facility using Illumina HiSeq 2500 (2 \times 250 bp).

To confirm the presence of proteorhodopsin, we analyzed a single-cell amplified genome from *Nitrosospira* collected from Lake Michigan and sequenced by the Joint Genome Institute. Quality filtered reads were downloaded from Joint Genome Institute (JGI) IMG/ER and normalized using `bbnorm.sh` with a target of 100 and a mindepth of 2. Normalized reads were assembled using SPAdes 3.1.11 in single-cell mode (93) with flags `-sc` and `-careful`. Resulting scaffolds were annotated identically to MAGs as described below.

Obtaining metagenome-assembled genomes. Raw reads for spring surface samples were quality controlled at the Joint Genome Institute, using `bbduk.sh` for adapter trimming (`ktrim = r`, `minlen = 40`, `minlenfraction = 0.6`, `mink = 11`, `tbo`, `tpe`, `k = 23`, `hdist = 1`, `hdist2 = 1`, `ftm = 5`) and quality filtering (`maq = 8`, `maxns = 1`, `minlen = 40`, `minlenfraction = 0.6`, `k = 27`, `hdist = 1`, `trimq = 12`, `qtrim = rl`). Raw reads for summer hypolimnion samples were adapter trimmed, quality filtered, and interleaved using `bbduk` (parameters: `ktrim = r`, `mink = 8`, `hdist = 2`, `k = 21`, `forcetrimleft = 10`, `forcetrimright = 199`, `minlen = 150`) using BBTools suite version 35.74 (<https://sourceforge.net/projects/bbmap/>). Separate assemblies of quality-filtered reads were carried out for each metagenome using metaSPAdes 3.1.11 `-meta` mode using default *k* sizes of 21, 33, and 55 (94). To enable binning based on sequence coverage, forward and reverse reads were merged using `bbmerge` in BBtools, using `qtrim2 = r trimq = 10,13,16` and `adapter = default`. Merged short reads were then mapped onto each assembly using `bowtie2` 2.2.9 in `-sensitive` mode (95), and this coverage information was used to bin assembled contigs. Binning was performed using MetaBAT2 2.12.1 (96), Binsanity 0.2.6.3 (97), and CONCOCT 1.0.0 (98) using default parameters. The resulting bins were scored, aggregated, and dereplicated using DAS_Tool 1.1.1 (99), followed by manual curation using Anvi'o 4.0 (100). We assessed genome completion and contamination of manually curated bins using CheckM 1.1.0 `lineage_wf` (101), and all new MAGs presented here are greater than 70% complete with less than 10% contamination (see Data Set S3 at <https://doi.org/10.6084/m9.figshare.15130350.v4>). Potential nitrifiers were screened by searching for ammonia monooxygenase, hydroxylamine oxidoreductase, and nitrite oxidoreductase genes within reconstructed genomes using `blastp` 2.5.0 (102). For bins where any of these genes were detected, we identified bacterial single-copy core genes (103) or archaeal single-copy core genes (104) using HMMER (105), as implemented in Anvi'o. Single-copy core genes were queried against proteins predicted from bacterial and archaeal genomes in RefSeq (NCBI) (106), and taxonomic identity of these core genes was ascertained based on a least common ancestor approach using a 0.1% window around the bit score of the best hit using KronaTools 2.7.1 (107). Taxonomic assignment was further validated using GTDB-tk 1.0.0 (108). Grouping of MAGs into clades and subclades based on ANI was carried out using `fastANI` 1.1.0 (109). Genome characteristics for each genome group were calculated as the median of those values for the group. Estimated complete genome size was calculated for MAGs and for references in the pangenome analysis using CheckM (101) completion and contamination, as follows: $\text{estimated} = \text{actual} \times [(1 - \text{contamination})/\text{completion}]$. To quantify the abundance of each clade/ecotype across samples, we used competitive mapping of merged short reads using `bowtie2` in sensitive mode against all nitrifier MAGs, summing up the mapped read counts across all MAGs in a given clade/ecotype and dividing by the total mapped nitrifier reads in a sample; these values are shown in bubble plots (Fig. 2 and 4).

Annotation and gene cluster analysis. Reference genomes were obtained from GenBank (accession numbers listed in Data Set S4 at <https://doi.org/10.6084/m9.figshare.15130350.v4>). The full pangenome analyses included all the genomes listed therein, but we only report results from the subset of genomes most closely related to our MAGs. This subset consists of 86 *Nitrosomonadaceae*, 5 "*Ca. Nitrotoga*," 78 *Nitrosopumilaceae* within *Thaumarchaeota*, and 75 *Nitrosospira* genomes that fall within

lineage II. Reference genomes were treated consistently with GL MAGs, with *de novo* gene calling by prodigal 2.6.3 (110) via Anvi'o. Unless otherwise noted, default settings were used for all software. Genes were annotated using InterProScan 5.30–69.0 (111), GhostKOALA (112), and eggNOG-mapper 1.0.3 against the bactNOG database (113). Gene cluster analysis was carried out using the Anvi'o pangenome pipeline (114), using blastp to determine sequence similarity, ITEP to eliminate weak similarity (115), and MCL to cluster, using a minbit of 0.5, MCL inflation of 2, and minimum gene occurrence of 1 (116). Sigma factors were tallied by identifying gene clusters annotated with the following PFAMs: PF00309, PF03979, PF00140, PF04542, PF04539, PF04545, and PF08281. Pseudogene counts were retrieved where available from NCBI PGAP annotated genomes (117). Paralog counts are reported as the number of gene clusters with more than one gene per genome. Intergenic spacers were calculated using bedtools complementBed function (118). Coding fraction is defined as the summed length of all protein-coding genes divided by the estimated total genome length. Prokka 1.14.5 (119) was used to generate GenBank format files from MAGs and SAGs, and genoPlotR 0.8.9 (120) was used to generate initial gene neighborhood maps.

Gene tree construction. The NspGL1 proteorhodopsin sequence was inserted into the MicRhoDE rhodopsin tree using pplacer (121) through the MicRhoDE Galaxy pipeline (56). We then constructed a more targeted phylogenetic tree using aligned reference sequences of supercluster III from MicRhoDE, filtered to exclude fragments shorter than 220 amino acids. To this alignment, we added NspGL1 sequences using MAFFT 7.310 (122) along with high-similarity sequences from NCBI nr that were not present in MicRhoDE. The tree was inferred using RaxML 8.2.12 with model PROTGAMMALG (123). The tree was visualized in iTOL (124), and more distant clusters were collapsed for clarity.

A cyanase phylogenetic tree was created using sequences drawn from querying NtogaGL cyanase against NCBI nr using blastp, as well as sequences from references 2, 25, and 125. Sequences were aligned using MAFFT (122), and the tree was inferred using RaxML 8.2.12 with model PROTGAMMALG (123). The tree was visualized in iTOL (124), and branches were colored based on the taxonomy of the parent genome.

Photolyase-related proteins in GL MAGs were identified by searching for the following features: KEGG Orthology K01669, NCBI Clusters of Orthologous Genes COG0415, Pfams PF03441, PF00875, PF04244, Superfamilies SSF48173, and SSF52425. Reference proteins ($n = 56$) spanning the previously defined families of photolyases and cryptochromes (126) were obtained from UniProt, along with aquatic bacterial sequences described by Maresca and colleagues (76). The reference sequences were aligned using MAFFT (122), and sequences from GL MAGs were added using the MAFFT –addfragments option. The tree was estimated using IQ-TREE 2 1.6.11 (127) and visualized using iTOL (124).

Phylogenomic tree construction. *Nitrospirae*, *Thaumarchaeota*, *Gallionellaceae*, and *Nitrosomonadaceae* genomes were downloaded from GenBank (NCBI) (128) and included in the phylogenomic trees for their respective family. Phylogenomic analyses were carried out within Anvi'o. Briefly, single-copy core genes were extracted as described above, individually aligned at the protein level using muscle (129), and concatenated for each genome. Concatenated alignments were trimmed using Gblocks 0.91b (130) and analyzed by RAXML 8.2.12 (123) to create a phylogenetic tree using the PROTGAMMALG model and 50 bootstraps. Trees were visualized in iTOL (124).

Proteorhodopsin assembly verification. We used several approaches to validate the presence of proteorhodopsin in assembled *Nitrosospira* genomes, to rule out the possibility of chimeric assemblies from different species. We note that proteorhodopsin-containing contigs were independently assembled and binned together with core *Nitrosospira* contigs from seven different samples (i.e., each sample was assembled and binned separately, rather than coassembled). In five of seven cases, proteorhodopsin and retinal biosynthesis genes were assembled together with core *Nitrosospira* genes on the same contig. To rule out a systematic reproducible error in assembly and/or binning, we compared these seven MAGs to a single-cell amplified *Nitrosospira* genome (SAG) from Lake Michigan, obtained as part of another project with the JGI. This SAG was processed through JGI's standard decontamination pipeline and manually investigated to ensure lack of contamination. We found no evidence of contaminating core genes, as all core genes had best hits to either *Nitrosospira* or more generally *Nitrosomonadaceae* in nr. SAG contigs were matched to homologous contigs from NspGL1 MAGs to determine if any SAG contigs were unique using FastANI 1.1.0 (109) with –visualize flag. All contigs from this *Nitrosospira* SAG were found within an NspGL1 MAG. Bandage 0.8.1 (131) was used to manually inspect the assembly graph around the contig that contained the NspGL1 *Nitrosospira* proteorhodopsin to ensure that the assembled contig did not represent a chimeric contig or inappropriate scaffolding. We verified that a single, unique path exists from the beginning to the end of the NspGL1 contig containing proteorhodopsin (Fig. 3). Further, we verified that consistent coverage across this contig existed by mapping short reads from the original sample using bowtie2 (95) and viewing results using Integrated Genomics Viewer 2.7.0 (132). A closely related assembly of the same genomic region from Lake Biwa did not show evidence of proteorhodopsin; to confirm this difference between the Lake Biwa and Great Lakes MAGs, we mapped reads from Lake Biwa (133) (BioProject PRJDB6644) onto the assembled contig described above using bowtie2 (95). This analysis demonstrated that while a large fraction of the NspGL1 contig in question recruited reads from Lake Biwa at high identity (98 to 99%), starting upstream of proteorhodopsin and retinal biosynthesis, this contig no longer recruited reads from Lake Biwa.

Manual identification of key nitrification genes. Despite recovery of 15 high-completion MAGs in NspGL1/2a/2b/3, many of these MAGs lacked key nitrification genes in *amo* and *hao* operons. This was largely due to the fact that *amo* and *hao* operons were often assembled on small contigs below the minimum size cutoff we imposed for binning contigs. Difficulty in assembling these contigs was likely due in part to the several *amo* and *hao* operons with extremely high identity to one another in each

genome, a phenomenon which has been observed in other *Nitrosospora* genomes (18). Manual assembly graph inspection with Bandage (131) supported this hypothesis, as did assessment of abundance of short reads associated with *amo* operons from NspGL and comparison of abundance of short reads associated with core gene *rpoB* from NspGL, using ROcker (134). Still, an exemplar MAG from at least one representative of each ecotype (NspGL1/2a/2b/3) was found with both *amo* and *hao* operons. Further, manual inspection of unbinned contigs confirmed that *amo* and *hao* operons existed on contigs in every sample from which a MAG for a particular ecotype was recovered. That is, for every time that an NspGL1 MAG was recovered from a sample, we were able to determine that an *amo* and *hao* operon which could be affiliated with NspGL1 existed, even if it was not correctly binned. Affiliation for these unbinned key nitrification genes was carried out by alignment of *amoAB* and *haoAB* sequences to *amoAB* and *haoAB* sequences correctly binned in NspGL ecotypes. This process was also carried out for two NtogaGL1a MAGs for *nrxAB*, which were poorly assembled in those two samples. Data Set S6 (at <https://doi.org/10.6084/m9.figshare.15130350.v4>) summarizes the presence of genes related to nitrification and nitrogen metabolism across all our MAGs.

Verification of gene absences. Metagenome-assembled genomes typically comprise tens or even hundreds of contigs, and this fragmented nature makes it impossible to say with certainty whether a particular gene is truly absent. To substantiate our claims of gene absence based on MAGs, we used several lines of evidence. First, we note that our MAGs have high estimated completion (median of 96.4%, mean of 94.3%), based on the presence of universal core gene markers. Second, for all new lineages described here except NspGL3, we assembled multiple similar MAGs independently from different samples, and we inferred gene absences only if the absence was replicated in multiple assemblies. Together, these two factors provide strong support for cases where a missing gene would be expected to occur in a region of predominantly core genes; however, these factors are less informative for cases where a missing gene might occur in a genomic island, because we have no way of assessing the completion of regions lacking core genes, and islands tend to have systematic poor assemblies across samples. A third line of evidence that we considered is chromosome organization: if a single gene is deleted from an otherwise conserved region of synteny, then this deletion should be apparent in a gene neighborhood diagram (e.g., see Fig. S5, S7, S8, and S10 in the supplemental material). Unfortunately, in many cases, our MAGs are too dissimilar from reference genomes and share little synteny with them, so this approach is not always informative.

We used a fourth approach based on quantitative analysis of short reads to verify gene absences. If a suspected missing gene were actually present in the population, but failed to assemble and/or bin with the rest of the genome, then it should be detectable in the unassembled short reads. The frequency of a gene in the population can be estimated from its abundance in the short reads, compared to the abundance of core marker genes in the short reads. We implemented this approach as follows. We searched unassembled short reads for each gene of interest that we identified as absent from MAGs (e.g., nitrosocyanin) using tblastn. Short reads with significant similarity were then filtered by best-hit taxonomy to the appropriate nitrifier group (i.e., *Nitrosomonadaceae*, "*Ca. Nitrotoga*," *Nitrospira*, *Thaumarchaea*). These filtered short reads were enumerated and length normalized [$1,000 \times (\text{number of short reads} / \text{length of the target gene of interest})$]. The same procedure was repeated for genes expected to be present in every cell (e.g., *amoAB*, *hao*, *nrxAB*, ribosomal protein genes) for comparison. If a putative missing gene (based on MAGs) has near-zero detection in the short reads, we can be confident that the gene is truly missing (or has undetectable sequence similarity, or was so recently acquired from another lineage that its best hit points to a different taxon). In contrast, if a putative missing gene is detected in the short reads, then the gene may be present in genomes related to our MAGs but was unassembled/unbinned, or the gene may be present in another lineage of nitrifiers that is not represented by our MAGs. Short read-based quantification of select genes is presented and described in Data Set S7 and Supplemental Text (both at <https://doi.org/10.6084/m9.figshare.15130350.v4>).

Statistical analysis and plots. All statistical comparisons were carried out in R version 3.5.3 (135), and plots were generated using ggplot2 3.2.0 (136). Code and data files are available at bitbucket.org/greatlakes/gl_nitrifiers.

Data availability. The metagenome-assembled genomes presented here are available via NCBI BioProject [PRJNA636190](https://www.ncbi.nlm.nih.gov/bioproject/PRJNA636190). 16S rRNA data are available at NCBI BioProject [PRJNA591360](https://www.ncbi.nlm.nih.gov/bioproject/PRJNA591360). Metagenomes sequenced by JGI are available at <https://genome.jgi.doe.gov> under project ID 1045056, 1045059, 1045062, 1045065, 1045068, and 1045071. The single-cell amplified genome is available at <http://img.jgi.doe.gov> under IMG Genome ID 3300033241. Raw reads are available in NCBI SRA ([SRR14240538](https://www.ncbi.nlm.nih.gov/sra/SRR14240538)–[SRR14240543](https://www.ncbi.nlm.nih.gov/sra/SRR14240543)) or through JGI with the project IDs listed above.

SUPPLEMENTAL MATERIAL

Supplemental material is available online only.

FIG S1, EPS file, 0.8 MB.

FIG S2, EPS file, 1.1 MB.

FIG S3, EPS file, 1 MB.

FIG S4, EPS file, 1.3 MB.

FIG S5, EPS file, 0.8 MB.

FIG S6, EPS file, 1.2 MB.

FIG S7, EPS file, 1.3 MB.

FIG S8, EPS file, 0.7 MB.

FIG S9, EPS file, 1.2 MB.

FIG S10, EPS file, 0.6 MB.

ACKNOWLEDGMENTS

The work conducted by the U.S. Department of Energy Joint Genome Institute, a DOE Office of Science User Facility, was supported under contract no. DE-AC02-05CH11231. Sequencing support was provided by the DOE JGI Community Sequencing Program (CSP no. 1565 and 503460). Funding for this work was provided by Illinois-Indiana Sea Grant (grant no. NA14OAR170095), the UChicago Women's Board, and the National Science Foundation (OCE-1830011 to M.L.C.).

Computational resources were provided by the UChicago Research Computing Center. We thank the science staff in the Great Lakes National Program Office of the U.S. EPA, and the captain and crew of the R/V *Lake Guardian*, for facilitating sample collection. We thank members of the Coleman and Waldbauer labs for assistance with sample collection and processing and for discussion and comments on the manuscript.

We declare no competing interests.

REFERENCES

- Canfield DE, Glazer AN, Falkowski PG. 2010. The evolution and future of Earth's nitrogen cycle. *Science* 330:192–196. <https://doi.org/10.1126/science.1186120>.
- Pachiadaki MG, Sintez E, Bergauer K, Brown JM, Record NR, Swan BK, Mathyer ME, Hallam SJ, Lopez-García P, Takaki Y, Nunoura T, Woyke T, Herndl GJ, Stepanauskas R. 2017. Major role of nitrite-oxidizing bacteria in dark ocean carbon fixation. *Science* 358:1046–1051. <https://doi.org/10.1126/science.aan8260>.
- Reinthal T, van Aken HM, Herndl GJ. 2010. Major contribution of autotrophy to microbial carbon cycling in the deep North Atlantic's interior. *Deep Sea Res II* 57:1572–1580. <https://doi.org/10.1016/j.dsr2.2010.02.023>.
- Swan BK, Martínez-García M, Preston CM, Sczyrba A, Woyke T, Lamy D, Reinthal T, Poulton NJ, Masland EDP, Gomez ML, Sieracki ME, DeLong EF, Herndl GJ, Stepanauskas R. 2011. Potential for chemolithoautotrophy among ubiquitous bacteria lineages in the dark ocean. *Science* 333:1296–1300. <https://doi.org/10.1126/science.1203690>.
- Baltar F, Herndl GJ. 2019. Ideas and perspectives: is dark carbon fixation relevant for oceanic primary production estimates? *Biogeosciences* 16:3793–3799. <https://doi.org/10.5194/bg-16-3793-2019>.
- Callieri C, Coci M, Eckert EM, Salcher MM, Bertoni R. 2014. Archaea and Bacteria in deep lake hypolimnion: in situ dark inorganic carbon uptake. *J Limnol* 73:31–38. <https://doi.org/10.4081/jlimnol.2014.937>.
- Santoro AE, Richter RA, Dupont CL. 2019. Planktonic marine Archaea. *Annu Rev Mar Sci* 11:131–158. <https://doi.org/10.1146/annurev-marine-121916-063141>.
- Schleper C. 2010. Ammonia oxidation: different niches for bacteria and archaea? *ISME J* 4:1092–1094. <https://doi.org/10.1038/ismej.2010.111>.
- Santoro AE, Dupont CL, Richter RA, Craig MT, Carini P, McIlvin MR, Yang Y, Orsi WD, Moran DM, Saito MA. 2015. Genomic and proteomic characterization of “*Candidatus Nitrosopelagicus brevis*”: an ammonia-oxidizing archaeon from the open ocean. *Proc Natl Acad Sci U S A* 112:1173–1178. <https://doi.org/10.1073/pnas.1416223112>.
- Auguet J-C, Triadó-Margarit X, Nomokonova N, Camarero L, Casamayor EO. 2012. Vertical segregation and phylogenetic characterization of ammonia-oxidizing Archaea in a deep oligotrophic lake. *ISME J* 6:1786–1797. <https://doi.org/10.1038/ismej.2012.33>.
- Herber J, Klotz F, Frommeyer B, Weis S, Straile D, Kolar A, Sikorski J, Egert M, Dannenmann M, Pester M. 2020. A single Thaumarchaeon drives nitrification in deep oligotrophic Lake Constance. *Environ Microbiol* 22:212–228. <https://doi.org/10.1111/1462-2920.14840>.
- Urbach E, Vergin KL, Young L, Morse A, Larson GL, Giovannoni SJ. 2001. Unusual bacterioplankton community structure in ultra-oligotrophic Crater Lake. *Limnol Oceanogr* 46:557–572. <https://doi.org/10.4319/lo.2001.46.3.0557>.
- Okazaki Y, Fujinaga S, Tanaka A, Kohzu A, Oyagi H, Nakano S. 2017. Ubiquity and quantitative significance of bacterioplankton lineages inhabiting the oxygenated hypolimnion of deep freshwater lakes. *ISME J* 11:2279–2293. <https://doi.org/10.1038/ismej.2017.89>.
- Mukherjee M, Ray A, Post AF, McKay RM, Bullerjahn GS. 2016. Identification, enumeration and diversity of nitrifying planktonic archaea and bacteria in trophic end members of the Laurentian Great Lakes. *J Great Lakes Res* 42:39–49. <https://doi.org/10.1016/j.jglr.2015.11.007>.
- Hugoni M, Etien S, Bourges A, Lepère C, Domaizon I, Mallet C, Bronner G, Debroas D, Mary I. 2013. Dynamics of ammonia-oxidizing Archaea and Bacteria in contrasted freshwater ecosystems. *Res Microbiol* 164:360–370. <https://doi.org/10.1016/j.resmic.2013.01.004>.
- Hayden CJ, Beman JM. 2014. High abundances of potentially active ammonia-oxidizing Bacteria and Archaea in oligotrophic, high-altitude lakes of the Sierra Nevada, California, USA. *PLoS One* 9:e111560. <https://doi.org/10.1371/journal.pone.0111560>.
- Bollmann A, Bar-Gilissen M-J, Laanbroek HJ. 2002. Growth at low ammonium concentrations and starvation response as potential factors involved in niche differentiation among ammonia-oxidizing bacteria. *Appl Environ Microbiol* 68:4751–4757. <https://doi.org/10.1128/AEM.68.10.4751-4757.2002>.
- Sedlacek CJ, McGowan B, Suwa Y, Sayavedra-Soto L, Laanbroek HJ, Stein LY, Norton JM, Klotz MG, Bollmann A. 2019. A physiological and genomic comparison of *Nitrosomonas* cluster 6a and 7 ammonia-oxidizing bacteria. *Microb Ecol* 78:985–994. <https://doi.org/10.1007/s00248-019-01378-8>.
- Alonso-Sáez L, Waller AS, Mende DR, Bakker K, Farnelid H, Yager PL, Lovejoy C, Tremblay J-É, Potvin M, Heinrich F, Estrada M, Riemann L, Bork P, Pedrós-Alió C, Bertilsson S. 2012. Role for urea in nitrification by polar marine Archaea. *Proc Natl Acad Sci U S A* 109:17989–17994. <https://doi.org/10.1073/pnas.1201914109>.
- Daims H, Lüscher S, Wagner M. 2016. A new perspective on microbes formerly known as nitrite-oxidizing bacteria. *Trends Microbiol* 24:699–712. <https://doi.org/10.1016/j.tim.2016.05.004>.
- Nowka B, Daims H, Spieck E. 2015. Comparison of oxidation kinetics of nitrite-oxidizing bacteria: nitrite availability as a key factor in niche differentiation. *Appl Environ Microbiol* 81:745–753. <https://doi.org/10.1128/AEM.02734-14>.
- Wegen S, Nowka B, Spieck E. 2019. Low temperature and neutral pH define “*Candidatus Nitrotoga* sp.” as a competitive nitrite oxidizer in co-culture with *Nitrospira defluvi*. *Appl Environ Microbiol* 85:e02569-18. <https://doi.org/10.1128/AEM.02569-18>.
- Koch H, Galushko A, Albertsen M, Schintmeister A, Gruber-Dorninger C, Lucker S, Pelletier E, Le Paslier D, Spieck E, Richter A, Nielsen PH, Wagner M, Daims H. 2014. Growth of nitrite-oxidizing bacteria by aerobic hydrogen oxidation. *Science* 345:1052–1054. <https://doi.org/10.1126/science.1256985>.
- Füssel J, Lüscher S, Yilmaz P, Nowka B, van Kessel MAHJ, Bourceau P, Hach PF, Littmann S, Berg J, Spieck E, Daims H, Kuypers MMM, Lam P. 2017. Adaptability as the key to success for the ubiquitous marine nitrite oxidizer *Nitrococcus*. *Sci Adv* 3:e1700807. <https://doi.org/10.1126/sciadv.1700807>.

25. Palatinszky M, Herbold C, Jehmlich N, Pogoda M, Han P, von Bergen M, Lagkouvardos I, Karst SM, Galushko A, Koch H, Berry D, Daims H, Wagner M. 2015. Cyanate as an energy source for nitrifiers. *Nature* 524:105–108. <https://doi.org/10.1038/nature14856>.
26. Koch H, Lückner S, Albertsen M, Kitzinger K, Herbold C, Spieck E, Nielsen PH, Wagner M, Daims H. 2015. Expanded metabolic versatility of ubiquitous nitrite-oxidizing bacteria from the genus *Nitrospira*. *Proc Natl Acad Sci U S A* 112:11371–11376. <https://doi.org/10.1073/pnas.1506533112>.
27. Boddicker AM, Mosier AC. 2018. Genomic profiling of four cultivated *Candidatus Nitrotoga* spp. predicts broad metabolic potential and environmental distribution. *ISME J* 12:2864–2882. <https://doi.org/10.1038/s41396-018-0240-8>.
28. Sterner RW. 2010. In situ-measured primary production in Lake Superior. *J Great Lakes Res* 36:139–149. <https://doi.org/10.1016/j.jglr.2009.12.007>.
29. Small GE, Bullerjahn GS, Sterner RW, Beall BFN, Brovold S, Finlay JC, McKay RML, Mukherjee M. 2013. Rates and controls of nitrification in a large oligotrophic lake. *Limnol Oceanogr* 58:276–286. <https://doi.org/10.4319/lo.2013.58.1.0276>.
30. Vollenweider RA, Munawar M, Stadelmann P. 1974. A comparative review of phytoplankton and primary production in the Laurentian Great Lakes. *J Fish Res Bd Can* 31:739–762. <https://doi.org/10.1139/f74-100>.
31. Clevinger CC, Heath RT, Bade DL. 2014. Oxygen use by nitrification in the hypolimnion and sediments of Lake Erie. *J Great Lakes Res* 40:202–207. <https://doi.org/10.1016/j.jglr.2013.09.015>.
32. Neilson MA, Stevens RJ. 1987. Spatial heterogeneity of nutrients and organic matter in Lake Ontario. *Can J Fish Aquat Sci* 44:2192–2203. <https://doi.org/10.1139/f87-269>.
33. Lean DRS, Knowles R. 1987. Nitrogen transformations in Lake Ontario. *Can J Fish Aquat Sci* 44:2133–2143. <https://doi.org/10.1139/f87-262>.
34. Paver SF, Newton RJ, Coleman ML. 2020. Microbial communities of the Laurentian Great Lakes reflect connectivity and local biogeochemistry. *Environ Microbiol* 22:433–446. <https://doi.org/10.1111/1462-2920.14862>.
35. Rozmarynowycz MJ, Beall BFN, Bullerjahn GS, Small GE, Sterner RW, Brovold SS, D'souza NA, Watson SB, McKay RML. 2019. Transitions in microbial communities along a 1600 km freshwater trophic gradient. *J Great Lakes Res* 45:263–276. <https://doi.org/10.1016/j.jglr.2019.01.004>.
36. Fujimoto M, Cavaletto J, Liebig JR, McCarthy A, Vanderploeg HA, Deneff VJ. 2016. Spatiotemporal distribution of bacterioplankton functional groups along a freshwater estuary to pelagic gradient in Lake Michigan. *J Great Lakes Res* 42:1036–1048. <https://doi.org/10.1016/j.jglr.2016.07.029>.
37. Hooper AB, Terry KR. 1974. Photoinactivation of ammonia oxidation in *Nitrosomonas*. *J Bacteriol* 119:899–906. <https://doi.org/10.1128/jb.119.3.899-906.1974>.
38. Horrigan SG, Springer AL. 1990. Oceanic and estuarine ammonium oxidation: effects of light. *Limnol Oceanogr* 35:479–482. <https://doi.org/10.4319/lo.1990.35.2.0479>.
39. Guerrero M, Jones R. 1996. Photoinhibition of marine nitrifying bacteria. I. Wavelength-dependent response. *Mar Ecol Prog Ser* 141:183–192. <https://doi.org/10.3354/meps141183>.
40. Merbt SN, Stahl DA, Casamayor EO, Martí E, Nicol GW, Prosser JI. 2012. Differential photoinhibition of bacterial and archaeal ammonia oxidation. *FEMS Microbiol Lett* 327:41–46. <https://doi.org/10.1111/j.1574-6968.2011.02457.x>.
41. Smith JM, Chavez FP, Francis CA. 2014. Ammonium uptake by phytoplankton regulates nitrification in the sunlit ocean. *PLoS One* 9:e108173. <https://doi.org/10.1371/journal.pone.0108173>.
42. Kumar S, Sterner RW, Finlay JC, Brovold S. 2007. Spatial and temporal variation of ammonium in Lake Superior. *J Great Lakes Res* 33:581–591. [https://doi.org/10.3394/0380-1330\(2007\)33\[581:SATVOA\]2.0.CO;2](https://doi.org/10.3394/0380-1330(2007)33[581:SATVOA]2.0.CO;2).
43. Dove A, Chapra SC. 2015. Long-term trends of nutrients and trophic response variables for the Great Lakes. *Limnol Oceanogr* 60:696–721. <https://doi.org/10.1002/lno.10055>.
44. Belisle BS, Steffen MM, Pound HL, Watson SB, DeBruyn JM, Bourbonniere RA, Boyer GL, Wilhelm SW. 2016. Urea in Lake Erie: organic nutrient sources as potentially important drivers of phytoplankton biomass. *J Great Lakes Res* 42:599–607. <https://doi.org/10.1016/j.jglr.2016.03.002>.
45. Rice MC, Norton JM, Valois F, Bollmann A, Bottomley PJ, Klotz MG, Laanbroek HJ, Suwa Y, Stein LY, Sayavedra-Soto L, Woyke T, Shapiro N, Goodwin LA, Huntemann M, Clum A, Pillay M, Kyrpides N, Varghese N, Mikhailova N, Markowitz V, Palaniappan K, Ivanova N, Stamatis D, Reddy TBK, Ngan CY, Daum C. 2016. Complete genome of *Nitrosospira briensis* C-128, an ammonia-oxidizing bacterium from agricultural soil. *Stand Genomic Sci* 11:46. <https://doi.org/10.1186/s40793-016-0168-4>.
46. Garcia JC, Urakawa H, Le VQ, Stein LY, Klotz MG, Nielsen JL. 2013. Draft genome sequence of *Nitrosospira* sp. strain APG3, a psychrotolerant ammonia-oxidizing bacterium isolated from sandy lake sediment. *Genome Announc* 1:e00930-13. <https://doi.org/10.1128/genomeA.00930-13>.
47. Norton JM, Klotz MG, Stein LY, Arp DJ, Bottomley PJ, Chain PSG, Hauser LJ, Land ML, Larimer FW, Shin MW, Starckenburg SR. 2008. Complete genome sequence of *Nitrosospira multiformis*, an ammonia-oxidizing bacterium from the soil environment. *Appl Environ Microbiol* 74:3559–3572. <https://doi.org/10.1128/AEM.02722-07>.
48. Giovannoni SJ, Thrash JC, Temperton B. 2014. Implications of streamlining theory for microbial ecology. *ISME J* 8:1553–1565. <https://doi.org/10.1038/ismej.2014.60>.
49. Bollmann A, Sedlacek CJ, Norton J, Laanbroek HJ, Suwa Y, Stein LY, Klotz MG, Arp D, Sayavedra-Soto L, Lu M, Bruce D, Detter C, Tapia R, Han J, Woyke T, Lucas SM, Pitluck S, Pennacchio L, Nolan M, Land ML, Huntemann M, Deshpande S, Han C, Chen A, Kyrpides F, Mavromatis K, Markowitz V, Szeto E, Ivanova N, Mikhailova N, Pagani I, Pati A, Peters L, Ovchinnikova G, Goodwin LA. 2013. Complete genome sequence of *Nitrosomonas* sp. Is79, an ammonia oxidizing bacterium adapted to low ammonium concentrations. *Stand Genomic Sci* 7:469–482. <https://environmentalmicrobiome.biomedcentral.com/articles/10.4056/signs.3517166>.
50. Bèjà O, Aravind L, Koonin EV, Suzuki MT, Hadd A, Nguyen LP, Jovanovich SB, Gates CM, Feldman RA, Spudich JL, Spudich EN, DeLong EF. 2000. Bacteroid rhodopsin: evidence for a new type of phototrophy in the sea. *Science* 289:1902–1906. <https://doi.org/10.1126/science.289.5486.1902>.
51. Bèjà O, Spudich EN, Spudich JL, Leclerc M, DeLong EF. 2001. Proteorhodopsin phototrophy in the ocean. *Nature* 411:786–789. <https://doi.org/10.1038/35081051>.
52. Sabehi G, Loy A, Jung K-H, Partha R, Spudich JL, Isaacson T, Hirschberg J, Wagner M, Bèjà O. 2005. New insights into metabolic properties of marine bacteria encoding proteorhodopsins. *PLoS Biol* 3:e273. <https://doi.org/10.1371/journal.pbio.0030273>.
53. Martinez A, Bradley AS, Waldbauer JR, Summons RE, DeLong EF. 2007. Proteorhodopsin photosystem gene expression enables photophosphorylation in a heterologous host. *Proc Natl Acad Sci U S A* 104:5590–5595. <https://doi.org/10.1073/pnas.0611470104>.
54. Reckel S, Gottstein D, Stehle J, Löhr F, Verhoeven M-K, Takeda M, Silvers R, Kainosho M, Glaubitz C, Wachtveitl J, Bernhard F, Schwalbe H, Güntert P, Dötsch V. 2011. Solution NMR structure of proteorhodopsin. *Angew Chem Int Ed Engl* 50:11942–11946. <https://doi.org/10.1002/anie.201105648>.
55. Kralj JM, Bergo VB, Amsden JL, Spudich EN, Spudich KJ. 2008. Protonation state of Glu142 differs in the green- and blue-absorbing variants of proteorhodopsin. *Biochemistry* 47:3447–3453. <https://doi.org/10.1021/bi7018964>.
56. Boeuf D, Audic S, Brillet-Guéguen L, Caron C, Jeanthon C. 2015. MicRhoDE: a curated database for the analysis of microbial rhodopsin diversity and evolution. *Database (Oxford)* 2015:bav080. <https://doi.org/10.1093/database/bav080>.
57. Hyman MR, Arp DJ. 1992. ¹⁴C₂H₂- and ¹⁴CO₂-labeling studies of the de novo synthesis of polypeptides by *Nitrosomonas europaea* during recovery from acetylene and light inactivation of ammonia monooxygenase. *J Biol Chem* 267:1534–1545. [https://doi.org/10.1016/S0021-9258\(18\)45979-0](https://doi.org/10.1016/S0021-9258(18)45979-0).
58. Yousef F, Shuchman R, Sayers M, Fahnenstiel G, Henareh A. 2017. Water clarity of the Upper Great Lakes: tracking changes between 1998–2012. *J Great Lakes Res* 43:239–247. <https://doi.org/10.1016/j.jglr.2016.12.002>.
59. Beaumont HJE, Hommes NG, Sayavedra-Soto LA, Arp DJ, Arciero DM, Hooper AB, Westerhoff HV, van Spanning RJM. 2002. Nitrite reductase of *Nitrosomonas europaea* is not essential for production of gaseous nitrogen oxides and confers tolerance to nitrite. *J Bacteriol* 184:2557–2560. <https://doi.org/10.1128/JB.184.9.2557-2560.2002>.
60. Watzel B, Forchhammer K. 2018. Cyanophycin synthesis optimizes nitrogen utilization in the unicellular cyanobacterium *Synechocystis* sp. strain PCC 6803. *Appl Environ Microbiol* 84:e01298-18. <https://doi.org/10.1128/AEM.01298-18>.
61. Badger MR, Bek EJ. 2008. Multiple Rubisco forms in proteobacteria: their functional significance in relation to CO₂ acquisition by the CBB cycle. *J Exp Bot* 59:1525–1541. <https://doi.org/10.1093/jxb/erm297>.
62. Rae BD, Long BM, Badger MR, Price GD. 2013. Functions, compositions, and evolution of the two types of carboxysomes: polyhedral microcompartments that facilitate CO₂ fixation in cyanobacteria and some proteobacteria. *Microbiol Mol Biol Rev* 77:357–379. <https://doi.org/10.1128/MMBR.00061-12>.
63. Cabello-Yeves PJ, Zemskaia TI, Rosselli R, Coutinho FH, Zakharenko AS, Blinov VV, Rodriguez-Valera F. 2018. Genomes of novel microbial

- lineages assembled from the sub-ice waters of Lake Baikal. *Appl Environ Microbiol* 84:e02132-17. <https://doi.org/10.1128/AEM.02132-17>.
64. Sterner RW, Smutka TM, McKay RML, Xiaoming Q, Brown ET, Sherrill RM. 2004. Phosphorus and trace metal limitation of algae and bacteria in Lake Superior. *Limnol Oceanogr* 49:495–507. <https://doi.org/10.4319/lo.2004.49.2.0495>.
 65. Ryals J, Hsu RY, Lipsett MN, Bremer H. 1982. Isolation of single-site *Escherichia coli* mutants deficient in thiamine and 4-thiouridine syntheses: identification of a *nuvC* mutant. *J Bacteriol* 151:899–904. <https://doi.org/10.1128/jb.151.2.899-904.1982>.
 66. Kitzinger K, Koch H, Lückner S, Sedlacek CJ, Herbold C, Schwarz J, Daebeler A, Mueller AJ, Lukumbuzya M, Romano S, Leisch N, Karst SM, Kirkegaard R, Albertsen M, Nielsen PH, Wagner M, Daims H. 2018. Characterization of the first “*Candidatus Nitrotoga*” isolate reveals metabolic versatility and separate evolution of widespread nitrite-oxidizing bacteria. *mBio* 9:e01186-18. <https://doi.org/10.1128/mBio.01186-18>.
 67. Ishii K, Fujitani H, Sekiguchi Y, Tsuneda S. 2020. Physiological and genomic characterization of a new “*Candidatus Nitrotoga*” isolate. *Environ Microbiol* 22:2365–2382. <https://doi.org/10.1111/1462-2920.15015>.
 68. Starckenburg SR, Chain PSG, Sayavedra-Soto LA, Hauser L, Land ML, Larimer FW, Malfatti SA, Klotz MG, Bottomley PJ, Arp DJ, Hickey WJ. 2006. Genome sequence of the chemolithoautotrophic nitrite-oxidizing bacterium *Nitrobacter winogradskyi* Nb-255. *Appl Environ Microbiol* 72:2050–2063. <https://doi.org/10.1128/AEM.72.3.2050-2063.2006>.
 69. Lückner S, Wagner M, Maixner F, Pelletier E, Koch H, Vacherie B, Rattei T, Damste JSS, Spieck E, Le Paslier D, Daims H. 2010. A *Nitrospira* metagenome illuminates the physiology and evolution of globally important nitrite-oxidizing bacteria. *Proc Natl Acad Sci U S A* 107:13479–13484. <https://doi.org/10.1073/pnas.1003860107>.
 70. Lückner S, Nowka B, Rattei T, Spieck E, Daims H. 2013. The genome of *Nitrospina gracilis* illuminates the metabolism and evolution of the major marine nitrite oxidizer. *Front Microbiol* 4:27. <https://doi.org/10.3389/fmicb.2013.00027>.
 71. Kitzinger K, Marchant HK, Bristow LA, Herbold CW, Padilla CC, Kidane AT, Littmann S, Daims H, Pjevac P, Stewart FJ, Wagner M, Kuypers MMM. 2020. Single cell analyses reveal contrasting life strategies of the two main nitrifiers in the ocean. *Nat Commun* 11:767. <https://doi.org/10.1038/s41467-020-14542-3>.
 72. Munding AB, Lawson CE, Jetten MSM, Koch H, Lückner S. 2019. Cultivation and transcriptional analysis of a canonical *Nitrospira* under stable growth conditions. *Front Microbiol* 10:1325. <https://doi.org/10.3389/fmicb.2019.01325>.
 73. Sancar A. 2003. Structure and function of DNA photolyase and cryptochrome blue-light photoreceptors. *Chem Rev* 103:2203–2238. <https://doi.org/10.1021/cr0204348>.
 74. von Zadow A, Ignatz E, Pokorny R, Essen L-O, Klug G. 2016. *Rhodobacter sphaeroides* CryB is a bacterial cryptochrome with (6–4) photolyase activity. *FEBS J* 283:4291–4309. <https://doi.org/10.1111/febs.13924>.
 75. Zhang F, Scheerer P, Oberpichler I, Lamparter T, Krauß N. 2013. Crystal structure of a prokaryotic (6–4) photolyase with an Fe-S cluster and a 6,7-dimethyl-8-ribityllumazine antenna chromophore. *Proc Natl Acad Sci U S A* 110:7217–7222. <https://doi.org/10.1073/pnas.1302377110>.
 76. Maresca JA, Keffer JL, Hempel PP, Polson SW, Shevchenko O, Bhavsar J, Powell D, Miller KJ, Singh A, Hahn MW. 2019. Light modulates the physiology of nonphototrophic Actinobacteria. *J Bacteriol* 201:e00740-18. <https://doi.org/10.1128/JB.00740-18>.
 77. Wink DA, Kasprzak KS, Maragos CM, Elespuru RK, Misra M, Dunams TM, Cebula TA, Koch WH, Andrews AW, Allen JS, et al. 1991. DNA deaminating ability and genotoxicity of nitric oxide and its progenitors. *Science* 254:1001–1003. <https://doi.org/10.1126/science.1948068>.
 78. Imlay JA. 2013. The molecular mechanisms and physiological consequences of oxidative stress: lessons from a model bacterium. *Nat Rev Microbiol* 11:443–454. <https://doi.org/10.1038/nrmicro3032>.
 79. Zhang T, Hansel CM, Voelker BM, Lamborg CH. 2016. Extensive dark biological production of reactive oxygen species in brackish and freshwater ponds. *Environ Sci Technol* 50:2983–2993. <https://doi.org/10.1021/acs.est.5b03906>.
 80. Khademian M, Imlay JA. 2017. *Escherichia coli* cytochrome *c* peroxidase is a respiratory oxidase that enables the use of hydrogen peroxide as a terminal electron acceptor. *Proc Natl Acad Sci U S A* 114:E6922–E6931. <https://doi.org/10.1073/pnas.1701587114>.
 81. O’Reilly CM, Sharma S, Gray DK, Hampton SE, Read JS, Rowley RJ, Schneider P, Lenters JD, McIntyre PB, Kraemer BM, Weyhenmeyer GA, Straile D, Dong B, Adrian R, Allan MG, Anneville O, Arvola L, Austin J, Bailey JL, Baron JS, Brookes JD, de Eyto E, Dokulil MT, Hamilton DP, Havens K, Hetherington AL, Higgins SN, Hook S, Izmeševa LR, Joehnk KD, Kangur K, Kasprzak P, Kumagai M, Kuusisto E, Leshkevich G, Livingstone DM, MacIntyre S, May L, Melack JM, Mueller-Navarra DC, Naumenko M, Noges P, Noges T, North RP, Plisnier P-D, Rigosi A, Rimmer A, Rogora M, Rudstam LG, Rusak JA, Salmaso N, Samal NR, Schindler DE, Schladow SG, Schmid M, Schmidt SR, Silow E, et al. 2015. Rapid and highly variable warming of lake surface waters around the globe. *Geophys Res Lett* 42:10773–10781.
 82. Sterner RW, Keeler B, Polasky S, Poudel R, Rhude K, Rogers M. 2020. Ecosystem services of Earth’s largest freshwater lakes. *Ecosyst Serv* 41:101046. <https://doi.org/10.1016/j.ecoser.2019.101046>.
 83. Barbiero RP, Lesht BM, Hinchey EK, Nettesheim TG. 2018. A brief history of the U.S. EPA Great Lakes National Program Office’s water quality survey. *J Great Lakes Res* 44:539–546. <https://doi.org/10.1016/j.jglr.2018.05.011>.
 84. U.S. Environmental Protection Agency. 2003. Sampling and analytical procedures for GLNPO’s open lake water quality survey of the Great Lakes. Great Lakes National Program Office, Chicago, IL.
 85. Holmes RM, Aminot A, Kérouel R, Hooker BA, Peterson BJ. 1999. A simple and precise method for measuring ammonium in marine and freshwater ecosystems. *Can J Fish Aquat Sci* 56:1801–1808. <https://doi.org/10.1139/f99-128>.
 86. Miranda KM, Espey MG, Wink DA. 2001. A rapid, simple spectrophotometric method for simultaneous detection of nitrate and nitrite. *Nitric Oxide* 5:62–71. <https://doi.org/10.1006/niox.2000.0319>.
 87. Revilla M, Alexander J, Glibert PM. 2005. Urea analysis in coastal waters: comparison of enzymatic and direct methods. *Limnol Oceanogr Methods* 3:290–299. <https://doi.org/10.4319/lom.2005.3.290>.
 88. 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>.
 89. Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, Peplies J, Glöckner FO. 2013. The SILVA ribosomal RNA gene database project: improved data processing and web-based tools. *Nucleic Acids Res* 41:D590–D596. <https://doi.org/10.1093/nar/gks1219>.
 90. Wang Q, Garrity GM, Tiedje JM, Cole JR. 2007. Naïve Bayesian classifier for rapid assignment of rRNA sequences into the new bacterial taxonomy. *Appl Environ Microbiol* 73:5261–5267. <https://doi.org/10.1128/AEM.00062-07>.
 91. Kozich JJ, Westcott SL, Baxter NT, Highlander SK, Schloss PD. 2013. Development of a dual-index sequencing strategy and curation pipeline for analyzing amplicon sequence data on the MiSeq Illumina sequencing platform. *Appl Environ Microbiol* 79:5112–5120. <https://doi.org/10.1128/AEM.01043-13>.
 92. Eren AM, Maignien L, Sul WJ, Murphy LG, Grim SL, Morrison HG, Sogin ML. 2013. Oligotyping: differentiating between closely related microbial taxa using 16S rRNA gene data. *Methods Ecol Evol* 4:1111–1119. <https://doi.org/10.1111/2041-210X.12114>.
 93. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin VM, Nikolenko SI, Pham S, Pribelski AD, Pyshkin AV, Sirotkin AV, Vyahhi N, Tesler G, Alekseyev MA, Pevzner PA. 2012. SPAdes: a new genome assembly algorithm and its applications to single-cell sequencing. *J Comput Biol* 19:455–477. <https://doi.org/10.1089/cmb.2012.0021>.
 94. Nurk S, Meleshko D, Korobeynikov A, Pevzner PA. 2017. metaSPAdes: a new versatile metagenomic assembler. *Genome Res* 27:824–834. <https://doi.org/10.1101/gr.213959.116>.
 95. Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nat Methods* 9:357–359. <https://doi.org/10.1038/nmeth.1923>.
 96. Kang DD, Froula J, Egan R, Wang Z. 2015. MetaBAT, an efficient tool for accurately reconstructing single genomes from complex microbial communities. *PeerJ* 3:e1165. <https://doi.org/10.7717/peerj.1165>.
 97. Graham ED, Heidelberg JF, Tully BJ. 2017. BinSanity: unsupervised clustering of environmental microbial assemblies using coverage and affinity propagation. *PeerJ* 5:e3035. <https://doi.org/10.7717/peerj.3035>.
 98. Alneberg J, Bjarnason BS, de Bruijn I, Schirmer M, Quick J, Ijaz UZ, Lahti L, Loman NJ, Andersson AF, Quince C. 2014. Binning metagenomic contigs by coverage and composition. *Nat Methods* 11:1144–1146. <https://doi.org/10.1038/nmeth.3103>.
 99. Sieber CMK, Probst AJ, Sharrar A, Thomas BC, Hess M, Tringe SG, Banfield JF. 2018. Recovery of genomes from metagenomes via a dereplication, aggregation and scoring strategy. *Nat Microbiol* 3:836–843. <https://doi.org/10.1038/s41564-018-0171-1>.

100. Eren AM, Esen ÖC, Quince C, Vineis JH, Morrison HG, Sogin ML, Delmont TO. 2015. Anvi'o: an advanced analysis and visualization platform for 'omics data. *PeerJ* 3:e1319. <https://doi.org/10.7717/peerj.1319>.
101. 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>.
102. 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>.
103. Campbell BJ, Yu L, Heidelberg JF, Kirchman DL. 2011. Activity of abundant and rare bacteria in a coastal ocean. *Proc Natl Acad Sci U S A* 108:12776–12781. <https://doi.org/10.1073/pnas.1101405108>.
104. 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, Kyrpides 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>.
105. Mistry J, Finn RD, Eddy SR, Bateman A, Punta M. 2013. Challenges in homology search: HMMER3 and convergent evolution of coiled-coil regions. *Nucleic Acids Res* 41:e121. <https://doi.org/10.1093/nar/gkt263>.
106. O'Leary NA, Wright MW, Brister JR, Ciuflo S, Haddad D, McVeigh R, Rajput B, Robbertse B, Smith-White B, Ako-Adjei D, Astashyn A, Badretdin A, Bao Y, Blinkova O, Brover V, Chetverin V, Choi J, Cox E, Ermolaeva O, Farrell CM, Goldfarb T, Gupta T, Haft D, Hatcher E, Hlavina W, Joardar VS, Kodali VK, Li W, Maglott D, Masterson P, McGarvey KM, Murphy MR, O'Neill K, Pujar S, Rangwala SH, Rausch D, Riddick LD, Schoch C, Shkeda A, Storz SS, Sun H, Thibaud-Nissen F, Tolstoy I, Tully RE, Vatsan AR, Wallin C, Webb D, Wu W, Landrum MJ, Kimchi A, et al. 2016. Reference sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and functional annotation. *Nucleic Acids Res* 44:D733–D745. <https://doi.org/10.1093/nar/gkv1189>.
107. Ondov BD, Bergman NH, Phillippy AM. 2011. Interactive metagenomic visualization in a Web browser. *BMC Bioinformatics* 12:385. <https://doi.org/10.1186/1471-2105-12-385>.
108. Chaumeil P-A, Mussig AJ, Hugenholtz P, Parks DH. 2019. GTDB-Tk: a toolkit to classify genomes with the Genome Taxonomy Database. *Bioinformatics* 36:1925–1927. <https://doi.org/10.1093/bioinformatics/btz848>.
109. Jain C, Rodriguez-R LM, Phillippy AM, Konstantinidis KT, Aluru S. 2018. High throughput ANI analysis of 90K prokaryotic genomes reveals clear species boundaries. *Nat Commun* 9:5114. <https://doi.org/10.1038/s41467-018-07641-9>.
110. Hyatt D, Chen G-L, LoCascio PF, Land ML, Larimer FW, Hauser LJ. 2010. Prodigal: prokaryotic gene recognition and translation initiation site identification. *BMC Bioinformatics* 11:119. <https://doi.org/10.1186/1471-2105-11-119>.
111. Jones P, Binns D, Chang H-Y, Fraser M, Li W, McAnulla C, McWilliam H, Maslen J, Mitchell A, Nuka G, Pesseat S, Quinn AF, Sangrador-Vegas A, Scheremetjew M, Yong S-Y, Lopez R, Hunter S. 2014. InterProScan 5: genome-scale protein function classification. *Bioinformatics* 30:1236–1240. <https://doi.org/10.1093/bioinformatics/btu031>.
112. Kanehisa M, Sato Y, Morishima K. 2016. BlastKOALA and GhostKOALA: KEGG tools for functional characterization of genome and metagenome sequences. *J Mol Biol* 428:726–731. <https://doi.org/10.1016/j.jmb.2015.11.006>.
113. Huerta-Cepas J, Forslund K, Coelho LP, Szklarczyk D, Jensen LJ, von Mering C, Bork P. 2017. Fast genome-wide functional annotation through orthology assignment by eggNOG-Mapper. *Mol Biol Evol* 34:2115–2122. <https://doi.org/10.1093/molbev/msx148>.
114. Delmont TO, Eren AM. 2018. Linking pangenomes and metagenomes: the *Prochlorococcus* metapangenome. *PeerJ* 6:e4320. <https://doi.org/10.7717/peerj.4320>.
115. Benedict MN, Henriksen JR, Metcalf WW, Whitaker RJ, Price ND. 2014. ITEP: an integrated toolkit for exploration of microbial pan-genomes. *BMC Genomics* 15:8. <https://doi.org/10.1186/1471-2164-15-8>.
116. van Dongen S, Abreu-Goodger C. 2012. Using MCL to extract clusters from networks. *Methods Mol Biol* 804:281–295. https://doi.org/10.1007/978-1-61779-361-5_15.
117. Tatusova T, DiCuccio M, Badretdin A, Chetverin V, Nawrocki EP, Zaslavsky L, Lomsadze A, Pruitt KD, Borodovsky M, Ostell J. 2016. NCBI prokaryotic genome annotation pipeline. *Nucleic Acids Res* 44:6614–6624. <https://doi.org/10.1093/nar/gkw569>.
118. 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>.
119. Seemann T. 2014. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* 30:2068–2069. <https://doi.org/10.1093/bioinformatics/btu153>.
120. Guy L, Roat Kultima J, Andersson SGE. 2010. genoPlotR: comparative gene and genome visualization in R. *Bioinformatics* 26:2334–2335. <https://doi.org/10.1093/bioinformatics/btq413>.
121. Matsen FA, Kodner RB, Armbrust EV. 2010. pplacer: linear time maximum-likelihood and Bayesian phylogenetic placement of sequences onto a fixed reference tree. *BMC Bioinformatics* 11:538. <https://doi.org/10.1186/1471-2105-11-538>.
122. Katoh K, Misawa K, Kuma K, Miyata T. 2002. MAFFT: a novel method for rapid multiple sequence alignment based on fast Fourier transform. *Nucleic Acids Res* 30:3059–3066. <https://doi.org/10.1093/nar/gkf436>.
123. 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>.
124. Letunic I, Bork P. 2016. Interactive tree of life (iTOL) v3: an online tool for the display and annotation of phylogenetic and other trees. *Nucleic Acids Res* 44:W242–W245. <https://doi.org/10.1093/nar/gkw290>.
125. Spang A, Poehlein A, Offre P, Zumbärgel S, Haider S, Rychlik N, Nowka B, Schmeisser C, Lebedeva EV, Rattei T, Böhm C, Schmid M, Galushko A, Hatzenpichler R, Weinmaier T, Daniel R, Schleper C, Spieck E, Streit W, Wagner M. 2012. The genome of the ammonia-oxidizing *Candidatus Nitrososphaera gargensis*: insights into metabolic versatility and environmental adaptations. *Environ Microbiol* 14:3122–3145. <https://doi.org/10.1111/j.1462-2920.2012.02893.x>.
126. Vechtomova YL, Telegina TA, Kritsky MS. 2020. Evolution of proteins of the DNA photolyase/cryptochrome family. *Biochemistry (Mosc)* 85 (Suppl):S131–S153. <https://doi.org/10.1134/S0006297920140072>.
127. Minh BQ, Schmidt HA, Chernomor O, Schrempf D, Woodhams MD, von Haeseler A, Lanfear R. 2020. IQ-TREE 2: new models and efficient methods for phylogenetic inference in the genomic era. *Mol Biol Evol* 37:1530–1534. <https://doi.org/10.1093/molbev/msaa015>.
128. Benson DA, Karsch-Mizrachi I, Lipman DJ, Ostell J, Wheeler DL. 2005. GenBank. *Nucleic Acids Res* 33:D34–D38. <https://doi.org/10.1093/nar/gki063>.
129. 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>.
130. 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>.
131. Wick RR, Schultz MB, Zobel J, Holt KE. 2015. Bandage: interactive visualization of de novo genome assemblies. *Bioinformatics* 31:3350–3352. <https://doi.org/10.1093/bioinformatics/btv383>.
132. Thorvaldsdóttir H, Robinson JT, Mesirov JP. 2013. Integrative Genomics Viewer (IGV): high-performance genomics data visualization and exploration. *Brief Bioinformatics* 14:178–192. <https://doi.org/10.1093/bib/bbs017>.
133. Okazaki Y, Nishimura Y, Yoshida T, Ogata H, Nakano S-I. 2019. Genome-resolved viral and cellular metagenomes revealed potential key virus-host interactions in a deep freshwater lake. *Environ Microbiol* 21:4740–4754. <https://doi.org/10.1111/1462-2920.14816>.
134. Orellana LH, Rodriguez-R LM, Konstantinidis KT. 2017. ROCKER: accurate detection and quantification of target genes in short-read metagenomic data sets by modeling sliding-window bitscores. *Nucleic Acids Res* 45:e14. <https://doi.org/10.1093/nar/gkw900>.
135. R Core Team. 2019. R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. <https://www.R-project.org>.
136. Wickham H. 2009. ggplot2: Elegant graphics for data analysis. Springer-Verlag, New York, NY.
137. Murphy TP. 1980. Ammonia and nitrate uptake in the lower Great Lakes. *Can J Fish Aquat Sci* 37:1365–1372. <https://doi.org/10.1139/f80-175>.
138. Gardner WS, Lavrentyev PJ, Cavaletto JF, McCarthy MJ, Eadie BJ, Johengen TH, Cotner JB. 2004. Distribution and dynamics of nitrogen and microbial plankton in southern Lake Michigan during spring transition 1999–2000. *J Geophys Res Oceans* 109:C03007. <https://doi.org/10.1029/2002JC001588>.