

Patterns of gene content and co-occurrence constrain the evolutionary path toward animal association in CPR bacteria

Alexander L. Jaffe¹, Christine He², Ray Keren³, Luis E. Valentin-Alvarado^{1,2}, Patrick Munk⁴, Keith Bouma-Gregson^{5,6}, Ibrahim F. Farag⁷, Yuki Amano^{8,9}, Rohan Sachdeva^{2,5}, Patrick T. West¹⁰ and Jillian F. Banfield^{*2,5,11,12}

¹Department of Plant and Microbial Biology, University of California, Berkeley, Berkeley, CA

²Innovative Genomics Institute, University of California, Berkeley, Berkeley, CA

³Department of Civil and Environmental Engineering, University of California, Berkeley, Berkeley, CA

⁴National Food Institute, Technical University of Denmark, Kongens Lyngby, Denmark

⁵Department of Earth and Planetary Science, University of California, Berkeley, Berkeley, CA

⁶Department of Integrative Biology, University of California, Berkeley, Berkeley, CA

⁷School of Marine Science and Policy, University of Delaware, Lewes, Delaware 19958, USA

⁸Nuclear Fuel Cycle Engineering Laboratories, Japan Atomic Energy Agency

⁹Horonobe Underground Research Center, Japan Atomic Energy Agency

¹⁰Department of Medicine (Hematology & Blood and Marrow Transplantation), Stanford University, Stanford, CA

¹¹Department of Environmental Science, Policy, and Management, University of California, Berkeley, Berkeley, CA

¹²Chan Zuckerberg Biohub, San Francisco, CA

*Corresponding author email: jbanfield@berkeley.edu

ABSTRACT: Candidate Phyla Radiation (CPR) bacteria are small, likely episymbiotic organisms found across Earth's ecosystems. Despite their prevalence, the distribution of CPR lineages across habitats and the genomic signatures of transitions amongst these habitats remain unclear. Here, we expand the genome inventory for Absconditabacteria (SR1), Gracilibacteria, and Saccharibacteria (TM7), CPR bacteria known to occur in both animal-associated and environmental microbiomes, and investigate variation in gene content with habitat of origin. By overlaying phylogeny with habitat information, we show that bacteria from these three lineages have undergone multiple transitions from environmental habitats into animal microbiomes. Based on co-occurrence analyses of hundreds of metagenomes, we extend the prior suggestion that certain Saccharibacteria have broad bacterial host ranges and constrain possible host relationships for Absconditabacteria and Gracilibacteria. Full-proteome analyses show that animal-associated Saccharibacteria have smaller gene repertoires than their environmental counterparts and are enriched in numerous protein families, including those likely functioning in amino acid metabolism, phage defense, and detoxification of peroxide. In contrast, some freshwater Saccharibacteria encode a putative rhodopsin. For protein families exhibiting the clearest patterns of differential habitat distribution, we compared protein and species phylogenies to estimate the incidence of lateral gene transfer and genomic loss occurring over the species tree. These analyses suggest that habitat transitions were likely not accompanied by large transfer or loss events, but rather were associated with continuous proteome remodeling. Thus, we speculate that CPR habitat transitions were driven largely by availability of suitable host taxa, and were reinforced by acquisition and loss of some capacities.

IMPORTANCE: Studying the genetic differences between related microorganisms from different environment types can indicate factors associated with their movement among habitats. This is particularly interesting for bacteria from the Candidate Phyla Radiation because their minimal metabolic capabilities require symbiotic associations with microbial hosts. We found that shifts of Absconditabacteria, Gracilibacteria, and Saccharibacteria between environmental ecosystems and mammalian mouths/guts probably did not involve major episodes of gene gain and loss; rather, gradual genomic change likely followed habitat migration. The results inform our understanding of how little-known microorganisms establish in the human microbiota where they may ultimately impact health.

Keywords: CPR bacteria; habitat transition, animal microbiome; comparative genomics; bacterial evolution

1 INTRODUCTION:

2 The Candidate Phyla Radiation (CPR) is a phylogenetically diverse clade of bacteria
 3 characterized by reduced metabolisms, potentially episymbiotic lifestyles, and ultrasmall cells.
 4 While the first high-quality CPR genomes were primarily from groundwater, sediment, and
 5 wastewater (1–3), subsequently genomes have been recovered from diverse environmental and
 6 animal-associated habitats, including humans. Intriguingly, from dozens of major CPR lineages,
 7 only three – Candidatus Absconditabacteria (formerly SR1), Gracilibacteria (formerly BD1-5 and
 8 GN02), and Saccharibacteria (formerly TM7) – are consistently associated with animal oral
 9 cavities and digestive tracts (4). The Saccharibacteria are perhaps the most deeply studied of
 10 all CPR lineages to date, likely due to their widespread presence in human oral microbiomes
 11 and association with disease states such as gingivitis and periodontitis (5, 6). On the other
 12 hand, Absconditabacteria and Gracilibacteria remain deeply undersampled, potentially due to
 13 their rarity in microbial communities or their use of an alternative genetic code that may
 14 confound some gene content analyses (1, 7, 8).

15 Absconditabacteria, Gracilibacteria, and Saccharibacteria are predicted to be obligate
 16 fermenters, dependent on other microorganisms (hosts) for components such as lipids, nucleic
 17 acids, and many amino acids (1, 3). Despite a generally reduced metabolic platform, CPR
 18 bacteria display substantial variation in their genetic capacities, even within lineages (9, 10). For
 19 example, some Gracilibacteria lack essentially all genes of the glycolysis and
 20 pentose-phosphate pathways and the TCA cycle (11). In contrast to many CPR, soil-associated
 21 Saccharibacteria encode numerous genes related to oxygen metabolism (12, 13). Pangenome
 22 analyses have shown genetic evidence for niche partitioning among Saccharibacteria from the
 23 same body site (14). However, the lack of comprehensive genomic sampling of these three CPR
 24 lineages across habitats, particularly from environmental biomes, has left unclear the full extent
 25 to which CPR gene inventories vary with habitat type, and, relatedly, the extent to which
 26 changes in metabolic capacities might have been reshaped during periods of environmental
 27 transition. Of particular interest is whether rapid gene acquisitions (e.g., via lateral gene
 28 transfer) or losses enabled habitat switches, or if these changes occurred gradually following
 29 habitat change.

30 The availability of suitable hosts may also drive the colonization of new environments by CPR
 31 bacteria (14). While there has been significant progress in characterizing the relationship

1 between Saccharibacteria and Actinobacteria in the oral habitat (15–17), other CPR-host
2 relationships remain unclear. Elucidation of environmental transitions among CPR lineages will
3 require both thorough analysis of functional repertoires as well as a more comprehensive
4 understanding of associations with other microorganisms. Here, we improve existing sampling
5 of CPR genomes and their surrounding communities to examine patterns of distribution,
6 abundance, and gene content in different microbiome types. We also make use of
7 whole-community co-occurrence patterns to shed light on the potential host range of the CPR
8 bacteria in their associated ecosystems. In combination, our analyses shed light on the
9 frequency of habitat shifts in three CPR lineages and the evolutionary processes likely
10 underlying them.

11

12 **RESULTS:**

13 *Environmental diversity, phylogenetic relationships, and abundance patterns*

14 We gathered an environmentally comprehensive set of Absconditabacteria, Gracilibacteria, and
15 Saccharibacteria by querying multiple databases for genomes assembled in previous studies
16 and assembling new genomes from several additional metagenomic data sources (Table S1,
17 Materials and Methods) (1–4, 7, 12–15, 18–82). Quality filtration of this curated genome set at \geq
18 70% completeness and \leq 10% contamination and subsequent de-replication at 99% average
19 nucleotide identity (ANI) yielded a non-redundant set of 389 genomes for downstream analysis
20 (Table S1). Absconditabacteria and Gracilibacteria were less frequently sampled relative to
21 Saccharibacteria, comprising only ~7.5% and ~10.8% of the total genome set, respectively. All
22 three lineages were distributed across a broad range of microbiomes, encompassing various
23 environmental habitats (freshwater, marine, soil, engineered, plant-associated, hypersaline) as
24 well as multiple animal-associated microbiomes (oral and gut) (Fig. 1). Unlike animal-associated
25 Gracilibacteria and Absconditabacteria genomes, which were recovered only from human and
26 animal oral samples, animal-associated Saccharibacteria were found in both oral and gut
27 samples.

28 We extracted 16 syntenic, phylogenetically informative ribosomal proteins from each genome to
29 construct a CPR species tree and evaluate how habitat of origin maps onto phylogeny.

30 Sequences from related CPR bacteria were used as outgroups for tree construction (Materials

1 and Methods). The resolved topology supports monophyly of all three lineages and a sibling
2 relationship between the two alternatively coded lineages, Absconditabacteria and
3 Gracilibacteria (Figure 1a, File S1), consistent with previous findings for the CPR (10). For the
4 Absconditabacteria, a single clade of organisms derived from animal-associated microbiomes
5 was deeply nested within genomes from the environment. On the other hand, Gracilibacteria
6 clearly formed two major lineages (GRA 1-2), each with a small subclade comprised of
7 animal-associated genomes. For Saccharibacteria, deeply-rooting lineages were also almost
8 exclusively of environmental origin (soil, water, sediment) and animal-associated genomes were
9 strongly clustered into at least three independent subclades (Fig. 1a). Two of these three
10 subclades were exclusively composed of animal-associated sequences whereas one (SAC 5),
11 was a mixture of animal-associated, wastewater (potentially of human origin) and a few aquatic
12 sequences. Intriguingly, for both Saccharibacteria and Gracilibacteria, a subset of organisms
13 from the dolphin mouth (22) did not affiliate with those from terrestrial mammals/humans and
14 instead fell within marine/environmental clades (indicated by asterisks in Fig. 1). In primarily
15 environmental clades (SAC 1 and 4), genomes from soil, freshwater, engineered, and halophilic
16 environments were phylogenetically interspersed, suggesting comparatively wide global
17 distributions for these lineages. One exception to this pattern were two clades representing
18 distinct hypersaline environments – a hypersaline lake and salt crust (65, 70).

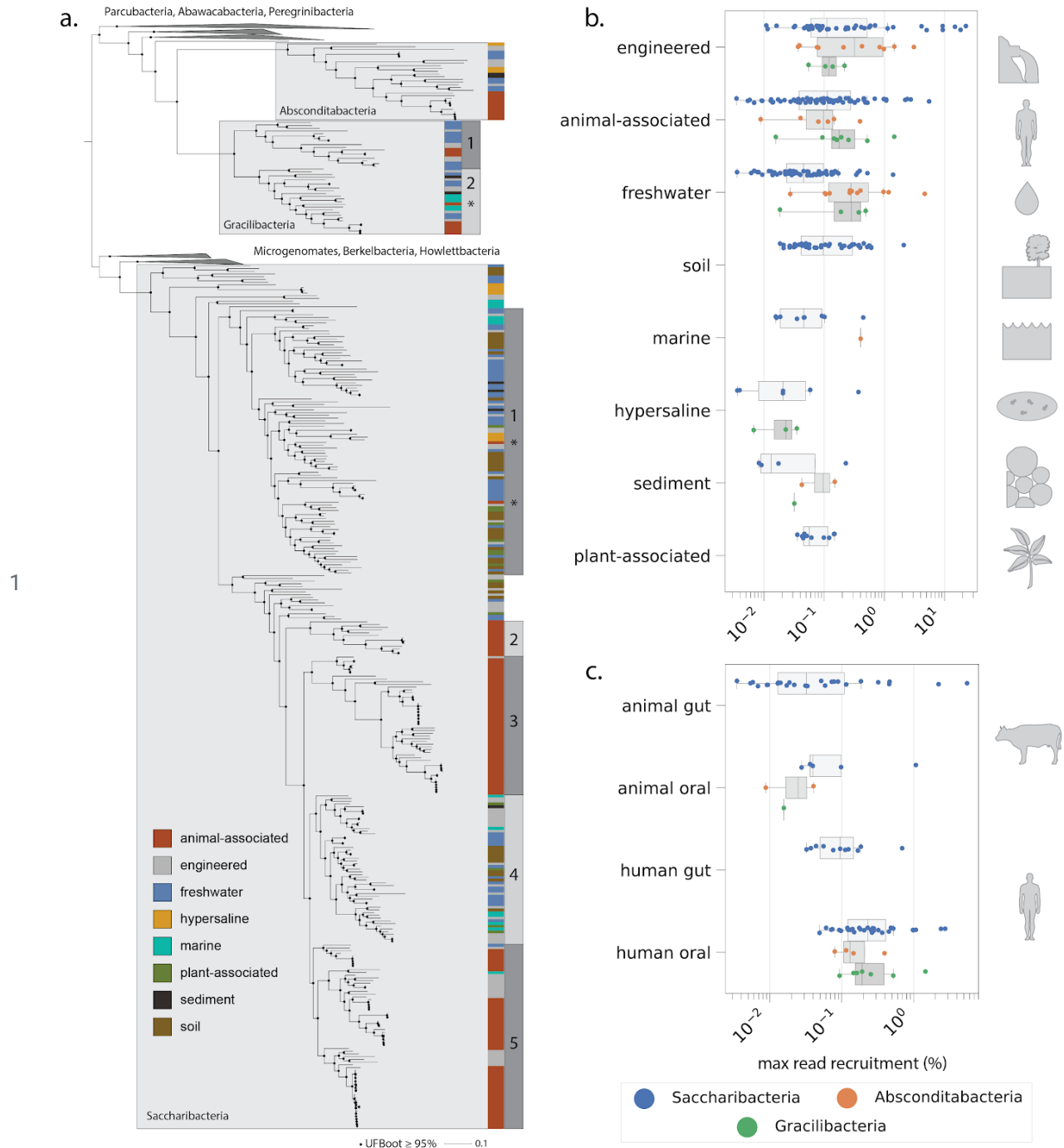


Fig. 1. Phylogenetic and environmental patterns for the Absconditabacteria, Gracilibacteria, and Saccharibacteria. **a)** Maximum-likelihood tree based on 16 concatenated ribosomal proteins (1976 amino acids, LG+R10 model). Scale bar represents the average number of substitutions per site. Asterisks indicate phylogenetic position of a subset of organisms derived from dolphin mouth metagenomes. Percentage of reads per metagenomic sample mapping to individual genomes across **b)** environments and **c)** body sites of humans and animals.

1 We used read mapping to assess the abundance of Absconditabacteria, Gracilibacteria, and
 2 Saccharibacteria genomes in the samples from which they were originally reconstructed.
 3 Generally, these CPR bacteria are not dominant members of microbial communities (<1% of
 4 reads). However, they were relatively abundant in some engineered, animal-associated, and
 5 freshwater environments (Fig. 1b). In rare cases, CPR taxa comprised >10% of reads (Fig. 1b),
 6 and in a bioreactor (engineered) reached a maximum of ~22% of reads. Gracilibacteria and
 7 Absconditabacteria attained comparable read recruitment to Saccharibacteria and were
 8 particularly abundant in some groundwater and animal-associated habitats. In contrast to
 9 Saccharibacteria, Gracilibacteria and Absconditabacteria have so far only been minimally
 10 detected in soil and plant-associated microbiomes. We also compared abundance patterns
 11 across animal body sites. As expected based on extensive prior work (5, 16, 44),
 12 Saccharibacteria exhibited highest read recruitment in the human oral microbiome. However,
 13 these bacteria can also comprise a significant fraction of the microbial community in exceptional
 14 gut/oral microbiomes from cows, pigs, and dolphins (Fig. 1c), in one case approaching 5% of
 15 sequenced reads (Table S2). When detected, Saccharibacteria in the human gut were relatively
 16 rare, comprising a median of ~0.1% of reads across samples.

17
 18

19 *Patterns of co-occurrence constrain CPR host range across environments*

20 Despite recent progress made in experimentally identifying bacterial host range for oral
 21 Saccharibacteria, little is known about associations in other habitats. Abundance pattern
 22 correlations can be informative regarding associations involving obligate symbionts and their
 23 microbial hosts (37, 55); however, such analyses often rely on highly resolved time-series for
 24 statistical confidence. Here, we instead examine patterns of co-occurrence within samples to
 25 probe potential relationships between CPR bacteria and their microbial hosts. Given recent
 26 experimental evidence demonstrating the association of multiple Saccharibacteria strains with
 27 various Actinobacteria in the human oral microbiome (15–17, 44, 83), we predicted that
 28 Actinobacteria may be common hosts of Saccharibacteria in microbiomes other than the mouth
 29 and asked to what extent co-occurrence data supported this relationship.

30 We first identified all ribosomal protein S3 (rps3) sequences from Actinobacteria and
 31 Saccharibacteria in the source metagenomes probed in this study for relative abundance
 32 patterns (Fig. 1bc). Rps3 sequences from all samples were clustered into ‘species groups’

1 (Materials and Methods). We observed that species groups from Actinobacteria and
 2 Saccharibacteria frequently co-occurred in the soil and plant-associated microbiomes as well as
 3 several hypersaline microbiomes (Fig. 2a). On the other hand, co-occurrence of the two
 4 lineages was less frequent in engineered and freshwater environments relative to other
 5 environments. Surprisingly, only ~78% of animal-associated samples containing
 6 Saccharibacteria also contained Actinobacteria at abundances high enough to be detected (Fig.
 7 2a). Assemblies with well-sampled Saccharibacteria yet no detectable Actinobacteria could
 8 suggest that Saccharibacteria have alternative hosts in these samples or are able to (at least
 9 periodically) live independently.

10 For samples where both Saccharibacteria and Actinobacteria marker genes were detectable,
 11 we computed a ‘relative richness’ metric describing the ratio of distinct Saccharibacteria species
 12 groups to Actinobacteria species groups. In most animal-associated microbiomes,
 13 Actinobacteria were more species rich (lower richness ratios), as expected if individual
 14 Saccharibacteria can associate with multiple hosts (Fig. 2a). Greater species richness of
 15 Actinobacteria compared to Saccharibacteria was also observed for many plant-associated, soil,
 16 engineered, and freshwater microbiomes. However, some engineered freshwater samples had
 17 richness ratios equal to (equal richness) or greater than 1 (i.e., Saccharibacteria more species
 18 rich) (Fig. 2a). Specifically, we observed that several metagenomes from engineered and
 19 freshwater environments contained anywhere from 1-11 Saccharibacteria species but only one
 20 detectable Actinobacteria species (Table S4). Thus, if Actinobacteria serve as hosts for
 21 Saccharibacteria in these habitats, there may be both exclusive associations and associations
 22 linking multiple Saccharibacteria species with a single Actinobacteria host species.

23 We next tested for more specific possible associations in the animal microbiome, reasoning that
 24 if Actinobacteria are common hosts for Saccharibacteria, then exclusive co-occurrence of a
 25 particular Saccharibacteria species with singular Actinobacteria species within a sample might
 26 suggest an interaction *in vivo*. We mapped all pairs of Saccharibacteria and Actinobacteria
 27 species that co-occurred within a single sample onto the trees constructed from the rpS3
 28 sequences (Fig. 2b), including 22 Saccharibacteria-Actinobacteria pairs reported in previous
 29 experimental studies (Table S3). In three additional cases, we found that individual
 30 metagenomic samples contained only one assembled Saccharibacteria species group and one
 31 Actinobacteria species group (“exclusive co-occurrence - species group”, Fig. 2b). Two of these
 32 cases involved Actinobacteria from the order Actinomycetales, from which multiple

1 Saccharibacteria hosts have already been identified. We also noted exclusive species-level
 2 co-occurrence of a Saccharibacteria species group from the human gut and an Actinobacteria
 3 species group from the order Coriobacteriales (Table S4). In an additional seven cases, one
 4 Saccharibacteria species group occurred with multiple Actinobacteria species groups of the
 5 same order-level classification based on rps3 gene profiling (“exclusive co-occurrence - order”,
 6 Fig. 2b). Five of the seven instances involved pairs of Saccharibacteria and Coriobacteriales
 7 from termite and swine gut metagenomes. Thus, unlike in human oral environments,
 8 Coriobacteriales may serve as hosts for Saccharibacteria in gut environments of multiple animal
 9 species. More generally, we also observed that Saccharibacteria from the same phylogenetic
 10 clade had predicted relationships to phylogenetically unrelated Actinobacteria (Fig. 2b),
 11 consistent with previous experimental observations for individual species (44).

12 Compared to Saccharibacteria, host relationships for Gracilibacteria and Absconditabacteria
 13 have received little attention. There are preliminary indications that Absconditabacteria may
 14 associate with members of the Fusobacteria or Firmicutes in the oral microbiome (44) or the
 15 gammaproteobacterium *Halochromatium* in certain salt lakes (84). We thus explored
 16 co-occurrence patterns in microbial communities containing Absconditabacteria and
 17 Gracilibacteria, attempting to further constrain possible host taxa. In animal and
 18 human-associated microbiomes, bacteria from several lineages, including Fusobacteria (Fig.
 19 2c), were relatively abundant in nearly all samples that contained Absconditabacteria. Members
 20 of the Chitinophagales, Pseudomonadales, and Acidimicrobiales were detected in high
 21 abundance in three wastewater samples from similar treatment plants (40) and one dairy pond
 22 sample containing Absconditabacteria. No clear patterns of potential host co-occurrence were
 23 observed for Gracilibacteria, with the exception of the proteobacterial order Campylobacteriales,
 24 which co-occurred in 8 of 10 groundwater samples where Gracilibacteria were found (Fig. 2c).
 25 Across all environments, only members of the order Burkholderiales (a large order of
 26 Gammaproteobacteria) consistently co-occurred with Gracilibacteria.

27

28

29

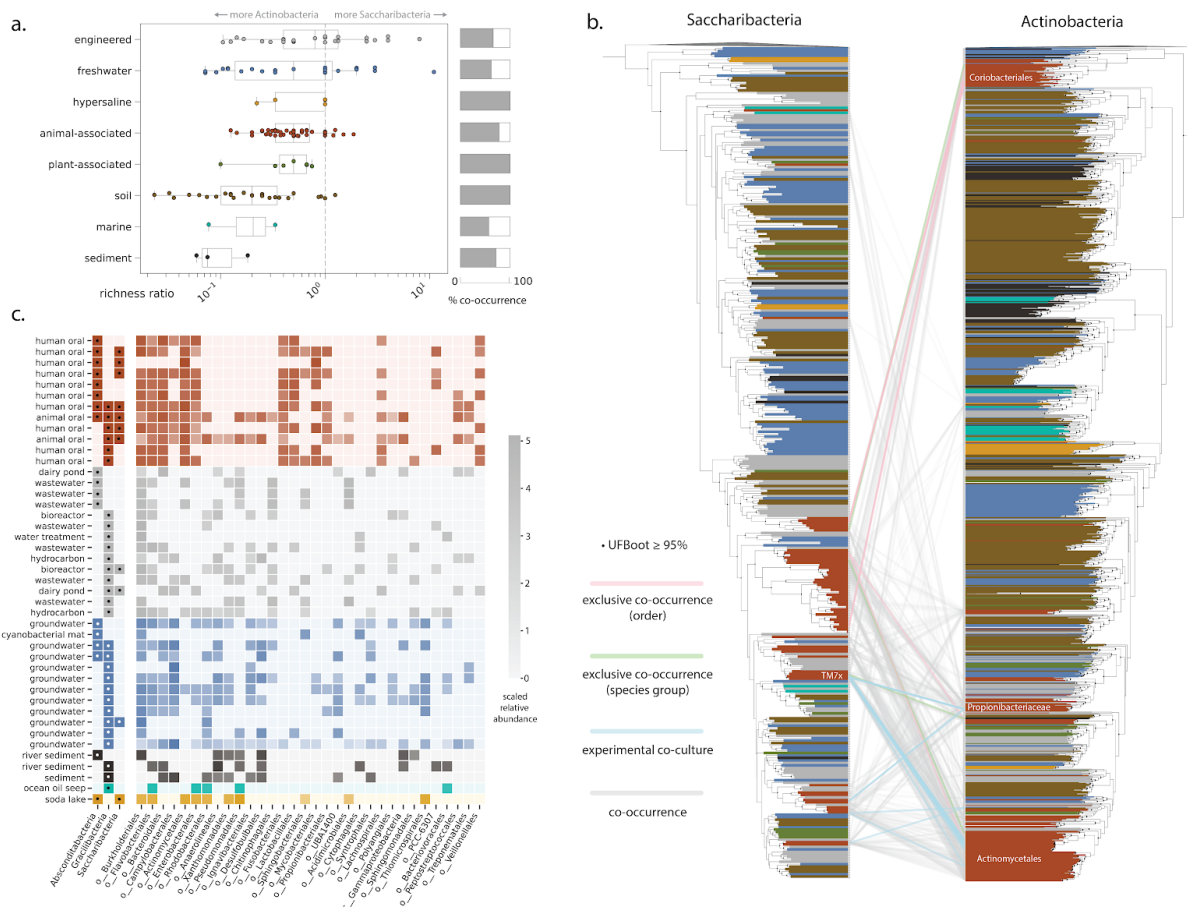


Fig. 2. Patterns of co-occurrence between CPR and potential host lineages across environments. **a)** Relative richness ratio, describing the ratio of distinct Saccharibacteria species groups to Actinobacteria species groups, for each sample and overall co-occurrence percentage across habitat categories. **b)** Maximum-likelihood trees for Saccharibacteria and Actinobacteria based on ribosomal protein S3 sequences extracted from all source metagenomes. Co-occurrence patterns are shown only for species groups derived from animal-associated metagenomes. **c)** Community composition for metagenomic samples containing Absconditabacteria and Gracilibacteria. Cells with dots indicate only presence, whereas those without dots convey quantitative relative abundance information. Only potential host lineages present in 8 or more samples are shown.

Among the least complex communities that contained Absconditabacteria were cyanobacterial mats from a California river network, where dominant cyanobacterial taxa accounted for ~60-98% relative abundance (32). To complement the above co-occurrence analyses, we re-analyzed 22 published metagenomes representing spatially separated mats and discovered that Absconditabacteria were detectable in 12 of them at varying degrees of coverage (0.12-37x). As noted previously, also present in the mats were members of the phyla Bacteroidetes, Betaproteobacteria, and Verrucomicrobia (32). Correlation of read coverage profiles across mats provided moderate support for the association of Absconditabacteria and

1 Bacteroidetes. Specifically, many of the strongest species-level correlations, including five of the
2 top ten, involved Bacteroidetes (Table S5).

3

4

5 *Gene content of Absconditabacteria, Gracilibacteria, and Saccharibacteria*

6

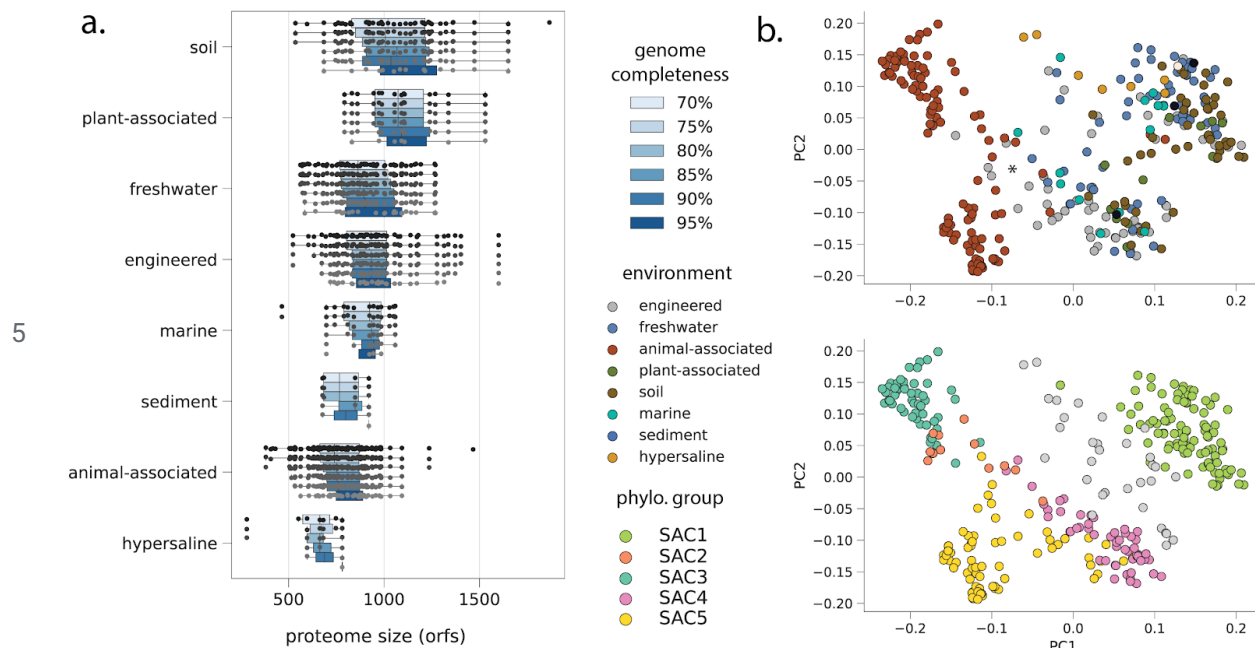
7 We next examined how gene content of these CPR lineages varied across environments. We
8 first compared the predicted proteome size of these bacteria across habitats, taking into
9 consideration differing degrees of genome completeness. This analysis revealed that genomes
10 from soil and the rhizosphere (plant-associated) have on average larger predicted proteomes
11 relative to those from animal-associated environments (Fig. 3b). Saccharibacteria from
12 hypersaline environments appear to have the smallest predicted proteomes, although the
13 limited number of high-quality genomes in this category currently limits a firm conclusion. We
14 observed some evidence for variance in predicted proteome sizes among Absconditabacteria
15 and Gracilibacteria, including potentially smaller predicted proteomes among animal-associated
16 Gracilibacteria (Fig. S1). Additional high quality genomes will be required to confirm this trend.

17

18 To examine overall proteome similarity as a function of habitat type, we employed a recently
19 developed protein-clustering approach that is agnostic to functional annotation (9) (Table S6,
20 Materials and Methods). Among Saccharibacteria, principal coordinates analysis (PCoA) of
21 presence/absence profiles for all protein families with 5 or more members yielded a primary axis
22 of variation (~12% variance explained) that distinguished animal-associated Saccharibacteria
23 from environmental or plant-associated ones and a secondary axis (~8 % variance explained)
24 that distinguished between phylogenetic clades (SAC1-3 vs 4-5). We did not observe strong
25 clustering of Saccharibacteria by specific environmental biome, consistent with the interspersed
26 nature of their phylogenetic relationships (Fig. 1a, 3b). Notably, several SAC5 genomes from
27 wastewater have protein family contents that are intermediate between those of
28 animal-associated Saccharibacteria and Saccharibacteria from the large environmental clade
29 (indicated by an asterisk in Fig. 3b). This finding may indicate selection within the engineered
30 environments for variants introduced from human waste. PCoAs of predicted proteome content
31 among Absconditabacteria and Gracilibacteria generally showed that, with the exception of
32 dolphin-derived genomes, animal-associated lineages are also distinct from their relatives from
33 environmental biomes (Fig. S2). Overall, our results indicate that the CPR lineages examined

1 here have predicted proteomes whose content and size vary substantially with their
2 environment. This is particularly evident for animal-associated Saccharibacteria, which are
3 notably dissimilar in their protein family content compared to environmental counterparts.

4



6 **Fig. 3.** Proteome characteristics for Saccharibacteria. **a)** Predicted proteome size (open reading frame
7 count) at increasing genome completeness thresholds. **b)** Overall proteome similarity among
8 Saccharibacteria from different environmental categories (top panel) and phylogenetic clades (bottom
9 panel). PCoAs were computed from presence/absence profiles of all protein clusters with 5 or more
10 member sequences. The primary (PC1) and secondary (PC2) axis of variation explained 12% and 8% of
11 variance, respectively.

12

13 To further examine the distinctions evident in the PCoA analysis, we arrayed presence/absence
14 information for each protein family and hierarchically clustered them based on their distribution
15 patterns across all three CPR phyla. This strategy allowed us to explore specific protein family
16 distributions and to test for groups of co-occurring protein families (modules) that are common
17 to bacteria from a single lineage or are shared by most bacteria from one or more CPR
18 lineages. We first observed one large module that is generally conserved across all genomes.
19 This module is comprised of families for essential cellular functions such as transcription,
20 translation, cell division, and basic energy generating mechanisms (Fig. 4a, “core”).

21

22 The protein family analysis also revealed multiple modules specific to Gracilibacteria,
23 Absconditabacteria, and modules shared by both lineages but not present in Saccharibacteria,

1 paralleling their phylogenetic relationships (Fig. 1a, 4a). Of the ~70 families shared only by
 2 Gracilibacteria and Absconditabacteria (M2, Fig. 4a), nearly half had no KEGG annotation at the
 3 thresholds employed. One family shared by these phyla but not in Saccharibacteria is the
 4 ribosomal protein L9, which supports prior findings on the composition of Saccharibacteria
 5 ribosomes (29). The remaining families also include two that were fairly confidently annotated
 6 as the DNA mismatch repair proteins, MutS and MutL (fam01378 and fam00753), nicking
 7 endonucleases involved in correction of errors made during replication (85) (Table S6). Despite
 8 the generally wide conservation of these proteins among Bacteria, we saw no evidence for the
 9 presence of either enzyme in Saccharibacteria, suggesting that aspects of DNA repair may vary
 10 in this group relative to other CPR. We recovered a module of approximately 60 proteins highly
 11 conserved among the Saccharibacteria and only rarely encoded in the other lineages (M5, Fig.
 12 4a). This module contained several protein families confidently annotated as core components
 13 of glycolysis and the pentose phosphate pathway, including three enzymes present in almost all
 14 CPR (10): glyceraldehyde 3-phosphate dehydrogenase, (GAPDH) triosephosphate isomerase
 15 (TIM), and phosphoglycerate kinase (PGK). These results indicate that Gracilibacteria and
 16 Absconditabacteria may have extremely patchy, if not entirely lacking, components of core
 17 carbon metabolism, even when a high-quality genome set is considered.

18
 19 For all three lineages of CPR, we also observed numerous small modules with narrow
 20 distributions. To test whether these modules represent functions differentially distributed among
 21 organisms from different habitats, we computed ratios describing the incidence of each protein
 22 family in one habitat compared to all others (Materials and Methods). Enriched families were
 23 defined as those with ratios ≥ 5 , whereas depleted families were defined as those that were
 24 encoded by $<10\%$ of genomes in a given habitat, but $\geq 50\%$ of genomes from other habitats. To
 25 account for the fact that small families might appear to be differentially distributed due to chance
 26 alone, we also stipulated that comparisons be statistically significant ($p \leq 0.05$, two-sided
 27 Fisher's exact test corrected for multiple comparisons).

28
 29
 30
 31
 32

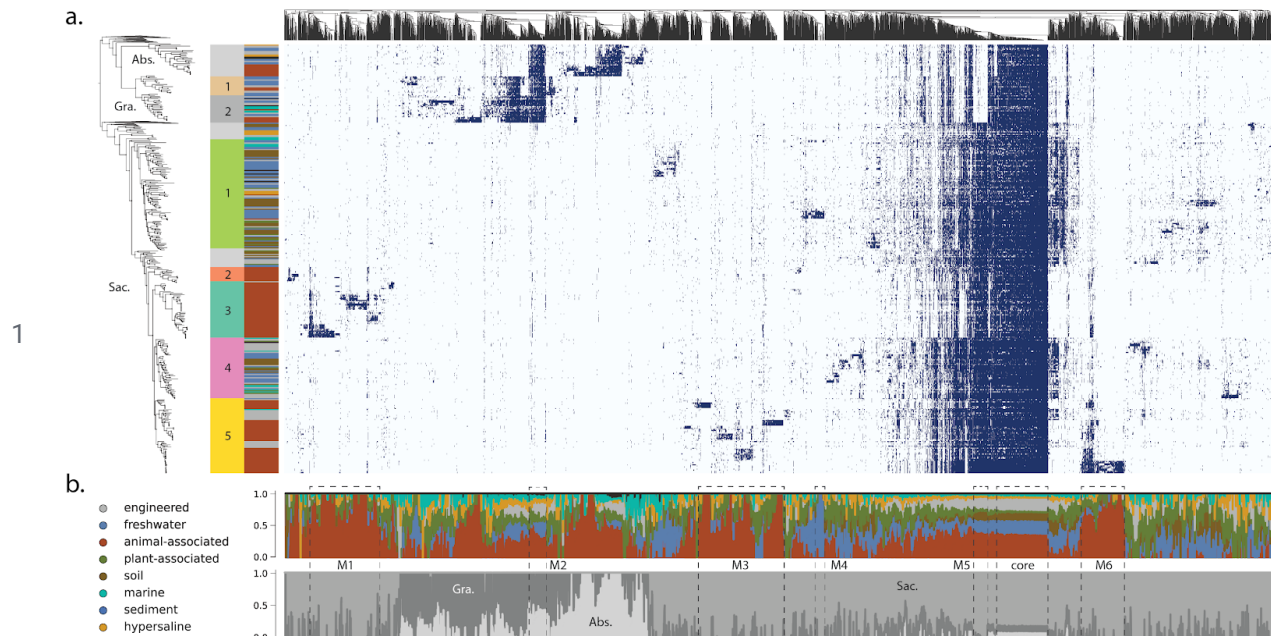


Fig. 4. Phylogenetic and environmental distribution of protein families recovered among CPR. **a)** Presence/absence profiles for protein families with 5 or more members, with shaded cells indicating presence, and light cells indicating absence. Heatmap columns represent protein families, hierarchically clustered by similarity in distribution across the genome set. Rows correspond to genomes, ordered by their phylogenetic position in the species tree (left). Abbreviations: Abs., Absconditabacteria; Gra., Gracilibacteria; Sac., Saccharibacteria. **b)** Percentage of genomes encoding individual protein families that belonged to broad habitat groups (top panel) or taxonomic groups (bottom panel). Modules of protein families indicated in the text are represented by dotted lines (M1-6 and 'core').

Using this approach, we identified 926 families that were either enriched ($n=872$) or depleted ($n=54$) in genomes from one or more broad habitat groups. We identified 45 families enriched in Absconditabacteria from animal-associated environments relative to those from environmental biomes. The majority of these families were either poorly functionally characterized or entirely without a functional annotation at the thresholds employed. Similarly, families enriched in animal-associated Gracilibacteria relative to environmental counterparts were primarily unannotated; among those families with confident annotations was a family likely encoding a phosphate: Na^+ symporter (fam04488) and a putative membrane protein (fam06579). Intriguingly, 6 families were co-enriched in both animal-associated Gracilibacteria and Absconditabacteria, suggesting that these sibling lineages might have acquired or retained a small complement of genes that are important in adaptation to animal habitats or their associated bacteria.

1 Animal-associated Saccharibacteria, on the other hand, encoded 417 unique families that were
 2 exclusive or highly enriched relative to those from other habitats. Enriched families largely fell
 3 into three major groups (M1, M3, M6; Fig. 4), and the large majority of them, particularly among
 4 modules with narrow, lineage-specific distributions, were without functional annotations.
 5 However, our analysis also revealed some protein families with broader distributions across
 6 multiple clades of animal-associated Saccharibacteria (Fig. 4). Here, among families with
 7 functional annotations, we found several apparently involved in the transport of amino acids and
 8 dicarboxylates that were highly enriched (ratios ranging from 10 to 112.9) in the majority of
 9 animal-associated Saccharibacteria (52-58% of genomes across clades) (Table S7). Two of
 10 these families, corresponding to a putative amino acid transport permease and
 11 substrate-binding protein (fam00393 and fam11477, respectively) were co-located in some
 12 genomes along with a ATP-binding protein (subset of fam00001), suggesting that they may
 13 function together to uptake amino acids. We also recovered several other functions that were
 14 previously predicted to be enriched based on analysis of a smaller set of animal-associated
 15 Saccharibacteria (4), including phosphoglycerate mutase, glycogen phosphorylase, and a
 16 uracil-DNA glycosylase (ratio 8.3-33.5). Lastly, we found that the CRISPR-associated protein
 17 cas9 was moderately enriched among animal-associated genomes (ratio ~5 among 33
 18 genomes), consistent with the suggestion that these Saccharibacteria likely acquired their viral
 19 defense systems after colonizing animals (Table S7) (4).

20

21 We identified multiple families that are either enriched or depleted in animal-associated
 22 Saccharibacteria that were functionally related to oxidative stress (Table S7). Among enriched
 23 families, one (fam00662) set was mostly annotated with low confidence as rubrerythrin, a family
 24 of iron containing proteins generally involved in detoxification of peroxide (86). Member
 25 sequences of this family were present in over a third of animal-associated Saccharibacteria and
 26 were highly enriched relative to environmental genomes (fold-enrichment ratio of 36.2),
 27 suggesting that acquisition may have conferred an adaptive benefit in the gut and/or oral cavity.
 28 In contrast, we also observed that animal-associated Saccharibacteria were significantly
 29 depleted in another family confidently annotated as a Fe-Mn family superoxide dismutase
 30 (fam01569) and likely involved in radical detoxification. Animal-associated lineages were also
 31 strongly depleted for the genes comprising the cytochrome o ubiquinol oxidase operon
 32 (fam00281, fam00112, fam01347, fam00624, and fam10494), with very few, if any,
 33 animal-associated genomes and more than 50% of environmental genomes encoding each of

1 the five genes. This operon has been previously suggested to confer an advantage in aerophilic
2 environments like soil through detoxification (3) or use of oxygen (12, 13).
3
4 Among genomes from environmental biomes, we identified a module of approximately 100
5 protein families, also primarily without functional annotation, that were associated with a
6 subclade of Saccharibacteria recently reconstructed from metagenomes of freshwater lakes and
7 glacier ice (M4, Fig. 4) (61, 87). Notably, among the most widespread families in this module
8 was one in which sequences were annotated as bacteriorhodopsin with low confidence
9 (fam11249). Further analysis indicated that these sequences fall within the bacterial/archaeal
10 Type 1 rhodopsin clade and contain both the retinal-binding lysine associated with light
11 sensitivity and a DTS motif (Fig. S3), suggesting that they may function as proton pumps (88,
12 89). Distinct rhodopsin sequences were also recovered in the genomes of environmental
13 Absconditabacteria (NDQ motif) and Gracilibacteria (DTE motif), although they were not
14 statistically enriched (Fig. S3). Genomes of soil-associated Saccharibacteria were enriched for
15 about 130 protein families largely without strong functional annotations (Fig. 3). Despite their
16 small proteome sizes, Saccharibacteria from hypersaline environments were only statistically
17 depleted in about 15 families at the thresholds employed here. Sequence files for all protein
18 families are provided in the Supplementary Materials (File S2).
19

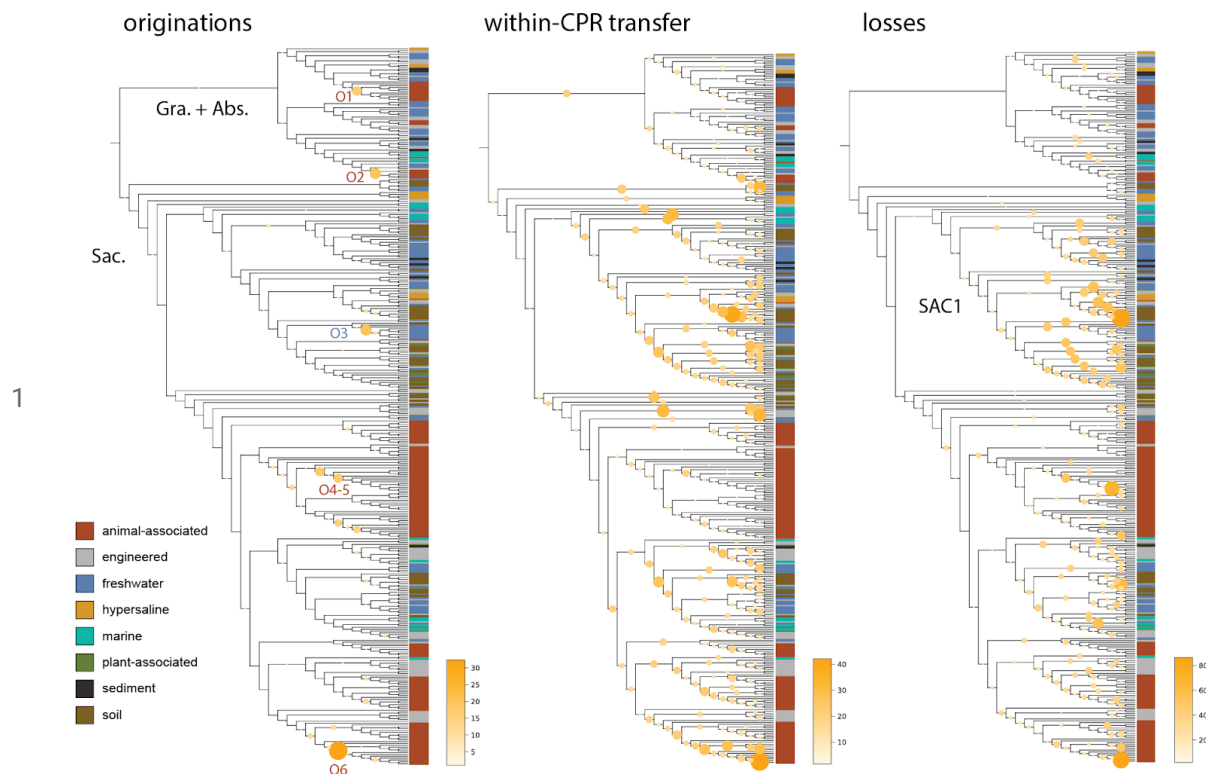


Fig. 5. Evolutionary processes shaping proteome evolution in three lineages of CPR bacteria. Each panel displays the species tree from Fig. 1 in cladogram format. The size and color of circles mapped onto interior branches represent the cumulative number of **a)** originations (defined as either lateral transfer from outside the lineages examined here, or *de novo* evolution) **b)** transfer among the three CPR lineages included here and **c)** genomic losses predicted to occur on that branch for all 902 differentially distributed families where gene-species tree reconciliation was possible. Abbreviations: Abs., Absconditabacteria; Gra., Gracilibacteria; Sac., Saccharibacteria. SAC1 indicates a monophyletic clade of Saccharibacteria referenced in the text.

Evolutionary processes shaping proteome evolution

The observation that some differentially distributed traits among CPR were apparently lineage specific, whereas others were more widespread, motivated us to examine the relative contributions of gene transfer and loss to proteome evolution. To do so, we first inferred unrooted, maximum-likelihood phylogenies for the sequences in each protein family that was differentially distributed, then compared these phylogenies to the previously reconstructed species tree (Materials and Methods). For each family, the likelihood of transfer and loss events on each branch of the species tree were then estimated using a probabilistic framework that takes into consideration genome incompleteness, variable rates of transfer and loss, and uncertainty in gene tree reconstruction (90, 91). The results of this analysis reveal relatively few

1 instances of originations, defined as lateral transfer from outside the three lineages of CPR or
2 *de novo* evolution ('originations', Fig. 5). In the Absconditabacteria and Gracilibacteria,
3 gene-species tree reconciliation revealed that small modules of families of mostly hypothetical
4 proteins were acquired near the base of animal-associated clades (O1-O2, Fig. 5). On the other
5 hand, in Saccharibacteria, originations were primarily associated with shallower subclades of
6 animal-associated (and in one case, freshwater) genomes (O3-O6, Fig. 5). These findings
7 generally corresponded with the distribution of small, highly enriched modules of largely
8 hypothetical proteins (Fig. 4) and suggest that the distribution of these modules is best
9 explained by lineage-specific acquisition events of relatively few genes at one time, rather than
10 large acquisition events at deeper nodes. Intriguingly, one subclade of animal-associated
11 Saccharibacteria had the highest incidence of originations of all groups in our analysis (O6),
12 suggesting that these genomes may be phylogenetic 'hotspots' for transfer.

13
14 While origination events were relatively infrequent in all three CPR lineages, instances of
15 within-CPR transfer and loss were very frequent and dispersed across most interior branches of
16 the tree (Fig. 5). Notably, we detected sporadic losses across internal branches, which is
17 inconsistent with a major gene loss event at the time of adaptation to animal-associated
18 habitats. Surprisingly, we noticed that genomes of non-animal associated Saccharibacteria,
19 particularly those from the SAC1 clade, displayed substantial patterns of loss despite their
20 relatively large proteome sizes. Thus, losses in these environmental lineages were possibly
21 balanced by lateral transfer events over the course of evolution.

22

23

24 **DISCUSSION:**

25

26 Here, we expand sampling of genomes from the Absconditabacteria, Gracilibacteria, and
27 Saccharibacteria, particularly from environmental biomes. The basal positioning of
28 environmental clades in phylogenetic reconstructions provides strong support for the hypothesis
29 that these lineages originated in the environment (Fig. 1a), and potentially migrated into humans
30 and terrestrial animals via consumption of groundwater (4, 41). Unlike the Absconditabacteria,
31 which appear to have transitioned only once into animal oral cavities and guts, our phylogenetic
32 evidence suggests that Gracilibacteria and Saccharibacteria have undergone multiple
33 transitions into the animal microbiome in unique phylogenetic clades. This includes

1 Saccharibacteria and Gracilibacteria from the dolphin mouth that appear to have colonized
2 these marine mammals from a distinct source compared to those that colonized the oral
3 environments of terrestrial animals (Fig. 1). In other clades, phylogenetically interspersed
4 environmental and oral/gut Saccharibacteria could reflect either independent migrations into the
5 animal environment or lineage-specific reversion back to environmental niches (SAC2-5, Fig. 1).
6
7 Currently, the mechanisms that enable environmental transition among CPR are unknown.
8 Several observations, including that CPR host-pairs may be taxonomically distinct between oral
9 and gut habitats, raise the question of whether habitat transitions among CPR involve
10 co-migration with their hosts or the acquisition of new hosts. The finding that single CPR
11 species co-occur with a single Actinobacteria species, or several closely related ones, in
12 multiple animal-associated metagenomes contributes further evidence that these associations
13 can be flexible and phylogenetically diverse rather than highly evolutionarily conserved (44).
14 Supporting this, some laboratory strains of oral Saccharibacteria can adapt to new hosts after
15 periods of living independently (92). The lack of evidence for lateral gene transfer between
16 experimentally profiled pairs (4) also suggests that some CPR-host pairs may have established
17 fairly recently.
18
19 Host associations for Absconditabacteria, Gracilibacteria, and environmental Saccharibacteria
20 are currently unknown. Previously, changes in abundance over a sample series from a
21 bioreactor system treating thiocyanate was used to suggest that *Microbacterium ginsengisoli*
22 may serve as a host for a co-occurring Saccharibacteria (55). One Absconditabacteria lineage
23 (*Vampirococcus*) has been predicted to have a host from the Gammaproteobacteria (84) and
24 one Gracilibacteria was suggested to have a *Colwellia* host based on a shared repeat protein
25 motif (11). Given the scant information on hand about possible hosts for these CPR, especially
26 for Absconditabacteria and Gracilibacteria, the patterns of co-occurrence we report for specific
27 organisms provide starting points for host identification via targeted co-isolation.
28
29 To evaluate to what extent changes in gene content are associated with habitat transition, we
30 first established core gene sets indicated that overall proteome size and content differed
31 between environmental and animal-associated Saccharibacteria, and to some extent
32 Gracilibacteria. Despite overall smaller proteome size, we identified a large number of protein
33 families that were highly enriched among animal-associated CPR from all three lineages. The

1 most striking capacities involve amino acid transport, oxidative stress tolerance, and viral
2 defense, which may enable use of habitat-specific resources or tolerance of habitat-specific
3 stressors. These findings complement previous suggestions that prophages are enriched in
4 animal-associated Saccharibacteria relative to environmental counterparts (14).

5
6 Only three lineages of CPR (of potentially dozens) have been consistently recovered in the
7 animal-associated microbiome. Given the enormous diversity of CPR bacteria in drinking water
8 (41), there has likely been ample opportunity for various taxa to disperse into the mouths of
9 terrestrial animals; however, establishment and persistence of these bacteria may have been
10 limited by the absence of a suitable host in oral and gut environments. Thus, we predict that
11 other CPR bacteria - including those from the large Microgenomates and Parcubacteria
12 lineages - have hosts that are infrequent or transient members of the animal microbiome, or
13 have insufficient ability to 'adapt' to new hosts upon contact. For example, formation of new
14 associations may be limited by the specificity of pili involved in host interaction or proteins
15 involved in attachment (14, 41, 93).

16
17 It is also interesting to compare processes of habitat transition in CPR with those proposed for
18 other bacteria and for archaea. Our results suggest that Saccharibacteria (and potentially
19 Gracilibacteria) from the human/animal microbiome have smaller genome sizes than related,
20 deeper-branching lineages of environmental origin. This pattern is also apparent for other,
21 free-living groups adapted to the animal microbiome from the environment, like the
22 Elusimicrobia (94) and intracellular symbionts of insects (95). Furthermore, in contrast to
23 findings for Elusimicrobia, where host-associated lineages have common patterns of loss of
24 metabolic capacities compared to relatives from non-host environments (94), patterns of gene
25 loss in animal-associated CPR appear to be heterogeneous and lineage-specific. One
26 possibility is that gene loss in CPR is primarily modulated by strong dependence on host
27 bacteria, whose capacities may vary substantially, rather than by adaptation to the relatively
28 stable, nutrient-rich animal habitat that likely shaped evolution of some non-CPR bacteria.

29
30 Changes in gene content could enable, facilitate, or follow habitat transitions. Our evolutionary
31 reconstructions revealed that habitat-specific differences in gene content are more likely the
32 product of combinations of intra-CPR transfer and loss rather than major acquisition events at
33 time of lineage divergence. Thus, modules enriched in specific lineages were probably acquired

1 via lateral transfer after habitat transition, suggesting that proteome remodeling has been
2 continuous in CPR over evolutionary time. As such, the processes shaping CPR lineage
3 evolution share both similarities and differences with those predicted for other microbes,
4 including Haloarchaeota (96) and ammonia-oxidizing lineages of Thaumarchaeota (90, 97),
5 where both large lateral transfer events and gradual patterns of gene loss, gain, and duplication
6 worked together to shape major habitat transitions.

7

8 **CONCLUSION:**

9

10 Overall, our findings highlight factors associated with habitat transitions in three CPR lineages
11 that occur in both human/animal and environmental microbiomes. We expand the evidence for
12 niche-based differences in protein content (4, 14) and identify a large set of protein families that
13 could guide future studies of CPR symbiosis. Furthermore, patterns of co-occurrence may
14 inform experiments aiming to co-cultivate CPR and their hosts. Our analyses point to a history
15 of continuous genome remodeling accompanying transition into human/animal habitats, rather
16 than rapid gene gain/loss around the time of habitat switches. Thus, habitat transitions in CPR
17 may have been primarily driven by the availability of suitable hosts and reinforced by acquisition
18 and/or loss of genetic capacities. These processes may be distinct from those shaping
19 transitions in other bacteria and archaea that are not obligate symbionts of other
20 microorganisms.

21

22 **MATERIALS AND METHODS:**

23

24 *Genome database preparation and curation*

25

26 To compile an environmentally comprehensive set of genomes from the selected CPR lineages,
27 we first queried four genomic information databases - GTDB (<https://gtdb.ecogenomic.org/>),
28 NCBI assembly (<https://www.ncbi.nlm.nih.gov/assembly>), PATRIC (<https://www.patricbrc.org/>),
29 and IMG (<https://img.jgi.doe.gov/>) - for records corresponding to the Absconditabacteria,
30 Gracilibacteria, and Saccharibacteria genomes. Genomes gathered from these databases were
31 combined with those drawn from several recent publications as well as genomes newly binned
32 from metagenomic samples of sulfidic springs, an advanced treatment system for potable reuse
33 of wastewater, human saliva, cyanobacterial mats, fecal material from primates, baboons, pigs,
34 goats, cattle, and rhinoceros, several deep subsurface aquifers, dairy-impacted groundwater
35 and associated enrichments, multiple bioreactors, soil, and sediment (Table S1). Assembly,
36 annotation, and binning procedures followed those from Anantharaman et al. 2016. In some
37 cases, manual binning of the alternatively coded Absconditabacteria was aided by a strategy in

1 which a known Absconditabacteria gene was blasted against predicted metagenome scaffolds
2 to find ‘seed’ scaffolds, whose coverage and GC profile were used to probe remaining scaffolds
3 for those with similar characteristics. For newly binned genomes, genes were predicted for
4 scaffolds > 1 kb using prodigal (meta mode) and annotated using usearch against the KEGG,
5 UniProt, and UniRef100 databases. Bins were ‘polished’ by removing potentially contaminating
6 scaffolds with phylogenetic profiles that deviated from consensus taxonomy at the phylum level.
7 One genome was further manually curated to remove scaffolding errors identified by read
8 mapping, following the procedures outlined in (98).

9
10 We removed exact redundancy from the combined genome set by identifying identical genome
11 records and selecting one representative for downstream analyses. We then computed
12 contamination and completeness for the genome set using a set of 43 marker genes sensitive to
13 described lineage-specific losses in the CPR (29, 31) using the custom workflow in checkm
14 (99). Results were used to secondarily filter the genome set to those with $\geq 70\%$ of the 43
15 marker genes present and $\leq 10\%$ of marker genes duplicated. The resulting genomes were then
16 de-replicated at 99% ANI using dRep (-sa 0.99 -comp 70 -con 10) (100), yielding a set of 389
17 non-redundant genomes. Existing metadata were used to assign both “broad” and “narrow”
18 habitat of origin for each non-redundant genome. The “engineered” habitat category was
19 defined to include human-made or industrial systems like wastewater treatment, bioreactors,
20 and water impacted by farming/mining. Curated metadata, along with accession/source
21 information for each genome in the final set, is available in Table S1. All newly binned genomes
22 are available through Zenodo (Data and Software Availability).

23 24 *Functional annotation and phylogenomics*

25
26 We predicted genes for each genome using prodigal (“single” mode), adjusting the translation
27 table (-g 25) for Gracilibacteria and Absconditabacteria, which are known to utilize an alternative
28 genetic code (7, 8). Predicted proteins were concatenated and functionally annotated using
29 kofamscan (101). Results with an e-value $\leq 1e-6$ were retained and subsequently filtered to
30 yield the highest scoring hit for each individual protein.

31
32 To create a species tree for the CPR groups of interest, functional annotations from kofamscan
33 were queried for 16 syntenic ribosomal proteins (rp16). Marker genes were combined with those
34 from a set of representative sequences of major, phylogenetically proximal CPR lineages (10).
35 Sequences corresponding to each ribosomal protein were separately aligned with MAFFT and
36 subsequently trimmed for phylogenetically informative regions using BMGE (-m BLOSUM30)
37 (102). We then concatenated individual protein alignments, retaining only genomes for which at
38 least 8 of 16 syntenic ribosomal proteins were present. A maximum-likelihood tree was then
39 inferred for the concatenated rp16 (1976 amino acids) set using ultrafast bootstrap and
40 IQTREE’s extended Free-Rate model selection (-m MFP -st AA -bb 1000) (103). The maximum
41 likelihood tree is available as File S1. The tree and associated metadata were visualized in iTol
42 (104) where well-supported, monophyletic subclades were manually identified within
43 Gracilibacteria and Saccharibacteria for use in downstream analysis.

Abundance analysis

To assess the global abundance of Absconditabacteria, Gracilibacteria, and Saccharibacteria, we manually compiled the original read data associated with each genome in the analysis set, where available. We included only those genomes from short-read, shotgun metagenomics of microbial entire communities (genomes derived from single cell experiments, stable isotope probing experiments, “mini” metagenomes, long-read sequencing experiments, and co-cultures were excluded). For each sequencing experiment, we downloaded the corresponding raw reads and, where appropriate, filtered out animal-associated reads by mapping to the host genome using bbdut (*qhdist=1*). Sequencing experiments downloaded from the NCBI SRA database were sub-sampled to the average number of reads across all compiled experiments (~36 million) using seqtk (*sample -s 7*) if the starting read pair count exceeded 100 million. We then removed Illumina adapters and other contaminants from the remaining reads and further quality trimmed them using Sickle. The filtered read set was then mapped against all genomes assembled (or co-assembled) from it using bowtie2 (default parameters). For mappings with a non-zero number of read alignments, relative abundance of each genome was calculated by counting the number of stringently mapped ($\geq 99\%$ identity) using coverm (*--min-read-percent-identity 0.99*) and dividing by the total number of reads in the quality-filtered read set. Genomes were considered present in a sample if at least 10% of sequence length was covered by reads. In cases where genomes were derived from co-assemblies of multiple sequencing experiments, we computed the abundance for each sample individually and then selected the one with the highest value as a ‘representative’ sample for downstream analyses.

Co-occurrence analyses

Each representative sample was then probed for co-occurrence patterns of CPR and potential host lineages. To account for across-study differences in binning procedures, quality-filtered read sets were re-assembled using megahit (*--min-contig-len 1000*) and subsequently analyzed using graftm (105) with a ribosomal protein S3 (*rps3*) gpackage custom built from GTDB (release 05-RS95). Recovered *rps3* protein sequences in each sample were clustered to form ‘species groups’ at 99% identity using usearch cluster_fast (*-sort length -id 0.99*). For all samples with >0 marker hits, we then performed three downstream analyses to examine patterns of co-occurrence for various taxa. First, we counted the number of unique species groups in each sample taxonomically annotated as Saccharibacteria (“*c__Saccharimonadia*”) and Actinobacteria (“*p__Actinobacteriota*”), dividing the former by the latter to compute a species ‘richness ratio’ for each sample (where *p__Actinobacteria* did not equal 0).

Second, to examine the co-occurrence of Saccharibacteria and Actinobacteria within a phylogenetic framework, we inferred maximum likelihood trees for the set of *rps3* marker genes recovered across samples. Species group sequences were clustered across samples to further reduce redundancy using usearch (as described above) and were combined with *rps3* sequences drawn from a taxonomically balanced set of bacterial reference genomes (10) as an outgroup. Saccharibacteria and Actinobacteria sequence sets were then aligned, trimmed, and

1 used to build trees as described above for the 16 ribosomal protein tree, with the exception of
 2 using trimal (*-gt 0.1*) (106) instead of BMGE. Species groups that co-occurred in one or more
 3 metagenomic samples were then linked. If a given Saccharibacteria species group exclusively
 4 co-occurred with an Actinobacteria species group in *at least* one sample, or Actinobacteria
 5 species groups belonging to the same order level in *all* samples, those linkages were labelled.
 6 Finally, experimental co-cultures of Saccharibacteria and Actinobacteria from previous studies
 7 were mapped onto the trees. To do this, we compiled a list of strain pairs and their
 8 corresponding genome assemblies (Table S3) and then used graftm to extract rps3 sequences
 9 from corresponding genome assemblies downloaded from NCBI. We then matched these rps3
 10 sequences to their closest previously defined species group using blastp (*-evaluate 1e-3*
 11 *-max_target_seqs 10 -num_threads 16 -sorthits 3 -outfmt 6*), prioritizing hits with the highest
 12 bitscore and alignment length. Reference rps3 sequences with no match at $\geq 99\%$ identity and
 13 $\geq 95\%$ coverage among the species groups were inserted separately into the tree. We then
 14 labelled all experimental pairs of species in the linkage diagram.

15
 16 Third, we profiled a subset of 43 metagenomes containing Gracilibacteria and
 17 Absconditabacteria for overall community composition. For each sample, we extracted all
 18 contigs bearing ribosomal S3 proteins and mapped the corresponding quality filtered read set to
 19 them using bowtie2. Mean coverage for each contig was then computed using coverm (*contig*
 20 *--min-read-percent-identity .99*) and a minimum covered fraction of 0.10 was again employed.
 21 Relative coverage for each order level lineage (as predicted by graftm) was computed by
 22 summing the mean coverage values for all rps3-bearing contigs belonging to that lineage.
 23 Where species groups did not have order-level taxonomic predictions, the lowest available rank
 24 was used. Finally, relative coverage values were scaled by first dividing by the lowest relative
 25 coverage observed across samples and then taking the base-10 log. For the re-analysis of 22
 26 cyanobacterial mat metagenomes (32), the same approach was taken, and coverage profiles for
 27 rps3-bearing scaffolds were correlated using the pearsonr function in the scipy.stats package.

28 29 *Proteome size, content, enrichment*

30
 31 We subjected all predicted proteins from the genome set to a two-part, de novo protein
 32 clustering pipeline recently applied to CPR genomes, in which proteins are first clustered into
 33 “subfamilies” and highly similar/overlapping subfamilies are merged using and HMM-HMM
 34 comparison approach (*--coverage 0.70*) (9) (<https://github.com/raphael-upmc/proteinClustering-Pipeline>). For each protein cluster, we recorded the most common KEGG annotation among its
 35 member sequences and the percent of sequences bearing this annotation (e.g. 69% of
 36 sequences in fam00095 were matched with K00852).

37
 38
 39 We then performed three subsequent analyses to describe broad proteome features of included
 40 CPR. First, we computed proteome size across habitats, defined as the number of predicted
 41 ORFs per genome when considering genomes at increasing thresholds of completeness in
 42 single copy gene inventories (75%, 80%, 85%, etc.). Second, we examined similarity between
 43 proteomes by generating a matrix describing the presence/absence patterns of protein families

1 with 5 or more member sequences. We then used this matrix to compute distance metrics
 2 between each genome based on protein content using the ecopy package in Python
 3 (method='jaccard', transform='1') and performing a principal coordinates analysis (PCoA) using
 4 the skbio package. The first two axes of variation were retained for visualization alongside
 5 environmental and phylogenetic metadata. Finally, we used the clustermap function in seaborn
 6 (metric='jaccard', method='average') to hierarchically cluster the protein families based on their
 7 distribution patterns and plot these patterns across the genome set. For each protein family, we
 8 also computed the proportion of genomes encoding at least one member sequence that
 9 belonged to each of the three CPR lineages and each broad environmental category (Fig. 4b)
 10 (see custom code linked in Data and Software Availability).

11
 12 We next identified protein families that were differentially distributed among genomes from
 13 broad environmental categories. For each protein family, we divided the fraction of genomes
 14 from a given habitat ('in-group') encoding the family by the same fraction for genomes from all
 15 other habitats ('out-group'). In cases where no 'out-group' genome encoded a member protein,
 16 the protein family was simply noted as 'exclusive' to the 'in-group' habitat. In all cases, we
 17 calculated the Fisher's exact statistic using the *fisher_exact* function in scipy.stats
 18 (alternative='two-sided'). To account for discrepancies in genome sampling among lineages, we
 19 determined ratios and corresponding statistical significance values separately for each lineage.
 20 All statistical comparisons for a given lineage were corrected for false discovery rate using the
 21 *multipletests* function in statsmodels.stats.multitest (method="fdr_bh"). Finally, we selected
 22 families that were predicted to be enriched or depleted in particular habitats. We considered
 23 enriched families to be those with ratios ≥ 5 , and depleted families as those that were encoded
 24 in 10% or fewer of genomes from a given habitat but present in 50% or more of genomes
 25 outside that environmental category. Retaining only those comparisons with corrected Fisher's
 26 statistics at 0.05 or below resulted in a set of 926 unique, differentially distributed protein
 27 families for downstream analysis.

28 29 *Analysis of putative rhodopsins*

30
 31 Protein sequences from the CPR (fam11249) were combined with a set of reference protein
 32 sequences spanning Type 1 bacterial/archaeal rhodopsin and heliorhodopsin (107). Sequences
 33 were then aligned using MAFFT (--auto) and a tree was inferred using FastTreeMP (default
 34 parameters). Alignment columns with 95% or more gaps were trimmed manually in Geneious
 35 for the purposes of visualization. Transmembrane domains were identified by BLASTp searches
 36 (<https://blast.ncbi.nlm.nih.gov/Blast.cgi>) and conserved residues were defined by manual
 37 comparison with an annotated alignment of previously published reference sequences (108).

38 39 *Processes driving protein family evolution*

40
 41 To examine the evolutionary processes shaping the differentially-distributed protein families, we
 42 next subjected each family to an automated gene-species tree reconciliation workflow adapted
 43 from (90). Briefly, for each family, truncated sequences (defined as those with lengths less than

1 2 standard deviations from the family mean) were removed and the remaining sequences
2 aligned with MAFFT (*--retree 2*). Resulting alignments were then trimmed using trimal (*-gt 0.1*)
3 and used to infer maximum-likelihood phylogenetic trees using IQTree with 1000 ultrafast
4 bootstrap replicates (*-bnni -m TEST -st AA -bb 1000 -nt AUTO*). We removed reference
5 sequences from the inferred species tree and rooted it on the branch separating
6 Saccharibacteria from the monophyletic clade containing Gracilibacteria and
7 Absconditabacteria. A random sample of 100 bootstrap replicates were then used to
8 probabilistically reconcile each protein family with the pruned species tree using the ALE
9 package (*ALE_undated*) (91). Estimates of missing gene fraction were derived from the checkM
10 genome completeness estimates described above. We then calculated the total number of
11 originations (horizontal gene transfer from non-CPR, or *de novo* gene formation), within-CPR
12 horizontal transfers, and losses over each non-terminal branch and mapped branch-wise counts
13 for each event to a species-tree cladogram in iTol (104).

14

SUPPLEMENTARY MATERIAL:

16

17 All supplementary figures, tables, and files are available through Zenodo
18 (<https://doi.org/10.5281/zenodo.4586041>).

19

20 **Fig. S1.** Proteome size as a function of genome completeness and habitat of origin for
21 Absconditabacteria and Gracilibacteria.

22

23 **Fig. S2.** Principal coordinates analysis based on all protein families with 5 or more members
24 among **ab**) all lineages, **c**) Absconditabacteria, and **d**) Gracilibacteria.

25

26 **Fig. S3.** Phylogenetic relationships and trimmed protein alignment among Type 1 rhodopsins
27 and CPR rhodopsin homologs. CPR sequences from the Absconditabacteria (NDQ motif),
28 Gracilibacteria (DTE motif), and Saccharibacteria (DTS motif) are highlighted in the tree.
29 Sequence conservation at each aligned site, the location of bacteriorhodopsin (BR) site 85, 89,
30 96, and the location of retinal-binding lysine (Schiff Base linkage) are also indicated.

31

32 **File S1.** Absconditabacteria, Gracilibacteria, and Saccharibacteria species tree based on 16
33 syntenic ribosomal proteins (newick format).

34

35 **File S2.** FASTA-formatted files for 3,787 protein families with 5 or more member sequences.

36

37 **Table S1.** Characteristics of genomes used in this study, including environmental metadata and
38 accession information.

39

40 **Table S2.** Read accession and mapping information for genomes included in the global
41 abundance analysis.

42

Table S3. Metadata and ribosomal protein S3 (rpS3) information for experimentally validated Saccharibacteria-Actinobacteria pairs. Only strains with publicly available reference genomes and detectable rpS3 sequences are listed.

Table S4. Pairs of Saccharibacteria-Actinobacteria species groups with exclusive co-occurrence at the species group or order level in the included metagenomic samples.

Table S5. Top 25 coverage correlations between Absconditabacteria and other organisms (as denoted by representative rpS3-bearing scaffolds) across 22 cyanobacterial mat metagenomes.

Table S6. Characteristics of the 3,787 protein families with 5 or more member sequences.

Table S7. Characteristics of protein families that are statistically enriched/depleted across habitat categories.

Acknowledgments:

We thank Shufei Lei, Lily Law, Alex Crits-Christoph, Tom Williams, Oded Béjà, Adair Borges, Raphaël Méheust, Alison Sharrar, Alexa Nicolas, Jett Liu, and Simonetta Gribaldo for informatics support, helpful discussions, and comments on the manuscript. We thank Alex Thomas, Ariel Amadio, Mircea Podar, Ramunas Stepanauskas, Connor Skennerton, Stefano Campanaro, Cédric Laczny, Paul Wilmes, Clara Chan, Scott E. Miller, Lauren C. Kennedy, Rose S. Kantor, Kara L. Nelson, Lauren Lui, Maliheh Mehrshad, Chris Greening, Mads Albertsen, and Sari Peura for permission to use genomic data that were unpublished at the time of writing.

Christine He was funded by a Camille & Henry Dreyfus Environmental Chemistry Postdoctoral Fellowship. Patrick Munk was supported by the Danish Veterinary and Food Administration and The Novo Nordisk Foundation (NNF16OC0021856). JAEA was funded by the Ministry of Economy, Trade and Industry of Japan, as “The Project for Validating Near-field System Assessment Methodology in Geological Disposal System”. Keith Bouma-Gregson and cyanobacterial mat sample collection was supported by the National Science Foundation’s Eel River Critical Zone Observatory [EAR-1331940], Department of Energy grant [DOE-SC10010566], NSF Division of Environmental Biology [1656009], and US EPA STAR Fellowship [91767101-0]. Rohan Sachdeva and thiocyanate reactor genome construction was supported by a grant from the National Science Foundation (USA) to JFB (EAR-1349278). We also thank the Innovative Genomics Institute at UC Berkeley.

Author’s contributions: A.L.J and J.F.B. compiled the dataset, performed genome curation and analysis, developed the project, and wrote the manuscript. C.H., R.K., L.E.V.A., P.M., K.B.G., I.F.F., Y.A., R.S., and P.T.W. generated data for the study and provided comments on the manuscript.

Data and software availability: All accession information for the genomes analyzed in this study are listed in Supplementary Table 1. Genomes as well as custom code for the described analyses are also available on GitHub: <https://github.com/alexanderjaffe/cpr-crossenv>. All supplementary figures, tables, and files are available through Zenodo (<https://doi.org/10.5281/zenodo.4586041>).

REFERENCES:

1. Wrighton KC, Thomas BC, Sharon I, Miller CS, Castelle CJ, VerBerkmoes NC, Wilkins MJ, Hettich RL, Lipton MS, Williams KH, Long PE, Banfield JF. 2012. Fermentation, hydrogen, and sulfur metabolism in multiple uncultivated bacterial phyla. *Science* 337:1661–1665.
2. Albertsen M, Hugenholtz P, Skarshewski A, Nielsen KL, Tyson GW, Nielsen PH. 2013. Genome sequences of rare, uncultured bacteria obtained by differential coverage binning of multiple metagenomes. *Nat Biotechnol* 31:533–538.
3. Kantor RS, Wrighton KC, Handley KM, Sharon I, Hug LA, Castelle CJ, Thomas BC, Banfield JF. 2013. Small genomes and sparse metabolisms of sediment-associated bacteria from four candidate phyla. *MBio* 4:e00708–13.
4. McLean JS, Bor B, Kerns KA, Liu Q, To TT, Solden L, Hendrickson EL, Wrighton K, Shi W, He X. 2020. Acquisition and Adaptation of Ultra-small Parasitic Reduced Genome Bacteria to Mammalian Hosts. *Cell Rep* 32:107939.
5. Bor B, Bedree JK, Shi W, McLean JS, He X. 2019. Saccharibacteria (TM7) in the Human Oral Microbiome. *J Dent Res* 98:500–509.
6. Abusleme L, Dupuy AK, Dutzan N, Silva N, Burleson JA, Strausbaugh LD, Gamonal J, Diaz PI. 2013. The subgingival microbiome in health and periodontitis and its relationship with community biomass and inflammation. *ISME J* 7:1016–1025.
7. Campbell JH, O'Donoghue P, Campbell AG, Schwientek P, Sczyrba A, Woyke T, Söll D, Podar M. 2013. UGA is an additional glycine codon in uncultured SR1 bacteria from the human microbiota. *Proc Natl Acad Sci U S A* 110:5540–5545.
8. Hanke A, Hamann E, Sharma R, Geelhoed JS, Hargesheimer T, Kraft B, Meyer V, Lenk S, Osmer H, Wu R, Others. 2014. Recoding of the stop codon UGA to glycine by a BD1-5/SN-2 bacterium and niche partitioning between Alpha- and Gammaproteobacteria in a tidal sediment microbial community naturally selected in a laboratory chemostat. *Front Microbiol* 5:231.
9. Méheust R, Burstein D, Castelle CJ, Banfield JF. 2019. The distinction of CPR bacteria from other bacteria based on protein family content. *Nat Commun* 10:4173.
10. Jaffe AL, Castelle CJ, Matheus Carnevali PB, Gribaldo S, Banfield JF. 2020. The rise of diversity in metabolic platforms across the Candidate Phyla Radiation. *BMC Biology*.
11. Sieber CMK, Paul BG, Castelle CJ, Hu P, Tringe SG, Valentine DL, Andersen GL, Banfield JF. 2019. Unusual metabolism and hypervariation in the genome of a Gracilibacteria

- 1 (BD1-5) from an oil degrading community. bioRxiv.
- 2 12. Starr EP, Shi S, Blazewicz SJ, Probst AJ, Herman DJ, Firestone MK, Banfield JF. 2018.
- 3 Stable isotope informed genome-resolved metagenomics reveals that Saccharibacteria
- 4 utilize microbially-processed plant-derived carbon. Microbiome 6:122.
- 5 13. Nicolas AM, Jaffe AL, Nuccio EE, Taga ME, Firestone MK, Banfield JF. 2020. Unexpected
- 6 diversity of CPR bacteria and nanoarchaea in the rare biosphere of rhizosphere-associated
- 7 grassland soil. Cold Spring Harbor Laboratory.
- 8 14. Shaiber A, Willis AD, Delmont TO, Roux S, Chen L-X, Schmid AC, Yousef M, Watson AR,
- 9 Lolans K, Esen ÖC, Lee STM, Downey N, Morrison HG, Dewhirst FE, Mark Welch JL, Eren
- 10 AM. 2020. Functional and genetic markers of niche partitioning among enigmatic members
- 11 of the human oral microbiome. Genome Biol 21:292.
- 12 15. Murugkar PP, Collins AJ, Chen T, Dewhirst FE. 2020. Isolation and cultivation of candidate
- 13 phyla radiation Saccharibacteria (TM7) bacteria in coculture with bacterial hosts. J Oral
- 14 Microbiol 12:1814666.
- 15 16. He X, McLean JS, Edlund A, Yooseph S, Hall AP, Liu S-Y, Dorrestein PC, Esquenazi E,
- 16 Hunter RC, Cheng G, Nelson KE, Lux R, Shi W. 2015. Cultivation of a human-associated
- 17 TM7 phylotype reveals a reduced genome and epibiotic parasitic lifestyle. Proc Natl Acad
- 18 Sci U S A 112:244–249.
- 19 17. Utter DR, He X, Cavanaugh CM, McLean JS, Bor B. 2020. The saccharibacterium TM7x
- 20 elicits differential responses across its host range. ISME J
- 21 <https://doi.org/10.1038/s41396-020-00736-6>.
- 22 18. Lui LM, Nielsen TN, Arkin AP. 2020. A method for achieving complete microbial genomes
- 23 and improving bins from metagenomics data. Cold Spring Harbor Laboratory.
- 24 19. Dueholm MS, Albertsen M, Stokholm-Bjerregaard M, McIlroy SJ, Karst SM, Nielsen PH.
- 25 2015. Complete Genome Sequence of the Bacterium Aalborg_AAW-1, Representing a
- 26 Novel Family within the Candidate Phylum SR1. Genome Announc 3.
- 27 20. Ornaghi M, Prado RM, Ramos TR, Catalano FR, Mottin C, Creevey CJ, Huws SA, Prado
- 28 IN. 2020. Natural plant-based additives can improve ruminant performance by influencing
- 29 the Rumen microbiome. Research Square.
- 30 21. Cabello-Yeves PJ, Zemskaya TI, Zakharenko AS, Sakirko MV, Ivanov VG, Ghai R,
- 31 Rodriguez-Valera F. 2020. Microbiome of the deep Lake Baikal, a unique oxic bathypelagic
- 32 habitat. Limnol Oceanogr 65:1471–1488.
- 33 22. Dudek NK, Sun CL, Burstein D, Kantor RS, Aliaga Goltsman DS, Bik EM, Thomas BC,
- 34 Banfield JF, Relman DA. 2017. Novel Microbial Diversity and Functional Potential in the
- 35 Marine Mammal Oral Microbiome. Curr Biol 27:3752–3762.e6.
- 36 23. Andersen VD, Aarestrup FM, Munk P, Jensen MS, de Knecht LV, Bortolaia V, Knudsen BE,
- 37 Lukjancenko O, Birkegård AC, Vigre H. 2020. Predicting effects of changed antimicrobial
- 38 usage on the abundance of antimicrobial resistance genes in finisher' gut microbiomes.

- 1 Prev Vet Med 174:104853.
- 2 24. Rehman ZU, Fortunato L, Cheng T, Leiknes T. 2020. Metagenomic analysis of sludge and
3 early-stage biofilm communities of a submerged membrane bioreactor. *Sci Total Environ*
4 701:134682.
- 5 25. Youssef NH, Farag IF, Hahn CR, Premathilake H, Fry E, Hart M, Huffaker K, Bird E,
6 Hambricht J, Hoff WD, Elshahed MS. 2019. *Candidatus Krumholzibacterium zodletense*
7 gen. nov., sp nov, the first representative of the candidate phylum Krumholzibacteriota phyl.
8 nov. recovered from an anoxic sulfidic spring using genome resolved metagenomics. *Syst*
9 *Appl Microbiol* 42:85–93.
- 10 26. Poghosyan L, Koch H, Frank J, van Kessel MAHJ, Cremers G, van Alen T, Jetten MSM, Op
11 den Camp HJM, Lückner S. 2020. Metagenomic profiling of ammonia- and
12 methane-oxidizing microorganisms in two sequential rapid sand filters. *Water Res*
13 185:116288.
- 14 27. Hermsdorf AW, Amano Y, Miyakawa K, Ise K, Suzuki Y, Anantharaman K, Probst A, Burstein
15 D, Thomas BC, Banfield JF. 2017. Potential for microbial H₂ and metal transformations
16 associated with novel bacteria and archaea in deep terrestrial subsurface sediments. *ISME*
17 *J* 11:1915–1929.
- 18 28. Rinke C, Schwientek P, Sczyrba A, Ivanova NN, Anderson IJ, Cheng J-F, Darling A, Malfatti
19 S, Swan BK, Gies EA, Dodsworth JA, Hedlund BP, Tsiamis G, Sievert SM, Liu W-T, Eisen
20 JA, Hallam SJ, Kyrpides NC, Stepanauskas R, Rubin EM, Hugenholtz P, Woyke T. 2013.
21 Insights into the phylogeny and coding potential of microbial dark matter. *Nature*
22 499:431–437.
- 23 29. Brown CT, Hug LA, Thomas BC, Sharon I, Castelle CJ, Singh A, Wilkins MJ, Wrighton KC,
24 Williams KH, Banfield JF. 2015. Unusual biology across a group comprising more than 15%
25 of domain Bacteria. *Nature* 523:208–211.
- 26 30. Parks DH, Chuvochina M, Waite DW, Rinke C, Skarshewski A, Chaumeil P-A, Hugenholtz
27 P. 2018. A standardized bacterial taxonomy based on genome phylogeny substantially
28 revises the tree of life. *Nat Biotechnol* 36:996–1004.
- 29 31. Anantharaman K, Brown CT, Hug LA, Sharon I, Castelle CJ, Probst AJ, Thomas BC, Singh
30 A, Wilkins MJ, Karaoz U, Brodie EL, Williams KH, Hubbard SS, Banfield JF. 2016.
31 Thousands of microbial genomes shed light on interconnected biogeochemical processes
32 in an aquifer system. *Nat Commun* 7:13219.
- 33 32. Bouma-Gregson K, Olm MR, Probst AJ, Anantharaman K, Power ME, Banfield JF. 2019.
34 Impacts of microbial assemblage and environmental conditions on the distribution of
35 anatoxin-a producing cyanobacteria within a river network. *ISME J* 13:1618–1634.
- 36 33. Engelberts JP, Robbins SJ, de Goeij JM, Aranda M, Bell SC, Webster NS. 2020.
37 Characterization of a sponge microbiome using an integrative genome-centric approach.
38 *ISME J* 14:1100–1110.
- 39 34. Zhou Z, Tran PQ, Kieft K, Anantharaman K. 2020. Genome diversification in globally

- 1 distributed novel marine Proteobacteria is linked to environmental adaptation. ISME J
2 14:2060–2077.
- 3 35. Pereira FC, Wasmund K, Cobankovic I, Jehmlich N, Herbold CW, Lee KS, Sziranyi B,
4 Vesely C, Decker T, Stocker R, Warth B, von Bergen M, Wagner M, Berry D. 2020. Rational
5 design of a microbial consortium of mucosal sugar utilizers reduces *Clostridiodes difficile*
6 colonization. Nat Commun 11:5104.
- 7 36. Parks DH, Rinke C, Chuvochina M, Chaumeil P-A, Woodcroft BJ, Evans PN, Hugenholtz P,
8 Tyson GW. 2017. Recovery of nearly 8,000 metagenome-assembled genomes substantially
9 expands the tree of life. Nat Microbiol 2:1533–1542.
- 10 37. Probst AJ, Ladd B, Jarett JK, Geller-McGrath DE, Sieber CMK, Emerson JB,
11 Anantharaman K, Thomas BC, Malmstrom RR, Stieglmeier M, Klingl A, Woyke T, Ryan MC,
12 Banfield JF. 2018. Differential depth distribution of microbial function and putative symbionts
13 through sediment-hosted aquifers in the deep terrestrial subsurface. Nat Microbiol
14 3:328–336.
- 15 38. Solden LM, Naas AE, Roux S, Daly RA, Collins WB, Nicora CD, Purvine SO, Hoyt DW,
16 Schückel J, Jørgensen B, Willats W, Spalinger DE, Firkins JL, Lipton MS, Sullivan MB,
17 Pope PB, Wrighton KC. 2018. Interspecies cross-feeding orchestrates carbon degradation
18 in the rumen ecosystem. Nat Microbiol 3:1274–1284.
- 19 39. Robbins SJ, Singleton CM, Chan CX, Messer LF, Geers AU, Ying H, Baker A, Bell SC,
20 Morrow KM, Ragan MA, Miller DJ, Forêt S, ReFuGe2020 Consortium, Voolstra CR, Tyson
21 GW, Bourne DG. 2019. A genomic view of the reef-building coral *Porites lutea* and its
22 microbial symbionts. Nat Microbiol 4:2090–2100.
- 23 40. Martínez Arbas S, Narayanasamy S, Herold M, Lebrun LA, Hoopmann MR, Li S, Lam TJ,
24 Kunath BJ, Hicks ND, Liu CM, Price LB, Laczny CC, Gillece JD, Schupp JM, Keim PS,
25 Moritz RL, Faust K, Tang H, Ye Y, Skupin A, May P, Muller EEL, Wilmes P. 2021. Roles of
26 bacteriophages, plasmids and CRISPR immunity in microbial community dynamics
27 revealed using time-series integrated meta-omics. Nature Microbiology 6:123–135.
- 28 41. He C, Keren R, Whittaker ML, Farag IF, Doudna JA, Cate JHD, Banfield JF. 2021.
29 Genome-resolved metagenomics reveals site-specific diversity of episymbiotic CPR
30 bacteria and DPANN archaea in groundwater ecosystems. Nat Microbiol
31 <https://doi.org/10.1038/s41564-020-00840-5>.
- 32 42. Woodcroft BJ, Singleton CM, Boyd JA, Evans PN, Emerson JB, Zayed AAF, Hoelzle RD,
33 Lamberton TO, McCalley CK, Hodgkins SB, Wilson RM, Purvine SO, Nicora CD, Li C,
34 Frolking S, Chanton JP, Crill PM, Saleska SR, Rich VI, Tyson GW. 2018. Genome-centric
35 view of carbon processing in thawing permafrost. Nature 560:49–54.
- 36 43. Nayfach S, Shi ZJ, Seshadri R, Pollard KS, Kyrpides NC. 2019. New insights from
37 uncultivated genomes of the global human gut microbiome. Nature 568:505–510.
- 38 44. Cross KL, Campbell JH, Balachandran M, Campbell AG, Cooper SJ, Griffen A, Heaton M,
39 Joshi S, Klingeman D, Leys E, Yang Z, Parks JM, Podar M. 2019. Targeted isolation and

- 1 cultivation of uncultivated bacteria by reverse genomics. *Nat Biotechnol* 37:1314–1321.
- 2 45. Almeida A, Nayfach S, Boland M, Strozzi F, Beracochea M, Shi ZJ, Pollard KS, Sakharova
3 E, Parks DH, Hugenholtz P, Segata N, Kyrpides NC, Finn RD. 2021. A unified catalog of
4 204,938 reference genomes from the human gut microbiome. *Nat Biotechnol* 39:105–114.
- 5 46. Fettweis JM, Serrano MG, Brooks JP, Edwards DJ, Girerd PH, Parikh HI, Huang B, Arodz
6 TJ, Edupuganti L, Glascock AL, Xu J, Jimenez NR, Vivadelli SC, Fong SS, Sheth NU, Jean
7 S, Lee V, Bokhari YA, Lara AM, Mistry SD, Duckworth RA 3rd, Bradley SP, Koparde VN,
8 Orenda XV, Milton SH, Rozycki SK, Matveyev AV, Wright ML, Huzurbazar SV, Jackson EM,
9 Smirnova E, Korlach J, Tsai Y-C, Dickinson MR, Brooks JL, Drake JI, Chaffin DO, Sexton
10 AL, Gravett MG, Rubens CE, Wijesooriya NR, Hendricks-Muñoz KD, Jefferson KK, Strauss
11 JF 3rd, Buck GA. 2019. The vaginal microbiome and preterm birth. *Nat Med* 25:1012–1021.
- 12 47. Bandla A, Pavagadhi S, Sridhar Sudarshan A, Poh MCH, Swarup S. 2020. 910
13 metagenome-assembled genomes from the phytobiomes of three urban-farmed leafy Asian
14 greens. *Sci Data* 7:278.
- 15 48. Gibson KM, Nguyen BN, Neumann LM, Miller M, Buss P, Daniels S, Ahn MJ, Crandall KA,
16 Pukazhenthi B. 2019. Gut microbiome differences between wild and captive black
17 rhinoceros - implications for rhino health. *Sci Rep* 9:7570.
- 18 49. Breister AM, Imam MA, Zhou Z, Ahsan MA, Noveron JC, Anantharaman K, Prabhakar P.
19 2020. Soil microbiomes mediate degradation of vinyl ester-based polymer composites.
20 *Communications Materials* 1:101.
- 21 50. Tully BJ, Graham ED, Heidelberg JF. 2018. The reconstruction of 2,631 draft
22 metagenome-assembled genomes from the global oceans. *Sci Data* 5:170203.
- 23 51. Clayton JB, Vangay P, Huang H, Ward T, Hillmann BM, Al-Ghalith GA, Travis DA, Long HT,
24 Van Tuan B, Van Minh V, Cabana F, Nadler T, Toddes B, Murphy T, Glander KE, Johnson
25 TJ, Knights D. 2016. Captivity humanizes the primate microbiome. *Proc Natl Acad Sci U S*
26 *A* 113:10376–10381.
- 27 52. Hu P, Dubinsky EA, Probst AJ, Wang J, Sieber CMK, Tom LM, Gardinali PR, Banfield JF,
28 Atlas RM, Andersen GL. 2017. Simulation of Deepwater Horizon oil plume reveals
29 substrate specialization within a complex community of hydrocarbon degraders. *Proc Natl*
30 *Acad Sci U S A* 114:7432–7437.
- 31 53. Schulze-Makuch D, Wagner D, Kounaves SP, Mangelsdorf K, Devine KG, de Vera J-P,
32 Schmitt-Kopplin P, Grossart H-P, Parro V, Kaupenjohann M, Galy A, Schneider B, Airo A,
33 Frösler J, Davila AF, Arens FL, Cáceres L, Cornejo FS, Carrizo D, Dartnell L, DiRuggiero J,
34 Flury M, Ganzert L, Gessner MO, Grathwohl P, Guan L, Heinz J, Hess M, Keppler F, Maus
35 D, McKay CP, Meckenstock RU, Montgomery W, Oberlin EA, Probst AJ, Sáenz JS, Sattler
36 T, Schirmack J, Sephton MA, Schlöter M, Uhl J, Valenzuela B, Vestergaard G, Wörmer L,
37 Zamorano P. 2018. Transitory microbial habitat in the hyperarid Atacama Desert. *Proc Natl*
38 *Acad Sci U S A* 115:2670–2675.
- 39 54. Munk P, Andersen VD, de Knecht L, Jensen MS, Knudsen BE, Lukjancenko O, Mordhorst H,
40 Clasen J, Agersø Y, Folkesson A, Pamp SJ, Vigre H, Aarestrup FM. 2017. A sampling and

- 1 metagenomic sequencing-based methodology for monitoring antimicrobial resistance in
2 swine herds. *J Antimicrob Chemother* 72:385–392.
- 3 55. Huddy RJ, Sachdeva R, Kadzinga F, Kantor R, Harrison STL, Banfield JF. 2020.
4 Thiocyanate and organic carbon inputs drive convergent selection for specific autotrophic
5 *Afpia* and *Thiobacillus* strains within complex microbiomes. Cold Spring Harbor Laboratory.
- 6 56. Kantor RS, van Zyl AW, van Hille RP, Thomas BC, Harrison STL, Banfield JF. 2015.
7 Bioreactor microbial ecosystems for thiocyanate and cyanide degradation unravelled with
8 genome-resolved metagenomics. *Environ Microbiol* 17:4929–4941.
- 9 57. Probst AJ, Castelle CJ, Singh A, Brown CT, Anantharaman K, Sharon I, Hug LA, Burstein
10 D, Emerson JB, Thomas BC, Banfield JF. 2017. Genomic resolution of a cold subsurface
11 aquifer community provides metabolic insights for novel microbes adapted to high CO₂
12 concentrations. *Environ Microbiol* 19:459–474.
- 13 58. Okazaki Y, Nishimura Y, Yoshida T, Ogata H, Nakano S-I. 2019. Genome-resolved viral and
14 cellular metagenomes revealed potential key virus-host interactions in a deep freshwater
15 lake. *Environ Microbiol* 21:4740–4754.
- 16 59. Lemos LN, Medeiros JD, Dini-Andreote F, Fernandes GR, Varani AM, Oliveira G, Pylro VS.
17 2019. Genomic signatures and co-occurrence patterns of the ultra-small *Saccharimonadia*
18 (phylum CPR/Patescibacteria) suggest a symbiotic lifestyle. *Mol Ecol* 28:4259–4271.
- 19 60. Sharrar AM, Crits-Christoph A, Méheust R, Diamond S, Starr EP, Banfield JF. 2020.
20 Bacterial Secondary Metabolite Biosynthetic Potential in Soil Varies with Phylum, Depth,
21 and Vegetation Type. *MBio* 11.
- 22 61. Zeng Y, Chen X, Madsen AM, Zervas A, Nielsen TK, Andrei A-S, Lund-Hansen LC, Liu Y,
23 Hansen LH. 2020. Potential Rhodopsin- and Bacteriochlorophyll-Based Dual Phototrophy in
24 a High Arctic Glacier. *MBio* 11.
- 25 62. Alteio LV, Schulz F, Seshadri R, Varghese N, Rodriguez-Reillo W, Ryan E, Goudeau D,
26 Eichorst SA, Malmstrom RR, Bowers RM, Katz LA, Blanchard JL, Woyke T. 2020.
27 Complementary Metagenomic Approaches Improve Reconstruction of Microbial Diversity in
28 a Forest Soil. *mSystems* 5.
- 29 63. Vavourakis CD, Mehrshad M, Balkema C, van Hall R, Andrei A-S, Ghai R, Sorokin DY,
30 Muyzer G. 2019. Metagenomes and metatranscriptomes shed new light on the
31 microbial-mediated sulfur cycle in a Siberian soda lake. *BMC Biol* 17:69.
- 32 64. Campanaro S, Treu L, Rodriguez-R LM, Kovalovszki A, Ziels RM, Maus I, Zhu X, Kougias
33 PG, Basile A, Luo G, Schlüter A, Konstantinidis KT, Angelidaki I. 2020. New insights from
34 the biogas microbiome by comprehensive genome-resolved metagenomics of nearly 1600
35 species originating from multiple anaerobic digesters. *Biotechnol Biofuels* 13:25.
- 36 65. Vavourakis CD, Andrei A-S, Mehrshad M, Ghai R, Sorokin DY, Muyzer G. 2018. A
37 metagenomics roadmap to the uncultured genome diversity in hypersaline soda lake
38 sediments. *Microbiome* 6:168.

- 1 66. Wang W, Hu H, Zijlstra RT, Zheng J, Gänzle MG. 2019. Metagenomic reconstructions of gut
2 microbial metabolism in weanling pigs. *Microbiome* 7:48.
- 3 67. Tian R, Ning D, He Z, Zhang P, Spencer SJ, Gao S, Shi W, Wu L, Zhang Y, Yang Y, Adams
4 BG, Rocha AM, Detienne BL, Lowe KA, Joyner DC, Klingeman DM, Arkin AP, Fields MW,
5 Hazen TC, Stahl DA, Alm EJ, Zhou J. 2020. Small and mighty: adaptation of superphylum
6 Patescibacteria to groundwater environment drives their genome simplicity. *Microbiome*
7 8:51.
- 8 68. Keren R, Lawrence JE, Zhuang W, Jenkins D, Banfield JF, Alvarez-Cohen L, Zhou L, Yu K.
9 2020. Increased replication of dissimilatory nitrate-reducing bacteria leads to decreased
10 anammox bioreactor performance. *Microbiome* 8:7.
- 11 69. Cao Y, Xu H, Li R, Gao S, Chen N, Luo J, Jiang Y. 2019. Genetic Basis of Phenotypic
12 Differences Between Chinese Yunling Black Goats and Nubian Goats Revealed by
13 Allele-Specific Expression in Their F1 Hybrids. *Front Genet* 10:145.
- 14 70. Finstad KM, Probst AJ, Thomas BC, Andersen GL, Demergasso C, Echeverría A,
15 Amundson RG, Banfield JF. 2017. Microbial Community Structure and the Persistence of
16 Cyanobacterial Populations in Salt Crusts of the Hyperarid Atacama Desert from
17 Genome-Resolved Metagenomics. *Frontiers in Microbiology*.
- 18 71. Kantor RS, Miller SE, Nelson KL. 2019. The Water Microbiome Through a Pilot Scale
19 Advanced Treatment Facility for Direct Potable Reuse. *Front Microbiol* 10:993.
- 20 72. Beam JP, Becraft ED, Brown JM, Schulz F, Jarett JK, Bezuidt O, Poulton NJ, Clark K,
21 Dunfield PF, Ravin NV, Spear JR, Hedlund BP, Kormas KA, Sievert SM, Elshahed MS,
22 Barton HA, Stott MB, Eisen JA, Moser DP, Onstott TC, Woyke T, Stepanauskas R. 2020.
23 Ancestral Absence of Electron Transport Chains in Patescibacteria and DPANN. *Front*
24 *Microbiol* 11:1848.
- 25 73. Tung J, Barreiro LB, Burns MB, Grenier J-C, Lynch J, Grieneisen LE, Altmann J, Alberts
26 SC, Blekhman R, Archie EA. 2015. Social networks predict gut microbiome composition in
27 wild baboons. *Elife* 4.
- 28 74. Hervé V, Liu P, Dietrich C, Sillam-Dussès D, Stiblik P, Šobotník J, Brune A. 2020.
29 Phylogenomic analysis of 589 metagenome-assembled genomes encompassing all major
30 prokaryotic lineages from the gut of higher termites. *PeerJ* 8:e8614.
- 31 75. UQ eSpace.
- 32 76. Espinoza JL, Harkins DM, Torralba M, Gomez A, Highlander SK, Jones MB, Leong P,
33 Saffery R, Bockmann M, Kuelbs C, Inman JM, Hughes T, Craig JM, Nelson KE, Dupont CL.
34 2018. Supragingival Plaque Microbiome Ecology and Functional Potential in the Context of
35 Health and Disease. *MBio* 9.
- 36 77. Stamps BW, Spear JR. 2020. Identification of Metagenome-Assembled Genomes
37 Containing Antimicrobial Resistance Genes, Isolated from an Advanced Water Treatment
38 Facility. *Microbiol Resour Announc* 9.

- 1 78. Zhou Z, Liu Y, Xu W, Pan J, Luo Z-H, Li M. 2020. Genome- and Community-Level
2 Interaction Insights into Carbon Utilization and Element Cycling Functions of
3 Hydrothermarchaeota in Hydrothermal Sediment. *mSystems* 5.
- 4 79. Mehrshad M, Lopez-Fernandez M, Sundh J, Bell E, Simone D, Buck M, Bernier-Latmani R,
5 Bertilsson S, Dopson M. 2020. Energy efficiency and biological interactions define the core
6 microbiome of deep oligotrophic groundwater. Cold Spring Harbor Laboratory.
- 7 80. Ortiz M, Leung PM, Shelley G, Van Goethem MW, Bay SK, Jordaan K, Vikram S, Hogg ID,
8 Makhalanyane TP, Chown SL, Grinter R, Cowan DA, Greening C. 2020. A genome
9 compendium reveals diverse metabolic adaptations of Antarctic soil microorganisms. Cold
10 Spring Harbor Laboratory.
- 11 81. Buck M, Garcia SL, Vidal LF, Martin G, Martinez Rodriguez GA, Saarenheimo J, Zopfi J,
12 Bertilsson S, Peura S. 2020. Comprehensive dataset of shotgun metagenomes from
13 stratified freshwater lakes and ponds. Cold Spring Harbor Laboratory.
- 14 82. Shaiber A, Eren AM. 2019. Composite Metagenome-Assembled Genomes Reduce the
15 Quality of Public Genome Repositories. *MBio*.
- 16 83. Soro V, Dutton LC, Sprague SV, Nobbs AH, Ireland AJ, Sandy JR, Jepson MA, Micaroni M,
17 Splatt PR, Dymock D, Jenkinson HF. 2014. Axenic culture of a candidate division TM7
18 bacterium from the human oral cavity and biofilm interactions with other oral bacteria. *Appl*
19 *Environ Microbiol* 80:6480–6489.
- 20 84. Moreira D, Zivanovic Y, López-Archilla AI, Iniesto M, López-García P. 2020. Reductive
21 evolution and unique infection and feeding mode in the CPR predatory bacterium
22 *Vampirococcus lugosii*. Cold Spring Harbor Laboratory.
- 23 85. Yamamoto T, Iino H, Kim K, Kuramitsu S, Fukui K. 2011. Evidence for ATP-dependent
24 structural rearrangement of nuclease catalytic site in DNA mismatch repair endonuclease
25 MutL. *J Biol Chem* 286:42337–42348.
- 26 86. Cardenas JP, Quatrini R, Holmes DS. 2016. Aerobic Lineage of the Oxidative Stress
27 Response Protein Rubrerythrin Emerged in an Ancient Microaerobic, (Hyper)Thermophilic
28 Environment. *Frontiers in Microbiology*.
- 29 87. Rissanen AJ, Saarela T, Jäntti H, Buck M, Peura S, Aalto SL, Ojala A, Pumpanen J, Tirola
30 M, Elvert M, Nykänen H. 2020. Vertical stratification patterns of methanotrophs and their
31 genetic controllers in water columns of oxygen-stratified boreal lakes. *FEMS Microbiol Ecol*
32 <https://doi.org/10.1093/femsec/fiaa252>.
- 33 88. Maliar N, Okhrimenko IS, Petrovskaya LE, Alekseev AA, Kovalev KV, Soloviov DV, Popov
34 PA, Rokitskaya TI, Antonenko YN, Zabelskii DV, Dolgikh DA, Kirpichnikov MP, Gordeliy VI.
35 2020. Novel pH-Sensitive Microbial Rhodopsin from *Sphingomonas paucimobilis*. *Dokl*
36 *Biochem Biophys* 495:342–346.
- 37 89. Béjà O, Lanyi JK. 2014. Nature's toolkit for microbial rhodopsin ion pumps. *Proc Natl Acad*
38 *Sci U S A*.

- 1 90. Sheridan PO, Raguideau S, Quince C, Holden J, Zhang L, Thames Consortium, Williams
2 TA, Gubry-Rangin C. 2020. Gene duplication drives genome expansion in a major lineage
3 of Thaumarchaeota. *Nat Commun* 11:5494.
- 4 91. Szöllösi GJ, Rosikiewicz W, Boussau B, Tannier E, Daubin V. 2013. Efficient Exploration of
5 the Space of Reconciled Gene Trees. *Syst Biol* 62:901–912.
- 6 92. Bor B, McLean JS, Foster KR, Cen L, To TT, Serrato-Guillen A, Dewhurst FE, Shi W, He X.
7 2018. Rapid evolution of decreased host susceptibility drives a stable relationship between
8 ultrasmall parasite TM7x and its bacterial host. *Proc Natl Acad Sci U S A*
9 115:12277–12282.
- 10 93. Luef B, Frischkorn KR, Wrighton KC, Holman H-YN, Birarda G, Thomas BC, Singh A,
11 Williams KH, Siegerist CE, Tringe SG, Downing KH, Comolli LR, Banfield JF. 2015. Diverse
12 uncultivated ultra-small bacterial cells in groundwater. *Nat Commun* 6:6372.
- 13 94. Méheust R, Castelle CJ, Matheus Carnevali PB, Farag IF, He C, Chen L-X, Amano Y, Hug
14 LA, Banfield JF. 2020. Groundwater Elusimicrobia are metabolically diverse compared to
15 gut microbiome Elusimicrobia and some have a novel nitrogenase paralog. *ISME J*
16 14:2907–2922.
- 17 95. McCutcheon JP, Moran NA. 2011. Extreme genome reduction in symbiotic bacteria. *Nat*
18 *Rev Microbiol* 10:13–26.
- 19 96. Martijn J, Schön ME, Lind AE, Vosseberg J, Williams TA, Spang A, Ettema TJG. 2020.
20 Hikarchaeia demonstrate an intermediate stage in the methanogen-to-halophile transition.
21 *Nat Commun* 11:5490.
- 22 97. Abby SS, Kerou M, Schleper C. 2020. Ancestral Reconstructions Decipher Major
23 Adaptations of Ammonia-Oxidizing Archaea upon Radiation into Moderate Terrestrial and
24 Marine Environments. *MBio* 11.
- 25 98. Chen L-X, Anantharaman K, Shaiber A, Eren AM, Banfield JF. 2020. Accurate and
26 complete genomes from metagenomes. *Genome Res* 30:315–333.
- 27 99. Parks DH, Imelfort M, Skennerton CT, Hugenholtz P, Tyson GW. 2015. CheckM: assessing
28 the quality of microbial genomes recovered from isolates, single cells, and metagenomes.
29 *Genome Res* 25:1043–1055.
- 30 100. Olm MR, Brown CT, Brooks B, Banfield JF. 2017. dRep: a tool for fast and accurate
31 genomic comparisons that enables improved genome recovery from metagenomes through
32 de-replication. *ISME J* 11:2864–2868.
- 33 101. Aramaki T, Blanc-Mathieu R, Endo H, Ohkubo K, Kanehisa M, Goto S, Ogata H. 2020.
34 KofamKOALA: KEGG Ortholog assignment based on profile HMM and adaptive score
35 threshold. *Bioinformatics* 36:2251–2252.
- 36 102. Criscuolo A, Gribaldo S. 2010. BMGE (Block Mapping and Gathering with Entropy): a
37 new software for selection of phylogenetic informative regions from multiple sequence

- 1 alignments. BMC Evol Biol 10:210.
- 2 103. Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ. 2015. IQ-TREE: a fast and effective
3 stochastic algorithm for estimating maximum-likelihood phylogenies. Mol Biol Evol
4 32:268–274.
- 5 104. Letunic I, Bork P. 2016. Interactive tree of life (iTOL) v3: an online tool for the display and
6 annotation of phylogenetic and other trees. Nucleic Acids Res 44:W242–5.
- 7 105. Boyd JA, Woodcroft BJ, Tyson GW. 2018. GraftM: a tool for scalable, phylogenetically
8 informed classification of genes within metagenomes. Nucleic Acids Res 46:e59.
- 9 106. Capella-Gutiérrez S, Silla-Martínez JM, Gabaldón T. 2009. trimAl: a tool for automated
10 alignment trimming in large-scale phylogenetic analyses. Bioinformatics 25:1972–1973.
- 11 107. Pushkarev A, Inoue K, Larom S, Flores-Urbe J, Singh M, Konno M, Tomida S, Ito S,
12 Nakamura R, Tsunoda SP, Philosof A, Sharon I, Yutin N, Koonin EV, Kandori H, Béjà O.
13 2018. A distinct abundant group of microbial rhodopsins discovered using functional
14 metagenomics. Nature 558:595–599.
- 15 108. Hasegawa M, Hosaka T, Kojima K, Nishimura Y, Nakajima Y, Kimura-Someya T,
16 Shirouzu M, Sudo Y, Yoshizawa S. 2020. A unique clade of light-driven proton-pumping
17 rhodopsins evolved in the cyanobacterial lineage. Sci Rep 10:16752.