

Functional and genetic markers of niche partitioning among enigmatic members of the human oral microbiome

Alon Shaiber^{1,2}, Amy D. Willis³, Tom O. Delmont⁴, Simon Roux⁵, Lin-Xing Chen⁶, Abigail
 5 C. Schmid⁷, Mahmoud Yousef⁸, Andrea R. Watson^{1,9}, Karen Lolans¹, Özcan C. Esen¹,
 Sonny T. M. Lee¹⁰, Nora Downey¹¹, Hilary G. Morrison¹¹, Floyd E. Dewhirst^{12,13}, Jessica
 L. Mark Welch^{11,*}, A. Murat Eren^{1,2,9,11,*}

¹ Department of Medicine, University of Chicago, Chicago, IL 60637, USA; ² Graduate Program in Biophysical Sciences, University of Chicago, Chicago, IL 60637, USA; ³ Department of Biostatistics, University of Washington, Seattle WA 98195, USA; ⁴ Génomique Métabolique, Genoscope, Institut François Jacob, CEA, CNRS, Univ Evry, Université Paris-Saclay, 91057 Evry, France; ⁵ Department of Energy Joint Genome Institute, Berkeley CA 94720, USA; ⁶ Department of Earth and Planetary Sciences, University of California, Berkeley, CA 94720, USA; ⁷ Undergraduate Student, Computational and Applied Mathematics, University of Chicago, Chicago, IL 60637, USA; ⁸ Undergraduate Student, Computer Science, University of Chicago, Chicago, IL 60637, USA; ⁹ Committee on Microbiology, University of Chicago, Chicago, IL 60637, USA; ¹⁰ Division of Biology, Kansas State University, Manhattan, KS 66506, USA; ¹¹ Josephine Bay Paul Center for Comparative Molecular Biology and Evolution, Marine Biological Laboratory, Woods Hole, MA 02543, USA; ¹² Department of Microbiology, The Forsyth Institute, Cambridge, MA 02142, USA; ¹³ Department of Oral Medicine, Infection and Immunity, Harvard School of Dental Medicine, Boston, MA 02115, USA

* Correspondence: meren@uchicago.edu and jmarkwelch@mbl.edu

Running Title

Niche partitioning in the human oral cavity

Keywords

Metagenomics, metapangenomics, niche partitioning, human oral cavity, Candidate Phyla Radiation, Saccharibacteria

Abstract

Microbial residents of the human oral cavity have long been a major focus of microbiology due to their influence on host health and their intriguing patterns of site specificity amidst the lack of dispersal limitation. Yet, the determinants of niche partitioning in this habitat are yet to be fully understood, especially among the taxa that belong to recently discovered branches of microbial life. Here we assembled metagenomes from daily tongue and dental plaque samples from multiple individuals and reconstructed 790 non-redundant genomes, 43 of which resolved to TM7 that formed six monophyletic clades distinctly associated either with plaque or with tongue. Both pangenomic and phylogenomic analyses grouped tongue-specific TM7 clades with other host-associated TM7 genomes. In contrast, plaque-specific TM7 grouped together with environmental TM7 genomes. Besides offering deeper insights into the ecology, evolution, and the mobilome of the cryptic members of the oral microbiome, our study reveals an intriguing resemblance between dental plaque and non-host environments indicated by the TM7 evolution, suggesting that plaque may have served as a stepping stone for environmental microbes to adapt to host environments for some clades of human associated microbes. Additionally, we report that prophages are widespread amongst oral-associated TM7, while absent from environmental TM7, suggesting that prophages may have played a role in adaptation of TM7 to the host environment.

Introduction

Since the inception of microbiology as a new discipline following Antoni van Leeuwenhoek's historical observation of the *animalcules* (Lane, 2015), the human mouth has remained a major focus among microbiologists. The oral cavity is a rich environment with multiple distinct niches in a relatively small space partially due to (1) its diverse anatomy with hard and soft tissue structures (German & Palmer, 2006), (2) the differential influence of the host immunity throughout the oral tissue types (Moutsopoulos & Konkel, 2018), and (3) its constant exposure to exogenous factors. Microbial residents of the oral cavity complement their environment with their own sophisticated lifestyles. Oral microbes form complex communities that show remarkable patterns of horizontal and vertical transmission across humans and animals (Song et al., 2013; Ferretti et al., 2018), temporal dynamism (Caporaso et al., 2011; Mark Welch et al., 2014; Hall et al., 2017), spatial organization (Mark Welch et al., 2016), and site-specificity (Dewhirst et al., 2010; Eren et al., 2014a; Mark Welch, Dewhirst & Borisy, 2019), where they influence the host health (Lamont, Koo & Hajishengallis, 2018) and the ecology of the gastrointestinal tract (Schmidt et al., 2019). Altogether, the oral cavity offers a powerful environment to study the ecology and evolution of microbial systems.

One of the fundamental pursuits of microbiology is to understand the determinants of microbial colonization and niche partitioning that govern the distribution of microbes in their natural habitats. Despite the low dispersal limitation in the human oral cavity that ensures everything could be everywhere, extensive site-specificity among oral microbes has been observed since the earliest studies that used microscopy and cultivation (Socransky & Manganiello, 1971), DNA-DNA hybridization (Mager et al., 2003) and cloning (Aas et al., 2005) strategies. Factors influencing microbial site-specificity include (1) the nature of the underlying substrate (permanent teeth vs. mucosal surfaces), (2) keratinization and other features of the surface topography, (3) proximity to sources of saliva, gingival crevicular fluid, and oxygen, (4) and ability of microbes to adhere both to the substrate and to one another (Socransky & Manganiello, 1971; Gibbons & Houte,

1975; Simón-Soro et al., 2013), overall creating a fascinating environment to study microbial colonization.

Our understanding of the ecology of oral microbes surged thanks to the Human Microbiome Project (HMP) (Human Microbiome Project Consortium, 2012), which generated extensive sequencing data from 9 oral sites sampled from 200 healthy individuals and over 300 reference genomes for bacteria isolated from the human oral cavity. Studies focused on the HMP data confirmed major taxonomic differences between microbial communities associated with dental plaque and mucosal sites in the mouth (Segata et al., 2012; Lloyd-Price et al., 2017). Recruiting metagenomic short reads using single-copy core genes, Donati et al. demonstrated that while some members of the genus *Neisseria* were predominantly found in tongue dorsum samples, others were predominant in plaque samples (Donati et al., 2016), and Eren et al. revealed that even populations of the same species that differed by as little as one nucleotide in 16S rRNA gene amplicons could show extensive site specificity (Eren et al., 2014a). Strong associations between oral sites and their microbial residents even at the finest levels of resolution raise questions regarding the drivers of such exclusiveness (Mark Welch, Dewhirst & Borisy, 2019). However, identifying genetic or functional determinants of site-specificity requires the investigation of microbial pangenomes.

The human oral cavity is one of the most well characterized microbial habitats of the human body. The Human Oral Microbiome Database (HOMD) (<http://www.homd.org>) (Chen et al., 2010; Escapa et al., 2018) describes more than 750 oral taxa based on full-length 16S rRNA gene sequences, 70% of which have cultured representatives, enabling genome-resolved analyses that cover a considerable fraction of oral metagenomes (Nayfach et al., 2016). Yet, one-third of the known oral taxa are missing or poorly represented in culture collections and genomic databases, and include some that are common in the oral cavity (Vartoukian et al., 2016), including members of the Candidate Phyla Radiation (CPR) (Brown et al., 2015), such as *Saccharibacteria* (TM7), *Absconditabacteria* (SR1), and *Gracilibacteria* (GN02). CPR bacteria form distinct branches in the Tree of Life both based on their phylogenetic origins (Hug et al., 2016) and functional makeup (Méheust et al., 2019); they lack many biological pathways that

are considered essential (Brown et al., 2015) and have been shown to rely on epibiotic lifestyles (Bor et al., 2019), with a complex and poorly understood relationship with a microbial host (Bor et al., 2018). Their unique lifestyle (He et al., 2015), diversity and prevalence in the oral cavity (Camanocha & Dewhirst, 2014), association with distinct oral sites (Bor et al., 2019), and potential role in disease (Brinig et al., 2003; Abusleme et al., 2013) make them important clades to characterize for a fuller understanding of the ecology of the oral cavity.

Successful efforts targeting these enigmatic members of the oral microbiome produced the first genomic evidence to better understand their functional potential and ecology. The first genomes for oral TM7 emerged from single-amplified genomics studies (Marcy et al., 2007) and were followed by He et al.'s pioneering work that brought the first TM7 population into culture (He et al., 2015), establishing a deeper understanding of its relationship with an *Actinomyces* host. Additional recent cultivation efforts are proving successful in providing access to a wider variety of oral TM7 (Collins, Murugkar & Dewhirst, 2019; Murugkar, Collins & Dewhirst, 2019; Cross et al., 2019). Recent genome-resolved and single-amplified genomics studies have also produced genomes for oral GN02 and SR1 (Campbell et al., 2013; Espinoza et al., 2018), and recently the first targeted isolation of oral SR1 strains has been reported, but genomes were not produced (Cross et al., 2019). Despite the promise of these studies, our understanding of the ecology and evolution of these fastidious oral clades is incomplete.

Here we investigated phylogenetic and functional markers of niche partitioning of enigmatic members of the oral cavity, with a focus on members of the candidate phylum TM7. We used a metagenomic assembly and binning approach to recover metagenome-assembled genomes (MAGs) from the supragingival plaque and tongue dorsum of healthy individuals, and used long-read sequencing to associate TM7 MAGs with previously identified phylotypes through 16S rRNA sequence comparison. Our genomes represent prevalent and abundant lineages that lack genomic representation in the HOMD and National Center for Biotechnology Information (NCBI) genomic databases, including members of the CPR. Using a multi-omics approach we show that oral TM7 species are split into plaque and tongue specialists, and that plaque TM7 phylogenetically

and functionally associate with environmental TM7, while tongue TM7 associate with TM7 from animal guts. To assess the generality of our results we carried out read recruitment from approximately 200 tongue and 200 plaque Human Microbiome Project (HMP) samples; which confirm that the genomes we identified are prevalent, abundant, and site-specific. Our findings suggest that at least for TM7, dental plaque resembles non-host habitats, while tongue- and gut-associated TM7s are more strongly shaped by the host. In addition, our results shed light on other understudied members of the oral cavity, and allow for better genomic insight into prevalent, yet poorly understood members of the oral microbiome.

Results and Discussion

Genome-resolved analysis of tongue and plaque metagenomes of seven individuals yield 790 non-redundant genomes that represent the majority of microbial DNA in samples

To create a genomic collection of oral microbes, we sampled supragingival plaque and tongue dorsum of seven individuals on four to six consecutive or nearly-consecutive days. Shotgun metagenomic sequencing of the resulting 71 samples yielded 1.7 billion high-quality short-reads (Supplementary Table 1a). We independently co-assembled plaque and tongue samples from each individual to improve our ability to detect rare organisms and to minimize errors associated with single-assemblies (Chen et al., 2019b). The resulting 14 co-assemblies (7 people x 2 sites) contained 267,456 contigs longer than 2,500 nts that described approximately 1,163 million nucleotides and 1,554,807 genes (Supplementary Table 1b). To reconstruct genomes from these metagenomes we used a combination of automatic and manual binning strategies that resulted in 2,463 genome bins. Independent assembly and binning of metagenomes from similar habitats can result in the recovery of multiple near-identical genomes (Raveh-Sadka et al., 2015; Delmont et al., 2018). To increase the accuracy of downstream analyses we employed only the 857 of 2,463 bins that were 0.5 Mbp or larger (Supplementary Table 2g), then removed

redundancy by selecting a single representative for each set of genomes that shared an average nucleotide identity (ANI) of greater than 99.8% (see Methods). This resulted in a final collection of 790 non-redundant genomes (Supplementary Tables 2a-b, 3).

Automatic binning approaches can yield composite genomes that suffer from contamination, influencing downstream ecological and evolutionary insights (Shaiber & Eren, 2019), even when single-copy core genes suggest the absence of an apparent contamination (Chen et al., 2019b). Here we sampled each subject on at least 4 separate days to improve the accuracy of automatic binning through differential coverage (Quince et al., 2017). To further minimize potential errors in automatic binning results, we used *anvi'o* to manually inspect, and when necessary, further refine key genomes in our study by (1) visualizing the change in GC-content and gene taxonomy of each contig, (2) performing *ad hoc* searches of sequences in public databases, and (3) ensuring the agreement across all contigs with respect to sequence composition signal and differential coverage, the coverage of contigs by reads recruited from our metagenomes as well as metagenomes from other studies. Our data report includes each genome bin for interactive inspection (see Data Availability).

After removal of human host DNA contamination, which accounted for 5%-45% of the reads per sample, competitive read recruitment revealed that the final list of genomes recruited nearly half of the reads from our metagenomes (mean 47%, with a range of 10%-74% per sample). Confidently assessing the origins of the remaining short reads is difficult as reconstructing genomes from metagenomes is a challenging task that often leads to incomplete genomic descriptions of complex environments such as the human oral cavity. Factors that influence the MAG recovery include the extent of residual eukaryotic host contamination, the poor assembly of strain mixtures, and mobile genetic elements such as viruses and plasmids which are often difficult to bin. A major driver of the variability we observed in the percentage of reads recruited by our MAGs across samples was the assembly quality, as we found a significant correlation (R^2 : 0.67, t-statistic: 11.9, p-value: $2e-18$) between the percent of reads recruited by the assembled contigs and the percent of reads recruited by MAGs for each metagenome (Supplementary Table 1a, Figure S1a). Interestingly, assemblies differed according to

sample type, where plaque assemblies recruited a significantly higher portion of reads as compared to tongue samples (Wilcoxon sum-rank test, W : 960, p -value: $9.882e-05$, Figure S1b). Additionally, the total number of expected genomes in assemblies, as estimated based on single copy core genes (SCGs), was higher in plaque as compared to tongue samples (Figure S1c). These differences between assemblies of plaque and tongue metagenomes could explain the fact that a larger number of our genomes were derived from plaque samples (463 vs 327), as well as the fact that our collection of 790 genomes recruited a significantly larger fraction of the reads in plaque metagenomes (51.6%) than in tongue metagenomes (38.3%) (z -score: 3.73, p -value: 0.0002). Overall, despite variation between samples, our analysis shows that MAGs encompassed most of the microbial genomic content estimated to be included in each assembly, and represent a large (near 50%) portion of the reads after removal of human DNA.

Metagenome-assembled genomes reveal new lineages including members of the Candidate Phyla Radiation

To assess how taxa represented by our MAGs are distributed relative to known oral taxa, we performed a phylogenomic analysis using our genomes as well as the 1,332 genomes from the HOMD (accessed on August 1st 2018) (Supplementary Table 6b). Our strict criterion of inclusion of genomes with at least 18 of the 37 ribosomal proteins that we used for phylogenomics removed 539 genomes from the analysis, including 492 low completion (<70%) and 23 high completion ($\geq 70\%$) MAGs from our samples, and 24 genomes from the HOMD collection. The 275 MAGs that passed this quality-control threshold covered much of the diversity at the abundant genera of the samples we collected, as evidenced by a comparison between the taxonomy of MAGs (Supplementary Table 2e-f) and the taxonomic composition of metagenomes estimated by short reads (Supplementary Table 4a-h) and 16S rRNA gene amplicon sequencing (Supplementary Table 5a-j and Supplementary Information).

Some lineages contained members exclusively from our collection and not in the HOMD (Figure 1), including 51 genomes that we identified as members of the CPR, which formed a distinct branch, as expected (Figure 1). Our MAGs also included novel genomes from

non-CPR lineages not represented in the HOMD (Figure 1). While some of these deeply
 225 branching MAGs clearly represent novel genomes, it is conceivable that others could
 have been due to contamination that mixes ribosomal proteins from distant populations
 in a single MAG. To guard against this possibility, we performed additional steps of
 manual refinement using public genomic and metagenomic resources (see Methods).

A large fraction of the CPR genomes in our collection belonged to the phylum *Ca.*
 230 *Saccharibacteria* (TM7; 43). The rest were affiliated with the phyla *Ca. Absconditabacteria*
 (*SR1*; 5) and *Ca. Gracilibacteria* (GN02; 3).

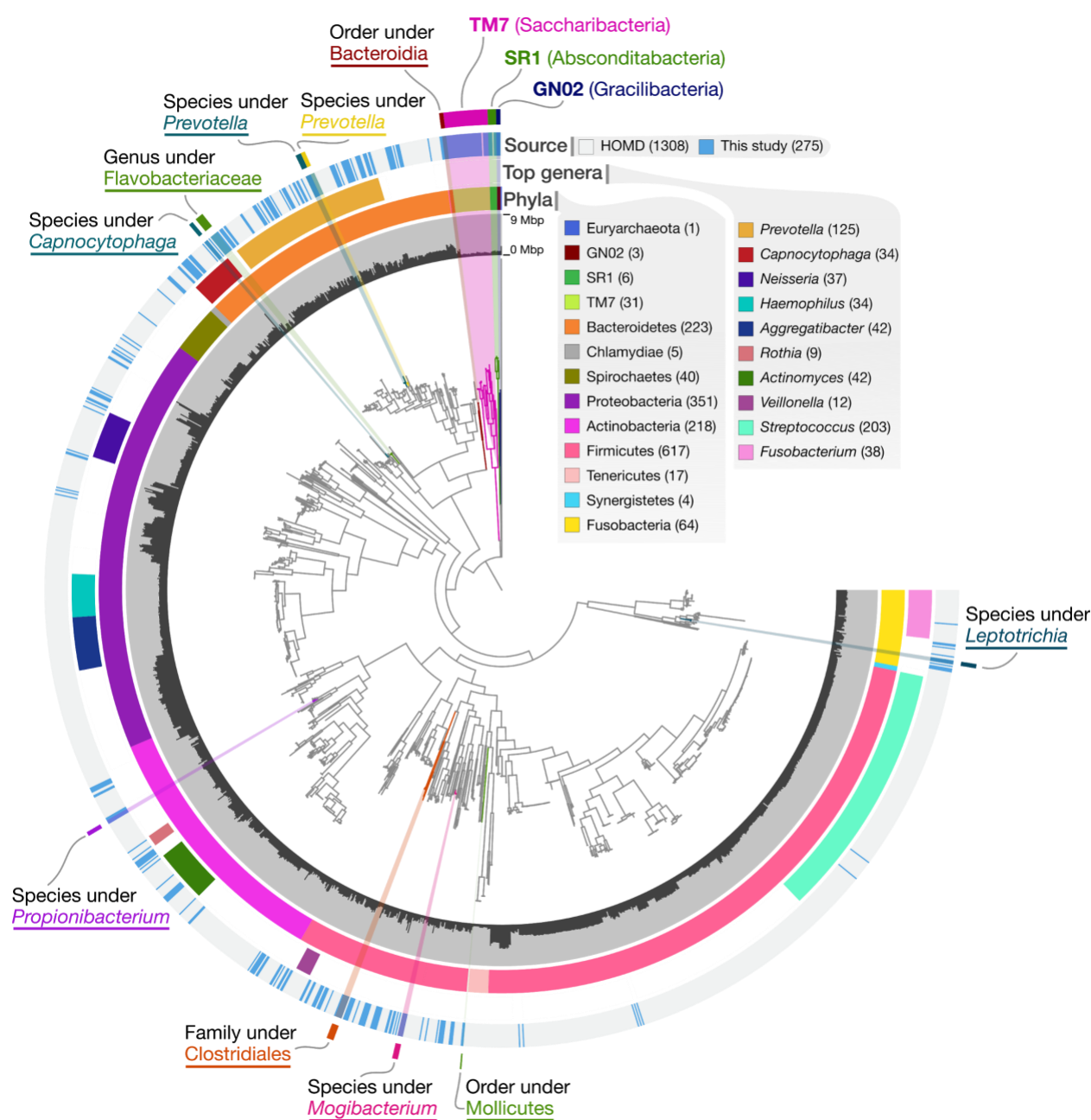


Figure 1: MAGs cover most of the abundant genera of the oral microbiome as well as represent lineages absent in public genomic databases. The dendrogram in the middle of the figure organizes 227 MAGs, 1582 genomes from the HOMD, and a single archaeon, which was used to root the tree, according to their phylogenomic organization based on our collection of ribosomal proteins. The bars in the innermost circular layer represent the length of each genome. The second layer shows the phylum affiliation of each genome. The third layer shows the 10 most abundant genera in our samples as estimated by KrakenUniq. The fourth layer shows the affiliation of genomes as either MAGs from our study (blue) or genomes from HOMD (grey). The outermost layer marks novel genomes of lineages that lack

representation in HMD and NCBI. The lowest taxonomic level that could be assigned using CheckM and sequence search (see Methods) is listed for each novel lineage.

TM7 phylogenomic clades correspond to site of recovery

Our collection included 43 non-redundant TM7 MAGs (Supplementary Table 2b),
 245 presenting an opportunity to investigate associations between their lifestyles (i.e.,
 cosmopolitan or site-specific) and their ancestral relationships. For this, we first examined
 the biogeography of TM7 populations by estimating their relative abundance in each of
 the 71 metagenomes through metagenomic read recruitment (Figure 2a, Supplementary
 Tables 7a-c). We defined a given TM7 population as *detected* in one of the 71 samples if
 250 at least 50% of the nucleotides of the genome were covered by at least one short read.
 We detected 42 of the 43 TM7 populations either only in plaque or only in tongue samples,
 but never in both (Figures 2a, S2a, S2b). The exception was T_C_M_Bin_00022, the
 entirety of which was covered in 4/6 tongue samples and 6/6 plaque samples from
 participant C_M, but not in any other participant (Figure 2a). To investigate the structure
 255 of this population in tongue and plaque habitats, we studied single-nucleotide variants
 (SNVs) at positions of T_C_M_Bin_00022 with a coverage of at least 20X across all
 samples. None of the 22,507 nucleotide positions passed the coverage criterion showed
 any variation in plaque metagenomes. In contrast, 449 of these positions (2%) showed
 notable variation in tongue samples (where the frequency of the minor allele was at least
 260 10%), yet the median frequency of the minor allele was 40% (Supplementary Table 7t).
 This demonstrates (1) the existence of at least two subpopulations where some of them
 are specific to tongue (i.e., those that are represented by minor alleles in tongue
 metagenomes), and (2) the occurrence of the monoclonal plaque population in tongue as
 the dominant member of the subpopulations associated with tongue. Other than this
 265 seemingly “cosmopolitan” population that was present in both tongue and plaque
 metagenomes, all TM7 genomes in our collection appeared to be specialists for plaque
 or tongue habitats (Figure 2).

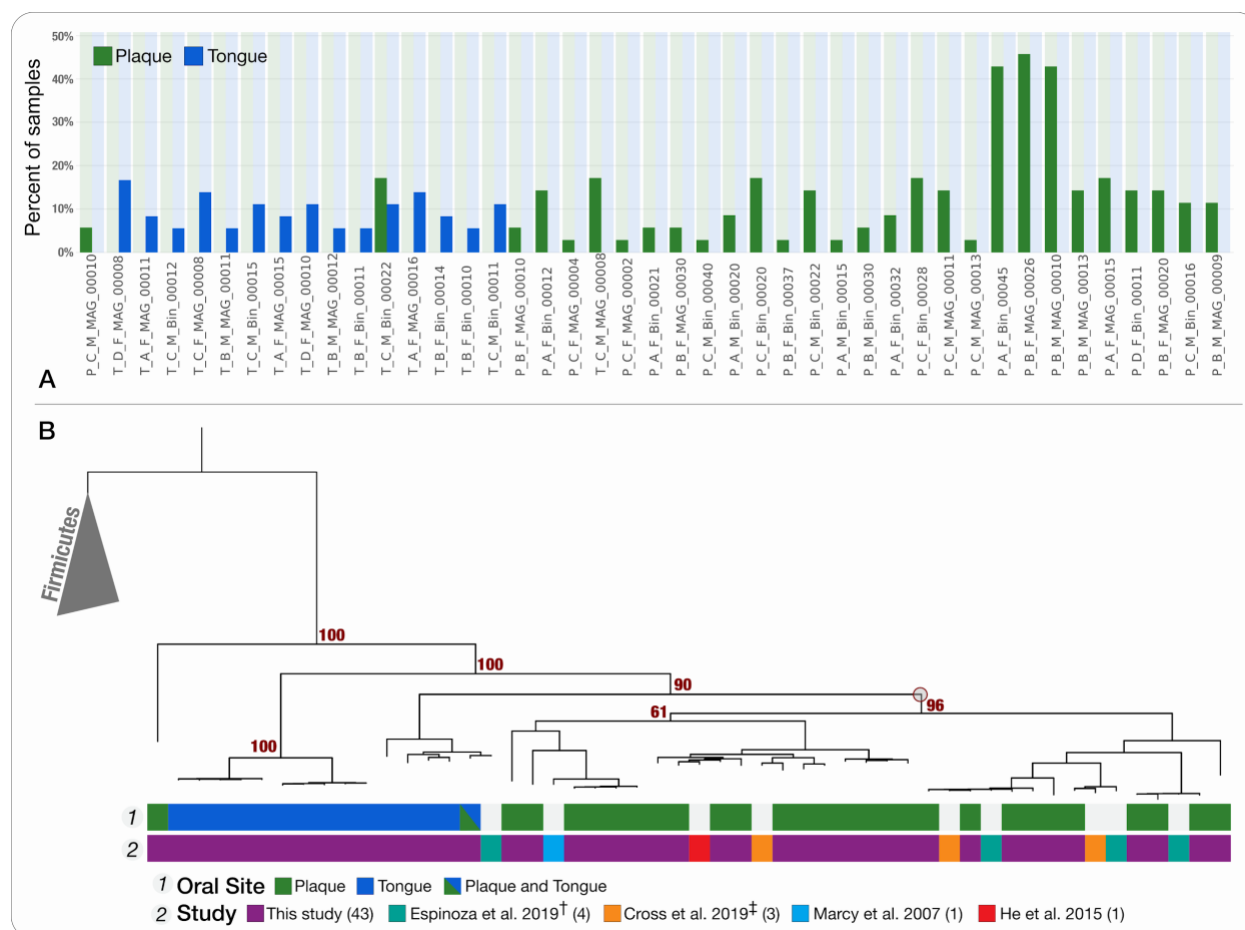


Figure 2: Detection of TM7 genomes across oral metagenomes and their phylogeny. (A) Most TM7 populations are exclusively detected in either tongue or plaque samples in our dataset. For each of the 43 MAGs (on the x-axis) the green and blue bars represent the portion of plaque and tongue samples, respectively, in which it is detected (detection > 0.5). **(B) Phylogenetic organization of TM7 genomes reveals niche-associated oral clades.** The phylogenetic tree includes the 52 oral TM7 genomes (9 of which were previously published), as well as 5 genomes of Firmicutes that root the tree. The layers below the tree describe (top to bottom): “Oral site” - the oral site to which each of our MAGs corresponded, where blue marks tongue dorsum, green marks supragingival plaque and a green-blue combination marks the “cosmopolitan” TM7; “Study” - the study associated with each genome: our MAGs (purple), Espinoza et al. 2018 (teal), Marcy et al. 2007 (blue), He et al. 2015 (red), and Cross et al. 2019 (orange). A red circle appears on the dendrogram and indicates the junction that separates the majority of plaque specialists from tongue specialists, and bootstrap values appear above branches that separate major clusters. † Refined versions of genomes, which we previously published (Shaiber & Eren, 2019). ‡ Genomes from IMG that we refined in this study, but for which accession numbers for refined versions are available in Cross et al. 2019.

We then sought to investigate whether the ancestral relationships among TM7 genomes could explain their intriguing site-specificity. For this, we combined our 43 MAGs with 9 human oral TM7 genomes from the literature. In addition to 3 single amplified genomes

that we downloaded from the Integrated Microbial Genomes and Microbiomes database (IMG/M) (Chen et al., 2019c) and refined (see Methods) and a MAG from Marcy et al. (Marcy et al., 2007), we included 4 MAGs from Espinoza et al. (Espinoza et al., 2018) after manually refining composite TM7 genomes (Shaiber & Eren, 2019), and the first cultivated strain of TM7, TM7x (He et al., 2015) (Supplementary Table 7d). The phylogenomic analysis of these 52 genomes separated tongue and plaque-associated genomes into distinct branches, where we could identify a single node on the tree that separated 41 of the 42 plaque associated genomes, suggesting that the site-specificity of TM7 is an ancestral trait. Another observation emerging from this analysis was that TM7x, which was cultivated from a saliva sample, clustered together with plaque-associated genomes, suggesting that its niche is most likely dental plaque rather than tongue (Figure 2b).

TM7s found in plaque and tongue share exclusive ancestry with environment- and host-associated TM7s

Previous studies have shown that the human associated members of TM7 are polyphyletic, and cluster together with TM7 genomes of environmental origin (Camanocha & Dewhirst; McLean et al., 2018b). Taking advantage of the large number of genomes we have reconstructed, we revisited this observation by performing a phylogenomic analysis using the 156 publicly available TM7 genomes in the NCBI's GenBank database as of 1/16/2019 (Figure 3). We identified six monophyletic human oral clades that were associated either with tongue (T1, T2) or plaque (P1, P2, P3, P4) (Figure 3). Using a pair-wise comparison of the ANI of oral TM7 genomes, we further identified subclades corresponding to genus and species level groups within the six monophyletic clades, including 12 species of TM7 represented each by at least 2 genomes in our collection (Figure 3, Supplementary Tables 7f-h, Supplementary Information).

MAGs reconstructed from metagenomic short-reads often lack ribosomal RNA operons and therefore cannot be identified with known 16S rRNA gene sequences. We surmounted this obstacle by carrying out nanopore long-read sequencing of samples from an additional volunteer (individual L in Supplementary Table 1d). Long reads that

contained both 16S rRNA genes and flanking sequences allowed us to compare our
 315 clades to the 6 previously described TM7 oral groups (G1-G6) based on 16S rRNA gene
 amplicons (Camanocha & Dewhirst). We determined that our monophyletic clades T1,
 T2, and P4 correspond to G3, G6, and G5, respectively (Supplementary Table 7e,i). In
 contrast, clades P1, P2, and P3 all correspond to group G1, showing that G1 is composed
 of at least 3 distinct monophyletic oral clades. We have not recovered any MAGs for TM7
 320 groups G2 and G4, which have been previously shown to have low prevalence as
 compared to other TM7 groups (Bor et al., 2019).

While tongue clades T1 and T2 clustered with genomes recovered from animal gut and
 together formed a deep monophyletic branch of an exclusively host-associated
 superclade shaded blue in Figure 3, plaque clades were interspersed with genomes from
 325 environmental sources (Figure 3). The exceptions to this clear distinction between plaque
 and tongue clades were T_C_M_Bin_00022, a cosmopolitan oral population that
 clustered within the clade T2, and the plaque-associated P_C_M_MAG_00010 (the only
 member of the clade P4) which was placed as a far outlier to all other oral TM7 and
 clustered together with genomes from animal gut (baboon feces). Animal gut-associated
 330 genomes that grouped within the host-associated superclade were recovered
 predominantly from sheep and cow rumen samples, but also included genomes from
 termite gut, mouse colon, and elephant feces, suggesting an ancient association for
 members of the host-associated superclade and their host habitats (Figure 3,
 Supplementary Table 7e). Similarly, the inclusion of genomes recovered from dolphin
 335 dental plaque together with human-plaque-associated TM7 suggests an ancient
 association for plaque-specialists with the dental plaque environment. The phylogenetic
 clustering of tongue-associated TM7 genomes with TM7 genomes from animal gut, to the
 exclusion of environmental TM7, suggests that tongue and gut microbiota share a higher
 degree of ancestral relationship compared to those that are associated with plaque and
 340 with environments outside of a host. We know from previous studies that even though
 microbial community structures and membership in the human oral cavity and gut
 microbiomes are different, the 'community types' observed at these habitats are predictive
 of each other (Ding & Schloss, 2014), suggesting a level of continuity for host influences
 at these distinct sites that shape microbial community succession. Ancestral similarity

between tongue- and gut-associated TM7s compared to those associated with non-host environments suggests that the host factors that influence microbial community succession may also have played a key role in the differentiation of host-associated and non-host-associated branches of TM7. We also know from previous studies that overall microbial community profiles in dental plaque dramatically differs from mucosal sites in the mouth with little overlap in membership (Segata et al., 2012; Eren et al., 2014a). The strong ancestral associations between TM7 clades of plaque and non-host environments, as well as the depletion of plaque specialists from the host-associated superclade, suggest that from a microbial point of view, at least in the context of TM7, dental plaque resembles a non-host environment.

What led to the divergence of TM7 populations? Since TM7 have highly reduced genomes and have been found to be epibionts of other bacteria, primarily Actinobacteria (Kantor et al., 2013; Bor et al., 2019), one reasonable hypothesis is that the bacterial hosts of each TM7 clade are the drivers of the link between TM7 ecology and evolution. Such an hypothesis would imply that the similarity between tongue TM7 and gut TM7 is driven by the colonization of the gut and tongue environments by closely related bacterial hosts that provide a niche for TM7. Furthermore, it would imply the exclusion of such suitable hosts from the plaque environment, and vice versa, it would imply that plaque-specialist TM7 are dependent on bacterial hosts that are absent from the tongue and gut environments. In this context, it is notable that human oral *Actinomyces* species show strong site-specificity and little overlap in membership of dental plaque vs. tongue dorsum inhabitants (Mark Welch et al. 2019) and that Actinobacteria are rare in the human gut (Segata et al., 2012). An alternative hypothesis is that the mechanisms by which TM7 adapt to distinct habitats and distinct bacterial hosts are shaped by independent evolutionary events. While the existence of suitable bacterial hosts is likely an important factor, under this hypothesis, TM7 may acquire “local” bacterial hosts as they adapt to new environments. Our data are not suitable to evaluate either of these hypotheses. Yet given the ancestral similarity between dental plaque TM7 and TM7 from soils and sediments, it is conceivable to hypothesize that the dental plaque environment was able to support environmental TM7, while tongue and gut environments forced a distinct evolutionary path as suggested by the nested monophyletic superclade that is exclusively

associated with host habitats. This depiction of TM7 evolution raises another question about the nature of dental plaque as a host habitat: why is dental plaque not as different from soil and sediment as tongue or gut? It is possible that fixed hard substrate of dental plaque renders it more similar to soils and sediments than to the constantly shedding
380 epithelial surfaces of tongue and gut habitats from a microbial point of view. Whether dental plaque may have served as a stepping stone for environmental microbes by offering them a relatively safe harbor on the human body for host adaptation for some clades of human associated microbes is an intriguing question that warrants further study.

In summary, our data reveal the existence of at least 6 monophyletic oral TM7 clades with
385 clear biogeography within the oral cavity, and a strong divide between the evolutionary history of host-associated and non-host-associated TM7 genomes. Additionally, our analysis reveals 12 species of TM7 that are represented by multiple genomes in our collection and lays the groundwork for definition of taxonomic groups within this candidate phylum. The phylogenomic organization of genomes corresponds to their niche
390 (tongue/plaque) in our dataset, suggesting a link between environmental distribution of these genomes and their evolutionary history in the context of ribosomal proteins.

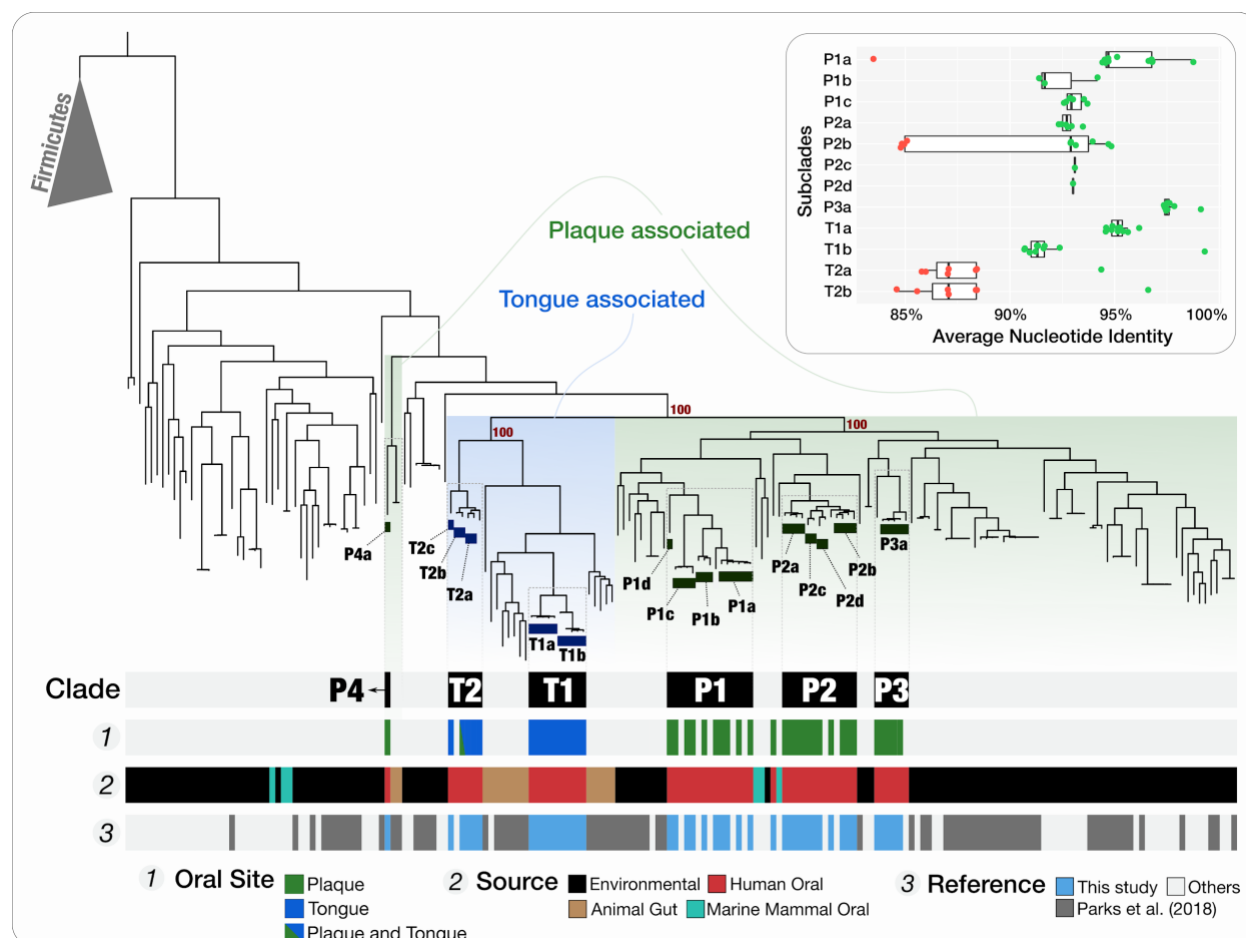


Figure 3: Phylogenetic analysis of human oral TM7 with all TM7 genomes on the NCBI's GenBank shows association of plaque TM7 with environmental genomes, and tongue TM7 with TM7 from animal stool. The phylogenetic tree at the top of the figure was computed using ribosomal proteins and includes 5 Firmicutes as an outgroup. Regions of the tree that are associated with either plaque or tongue clades from Figure 2 are marked with green or blue shaded backgrounds respectively. Bootstrap support values are shown next to branches separating major clusters of oral clades. Subclades are marked with rectangles below the branches they represent. The layers below the tree provide additional information for each genome. From top to bottom: **Clade:** the clade association is shown for each cluster of oral genomes. **Oral Site:** the oral site with which the genome is associated is shown for our MAGs in accordance with Figure 2. **Source:** the source of the genome, where red: human oral, brown: animal gut, cyan: dolphin oral, black: environmental samples. **Reference:** the genomes from this study in blue, and genomes from Parks et al. in grey (Parks et al., 2017). The majority of the rest of the genomes originate from various publications from the Banfield Lab at UC Berkeley. The insert at the top right of the figure shows boxplots for ANI results for genomes in each subclade against all other genomes. Data points represent the ANI score for comparisons in which the alignment coverage was at least 25%. Within-subclade comparisons appear in green and between-subclades comparisons appear in red.

Prevalence of TM7 across individuals is associated with TM7 clades, linking TM7 ecology and evolution

Since the samples we used to generate our 43 *Saccharibacteria* MAGs represent only 7 individuals we next sought to identify whether these patterns were representative of the distribution of TM7 among a wider cohort of healthy individuals. To assess the occurrence of these oral TM7 populations in a larger cohort of healthy individuals, we used a metagenomic short-read recruitment strategy to characterize the distribution of 52 oral TM7 genomes within 413 HMP oral metagenomes (with 30,005,746,488 pairs of reads) that included 196 samples from supragingival plaque and 217 tongue dorsum samples and were sampled from 131 individuals (Supplementary Tables 7j-k). We conservatively defined a genome to be present in a metagenome only if at least 50% of it was covered by at least one short read (see Methods). In addition to oral genomes, we also included three circular TM7 MAGs that were reconstructed from environmental samples and manually curated to circularity (Albertsen et al., 2013; Kantor et al., 2013; Brown et al., 2015). As expected, these 3 environmentally derived genomes (RAAC3, GWC2, and S_aal) were not detected in any oral metagenome (Figure 4, Supplementary Tables 7l-n). The occurrence pattern of TM7 genomes across the HMP individuals matched their occurrence in our seven participants, where all populations except the two genomes of subclade T2_b (T_C_M_Bin_00022, and TM7_MAG_III_B_1) were strongly associated with either tongue or plaque (Figure 4). Members of subclade T2_b indeed appeared to be cosmopolitan and were detected in both plaque and tongue samples (Figure 4). The most prevalent tongue-associated genome and plaque-associated genome were detected in samples from 45% and 50% of the HMP individuals, respectively (Figure 4). In contrast, TM7x, the first cultured strain of TM7, was detected in only 5% of the HMP individuals. While the majority of the samples in the HMP dataset were taken from the tongue dorsum and supragingival plaque, there are additional oral sample types. Our analysis of these additional sample types suggested that certain TM7 populations have a preferential association with oral sites other than the tongue and supragingival plaque (Supplementary Table 7o, Supplementary Information). Of particular notice, the single MAG of clade P4 (group G5), which was previously suggested to associate with

periodontitis (Abusleme et al., 2013) appeared to associate with subgingival plaque, but occurred similarly in subgingival plaque metagenomes of patients with periodontitis and healthy individuals (Supplementary Table 7p-s). These results confirm that the exclusive association of most TM7 oral populations with either plaque or tongue is a general feature and not restricted to the participants of our study and reveal prevalent and abundant tongue and plaque specialists.

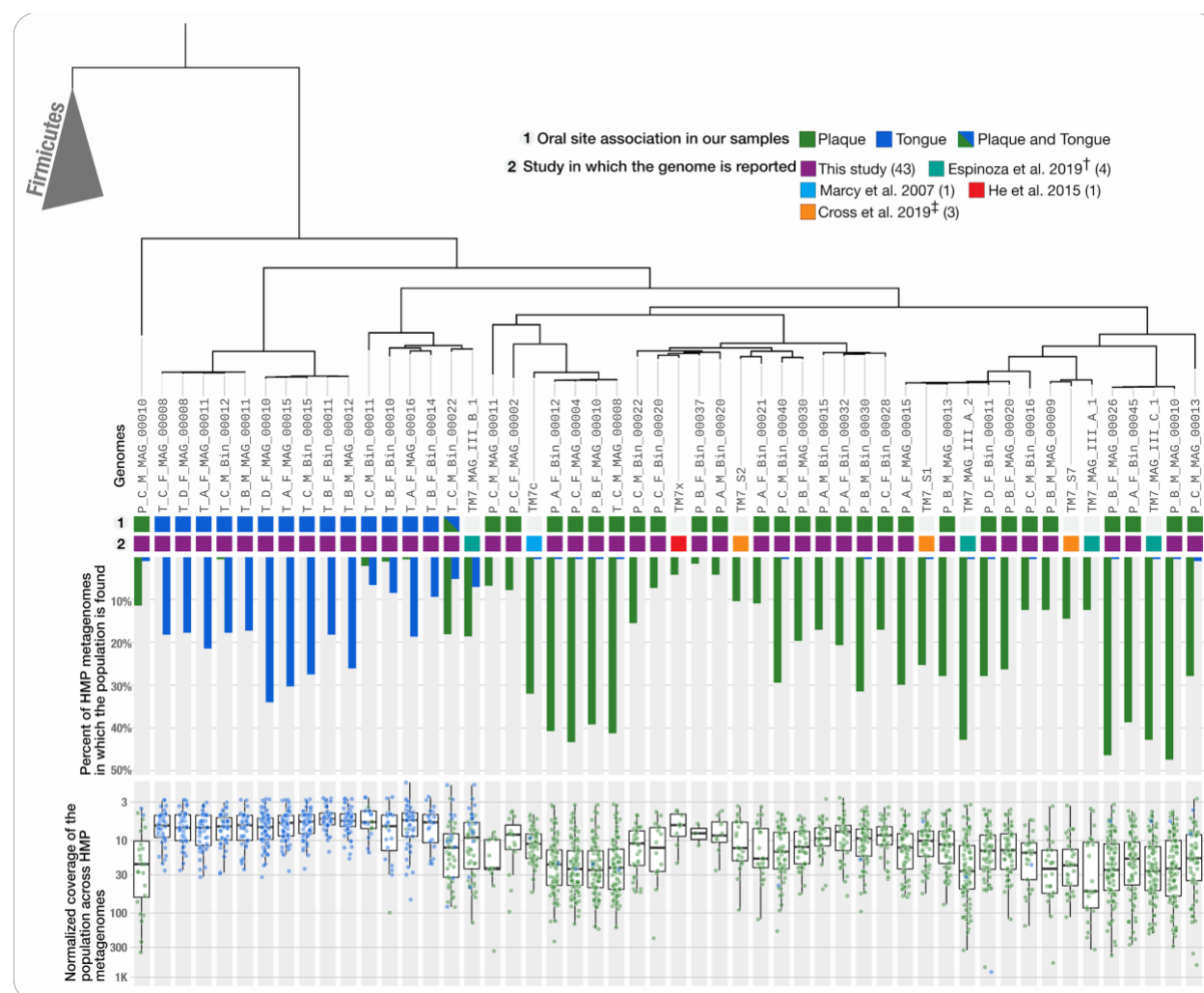


Figure 4: Detection and coverage of TM7 populations in the HMP plaque and tongue samples reveals abundant populations and niche specificity. The tree at the top of the figure and the two layers of information below it are identical to the one in Figure 2. Barplots below the tree show the portion of plaque (green) and tongue (blue) HMP samples in which each TM7 was detected, using a detection threshold of 0.5. Boxplots at the bottom of the figure show the normalized coverages of each TM7 in plaque (green) and tongue (blue) HMP samples in which it was detected.

TM7 pangenome reveals functional markers of niche specificity

Our read recruitment and phylogenomic analyses demonstrated a clear distinction between tongue and plaque specialists in their ecological distribution and evolutionary history, but what are the drivers of this niche specificity? In order to identify functional determinants of niche specificity and investigate the functional differences between members of the various TM7 clades and subclades we utilized a pangenomic approach. Our analysis organized the total 40,832 genes across 55 genomes into 9,117 gene clusters (GCs), which describe one or more homologous genes across genomes based on translated DNA sequences (not to be confused with biosynthetic GCs). Of all GCs, 4,045 were non-singletons (i.e., occurred in at least 2 genomes) and included up to 162 homologous genes from the collection of 55 TM7 genomes described above (Figure 5a, Supplementary Tables 8a-b). The GCs can themselves be clustered into groups that show similar distribution across genomes. By computing the hierarchical clustering of GCs based on their distribution across genomes we identified a collection of 205 core GCs that are found in nearly all genomes, as well as clusters of accessory GCs, many of which were exclusively associated with oral habitats or phylogenomic clades (Figure 5a), confirming that the agreement between phylogenomics and ecology of these genomes was also represented by differentially occurring GCs. The proportion of genes with functional hits varied dramatically between the core and accessory TM7 genes. While more than 90% of core GCs had functional annotations, NCBI's Clusters of Orthologous Groups (COGs) database (Tatusov et al., 2003) annotated only 29% of singletons, and 22% to 88% of other clusters of GCs (Supplementary Table 8c), revealing a vast number of unknown genes.

Whereas phylogenomics infers associations among genomes based on ancestral relationships, pangenomics reveals associations based on gene content (Dutilh et al., 2004), which can emphasize ecological similarities between genomes (Delmont & Eren, 2018), primarily due to the fact that non-singleton accessory genes are the only drivers of hierarchical clustering based on gene content. The hierarchical clustering of TM7 genomes based on GCs predominantly matched their phylogenomic organization (Figure S4b); however, it recapitulated their niche-association better than phylogenomics (Figure

S4b). Specifically, the plaque-associated genome P_C_M_MAG_00010 of the clade P4
 480 (group 'G5'), which is a distant outlier to all other oral TM7 according to phylogeny (Figure
 2b), was placed together with all other plaque-associated TM7 as well as environmental
 TM7 (Figure S4b). GCs driving this organization belong to the 'Extended Core 2' cluster
 of GCs (Figure 5). These GCs are generally characteristic of plaque and environmental
 TM7, but absent from tongue-associated TM7 of clades T1 and T2 (Figure 5,
 485 Supplementary Table 8c). The 31 of 123 GCs of 'Extended Core 2' that are present in
 P_C_M_MAG_00010 and appear to drive the grouping of P_C_M_MAG_00010 with
 other plaque genomes include proteins involved in a variety of functions including stress
 response, metabolism, and cell division, but are particularly enriched for membrane
 proteins and proteins involved in transport across the cell membrane, suggesting that
 490 proteins involved in interaction with the environment play a key role in grouping plaque
 and environmental genomes together. In summary, these data show that phylogenetically
 distinct clades of plaque-associated TM7s as well as environmental genomes are
 grouped together according to their GC composition and that proteins with a potential role
 in interaction with the environment are driving this organization. While an analysis of a
 495 larger variety of environmental and oral genomes would be required for assertion, these
 findings imply an ecological similarity between plaque and non-host environments, at
 least for TM7.

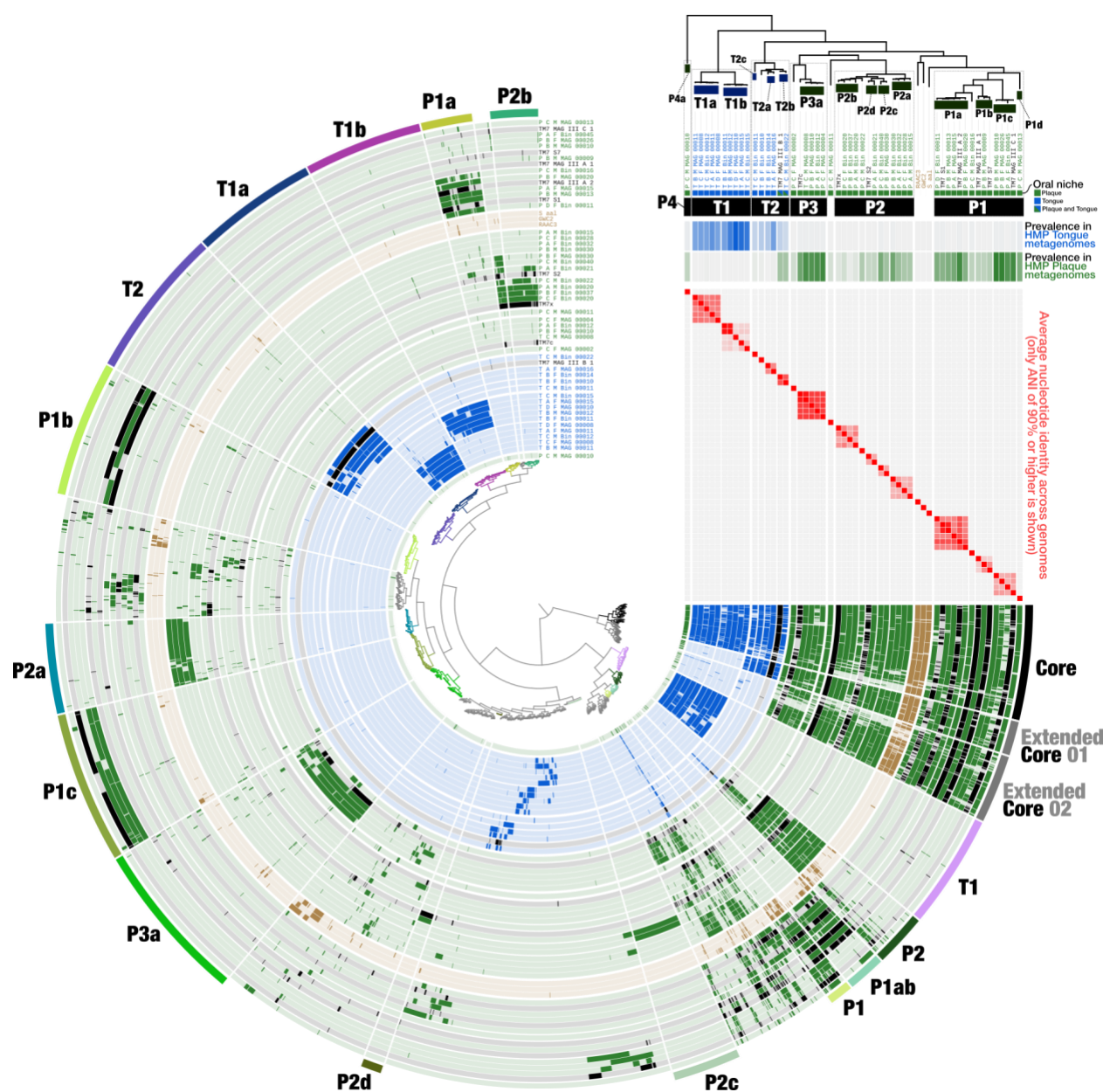


Figure 5: Pangenome of TM7 - Accessory gene-clusters include clade-specific and niche-specific markers. The dendrogram in the center of the figure organizes the 4,045 gene-clusters that occurred in more than one genome according to their frequency of occurrence in the 55 TM7 genomes. The 55 inner layers correspond to the 55 genomes, where our MAGs that were associated with tongue and plaque are shown in blue and green, respectively. Previously published oral and environmental genomes are shown in black and brown, respectively. The data points in the 55 concentric layers show the presence of a gene-cluster in a given genome, and the outermost circular layer highlights clusters of GCs that correspond to the core or to group-specific GCs. Genomes in this figure are ordered according to their phylogenomic organization which is shown at the top-right corner. The three top horizontal layers underneath the phylogenomic tree represent subclade, clade, and oral-site associations of genomes. The next three layers include statistics of coverage for each genome in the HMP oral metagenomes and show (from top to bottom) (1) the maximum interquartile mean coverage, (2) occurrence in tongue samples, and (3) occurrence in supragingival plaque samples.

510 The large number of TM7 genomes we recovered affords the opportunity to investigate key functional properties shared by all TM7s by examining the functions encoded by core GCs. As expected, the TM7 core GCs included many genes involved in translation, replication, and housekeeping (Supplementary Table 8d). The core GCs also included genes involved in amino acid transport. Since TM7s lack the genes to produce their own

515 amino acids (Bor et al., 2019), these genes likely play an important role in scavenging amino acids from the environment or from the bacterial host. The core GCs also included several genes with potential roles in binding to the host, including components of a Type IV pilus system that was identified in all genomes. Oral-associated TM7 have been shown to have a parasitic lifestyle in which they attach to the surface of their bacterial host (He

520 et al., 2015; Cross et al., 2019), but the mechanism utilized for this attachment is unknown. Type IV pilus systems have been found to be enriched in CPR genomes as compared to other bacteria (Méheust et al., 2019) and were also specifically noted in TM7 genomes (Marcy et al., 2007). Type IV pilus systems are involved in many functions, including adherence (Craig, Forest & Maier, 2019), and could potentially be utilized by

525 TM7 to attach to the bacterial host. Most of the components of the type IV pilus system we detected in the TM7 genomes occurred in a single operon with conserved gene synteny (Figure 6a). Additional copies of some of the type IV pilus proteins appear in various loci of the genome (Supplementary Table 8a). We found that while the cytosolic components of the type IV pilus system (PilT, PilB, PilC, PilM) were highly conserved

530 across all genomes, components involved in the alignment of the system in the peptidoglycan (PilN) and the major and minor pilin proteins (PilE, and PilV) appeared in clade or subclade-specific gene clusters and were completely absent from all genomes of clade T1 and from the single genome of clade P4 (Figure 6a, Supplementary Table 8d). Variability in PilV has been shown in the past to confer binding specificity (Ishiwa &

535 Komano, 2003) and in the case of TM7, the clade-specific nature of PilV and PilE sequences could be driven by host-specificity. While T1 genomes were lacking the components of the pilus system with known adhesive roles, they were highly enriched in proteins with a Leucine-rich repeat (LRR) (COG4886), which are often found in membrane bound proteins that are involved in adherence (Bella et al., 2008). 104 of the

540 207 proteins that were annotated with an LRR belonged to a single gene-cluster

(GC_00000003) which was exclusively associated with T1 genomes, and each T1 genome had a total of 12-24 LRR proteins (COG4886, Supplementary Table 8a). In summary, our analyses suggest that the diversity of pilin proteins could be driven by the host-specificity of TM7 species, and that TM7 species that lack pilin proteins could rely on alternative mechanisms such as LRR proteins for adherence.

Additional proteins that we identified with a potential role in host attachment included proteins with a LysM repeat, which is a motif found in a wide range of proteins that are involved in binding to peptidoglycans (Buist et al., 2008). So far, the identified hosts of TM7 are all Gram-positive bacteria, and hence peptidoglycan binding could be a mechanism in which TM7 attach to their hosts. We found 33 GCs associated each with one of four COG functions that included LysM repeats and comprised a total of 169 genes (91 with COG0739, 6 with COG0741, 71 with COG1388, 1 with COG1652). We identified a Murein DD-endopeptidase MepM with a LysM domain (COG0739) in most genomes within a conserved operon, which included components of a Type IV Secretion system including VirB4 and VirB6 (Supplementary Table 8a). Similar to what we observed for the type IV pilus system, the cytosolic component Virb4 was highly conserved across all genomes, while the membrane bound Virb6 varied and appeared to be clade (and even subclade) -specific. This secretion system is also associated with motility in gram-positive bacteria (Marcy et al., 2007), and could potentially be used by TM7 for motion, and/or translocation from one host to another. We detected an additional protein with a LysM repeat (COG1388) in nearly all genomes. While in most genomes this proteins was flanked by genes involved in cell division, in the genomes of Clade T1_b, this locus included an insertion of 1-3 copies of a Leucine-rich repeat (LRR) protein, which as we mentioned above, also has a potential role in adherence. Overall, proteins with a LysM domain are common amongst oral TM7 and could provide another mechanism for attachment to the host surface.

The occurrences of functions across phylogenetic clades could reveal lifestyle differences that are not necessarily highlighted by the occurrences of gene-clusters. Since gene-clusters in a pangenome describe genes that are highly conserved in sequence space, identical functions can occur in distinct gene-clusters, rendering it difficult to describe core

and accessory functions in a pangenome based on core and accessory gene-clusters. Here we developed a statistical approach for functional enrichment analyses in pangenomes (see Methods) that employs logistic regression to reveal enriched functions in any given subset of genomes (e.g., those that match to a phylogenomic clade).

Of the 972 unique functions, we identified 320 (34%) as the functional core, which included genes predominantly identified in all genomes, and 300 that were significantly enriched (q -value < 0.05) in specific clades (Figure S4c and Supplementary Table 8q-r). While there was a wide overlap between core functions and core GCs, 131 core functions occurred in clade-specific GCs, suggesting that these core functions have undergone selective pressure in a clade-specific manner. Of these 131 core functions, 21 were exclusively associated with one GC from the 'Extended Core 1' cluster and one GC from the 'T1' cluster. The large number of core functions with a T1-specific variant further demonstrates the uniqueness of clade T1 amongst the oral TM7 genomes (Figure S4c, supplementary table 8a). Other cases revealed additional functions that may have undergone selective pressure in a clade-specific manner. For example, a single copy of an RTX toxin-related Ca^{2+} -binding protein, was highly conserved in nearly all genomes (gene cluster GC_00000221), but appeared to have a slightly different variant in genomes of the P1_c (GC_00001826) and T2 (GC_00001332) clades. In addition to clade-specific variation in sequences of core functions, we identified accessory functions associated with specific clades. Our examination of the top 100 most enriched functions revealed many membrane associated genes, including, but not limited to, functions that were highlighted above by our examination of GCs (Supplementary Table 8f). For example, tongue and plaque clades appeared to be differentially enriched for transporters of ions and metals. Genes involved in respiration as well as genes involved in translation and stress response were also differentially enriched for tongue and plaque clades. Overall, our analysis of the functional composition of oral TM7 shows that along with differences in accessory functions, sequence divergence of particular core genes distinguishes various clades, and in particular highlights members of clade T1 as outliers amongst the TM7 oral clades, matching their deep phylogenetic position. In addition, we identified functions that characterize tongue and plaque clades and could provide targets for future endeavors to understand the unique biological features of members of each clade.

Overall, our data show that both accessory functions and core functions distinguish plaque and tongue specialists. While the core genome includes many functions common to all bacteria, it also includes many functions that are known to be enriched in CPR
605 genomes. In particular, our data reveal proteins with potential roles in adherence, and suggests that while cytosolic components are highly conserved, extracellular proteins appear to be clade-specific, suggesting that interaction with the host and with the environment are important drivers in differentiating between TM7 oral clades. While members of clade T1 appear as outliers that differ both in functional composition and in
610 the sequence divergence of many core functions as compared to other oral TM7, the functional composition of members of clade T2, which includes the cosmopolitan T2_b genomes, appears to represent an intermediate between the strictly host-associated group and the plaque/environmental group. In addition, plaque-specialists that are phylogenetically distinct are functionally related and group together with environmental
615 genomes based on GCs, while tongue-specialists form separate clusters, suggesting that a stronger ecological similarity exists between plaque and non-host environments as compared to tongue.

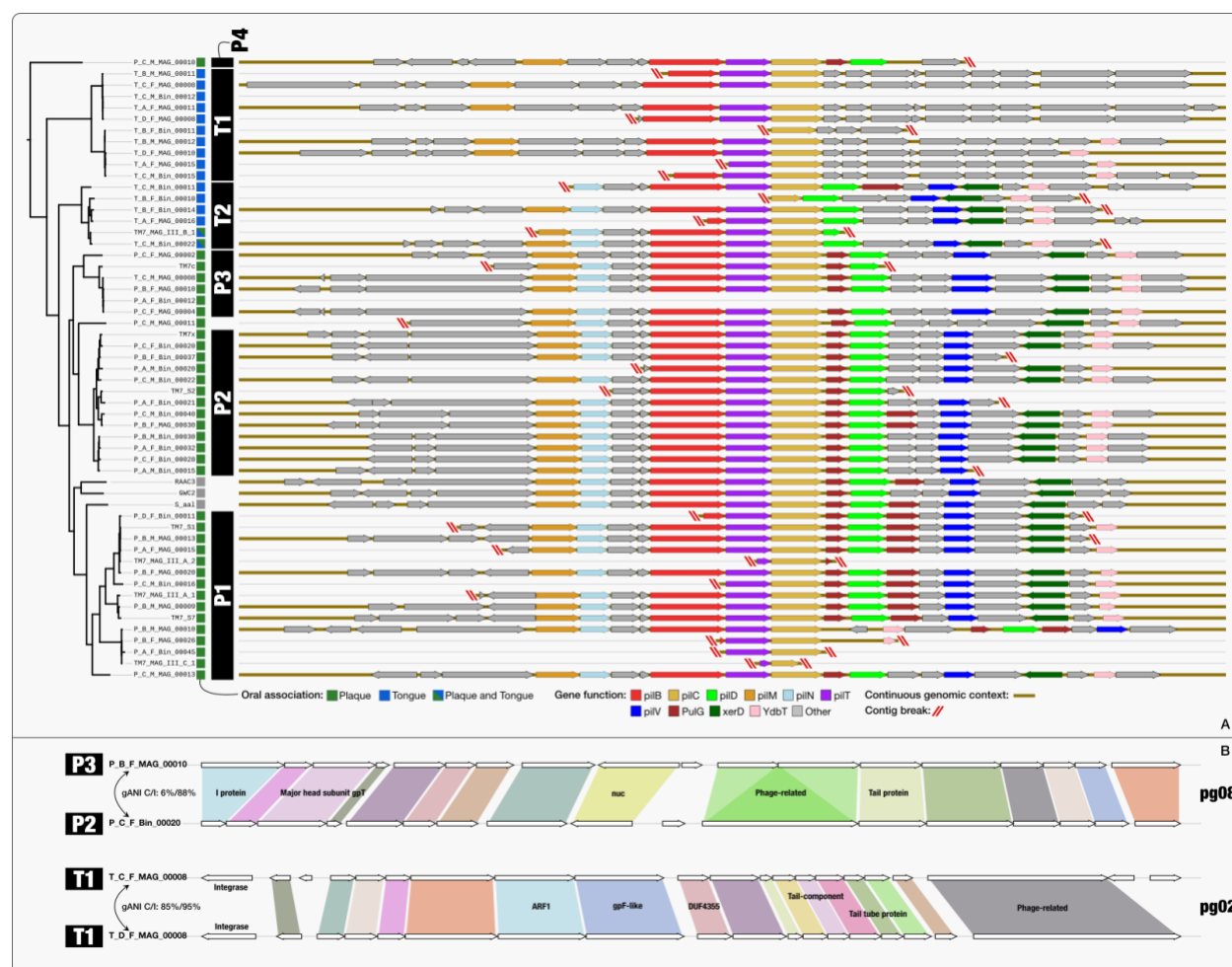


Figure 6: A) Type IV pilus operon is highly conserved in TM7 genomes, but missing many components in genomes of the tongue-associated clade T1. Type IV pili operons from 52 of the 55 TM7 that included pilC are aligned according to pilC (yellow). Genomes are organized according to their phylogenetic organization shown in Figure 5. The top 10 functions identified in these operons appear with color filling, while the rest of the functions appear in grey. Contig breaks are marked with red lines for contigs that include less than 9 genes either upstream or downstream from pilC. **B) Some phage groups span phylogenetic clades, while other phage groups associate with specific clades.** At the top of the panel the two prophages of phage group pg08 are compared and on the bottom of the panel the two prophages of the phage group pg02 are compared. White arrows signify genes as identified by Prodigal. Homologous genes, identified as belonging to the same gene-cluster, are connected by colored areas. A function name assigned by KEGG, COG or Pfam functional annotation source appears for genes for which it was available. On the left the phylogenetic clade of the TM7 host of each prophage is listed next to the host genome name. The genome-wide average nucleotide identity (gANI) appears for each pair of the host genomes, where C/I stands for alignment coverage / alignment identity.

Mobile elements and prophages in TM7 genomes

Little evidence for phage association with members of the CPR has been found so far (Chen et al., 2019a). Dudek et al. recovered a phage associated with a TM7 genome from a dolphin plaque metagenome (Dudek et al., 2017) and Paez-Espino et al. identified phages with a predicted SR1 host (Paez-Espino et al., 2016) in human oral metagenomes. A smaller genome size has been shown to correlate with the lack of lysogenic phages (Touchon, Bernheim & Rocha, 2016), and a lack of prophages in CPR genomes would fit this trend. To evaluate whether oral TM7 were indeed devoid of integrated prophages, we used an automatic approach based on VirSorter (Roux et al., 2015) and the recently published “inovirus detector” (Roux et al., 2019), along with a manual approach (see Supplementary Information), to identify 9 “phage groups” each composed of closely related prophages that were recovered from multiple TM7 genomes spanning all oral clades (Supplementary Table 8g). We did not identify any prophages in the three environmental genomes. Phage groups generally associated with closely related hosts but were not restricted to hosts of the same TM7 species, or even the same oral clades (Figure 6b, Supplementary Table 8g). A BLAST search of prophage nucleotide sequences against the NCBI’s nr nucleotide collection returned no significant hits, confirming the novelty of these phage sequences. Using CRISPRCasFinder (Couvin et al., 2018) we identified CRISPR spacers targeting prophages of two “phage groups” in closely related hosts, validating the association of these prophages with their corresponding TM7 hosts. We identified CRISPR spacers and CRISPR related proteins in genomes representing clades P1, P2, P3, P4, and T2, but not in T1 nor in the three environmental genomes. The lack of CRISPR systems in the environmental TM7, despite their close affiliation with plaque TM7, would be consistent with a recent acquisition of these systems by oral clades. To investigate this hypothesis, we BLASTed cas9 proteins from 6 genomes representing all 5 CRISPR-containing clades, and found that these best matched cas9 protein from a variety of oral TM7 and Firmicutes, but no environmental TM7 nor any other CPR (Supplementary Table 8p). While some TM7 clades appear to lack CRISPR systems, we identified restriction modification (RM) systems in genomes representing all oral clades, including clade T1, as well as in the environmental genomes

GWC2 and RAAC3 (Supplementary Table 8a). These RM systems could serve as alternative measures against foreign DNA for TM7 that lack CRISPR systems. Overall our data show that prophages are common amongst oral TM7, and appear to be a unique
 665 feature of oral TM7, while absent from environmental TM7. In addition, CRISPR systems appear to be common amongst specific clades of oral TM7, but not a common feature of all TM7. While additional analyses that include a larger collection of environmental genomes will be required to verify this observation, a specific association of prophages with host-associated TM7 suggests that prophages may have played a role in the
 670 adaptation of TM7 to the host environment, perhaps by facilitating horizontal gene transfer.

In search of other mobile genetic elements, we identified transposases in 18 TM7 genomes representing all oral clades and environmental genomes (Supplementary Table 8n). The varying location of the highly conserved transposases we identified in genomes
 675 of subclade T1_a suggests recent mobility, and that at least some of these elements are indeed active transposons (Supplementary Tables 8a,o). BLAST search of genes annotated as transposases revealed that while the majority appear to be strongly associated with members of the CPR, two transposases had more close hits to those from non-CPR bacteria.

680 Additional members of the CPR are prevalent in the oral cavity, including a tongue-associated SR1

In addition to TM7, other members of the CPR have been commonly found in the human oral cavity, specifically members of the candidate phyla SR1 and GN02 (Camanocha & Dewhirst). Using full length 16S rRNA sequences from clone libraries, Camanocha and
 685 Dewhirst identified three distinct oral taxa within SR1 (HMT-345, HMT-874, and HMT-875) and three within GN02 (HMT-871, HMT-872, and HMT-873). Genomes have been previously published for all of these taxa except SR1 HMT-875 (Camanocha & Dewhirst; Campbell et al., 2013). While none of the GN02 and SR1 MAGs in our collection included 16S rRNA, which would allow a direct match to the Human Microbial Taxon (HMT)
 690 designation, using a pangenomic analysis along with ANI statistics we were able to match

MAGs to genomes representing HMT-871, HMT-873, HMT-345, and HMT-874 (Figures S6a, S7a, Supplementary Tables 9a-h). Only a single tongue-associated SR1 (T_B_F_MAG_00004) did not match any previously published genome. A recent study presented the successful isolation of an SR1 HMT-875, but a genome has not been sequenced (Cross et al., 2019).

To investigate the niche association of these CPR genomes, we characterized their distribution across HMP metagenomes through read recruitment. While SR1 HMT-874 and HMT-345 were enriched in plaque samples, T_B_F_MAG_00004 was highly enriched in tongue samples (detected in 37% of tongue and 9% of plaque metagenomes), and it recruited up to 2.09% of all metagenomic reads from tongue samples (Figures S6b-c, Supplementary Tables 9l-n). Oral GN02 were all associated with plaque, and nearly absent from tongue samples (Figures S6b-c, Supplementary Tables 9i-l). Our ANI analysis suggests that HMT-871 and HMT-872 represent the same genus as genomes from both of these lineages match with ANI>85% (alignment coverage>30%), while HMT-873 represents a separate genus and likely a separate family or order, as suggested by Camanocha & Dewhirst (Camanocha & Dewhirst) (Supplementary Tables 9e-f). Overall our GN02 and SR1 MAGs extend the collection of genomes available for these understudied members of the oral microbiome, and our analysis demonstrates their niche partitioning and reveals the prevalence of a tongue-associated SR1.

Novel non-CPR lineages represent prevalent members of the oral microbiome

Our collection included 34 MAGs that based on phylogenomics and BLAST sequence search represent 11 lineages with no representation at NCBI (from here on referred to as “novel MAGs”). These appear to include two unnamed species of the genus *Prevotella*, single unnamed species of the genera *Mogibacterium*, *Propionibacterium*, *Leptotrichia*, and *Capnocytophaga* each, as well as an unnamed genus in the family *Flavobacteriaceae*, an unnamed family within the class *Clostridia*, and unnamed families (and potentially unnamed orders) within the classes *Bacteroidia* and *Mollicutes* (Figure 1, Supplementary Table 10a-d, Supplementary Information file). Populations represented

by these novel MAGs were absent from skin and gut samples, and of our 790 MAGs, we found only two that were consistently detected in gut samples. Both of these MAGs belong to the species *Dialister invisus*, which were previously found to be the only abundant gut-associated microbes that were detected with considerable abundance in the oral cavity (Franzosa et al. 2014, Eren et al. 2014).

The oral microbiome is highly represented in genomic databases (Vartoukian et al., 2016; Nayfach et al., 2016), hence we next sought to check if the lack of genomic representation for these novel MAGs is due to low prevalence. We mapped short reads from the HMP metagenomes to these MAGs to estimate their prevalence and abundance across oral sites. Overall, the organisms represented by these novel genomes presented strong tropism for either tongue or plaque, with the exception of three populations that appear to consistently recruit reads from both plaque and tongue samples, represented by the Flavobacteriaceae MAGs, T_A_M_MAG_00009 (Clostridiales), and three Capnocytophaga MAGs (Figure S7a). While we found some populations to be rare, which could explain their lack of genomic representation in databases, other populations were extremely prevalent (Figures S7a-c, Supplementary Table 10e-h). In addition to their high prevalence, some of these novel MAGs were highly abundant. P_B_M_MAG_00008 (Capnocytophaga) recruited on average 1% of the reads of plaque samples and two of the Propionibacterium MAGs recruited up to 18% of the reads of a single plaque metagenome, and on average 0.7% for plaque metagenomes (Supplementary Table 10h).

The most prevalent novel MAGs were five closely related MAGs of the family Flavobacteriaceae, which we detected in approximately 98.5% and 80% of HMP plaque and tongue samples, respectively, and reached high relative abundance, recruiting up to 2.98% of the reads of a single metagenome, and on average 0.19%, 0.62% of tongue, and plaque samples respectively (Supplementary Tables 10e,g). ANI comparison of these MAGs to each other and to representatives of all Flavobacteriaceae species on RefSeq suggested they represent a single new species in an unnamed genus, as within group ANI was >93.8% (with >80% alignment coverage), while they had no significant alignment with any other Flavobacteriaceae genome (Supplementary Table 10i-j). A

phylogenomic analysis placed these MAGs in a subgroup of Flavobacteriaceae together with *Cloacibacterium*, *Chryseobacterium*, *Bergeyella*, *Riemerella*, *Cruoricaptor*, *Elithabetkingia*, and *Soonwooa* (Figure S8). While all five Flavobacteriaceae MAGs had high sequence similarity, both ANI results and the phylogenetic analysis clustered these genomes according to the site of recovery, suggesting the existence of a plaque and tongue-specific sub-population. Three of our Flavobacteriaceae genomes were highly complete according to estimation by SCGs and were of length 1.7-1.8Mbp, considerably shorter than other Flavobacteriaceae genomes, as well as other commonly found oral microbes. The short length of these genomes as compared to other Flavobacteriaceae suggests a recent genomic reduction and possibly strong host-association. A strong host-association could lead to many auxotrophies and could explain why this species has never been isolated despite being an abundant and ubiquitous member of the oral microbiome. The recovery of novel genomes for these prevalent members of the oral microbiome could help shed light on their role and could assist future cultivation efforts.

Conclusions

Using genome resolved metagenomics, we have recovered much of the known diversity of the human oral cavity using samples from only 7 individuals, providing genomes for prevalent, yet uncultivated members of the microbiome, and highlighting phylogenetic and functional markers of niche partitioning of the cryptic candidate phylum TM7. Our findings group TM7 from the supragingival plaque with environmental TM7, both functionally and phylogenetically, while tongue-associated TM7 group together with lineages associated with animal gut, suggesting that at least for TM7, the supragingival plaque resembles non-host environments, while the tongue and gut TM7s are more strongly shaped by the host. Drivers of differentiation between the various microbial niches within the oral cavity are largely unknown, and could be revealed by applying similar approaches to study additional members of the oral microbiome.

Material and methods

Sampling. We recruited human subjects and collected samples according to protocol #15-247 as approved by New England IRB (Newton, Massachusetts, USA). Seven healthy subjects, 4 female, 3 male, in the age range 21 to 55 years, contributed to the study. The seven subjects included three male-female married couples and a single individual; the couples had been married for 10 to 22 years at the time of sampling. All participants gave informed consent prior to sampling. Subjects were sampled 5 to 6 times over a 5- to 9-day period. We instructed subjects to refrain from using mouthwash during the sampling period, to refrain from eating and from oral hygiene on each morning of sampling until after the samples had been collected, and to refrain from flossing teeth the evening before sampling. Eating, drinking, and oral hygiene were otherwise permitted as was customary for the subject. Subjects collected tongue material by passing a ridged plastic tongue scraper (BreathRx Gentle Tongue Scraper, Discus Dental, Culver City, CA) with gentle pressure over the surface of the tongue from back to front; three subjects (T-A-F, T-A-M, and T-D-F) also collected tongue material by swabbing the tongue dorsum with a sterile Catch-All™ Specimen Collection Swab (Epicentre Biotechnologies, Madison, WI). Subjects collected dental plaque samples using multiple toothpicks to collect plaque from the gingival margin, from the surface of teeth on the buccal (cheek) side, and from between the teeth. Subjects transferred collected material directly into the bead tube of a PowerSoil DNA Isolation kit (Qiagen). We extracted DNA following the manufacturer's protocol, processed the extracted DNA with the NEBNext Microbiome DNA Enrichment Kit (New England Biolabs, Ipswich, MA), and quantified the resulting DNA using Picogreen (Invitrogen).

Shotgun metagenomic library preparation and sequencing. For short-read library preparation, we sheared DNA using a Covaris acoustic platform. We visualized amplified libraries on an Agilent Bioanalyzer DNA1000 chip, pooled at equimolar concentrations based on these results, and size-selected to an insert size of 600 bp using a PippinPrep 2% cassette (Sage Biosciences). To quantify library pools we used a Kapa Biosystems

qPCR library quantification protocol and then performed sequencing on the Illumina NextSeq in a 2x150 paired-end sequencing run using dedicated read indexing.

Metagenomic assembly and processing of contigs. We used illumina-utils (Eren et al., 2013) for quality filtering of short reads from the 71 metagenomes with the 'iu-filter-quality-minoche' program using default parameters, which removes noisy reads using the method described in (Minoche, Dohm & Himmelbauer, 2011). We then used MEGAHIT (Li et al., 2015) v1.0.6 to co-assemble the set of all quality filtered metagenomes originating from one oral site (either plaque or tongue) of one donor, for a total of 14 co-assemblies. We used anvi-display-contigs-stats to get a summary of contigs statistics for each co-assembly. To process FASTA files for each of the 14 assemblies we used the contigs workflow implemented in anvi'o (Eren et al., 2015), v5.5.1, which (1) generated an anvi'o contigs database, (2) identified open reading frames using Prodigal (Hyatt et al., 2010) v2.60, (3) predicted gene-level taxonomy using Centrifuge (Kim et al., 2016), (4) annotated genes with functions using the NCBI's Clusters of Orthologous Groups (COG) (Tatusov et al., 2003), and (5) identified single-copy core genes using HMMER (Eddy, 2011) v3.2.1 and a collection of built-in HMM profiles for bacteria and archaea.

Metagenomic read recruitment, and initial automatic binning. In our metagenomic workflow we used Bowtie2 v2.3.4.3 (Langmead & Salzberg, 2012) to recruit short reads from the set of metagenomes used for co-assembly to the assembly product; we used samtools (Li et al., 2009) to sort the output sam files into bam files; and we used anvi'o to profile the bam files and compute coverage and detection statistics, and merge the profiles of each metagenomic sets. We then used CONCOCT (Aineberg et al., 2013) to create preliminary clusters of contigs. In short, CONCOCT uses differential coverage and sequence composition of contigs to bin contigs together. For each co-assembly, we constrained CONCOCT to generate 10 superclusters to maximize explained patterns while minimizing fragmentation error (where contigs that belong to the same population distribute into more than one bin). We then used the anvi'o interactive interface to manually refine the superclusters generated by CONCOCT using the method described below. Finally, we retained all MAGs of length greater than 0.5Mbp, and redundancy in SCGs below 10% for the rest of the analysis.

Manual bin refinement. We used the anvi'o interactive interface to refine our MAGs, as well as TM7 genomes we downloaded from the IMG, which as previously reported (McLean et al., 2018a), include contamination. In our refinement approach, we utilized the different clustering organizations available on the anvi'o interactive interface, which rely on sequence composition and differential coverage across multiple metagenomes. In addition, to assist our refinement we utilized taxonomic assignments of contigs assigned based on Centrifuge annotation of genes. In cases in which we could not confidently distinguish contamination based on the clustering organizations, we used BLAST of specific sequences to assist us in making refinement decisions. We performed two to three rounds of refinement per MAG: (1) Refinement using the coverage information in the 4-6 samples used to assemble each MAG, (2) Refinement of 63 MAGs which we identified as contaminated based on their coverage across our full collection of 71 metagenomes, and then used this coverage profile for refinement, and (3) Refinement of CPR MAGs and novel MAGs based on their coverage patterns in the HMP samples. Some contamination may remain despite these efforts. Refinement of TM7 genomes downloaded from IMG was done using coverage of their contigs across the HMP samples.

Naming scheme of MAGs. We named final MAGs according to the following scheme: names of the final MAGs included the prefix "ORAL", followed by a single letter to specify the type of samples used for the assembly of the MAG ("P" or "T" for plaque or tongue), followed by the ID of the individual (for example "C_M", which stands for "couple 'C', male"), followed by either "Bin" or "MAG" if the MAG had completion below or above 70% as estimated using the Campbell et al. collection of single-copy core genes (SCGs) (Campbell et al., 2013), and followed by a number, where for each co-assembly the MAGs had a series of numbers from "00001" to the maximum number of MAGs that were retained from that co-assembly.

Removing redundancy and analysis of the non-redundant collection of MAGs. In order to identify near-identical MAGs, we used NUCmer (Delcher et al., 2002) to calculate the average nucleotide identity (ANI) between each pair of MAGs that were estimated by CheckM to belong to the same phylum. To assign phylum affiliations to MAGs that had

no phylum designation from CheckM we used phylogenomics (see below) and complimented phylogenomics with BLAST of protein sequences against the NCBI's non-redundant database. We determined that a pair of MAGs are redundant if their ANI was 99.8% with the alignment length covering at least 50% of the shorter of the two genomes. For each group of redundant genomes, we chose the genome with the highest 'completion minus redundancy' as the representative of the group, where completion and redundancy were calculated by anvi'o based on single-copy core genes. If multiple redundant genomes had the same 'completion minus redundancy' then we selected the longest genome as the representative genome. We merged the sequences of the collection of non-redundant bins into one FASTA file, and processed this FASTA file using the anvi'o contigs workflow as mentioned above. We then also used this FASTA file to recruit reads from all 71 metagenomes, and used the anvi'o metagenomics workflow as mentioned above to generate a merged profile database. To rapidly inspect contamination that may be missed by SCGs, we generated images that visually describe the coverage of each contig in each MAG across metagenomes. For this, we used (1) anvi-split to split each MAG into their stand-alone database files, (2) anvi-interactive with the flag --export-svg, which stores the display as an SVG without user interaction, and (3) Inkscape from terminal for SVG to PNG conversion for quick screening. We used these images to identify MAGs that required additional refinement. For the novel Flavobacteriaceae population genome ORAL_T-B-M_MAG_00001, we followed the previously explained scaffold extension and gap closing strategies (Chen et al., 2019b) to reduce the number of contigs from 48 to 8.

Sequence searches in public databases. We used the NCBI nucleotide collection to search for nucleotide sequences, and the NCBI non-redundant protein sequences database to search for protein sequences. For 16S rRNA sequences, we used the 16S rRNA RefSeq Version 15.2 (starts at position 28) through the online search tool of HOMD (<http://www.homd.org>) with default settings.

Read recruitment from public metagenomes. To recruit reads from Human Microbiome Project (HMP) (Human Microbiome Project Consortium, 2012) oral metagenomes we used the program 'anvi-run-workflow' with '--workflow metagenomics',

which uses Snakemake (Köster & Rahmann, 2012) to execute the steps described above for our metagenomic read recruitment analysis. We used the same approach to also recruit reads from previously published metagenomes from periodontitis patients (Califf et al., 2017) to the TM7 pangenome.

Quantifying human contamination in metagenomes. We ran the aforementioned metagenomics workflow using anvi-run-workflow and used the human genome build 38 (GRCh38) from NCBI to quantify the number of reads matching the human genome in each sample. We estimated the number of reads that originate from microbes (or “non-human” reads) in each sample as the total number of reads minus the number of reads that mapped to the human genome.

Relative abundance estimations of MAGs. For each MAG we used the number of reads that mapped to it, divided by the total number of non-human reads as the unnormalized abundance. All unmapped reads were counted as an UNKNOWN bin. In order to account for different genome lengths, which is expected to impact the number of reads expected from each population at a given true abundance, we divided each normalized abundance by the genome length. Since the genome length is unknown for the UNKNOWN bin, as it represents an agglomeration of whole genomes and the portion of genomes that we did not recover, we used an arbitrary choice of 2Mbp as the normalization factor. The choice of this arbitrary factor changes the overall estimation of the portion of unknown reads, but not the observed trends.

Taxonomic profiles of metagenomes based on short reads. We used KrakenUniq (Breitwieser & Salzberg, 2018) to generate taxonomic profiles for all metagenomes. Briefly, KrakenUniq uses counts of unique k-mers to estimate the relative abundance of taxa in a sample, based on short-reads.

Phylogenomic analyses. For phylogenomic analyses we used our collection of 37 ribosomal proteins that occurred both in bacterial (Campbell et al., 2013) and archaeal (Rinke et al., 2013) single-copy core gene collections: Ribosom_S12_S23, Ribosomal_L1, Ribosomal_L10, Ribosomal_L11, Ribosomal_L11_N, Ribosomal_L13, Ribosomal_L14, Ribosomal_L16, Ribosomal_L18e, Ribosomal_L18p, Ribosomal_L19,

Ribosomal_L2, Ribosomal_L21p, Ribosomal_L22, Ribosomal_L23, Ribosomal_L29,
 925 Ribosomal_L2_C, Ribosomal_L3, Ribosomal_L32p, Ribosomal_L4, Ribosomal_L5,
 Ribosomal_L5_C, Ribosomal_L6, Ribosomal_S11, Ribosomal_S13, Ribosomal_S15,
 Ribosomal_S17, Ribosomal_S19, Ribosomal_S2, Ribosomal_S3_C, Ribosomal_S4,
 Ribosomal_S5, Ribosomal_S5_C, Ribosomal_S6, Ribosomal_S7, Ribosomal_S8,
 Ribosomal_S9. To compute phylogenetic trees we used the program 'anvi-run-workflow'
 930 with '--workflow phylogenomics' parameter, which runs 'anvi-get-sequences-for-hmm-
 hits' using parameters (1) '--align-with famsa' to perform alignment of protein sequences
 using FAMSA (Deorowicz, Debudaj-Grabysz & Gudyś, 2016), (2) '--concatenate-genes'
 to concatenate separately aligned and concatenated ribosomal proteins, (3) '--return-
 best-hit' to return only the most significant hit when a single HMM profile had multiple hits
 935 in one genome, (4) '--get-aa-sequences' to output amino-acid sequence, and (4) '--hmm-
 sources Campbell_et_al' to use the Campbell_et_al HMM source (Campbell et al., 2013)
 to search for genes. For Figure 1 we also included the parameter '--max-num-genes-
 missing-from-bin 19' to only include genomes that contain at least 18 of the 37 ribosomal
 proteins. For the rest of the phylogenomics analyses we used '--min-num-bins-gene-
 940 occurs' to ensure that only ribosomal proteins that occur in at least 50% of the genomes
 are used for the analysis. We trimmed alignments using trimAl (Capella-Gutiérrez, Silla-
 Martínez & Gabaldón, 2009) with the setting '-gt 0.5' to remove all positions that were
 gaps in more than 50% of sequences, and computed maximum likelihood phylogenetic
 trees using IQ-TREE (Nguyen et al., 2015) with the 'WAG' general matrix model (Whelan
 945 & Goldman, 2001). We computed the phylogeny of CPR genomes using only 36 of the
 37, excluding Ribosomal_L32p since it was absent from all TM7 genomes. To root
 phylogenetic trees we used an outlier genome in each analysis: for Figure 1 we used a
 genome of the archeal *Methanobrevibacter oralis*, and for all other phylogenomic
 analyses we used a collection of five members of the Firmicutes: *Acidaminococcus*
 950 *intestini*, *Eubacterium rectale*, *Staphylococcus aureus*, *Streptococcus pneumoniae*,
Veillonella parvula. To remove the Firmicutes from the trees in Figure 5, S7a, S7d we
 used the python package ete3 version 3.1.1 (Huerta-Cepas, Serra & Bork, 2016).

Processing publicly available genomes. To process FASTA files, we used the program
 'anvi-run-workflow' with '--workflow contigs' parameter, which includes the steps of the

anvi'o contigs workflow as described above. To generate the data in Supplementary Table 8a, our workflow also included running 'anvi-run-pfams' to annotate functions with Pfams (El-Gebali et al. 2019). We used 'anvi-get-sequences-for-gene-calls' to get all protein sequences and used GhostKoala (<https://www.kegg.jp/ghostkoala/>) to annotate genes with KEGG functions (Kanehisa, Sato, and Morishima 2016).

Assessing the occurrence of populations in metagenomes. We used anvi-mcg-classifier with the settings '--get-samples-stats-only', '--alpha 0.1', which determines a threshold of 0.6 detection value for to determine occurrence, '--zeros-are-outliers', which considers positions with zero coverage as outlier coverage values when computing the non-outlier mean coverage. We used the anvi-mcg-classifier output to determine the occurrence of TM7 populations in our collection of 71 metagenomes. In order to account for the different number of reads per sample when comparing non-outlier mean coverage values, we normalized these values. To compute the normalization factor, we first divided the number of reads in each sample by the maximum number of reads in the biggest sample (so that the normalization factor would be ≤ 1 for all samples). We then divided the non-outlier mean coverage values in each sample by the normalization factor.

Pangenomic analyses. To compute pangenomes in our study we used the program 'anvi-run-workflow' with '--workflow pangenomics'. Anvi'o pangenomics workflow is detailed elsewhere (Delmont & Eren, 2018), but briefly the pangenomic analysis used the NCBI's BLAST (Altschul et al., 1990) to quantify similarity between each pair of genes, and the Markov Cluster algorithm (MCL) (Enright, Van Dongen & Ouzounis, 2002) (with inflation parameter of 2) to resolve clusters homologous genes. The program 'anvi-summarize' created summary tables for pangenomes and 'anvi-display-pan' provided interactive visualizations of pangenomes. To simplify visualizations of complex pangenomes we removed singleton gene clusters using the parameter '--min-occurrence 2'.

Average nucleotide identity (ANI). We used 'anvi-compute-ani' with parameters '--method ANIm' to align genomes using MUMmer (Marçais et al., 2018) and '--min-alignment-fraction 0.25' to only keep scores if the alignment fraction covers at least 25% percent of both genomes. For the ANI data presented in Figures 5 and 6, we first

computed ANI without the flag '--min-alignment-fraction' to get all alignment statistics, and then we imported ANI values only for pairs of genomes with alignment coverage of at least 25%.

Long-read sequencing, analysis, and extraction of 16S rRNA sequences.

We collected additional samples for long-read sequencing from two of the initial seven

individuals (C-M and C-F) as well as a new female participant (L). To increase microbial biomass and the likelihood of getting reads from low abundance members, we pooled four daily tongue dorsum scrapings from individuals C_F and C_M. From the individual L, we collected 13 daily tongue dorsum scrapings (Supplementary Table 1d). To obtain

high-molecular weight (HMW) DNA from these low-biomass samples, we extracted total

genomic DNA using the Qiagen Genomic Tip 20/G gravity flow columns (Qiagen, Germantown, MD) using the manufacturer's protocol. We modified the lytic enzyme

cocktail to include lysostaphin (by Sigma-Aldrich, final concentration: 24U/mL) and mutanolysin (by Sigma-Aldrich; final concentration: 0.3KU/mL), and extended each incubation step to 2-hrs to increase the DNA yield from Gram-positive bacteria. During all

laboratory steps, we sought to maximize read lengths by implementing best practices for handling HMW DNA, including (1) smooth and slow pipetting with wide bore pipette tips and when possible, (2) replacing centrifugation/vortexing with end-over-end rotations to minimize velocity gradients and avoid further shearing to DNA molecules. For library

synthesis we used the 1D Native barcoding genomic DNA protocol (Oxford Nanopore

Technologies, UK), but our procedure differed slightly as (1) we 'padded' a given sample

with linear double-stranded lambda DNA (New England Biolabs) if the sample did not

meet the manufacturer DNA input mass recommendations (1000 ng) and (2) we changed

the incubation time at the end-prep step of the library preparation 30 min at 20°C and 30

min at 65°C to minimize contamination with short reads. We used 1× Agencourt AMPure

XP beads (A63882, Beckman Coulter) for sample clean-up and concentration of pooled

barcoded samples. The incubation times for the DNA binding and elution steps were

modified to 20-min at room temperature and 20-min at 37°C, respectively. We quantified

DNA yield on a Qubit® 1.0 Fluorometer (Thermo Scientific), using the dsDNA HS (High

Sensitivity) Assay kit. Two R9.4/FLO-MIN106D MinION flow cells (Oxford Nanopore

Technologies) sequenced the resulting libraries with a starting voltage of -180 mV and

run times of 40 and 48 hr. Runs were stopped when the number of active pores fell below 10. We processed the raw sequencing data with ONT MinKNOW software (v.1.15.4-3.3.2), removed sequences with Q-score <7, and called bases using Guppy (version 3.2.1, Oxford Nanopore Technologies). Sequences then underwent demultiplexing, barcode trimming, and conversion of raw FAST5 files to FASTQ files. To filter human contamination we mapped final long-read sequences to the human genome using minimap2 (Li, 2018). We used the remaining contigs to generate anvi'o contigs databases as described above, extracted 16S rRNA gene sequences using 'anvi-get-sequences-for-hmm-hits' program with '--hmm-sources Ribosomal_RNAs' parameter, and used their HOMD matches to assign group affiliation to TM7 genomes.

Group affiliation of TM7 based on 16S rRNA gene sequences. We exported ribosomal RNA sequences from all TM7 genomes, including ones downloaded from NCBI. We then searched 16S rRNA sequences against the HOMD as explained above. For each genome, we identified the group affiliation (G-1, G-2, etc.) of the closest hit on HOMD. In addition, we searched for nanopore reads that matched to TM7 against the collection of oral TM7 genomes. We used search results to associate TM7 MAGs with a 16S rRNA group affiliation. The 16S rRNA group affiliations are summarized in Supplementary Table 7i for oral genomes, and in Supplementary Table 7e for all TM7 downloaded from NCBI.

Functional enrichment analysis. We employed a statistical approach to identify functions enriched within a phylogenomic clade. This approach fits a logistic regression (binomial GLM) to the occurrence of each gene function using clade affiliation as the explanatory variable using R (R Development Core Team, 2011). We test for equality of proportions across clade affiliation using a Rao score test, which gives a test statistic ("enrichment score") and a p-value. As this test is performed independently for each function, we computed q-values from p-values to account for multiple testing using the R package "qvalue" (Storey & Tibshirani, 2003). To apply functional enrichment analysis to our pangemome we used the program 'anvi-get-enriched-functions-per-pan-group' using COG functions across genomes and using clade affiliation as the explanatory variable. We considered a function to be enriched if the q-value was below 0.05; this controls the

expected proportion of false positives at 0.05. The URL <http://merenlab.org/p> provides details on how to use this method.

Identifying prophages in TM7 genomes. We used Virsorter (Roux et al., 2015) and the “Inovirus detector” (Roux et al., 2019) to identify contigs that include putative phage sequences. We manually inspected contigs predicted as viral, and all contigs which gene content was also consistent with a plasmid or another mobile genetic element, i.e. did not include either a viral hallmark gene or capsid-related gene(s) were excluded. We further examined all remaining contigs to verify their placement in the prospective genomes, using the data in Supplementary Table 8a, as well as BLAST searches of protein sequences (see the notes in Supplementary Table 8g for more details). We used functional annotations to identify additional contigs containing phage-related functions that were not identified by VirSorter/Inovirus detector. In addition, we identified additional phages by searching for contigs with many homologs (according to GC occurrence) to the identified phages. We repeated this process recursively and identified 11 more contigs that contain partial or complete prophages. To identify start and end positions of prophages, we relied on identifying genes that appear to be TM7 genes as per their association with GCs. When possible, we used closely related TM7 genomes that lacked the prophage genes, to identify the position of the genes flanking the prophage, and hence confirming the insertion site of the prophage.

Identifying CRISPRs. We used the web service CRISPRCasFinder (Couvin et al. 2018) to search for CRISPR spacers in the 55 TM7 genomes. Along with a summary of the results (Supplementary Table 8l), the web application allows the direct download of a FASTA file of all high confidence spacers (evidence level 3 or 4 as defined by Couvin et al).

16S rRNA gene amplicon library preparation and sequencing. We amplified the V4-V5 hypervariable regions of the bacterial SSU rRNA gene using degenerate primers. 518F (CCAGCAGCYGCGGTAAN) and 926R (CCGTCAATTCNTTTTRAGT CCGTCAATTTCTTTGAGT CCGTCTATTCCTTTGANT). Amplification was done with fusion primers containing the 16S-only sequences fused to Illumina adapters. The forward primers included a 5 nt multiplexing barcode and the reverse a 6 nt index. We generated

PCR amplicons in triplicate 33 uL reaction volumes with an amplification cocktail containing 0.67 U SuperFi Taq Polymerase (Invitrogen, Carlsbad, CA), 1X enzyme buffer (includes MgCl₂), 200 uM dNTP PurePeak DNA polymerase mix (ThermoFisher), and 0.3 uM of each primer. We added approximately 10-25 ng template DNA to each PCR and ran a no-template control for each primer pair. Amplification conditions were: initial 94C, 3 minute denaturation step; 30 cycles of 94C for 30s, 57C for 45s, and 72C for 60s; final 2 minute extension at 72C. The triplicate PCR reactions were pooled after amplification, visualized with the negative controls on a Caliper LabChipGX or Agilent TapeStation 4200, and purified using Ampure followed by PicoGreen quantitation and Ampure size selection. We used Minimum Entropy Decomposition (Eren et al., 2014b) to identify amplicon sequence variants (ASVs) across samples and determine the microbial community structure, and Global Alignment for Sequence Taxonomy (GAST) (Huse et al., 2008) to assign taxonomic affiliation to each ASV.

Statistics and visualization. We used ggplot2 version 3.2.1 to generate boxplots and barplots of abundances, as well as barplots of occurrences across metagenomes. To compare the number of reads recruited by our MAGs from our plaque and tongue metagenomes, we ran a two-sided Z-test, using the Python package statsmodels (Seabold & Perktold, 2010). We finalized figures for publication using the open-source vector graphics editor Inkscape (<https://inkscape.org/>).

Access to previously published sequences. We downloaded all oral genomes from the HOMD FTP site (ftp://ftp.homd.org/HOMD_annotated_genomes/, and ftp://ftp.homd.org/NCBI_annotated_genomes/); accession numbers are available in Supplementary Table 6b, which also includes the archeon *Methanobrevibacter oralis* that was used to root the phylogeny in Figure 1. While the TM7 genomes we downloaded from IMG had no accession numbers, refined versions of these genomes have recently been published (Cross et al., 2019). To download genomes from GenBank we used 'ncbi-genome-download' (<https://github.com/kblin/ncbi-genome-download>) and processed them with 'anvi-script-process-genbank-metadata' to generate input files for the anvi'o contigs workflow. We downloaded TM7 genomes from GenBank on 1/16/2019 (accession numbers provided in Supplementary Table 7e); GN02 and SR1 on 12/17/2018 (accession

numbers provided in Supplementary Table 9a,b); and Flavobacteriaceae on 9/20/2019. We obtained the raw metagenomes from Califf et al. 2017 directly from the authors of Califf et al. since the FASTQ files published by Califf et al. included only a single read for each pair of raw reads. Firmicutes genomes used to root the trees in Figures 2, 3, and 4 were downloaded from RefSeq (Accessions: GCF_000147095.1, GCF_000210315.1, GCF_000024945.1, GCF_000020605.1, GCF_000230275.1).

Data availability

We deposited the short-read sequencing data for amplicons and shotgun metagenomes under the NCBI BioProject PRJNA625082 (Supplementary Table lists accession numbers for each sample individually). We deposited the long-read sequencing data in the NCBI's SRA database for individuals C-F (SRR11547007), C-M (SRR11547005, SRR11547006), and L (SRR11547004). For reproducibility, reusability, and transparency, we also have made available for each individual the FASTA files for co-assembled metagenomes (doi:10.6084/m9.figshare.12217799) and anvi'o merged profile databases (doi:10.6084/m9.figshare.12217802) in addition to anvi'o split profiles for each MAG (doi:10.6084/m9.figshare.12217805) and the TM7 pangenome (doi:10.6084/m9.figshare.12217811).

Acknowledgements

We thank Gary Borisy (ORCID:0000-0002-0266-8018), Emily Fogarty (0000-0002-8957-9922), Evan Kiefl (0000-0002-6473-0921) for helpful discussions, and Elaina Graham (0000-0001-5036-1929) for her help with GhostKOALA. ADW was funded by NIH R35GM133420; SR was supported by the Office of Science of the U.S. Department of Energy under contract no. DE-AC02-05CH11231; FED was funded by NIH NIDCR DE016937 and DE024468; JMW was funded by NIH NIDCR DE022586; AME and JMW were supported by a G. Unger Vetlesen Foundation grant to the Marine Biological Laboratory. This research was funded by the Frank R. Lillie Research Innovation Award

and start-up funds from the University of Chicago to A.M.E, Gastro-Intestinal Research Foundation (GIRF), and the Mutchnik Family Fund.

Author contributions

TOD, JLMW, and AME conceived the study. AS, TOD, and AME performed the primary analysis of data. AS and AME prepared figures and tables. AS, JLMW, and AME wrote the manuscript. AS, MY, OCE, AME developed computational analysis tools. ADW developed statistical analysis tools. SR, LXC, ACS, ARW, STM, and FED analyzed data. ND and HGM performed short-read sequencing. KL performed long-read sequencing. FED, JLMW, AME supervised research. All authors read and edited the drafts of the manuscript.

Competing interests

The authors declare they have no competing interests.

References

- 1145 Aas JA, Paster BJ, Stokes LN, Olsen I, Dewhirst FE. 2005. Defining the normal bacterial flora of the oral cavity. *Journal of clinical microbiology* 43:5721–5732.
- 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. *The ISME journal* 7:1016–1025.
- 1150 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. *Nature biotechnology* 31:533–538.
- Alneberg J, Bjarnason BS, de Bruijn I, Schirmer M, Quick J, Ijaz UZ, Loman NJ, Andersson AF, Quince C. 2013. CONCOCT: Clustering cONTigs on COverage and ComposiTion. *arXiv [q-bio.GN]*.
- 1155 Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search tool. *Journal of molecular biology* 215:403–410.
- Bella J, Hindle KL, McEwan PA, Lovell SC. 2008. The leucine-rich repeat structure. *Cellular and molecular life sciences: CMLS* 65:2307–2333.
- 1160 Bor B, Bedree JK, Shi W, McLean JS, He X. 2019. Saccharibacteria (TM7) in the Human Oral Microbiome. *Journal of dental research* 98:500–509.
- Bor B, McLean JS, Foster KR, Cen L, To TT, Serrato-Guillen A, Dewhirst FE, Shi W, He X. 2018. Rapid evolution of decreased host susceptibility drives a stable relationship between ultrasmall parasite TM7x and its bacterial host. *Proceedings of the National Academy of Sciences of the United States of America* 115:12277–12282.
- 1165 Breitwieser FP, Salzberg SL. 2018. KrakenHLL: Confident and fast metagenomics classification using unique k-mer counts. *bioRxiv*:262956. DOI: 10.1101/262956.
- Brinig MM, Lepp PW, Ouverney CC, Armitage GC, Relman DA. 2003. Prevalence of bacteria of division TM7 in human subgingival plaque and their association with disease. *Applied and environmental microbiology* 69:1687–1694.
- 1170 Brown CT, Hug LA, Thomas BC, Sharon I, Castelle CJ, Singh A, Wilkins MJ, Wrighton KC, Williams KH, Banfield JF. 2015. Unusual biology across a group comprising more than 15% of domain Bacteria. *Nature* 523:208–211.
- Buist G, Steen A, Kok J, Kuipers OP. 2008. LysM, a widely distributed protein motif for binding to (peptido)glycans. *Molecular microbiology* 68:838–847.
- 1175 Califf KJ, Schwarzberg-Lipson K, Garg N, Gibbons SM, Gregory Caporaso J, Slots J, Cohen C, Dorrestein PC, Kelley ST. 2017. Multi-omics Analysis of Periodontal Pocket Microbial Communities Pre- and Posttreatment. *mSystems* 2. DOI: 10.1128/msystems.00016-17.
- 1180 Camanocha A, Dewhirst FE. 2014. Host-associated bacterial taxa from Chlorobi, Chloroflexi, GN02, Synergistetes, SR1, TM7, and WPS-2 Phyla/candidate divisions. *Journal of oral microbiology* 6. DOI: 10.3402/jom.v6.25468.
- Camanocha A, Dewhirst FE. *Host-associated bacterial taxa from Chlorobi, Chloroflexi, GN02, Synergistetes, SR1, TM7, and WPS-2 Phyla/candidate divisions. J Oral Microbiol.* 2014; 6. Epub 2014/10/16. doi: 10.3402/jom.v6.25468 PMID: 25317252.
- 1185 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. *Proceedings of the National Academy of Sciences of the United States of America* 110:5540–5545.
- Capella-Gutiérrez S, Silla-Martínez JM, Gabaldón T. 2009. trimAl: a tool for automated alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* 25:1972–1973.

- 1190 Caporaso JG, Lauber CL, Costello EK, Berg-Lyons D, Gonzalez A, Stombaugh J, Knights D, Gajer P, Ravel J, Fierer N, Gordon JI, Knight R. 2011. Moving pictures of the human microbiome. *Genome biology* 12:R50.
- Chen L-X, Al-Shayeb B, Méheust R, Li W-J, Doudna JA, Banfield JF. 2019a. Candidate Phyla Radiation Roizmanbacteria From Hot Springs Have Novel and Unexpectedly Abundant CRISPR-Cas Systems. *Frontiers in microbiology* 10:928.
- 1195 Chen L-X, Anantharaman K, Shaiber A, Murat Eren A, Banfield JF. 2019b. Accurate and Complete Genomes from Metagenomes. *bioRxiv*:808410. DOI: 10.1101/808410.
- Chen I-MA, Chu K, Palaniappan K, Pillay M, Ratner A, Huang J, Huntemann M, Varghese N, White JR, Seshadri R, Smirnova T, Kirton E, Jungbluth SP, Woyke T, Elie-Fadrosh EA, Ivanova NN, Kyrpides NC. 2019c. IMG/M v.5.0: an integrated data management and comparative analysis system for microbial genomes and microbiomes. *Nucleic acids research* 47:D666–D677.
- 1200 Chen T, Yu W-H, Izard J, Baranova OV, Lakshmanan A, Dewhirst FE. 2010. The Human Oral Microbiome Database: a web accessible resource for investigating oral microbe taxonomic and genomic information. *Database: the journal of biological databases and curation* 2010:baq013.
- 1205 Collins AJ, Murugkar PP, Dewhirst FE. 2019. Complete Genome Sequence of Strain AC001, a Novel Cultured Member of the Human Oral Microbiome from the Candidate Phylum Saccharibacteria (TM7). *Microbiology Resource Announcements* 8. DOI: 10.1128/mra.01158-19.
- Couvin D, Bernheim A, Toffano-Nioche C, Touchon M, Michalik J, Néron B, Rocha EPC, Vergnaud G, Gautheret D, Pourcel C. 2018. CRISPRCasFinder, an update of CRISPRFinder, includes a portable version, enhanced performance and integrates search for Cas proteins. *Nucleic acids research* 46:W246–W251.
- 1210 Craig L, Forest KT, Maier B. 2019. Type IV pili: dynamics, biophysics and functional consequences. *Nature reviews. Microbiology* 17:429–440.
- Cross KL, Campbell JH, Balachandran M, Campbell AG, Cooper SJ, Griffen A, Heaton M, Joshi S, Klingeman D, Leys E, Yang Z, Parks JM, Podar M. 2019. Targeted isolation and cultivation of uncultivated bacteria by reverse genomics. *Nature biotechnology* 37:1314–1321.
- 1215 Delcher AL, Phillippy A, Carlton J, Salzberg SL. 2002. Fast algorithms for large-scale genome alignment and comparison. *Nucleic acids research* 30:2478–2483.
- Delmont TO, Eren AM. 2018. Linking pangenomes and metagenomes: the Prochlorococcus metapangenome. *PeerJ* 6:e4320.
- 1220 Delmont TO, Quince C, Shaiber A, Esen ÖC, Lee ST, Rappé MS, McLellan SL, Lückner S, Eren AM. 2018. Nitrogen-fixing populations of Planctomycetes and Proteobacteria are abundant in surface ocean metagenomes. *Nature microbiology* 3:804–813.
- Deorowicz S, Debudaj-Grabysz A, Gudyś A. 2016. FAMSA: Fast and accurate multiple sequence alignment of huge protein families. *Scientific reports* 6:33964.
- 1225 Dewhirst FE, Chen T, Izard J, Paster BJ, Tanner ACR, Yu W-H, Lakshmanan A, Wade WG. 2010. The human oral microbiome. *Journal of bacteriology* 192:5002–5017.
- Ding T, Schloss PD. 2014. Dynamics and associations of microbial community types across the human body. *Nature* 509:357–360.
- 1230 Donati C, Zolfo M, Albanese D, Tin Truong D, Asnicar F, Iebba V, Cavalieri D, Jousson O, De Filippo C, Huttenhower C, Segata N. 2016. Uncovering oral Neisseria tropism and persistence using metagenomic sequencing. *Nature microbiology* 1:16070.
- Dudek NK, Sun CL, Burstein D, Kantor RS, Aliaga Goltsman DS, Bik EM, Thomas BC, Banfield JF, Relman DA. 2017. Novel Microbial Diversity and Functional Potential in the Marine Mammal Oral Microbiome. *Current biology: CB* 27:3752–3762.e6.
- 1235 Dutilh BE, Huynen MA, Bruno WJ, Snel B. 2004. The consistent phylogenetic signal in genome trees

- revealed by reducing the impact of noise. *Journal of molecular evolution* 58:527–539.
- Eddy SR. 2011. Accelerated Profile HMM Searches. *PLoS computational biology* 7:e1002195.
- Enright AJ, Van Dongen S, Ouzounis CA. 2002. An efficient algorithm for large-scale detection of protein families. *Nucleic acids research* 30:1575–1584.
- 1240 Eren AM, Borisy GG, Huse SM, Mark Welch JL. 2014a. Oligotyping analysis of the human oral microbiome. *Proceedings of the National Academy of Sciences of the United States of America* 111:E2875–84.
- Eren AM, Esen ÖC, Quince C, Vineis JH, Morrison HG, Sogin ML, Delmont TO. 2015. Anvi'o: an advanced analysis and visualization platform for 'omics data. *PeerJ* 3:e1319.
- 1245 Eren AM, Morrison HG, Lescault PJ, Reveillaud J, Vineis JH, Sogin ML. 2014b. Minimum entropy decomposition: Unsupervised oligotyping for sensitive partitioning of high-throughput marker gene sequences. *The ISME journal* 9:968–979.
- Eren AM, Vineis JH, Morrison HG, Sogin ML. 2013. A filtering method to generate high quality short reads using illumina paired-end technology. *PloS one* 8:e66643.
- 1250 Escapa IF, Chen T, Huang Y, Gajare P, Dewhirst FE, Lemon KP. 2018. New Insights into Human Nostril Microbiome from the Expanded Human Oral Microbiome Database (eHOMD): a Resource for the Microbiome of the Human Aerodigestive Tract. *mSystems* 3. DOI: 10.1128/mSystems.00187-18.
- Espinoza JL, Harkins DM, Torralba M, Gomez A, Highlander SK, Jones MB, Leong P, Saffery R, Bockmann M, Kuelbs C, Inman JM, Hughes T, Craig JM, Nelson KE, Dupont CL. 2018.
- 1255 Supragingival Plaque Microbiome Ecology and Functional Potential in the Context of Health and Disease. *mBio* 9. DOI: 10.1128/mBio.01631-18.
- Ferretti P, Pasolli E, Tett A, Asnicar F, Gorfer V, Fedi S, Armanini F, Truong DT, Manara S, Zolfo M, Beghini F, Bertorelli R, De Sanctis V, Bariletti I, Canto R, Clementi R, Cologna M, Crifò T, Cusumano G, Gottardi S, Innamorati C, Masè C, Postai D, Savoi D, Duranti S, Lugli GA, Mancabelli L, Turrone F, Ferrario C, Milani C, Mangifesta M, Anzalone R, Viappiani A, Yassour M, Vlamakis H, Xavier R, Collado CM, Koren O, Tateo S, Soffiati M, Pedrotti A, Ventura M, Huttenhower C, Bork P, Segata N. 2018. Mother-to-Infant Microbial Transmission from Different Body Sites Shapes the Developing Infant Gut Microbiome. *Cell host & microbe* 24:133–145.e5.
- 1260 German RZ, Palmer JB. 2006. Anatomy and development of oral cavity and pharynx. *GI Motility online*. DOI: 10.1038/gimo5.
- 1265 Gibbons RJ, Houte JV. 1975. Bacterial adherence in oral microbial ecology. *Annual review of microbiology* 29:19–44.
- Hall MW, Singh N, Ng KF, Lam DK, Goldberg MB, Tenenbaum HC, Neufeld JD, G Beiko R, Senadheera DB. 2017. Inter-personal diversity and temporal dynamics of dental, tongue, and salivary microbiota in the healthy oral cavity. *NPJ biofilms and microbiomes* 3:2.
- 1270 He X, McLean JS, Edlund A, Yooseph S, Hall AP, Liu S-Y, Dorrestein PC, Esquenazi E, Hunter RC, Cheng G, Nelson KE, Lux R, Shi W. 2015. Cultivation of a human-associated TM7 phylotype reveals a reduced genome and epibiotic parasitic lifestyle. *Proceedings of the National Academy of Sciences of the United States of America* 112:244–249.
- 1275 Huerta-Cepas J, Serra F, Bork P. 2016. ETE 3: Reconstruction, Analysis, and Visualization of Phylogenomic Data. *Molecular biology and evolution* 33:1635–1638.
- Hug LA, Baker BJ, Anantharaman K, Brown CT, Probst AJ, Castelle CJ, Butterfield CN, Hernsdorf AW, Amano Y, Ise K, Suzuki Y, Dudek N, Relman DA, Finstad KM, Amundson R, Thomas BC, Banfield JF. 2016. A new view of the tree of life. *Nature Microbiology* 1:16048.
- 1280 Human Microbiome Project Consortium. 2012. Structure, function and diversity of the healthy human microbiome. *Nature* 486:207–214.
- Huse SM, Dethlefsen L, Huber JA, Mark Welch D, Welch DM, Relman DA, Sogin ML. 2008. Exploring

- microbial diversity and taxonomy using SSU rRNA hypervariable tag sequencing. *PLoS genetics* 4:e1000255.
- 1285 Hyatt D, Chen G-L, Locascio PF, Land ML, Larimer FW, Hauser LJ. 2010. Prodigal: prokaryotic gene recognition and translation initiation site identification. *BMC bioinformatics* 11:119.
- Ishiwa A, Komano T. 2003. Thin pilus PilV adhesins of plasmid R64 recognize specific structures of the lipopolysaccharide molecules of recipient cells. *Journal of bacteriology* 185:5192–5199.
- 1290 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.
- Kim D, Song L, Breitwieser FP, Salzberg SL. 2016. Centrifuge: rapid and sensitive classification of metagenomic sequences. *Genome research* 26:1721–1729.
- 1295 Köster J, Rahmann S. 2012. Snakemake--a scalable bioinformatics workflow engine. *Bioinformatics* 28:2520–2522.
- Lamont RJ, Koo H, Hajishengallis G. 2018. The oral microbiota: dynamic communities and host interactions. *Nature reviews. Microbiology* 16:745–759.
- Lane N. 2015. The unseen world: reflections on Leeuwenhoek (1677) “Concerning little animals.” *Philosophical transactions of the Royal Society of London. Series B, Biological sciences* 370:20140344.
- 1300 Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nature methods* 9:357–359.
- Li H. 2018. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* 34:3094–3100.
- Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R, 1000 Genome Project Data Processing Subgroup. 2009. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* 25:2078–2079.
- 1305 Li D, Liu C-M, Luo R, Sadakane K, Lam T-W. 2015. MEGAHIT: an ultra-fast single-node solution for large and complex metagenomics assembly via succinct de Bruijn graph. *Bioinformatics* 31:1674–1676.
- Lloyd-Price J, Mahurkar A, Rahnavard G, Crabtree J, Orvis J, Hall AB, Brady A, Creasy HH, McCracken C, Giglio MG, McDonald D, Franzosa EA, Knight R, White O, Huttenhower C. 2017. Strains, functions and dynamics in the expanded Human Microbiome Project. *Nature* 550:61–66.
- 1310 Mager DL, Ximenez-Fyvie LA, Haffajee AD, Socransky SS. 2003. Distribution of selected bacterial species on intraoral surfaces. *Journal of clinical periodontology* 30:644–654.
- Marçais G, Delcher AL, Phillippy AM, Coston R, Salzberg SL, Zimin A. 2018. MUMmer4: A fast and versatile genome alignment system. *PLoS computational biology* 14:e1005944.
- 1315 Marcy Y, Ouverney C, Bik EM, Lösekann T, Ivanova N, Martin HG, Szeto E, Platt D, Hugenholtz P, Relman DA, Quake SR. 2007. Dissecting biological “dark matter” with single-cell genetic analysis of rare and uncultivated TM7 microbes from the human mouth. *Proceedings of the National Academy of Sciences of the United States of America* 104:11889–11894.
- 1320 Mark Welch JL, Dewhirst FE, Borisy GG. 2019. Biogeography of the Oral Microbiome: The Site-Specialist Hypothesis. *Annual review of microbiology* 73:335–358.
- Mark Welch JL, Rossetti BJ, Rieken CW, Dewhirst FE, Borisy GG. 2016. Biogeography of a human oral microbiome at the micron scale. *Proceedings of the National Academy of Sciences of the United States of America* 113:E791–800.
- 1325 Mark Welch JL, Utter DR, Rossetti BJ, Mark Welch DB, Eren AM, Borisy GG. 2014. Dynamics of tongue microbial communities with single-nucleotide resolution using oligotyping. *Frontiers in microbiology* 5:568.
- McLean JS, Bor B, To TT, Liu Q, Kearns KA, Solden LM. 2018a. Evidence of independent acquisition and adaption of ultra-small bacteria to human hosts across the highly diverse yet reduced genomes of the

phylum *bioRxiv*.

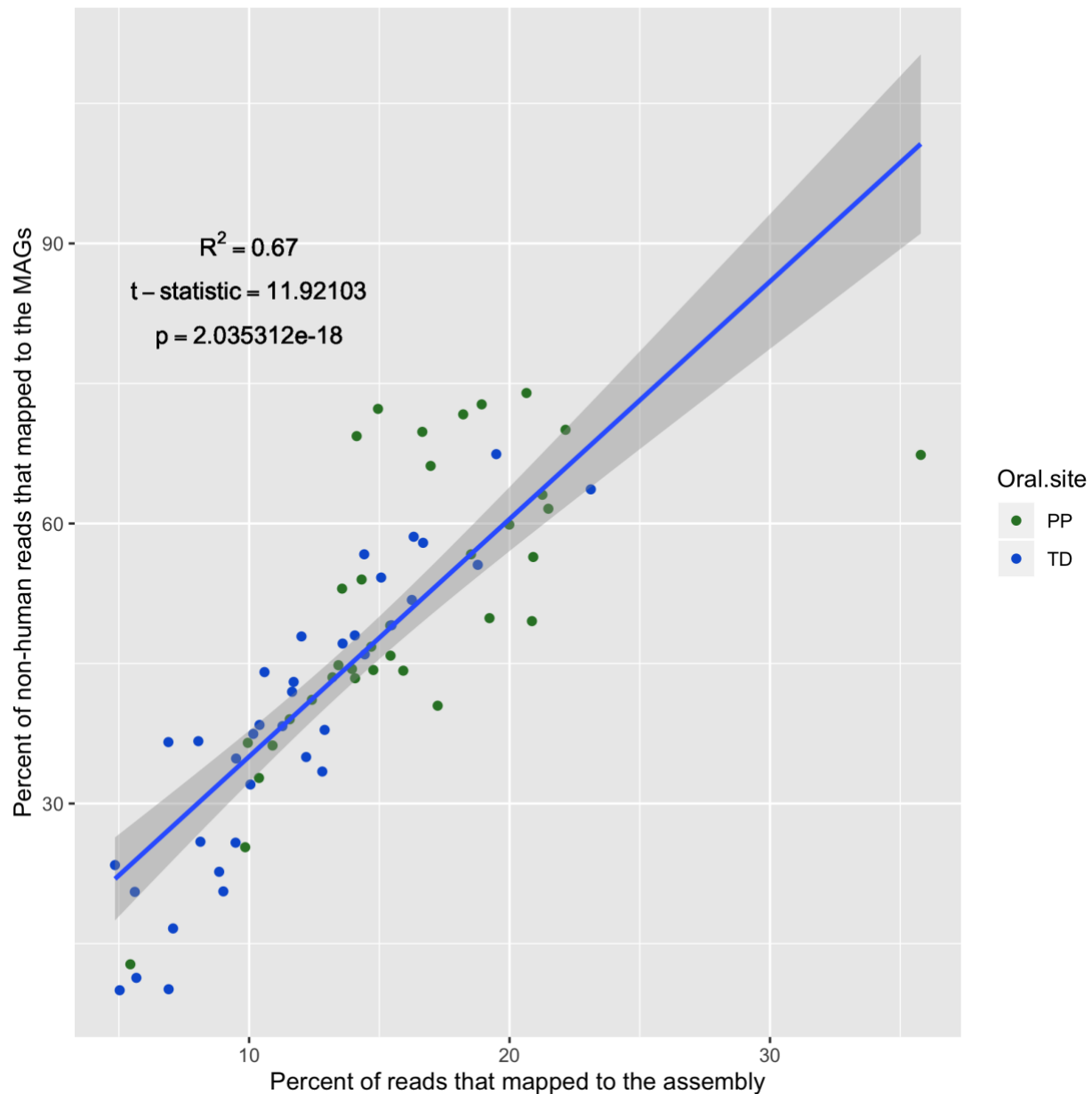
- 1330 McLean JS, Bor B, To TT, Liu Q, Kearns KA, Solden LM, Wrighton KC, He X, Shi W. 2018b. Evidence of independent acquisition and adaptation of ultra-small bacteria to human hosts across the highly diverse yet reduced genomes of the phylum Saccharibacteria. DOI: 10.1101/258137.
- Méheust R, Burstein D, Castelle CJ, Banfield JF. 2019. The distinction of CPR bacteria from other bacteria based on protein family content. *Nature communications* 10:4173.
- 1335 Minoche AE, Dohm JC, Himmelbauer H. 2011. Evaluation of genomic high-throughput sequencing data generated on Illumina HiSeq and genome analyzer systems. *Genome biology* 12:R112.
- Moutsopoulos NM, Konkel JE. 2018. Tissue-Specific Immunity at the Oral Mucosal Barrier. *Trends in immunology* 39:276–287.
- 1340 Murugkar PP, Collins AJ, Dewhirst FE. 2019. Complete Genome Sequence of Strain PM004, a Novel Cultured Member of the Human Oral Microbiome from the Candidate Phylum Saccharibacteria (TM7). *Microbiology resource announcements* 8. DOI: 10.1128/MRA.01159-19.
- Nayfach S, Rodriguez-Mueller B, Garud N, Pollard KS. 2016. An integrated metagenomics pipeline for strain profiling reveals novel patterns of bacterial transmission and biogeography. *Genome research* 26:1612–1625.
- 1345 Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ. 2015. IQ-TREE: a fast and effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Molecular biology and evolution* 32:268–274.
- Paez-Espino D, Elie-Fadrosh EA, Pavlopoulos GA, Thomas AD, Huntemann M, Mikhailova N, Rubin E, Ivanova NN, Kyrpides NC. 2016. Uncovering Earth's virome. *Nature* 536:425–430.
- 1350 Parks DH, Rinke C, Chuvochina M, Chaumeil P-A, Woodcroft BJ, Evans PN, Hugenholtz P, Tyson GW. 2017. Recovery of nearly 8,000 metagenome-assembled genomes substantially expands the tree of life. *Nature microbiology* 2:1533–1542.
- Quince C, Walker AW, Simpson JT, Loman NJ, Segata N. 2017. Shotgun metagenomics, from sampling to analysis. *Nature biotechnology* 35:833–844.
- 1355 Raveh-Sadka T, Thomas BC, Singh A, Firek B, Brooks B, Castelle CJ, Sharon I, Baker R, Good M, Morowitz MJ, Banfield JF. 2015. Gut bacteria are rarely shared by co-hospitalized premature infants, regardless of necrotizing enterocolitis development. *eLife* 4. DOI: 10.7554/eLife.05477.
- R Development Core Team R. 2011. *R: A Language and Environment for Statistical Computing*. DOI: 10.1007/978-3-540-74686-7.
- 1360 Rinke C, Schwientek P, Sczyrba A, Ivanova NN, Anderson IJ, Cheng J-F, Darling A, Malfatti S, Swan BK, Gies EA, Dodsworth JA, Hedlund BP, Tsiamis G, Sievert SM, Liu W-T, Eisen JA, Hallam SJ, Kyrpides NC, Stepanauskas R, Rubin EM, Hugenholtz P, Woyke T. 2013. Insights into the phylogeny and coding potential of microbial dark matter. *Nature* 499:431–437.
- Roux S, Enault F, Hurwitz BL, Sullivan MB. 2015. VirSorter: mining viral signal from microbial genomic data. *PeerJ* 3:e985.
- 1365 Roux S, Krupovic M, Daly RA, Borges AL, Nayfach S, Schulz F, Sharrar A, Matheus Carnevali PB, Cheng J-F, Ivanova NN, Bondy-Denomy J, Wrighton KC, Woyke T, Visel A, Kyrpides NC, Elie-Fadrosh EA. 2019. Cryptic inoviruses revealed as pervasive in bacteria and archaea across Earth's biomes. *Nature microbiology*. DOI: 10.1038/s41564-019-0510-x.
- 1370 Schmidt TS, Hayward MR, Coelho LP, Li SS, Costea PI, Voigt AY, Wirbel J, Maistrenko OM, Alves RJ, Bergsten E, de Beaufort C, Sobhani I, Heintz-Buschart A, Sunagawa S, Zeller G, Wilmes P, Bork P. 2019. Extensive transmission of microbes along the gastrointestinal tract. *eLife* 8. DOI: 10.7554/eLife.42693.
- 1375 Seabold S, Perktold J. 2010. Statsmodels: Econometric and statistical modeling with python. In: *Proceedings of the 9th Python in Science Conference*. Scipy, 61.

- Segata N, Haake SK, Mannon P, Lemon KP, Waldron L, Gevers D, Huttenhower C, Izard J. 2012. Composition of the adult digestive tract bacterial microbiome based on seven mouth surfaces, tonsils, throat and stool samples. *Genome biology* 13:R42.
- 1380 Shaiber A, Eren AM. 2019. Composite Metagenome-Assembled Genomes Reduce the Quality of Public Genome Repositories. *mBio* 10. DOI: 10.1128/mBio.00725-19.
- Simón-Soro A, Tomás I, Cabrera-Rubio R, Catalan MD, Nyvad B, Mira A. 2013. Microbial geography of the oral cavity. *Journal of dental research* 92:616–621.
- Socransky SS, Manganiello SD. 1971. The oral microbiota of man from birth to senility. *Journal of periodontology* 42:485–496.
- 1385 Song SJ, Lauber C, Costello EK, Lozupone CA, Humphrey G, Berg-Lyons D, Caporaso JG, Knights D, Clemente JC, Nakielny S, Gordon JI, Fierer N, Knight R. 2013. Cohabiting family members share microbiota with one another and with their dogs. *eLife* 2:e00458.
- Storey JD, Tibshirani R. 2003. Statistical significance for genomewide studies. *Proceedings of the National Academy of Sciences of the United States of America* 100:9440–9445.
- 1390 Tatusov RL, Fedorova ND, Jackson JD, Jacobs AR, Kiryutin B, Koonin EV, Krylov DM, Mazumder R, Mekhedov SL, Nikolskaya AN, Rao BS, Smirnov S, Sverdlov AV, Vasudevan S, Wolf YI, Yin JJ, Natale DA. 2003. The COG database: an updated version includes eukaryotes. *BMC bioinformatics* 4:41.
- 1395 Touchon M, Bernheim A, Rocha EP. 2016. Genetic and life-history traits associated with the distribution of prophages in bacteria. *The ISME journal* 10:2744–2754.
- Vartoukian SR, Adamowska A, Lawlor M, Moazzez R, Dewhirst FE, Wade WG. 2016. In Vitro Cultivation of “Unculturable” Oral Bacteria, Facilitated by Community Culture and Media Supplementation with Siderophores. *PloS one* 11:e0146926.
- 1400 Whelan S, Goldman N. 2001. A general empirical model of protein evolution derived from multiple protein families using a maximum-likelihood approach. *Molecular biology and evolution* 18:691–699.

Supplementary Material

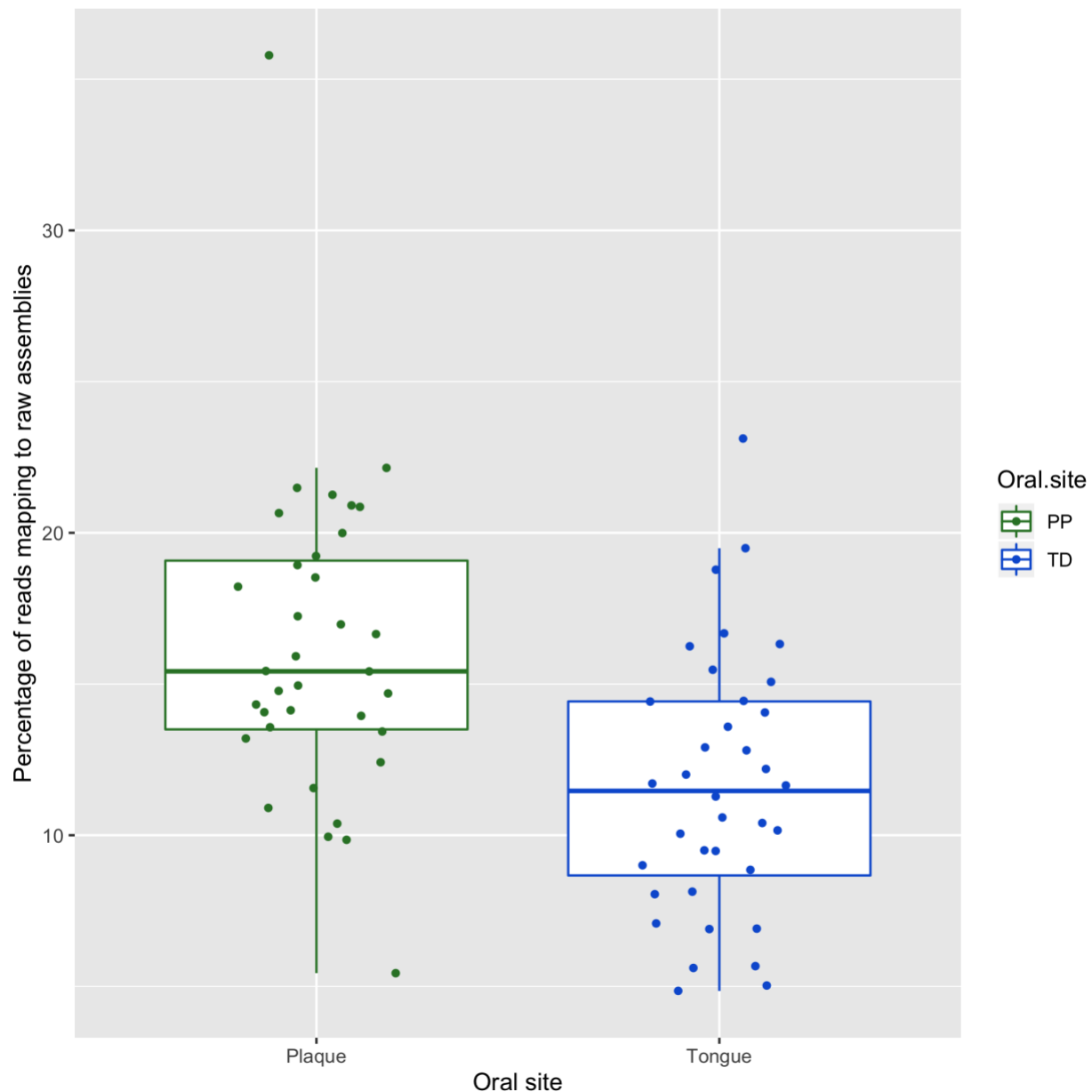
[doi:10.6084/m9.figshare.11634321](https://doi.org/10.6084/m9.figshare.11634321) gives access to all supplementary tables and the Supplementary Information.

Supplementary Figures

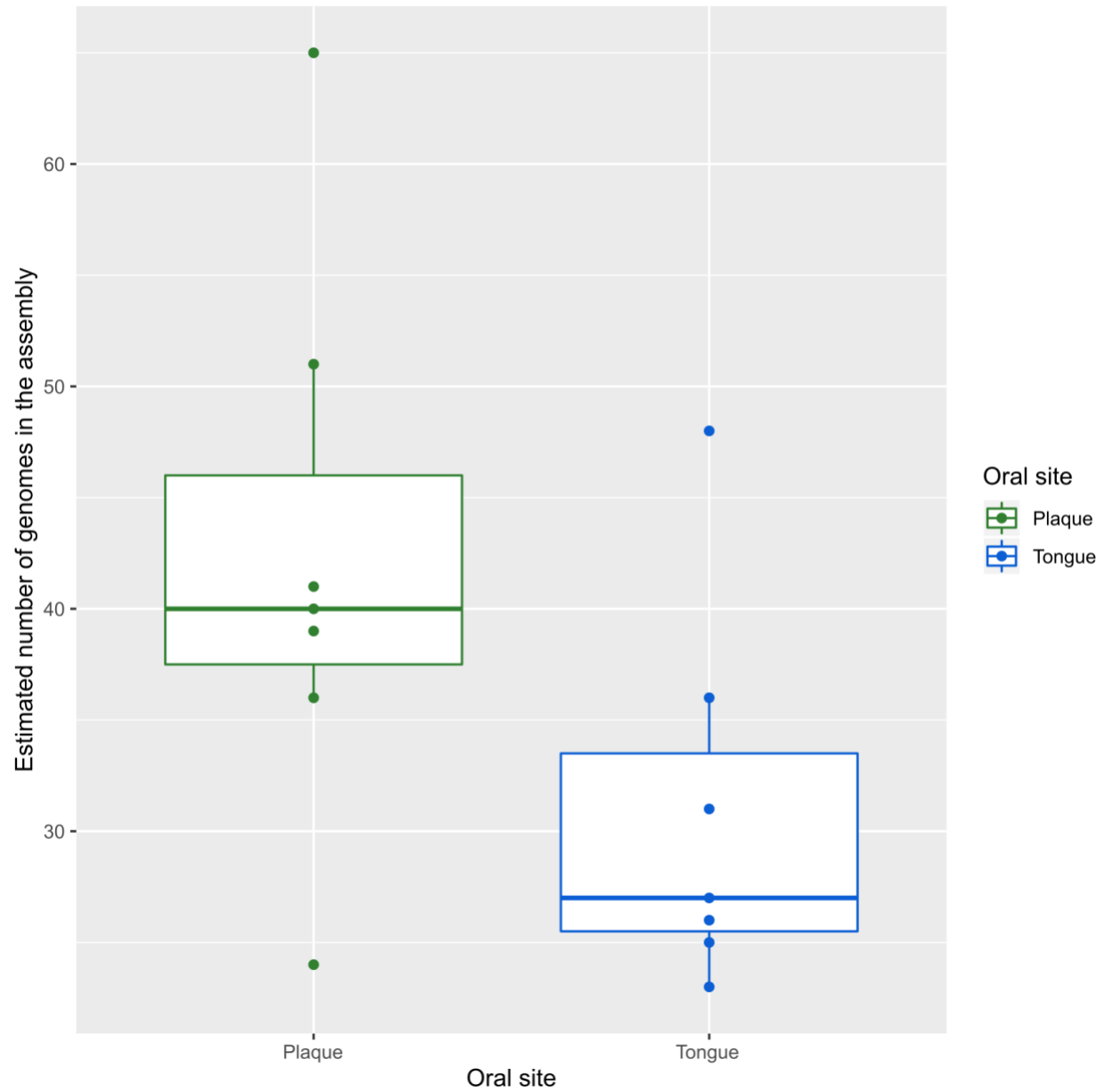


1405 **Supplementary Figure S1a: The percent of reads that map to MAGs is correlated with the quality of the assembly.** The percent of reads that mapped to the non-redundant collection of MAGs out of the total number of reads, excluding reads that mapped to the human genome is presented for each of the 35 plaque (green) and 36 tongue (blue) metagenomes as a function of the percent of reads that mapped to all contigs in the assembly. Blue curve represents

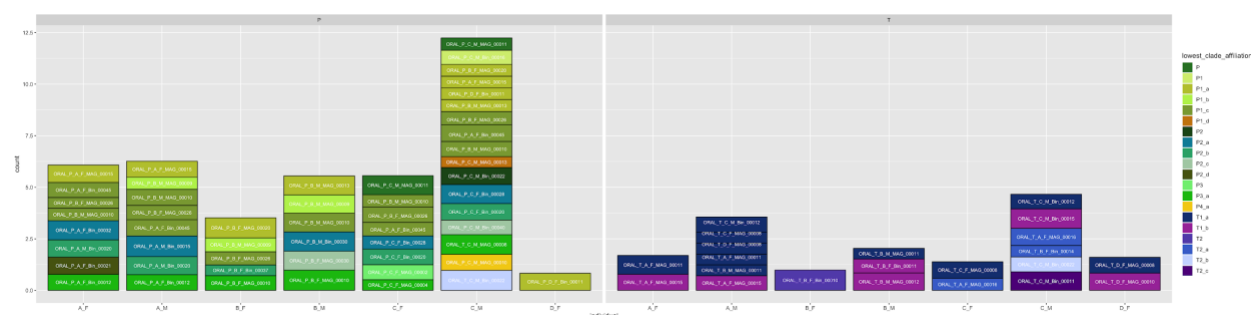
a linear regression model with the grey shaded area marking the 95% confidence intervals. R-squared value and p-value for the linear regression appear above the curve.



Supplementary Figure S1b: A significantly higher percentage of reads map to plaque assemblies as compared to tongue assemblies. Box and whisker plots for the percentage of reads that map to the raw (i.e. without binning) assembly contigs for plaque (green) and tongue (blue) metagenomes.



Supplementary Figure S1c: Higher number of estimated genomes in assemblies of plaque metagenomes as compared to assemblies of tongue metagenomes. Box and whisker plots for the number of estimated genomes, based on single copy core genes, in assemblies of plaque (green) and tongue (blue) metagenomes.



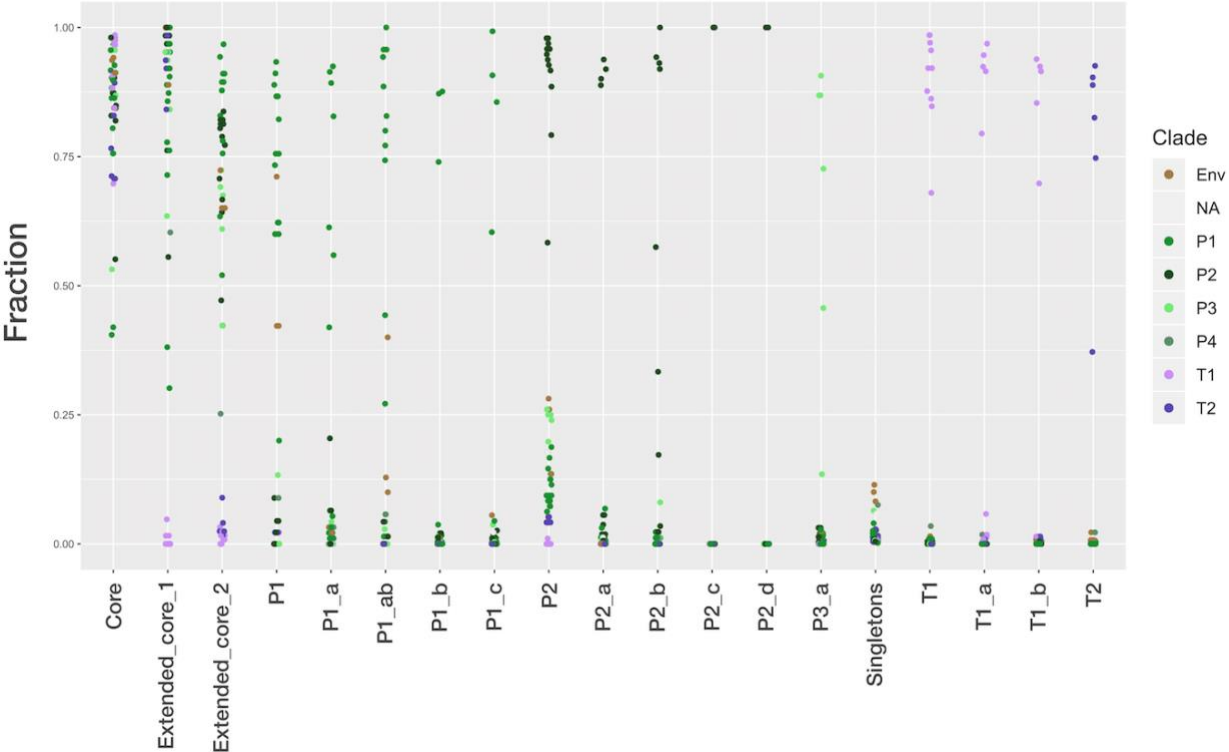


Figure S4a - Clade-specific clusters of GCs. Data points represent the fraction of the GCs in a given cluster of GCs in each genome.

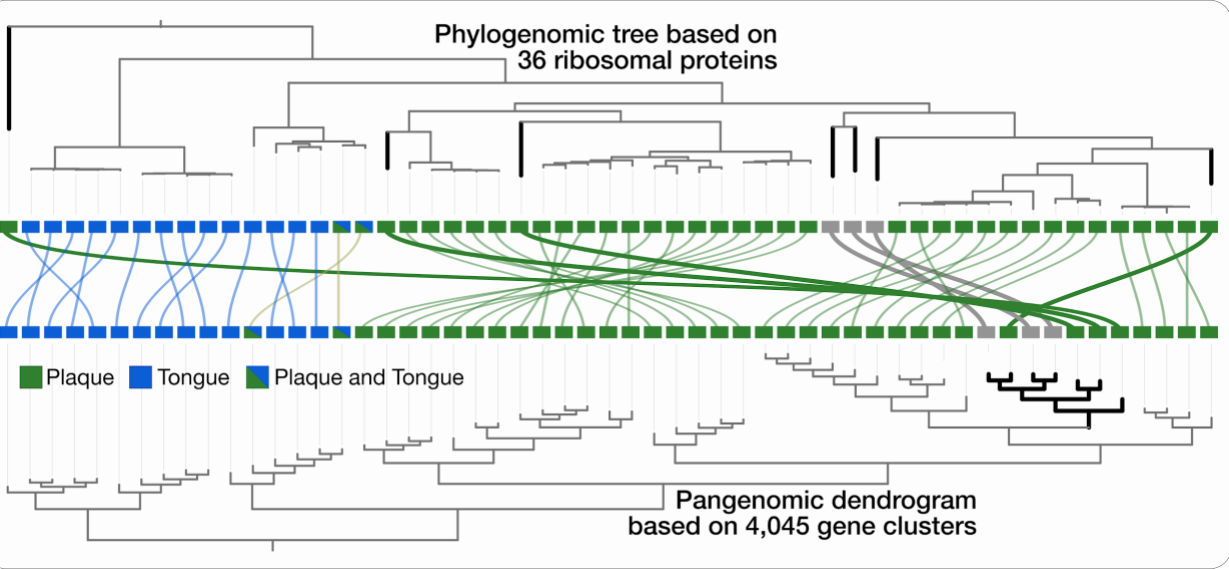
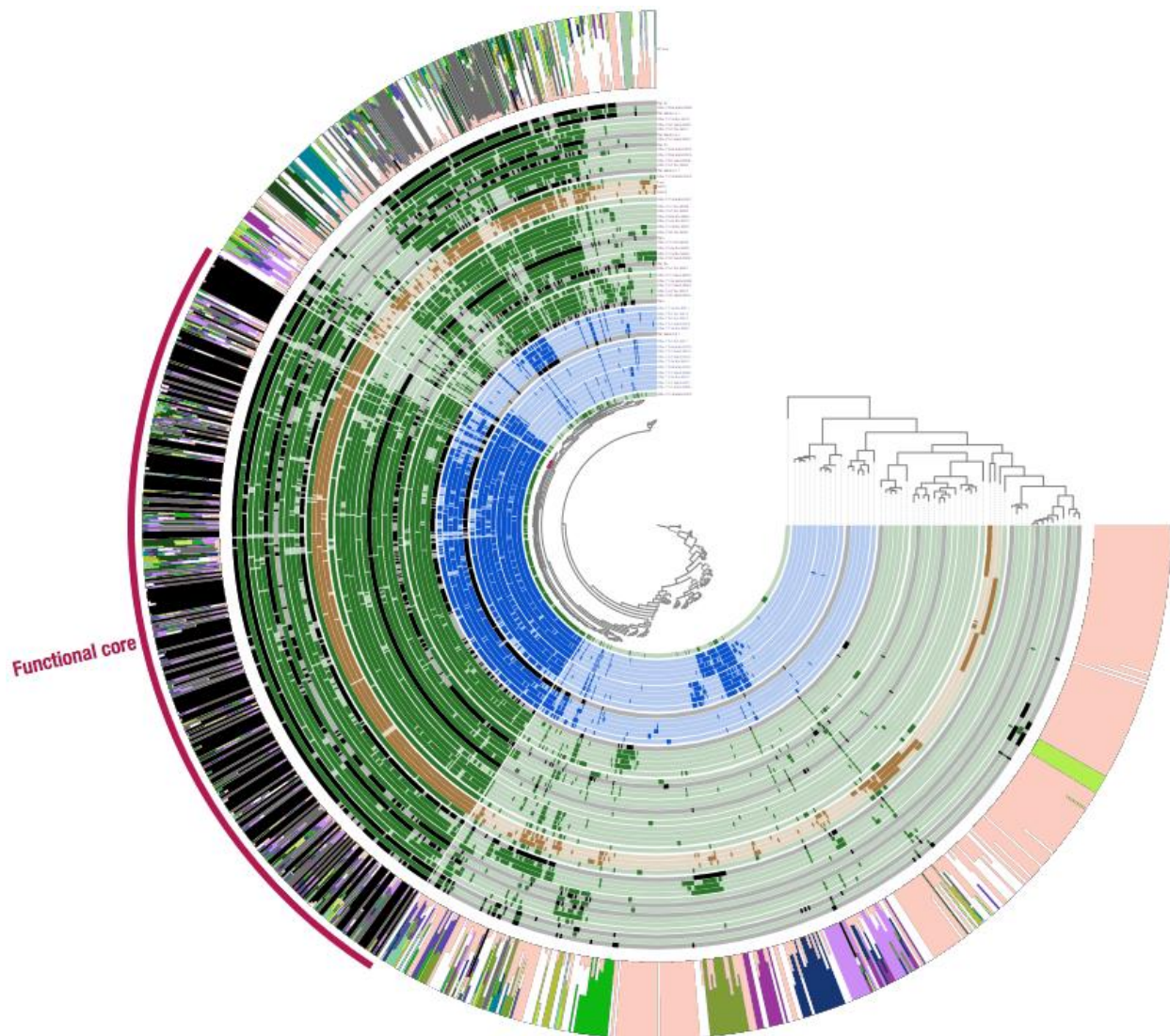


Figure S4b - Organization of TM7 genomes according to the occurrence of gene-clusters clusters oral genomes according to oral site affiliation. The dendrogram at the top represents the phylogenetic organization based on ribosomal proteins, while the

dendrogram on the bottom represents the hierarchical organization of genomes based on the gene-cluster frequency of occurrence across genomes using euclidean distance and ward ordination. The information at the center of the figure shows the site affiliation of each oral TM7 in accordance with Figure 4. Branches that appear in bold black color represent environmental and plaque-associated genomes that are phylogenetically distinct, but that are grouped together based on their gene content, and nested together with plaque-associated genomes.



FigureS4c - Functional core includes mostly core GCs, but also many clade specific GCs. Each of the 970 functions are organized in the tree in the center of the figure according to their occurrence in the 55 genomes (using Euclidean distance and

Ward's method) . The first 55 layers correspond to the TM7 genomes, where layers corresponding to tongue MAGs are blue, plaque MAGs are green, and previously published genomes are black. Bars in these 55 layers represent the presence of a function in the genome. The layers are ordered using the phylogenetic tree from Figure2b. The next layer includes a stacked bar representing the portion of GC bin affiliation of each gene associated with a function. The red arc in the outermost layer marks the functions that were defined as part of the core for this TM7 pangenome. Notice that while the majority of the core functions are associated with core GCs, there are many that are associated with clade-specific GCs.

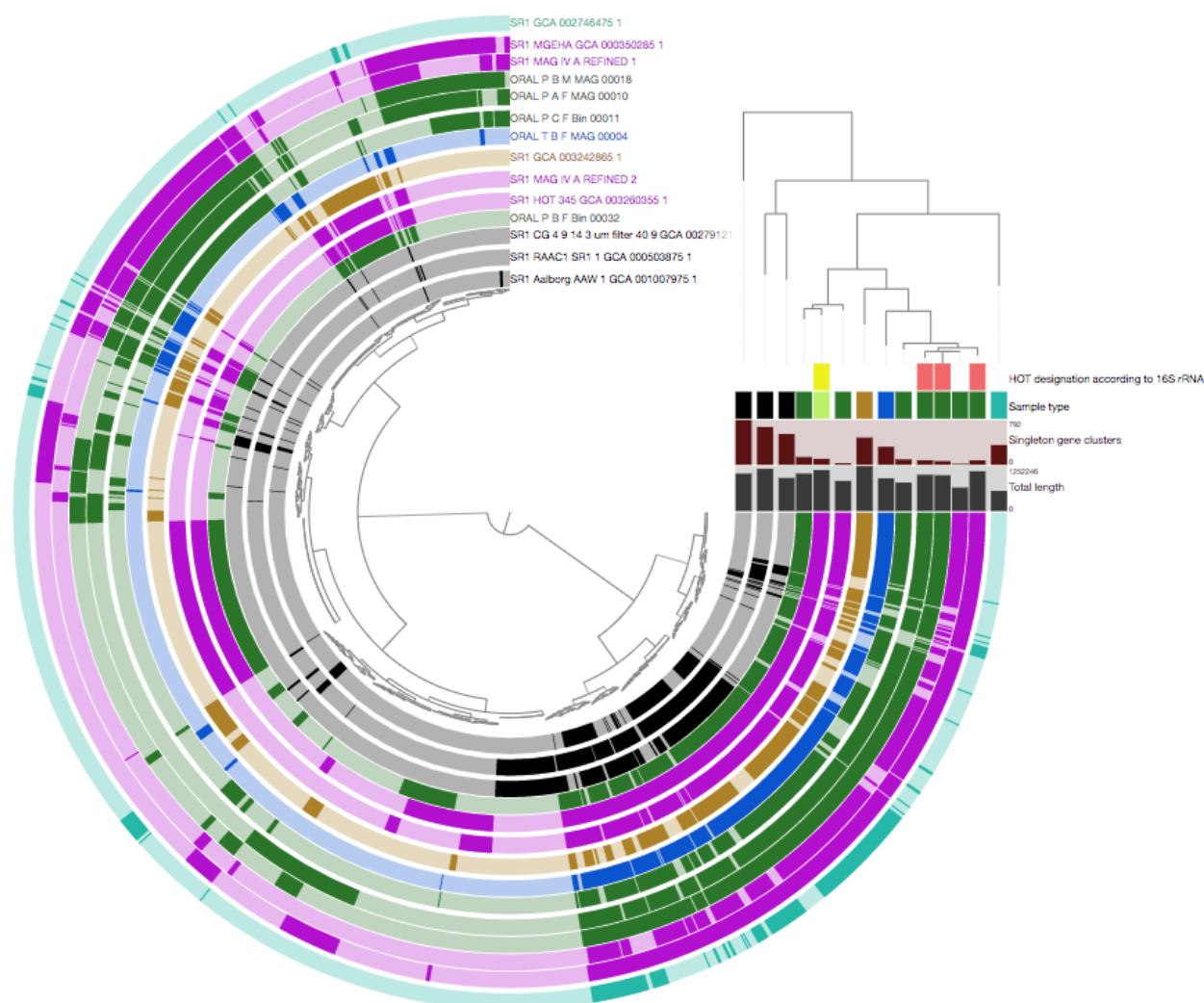


Figure S5a - pangenomic analysis of SR1 genomes. The dendrogram at the center of the figure organizes gene-clusters according to their occurrence across the 14 SR1 genomes. The circular layers correspond to the 14 SR1 genomes and are ordered according to their phylogenetic organization. In these circular layers, colored sections mark the presence of gene-clusters in the corresponding genome. On the top right, the phylogenetic tree is shown and below it, the four horizontal layers correspond to (top to bottom) 1) Human Oral Taxon designation according to 16S rRNA sequences 2) Sample type (environmental: black, plaque: dark green, saliva: light green, canine supragingival plaque: brown, tongue: blue, dolphin gingival sulcus: cyan) 3) Number of singleton gene-clusters 4) Total length of the genome.

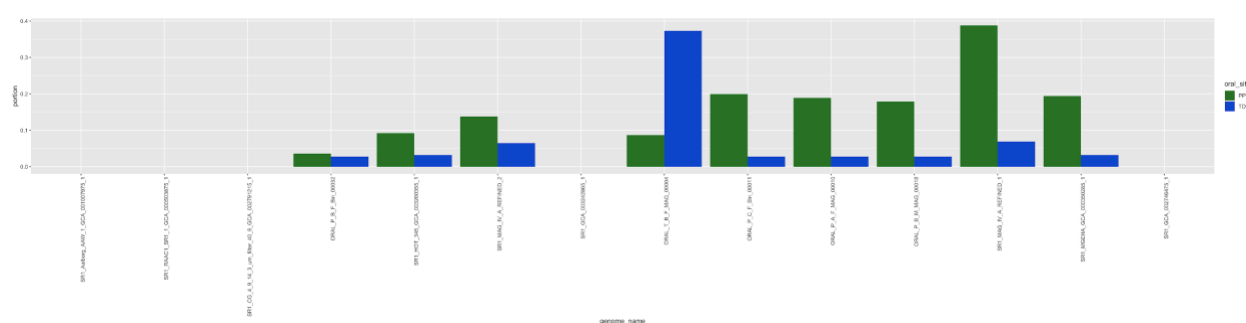


Figure S5b - Detection of SR1 populations in the HMP plaque and tongue samples reveals prevalent populations and niche specificity. Barplots showing the portion of plaque (green) and tongue (blue) HMP samples in which each SR1 was detected, using a detection threshold of 0.5.

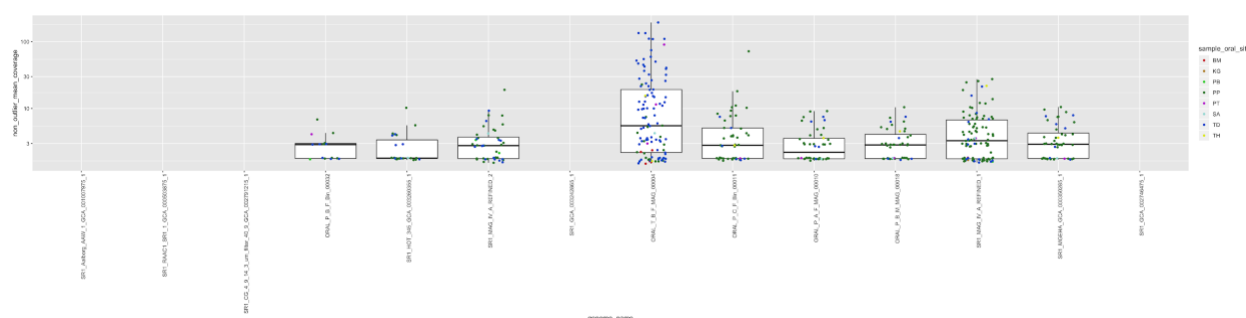


Figure S5c - Normalized coverage of SR1 populations in HMP oral samples according to sample type. Boxplots showing the normalized coverages of each SR1 in

plaque (green) and tongue (blue) HMP. For each genome, data is only shown for samples in which it was detected, according to the same criteria of detection used in Figure S5b.

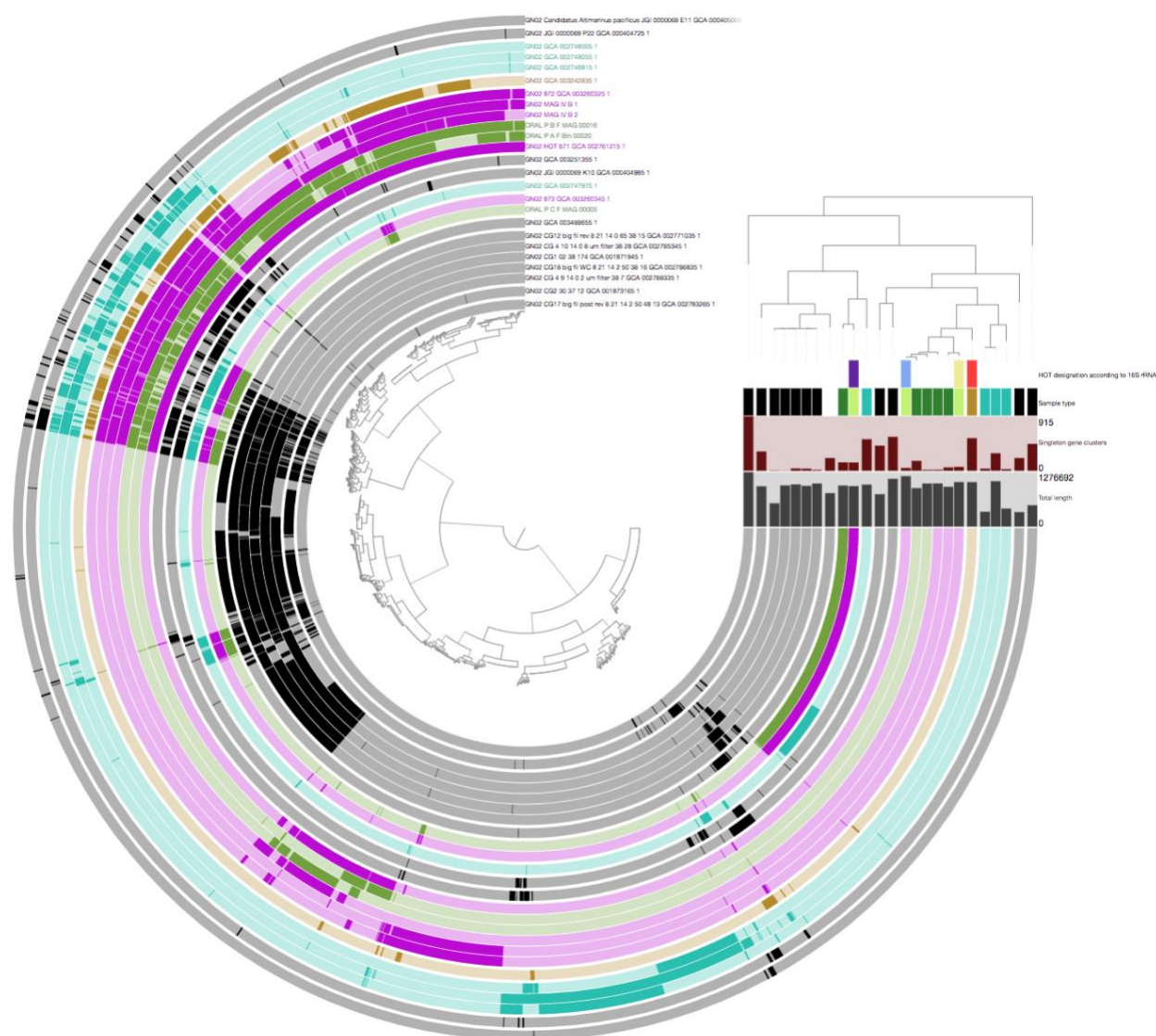


Figure S6a - pangenomic analysis of GN02 genomes. The dendrogram at the center of the figure organizes gene-clusters according to their occurrence across the 25 SR1 genomes. The circular layers correspond to the 25 SR1 genomes and are ordered according to their phylogenetic organization. In these circular layers, colored sections mark the presence of gene-clusters in the corresponding genome. On the top right, the phylogenetic tree is shown and below it, the four horizontal layers correspond to (top to bottom) 1) Human Oral Taxon designation according to 16S rRNA sequences 2) Sample type (environmental: black, plaque: dark green, saliva: light green, canine supragingival

plaque: brown, tongue: blue, dolphin gingival sulcus: cyan) 3) Number of singleton gene-clusters 4) Total length of the genome.

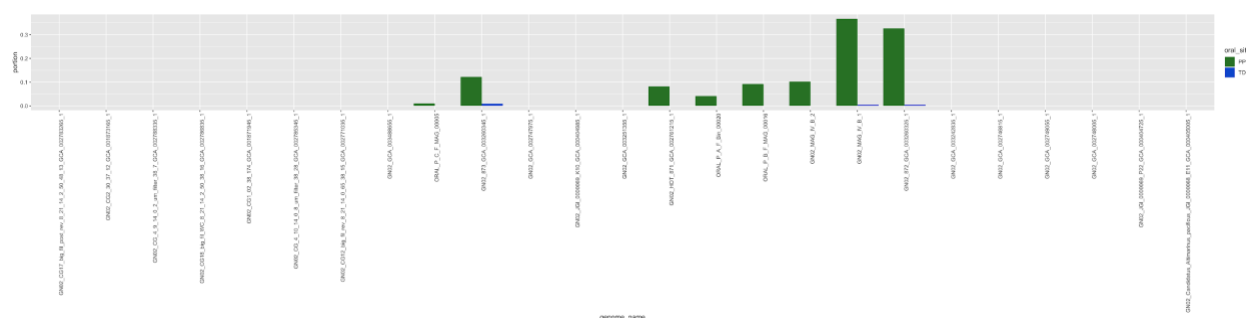


Figure S6b - Detection of GN02 populations in the HMP plaque and tongue samples reveals the plaque specificity of oral members of this candidate phylum. Barplots showing the portion of plaque (green) and tongue (blue) HMP samples in which each GN02 was detected, using a detection threshold of 0.5.

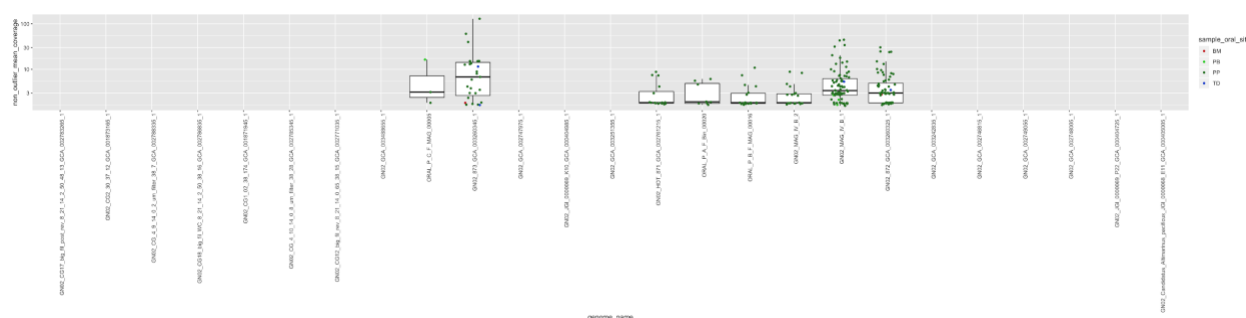


Figure S6c - Normalized coverage of GN02 populations in HMP oral samples according to sample type. Boxplots showing the normalized coverages of each GN02 in plaque (green) and tongue (blue) HMP. For each genome, data is only shown for samples in which it was detected, according to the same criteria of detection used in Figure S6b.

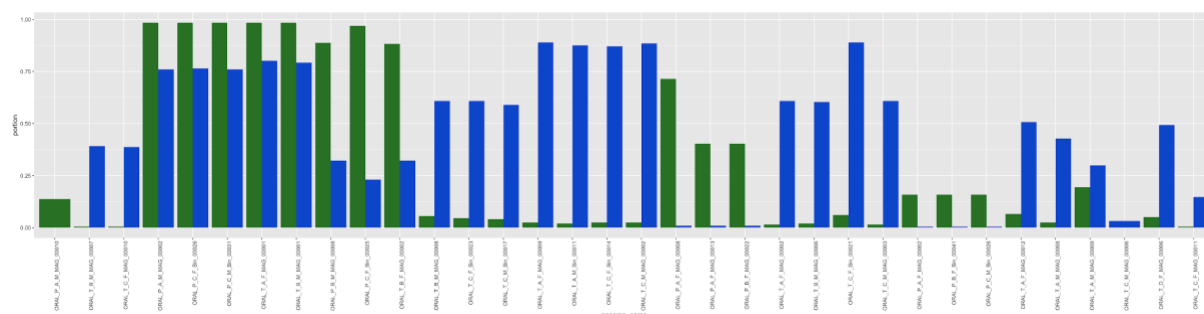


Figure S7a - Presence of the novel populations in HMP tongue and plaque samples.

Barplots of the portion of plaque (green) and tongue (blue) samples in which each of the novel genomes occur. The presence of a population in a sample was determined according to a threshold of 0.5 detection value.

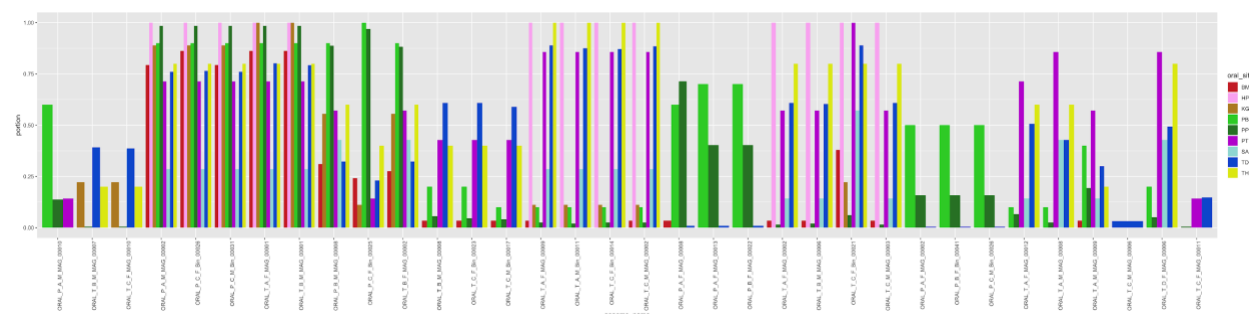


Figure S7b - Presence of the novel populations in HMP oral samples by sample type. Barplots of the portion of samples in which each of the novel genomes occur, plotted by sample type for all 9 HMP sample types in which at least one novel population was detected. The presence of a population in a sample was determined according to a threshold of 0.5 detection value.

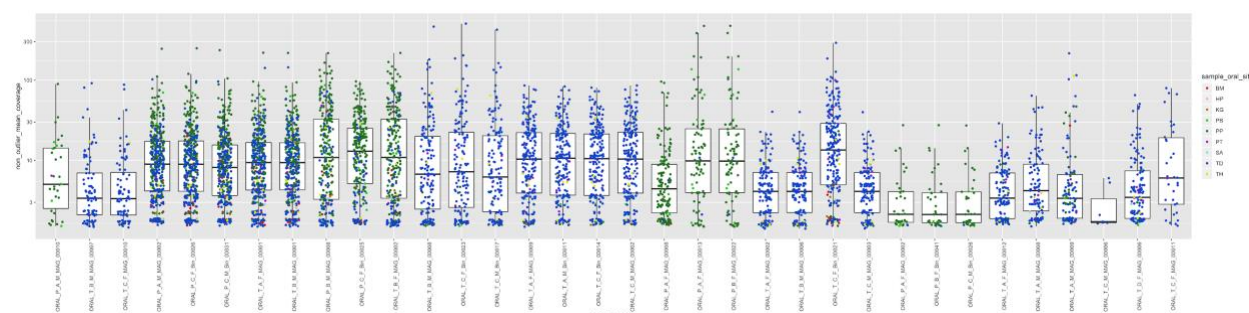


Figure S7c - Normalized coverage of the novel populations in HMP oral samples according to sample type. Boxplots of the normalized coverage of the novel population.

Color of data-points are according to the sample type. For each genome, data points are only shown for samples in which the genome was detected, according to the same detection threshold used in Figure S7b.

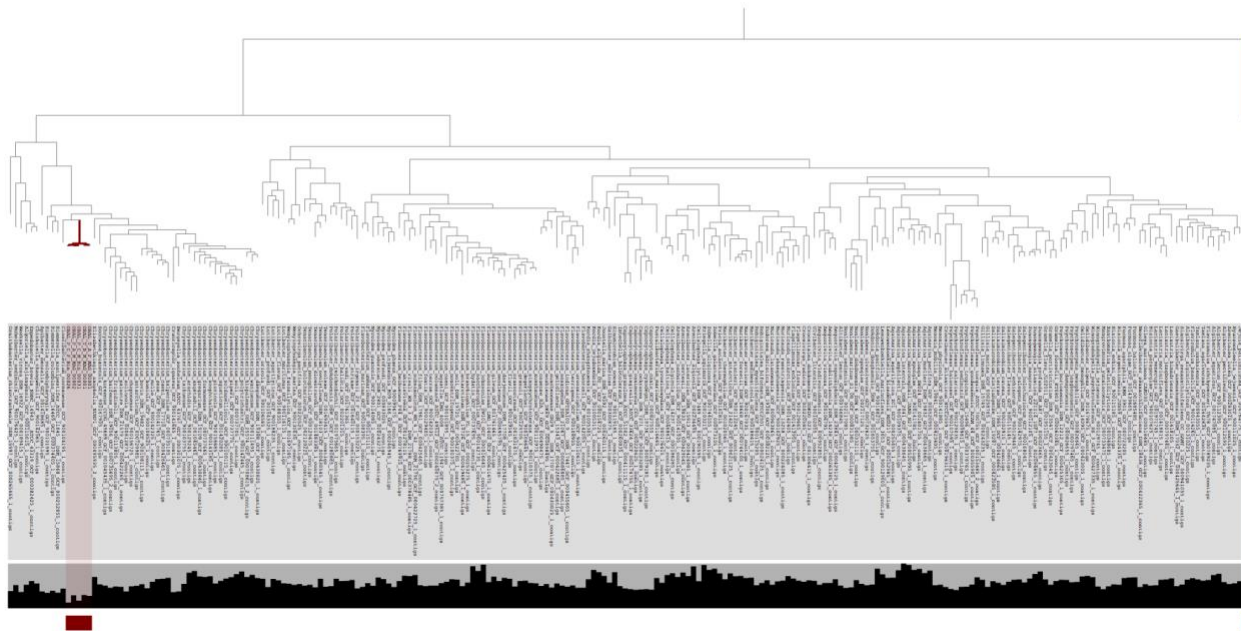


Figure S8 - Phylogenomic analysis of Flavobacteriaceae genomes indicates oral
MAGs represent an unnamed species in an unnamed genus within
Flavobacteriaceae. Below the dendrogram, layers include the name and length of each
genome. The 5 novel Flavobacteriaceae MAGs are indicated with red color and the
Prevotella genome that was used to root the tree is indicated with blue color.

Supplementary Tables

Supplementary Table 1: Details of collected samples, including statistics for
metagenomes and 16S rRNA amplicon data. (a) Summary of all 71 samples along with
summary of QC and assembly statistics for metagenomic short reads. (b) Summary of
statistics of contigs resulting from assemblies of metagenomes. (c) Number of reads
before and after QC for 16S rRNA gene amplicons. (d) Description of human subjects
and sample collection dates.

Supplementary Table 2: Details for MAGs. (a) Summary of sequence statistics,
taxonomy, and redundancy status for all 857 MAGs. (b) Summary of sequence statistics,
taxonomy, and redundancy status for the 790 non-redundant MAGs. (c) Mean coverage;
(d) Detection values; (e) Relative abundance estimations across our collection of 71
metagenomes for the 790 non-redundant MAGs. (f) Relative abundance estimations for

the 16 most abundant genera in our metagenomes based on our MAGs, in addition to relative abundance estimations for TM7 and the collection of contigs not assigned to any genomic bin (the 'UNKNOWN' bin). (g) Summary of sequence statistics for all 2,463 genomic bins.

1540 **Supplementary Table 3:** Dereplication statistics of metagenome-assembled genomes. (a) Phylum affiliation for all 857 MAGs. (b) Length and completion/redundancy of SCGs for the 857 MAGs. (c) Average nucleotide identity results for each pair of MAGs (only between pairs of MAGs that were classified to the same phylum). (d) Pairwise data for MAGs; this was used as the input for the redundancy analysis. (e) Affiliations for the 123
1545 redundant MAGs to one of 59 redundant groups.

Supplementary Table 4: Taxonomic annotation of metagenomes. Taxonomic annotation of short reads using KrakenUniq for each taxonomic level: (a) Domain. (b) Phylum. (c) Class. (d) Order. (e) Family. (f) Genus. (g) Species. (h) Relative abundance estimations for the 16 most abundant genera in our metagenomes based on KrakenUniq counts, in
1550 addition to relative abundance estimations for TM7 and the collection of reads not assigned any genus affiliation.

Supplementary Table 5: 16S ribosomal RNA amplicon sequencing results. Taxonomic annotations of amplicons in each of the 71 samples using GAST for each taxonomic level: (a) Phylum. (b) Class. (c) Order. (d) Family. (e) Genus. (f) Species. (g) Counts of amplicon
1555 sequence variants (ASVs) generated by the MED pipeline in each of the 71 samples. (h) Percent of counts of ASVs in the 71 samples. (i) Taxonomic affiliation of ASVs. (j) Relative abundance estimations for the 16 most abundant genera in our metagenomes based on MED/GAST results, in addition to relative abundance estimations for TM7 and the collection of ASVs not assigned any genus affiliation.

1560 **Supplementary Table 6:** Description of Human Oral Microbiome Database (HOMD) genomes used in this study. (a) Summary statistics for contigs for each of 1,334 HOMD genomes used in this study, including total length, number of contigs, and count of SCGs. (b) accession number information for each of the 1,334 HOMD genomes.

Supplementary Table 7: Metagenomic analysis of the TM7 genomes. (a) Detection of

our 43 non-redundant TM7 MAGs in our 71 samples as computed by anvi-mcg-classifier.

(b) Non-outlier mean coverage of our 43 non-redundant TM7 MAGs in our 71 samples

as computed by anvi-mcg-classifier. (c) Occurrence statistics for our 43 non-redundant

TM7 MAGs in our 71 samples as computed by anvi-mcg-classifier including normalized

coverage. (d) Accession and reference information for the 12 previously published TM7

genomes used in the pangenomic and metagenomic analyses. (e) Metadata, including

accession numbers, for the 156 TM7 genomes from NCBI used for the phylogenomic

analysis in Figure 3. (f) Average nucleotide identity (ANI) pairwise coverage values for

the 55 TM7 genomes (our 43 in addition to the previously published 12). (g) ANI percent

identity values. (h) ANI summary information for pairs of TM7 genomes with alignment

coverage greater than 0.25. (i) Clade affiliation (as defined in this study) and group

affiliation (as defined by Camanocha & Dewhirst, 2014) for the 52 oral TM7 genomes. (j)

Summary of mapping results of short reads from 481 oral metagenomes from the Human

Microbiome Projects (HMP) to the 55 TM7 genomes, including the total number of reads

that mapped from each sample, the total number of reads in each sample, and the

percentage of reads that mapped, as well as the oral site affiliation of each metagenomic

sample. (k) Metadata for HMP samples used in this study, including accession numbers.

(l) Detection of the 55 TM7 genomes in the HMP samples as computed by anvi-mcg-

classifier. (m) Non-outlier mean coverage of the 55 TM7 genomes in the HMP samples

as computed by anvi-mcg-classifier. (n) Summary of occurrence of the 55 TM7 genomes

in the HMP samples including the mean value of the non-outlier mean coverage for each

genome, number of samples in which a genome was detected above threshold of 0.5 with

break down for plaque and tongue samples, the percent of HMP samples in which the

genome was above detection threshold, and a the results from a chi-squared test to

compare occurrence of each genome between tongue and plaque samples. (o) Percent

detection (above 0.5 threshold) for each TM7 genome broken down by the 9 HMP sample

types. (p) Detection values as computed by anvi-mcg-classifier in subgingival plaque

metagenomes from the HMP and from Califf et al. 2017. (q) Non-outlier mean coverage

values subgingival plaque metagenomes (r) occurrence summary for the TM7 genomes

in subgingival plaque metagenomes. (s) Percent of reads from each subgingival plaque

sample that mapped to each TM7 genome. (t) Variability information as computed by anvi-gen-variability-profile for the “cosmopolitan” TM7 MAG (T_C_M_Bin_00022).

Supplementary Table 8: The TM7 pangenome and the functional enrichment analysis.

(a) Extended summary information for each of the 40,505 genes that were analyzed in the TM7 pangenomic analysis; including gene-cluster (GC) affiliation, functional annotations, contig name, start and stop positions, and functional enrichment information. (b) GC frequencies (c) Manually selected bins of GCs. (d) COG functions associated with each GC bin. (e) TM7 genes associated with type IV pilus. (f) Functional enrichment output. (g) Summary of prophage predictions using Virsorter, the inovirus detector or manual approaches. (h) Information on contigs that are associated with prophages. (i) start and stop nucleotide positions for prophages. (j) Genomic information on phage integrases. (k) Genomic information on phage terminases. (l) Output of CRISPRCasFinder. (m) Results of sequence search (using blastn) for CRISPR spacers. (n) Information on terminases identified in TM7 genomes. (o) Blast results of terminases against the NCBI protein database. (p) Blast results for CAS9 sequences. (q) Functional occurrence table for the TM7 pangenome. (r) Functions that are associated with the ‘Extended Core 2’ GC bin and occur in the clade P4 genome (ORAL_P_C_M_MAG_00010).

Supplementary Table 9: The SR1 and GN02 metapangenomes. (a) Reference

information on GN02 genomes used for analysis. (b) Reference information on SR1 genomes used for analysis. (c) Summary information of GN02 pangenome. (d) Summary information of SR1 pangenome. (e) Average nucleotide identity (ANI) pairwise coverage values for the 25 GN02 genomes. (f) ANI percent identity values for the GN02 genomes. (g) ANI pairwise coverage values for the 14 SR1 genomes. (h) ANI percent identity values for the SR1 genomes. (i) Occurrence summary for the GN02 genomes in the HMP oral metagenomes. (j) Non-outlier mean coverage of the GN02 genomes in the HMP samples as computed by anvi-mcg-classifier. (k) The portion of reads that mapped to each GN02 genome from each of the HMP samples. (l) Occurrence summary for the SR1 genomes in the HMP oral metagenomes. (m) Non-outlier mean coverage of the SR1 genomes in

the HMP samples as computed by anvi-mcg-classifier. (n) The portion of reads that mapped to each GN02 genome from each of the HMP samples.

Supplementary Table 10: Novel MAGs. (a) Group affiliations for the novel MAGs. (b) Blast results for ribosomal proteins of novel MAGs. (c) Summary of the blast results of ribosomal proteins of novel MAGs. (d) Summary information on the groups of novel MAGs. (e) Occurrence summary for the novel MAGs in the HMP oral metagenomes. (f) Detection of the novel MAGs in the HMP samples as computed by anvi-mcg-classifier. (g) Non-outlier mean coverage of the novel MAGs in the HMP samples as computed by anvi-mcg-classifier. (h) The portion of reads that mapped to each novel MAG from each of the HMP samples. (i) ANI pairwise coverage values for 41 Flavobacteriaceae (our 5 f__Flavobacteriaceae MAGs and 36 representatives of Flavobacteriaceae species). (j) ANI pairwise identity values for the 41 Flavobacteriaceae genomes.