



# Microbial Populations Are Shaped by Dispersal and Recombination in a Low Biomass Subseafloor Habitat

 Rika E. Anderson,<sup>a</sup>  Elaina D. Graham,<sup>b</sup>  Julie A. Huber,<sup>c</sup>  Benjamin J. Tully<sup>b,d,e</sup>

<sup>a</sup>Biology Department, Carleton College, Northfield, Minnesota, USA

<sup>b</sup>Department of Biological Sciences, University of Southern California, Los Angeles, California, USA

<sup>c</sup>Marine Chemistry and Geochemistry, Woods Hole Oceanographic Institution, Woods Hole, Massachusetts, USA

<sup>d</sup>Center for Dark Energy Biosphere Investigations, University of Southern California, Los Angeles, California, USA

<sup>e</sup>Wrigley Institute for Environmental Studies, University of Southern California, Los Angeles, California, USA

Rika E. Anderson and Benjamin J. Tully contributed equally to this article. Author order was determined based on the origination of the concept for the analysis by Benjamin J. Tully.

**ABSTRACT** The subseafloor is a vast habitat that supports microorganisms that have a global scale impact on geochemical cycles. Many of the endemic microbial communities inhabiting the subseafloor consist of small populations under growth-limited conditions. For small populations, stochastic evolutionary events can have large impacts on intra-specific population dynamics and allele frequencies. These conditions are fundamentally different from those experienced by most microorganisms in surface environments, and it is unknown how small population sizes and growth-limiting conditions influence evolution and population structure in the subsurface. Using a 2-year, high-resolution environmental time series, we examine the dynamics of microbial populations from cold, oxic crustal fluids collected from the subseafloor site North Pond, located near the mid-Atlantic ridge. Our results reveal rapid shifts in overall abundance, allele frequency, and strain abundance across the time points observed, with evidence for homologous recombination between coexisting lineages. We show that the subseafloor aquifer is a dynamic habitat that hosts microbial metapopulations that disperse frequently through the crustal fluids, enabling gene flow and recombination between microbial populations. The dynamism and stochasticity of microbial population dynamics in North Pond suggest that these forces are important drivers in the evolution of microbial populations in the vast subseafloor habitat.

**IMPORTANCE** The cold, oxic subseafloor is an understudied habitat that is difficult to access, yet important to global biogeochemical cycles and starkly different compared to microbial habitats on the surface of the Earth. Our understanding of microbial evolution and population dynamics is largely molded by studies of microbes living in surface habitats that can host 10 to 1,000 times more microbial biomass than is frequently observed in the subsurface. This study provides an opportunity to observe population dynamics within a low biomass, growth-limited environment and reveals that microbial populations in the subseafloor are influenced by changes in selection pressure and gene sweeps. In addition, recombination between strains that have dispersed from elsewhere within the aquifer has an important impact on the evolution of microbial populations. Much of the microbial life on the planet exists under growth-limited conditions, and the subseafloor provides a natural laboratory to explore how life evolves in such environments.

**KEYWORDS** microbial evolution, subseafloor, allele frequency, dispersal

The marine subsurface is home to more than half of the archaea and bacteria inhabiting the oceans (1, 2). Marine subsurface habitats are frequently energy limited, which impacts the total biomass that can be supported in these environments (3). As a

**Editor** Colleen M. Cavanaugh, Harvard University

**Copyright** © 2022 Anderson 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 Rika E. Anderson, [randerson@carleton.edu](mailto:randerson@carleton.edu), or Benjamin J. Tully, [tully.bj@gmail.com](mailto:tully.bj@gmail.com).

The authors declare no conflict of interest.

**Received** 6 March 2022

**Accepted** 1 July 2022

**Published** 1 August 2022

result, microbial communities in these habitats are often composed of relatively low biomass (4–7) ( $<10^4$  cells per unit volume). Population size can have important consequences for how microbial populations evolve over time, and factors such as genetic drift can have an exaggerated impact on small microbial populations compared to habitats in which microbial populations are larger. Most previous studies of adaptation in natural microbial communities have focused on habitats with high biomass and thus larger population sizes, such as the surface ocean, deep-sea hydrothermal vents, and freshwater lakes (8–11). It is unclear how natural selection and genetic drift shape the evolution of microbial communities in low-biomass habitats, despite the fact that a substantial fraction of global microbial biomass is found in such habitats.

North Pond, located near the mid-Atlantic ridge (12), offers a unique opportunity to examine the dynamics of microbial adaptation in natural populations in a low-biomass, growth-limited habitat. The North Pond site consists of exposed crustal ridges with a depression that has accumulated  $<100$  m of low permeability sediment, causing the bulk of fluid exchange and transport to occur beneath the sediment basin (13). This subseafloor aquifer connects the deep waters of the oceans through a porous network of crustal rock, which globally contains about 2% of the oceans' volume (14–16). The crustal aquifer beneath the sediments of North Pond is accessible through two circulation obviator retrofit kits (CORKs) installed by the International Ocean Drilling Program (IODP) in 2011, providing direct access to the fluids circulating within the aquifer, while preventing exchange with the overlying water and sediment environments (17, 18). Repeated sampling in 2012, 2014, and 2017, including an eighth sample *in situ* time series collected from 2012 to 2014, revealed a dynamic microbial community with a high degree of functional redundancy in a complex hydrological system depleted in dissolved organic carbon (6, 19–21). The aquifer microbial communities are distinct from bottom seawater, and there are multiple lines of genetic evidence to suggest that the microbial community within the aquifer is motile and metabolically flexible, with the ability to carry out both autotrophic and organotrophic pathways (6, 19, 22). The North Pond subseafloor microbial community has low biomass ( $10^3$  to  $10^4$  cells  $\text{mL}^{-1}$ ), and incubation experiments suggest the North Pond microbial community is growth limited and that increasing temperature and enrichment with carbon sources result in rapid growth (6, 7, 22).

An open question in subsurface habitats is what the major drivers of evolution are in energy- and growth-limited environments such as North Pond. It has been hypothesized that in energy- and growth-limited environments, such as low-deposition deep-sea sediments, adaptive mutations would be unable to sweep through populations and mutations would accumulate as a result of energy limitation (23). However, many crustal habitats like North Pond do not reflect the limited fluid flow and isolated cellular environments in low-energy sediments and instead may reflect the population dynamics of soils with high fluid connectivity (24). North Pond fluids are capable of traversing distances of approximately (21) 2 to 40  $\text{m day}^{-1}$ , although this fluid flux is spatially and temporally complex, resulting in convective and oscillatory fluid movement rather than linear flow along the North-South axis (13). Fluid movement appears to impact the microbial community, as previous analysis of the system revealed discrete ecological units within the time series that would oscillate between presence and absence with extended time periods (e.g., hundreds of days) between reoccurrences (6). However, it is unclear how this dynamic environment impacts the endemic microbial populations.

Intraspecific diversity results from the introduction of variation through mutation or gene flow, and that variation is removed through selection and genetic drift. Thus, examination of patterns of genomic variation within populations can provide insights into evolutionary dynamics within a system. Here, we present an analysis of the intraspecific population dynamics of the North Pond aquifer microbial community from a metagenomic time series spanning 10 time points over a 2-year period. We sought to characterize the dynamics of population variation over time in order to determine what processes constrain variation within microbial populations at North Pond and to

**TABLE 1** Information for the MAGs of interest<sup>a</sup>

	Calendar Date	25-Apr-12	6-Aug-12	5-Oct-12	4-Dec-12	2-Feb-13	3-Apr-13	2-Jun-13	1-Aug-13	30-Oct-13	5-Apr-14
	Julian Days	116	219	279	339	399	459	519	579	669	826
MAG ID	GTDB Taxonomy (R95)	TP0	TP1	TP2	TP3	TP4	TP5	TP6	TP7	TP8	TP9
NORP83	(o) Rhizobiales	0.02	16.62	0.23	13.69	11.27	2.73	9.95	1.04	0.38	0.00
NORP139	(f) Parvibaculaceae	2.13	6.84	6.90	4.37	9.78	8.50	6.97	1.27	1.21	0.01
NORP147	(f) Rhizobiaceae	0.00	36.41	0.08	20.33	12.36	0.11	0.30	0.02	0.02	0.03
NORP163	(c) Gammaproteobacteria	20.86	16.51	48.66	3.43	2.92	1.12	0.24	19.77	20.95	2.97
NORP169	(g) Nitratireductor_B	2.12	6.53	5.09	0.82	0.52	1.43	8.38	4.84	4.14	0.08
NORP246	(g) Shewanella	0.02	11.60	0.10	8.34	5.24	1.74	0.12	0.10	0.13	0.00
NORP6	(g) Desulforhopalus	0.00	0.00	0.00	0.00	88.61	0.00	0.02	0.00	0.00	0.00
NORP57	(g) Halomonas	0.04	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	38.25
NORP100	(o) Thiohalomonadales	1.05	0.10	1.97	0.08	0.16	0.11	0.01	2.46	2.77	38.98
NORP167	(f) Flavobacteriaceae	0.03	30.13	0.07	0.13	0.01	0.00	0.00	0.01	0.00	0.00

<sup>a</sup>Values at each time point represent RPKM. MAGs in the top half of the table occurred in at least three samples with  $\geq 5$  RPKM. MAGs in the bottom half of the table represent MAGs that demonstrated a sudden increase in abundance at a single time point.

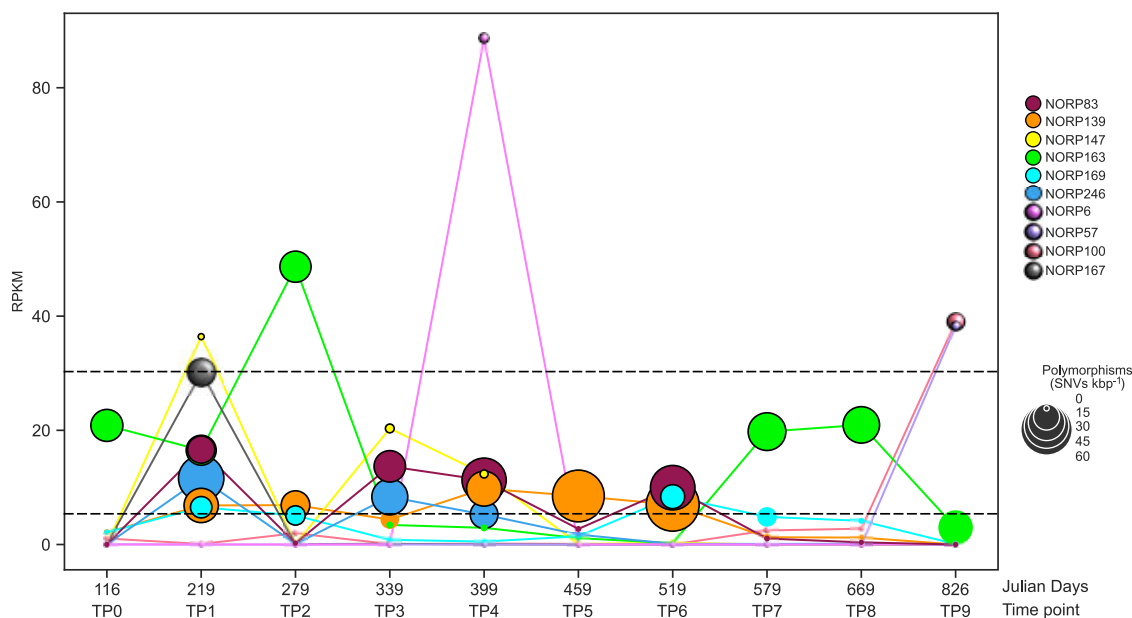
assess the balance between deterministic processes such as selection and more stochastic processes such as dispersal. We show that abundant microorganisms within the crustal fluids have patterns of genome diversity suggesting that large-scale genomic sweeps are infrequent and small gene sweeps, recombination, and dispersal of populations within the aquifer may have a more pronounced impact on genome diversity and therefore evolution in this habitat. These results provide new understanding of microbial population dynamics over time and how microorganisms adapt and persist in the largest habitat on the planet.

## RESULTS AND DISCUSSION

**Improving North Pond metagenome-assembled genomes.** By comparing and combining the total set of new metagenome-assembled genomes (MAGs) to the original 2018 North Pond (NORP) MAGs (6) ( $n = 195$  MAGs), it was possible to replace 22 2018 NORP MAGs with newly reconstructed MAGs with higher completion and/or lower redundancy statistics and to produce 137 completely new MAGs, previously unrecovered from the metagenomic data sets (see Data set S1 at <https://doi.org/10.6084/m9.figshare.17698631>). Additional results and discussion of these MAG results can be found in the supplemental material (Text S1).

It is important to note that MAGs represent the overall diversity of a mixture of strains in a microbial population at a specific site at a specific point in time. Therefore, the terms “MAG” and “population” are used interchangeably or in conjunction with each other for the remainder of this article. Many of the methods described below require a robust signal over many genomic sites to allow for interpretation, such that the analyses conducted here should not have been strongly affected by incompleteness or small amounts of redundancy in the MAGs.

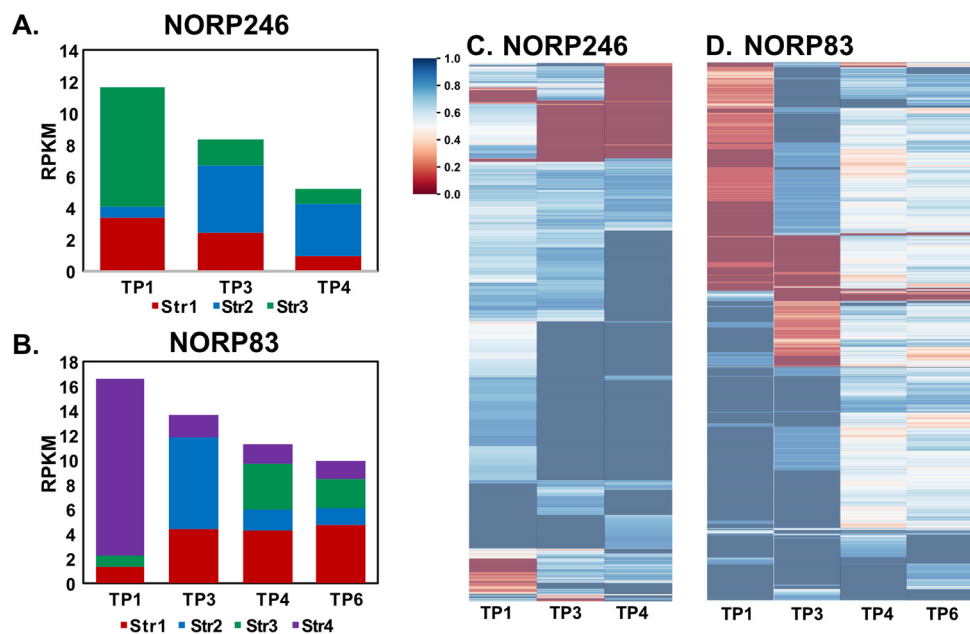
**Targeting of specific populations in North Pond fluids.** To conduct a robust analysis of genome dynamics at borehole U1382A over the 2-year time series, we focused on MAGs that occurred in at least three samples with  $\geq 5$  reads per kbp MAG per Mbp metagenomic sample (reads per kilobase per million [RPKM]) (Table 1; Fig. 1; see Data set S2 at <https://doi.org/10.6084/m9.figshare.17698631>). In total, six MAGs met these criteria. In addition, we identified another four MAGs with high coverage ( $>30$  RPKM) in only one sample. In all four instances, these MAGs had low coverage ( $<3$  RPKM) in the other time points, with three MAGs having extremely low coverage at other time points ( $<0.15$  RPKM) (Table 1; Fig. 1). We conducted an in-depth analysis of these 10 MAGs to understand the impacts of dispersal, selection, and recombination on microbial populations in the aquifer. Due to the complexity of fluid movement within the aquifer, it is challenging to determine whether rapid changes in cell abundance result from a dispersal event, whereby cells are flushed from the sampling location, or the result of “rapid” death and growth over the measured time points. Given the dynamism of the aquifer, we assumed that rapid shifts in population structure were likely



**FIG 1** Scatterplot of the abundance for each of the 10 MAGs of interest sampled from North Pond CORK U1382A. The size of each marker is scaled to the number of polymorphisms detected in the sample. Dotted lines highlighted RPKM values of 30 and 5, which were used as cutoffs to determine time points of interests. Markers with  $\geq 5$  RPKM are outlined with black lines.

associated with dispersal; however, given the small population sizes within these microbial communities, drift during bottleneck events remains likely.

**MAGs represent extant populations in the aquifer.** In general, the planktonic component of the microbial community within the aquifer has low cell density ( $10^3$  to  $10^4$  cell  $\text{mL}^{-1}$ ) (6, 7, 22). Under these low biomass conditions, we hypothesized that some of the microbial populations represented by the MAGs may be derived from organisms that had undergone a recent clonal expansion. A clonal expansion occurs when a population undergoes rapid growth, usually due to expansion into a new niche, and subsequently does not have enough time to accumulate mutations within the population (25). Few previous studies have defined metrics for what constitutes a clonal expansion in natural populations, and metrics will vary depending on methodology for assembly, binning, mapping, and defining single nucleotide variants (SNVs). Here, we assessed all MAGs for features indicative of a recent clonal expansion by searching for rapid shifts in abundance accompanied by extremely low SNV density relative to other populations. Each of the MAGs had rapid shifts in abundance, and we were particularly interested in the four MAGs (NORP6, NORP57, NORP100, and NORP167) that demonstrated a sudden increase in abundance at a single time point, with extremely low abundance at other time points (Table 1; Fig. 1). The number of detected polymorphisms, or single nucleotide variants per kilobase pair (SNVs  $\text{kbp}^{-1}$ ), nucleotide diversity, and evidence of active replication were used as metrics to determine if any of the populations were candidates for clonal expansion (see Data set S3 and 4 at <https://doi.org/10.6084/m9.figshare.17698631>). We calculated trough-to-peak ratios (T:P), which can be used as an estimate of genome replication from MAGs (26). It should be noted that while T:P ratios have been validated with pure cultures grown in chemostats (27), T:P ratios do not always correlate with cell growth rates in mixed communities (28). Based on genome coverage T:P ratios calculated for the North Pond MAGs, most MAGs in the time points of interest had evidence of active replication, with between 1.3 and 6.8. Despite drastic changes in abundance between time points (increasing from  $\sim 0$  to 38.3 and 88.6 RPKM, respectively; Table 1), a T:P ratio for NORP6 and NORP57 could not be determined by the tool iRep. Clonal expansions have been rarely observed in nature. The only defined clonal population from a previous study,



**FIG 2** Strain relative abundance as determined using DESMAN for NORP246 (A) and NORP83 (B). Major allele frequency for SNVs detected in time points of interest for NORP246 (C) and NORP83 (D). The major allele frequencies have been hierarchically clustered and scaled from 0 to 1. Both 0.0 and 1.0 represent a fixed allele; 0.5 represents a split allele.

which was also based on MAGs, observed a population with 0.003 SNVs  $\text{kbp}^{-1}$  polymorphisms (29) in a freshwater environment. Three of the North Pond MAGs (NORP6, NORP57, and NORP147) had notably few polymorphisms, ranging from 0.2 to 1.3 SNVs  $\text{kbp}^{-1}$ , but all others had more (5.7 to 63.8 SNVs  $\text{kbp}^{-1}$ ; Fig. 1; Data set S3). This level of population diversity is approximately equivalent to that observed for soils in grassland meadows (ranging between 4.7 and 43.0 SNVs  $\text{kbp}^{-1}$ ) and higher than single species-dominated photosynthetic mats in a connected watershed (ranging between 0.01 and 8.7 SNVs  $\text{kbp}^{-1}$ ) (11, 24). Based on the observed nucleotide diversity of the North Pond MAGs across all relevant time points ( $\pi = 0.002$  to 0.01), each population is diverse and shares minor alleles. The accumulation of diversity within the populations of interest, as measured through nucleotide diversity, suggests that there were no large-scale genome sweeps shortly before sampling. Collectively, this evidence suggests that none of the populations that met the criteria for analysis were candidates of recent clonal expansion. Instead, we observed high underlying diversity indicative of a complex system with coexisting subpopulations, approximately equivalent to subspecies or strains.

**Nucleotide-level heterogeneity reveals coexistence of multiple strains within populations.** For the six MAGs with high abundance across multiple time points, it was possible to compare patterns in allele frequency and reconstruct strain abundance patterns, revealing the coexistence of multiple subpopulations/strains. The MAG populations represented by NORP169 and NORP246 are estimated to contain approximately three strains each (Fig. 2; Fig. S1 in the supplemental material). There were pronounced shifts in the relative abundance of the three strains between time points. For example, in the NORP246 population, strain Str1 is  $\sim 65\%$  of the population in TP1, while in TP3 strain Str2 becomes the dominant strain at  $\sim 51\%$  of the population and reaches  $\sim 62\%$  in TP4 (Fig. 2A).

The other four MAG populations (NORP83, NORP139, NORP147, and NORP163) differed in that during the time points of interest, at least one strain enters the population that was not present at other time points. For example, the NORP139 population consisted of six strains (Str1 to Str6; Fig. S1). The NORP139 strain Str1 was dominant in TP2

(~62% of the population) but essentially absent (<0.7%) for other time points. Distinctly, NORP147, in the first time point of interest, was dominated by a single strain, with two different strains appearing as a minor component (<9%) of the population over subsequent time points (Fig. S1). For all four of these populations, the introduction of these strains supports the existence of a larger metapopulation, or spatially delimited local populations connected by migration (30), for which specific members were observed at North Pond at discrete time points. We observed shifts in the major allele frequency corresponding to the appearance of these new strains. This was exemplified by the transition in strain membership for NORP83 from TP1 to TP3, where the dominant strain Str4 (~87% of the population) was supplanted by the introduction of strain Str2 (~54% of the population), and more than half of the alleles shift from being nearly fixed at one allele to another different allele state (Fig. 2B and D). Moreover, when all four strains are present in TP4 and TP6 in a plurality (14% to 48% of the population), most of the major allele frequencies are close to 0.5 (Fig. 2D). The shifts in relative strain abundance were reflected in the major allele frequency, where time points dominated by a single strain have more sites that are fixed at a single allele (TP1 and TP3 major allele frequency 83.1% and 83.4%, respectively), while time points with multiple strains in similar proportions tend to have major allele frequencies closer to 0.5 (TP4 and TP6 major allele frequency 62.5% and 64.6%, respectively; Fig. 2). The ecology of these populations, with multiple strains disappearing and appearing across time points and changes in major allele frequencies corresponding to these shifts, and the high underlying genetic diversity, without indications of recent population-wide removal of diversity, supports an argument that the larger aquifer metapopulation is the source of genetic diversity observed at the North Pond site. This high observed diversity suggests that these strains originate through evolutionary selection and drift elsewhere in the aquifer over time periods longer than those observed in this study.

**Homologous recombination is common and varies by population.** If population differentiation occurs elsewhere and throughout the aquifer, the degree of recombination between coexisting strains has important implications for speciation and gene transfer. Low rates of recombination could lead to sympatric speciation or more frequent genome-wide sweeps, whereas high rates of recombination between strains could hold lineages together and also lead to gene-specific sweeps and facilitate gene transfer. We used the decay rate of linkage disequilibrium and the four-gamete test as a proxy for detecting recombination within the populations represented by each MAG. Linkage disequilibrium is a measurement that uses the distance between any two allelic variants as a proxy for recombination; variants in close proximity are less likely to become unlinked via recombination. As the recombination rate increases, the number of paired variants decreases as distance between variants increases (see Materials and Methods). Similarly, the four-gamete test is a measure of recombination that searches for pairs of segregating sites (see Materials and Methods). All of the populations had four alleles (denoted as H4; the four-gamete test assumes that the presence of the AB, ab, Ab, and aB haplotypes can only occur through homologous recombination) at some of the linked biallelic SNV sites. Low H4 frequencies indicate low recombination rates. The mean H4 frequencies for each population ranged from 0.44% to 14.84%. Several of the MAG populations (NORP83, NORP147, and NORP169) had H4 frequencies (0.44%, 1.17%, and 1.45%, respectively) below the lowest values estimated for single species photosynthetic mats (3.4% H4 frequency [11]; Table 2). NORP147 did not display decay of the linkage disequilibrium values as genetic distance increased for any of the mutation type pairs, indicating a lack of recombination within the population (S-S, N-S, or N-N; Fig. 3A). NORP83 displayed linkage disequilibrium decay (Fig. 3B), although the rate of decay was substantially lower than in the other North Pond populations (e.g., NORP246; Fig. 3C; Fig. S2), which is an indication of lower recombination rates. NORP57 did not display linkage decay (Fig. 3D) but had a slightly larger H4 frequency (5.40%). Both NORP83 and NORP169 had sufficient data to produce an estimate of recombination-to-mutation ( $\gamma/\mu$ ) from the tool mcorr, which had values



**TABLE 2** Percent mean four-gamete test results for all time points of interest<sup>d</sup>

MAG ID	H1 (%)	H2 (%)	H3 (%)	H4 (%)	Mean no. of biallelic sites
NORP83	0.11	84.61	14.85	0.44	25,653
NORP139	0.11	42.96	51.43	5.49	15,355
NORP147	0.00	76.48	22.35	1.17	452
NORP163	0.30	53.04	43.14	3.52	22,613
NORP169	0.13	76.76	21.66	1.45	5,132
NORP246	0.07	46.59	49.80	3.54	43,976
NORP6	0.45	23.58	61.13	14.84	3,558
NORP57	1.03	62.82	30.75	5.40	1,353
NORP100	2.83	63.73	28.14	5.29	11,680
NORP167	0.17	62.62	31.86	5.35	21,150

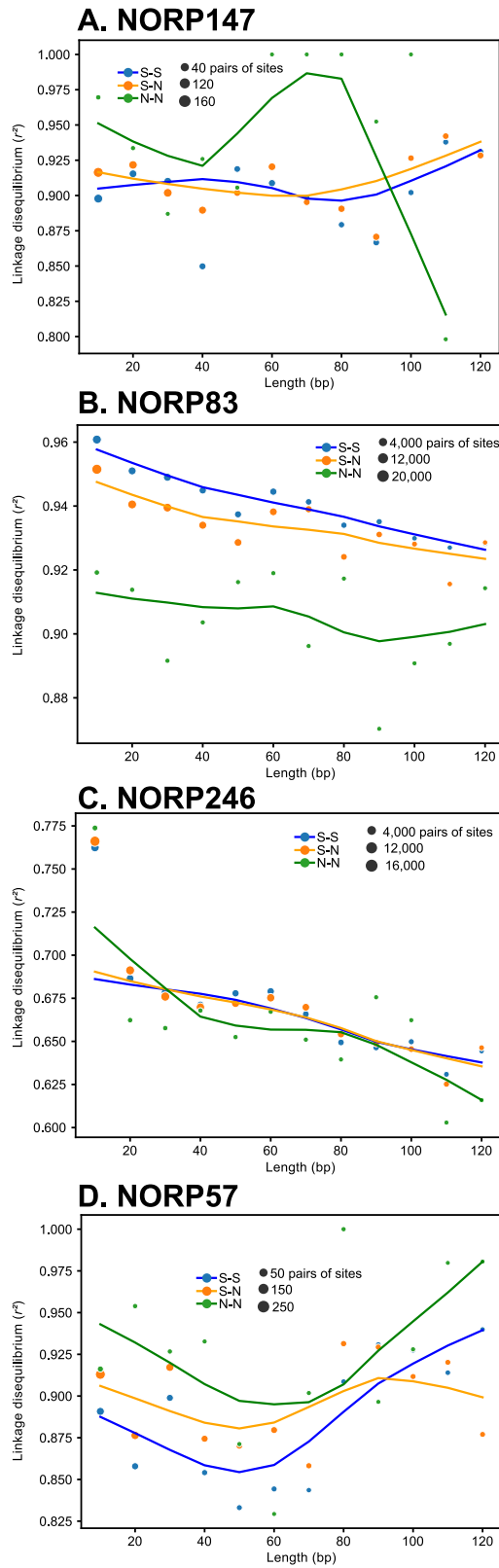
<sup>d</sup>Four-gamete tests for NORP6, NORP57, NORP100, and NORP167 were only computed for one time point.

(NORP83<sub>γ/μ</sub>: 1.10 to 4.59 and NORP169<sub>γ/μ</sub>: 4.71 to 8.25) similar to those of human pathogens (24, 31) (see data set S5 at <https://doi.org/10.6084/m9.figshare.17698631>). Collectively, this evidence suggests the NORP83, NORP57, NORP147, and NORP169 populations did not undergo high rates of recombination relative to other North Pond populations. For NORP147, the limited number of coexisting strains would reduce the opportunity for recombination, but for NORP83 and NORP169, which have multiple coexisting strains, this may imply that intrinsic biological characteristics prevented high rates of recombination within these populations. Diminished rates of recombination could increase the potential for sympatric speciation and genome-wide sweeps.

Conversely, the other six MAG populations had distinct patterns of linkage decay and higher frequencies of the H4 allele, suggesting likely recombination; however, the degree of linkage disequilibrium varied between them (Table 2; Fig. S2). For all of these populations, homologous recombination could maintain cohesion among the coexisting strains, preventing sympatric speciation while potentially facilitating gene-specific sweeps and possibly gene exchange among closely related strains. The balance between recombination and selection dictates the degree to which coexisting microbial lineages can speciate in the absence of geographical barriers (32–34). Recombination can play a cohesive role by maintaining sequence clusters and preventing speciation, particularly if the effects of recombination are higher than that of mutation (35). If recombination rates are low, clusters of organisms can form from selective sweeps, resulting in “ecotypes” that are adapted to specific niches (36). These ecotypes are maintained by periodic selection for or against specific mutations. However, if recombination rates are high enough, specific genes could sweep through populations without initiating a genome-wide sweep (29, 32). Our results suggest that recombination is common but populations likely do not reach panmictic equilibrium in the aquifer, in which recombination occurs frequently enough such that alleles are not linked across microbial genomes. Thus, both genome-wide and gene-specific sweeps are possible among these microbial populations.

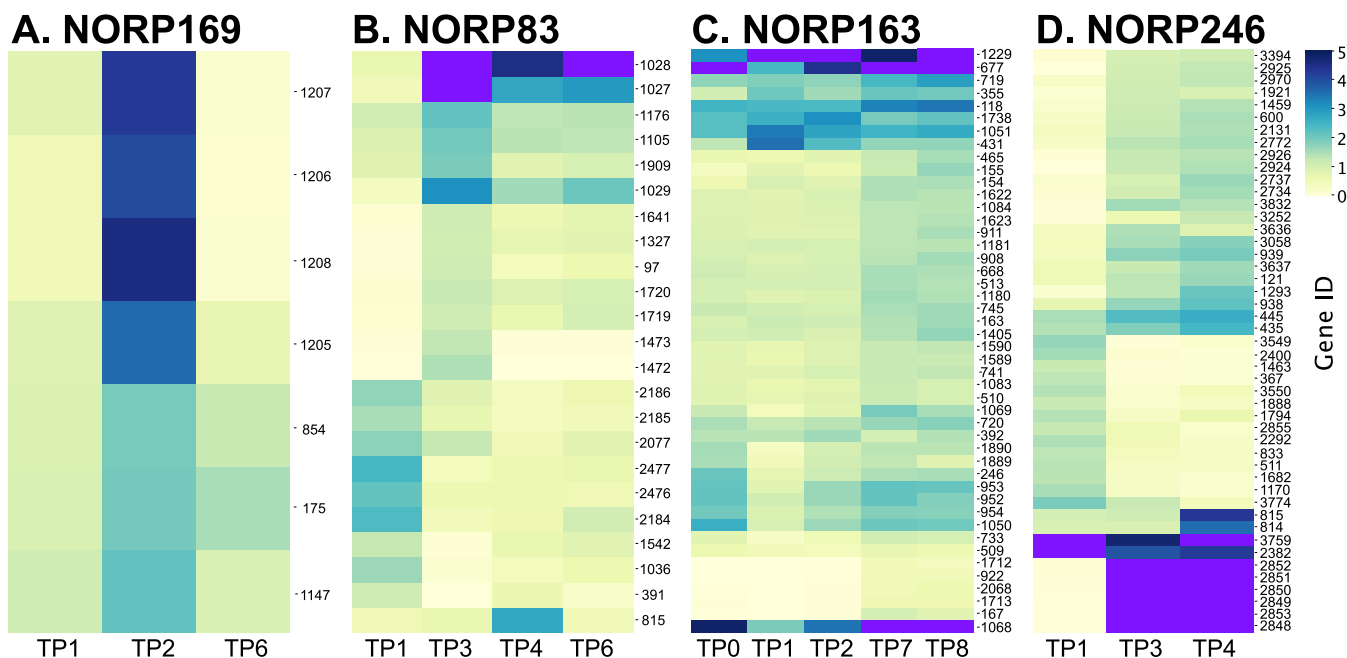
**Gene content variation across strains.** Genes that are unique to individual strains may provide a selective advantage. Genes limited to specific strains were identified by calculating changes in gene copy number based on observed changes in coverage over time ( $\geq 1 \times$  change in coverage between time points; see data set S6 at <https://doi.org/10.6084/m9.figshare.17698631>). In order to identify genes with abnormal coverage relative to each other, we used a minimum contig length of 2.5 kb to decrease the likelihood that any contig contained a single open read framing (ORF). As such, genes with variable coverage tend to be embedded on contigs between genes that were near the median genomic value.

The NORP147 population, dominated by a single strain at all three time points of interest, did not have genes with variable copy number, and NORP139, which had higher comparative rates of recombination, only had five genes with variable gene copy number, for which four had no KEGG annotations. We examined the four other populations



**FIG 3** Linkage disequilibrium of  $r^2$  for linked SNVs pairs for NORP147 (A), NORP83 (B), NORP246 (C), and NORP57 (D). Each circle is the mean  $r^2$  for pairs of linked SNVs at that distance range (e.g., 1–10 bp, 11–20 bp, etc.), with the area proportional to the number of linked SNVs in the mean. Linked SNVs are denoted by their predicted mutation type (nonsynonymous [N] and synonymous [S]).





**FIG 4** Hierarchically clustered heat map of gene frequencies for genes  $\geq 1\times$  change in coverage over the time points of interest for NORP169 (A), NORP83 (B), NORP163 (C), and NORP246 (D). Scaled range from 0 to  $5\times$  copies. Values that exceed  $5\times$  are denoted as purple in color. Maximum value of gene frequency for NORP246 (D) is  $21.7\times$  copies. Data are available in Data set S6 at <https://doi.org/10.6084/m9.figshare.17698631>.

to understand the putative functions of the genes that had variable copy numbers over time.

For the NORP169 population, none of the seven variable genes had KEGG annotations, but several of the functional assignments from GenBank implicated transcriptional regulation possibly related to plasmid replication (Fig. 4A). Four of the genes were colocalized within the MAG (NORP169 Gene IDs 1205 to 1208) and shifted from  $<1\times$  copies per genome in TP1 and TP6 to  $\sim 4\times$  copies per genome in TP2. This small shift corresponded with an increase in strain Str1 and may reflect an extrachromosomal element.

In NORP83, 13 of the 23 variable genes (57%) had no KEGG annotation. We observed shifts in gene abundance likely associated with the gradual loss of strain Str4 and the introduction of strain Str2 (Fig. 2B and C and Fig. 4B). Two sets of genes (NORP83 Gene IDs 2184 to 2186 and 2476 to 2477) decreased from  $\sim 2\times$  copies per genome in TP1 to  $<1\times$  copies in TP3 and were annotated as part of a type IV secretion system (IDs 2184 to 2186) and chromosome partitioning (ID 2477), which may indicate a role in horizontal gene transfer (37). Conversely, two sets of genes were absent in TP1 (IDs 1472 to 1473, a putative sodium:proton antiporter and serine hydrolase; IDs 1719 to 1720, sugar permease) and had  $>1\times$  copies per genome in TP3 (Fig. 4B). A set of genes (IDs 1027 to 1029) present at  $<0.5\times$  copies per genome in TP1 and annotated as putative transposases increased in gene copy in TP3 (5 to  $9\times$  copies per genome) before dropping to a lower, but still elevated level compared to TP1 in TP4 and TP6 (2 to  $6\times$  copies per genome; Fig. 4B). These patterns of gene frequency suggest that either these transposases were introduced into the North Pond NORP83 population with the arrival of strain Str2, or these transposases were introduced into other strains, resulting in the observed increase in copy number.

The NORP163 population had 45 putative genes that changed in gene copy number over time (Fig. 4C). While 29 of the variable genes (64%) had no KEGG annotations, 12 of these genes were putative transposases, with several exceeding  $>4\times$  copies per genome in all of the time points (Fig. 4C). There were five genes that were absent in TP0-TP2 and had  $>0\times$  copies per genome in TP7 and TP8 that were annotated as

ribosomal proteins (IDs 167 and 2068), hypothetical proteins (IDs 922 and 1712), and a sulfite reductase (ID 1714). A larger set of genes changed from  $<1\times$  copies per genome to  $>1\times$  copies from TP0 to TP8 (Fig. 4C). Colocalized genes of this larger set had functions related to core carbon metabolism (Gene IDs 1083 to 1084), heavy metal sensing and transcriptional response (IDs 1589 to 1590), and osmolyte transport (IDs 1622 to 1623). These genes that increased in copy number represent putative functions associated with environmental sensing and response through regulation and substrate transport.

Shifts in gene abundance for the population represented by NORP246 were highly dynamic over time and accompanied by a change in strain abundance, most notably the increase of strain Str2 over time (Fig. 2A and Fig. 4D). Forty-seven putative genes changed substantially in frequency over three time points in the NORP246 population (Fig. 4D). These genes can be divided into three broad groups: (i) present in TP1 and absent in TP3 and TP4, (ii) absent in TP1 and present in TP3 and TP4, and (iii) absent in TP1 and present in high coverage in TP3 and TP4. These most likely reflect genes that were either present or absent in strain Str2. The 13 putative genes in group one were present in TP1 ( $\sim 1\times$  copies per genome) and absent in TP3 and TP4, likely representing genes that were absent in strain Str2. These were predominantly annotated as hypothetical proteins with some conserved domains ( $n = 7$ ; Fig. 4D). Three of the putative genes had roles in transport: xanthine permease (ID 1682), sodium/glutamate symporter (ID 833), and cobalt ATP binding cassette-type permease (ID 2292). Two colocalized genes in this group were annotated as a transposase and reverse transcriptase (IDs 3549 to 3550). Conversely, the 23 putative genes in group two were absent in TP1 ( $<1\times$  copies per genome) and present in TP3 and TP4 (1 to  $2\times$  copies per genome), likely reflecting genes that were present in strain Str2 and absent in strains Str1 and Str3. These were predominantly annotated as hypothetical proteins with conserved domains ( $n = 18$ ; Fig. 4D). Genes in group three, which were absent in TP1 and present in high coverage in TP3 and TP4 (15 to  $20\times$  copies per genome), likely representing genes in high abundance in strain Str2, were colocalized (IDs 2848 to 2853) and lack annotations, except for gene ID 2851 with homology to proteins in GenBank annotated as DNA replication protein (Fig. 4D). This segment appeared to contain a fragment of DNA that in TP3 and TP4 could be independently replicating relative to the NORP246 genome, possibly a plasmid or similar extrachromosomal element. The lack of annotations makes it impossible to determine the exact roles of these genes within the TP3 and TP4 populations, but it is possible that selection favored the functional attributes of strain Str2 over that of strains Str1 and Str3.

There were three general trends in the types of genes that had gene copy variation: (i) genes lacking KEGG annotations, for which function remains unknown; (ii) genes annotated as and associated with extrachromosomal replication (e.g., plasmids, integrases, etc.) and horizontal gene transfer; and (iii) genes annotated as transporters. While variation in gene copy numbers over time for genes with functional annotations were identified, we cannot state with certainty that these genes provide a selective advantage for a particular variant, but some functions could provide specialization at the intraspecific level. For example, transporters play an important role in substrate availability and acquisition. The correlation between changes in gene frequency and the introduction/removal of a strain from the North Pond populations reinforces the role of dispersal as the likely mechanism for introducing genomic variation between strains of the same population. These populations may compete for and/or maintain partially overlapping ecological niches, such that changes in gene content may reflect the generation of "ecotypes" among the cooccurring strains. In general, the introduction of new strains did not correspond to an increase in total population abundance, perhaps reinforcing the idea that the intraspecific populations shared an ecological niche rather than directly competing for it (38). NORP83, NORP163, and NORP246 had multiple transposases that shifted in copy number over time. An increase in transposase activity has previously been identified as a signature of genetic drift (39). In this

**TABLE 3** Number of genes with decreased nucleotide diversity and regions detected with variable  $F_{ST}$ <sup>a</sup>

Genome	No. of CDS	No. of low nucleotide diversity CDS	Low dN/dS (<mean genome value) <sup>a</sup>	dN/dS > 1 <sup>b</sup>	No dN/dS <sup>b</sup>	Longest consecutive low nucleotide diversity CDS	No. of elevated $F_{ST}$ regions (>1 std dev)	No. of elevated $F_{ST}$ regions (>2 std dev)
NORP83	2,657	60 (2%)	4 (7%)	9 (15%)	43 (72%)	4	8	0
NORP139	1,538	32 (2%)	6 (19%)	12 (38%)	10 (31%)	1	0	0
NORP147	2,650	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0	0	0
NORP163	2,182	190 (9%)	37 (19%)	11 (6%)	117 (62%)	5	14	1
NORP169	2,842	238 (8%)	19 (8%)	19 (8%)	175 (74%)	3	3	0
NORP246	3,851	22 (1%)	5 (23%)	3 (14%)	12 (55%)	2	0	0

<sup>a</sup>Percent calculated as no. of low nucleotide diversity CDS/total CDS.

<sup>b</sup>Percent calculated as fraction of CDS with low nucleotide diversity.

instance, we observed an increase in the number of reads mapping to genes identified as transposases, although it is unclear whether this represented the spreading of transposases to different strains or an increase in transposase copy number within a strain. Other genes that changed in abundance, such as putative plasmids and the type IV secretion system, represent “selfish gene” elements that can propagate between strains during mixing and recombination events. Potentially in a system like North Pond, where growth is slow and infrequent, the actions and dispersal of self propagating elements may play an important role in gene transfer events as they occur outside canonical windows of genomic replication.

**Signatures of gene-specific sweeps.** While differences in gene content may result in the evolution of overlapping ecotypes or eventually to sympatric speciation, the higher rates of recombination in some of the North Pond populations should produce evidence of gene-specific sweeps. Genomic regions containing gene-specific sweeps are characterized by low SNV density. Likely due to the increased cutoff values used to determine the location of SNVs, compared to Crits-Christoph et al. (24), we did not observe a correlation between nucleotide diversity and gene coverage or purifying selection. As such, we used nucleotide diversity from the six MAG populations present over multiple time points as a proxy for potential gene sweep events based on a statistically significant decrease in nucleotide diversity for individual genes compared to the genomic mean in all time points (Table 3; see data set S7 to S9 at <https://doi.org/10.6084/m9.figshare.17698631>).

Four of the MAG populations had a small percentage of putative coding sequences (0% to 2%) that could represent gene-sweep events (Table 3). Both NORP83 (2%) and NORP147 (0%) were previously identified as having low predicted rates of recombination. The low number of putative gene sweeps is likely a direct result of low recombination rates, as low rates of recombination would increase gene linkage and thus prevent individual genes from sweeping through the population. However, both NORP139 (2%) and NORP246 (1%) had small percentages of their respective coding sequences that may have represented gene-specific sweeps, but estimates of recombination rates for these populations were higher compared to NORP83, NORP147, and NORP169. Conversely, NORP163 (9%) and NORP169 (8%) had a larger percentage of genes involved in putative gene-sweep events, although NORP163 was identified as having a relatively low recombination rate. The contradiction between a low expected recombination rate and a larger number of putative gene sweeps is likely the result of how the four-gamete test and inStrain recombination rates are calculated. Both approaches use linked biallelic SNVs to estimate recombination rates, but recent gene sweeps, especially sweeps through multiple strains, would not have many SNVs. With the exception of NORP139, most of the putative gene sweeps (55% to 74%) did not produce a dN/dS value from the Gretel approach analyzing all potential haplotypes and isoforms (Table 3). A smaller percentage of putative gene sweeps were attributed to coding regions under purifying selection (7% to 23%;  $gene_{dN/dS} < mean\ genome_{dN/dS}$ ) and relaxed selection (6% to 38%;  $dN/dS > 1$ ) (Table 3). While this might suggest that most genes are not under selection, this methodology filtered out putative coding regions without multiple haplotypes/

isoforms, lacking mutations ( $dN = dS = 0$ ), and with short timescales since divergence ( $dS < 0.01$ ). These data could support the interpretation that the identified putative gene sweeps are recent and have yet to accumulate mutations or been acted upon by selection. Given the inferred slow growth rates among the North Pond microbial community, longer periods of time may be required to observe such mutations and selection in a population.

We also searched for evidence of genetic differentiation among subpopulations using the pairwise fixation index  $F_{ST}$  by identifying genomic regions of 5+ genes that had variable allele frequencies across populations (see Data set S10 at <https://doi.org/10.6084/m9.figshare.17698631>). Generally, this metric followed the same pattern observed using nucleotide diversity. NORP163 and NORP169 had several regions with elevated  $F_{ST}$  when considering a deviation greater than one standard deviation from the mean  $F_{ST}$  value. Despite its low estimated recombination rate, NORP83 had the second highest number of regions with increased  $F_{ST}$  ( $n = 8$ ). For all three genomes, none of the elevated  $F_{ST}$  regions shared gene content with the putative gene-sweep regions identified with nucleotide diversity. However, increasing the required deviation to greater than two standard deviations from the mean reduced the number of differentiated regions to a single region in the NORP163 population. Due to the constraints placed on calculating the  $F_{ST}$  metric, which requires genes in two time points to have corresponding biallelic SNVs, many regions were likely excluded from this  $F_{ST}$  analysis. Additionally, the  $F_{ST}$  calculations required at least five genes in a variable region, although results from the nucleotide diversity metric demonstrated that the longest stretch in all genomes of decreased diversity was five sequential genes (Table 3). Furthermore, the nucleotide diversity metric could be calculated for all genes and the corresponding reads for all time points that met the inStrain cutoffs. Very few genomic regions had significantly elevated  $F_{ST}$  results, especially when filtered using a more stringent cutoff, supporting that even though some putative gene-specific sweeps were detected in the populations, the observed low recombination rates likely limit the total number of events between strains.

Collectively, for all six of these populations across all time points, there were several KEGG functional BRITE categories that were over represented within the identified putative gene-sweep regions compared to the genomic averages, including tRNA biogenesis (BR03016), Ribosome (BR03011), Ribosome biogenesis (BR03009), Transcription factors (BR03000), and Chaperones and folding catalysts (BR03110) (see Data set S11 at <https://doi.org/10.6084/m9.figshare.17698631>). The low number of SNVs in these regions may arise from the fact that these are often highly conserved core genes, which is consistent with previous studies (40). However, on average, all of these categories, except Transcription factors, also had increased mean  $dN/dS$  values in the putative gene-sweep regions (0.67 to 1.69) compared to the genomic average (0.40 to 0.48). Although the number of putative gene-sweep regions only reflects a small percentage of the total genes in the six MAGs, it is interesting that several core functional categories have comparatively relaxed selection. This could imply that the accumulation of mutations within core functional genes does not impart a fitness cost, which could again be a result of the slow growth conditions experienced by these populations. From the perspective of all six populations, the system dynamics appear to have been molded by dispersal of strains across the aquifer, and the low to moderate rates of homologous recombination may have facilitated species cohesion and gene transfer among closely related populations.

**Conclusion.** The MAGs described here represent subseafloor microbial populations that exhibited substantial changes over the  $\sim 825$ -day sampling period. Thus, these data sets provide evidence that the marine crustal aquifer hosts a dynamic habitat in which microbial populations grow and decay and can be dispersed over monthly timescales. While rapid allele frequency shifts can be linked to different types of population interactions, such as sweeps, dispersal, and clonal expansion, dispersal appears to play an important role in structuring the most abundant populations in the crustal fluid samples. The shifts in allele frequencies and nucleotide diversity suggest that stochastic events, such as dispersal and the mixing of populations throughout the aquifer, also mold the evolutionary trajectories

of microbial populations in this habitat and that selection and drift occur on timescales that were not captured as part of this study (41). Though classically structured as the introduction or removal of mutations through selection and drift, evolution, broadly defined, reflects changes in allele frequency within a population. Thus, dispersal can introduce extant populations from elsewhere in the environment and the interaction of these newly introduced alleles with the environment reflects an evolutionary response. As such, we are able to observe microbial evolution in action in the fluids moving through the oceanic crust. The observed population shifts occur in the span of months, reflecting the highly dynamic nature of the system, both biologically and physically. The accumulation of extrachromosomal insertions, such as transposable elements and plasmids, suggests that relaxed selection pressures and slow generation times can impact these populations in a manner distinct from dispersal. In summary, the subseafloor aquifer of North Pond represents a highly dynamic habitat where evolution may be governed largely by the stochastic forces of dispersal.

## MATERIALS AND METHODS

**Sample collection and sequencing.** As has been described elsewhere (6, 19, 22), North Pond water samples were collected from two CORKS at holes U1382A and U1383C in 2012 and 2014. For sample time points (TP) during expeditions (TP0 in April 2012 and TP9 in April 2014), the Mobile Pumping System (MPS) attached to the ROV Jason II was used to collect samples through umbilical lines connecting the CORK platform to fluids in the aquifer. For samples collected between expeditions (TP1 to TP8), the battery-powered GEOMICROBESLED (42, 43) deployed on the CORK collected samples approximately every 2 months. In all instances, lines were purged to remove stagnant water prior to sampling and all samples were fixed with RNALater *in situ*. A complete set ( $n = 8$ ) of *in situ* samples were collected from U1382A and provided the necessary samples and resolution for downstream analysis (as described below). Previously described in Tully et al. (6), DNA was extracted from filters using a phenol chloroform method (44) and paired-end sequencing was performed on an Illumina HiSeq 1000 at the Marine Biological Laboratory. Raw sequences were quality controlled using Cutadapt (45) v1.7.1 (parameters: `-e 0.08 --discard-trimmed --overlap = 3`) and Trimmomatic (46) v0.33 (parameters: `PE SLIDINGWINDOW:10:28 MINLEN:75`).

**Metagenome-assembled genomes and manual curation.** To build upon the set of MAGs generated in Tully et al. (6), the quality controlled paired-end reads were assembled using three different approaches to maximize recovery of MAGs that had not previously been assembled into MAGs (Table S1 and S2 in the supplemental material). The MAGs included in this specific study had a minimum completion of 58.25% and a maximum redundancy of 4.42%. A detailed methodology has been provided in the supplemental material, and a comprehensive assessment of the final set of MAGs can be found in data set S1 (see data set S1 at <https://doi.org/10.6084/m9.figshare.17698631>).

**SNV identification and analysis.** *anvi'o* (47, 48) v5.0 contig databases were generated for all MAGs. FASTA files of the contigs of each MAG were converted to an *anvi'o* database (*anvi-gen-contigs-database --skip-mindful-splitting*). *Bowtie2* (49) v2.3.4.1 (parameters: `--no-unal`) was used to create recruitment profiles for each MAG in each sample. The output SAM format file was converted to a sorted BAM format file using *samtools* (50) v1.9 and reads with  $<95\%$  sequence identity over 75% of the alignment were filtered from the recruitment profile using *BamM* v1.7.3 (parameters: `--percentage_id 0.95 --percentage_aln 0.75`; <https://github.com/Ecogenomics/BamM>). Recruitment profiles were merged with *anvi-merge*. During the *anvi-profile* step, the flag `--profile-SCVs` was set to include identification of single nucleotide variants (SNVs) and single codon variants (SCVs). Each individual MAG *anvi'o* database had the collection value 'DEFAULT' and bin value 'EVERYTHING' added using *anvi-script-add-default-collection*. Using those collection and bin values, *anvi-gen-variability-profile* (default parameters; `--engine AA --include-contig-names --engine CDN`) to extract SNV, single amino acid variants (SAAVs), and SCVs, from the profiles of each MAG in each North Pond sample (51). The SCV and SAAV values were used to calculate pN/pS ratios using *anvi-script-calculate-pn-ps-ratio* (*anvi'o* v6.2). The sorted, filtered BAM format recruitment profiles described above were used to determine read counts for each contig using *featureCounts* (52) v1.5.3 (default parameters) as implemented within *Binsanity-profile* (53) v0.3.3 (default parameters). Read counts were converted to the normalized unit reads per kilobase pair MAG per Megabase pair of metagenomic sample (RPKM). Read counts were also used to determine the relative fraction of each MAG in each metagenomic sample.

Hole U1382A provided a high-resolution sample data set (10 time points over  $\sim 825$  days), and MAGs detected in these samples were selected for downstream analysis based on two criteria from the time series: (i) sufficient coverage ( $\geq 5$  RPKM) in  $\geq 3$  time points or (ii) high coverage ( $\geq 30$  RPKM) in at least one time point. Only time points that had  $\geq 5$  RPKM were considered for SNV calculations. This criterion was set to ensure that most base pairs in each MAG had at least  $20\times$  read coverage, which was the minimum coverage value to determine SNVs in all downstream analyses (see below).

As has been reported previously (29, 54, 55), MAGs tend to represent a cohesive population represented by a  $>95\%$  nucleotide identity boundary. To confirm that a majority of the signal for each population came from closely related organisms, RPKM values were recalculated and compared using recruited reads filtered to 99% identity over 75% of the length of the alignment using *BamM* (parameters: `--percentage_id 0.99 --percentage_aln 0.75`). Additionally, results from *inStrain* (56) were used to track what fraction of reads recruited at 95 to 100% identity (see below).



Two of the MAGs (NORP143 and NORP246) were taxonomically assigned to the genus *Shewanella* and shared 98.2% ANI, as determined with FastANI (54). For consideration of this analysis these MAGs were deemed to represent the same species-level population and NORP143 (67.24% complete, 9.87% redundancy) was removed from downstream analysis while NORP246 was retained (82.57% complete, 2.28% redundancy; Table 1). The features identified in NORP246 were confirmed by creating a new anvi'o database with recruited reads filtered to 99% identity over 75% of the length of the alignment (as above). To ensure consistency, analyses using anvi'o calculated values below were based on recruited reads filtered to 95% identity over at least 75% of the length of the alignment.

**Functional annotation and taxonomic assignment.** Protein sequences as determined within anvi'o were submitted to the GhostKOALA KEGG annotation and mapping service (57) (submitted October 23, 2018; parameters: genus\_prokaryotes + family\_eukaryotes). KEGG annotations were added to the anvi'o contig databases for each MAG using the Python script KEGG-to-anvio (<https://github.com/edgraham/GhostKoalaParser>) and anvi-import-functions (<http://merenlab.org/2018/01/17/importing-ghostkoala-annotations/>).

GTDB-Tk (58) v1.3.0 (classify\_wf default parameters) using database R95 was used to determine a taxonomic assignment for each MAG.

**Data transformations for analysis.** The SNV occurrence table generated by anvi'o was filtered by retaining all entries with  $>0$  entropy value, SNV positions  $\geq 20\times$  coverage in the time point(s) of interest, and if the departure from consensus was  $\geq 10\%$ . Putative SNVs at each time point were converted to single nucleotide variants per kilobase pair (SNVs  $\text{kb}^{-1}$ ) (23), a measure of polymorphisms within the population, based on the full length of the MAG. Time points for which a MAG was present at  $<5$  RPKM were not used to calculate SNVs  $\text{kb}^{-1}$  or considered in further analysis.

We used anvi'o to obtain metrics of microbial population dynamics based on MAGs. anvi'o is a widely used platform for microbial genomics and provides a baseline to determine polymorphisms in the underlying metagenomic reads. The criteria detailed here provide a conservative method to compare features such as a SNV frequency, major allele frequency, and gene coverage. For data processed with anvi'o, major allele frequency was determined by counting the frequency of nucleotides in the aligned reads, identifying the consensus allele, and dividing the number of occurrences of that allele by the total coverage for that site (i.e., the nucleotide in the reference MAG sequence was not considered in this calculation). For each MAG, mean major allele frequency was determined by averaging all major allele frequencies for each SNV across the entire MAG. For comparing time points, if a SNV was detected in one time point, but not another (and coverage  $\geq 20\times$ ), the positions for which no SNV was recorded were assumed to be fixed at the consensus allele. Hierarchical clustering was performed and plotted using seaborn (<https://zenodo.org/record/3767070#.Ys1u8XbMI6Y>) with the default clustermap settings (average distance linkage method and Euclidean distance metric).

**Strain identification and abundance.** DESMAN (59) was used to identify strains of MAGs based on SNV proportions in core genes. DESMAN offers insight into the estimated complexity of the populations captured at a time point for a particular MAG. This analysis calculates the ratio of alleles at each time point in a set of core genes in order to identify discrete haplotypes (referred to throughout as "strains") and to define an estimate of relative strain contribution to the total population. This disentangles the bulk signal provided by the allele frequencies and offers a hypothesis as to how and why specific alleles are changing. We used anvi'o to first identify clusters of orthologous genes (COGs) in the MAGs, as well as to identify SNV variants in a set of 36 universal single-copy COGs as defined by Alneberg et al. (60). The DESMAN haplotype identifier was computed five times with a varying number of haplotypes (2 to 8), and the posterior mean deviance was plotted to select the optimal haplotype number and best replicate run. Results from Gretel (61) were used to support the optimal haplotype number for each MAG when running DESMAN (see below). Only time points with sufficient coverage (as described above) were included in depictions of haplotype abundance across samples.

**Recombination.** We used both mcorr (31) and inStrain (56) to compile further metrics of population dynamics and to assess recombination rates. The software package mcorr was used to determine the recombination rate for the populations represented by the MAGs. mcorr utilizes a correlation matrix to link recombination events at the 300-bp scale and fits these events to a model in order to estimate recombination and mutation rates. The correlation matrix requires polymorphisms in close proximity to correlate in an expected manner, which can be confounded as the complexity of the underlying subpopulation increases and rate of recombination increases between subpopulations. Additionally, should the population in question not adhere to the model trained from laboratory experiments, estimates of recombination and mutation may not be accurate or recoverable. An mcorr correlation profile for each time point was constructed using filtered recruitment profiles generated with the script filter\_reads.py (parameters: -m 0.96 -q 2 -l 1500; [https://github.com/alexcritschristoph/soil\\_popgen/blob/de1e09dd8416131a6b5feba74f0c3a5747d38da1/inStrain\\_lite/inStrain\\_lite/filter\\_reads.py](https://github.com/alexcritschristoph/soil_popgen/blob/de1e09dd8416131a6b5feba74f0c3a5747d38da1/inStrain_lite/inStrain_lite/filter_reads.py)) and a GFF formatted gene prediction file generated by Prodigal (62) v2.6.3 (parameters: -p m -m) using mcorr-bam (default parameters). The correlation profiles were fit to the mcorr model using mcorr-fit (default parameters). The correlation profiles and calculated fit were assessed for clear evidence of monotonic decay and normalized distribution of residuals. The estimated ratio of recombination-to-mutation ( $\gamma/\mu$  [ $[\gamma/\mu]$ ]) was incorporated into further analysis if the bootstrapping mean was  $<2\times$  the estimated value (24). Values of  $\gamma/\mu$  that failed this cut-off were interpreted as evidence of recombination among multiple haplotypes that could not be resolved using the correlation profile approach.

In contrast to mcorr, the inStrain profile links polymorphisms over the length of the metagenomic read. Thus, complex populations with extensive recombination, which may have been missed through the use of the mcorr correlation matrix, could be resolved but limited by the length of the read. Along with the accompanying genome statistics, inStrain provides highly precise measures of the metapopulation.



inStrain was used to profile the underlying diversity on a genome-wide basis for each MAG and on a gene-by-gene basis using Prodigal predicted ORFs from each time point (inStrain profile; parameters:  $-l$  0.95  $-\text{min\_mapq}$  1  $-c$  20  $-f$  0.1). Parameter settings, including read percent identity (95%), coverage (20 $\times$ ), and minimum allele frequency (10%), were selected to ensure that SNV detection by inStrain and anvio were conducted with the same level of stringency. inStrain includes the ability to filter reads that mapped to multiple locations during recruitment and these reads were removed, if present. Output from inStrain was assessed for the presence of linkage disequilibrium ( $r^2$ ), the distribution of read recruitment identity, per gene and genome average nucleotide diversity ( $\pi$ ), which provides a probability for two reads having the same nucleotide at a particular position, and estimate of replication rate (56). Linkage decay, a signal of recombination, was observed by plotting  $r^2$  against the distance between linked SNVs; where distance between linked SNVs were grouped into 10-bp ranges and separated by the type of mutation within the coding region (synonymous-synonymous [S-S], synonymous-nonsynonymous [S-N], and nonsynonymous-nonsynonymous [N-N]). Additionally, the four-gamete test, another signal of recombination, was performed for all biallelic (AB and ab haplotypes) linked SNV locations (11). Under the infinite-site model, the presence of all four haplotypes (AB, ab, Ab, aB) can only be explained by at least one recombination event (63, 64). The frequency of the occurrence of the number of haplotypes (one, two, three, or four haplotypes present) between two sites was calculated. Additionally, to identify putative genes that had undergone a selective sweep by detecting genes with statistically significantly lower nucleotide diversity, the mean MAG nucleotide diversity was compared to the gene mean for all time point of interest with the Welch's  $t$  test. A corresponding dN/dS value was calculated, when applicable, as computed by Gretel (see below).

The SNV and linkage output from inStrain was used to compute  $F_{ST}$ , a measure of allele frequency differences between two populations. For all linked biallelic sites within putative coding regions that had a complete start and stop codon, a moving Hudson  $F_{ST}$  calculation (65, 66) was performed as implemented in the scikit-allele package v1.3.3 (<https://scikit-allele.readthedocs.io/>) consistent with Crits-Christoph et al. (24). Genes with coverage outside of two standard deviations of the genome mean were excluded. A pairwise calculation was performed between all time points of interest and a genome mean was reported. All remaining putative coding sequences were screened in a five gene window for elevated  $F_{ST}$  signal that exceeds either one or two standard deviations from the genome mean.

**Gene frequency.** For each MAG, gene coverage values were exported from the anvio databases using anvio-export-gene-coverage-and-detection. Gene frequency in a sample was calculated by dividing each gene coverage value by the median coverage value of all genes. Genes of interest were identified based on changes in gene frequency such that the maximum frequency value minus the minimum frequency value in the time points of interest was  $\geq 1$ . Genes identified in this manner lacking a KEGG annotation were compared against GenBank nr (67) using BLASTP (68) with an e-value cutoff of  $10^{-5}$ , focusing only on the top two hits. Hierarchical clustering was performed and plotted using seaborn (as above).

**Gene variants.** The software package Gretel (61) was used to reconstruct haplotypes on a gene-by-gene basis (in contrast to DESMAN, which recovers haplotypes or strains for whole genomes based on SNVs in core genes). Gretel greedily uses the polymorphism patterns of the recruited reads to reconstruct all possible haplotypes for the genes within a MAG. This approach can lead to chimeric reconstructed genes, as clearly demonstrated by the presence of haplotypes within internal stop codons but also offers the full potential space of haplotypes for a gene within a time point. We simplified that potentially inflated value by converting the haplotypes to isoforms; this constrained estimate and describes a reduced space of environmental proteins. The combined mean dN/dS for each gene extended our ability to assess selection within the population. The filtered recruitment profiles for each MAG and time of interest was used to construct a variant call format (VCF) file using gretel-snpper (default parameters), compressed with bgzip (default parameters), and indexed with tabix (default parameters). Gene positions on each contig were determined from the Prodigal gene predictions and used as inputs for the analysis by gretel (default parameters). If a gene was determined to be in the reverse direction, putative haplotypes were reverse complemented and all correctly oriented genes were translated to proteins using seqmagick v0.8.4 (<https://fhrcr.github.io/seqmagick/>). Haplotypes that translated to isoforms with internal stop codons were excluded from further analysis. Haplotypes and corresponding isoforms for each putative gene from all time points of interest were combined. CD-HIT (69) was used to cluster haplotypes and isoforms with 100% identity (parameters:  $-c$  1  $-d$  0  $-g$  1). The combined isoforms were aligned using MUSCLE v3.8.31 (70) (default parameters) and the alignments were used as inputs to PAL2NAL (71) (default parameters) to create a codon alignment of the corresponding haplotypes. dN, dS, and dN/dS was calculated (72) for all genes with detected haplotypes using PAML (73) (See example control file at [https://github.com/bjtully/northpond\\_evolution/tree/main/dNdS-calc](https://github.com/bjtully/northpond_evolution/tree/main/dNdS-calc)). The output was filtered to include only  $dN > 0$ ,  $0.01 < dS < 1$ , and  $dN/dS < 5$ . dN/dS and dS for all pairwise haplotype comparisons were averaged. The distribution of the number of isoforms predicted for each gene was used to inform the selection of target strains in the DESMAN analysis.

**Data availability.** Data for this project can be found deposited at DDBJ/ENA/GenBank under the BioProject accession no. PRJNA391950. Raw sequence reads are available through the Short Read Archive with the accession no. SRX3143886-SRX3143902. Raw sequence reads from Meyer et al. (22) constituting the metagenomic samples from 2012, are available under the BioProject accession no. PRJNA280201. Files used to perform this research are available through figshare, including the contigs, anvio protein calls, Prodigal protein calls, anvio database profiles for the MAGs, plus an additional 11 supplemental data files (<https://doi.org/10.6084/m9.figshare.17698631>), and the filtered BAM files used in this analysis (<https://doi.org/10.6084/m9.figshare.17701254>). As detailed in Data set S1, all new/updated MAGs have been submitted to NCBI under the BioProject accession no. PRJNA391950. Scripts created to

process data, perform analyses, and create plots have been made accessible through GitHub ([https://github.com/bjtully/northpond\\_evolution](https://github.com/bjtully/northpond_evolution)).

## SUPPLEMENTAL MATERIAL

Supplemental material is available online only.

**TEXT S1**, PDF file, 0.3 MB.

**FIG S1**, PDF file, 2.3 MB.

**FIG S2**, PDF file, 0.1 MB.

**TABLE S1**, PDF file, 0.03 MB.

**TABLE S2**, PDF file, 0.03 MB.

## ACKNOWLEDGMENTS

We thank the crews of the *R/V Merian* and *ROV Jason II*, Wolfgang Bach, Peter Girguis, Chih-Chiang Hsieh, Ulrike Jaekel, Beate Kraft, Huei-Ting Lin, Beth Orcutt, Keir Becker, Stephanie Carr, and Heiner Villinger for support in accomplishing the 2012 and 2014 field programs. Ship time was provided by the German Science Foundation (DFG). Both Katrina Edwards and James Cowen's efforts were critical to the field component and success of this project. They are missed immensely. Leslie Murphy and Emily Reddington provided laboratory support at the WM Keck sequencing facility at the Marine Biological Laboratory.

This work was supported by NSF OCE-1062006, OCE-1745589, and OCE-1635208 to J.A.H. The Gordon and Betty Moore Foundation sponsored observatory components at North Pond through grant GBMF1609. The Center for Dark Energy Biosphere Investigations (C-DEBI) (OCE-0939564) supported J.A.H. and B.J.T. This is C-DEBI contribution 598.

Analyses were conducted by R.E.A., E.D.G., and B.J.T. E.D.G. performed the multitiered binning approach that generated the new MAGs. R.E.A. and B.J.T. performed the population analysis. R.E.A. and B.J.T. wrote the manuscript. All authors reviewed and edited the manuscript. J.A.H. was the lead researcher on establishing North Pond as a long-term research site and supported the collection and processing of all samples. The study was conceived by R.E.A. and B.J.T.

We declare no conflict of interest.

## REFERENCES

- Kallmeyer J, Pockalny R, Adhikari RR, Smith DC, D'Hondt S. 2012. Global distribution of microbial abundance and biomass in subseafloor sediment. *Proc Natl Acad Sci U S A* 109:16213–16216. <https://doi.org/10.1073/pnas.1203849109>.
- Parkes RJ, Cragg B, Roussel E, Webster G, Weightman A, Sass H. 2014. A review of prokaryotic populations and processes in sub-seafloor sediments, including biosphere:geosphere interactions. *Marine Geol* 352: 409–425. <https://doi.org/10.1016/j.margeo.2014.02.009>.
- Bradley JA, Arndt S, Amend JP, Burwicz E, Dale AW, Egger M, LaRowe DE. 2020. Widespread energy limitation to life in global subseafloor sediments. *Sci Adv* 6:eaba0697. <https://doi.org/10.1126/sciadv.aba0697>.
- Inagaki F, Hinrichs KU, Kubo Y, Bowles MW, Heuer VB, Hong W-L, Hoshino T, Ijiri A, Imachi H, Ito M, Kaneko M, Lever MA, Lin Y-S, Methe BA, Morita S, Morono Y, Tanikawa W, Bihan M, Bowden SA, Elvert M, Glombitza C, Gross D, Harrington GJ, Hori T, Li K, Limmer D, Liu C-H, Murayama M, Ohkouchi N, Ono S, Park Y-S, Phillips SC, Prieto-Mollar X, Purkey M, Riedinger N, Sanada Y, Sauvage J, Snyder G, Susilawati R, Takano Y, Tasumi E, Terada T, Tomaru H, Trembath-Reichert E, Wang DT, Yamada Y. 2015. Exploring deep microbial life in coal-bearing sediment down to 2.5 km below the ocean floor. *Science* 349:420–424. <https://doi.org/10.1126/science.aaa6882>.
- D'Hondt S, Inagaki F, Zarikian CA, Abrams LJ, Dubois N, Engelhardt T, Evans H, Ferdelman T, Gribsholt B, Harris RN, Hoppie BW, Hyun J-H, Kallmeyer J, Kim J, Lynch JE, McKinley CC, Mitsunobu S, Morono Y, Murray RW, Pockalny R, Sauvage J, Shimon T, Shiraishi F, Smith DC, Smith-Duque CE, Spivack AJ, Steinsbu BO, Suzuki Y, Zpak M, Toffin L, Uramoto G, Yamaguchi YT, Zhang G-I, Zhang X-H, Ziebis W. 2015. Presence of oxygen and aerobic communities from sea floor to basement in deep-sea sediments. *Nature Geosci* 8:299–304. <https://doi.org/10.1038/ngeo2387>.
- Tully BJ, Wheat CG, Glazer BT, Huber JA. 2018. A dynamic microbial community with high functional redundancy inhabits the cold, oxic subseafloor aquifer. *ISME J* 12:1–16. <https://doi.org/10.1038/ismej.2017.187>.
- Trembath-Reichert E, Walter SRS, Ortiz MAF, Carter PD, Girguis PR, Huber JA. 2021. Multiple carbon incorporation strategies support microbial survival in cold subseafloor crustal fluids. *Sci Adv* 7:eabg0153. <https://doi.org/10.1126/sciadv.abg0153>.
- Kashtan N, Roggensack SE, Rodrigue S, Thompson JW, Biller SJ, Coe A, Ding H, Martinen P, Malmstrom RR, Stocker R, Follows MJ, Stepanauskas R, Chisholm SW. 2014. Single-cell genomics reveals hundreds of coexisting subpopulations in wild *Prochlorococcus*. *Science* 344:416–420. <https://doi.org/10.1126/science.1248575>.
- Denef VJ, Banfield JF. 2012. In situ evolutionary rate measurements show ecological success of recently emerged bacterial hybrids. *Science* 336: 462–466. <https://doi.org/10.1126/science.1218389>.
- Anderson RE, Reveillaud J, Reddington E, Delmont TO, Eren AM, McDermott JM, Seewald JS, Huber JA. 2017. Genomic variation in microbial populations inhabiting the marine subseafloor at deep-sea hydrothermal vents. *Nat Commun* 8:1114. <https://doi.org/10.1038/s41467-017-01228-6>.
- BoumaGregson K, CritsChristoph A, Olm MR, Power ME, Banfield JF. 2021. Microcoleus (Cyanobacteria) form watershedwide populations without strong gradients in population structure. *Mol Ecol* 31:86–103. <https://doi.org/10.1111/mec.16208>.
- Edwards KJ, Bach W, Klaus A. 2010. Mid-Atlantic Ridge flank microbiology: initiation of long-term coupled microbiological, geochemical, and hydrological experimentation within the seafloor at North Pond, western flank of the Mid-Atlantic Ridge, p 1–62. *In* IODP scientific prospectus 336. International Ocean Discovery Program, College Station, TX.

13. Price AN, Fisher AT, Stauffer PH, Gable CW. 2021. Numerical simulation of cool hydrothermal processes in the upper volcanic crust beneath a marine sediment pond: North Pond, North Atlantic Ocean. *J Geophysical Res Solid Earth* 127:e2021JB023158. <https://doi.org/10.1029/2021JB023158>.
14. Sclater JG, Jaupart C, Galson D. 1980. The heat flow through oceanic and continental crust and the heat loss of the Earth. *Rev Geophys* 18:269–311. <https://doi.org/10.1029/RG018i001p00269>.
15. Stein CA, Stein S. 1994. Constraints on hydrothermal heat flux through the oceanic lithosphere from global heat flow. *J Geophys Res* 99:3081–3095. <https://doi.org/10.1029/93JB02222>.
16. Johnson HP, Pruis MJ. 2003. Fluxes of fluid and heat from the oceanic crustal reservoir. *Earth and Planetary Science Lett* 216:565–574. [https://doi.org/10.1016/S0012-821X\(03\)00545-4](https://doi.org/10.1016/S0012-821X(03)00545-4).
17. Wheat CG, Jannasch HW, Kastner M, Hulme S, Cowen J, Edwards KJ, Orcutt BN, Glazer B. 2011. Fluid sampling from oceanic borehole observatories: design and methods for CORK activities (1990–2010). In *Proceedings of the IODP. International Ocean Discovery Program, College Station, TX*.
18. Davis EE, Becker K, Pettigrew T, Carson B, MacDonald R. 1992. CORK: a hydrologic seal and downhole observatory for deep-ocean boreholes, p 43–53. In *Proceedings of the Ocean Drilling Program, initial report 139. International Ocean Discovery Program, College Station, TX*.
19. Seyler LM, Trembath-Reichert E, Tully BJ, Huber JA. 2020. Time-series transcriptomics from cold, oxic subseafloor crustal fluids reveals a motile, mixotrophic microbial community. *ISME J* 15:1192–1206. <https://doi.org/10.1038/s41396-020-00843-4>.
20. Walter SRS, Jaekel U, Osterholz H, Fisher AT, Huber JA, Pearson A, Dittmar T, Girguis PR. 2018. Microbial decomposition of marine dissolved organic matter in cool oceanic crust. *Nature Geoscience* 11:1–8. <https://doi.org/10.1038/s41561-018-0109-5>.
21. Wheat CG, Becker K, Villinger H, Orcutt BN, Fournier T, Hartwell A, Paul C. 2020. Subseafloor crosshole tracer experiment reveals hydrologic properties, heterogeneities, and reactions in slowspreading oceanic crust. *Geochem Geophys Geosyst* 21:2900–2915. <https://doi.org/10.1029/2019GC008804>.
22. Meyer JL, Jaekel U, Tully BJ, Glazer BT, Wheat CG, Lin H-T, Hsieh C-C, Cowen JP, Hulme SM, Girguis PR, Huber JA. 2016. A distinct and active bacterial community in cold oxygenated fluids circulating beneath the western flank of the Mid-Atlantic ridge. *Sci Rep* 6:22541. <https://doi.org/10.1038/srep22541>.
23. Starnawski P, Bataillon T, Ettema TJG, Jochum LM, Schreiber L, Chen X, Lever MA, Polz MF, Jorgensen BB, Schramm A, Kjeldsen KU. 2017. Microbial community assembly and evolution in subseafloor sediment. *Proc Natl Acad Sci U S A* 114:2940–2945. <https://doi.org/10.1073/pnas.1614190114>.
24. Crits-Christoph A, Olm MR, Diamond S, Bouma-Gregson K, Banfield JF. 2020. Soil bacterial populations are shaped by recombination and gene-specific selection across a grassland meadow. *ISME J* 14:1834–1846. <https://doi.org/10.1038/s41396-020-0655-x>.
25. Shapiro BJ, Polz MF. 2014. Ordering microbial diversity into ecologically and genetically cohesive units. *Trends Microbiol* 22:235–247. <https://doi.org/10.1016/j.tim.2014.02.006>.
26. Brown CT, Olm MR, Thomas BC, Banfield JF. 2016. Measurement of bacterial replication rates in microbial communities. *Nat Biotechnol* 34:1256–1263. <https://doi.org/10.1038/nbt.3704>.
27. Korem T, Zeevi D, Suez J, Weinberger A, Avnit-Sagi T, Pompan-Lotan M, Matot E, Jona G, Harmelin A, Cohen N, Sirota-Madi A, Thaiss CA, Pevsner-Fischer M, Sorek R, Xavier RJ, Elinav E, Segal E. 2015. Growth dynamics of gut microbiota in health and disease inferred from single metagenomic samples. *Science* 349:1101–1106. <https://doi.org/10.1126/science.aac4812>.
28. Long AM, Hou S, Ignacio-Espinoza JC, Fuhrman JA. 2021. Benchmarking microbial growth rate predictions from metagenomes. *ISME J* 15:183–195. <https://doi.org/10.1038/s41396-020-00773-1>.
29. Bendall ML, Stevens SL, Chan L-K, Malfatti S, Schwientek P, Tremblay J, Schackwitz W, Martin J, Pati A, Bushnell B, Froula J, Kang D, Tringe SG, Bertilsson S, Moran MA, Shade A, Newton RJ, McMahon KD, Malmstrom RR. 2016. Genome-wide selective sweeps and gene-specific sweeps in natural bacterial populations. *ISME J* 10:1589–1601. <https://doi.org/10.1038/ismej.2015.241>.
30. Hanski I, Gaggiotti O. 2004. Part I: perspectives on spatial dynamics, p 3–22. In *Ecology, genetics and evolution of metapopulations*. Academic Press, New York.
31. Lin M, Kussell E. 2019. Inferring bacterial recombination rates from large-scale sequencing datasets. *Nat Methods* 16:199–204. <https://doi.org/10.1038/s41592-018-0293-7>.
32. Cadillo-Quiroz H, Didelot X, Held NL, Herrera A, Darling A, Reno ML, Krause DJ, Whitaker RJ. 2012. Patterns of gene flow define species of thermophilic archaea. *PLoS Biol* 10:e1001265. <https://doi.org/10.1371/journal.pbio.1001265>.
33. Mallet J. 2008. Hybridization, ecological races and the nature of species: empirical evidence for the ease of speciation. *Philos Trans R Soc Lond B Biol Sci* 363:2971–2986. <https://doi.org/10.1098/rstb.2008.0081>.
34. Felsenstein J. 1981. Skepticism towards santa rosalia, or why are there so few kinds of animals? *Evolution* 35:124–138. <https://doi.org/10.1111/j.1558-5646.1981.tb04864.x>.
35. Fraser C, Hanage WP, Spratt BG. 2007. Recombination and the nature of bacterial speciation. *Science* 315:476–480. <https://doi.org/10.1126/science.1127573>.
36. Cohan FM, Koeppel AF. 2008. The origins of ecological diversity in prokaryotes. *Curr Biol* 18:R1024–R1034. <https://doi.org/10.1016/j.cub.2008.09.014>.
37. Wallden K, Rivera-Calzada A, Waksman G. 2010. Microreview: type IV secretion systems: versatility and diversity in function. *Cell Microbiol* 12:1203–1212. <https://doi.org/10.1111/j.1462-5822.2010.01499.x>.
38. Garcia SL, Stevens SLR, Crary B, Martinez-Garcia M, Stepanauskas R, Woyke T, Tringe SG, Andersson SGE, Bertilsson S, Malmstrom RR, McMahon KD. 2018. Contrasting patterns of genome-level diversity across distinct co-occurring bacterial populations. *ISME J* 12:742–714. <https://doi.org/10.1038/s41396-017-0001-0>.
39. Escobar-Páramo P, Ghosh S, DiRuggiero J. 2005. Evidence for genetic drift in the diversification of a geographically isolated population of the hyperthermophilic archaeon *Pyrococcus*. *Mol Biol Evol* 22:2297–2303. <https://doi.org/10.1093/molbev/msi227>.
40. Moulana A, Anderson RE, Fortunato CS, Huber JA. 2020. Selection is a significant driver of gene gain and loss in the pangenome of the bacterial genus *Sulfurovum* in geographically distinct deep-sea hydrothermal vents. *mSystems* 5:e00673-19. <https://doi.org/10.1128/mSystems.00673-19>.
41. Schloissnig S, Arumugam M, Sunagawa S, Mitreva M, Tap J, Zhu A, Waller A, Mende DR, Kultima JR, Martin J, Kota K, Sunyaev SR, Weinstock GM, Bork P. 2012. Genomic variation landscape of the human gut microbiome. *Nature* 493:1–6. <https://doi.org/10.1038/nature11711>.
42. Cowen JP, Copson DA, Jolly J, Hsieh C-C, Lin H-T, Glazer BT, Wheat CG. 2012. Advanced instrument system for real-time and time-series microbial geochemical sampling of the deep (basaltic) crustal biosphere. *Deep Sea Res Part I* 61:43–56. <https://doi.org/10.1016/j.dsr.2011.11.004>.
43. Lin H-T, Hsieh C-C, Repeta DJ, Rappé MS. 2020. Sampling of basement fluids via Circulation Obviation Retrofit Kits (CORKs) for dissolved gases, fluid fixation at the seafloor, and the characterization of organic carbon. *MethodsX* 7:101033. <https://doi.org/10.1016/j.mex.2020.101033>.
44. Sogin ML, Morrison HG, Huber JA, Welch DM, Huse SM, Neal PR, Arrieta JM, Herndl GJ. 2006. Microbial diversity in the deep sea and the underexplored rare biosphere. *Proc Natl Acad Sci U S A* 103:12115–12120. <https://doi.org/10.1073/pnas.0605127103>.
45. Martin M. 2011. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet j* 17:10. <https://doi.org/10.14806/ej.17.1.200>.
46. Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics* 30:2114–2120. <https://doi.org/10.1093/bioinformatics/btu170>.
47. Eren AM, Esen OC, 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>.
48. Eren AM, Kiefl E, Shaiber A, Veseli I, Miller SE, Schechter MS, Fink I, Pan JN, Yousef M, Fogarty EC, Trigodet F, Watson AR, Esen OC, Moore RM, Clayssen Q, Lee MD, Kivenson V, Graham ED, Merrill BD, Karkman A, Blankenberg D, Eppley JM, Sjödin A, Scott JJ, Vázquez-Campos X, McKay LJ, McDaniel EA, Stevens SLR, Anderson RE, Fuessel J, Fernández-Guerra A, Maignien L, Delmont TO, Willis AD. 2020. Community-led, integrated, reproducible multi-omics with anvi'o. *Nat Microbiol* 6:1–4. <https://doi.org/10.1038/s41564-020-00834-3>.
49. Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nat Methods* 9:357–359. <https://doi.org/10.1038/nmeth.1923>.
50. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R, 1000 Genome Project Data Processing Subgroup. 2009. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* 25:2078–2079. <https://doi.org/10.1093/bioinformatics/btp352>.
51. Delmont TO, Kiefl E, Kilinc O, Esen OC, Uysal I, Rappé MS, Giovannoni S, Eren AM. 2019. Single-amino acid variants reveal evolutionary processes that shape the biogeography of a global SAR11 subclade. *Elife* 8:e46497. <https://doi.org/10.7554/eLife.46497>.

52. Liao Y, Smyth GK, Shi W. 2014. featureCounts: an efficient general purpose program for assigning sequence reads to genomic features. *Bioinformatics* 30:923–930. <https://doi.org/10.1093/bioinformatics/btt656>.
53. Graham ED, Heidelberg JF, Tully BJ. 2017. BinSanity: unsupervised clustering of environmental microbial assemblies using coverage and affinity propagation. *PeerJ* 5:e3035-19. <https://doi.org/10.7717/peerj.3035>.
54. 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>.
55. Rodriguez-R LM, Tsementzi D, Luo C, Konstantinidis KT. 2020. Iterative subtractive binning of freshwater chronoserics metagenomes identifies over 400 novel species and their ecologic preferences. *Environ Microbiol* 22:3394–3412. <https://doi.org/10.1111/1462-2920.15112>.
56. Olm MR, Crits-Christoph A, Bouma-Gregson K, Firek BA, Morowitz MJ, Banfield JF. 2021. inStrain profiles population microdiversity from metagenomic data and sensitively detects shared microbial strains. *Nat Biotechnol* 25:1–15. <https://doi.org/10.1038/s41587-020-00797-0>.
57. 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>.
58. 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>.
59. Quince C, Delmont TO, Raguideau S, Alneberg J, Darling AE, Collins G, Eren AM. 2017. DESMAN: a new tool for de novo extraction of strains from metagenomes. *Genome Biol* 18:1–22. <https://doi.org/10.1186/s13059-017-1309-9>.
60. Alneberg J, Bjarnason BS, Bruijn I. d, 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>.
61. Nicholls SM, Aubrey W, Grave KD, Schietgat L, Creevey CJ, Clare A. 2020. On the complexity of haplotyping a microbial community. *Bioinformatics* 37:1360–1366. <https://doi.org/10.1093/bioinformatics/btaa977>.
62. Hyatt D, LoCascio PF, Hauser LJ, Uberbacher EC. 2012. Gene and translation initiation site prediction in metagenomic sequences. *Bioinformatics* 28:2223–2230. <https://doi.org/10.1093/bioinformatics/bts429>.
63. Garud NR, Pollard KS. 2020. Population genetics in the human microbiome. *Trends Genet* 36:53–67. <https://doi.org/10.1016/j.tig.2019.10.010>.
64. Hudson RR, Kaplan NL. 1985. Statistical properties of the number of recombination events in the history of a sample of DNA sequences. *Genetics* 111:147–164. <https://doi.org/10.1093/genetics/111.1.147>.
65. Hudson RR, Slatkin M, Maddison WP. 1992. Estimation of levels of gene flow from DNA sequence data. *Genetics* 132:583–589. <https://doi.org/10.1093/genetics/132.2.583>.
66. Bhatia G, Patterson N, Sankararaman S, Price AL. 2013. Estimating and interpreting FST: the impact of rare variants. *Genome Res* 23:1514–1521. <https://doi.org/10.1101/gr.154831.113>.
67. 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>.
68. Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W, Lipman DJ. 1997. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res* 25:3389–3402. <https://doi.org/10.1093/nar/25.17.3389>.
69. Fu L, Niu B, Zhu Z, Wu S, Li W. 2012. CD-HIT: accelerated for clustering the next-generation sequencing data. *Bioinformatics* 28:3150–3152. <https://doi.org/10.1093/bioinformatics/bts565>.
70. 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>.
71. Suyama M, Torrents D, Bork P. 2006. PAL2NAL: robust conversion of protein sequence alignments into the corresponding codon alignments. *Nucleic Acids Res* 34:W609–W612. <https://doi.org/10.1093/nar/gkl315>.
72. Goldman N, Yang Z. 1994. A codon-based model of nucleotide substitution for protein-coding DNA sequences. *Mol Biol Evol* 11:725–736. <https://doi.org/10.1093/oxfordjournals.molbev.a040153>.
73. Yang Z. 2007. PAML 4: phylogenetic analysis by maximum likelihood. *Mol Biol Evol* 24:1586–1591. <https://doi.org/10.1093/molbev/msm088>.