## RESEARCH

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

Alon Shaiber[1,2], Amy D. Willis[3], Tom O. Delmont[4], Simon Roux[5], Lin-Xing Chen[6], Abigail C. Schmid[7], Mahmoud Yousef[8], Andrea R. Watson[1,9], Karen Lolans[1], Özcan C. Esen[1], Sonny T. M. Lee[10], Nora Downey[11], Hilary G. Morrison[11], Floyd E. Dewhirst[12,13], Jessica L. Mark Welch[11*] and A. Murat Eren[1,2,9,11*]

* Correspondence: meren@
uchicago.edu; jmarkwelch@mbl.edu
[1]Department of Medicine, University
of Chicago, Chicago, IL 60637, USA
[11]Josephine Bay Paul Center for
Comparative Molecular Biology and
Evolution, Marine Biological
Laboratory, Woods Hole, MA 02543,
USA
Full list of author information is
available at the end of the article

## Abstract

**Introduction:** Microbial residents of the human oral cavity have long been a major focus of microbiology due to their influence on host health and intriguing patterns of site specificity amidst the lack of dispersal limitation. However, the determinants of niche partitioning in this habitat are yet to be fully understood, especially among taxa that belong to recently discovered branches of microbial life.

**Results:** Here, we assemble metagenomes from tongue and dental plaque samples from multiple individuals and reconstruct 790 non-redundant genomes, 43 of which resolve to TM7, a member of the Candidate Phyla Radiation, forming six monophyletic clades that distinctly associate with either plaque or tongue. Both pangenomic and phylogenomic analyses group tongue-specific clades with other host-associated TM7 genomes. In contrast, plaque-specific TM7 group with environmental TM7 genomes. Besides offering deeper insights into the ecology, evolution, and mobilome of 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 microbes. Additionally, we report that prophages are widespread among oral-associated TM7, while absent from environmental TM7, suggesting that prophages may have played a role in adaptation of TM7 to the host environment.

**Conclusions:** Our data illuminate niche partitioning of enigmatic members of the oral cavity, including TM7, SR1, and GN02, and provide genomes for poorly characterized yet prevalent members of this biome, such as uncultivated Flavobacteriaceae.

**Keywords:** Metagenomics, Metapangenomics, Niche partitioning, Human oral cavity, Candidate phyla radiation, Saccharibacteria

## Background

Since the inception of microbiology as a new discipline following Antonie van Leeuwen-hoek's historical observation of the animalcules [1], 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 [2], (2) the differential influence of the host immunity throughout the oral tissue types [3], 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 [4, 5], temporal dynamism [6–8], spatial organization [9], and site specificity [10–12], where they influence the host health [13] and the ecology of the gastrointestinal tract [14]. 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 [15], DNA-DNA hybridization [16], and cloning [17] 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, and (4) ability of microbes to adhere both to the substrate and to one another [15, 18, 19], 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) [20], 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 con-firmed major taxonomic differences between microbial communities associated with dental plaque and mucosal sites in the mouth [21, 22]. Recruiting metagenomic short reads using single-copy core genes, Donati et al. demonstrated that while some mem-bers of the genus *Neisseria* were predominantly found in tongue dorsum samples, others were predominant in plaque samples [23], 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 [11]. Strong associations between oral sites and their microbial residents even at the finest levels of resolution raise ques-tions regarding the drivers of such exclusiveness [12]. 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) [24, 25] 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 [26]. 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 [27], including

Shaiber *et al. Genome Biology* (2020) 21:292

Page 3 of 35

members of the Candidate Phyla Radiation (CPR) [28], 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 [29] and functional makeup [30]; they lack many biological pathways that are considered essential [28] and have been shown to rely on epibiotic lifestyles [31], with a complex and poorly understood relationship with a microbial host [32]. Their unique lifestyle [33], diversity and prevalence in the oral cavity [34], association with distinct oral sites [31], and potential role in disease [35, 36] 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 [37] and were followed by He et al.'s pioneering work that brought the first TM7 population into culture [33], 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 [38–40]. Recent genome-resolved and single-amplified genomics studies have also produced genomes for oral GN02 and SR1 [41, 42], and recently, the first targeted isolation of oral SR1 strains has been reported, but genomes were not produced [38]. 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 (Additional file 1: Table S1a). 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 [43]. The resulting 14 co-assemblies (7 people × 2 sites) contained 267,456 contigs longer than 2500 nts that described approximately 1163 million nucleotides and 1,554,807 genes (Additional file 1: Table S1b). To reconstruct genomes from these metagenomes, we used a combination of automatic and manual binning strategies that resulted in 2463 genome bins. Independent assembly and binning of metagenomes from similar habitats can result in the recovery of multiple near-identical genomes [44, 45]. To increase the accuracy of downstream analyses, we employed only the 857 of 2463 bins that were 0.5 Mbp or larger (Additional file 2: Table S2g), 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 "Material and methods"). This resulted in a final collection of 790 non-redundant genomes (Additional file 2: Table S2a-b, Additional file 3: Table S3).

Automatic binning approaches can yield composite genomes that suffer from contamination, influencing downstream ecological and evolutionary insights [46], even when single-copy core genes suggest the absence of an apparent contamination [43]. Here we sampled each subject on at least 4 separate days to improve the accuracy of automatic binning through differential coverage [47]. 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 "Availability of data and materials").

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 (Additional file 1: Table S1a, Additional file 4: Fig. 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, Additional file 4: Fig. 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 (Additional file 4: Fig. 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 1332 genomes from the HOMD (accessed on August 1, 2018) (Additional file 7: Table S6b). 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 (≥ 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 (Table S2e-f) and the taxonomic composition of metagenomes estimated by short reads (Additional file 5: Table S4a-h) and 16S rRNA gene amplicon sequencing (Additional file 6: Table S5a-j and Additional file 8: Fig. SI1-SI3).
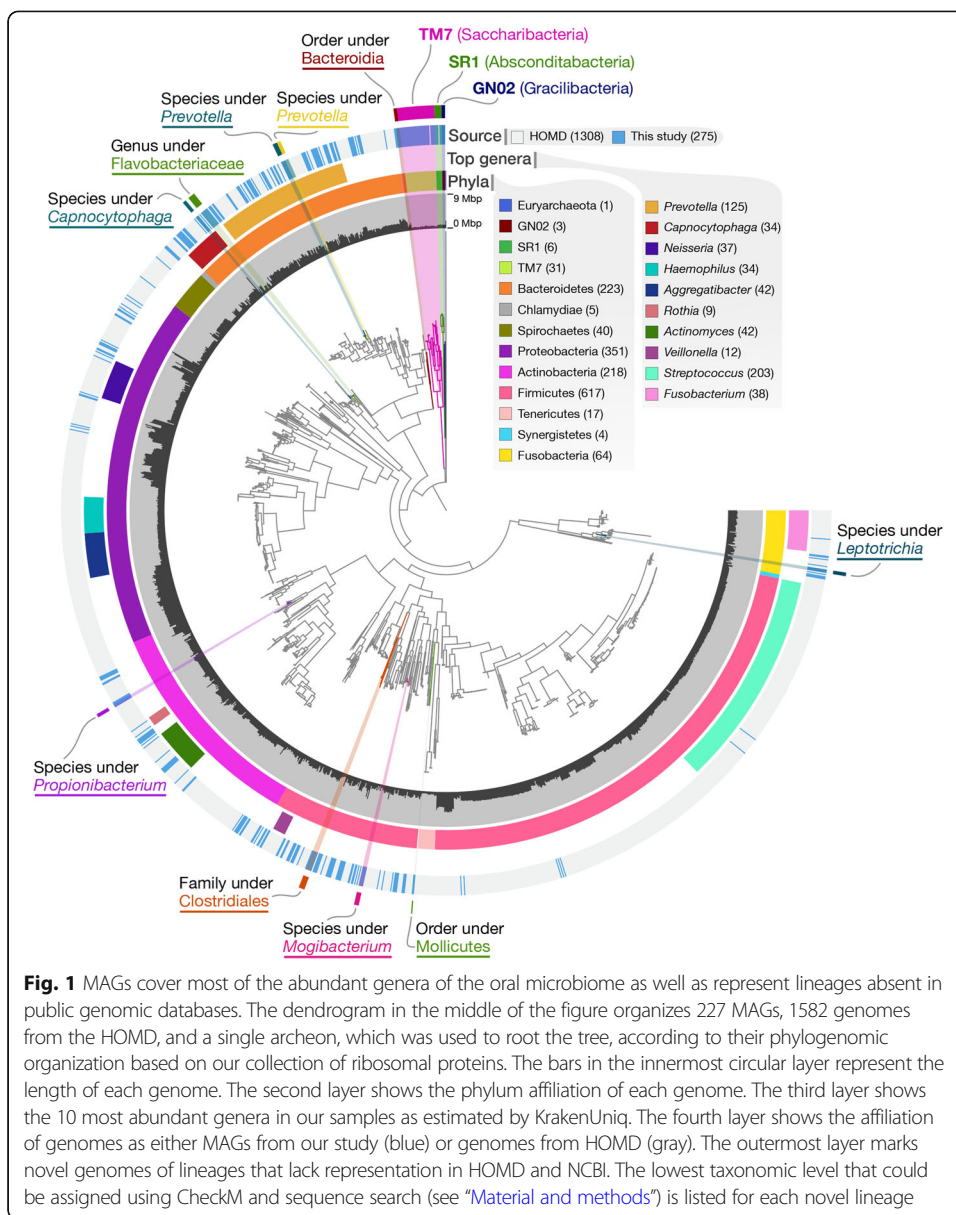
Some lineages contained members exclusively from our collection and not in the HOMD (Fig. 1), including 51 genomes that we identified as members of the CPR, which formed a distinct branch, as expected (Fig. 1). Our MAGs also included novel genomes from non-CPR lineages not represented in the HOMD (Fig. 1). While some of these deeply 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 "Material and methods"), during which we noticed and corrected binning errors in approximately 9% of our genomes (data not shown).

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

## TM7 phylogenomic clades are site-specific

Assessing the abundance of each TM7 MAG in the environment provides 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 (Fig. 2a, Additional file 9: Table S7a-c). We defined a given TM7 population as detected in one of the 71 samples if at least 50% of the nucleotides of the genome were covered by at least one short read.

**Fig. 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 archeon, 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 (gray). The outermost layer marks novel genomes of lineages that lack representation in HOMD and NCBI. The lowest taxonomic level that could be assigned using CheckM and sequence search (see "Material and methods") is listed for each novel lineage

Of the 43 TM7 populations represented by our MAGs, we detected 42 either only in plaque or only in tongue samples, but never in both (Fig. 2a, Additional file 4: Fig. S2a, Additional file 4: Fig. S2b). Thus, the vast majority of TM7 populations were site-specific. An exception was T_C_M_Bin_00022, which was detected in both tongue and plaque of a single subject; this entire MAG was covered in 4 of 6 tongue samples and 6 of 6 plaque samples from participant C_M, but not in any other participant (Fig. 2a). Thinking that T_C_M_Bin_00022 may represent multiple subpopulations each of which is in fact site-specific, we investigated its population structure in tongue and plaque using single-nucleotide variants (SNVs) at nucleotide positions with sufficient coverage for accurate assessment, which we defined as > 20× coverage across all samples. None of the 22,507 nucleotide positions that passed this coverage criterion showed any variation in plaque metagenomes. In contrast, 449 of these positions (2%)

**Fig. 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 turquoise marks the "cosmopolitan" TM7; "Study"—the study associated with each genome: our MAGs (purple), Espinoza et al. [41] (teal), Marcy et al. [37] (blue), He et al. [33] (red), and Cross et al. [38] (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 groups. † Refined versions of genomes, which we previously published [46]. ‡ Genomes from IMG that we refined in this study, but for which accession numbers for refined versions are available in Cross et al. [38]

showed notable variation in tongue samples, defined as a frequency of the minor allele (the less-common base) of at least 10%. The median frequency of the minor allele was 40% (Table S7t). This demonstrates (1) the existence of at least one subpopulation specific to tongue (represented by minor alleles in tongue metagenomes) and (2) that the monoclonal plaque population also occurs in tongue as the dominant member of the subpopulations associated with tongue. Other than this 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 (Fig. 2).

The site specificity of our TM7 MAGs led us to ask whether entire clades, or only individual genomes, showed a particular site specificity. For this analysis we combined our 43 MAGs with 9 human oral TM7 genomes from the literature. These genomes included the first cultivated strain of TM7, TM7x [33], and a MAG from Marcy et al. [37]. In addition, we used 3 single-amplified genomes (SAGs) that we downloaded from the Integrated Microbial Genomes and Microbiomes database (IMG/M) [48] and manually refined (see "Material and methods"), and 4 TM7 MAGs we manually refined [46] from composite Espinoza et al. MAGs [41] (Additional file 9: Table S7d).

The phylogenomic analysis of these 52 genomes separated tongue- and plaque-associated genomes into distinct branches. A single clade on the tree contained 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, grouped together with plaque-associated genomes, suggesting that its niche is most likely dental plaque rather than tongue (Fig. 2b).

### TM7s found in plaque and tongue share exclusive ancestry with environment- and host-associated TM7s
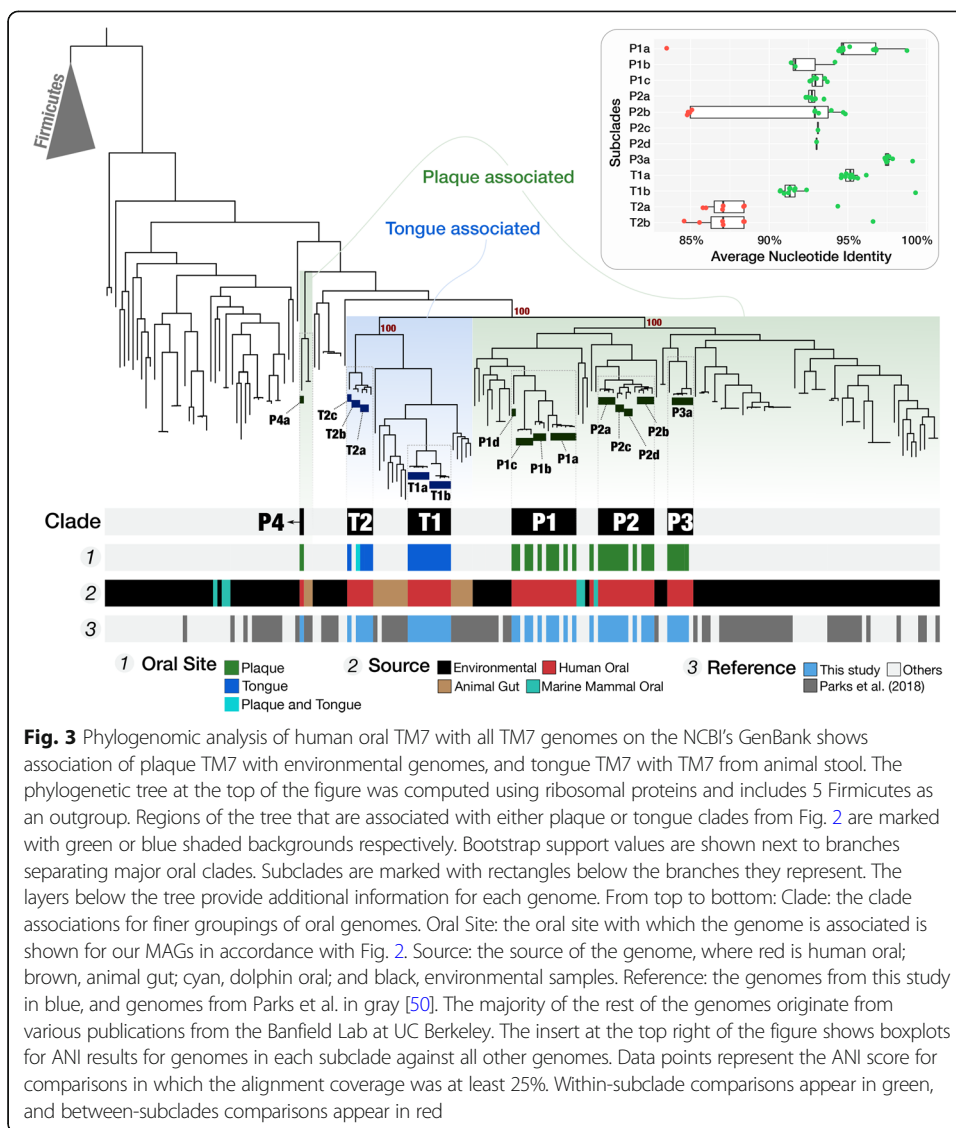
Previous studies have shown that the human-associated members of TM7 are polyphyletic, clustering phylogenetically together with TM7 genomes of environmental origin [34, 49]. Taking advantage of the large number of genomes we have reconstructed, we revisited this observation to determine the phylogenetic relationships between plaque-associated, tongue-associated, and environmental TM7s. We carried out a phylogenomic analysis, based on a set of ribosomal protein genes, using our 43 MAGs in addition to the 150 genomes of TM7 that were publicly available in the NCBI's GenBank database as of 1/16/2019 (Fig. 3).

On the resulting phylogenomic tree, we identified six monophyletic clades that contained exclusively human oral sequences (Fig. 3). Each of these was associated either with tongue (T1, T2) or with plaque (P1, P2, P3, P4) (Fig. 3). We then used a pairwise comparison of the ANI of oral TM7 genomes to identify more refined subclades that correspond to genus- and species-level groups, including 12 species of TM7 that are each represented by at least 2 genomes in our collection (Fig. 3, Additional file 9:Table S7f-h, Additional file 8).

Previous studies based on 16S rRNA gene sequences divided TM7 into 6 groups [34], leading us to ask whether these 6 groups corresponded to our 6 phylogenomic clades. Genomes reconstructed from metagenomic short reads often lack ribosomal RNA operons; hence, associating new genomes with populations only known through their 16S rRNA gene sequences poses a significant challenge. We surmounted this obstacle by carrying out nanopore long-read sequencing of samples from an additional volunteer (individual L in Additional file 1: Table S1d). Long reads that contained both 16S rRNA genes and flanking sequences allowed us to compare our clades to the 6 groups (G1-G6) of oral TM7 that were previously described based on 16S rRNA gene amplicons [34]. We determined that our monophyletic clades T1, T2, and P4 correspond to G3, G6, and G5, respectively (Additional file 9: Table S7e,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 did not recover any MAGs for TM7 groups G2 and G4, which have been previously shown to have low prevalence as compared to other TM7 groups [31].

The six clades of human-associated TM7 genomes in the phylogenomic tree each had sister groups composed of genomes isolated from non-human sources. Tongue clades T1 and T2 grouped with genomes recovered from animal gut and together formed a deep monophyletic branch of an exclusively host-associated superclade (Fig. 3). In contrast, plaque clades were interspersed with genomes from environmental sources (Fig. 3). The exceptions to this clear distinction between plaque and tongue clades were T_C_M_Bin_00022, a cosmopolitan oral population that

**Fig. 3** Phylogenomic 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 Fig. 2 are marked with green or blue shaded backgrounds respectively. Bootstrap support values are shown next to branches separating major 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 associations for finer groupings of oral genomes. Oral Site: the oral site with which the genome is associated is shown for our MAGs in accordance with Fig. 2. Source: the source of the genome, where red is human oral; brown, animal gut; cyan, dolphin oral; and black, environmental samples. Reference: the genomes from this study in blue, and genomes from Parks et al. in gray [50]. 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

was placed within the clade T2, and the plaque-associated P_C_M_MAG_00010 (the only member of the clade P4) that was placed as a far outlier to all other oral TM7 and grouped together with genomes from animal gut (baboon feces). Beyond these exceptions, the animal gut-associated genomes that were placed within the host-associated superclade originated 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 super-clade and their host habitats (Fig. 3, Additional file 9: Table S7e). Similarly, the inclusion of genomes recovered from dolphin dental plaque (shown between P1 and P2 on Fig. 3) together with human-plaque-associated TM7 suggests an ancient association for this clade of TM7 with the dental plaque environment.

The phylogenomic grouping of tongue-associated TM7 genomes with TM7 genomes from animal gut, to the exclusion of environmental TM7, suggests that tongue and gut TM7 share a higher degree of ancestral relationship compared to those that are

associated with plaque and with environments outside of a host. The selective forces imposed by animal hosts likely play a significant role in the evolution of the host-associated microbes [51] and can explain the emergence of a deep monophyletic branch of an exclusively host-associated superclade of TM7. In parallel, the absence of plaque TM7 from this host-associated superclade and the sister-group associations between plaque TM7 and TM7 clades of non-host environments suggest that from a microbial point of view, at least in the context of TM7, dental plaque resembles a non-host environment.

But what led to the divergence of tongue and plaque TM7? TM7 have highly reduced genomes and have been found to be epibionts of other bacteria, primarily Actinobacteria [31, 52]. In consideration of the TM7 lifestyle that has so far been observed, one reasonable hypothesis is that the bacterial hosts of each TM7 clade are the drivers of the link between TM7 ecology and evolution. Such a hypothesis would imply that the similarity between tongue 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, this hypothesis would imply the exclusion of such suitable hosts from the plaque environment, and vice versa, and 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 [11, 53] and that Actinobacteria are rare in the human gut [21]. 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 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 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 (tongue/plaque) in our dataset, suggesting a

link between environmental distribution of these genomes and their evolutionary history in the context of ribosomal proteins.

### Prevalence of TM7 across individuals is associated with TM7 clades, linking TM7 ecology and evolution
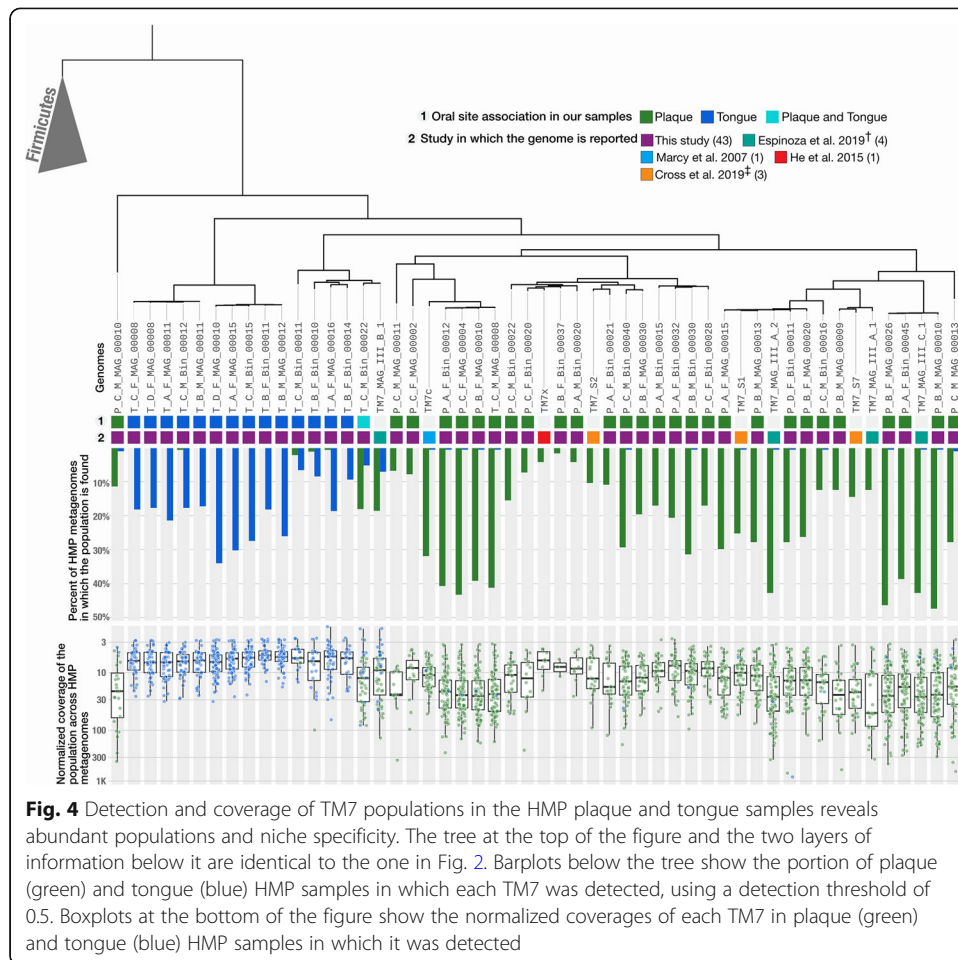
Since the samples we used to generate our 43 TM7 MAGs represent only 7 individuals, we next sought to identify whether the patterns of site specificity in these samples were representative of the distribution of TM7 among a wider cohort of healthy individuals.

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 paired-end reads) that included 196 samples from supragingival plaque and 217 tongue dorsum samples and were sampled from 131 individuals (Additional file 9: Table S7j-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 "Material and methods"). In addition to oral genomes, we also included three TM7 genomes (RAAC3, GWC2, and S_aal) that were reconstructed from environmental samples and manually curated to circularity [28, 52, 54].

The occurrence pattern of oral TM7 genomes across the HMP individuals matched their occurrence in our seven participants. 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 (Fig. 4). Members of subclade T2_b appeared to be cosmopolitan also in these data and were detected in both plaque and tongue samples (Fig. 4). As expected, the three environmentally derived TM7 genomes were not detected in any oral metagenome (Fig. 4, Additional file 9: Table S7l-n).

TM7 populations varied widely in their prevalence. The most prevalent tongue-associated genome and plaque-associated genome in our collection were detected in 45% and 50% of the HMP individuals, respectively (Fig. 4). In contrast, TM7x, the first cultured strain of TM7, was detected in only 5% of the HMP individuals (Fig. 4).

TM7 are commonly found in additional sites in the oral cavity other than the tongue and dental plaque [31]. While the majority of metagenome samples in the HMP dataset were from the tongue dorsum and supragingival plaque, HMP sampled additional oral sample types including subgingival plaque, the gingiva, buccal mucosa, hard palate, saliva, tonsils, and throat. 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 (Additional file 9: Table S7o, Additional file 8: Fig. SI4-SI8). TM7 populations that were abundant in supragingival plaque samples tended to be abundant in subgingival plaque samples and vice versa, and TM7 populations that were abundant on the tongue dorsum also tended to be abundant in saliva and in samples from tonsils and throat, results consistent with known patterns of microbial distribution within the mouth [11, 12]. Of particular note, the single MAG of clade P4 (group G5), which was previously suggested to associate with periodontitis [36] appeared to associate with subgingival plaque, but occurred similarly in subgingival plaque metagenomes of patients with periodontitis and healthy individuals (Additional file 9: Table S7p-s).
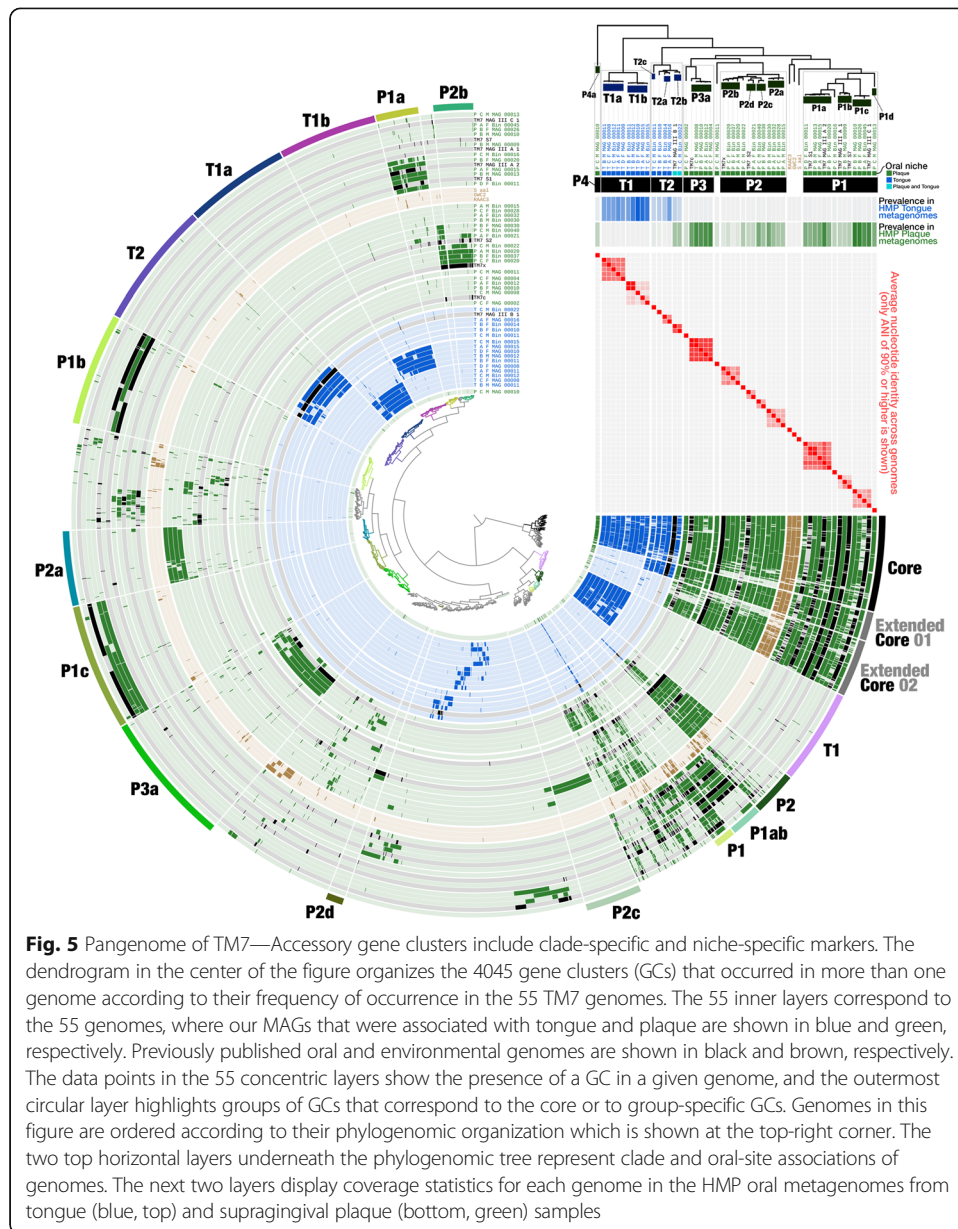
**Fig. 4** Detection and coverage of TM7 populations in the HMP plaque and tongue samples reveals abundant populations and niche specificity. The tree at the top of the figure and the two layers of information below it are identical to the one in Fig. 2. Barplots below the tree show the portion of plaque (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

Overall, these results confirm that the association of most TM7 oral populations with either dental plaque or tongue is a general feature and not restricted to the participants of our study. The MAGs that we constructed from 7 subjects represent prevalent and abundant tongue and plaque specialists.

### 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? We used a pangenomic approach to identify shared and unique genes and functions across the various TM7 clades and subclades in order to identify functional determinants of niche specificity. Our analysis organized the total of 40,832 genes across 55 genomes into 9117 gene clusters (GCs), each of which contains one or more homologous genes grouped on the basis of their amino acid similarity across genomes as judged by translated DNA sequences (not to be confused with operons or biosynthetic GCs). Of all GCs, 4045 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 (Fig. 5,

**Fig. 5** Pangenome of TM7—Accessory gene clusters include clade-specific and niche-specific markers. The dendrogram in the center of the figure organizes the 4045 gene clusters (GCs) 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 GC in a given genome, and the outermost circular layer highlights groups 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 two top horizontal layers underneath the phylogenomic tree represent clade and oral-site associations of genomes. The next two layers display coverage statistics for each genome in the HMP oral metagenomes from tongue (blue, top) and supragingival plaque (bottom, green) samples

Additional file 10: Table S8a-b). We used hierarchical clustering to group together the GCs that show similar distribution patterns across genomes (shown as the inner dendrogram in Fig. 5).

Our pangenomic analysis identified a collection of 205 core GCs that occurred in nearly all genomes, but we also identified groups of accessory GCs, many of which were exclusively associated with particular oral habitats or phylogenomic clades (Fig. 5, Additional file 4: Fig. S3a). Some groups of GCs were present only in plaque-associated clades, and some were present only in tongue-associated clades. Thus, the agreement between ecology and phylogenomics of these genomes was mirrored by the agreement between ecology and groups of differentially occurring GCs. The proportion of genes with functional hits varied dramatically between the core and accessory TM7 genes.
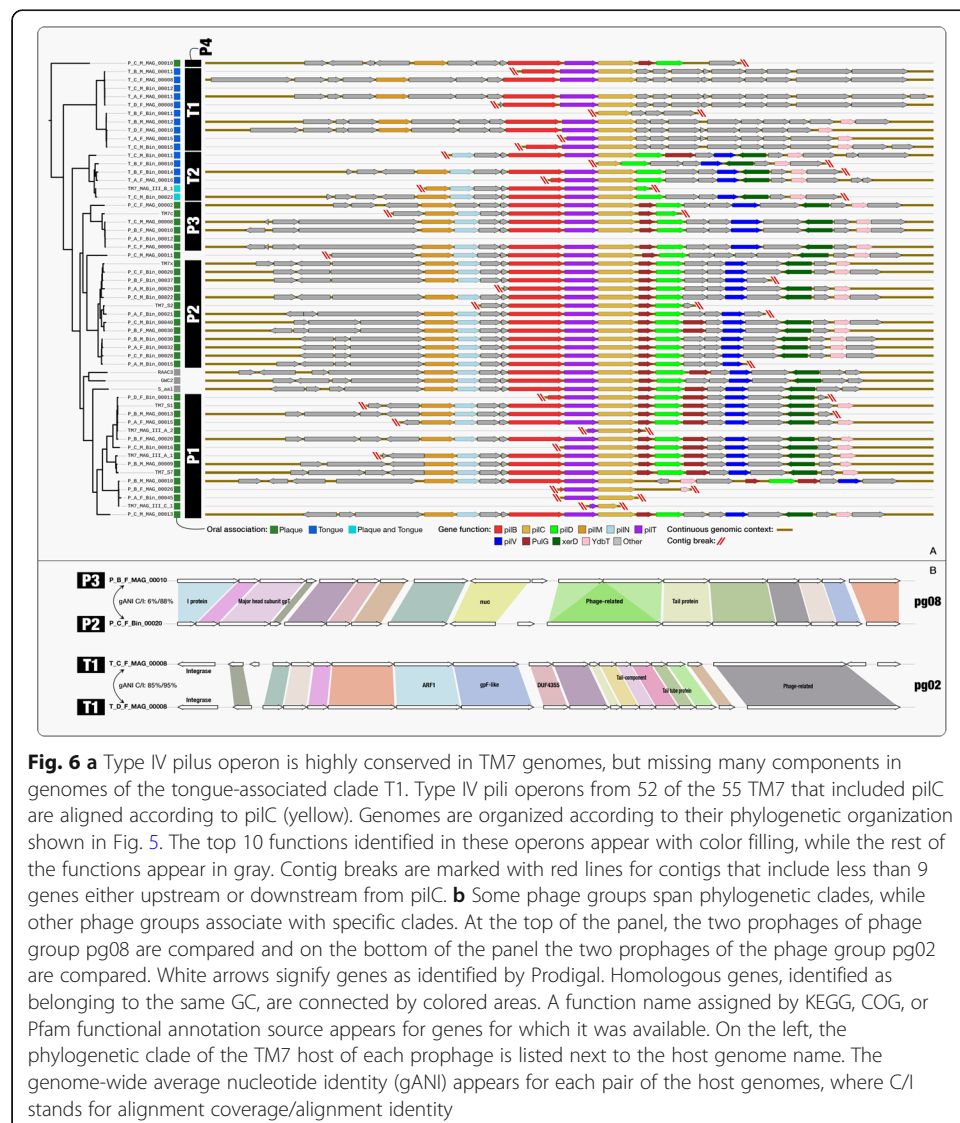
While more than 90% of core GCs had functional annotations, NCBI's Clusters of Orthologous Groups (COGs) database [55] annotated only 29% of singletons, and 22 to 88% of other groups of accessory GCs (Additional file 10: Table S8c), revealing a vast number of genes with unknown functions. Therefore, the nature and importance of many of these clade- and habitat-associated GCs is currently unknown.

Phylogenomics infers ancestral relationships as revealed by the evolution of conserved genes. Yet, genome phylogeny as predicted by conserved genes may not uniformly explain the distribution of all accessory genes across genomes, including genes that may be critical determinants of fitness against particular selective environmental pressures and influence the ecology of extant populations. In contrast to phylogenomics, pangenomics reveals associations between genomes based on the presence or absence of all genes across genomes [56]. Thus, pangenomics can more effectively capture ecological rather than phylogenetic similarities between genomes due to the strong influence of accessory genes on the estimation of associations [57]. To ask whether the distribution of GCs across genomes in the TM7 pangenome primarily reflected phylogeny or habitat, we also carried out a hierarchical clustering of genomes based on GCs and compared the resulting branching patterns (Additional file 4: Fig. S3b). The organization of TM7 genomes based on GCs predominantly matched their phylogenomic organization (Additional file 4: Fig. S3b); however, it recapitulated their niche association better than phylogenomics (Additional file 4: Fig. S3b). Specifically, the plaque-associated genome P_C_M_MAG_00010 of the clade P4 (group "G5"), which is a distant outlier to all other oral TM7 according to phylogeny (Fig. 2b), was placed together with all other plaque-associated TM7 as well as environmental TM7 (Additional file 4: Fig. S3b). GCs driving this organization belong to the "Extended Core 2" (Fig. 5). These GCs are generally characteristic of plaque and environmental TM7, but absent from tongue-associated TM7 of clades T1 and T2 (Fig. 5, Additional file 4: Fig. S3a, Additional file 10: Table S8c). 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 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 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.

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 (Additional file 10: Table S8d). The core GCs also included genes involved in amino acid transport. Since TM7s lack the genes to produce their own amino acids [31], 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 [33, 38], 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 [30] and were also specifically noted in TM7 genomes [37]. Type IV pilus systems are involved in many functions, including adherence [58], and could potentially be utilized by 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 (Fig. 6a). Additional copies of some of the type IV pilus proteins appear in various loci of the genome (Additional file 10: Table S8a). We found that while the cytosolic components of the type IV pilus system (PilT, PilB, PilC, PilM) were highly conserved across all genomes, components involved in the alignment of the system in the peptidoglycan (PilN) and the major and minor pilin proteins (PilE,



**Fig. 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 Fig. 5. The top 10 functions identified in these operons appear with color filling, while the rest of the functions appear in gray. 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 GC, 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

and PilV) appeared in clade or subclade-specific GCs and were completely absent from all genomes of clade T1 and from the single genome of clade P4 (Fig. 6a, Additional file 10: Table S8d). Variability in PilV has been shown in the past to confer binding specificity [59] 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 [60]. In total, 104 of the 207 proteins that were annotated with an LRR belonged to a single GC (GC_00000003) which was exclusively associated with T1 genomes, and each T1 genome had a total of 12–24 LRR proteins (COG4886, Table S8a). 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 [61]. 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 (Additional file 10: Table S8a). 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- or subclade-specific. This secretion system is also associated with motility in gram-positive bacteria [37] 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 52 out of 55 genomes. While this protein was flanked by genes involved in cell division in most genomes, in 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 among 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 GCs. Since GCs in a pangenome describe genes that are highly conserved in sequence space, identical functions can occur in distinct GCs, rendering it difficult to describe core and accessory functions in a pangenome based on core and accessory GCs. To address this issue, we developed a statistical approach for functional enrichment analyses in pangenomes (see "Material and methods") that employs logistic regression to reveal enriched functions in any given subset of genomes (e.g., those that match to a phylogenomic clade).

Applying this method to the TM7 pangenome, we identified 972 unique COGs assigned to gene clusters in the pangenome. Of these, 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 (Additional file 4: Fig. S3c and Additional file 10: Table S8q-r). While there was a wide overlap between core functions and core GCs, 131 core functions occurred in clade-specific GCs, suggesting that the genes encoding 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 "Extended Core 1" (the group of GCs characteristic of all genomes except tongue group 1) and one GC from "T1" (the group of GCs characteristic of tongue group 1 only) (Fig. 5). The large number of core functions with a T1-specific variant further demonstrates the uniqueness of clade T1 among the oral TM7 genomes (Additional file 4: Fig. S3c, Additional file 10: Table S8a). Other cases revealed additional functions that may have undergone selective pressures in a clade-specific manner. For example, a single copy of an RTX toxin-related Ca2+-binding protein (GC_00000221) was highly conserved in nearly all genomes 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. Among the 100 functions that were most enriched in one clade relative to other clades were many membrane-associated genes, including, but not limited to, functions that were highlighted above by our examination of GCs (Additional file 10: Table S8f). 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 among 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 genomes. In particular, our data reveal proteins with potential roles in adherence. Thus, 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 encoded functions and in core gene sequence composition as compared to other oral TM7, the functions encoded by members of clade T2, which includes the cosmopolitan T2_b genomes, appear to represent an intermediate between the strictly host-associated group and the plaque/environmental group. Finally, plaque specialists that are phylogenetically distinct resemble one another functionally and group together with environmental genomes based on GCs. By contrast, tongue specialists form a single host-associated superclade, suggesting that a stronger ecological similarity exists between teeth and non-host environments as compared to tongue.

## Mobile elements and prophages in TM7 genomes

Little evidence for phage association with members of the CPR has been found so far [62]. Dudek et al. recovered a phage associated with a TM7 genome from a dolphin plaque metagenome [63] and Paez-Espino et al. identified phages with a predicted SR1 host [64] in human oral metagenomes. A smaller genome size has been shown to correlate with the lack of lysogenic phages [65], 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 [66] and the recently published "inovirus detector" [67], along with a manual screening of gene functions to identify viral genes (see Additional file 8).

We identified 9 "phage groups" each composed of closely related prophages that were recovered from multiple TM7 genomes spanning all oral clades (Additional file 10: Table S8g). We did not identify any prophages in the three environmental genomes. Phage groups were generally associated with closely related hosts but were not restricted to hosts of the same TM7 species, or even the same oral clades (Fig. 6b, Additional file 10: Table S8g). 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 [68], 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 BLAST-searched 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 member of the CPR (Additional file 10: Table S8p). 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 (Additional file 10: Table S8a). 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 among oral TM7 and appear to be a unique feature of oral TM7, while absent from environmental TM7. In addition, CRISPR systems appear to be common among specific clades of oral TM7, but are 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 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 (Additional file 10: Table S8n). The varying location of the highly conserved transposases we identified in genomes of subclade T1_a suggests recent mobility, and that at least some of these elements are indeed active transposons (Additional file 10: Table S8a,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 closer hits to those from non-CPR bacteria.

### 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 [34]. Using full-length 16S rRNA sequences from clone libraries, Camanocha and 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) [34]. Genomes have been previously published for all of these taxa except SR1 HMT-875 [34, 42]. 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) 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 (Additional file 4: Fig. S4a, Additional file 4: Fig. S5a, Additional file 11: Table S9a-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 [38].

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 (Additional file 4: Fig. S4b-c, Additional file 11: Table S9l-n). Oral GN02 were all associated with plaque, and nearly absent from tongue samples (Additional file 4: Fig. S5b-c, Additional file 11: Table S9i-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 and Dewhirst [34] (Additional file 11: Table S9e-f). Overall, our GN02 and SR1 MAGs extend the collection of genomes available for these under-studied 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 10 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 each of the genera Mogibacterium, Propionibacterium, Leptotrichia, and Capnocytophaga, an unnamed genus in the family Flavobacteriaceae, an unnamed family within the class Clostridia, and finally unnamed families (potentially unnamed orders) within the classes Bacteroidia and Mollicutes (Fig. 1, Additional file 12: Table S10a-d, Additional file 8). Populations represented by these novel MAGs were

absent from skin and gut samples. Indeed, 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 [11, 69].

As the oral microbiome is highly represented in genomic databases [26, 27], the discovery of 10 unrepresented lineages seemed surprising. Hence, we next sought to check whether the absence of these novel MAGs from databases was 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 consistently recruited reads from both plaque and tongue samples, represented by the Flavobacteriaceae MAGs, T_A_M_MAG_00009 (Clostridiales), and three Capnocytophaga MAGs (Additional file 4: Fig. S6a). While we found some populations to be rare, which could explain their lack of genomic representation in databases, other populations were extremely prevalent; 22 of the 34 MAGs had at least 50% prevalence in at least one oral site and all but 1 had at least 10% prevalence (Additional file 4: Fig. S6a-c, Additional file 12: Table S10e-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 (Additional file 12: Table S10h).

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. These MAGs reached high relative abundance, recruiting up to 2.98% of the reads of a single metagenome, and recruiting on average 0.19% and 0.62% of tongue and plaque samples respectively (Additional file 12: Table S10e,g). To assess the likely taxonomic rank of these MAGs, we compared their ANI to one another and to representatives of all Flavobacteriaceae species in RefSeq. This analysis suggested that 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 (Additional file 12: Table S10i-j). A phylogenomic analysis placed these MAGs in a subgroup of Flavobacteriaceae together with Cloacibacterium, Chryseobacterium, Bergeyella, Riemerella, Cruoricaptor, Elizabethkingia, and Soonwooa (Additional file 4: Fig. S7). While all five Flavobacteriaceae MAGs had high sequence similarity, both ANI results and the phylogenetic analysis grouped these genomes according to the site of recovery, suggesting the existence of a plaque and tongue-specific subpopulation.

Approximately one third of human oral bacteria are not yet cultivable in the laboratory using standard approaches, perhaps because of dependence on other bacteria for metabolites or signals [26]. 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 may suggest a recent genomic reduction and possibly stronger 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

Our application of genome-resolved metagenomics to tongue and plaque metagenomes of seven individuals has provided new genomes for prevalent yet uncultivated members of the microbiome. In addition to making accessible 790 MAGs and their coverages across oral metagenomes in the form of reproducible anvi'o profiles (doi:https://doi.org/10.6084/m9.figshare.12217805), our study yields a set of 50 genomes (available via the NCBI accession PRJNA625082) that are of special interest to the oral microbiome community due to their novelty. Some of these genomes can inform future cultivation efforts as they suggest the existence of novel lineages with no cultured representatives within relatively well-studied microbial branches, while others, including those that resolve to TM7, SR1 and GN02, can enrich future comparative genomics and metagenomic read recruitment studies as they increase the known microbial gene pool, diversity, and mobilome of the more enigmatic branches of the human oral cavity.

Drivers of differentiation between the microbial members of distinct microbial niches within the human oral cavity remain largely unknown. However, integrated 'omics approaches that exploit rapidly increasing number of genomes promise unique opportunities. Our study demonstrates an integrated application of pangenomics and phylogenomics to a large number of new TM7 genomes and offers new insights into the diversity and evolution of this group. Our findings show that TM7 from the supragingival plaque group together with environmental TM7, while tongue-associated TM7 group together with lineages associated with animal gut. These results suggest that at least for TM7, supragingival plaque resembles non-host environments, while the tongue and gut TM7s are more strongly shaped by the host.

## 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 2 × 150 paired-end sequencing run using dedicated read indexing.

### Metagenomic assembly and processing of contigs

We used illumina-utils [70] 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 by Minoche et al. [71]. We then used MEGAHIT [72] 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 [73], v5.5.1, which (1) generated an anvi'o contigs database, (2) identified open reading frames using Prodigal [74] v2.60, (3) predicted gene-level taxonomy using Centrifuge [75], (4) annotated genes with functions using the NCBI's Clusters of Orthologous Groups (COG) [55], and (5) identified single-copy core genes using HMMER [76] 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 [77] to recruit short reads from the set of metagenomes used for co-assembly to the assembly product; we used samtools [78] 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 [79] 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 0.5 Mbp 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 [80], include contamination. In our refinement approach, we utilized the different clustering organizations available on the anvi'o interactive interface that rely on sequence composition and differential coverage across multiple metagenomes, taxonomic annotations of individual genes, and BLAST search results for gene sequences against public databases when necessary. For each MAG we have refined, we first identified contigs that showed anomalous coverage patterns across samples compared to the rest of the MAG using the program anvi-refine and visually inspected the nucleotide-level coverage and gene content of such contigs. If we could not confidently assign coverage variation to the presence of hypervariable genomic islands or prophages with matching taxonomy to the MAG, we conservatively eliminated them with the assumption that their inclusion in the MAG was likely due to binning errors. We performed two to three rounds of refinement per MAG: (1) initial refinement using the coverage information given the 4 to 6 metagenomes we have used to assemble our MAGs from each individual, (2) further refinement of 63 MAGs in which we identified likely contaminants based on their coverage across the full collection of 71 oral metagenomes in our study and then used the same coverage profiles to remove contaminating contigs, and (3) further refinement of 4 CPR and 3 non-CPR novel MAGs based on their coverage patterns across the HMP oral metagenomes, which we also used to refine TM7 genomes we have downloaded from the IMG. We carefully examined HOMD genomes and MAGs that formed deep phylogenetic branches and employed further steps of refinement when necessary to minimize the possibility of contamination-driven and artificially long phylogenetic branches. Some contamination may remain despite these efforts.

### 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) [42], and followed by a number, where for each co-assembly the MAGs had a series of numbers that start with "00001" and increment 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 [81] to calculate the average nucleotide identity (ANI) between each pair of MAGs that CheckM [82] assigned to the same phylum. To assign phylum affiliations to MAGs that had no phylum designation from CheckM, we used phylogenomics (see below) and complemented

phylogenomics with BLAST of protein sequences against the NCBI's non-redundant database. We assumed that a pair of MAGs were redundant if their ANI was at least 99.8% over the alignment between them that covers at least 50% of the shorter genome. 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 used the merged 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 [43] 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) [20] oral metagenomes, we used the program "anvi-run-workflow" with "--workflow metagenomics," which uses Snakemake [83] 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 [84] 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 2 Mbp 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 [85] 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 the anvi'o phylogenomics workflow with a set of 37 ribosomal proteins that occurred both in bacterial [42] and archaeal [86] 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, 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" 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 [87], (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 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 [42] to search for genes. For Fig. 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-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 [88] 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 [89] with the "WAG" general matrix model [90]. 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 our phylogenetic trees, we used an outlier genome in each analysis: for Fig. 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 intestini, Eubacterium rectale, Staphylococcus aureus, Streptococcus pneumoniae, Veillonella parvula.* In

Fig. 5, Additional file 4: Fig. S4a, Additional file 4: Fig. S5a, we removed Firmicutes after rooting using the Python package ete3 [91] version 3.1.1.

### 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 Additional file 10: Table S8a, our workflow also included running "anvi-run-pfams" to annotate functions with Pfams [92]. 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 [93].

### 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 to determine occurrence, and "--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 $\leq 1$ 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 [57], but briefly the pangenomic analysis used the NCBI's BLAST [94] to quantify similarity between each pair of genes, and the Markov Cluster algorithm (MCL) [95] (with inflation parameter of 2) to resolve clusters of 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."

### Computing average nucleotide identity (ANI)

We used "anvi-compute-ani" with parameters "--method ANIm" to align genomes using MUMmer [96] and "--min-alignment-fraction 0.25" to only keep scores if the alignment fraction covers at least 25% of both genomes. For the ANI data presented in Figs. 3 and 5, we first computed ANI without the flag "--min-alignment-fraction" to get all alignment statistics (which we report in Additional file 9: Table S7f-g), and then we imported ANI values into the anvi'o profile database 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 an additional 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 each of the individuals C_F and C_M. From individual L, we collected 13 daily tongue dorsum scrapings (Table S1d). 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: 24 U/mL) and mutanolysin (by Sigma-Aldrich; final concentration: 0.3KU/mL) and extended each incubation step to 2 h 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 to 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 h. 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 a $Q$-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 [97]. 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. We then searched 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 Additional file 9: Table S7i for oral genomes and in Table S7e for all TM7 downloaded from NCBI.

### Functional enrichment analysis

We developed 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 [98]. 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" [99]. To apply functional enrichment analysis to our pangenome, 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 [66] and the "Inovirus detector" [67] to identify contigs that include putative phage sequences. We manually inspected contigs predicted as viral and excluded 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). We further examined all remaining contigs to verify their placement in the prospective genomes, using the data in Additional file 10: Table S8a, as well as BLAST searches of protein sequences (see the notes in Table S8g for more details). We used functional annotations (from NCBI COGs, Pfam, and GhostKOALA) 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 contigs that we already identified as phage. 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, hence confirming the insertion site of the prophage.

### Identifying CRISPRs

We used the web service CRISPRCasFinder [68] to search for CRISPR spacers in the 55 TM7 genomes. Along with a summary of the results (Additional file 10: Table S8l), 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. [68]).

### 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 (CCGTCAATTC NTTTRAGT 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 6-

nt index. We generated PCR amplicons in triplicate 33 μL reaction volumes with an amplification cocktail containing 0.67 U SuperFi Taq Polymerase (Invitrogen, Carlsbad, CA), 1× enzyme buffer (includes MgCl2), 200 μM dNTP PurePeak DNA polymerase mix (Thermo Fisher), and 0.3 μM 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 as follows: initial 94 °C, 3 min denaturation step; 30 cycles of 94 °C for 30 s, 57 °C for 45 s, and 72 °C for 60 s; final 2 min extension at 72 °C. 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 [100] to identify amplicon sequence variants (ASVs) across samples and determine the microbial community structure, and Global Alignment for Sequence Taxonomy (GAST) [101] to assign taxonomic affiliation to each ASV.

### Statistics and visualization

We used ggplot2 [102] v3.2.1 to generate boxplots and barplots. 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 [103]. 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 Additional file 7: Table S6b, which also includes the archeon *Methanobrevibacter oralis* that was used to root the phylogeny in Fig. 1. While the TM7 genomes we downloaded from IMG had no accession numbers, refined versions of these genomes have recently been published [38]. 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 Additional file 9: Table S7e); GN02 and SR1 on 12/17/2018 (accession numbers provided in Additional file 11:Table S9a,b); and Flavobacteriaceae on 9/20/2019. We obtained the raw metagenomes from Califf et al. [84] directly from the authors since the published FASTQ files were missing the second pair. We downloaded Firmicutes genomes to root the trees in Figs. 2, 3, and 4 from the NCBI's RefSeq database (GCF_000147095.1, GCF_000210315.1, GCF_000024945.1, GCF_000020605.1, GCF_000230275.1).

### Supplementary information

The online version contains supplementary material available at https://doi.org/10.1186/s13059-020-02195-w.

---

**Additional file 1: Table S1.** Details of collected samples, including statistics and accession numbers 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.

**Additional file 2: Table S2.** 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 2463 genomic bins.

**Additional file 3: Table S3.** 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 redundant MAGs to one of 59 redundant groups.

**Additional file 4.** Supplementary figures S1, S2, S3, S4, S5, S6, and S7.

**Additional file 5. Table S4.** 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 addition to relative abundance estimations for TM7 and the collection of reads not assigned any genus affiliation.

**Additional file 6: Table S5.** 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 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.

**Additional file 7: Table S6.** Description of Human Oral Microbiome Database (HOMD) genomes used in this study. (a) Summary statistics for contigs for each of 1334 HOMD genomes used in this study, including total length, number of contigs, and count of SCGs. (b) accession number information for each of the 1334 HOMD genomes.

**Additional file 8.** Additional discussion and supplementary figures for (A) taxonomy based on metagenome-assembled genomes, metagenomics short reads, and 16S rRNA gene amplicons; (B) Average Nucleotide Identity (ANI) of oral TM7; (C) occurrence of TM7 across additional oral sample types, other than supragingival plaque and tongue dorsum, and including samples from patients with periodontitis; (D) mobile elements and prophages in TM7 genomes; (E) novel non-CPR MAGs; (F) a novel MAG that represents a member of the Mollicutes; (G) novel Clostridiales MAGs represent prevalent tongue-associated populations; and (H) novel Bacteroidia MAGs that include a tongue-specialist and a subgingival plaque specialist.

**Additional file 9: Table S7.** 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 150 TM7 genomes from NCBI used for the phylogenomic analysis in Fig. 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 [34]) 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 and (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 breakdown for plaque and tongue samples, the percent of HMP samples in which the genome was above detection threshold, and 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 [84]. (q) Non-outlier mean coverage and (r) detection summary statistics for 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).

**Additional file 10: Table S8.** 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).

**Additional file 11: Table S9.** 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.

**Additional file 12: Table S10.** 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 and (g) non-outlier mean coverage of the novel MAGs in the HMP samples as computerd 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.

**Additional file 13.** Review history.

## Review history

The review history is available as Additional file 13.

## Peer review information

Kevin Pang was the primary editor of this article and managed its editorial process and peer review in collaboration with the rest of the editorial team.

## Authors' contributions

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

## Authors' information

Twitter handles:
• Alon Shaiber @alon_shaiber
• Amy D. Willis @AmyDWillis
• Tom O. Delmont @tomodelmont
• Simon Roux @simroux_virus
• Abigail C. Schmid @SchmidAbby
• Andrea R. Watson @awatsontxt
• Sonny T. M. Lee @sonnytmlee
• Jessica L. Mark Welch @JMarkWelch
• A. Murat Eren @merenbey

## Availability of data and materials

We deposited the short-read sequencing data for amplicons and shotgun metagenomes as well as key MAGs that emerged from our study under the NCBI BioProject PRJNA625082 (Additional file 1: Table S1 lists accession numbers for each sample individually) [104]. 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) [104]. For reproducibility, reusability, and transparency, we also have made available the FASTA files for co-assembled metagenomes (doi:https://doi.org/10.6084/m9.figshare.12217799) and the anvi'o merged profile databases (doi:https://doi.org/10.6084/m9.figshare.12217802) for each individual in our study [105]. Anvi'o split profiles for each of the 790 MAGs (doi:https://doi.org/10.6084/m9.figshare.12217805) and the TM7 pangenome (doi:https://doi.org/10.6084/m9.figshare.12217811) are also

Shaiber *et al. Genome Biology*        (2020) 21:292

Page 32 of 35

publicly available [105]. Finally, doi:https://doi.org/10.6084/m9.figshare.11634321 gives access to all Supplementary Tables, Supplementary Figures, and the Supplementary Information file.

### Ethics approval and consent to participate
We recruited human subjects and collected samples according to protocol #15-247 as approved by New England IRB (Needham, Massachusetts, USA) and in accordance with the Declaration of Helsinki. All subjects provided informed consent prior to sampling.

### Consent for publication
Not applicable.

### Competing interests
The authors declare that they have no competing interests.

### Author details
[1]Department of Medicine, University of Chicago, Chicago, IL 60637, USA. [2]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]Computational and Applied Mathematics, University of Chicago, Chicago, IL 60637, USA. [8]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.

### References
1. Lane N. The unseen world: reflections on Leeuwenhoek (1677) "Concerning little animals.". Philos Trans R Soc Lond B Biol Sci R Soc. 2015;370:20140344.
2. German RZ, Palmer JB. Anatomy and development of oral cavity and pharynx. GI Motility online [Internet]: Nature Publishing Group; 2006. [cited 2019 Nov 28]. Available from: https://doi.org/10.1038/gimo5.
3. Moutsopoulos NM, Konkel JE. Tissue-specific immunity at the oral mucosal barrier. Trends Immunol. 2018;39:276–87.
4. Song SJ, Lauber C, Costello EK, Lozupone CA, Humphrey G, Berg-Lyons D, et al. Cohabiting family members share microbiota with one another and with their dogs. Elife. 2013;2:e00458.
5. Ferretti P, Pasolli E, Tett A, Asnicar F, Gorfer V, Fedi S, et al. Mother-to-infant microbial transmission from different body sites shapes the developing infant gut microbiome. Cell Host Microbe. 2018;24:133–45.e5.
6. Caporaso JG, Lauber CL, Costello EK, Berg-Lyons D, Gonzalez A, Stombaugh J, et al. Moving pictures of the human microbiome. Genome Biol. 2011;12:R50.
7. Mark Welch JL, Utter DR, Rossetti BJ, Mark Welch DB, Eren AM, Borisy GG. Dynamics of tongue microbial communities with single-nucleotide resolution using oligotyping. Front Microbiol. 2014;5:568.
8. Hall MW, Singh N, Ng KF, Lam DK, Goldberg MB, Tenenbaum HC, et al. Inter-personal diversity and temporal dynamics of dental, tongue, and salivary microbiota in the healthy oral cavity. NPJ Biofilms Microbiomes. 2017;3:2.
9. Mark Welch JL, Rossetti BJ, Rieken CW, Dewhirst FE, Borisy GG. Biogeography of a human oral microbiome at the micron scale. Proc Natl Acad Sci U S A. 2016;113:E791–800.
10. Dewhirst FE, Chen T, Izard J, Paster BJ, Tanner ACR, Yu W-H, et al. The human oral microbiome. J Bacteriol. 2010;192: 5002–17.
11. Eren AM, Borisy GG, Huse SM, Mark Welch JL. Oligotyping analysis of the human oral microbiome. Proc Natl Acad Sci U S A. 2014;111:E2875–84.
12. Mark Welch JL, Dewhirst FE, Borisy GG. Biogeography of the oral microbiome: the site-specialist hypothesis. Annu Rev Microbiol. 2019;73:335–58.
13. Lamont RJ, Koo H, Hajishengallis G. The oral microbiota: dynamic communities and host interactions. Nat Rev Microbiol. 2018;16:745–59.
14. Schmidt TS, Hayward MR, Coelho LP, Li SS, Costea PI, Voigt AY, et al. Extensive transmission of microbes along the gastrointestinal tract. Elife. 2019;8 Available from: https://doi.org/10.7554/eLife.42693.
15. Socransky SS, Manganiello SD. The oral microbiota of man from birth to senility. J Periodontol. 1971;42:485–96.
16. Mager DL, Ximenez-Fyvie LA, Haffajee AD, Socransky SS. Distribution of selected bacterial species on intraoral surfaces. J Clin Periodontol. 2003;30:644–54.
17. Aas JA, Paster BJ, Stokes LN, Olsen I, Dewhirst FE. Defining the normal bacterial flora of the oral cavity. J Clin Microbiol. 2005;43:5721–32.
18. Gibbons RJ, Houte JV. Bacterial adherence in oral microbial ecology. Annu Rev Microbiol. 1975;29:19–44.
19. Simón-Soro A, Tomás I, Cabrera-Rubio R, Catalan MD, Nyvad B, Mira A. Microbial geography of the oral cavity. J Dent Res. 2013;92:616–21.
20. Human Microbiome Project Consortium. Structure, function and diversity of the healthy human microbiome. Nature. 2012;486:207–14.
21. Segata N, Haake SK, Mannon P, Lemon KP, Waldron L, Gevers D, et al. Composition of the adult digestive tract bacterial microbiome based on seven mouth surfaces, tonsils, throat and stool samples. Genome Biol. 2012;13:R42.

22. Lloyd-Price J, Mahurkar A, Rahnavard G, Crabtree J, Orvis J, Hall AB, et al. Strains, functions and dynamics in the expanded Human Microbiome Project. Nature. 2017;550:61–6.
23. Donati C, Zolfo M, Albanese D, Tin Truong D, Asnicar F, Iebba V, et al. Uncovering oral Neisseria tropism and persistence using metagenomic sequencing. Nat Microbiol. 2016;1:16070.
24. Chen T, Yu W-H, Izard J, Baranova OV, Lakshmanan A, Dewhirst FE. The Human Oral Microbiome Database: a web accessible resource for investigating oral microbe taxonomic and genomic information. Database. 2010;2010:baq013.
25. Escapa IF, Chen T, Huang Y, Gajare P, Dewhirst FE, Lemon KP. 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. 2018;3 Available from: https://doi.org/10.1128/mSystems.00187-18.
26. Nayfach S, Rodriguez-Mueller B, Garud N, Pollard KS. An integrated metagenomics pipeline for strain profiling reveals novel patterns of bacterial transmission and biogeography. Genome Res. 2016;26:1612–25.
27. Vartoukian SR, Adamowska A, Lawlor M, Moazzez R, Dewhirst FE, Wade WG. In vitro cultivation of "unculturable" oral bacteria, facilitated by community culture and media supplementation with siderophores. PLoS One. 2016;11:e0146926.
28. Brown CT, Hug LA, Thomas BC, Sharon I, Castelle CJ, Singh A, et al. Unusual biology across a group comprising more than 15% of domain Bacteria. Nature. 2015;523:208–11.
29. Hug LA, Baker BJ, Anantharaman K, Brown CT, Probst AJ, Castelle CJ, et al. A new view of the tree of life. Nat Microbiol. 2016;1:16048. Macmillan Publishers Limited.
30. Méheust R, Burstein D, Castelle CJ, Banfield JF. The distinction of CPR bacteria from other bacteria based on protein family content. Nat Commun. 2019;10:4173.
31. Bor B, Bedree JK, Shi W, McLean JS, He X. Saccharibacteria (TM7) in the human oral microbiome. J Dent Res. 2019;98:500–9.
32. Bor B, McLean JS, Foster KR, Cen L, To TT, Serrato-Guillen A, et al. Rapid evolution of decreased host susceptibility drives a stable relationship between ultrasmall parasite TM7x and its bacterial host. Proc Natl Acad Sci U S A. 2018;115:12277–82.
33. He X, McLean JS, Edlund A, Yooseph S, Hall AP, Liu S-Y, et al. Cultivation of a human-associated TM7 phylotype reveals a reduced genome and epibiotic parasitic lifestyle. Proc Natl Acad Sci U S A. 2015;112:244–9.
34. 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 Available from: https://doi.org/10.3402/jom.v6.25468. Epub 2014/10/16. PMID: 25317252.
35. Brinig MM, Lepp PW, Ouverney CC, Armitage GC, Relman DA. Prevalence of bacteria of division TM7 in human subgingival plaque and their association with disease. Appl Environ Microbiol. 2003;69:1687–94.
36. Abusleme L, Dupuy AK, Dutzan N, Silva N, Burleson JA, Strausbaugh LD, et al. The subgingival microbiome in health and periodontitis and its relationship with community biomass and inflammation. ISME J. 2013;7:1016–25.
37. Marcy Y, Ouverney C, Bik EM, Lösekann T, Ivanova N, Martin HG, et al. Dissecting biological "dark matter" with single-cell genetic analysis of rare and uncultivated TM7 microbes from the human mouth. Proc Natl Acad Sci U S A. 2007;104: 11889–94.
38. Cross KL, Campbell JH, Balachandran M, Campbell AG, Cooper SJ, Griffen A, et al. Targeted isolation and cultivation of uncultivated bacteria by reverse genomics. Nat Biotechnol. 2019;37:1314–21.
39. Murugkar PP, Collins AJ, Dewhirst FE. Complete genome sequence of strain PM004, a novel cultured member of the human oral microbiome from the candidate phylum Saccharibacteria (TM7). Microbiol Resour Announc. 2019;8 Available from: https://doi.org/10.1128/MRA.01159-19.
40. Collins AJ, Murugkar PP, Dewhirst FE. Complete genome sequence of strain AC001, a novel cultured member of the human oral microbiome from the candidate phylum Saccharibacteria (TM7) [Internet]. Microbiol Res Announcements. 2019; Available from: https://doi.org/10.1128/mra.01158-19.
41. Espinoza JL, Harkins DM, Torralba M, Gomez A, Highlander SK, Jones MB, et al. Supragingival plaque microbiome ecology and functional potential in the context of health and disease. MBio. 2018;9 Available from: https://doi.org/10.1128/mBio.01631-18.
42. Campbell JH, O'Donoghue P, Campbell AG, Schwientek P, Sczyrba A, Woyke T, et al. UGA is an additional glycine codon in uncultured SR1 bacteria from the human microbiota. Proc Natl Acad Sci U S A. 2013;110:5540–5.
43. Chen L-X, Anantharaman K, Shaiber A, Murat Eren A, Banfield JF. Accurate and complete genomes from metagenomes [Internet]. bioRxiv. 2019:808410. [cited 2019 Dec 8]. Available from: https://www.biorxiv.org/content/10.1101/808410v1.
44. Raveh-Sadka T, Thomas BC, Singh A, Firek B, Brooks B, Castelle CJ, et al. Gut bacteria are rarely shared by co-hospitalized premature infants, regardless of necrotizing enterocolitis development. Elife. 2015;4 Available from: https://doi.org/10.7554/eLife.05477.
45. Delmont TO, Quince C, Shaiber A, Esen ÖC, Lee ST, Rappé MS, et al. Nitrogen-fixing populations of Planctomycetes and Proteobacteria are abundant in surface ocean metagenomes. Nat Microbiol. 2018;3:804–13.
46. Shaiber A, Eren AM. Composite metagenome-assembled genomes reduce the quality of public genome repositories. MBio. 2019; Available from: https://doi.org/10.1128/mBio.00725-19.
47. Quince C, Walker AW, Simpson JT, Loman NJ, Segata N. Shotgun metagenomics, from sampling to analysis. Nat Biotechnol. 2017;35:833–44.
48. Chen I-MA, Chu K, Palaniappan K, Pillay M, Ratner A, Huang J, et al. IMG/M v.5.0: an integrated data management and comparative analysis system for microbial genomes and microbiomes. Nucleic Acids Res. 2019;47:D666–77.
49. McLean JS, Bor B, To TT, Liu Q, Kearns KA, Solden LM, et al. Evidence of independent acquisition and adaption of ultra-small bacteria to human hosts across the highly diverse yet reduced genomes of the phylum Saccharibacteria [Internet]. 2018. Available from: https://doi.org/10.1101/258137.
50. Parks DH, Rinke C, Chuvochina M, Chaumeil P-A, Woodcroft BJ, Evans PN, et al. Recovery of nearly 8,000 metagenome-assembled genomes substantially expands the tree of life. Nat Microbiol. 2017;2:1533–42.
51. Schluter J, Foster KR. The evolution of mutualism in gut microbiota via host epithelial selection. PLoS Biol. 2012;10:e1001424.
52. Kantor RS, Wrighton KC, Handley KM, Sharon I, Hug LA, Castelle CJ, et al. Small genomes and sparse metabolisms of sediment-associated bacteria from four candidate phyla. MBio. 2013;4:e00708–13.
53. Wilbert SA, Mark Welch JL, Borisy GG. Spatial ecology of the human tongue dorsum microbiome. Cell Rep. 2020;30: 4003–15.e3.

54. Albertsen M, Hugenholtz P, Skarshewski A, Nielsen KL, Tyson GW, Nielsen PH. Genome sequences of rare, uncultured bacteria obtained by differential coverage binning of multiple metagenomes. Nat Biotechnol. 2013;31:533–8.
55. Tatusov RL, Fedorova ND, Jackson JD, Jacobs AR, Kiryutin B, Koonin EV, et al. The COG database: an updated version includes eukaryotes. BMC Bioinformatics. 2003;4:41.
56. Dutilh BE, Huynen MA, Bruno WJ, Snel B. The consistent phylogenetic signal in genome trees revealed by reducing the impact of noise. J Mol Evol. 2004;58:527–39.
57. Delmont TO, Eren AM. Linking pangenomes and metagenomes: the Prochlorococcus metapangenome. PeerJ. 2018;6:e4320.
58. Craig L, Forest KT, Maier B. Type IV pili: dynamics, biophysics and functional consequences. Nat Rev Microbiol. 2019;17: 429–40.
59. Ishiwa A, Komano T. Thin pilus PilV adhesins of plasmid R64 recognize specific structures of the lipopolysaccharide molecules of recipient cells. J Bacteriol. 2003;185:5192–9.
60. Bella J, Hindle KL, McEwan PA, Lovell SC. The leucine-rich repeat structure. Cell Mol Life Sci. 2008;65:2307–33.
61. Buist G, Steen A, Kok J, Kuipers OP. LysM, a widely distributed protein motif for binding to (peptido)glycans. Mol Microbiol. 2008;68:838–47.
62. Chen L-X, Al-Shayeb B, Méheust R, Li W-J, Doudna JA, Banfield JF. Candidate phyla radiation Roizmanbacteria from hot springs have novel and unexpectedly abundant CRISPR-Cas systems. Front Microbiol. 2019;10:928.
63. Dudek NK, Sun CL, Burstein D, Kantor RS, Aliaga Goltsman DS, Bik EM, et al. Novel microbial diversity and functional potential in the marine mammal oral microbiome. Curr Biol. 2017;27:3752–62.e6.
64. Paez-Espino D, Eloe-Fadrosh EA, Pavlopoulos GA, Thomas AD, Huntemann M, Mikhailova N, et al. Uncovering Earth's virome. Nature. 2016;536:425–30.
65. Touchon M, Bernheim A, Rocha EP. Genetic and life-history traits associated with the distribution of prophages in bacteria. ISME J. 2016;10:2744–54.
66. Roux S, Enault F, Hurwitz BL, Sullivan MB. VirSorter: mining viral signal from microbial genomic data. PeerJ. 2015;3:e985.
67. Roux S, Krupovic M, Daly RA, Borges AL, Nayfach S, Schulz F, et al. Cryptic inoviruses revealed as pervasive in bacteria and archaea across Earth's biomes. Nat Microbiol. 2019; Available from: https://doi.org/10.1038/s41564-019-0510-x.
68. Couvin D, Bernheim A, Toffano-Nioche C, Touchon M, Michalik J, Néron B, et al. CRISPRCasFinder, an update of CRIS RFinder, includes a portable version, enhanced performance and integrates search for Cas proteins. Nucleic Acids Res. 2018;46:W246–51.
69. Franzosa EA, Morgan XC, Segata N, Waldron L, Reyes J, Earl AM, et al. Relating the metatranscriptome and metagenome of the human gut. Proc Natl Acad Sci U S A. 2014;111:E2329–38.
70. Eren AM, Vineis JH, Morrison HG, Sogin ML. A filtering method to generate high quality short reads using illumina paired-end technology. PLoS One. 2013;8:e66643.
71. Minoche AE, Dohm JC, Himmelbauer H. Evaluation of genomic high-throughput sequencing data generated on Illumina HiSeq and genome analyzer systems. Genome Biol. 2011;12:R112.
72. Li D, Liu C-M, Luo R, Sadakane K, Lam T-W. MEGAHIT: an ultra-fast single-node solution for large and complex metagenomics assembly via succinct de Bruijn graph. Bioinformatics. 2015;31:1674–6.
73. Eren AM, Esen ÖC, Quince C, Vineis JH, Morrison HG, Sogin ML, et al. Anvi'o: an advanced analysis and visualization platform for 'omics data. PeerJ. 2015;3:e1319.
74. Hyatt D, Chen G-L, Locascio PF, Land ML, Larimer FW, Hauser LJ. Prodigal: prokaryotic gene recognition and translation initiation site identification. BMC Bioinformatics. 2010;11:119.
75. Kim D, Song L, Breitwieser FP, Salzberg SL. Centrifuge: rapid and sensitive classification of metagenomic sequences. Genome Res. 2016;26:1721–9.
76. Eddy SR. Accelerated profile HMM searches. PLoS Comput Biol. 2011;7:e1002195.
77. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. Nat Methods. 2012;9:357–9.
78. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence Alignment/Map format and SAMtools. Bioinformatics. 2009;25:2078–9.
79. Alneberg J, Bjarnason BS, de Bruijn I, Schirmer M, Quick J, Ijaz UZ, et al. CONCOCT: Clustering cONtigs on COverage and ComposiTion [Internet]. arXiv [q-bio.GN]. 2013; Available from: http://arxiv.org/abs/1312.4038.
80. McLean JS, Bor B, To TT, Liu Q, Kearns KA, Solden LM. Evidence of independent acquisition and adaption of ultra-small bacteria to human hosts across the highly diverse yet reduced genomes of the phylum. bioRxiv [Internet]. biorxiv.org. 2018; Available from: https://www.biorxiv.org/content/early/2018/02/02/258137.abstract.
81. Delcher AL, Phillippy A, Carlton J, Salzberg SL. Fast algorithms for large-scale genome alignment and comparison. Nucleic Acids Res. 2002;30:2478–83.
82. Parks DH, Imelfort M, Skennerton CT, Hugenholtz P, Tyson GW. CheckM: assessing the quality of microbial genomes recovered from isolates, single cells, and metagenomes. Genome Res. 2015;25:1043–55.
83. Köster J, Rahmann S. Snakemake--a scalable bioinformatics workflow engine. Bioinformatics. 2012;28:2520–2.
84. Califf KJ, Schwarzberg-Lipson K, Garg N, Gibbons SM, Gregory Caporaso J, Slots J, et al. Multi-omics analysis of periodontal pocket microbial communities pre- and posttreatment [Internet]. mSystems. 2017; Available from: https://doi.org/10.1128/msystems.00016-17.
85. Breitwieser FP, Salzberg SL, Kraken HLL. Confident and fast metagenomics classification using unique k-mer counts [Internet]. bioRxiv. 2018:262956. [cited 2018 Aug 9]. Available from: https://www.biorxiv.org/content/early/2018/02/09/262956.
86. Rinke C, Schwientek P, Sczyrba A, Ivanova NN, Anderson IJ, Cheng J-F, et al. Insights into the phylogeny and coding potential of microbial dark matter. Nature. 2013;499:431–7.
87. Deorowicz S, Debudaj-Grabysz A, Gudyś A. FAMSA: Fast and accurate multiple sequence alignment of huge protein families. Sci Rep. 2016;6:33964.
88. Capella-Gutiérrez S, Silla-Martínez JM, Gabaldón T, et al. Bioinformatics. 2009;25:1972–3.
89. Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ. IQ-TREE: a fast and effective stochastic algorithm for estimating maximum-likelihood phylogenies. Mol Biol Evol. 2015;32:268–74.
90. Whelan S, Goldman N. A general empirical model of protein evolution derived from multiple protein families using a maximum-likelihood approach. Mol Biol Evol. 2001;18:691–9.

91.  Huerta-Cepas J, Serra F, Bork P. ETE 3: reconstruction, analysis, and visualization of phylogenomic data. Mol Biol Evol. 2016;33:1635–8.
92.  El-Gebali S, Mistry J, Bateman A, Eddy SR, Luciani A, Potter SC, et al. The Pfam protein families database in 2019. Nucleic Acids Res. 2019;47:D427–32.
93.  Kanehisa M, Sato Y, Morishima K. BlastKOALA and GhostKOALA: KEGG tools for functional characterization of genome and metagenome sequences. J Mol Biol. 2016;428:726–31.
94.  Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool. J Mol Biol. 1990;215:403–10.
95.  Enright AJ, Van Dongen S, Ouzounis CA. An efficient algorithm for large-scale detection of protein families. Nucleic Acids Res. 2002;30:1575–84.
96.  Marçais G, Delcher AL, Phillippy AM, Coston R, Salzberg SL, Zimin A. MUMmer4: A fast and versatile genome alignment system. PLoS Comput Biol. 2018;14:e1005944.
97.  Li H. Minimap2: pairwise alignment for nucleotide sequences. Bioinformatics. 2018;34:3094–100.
98.  R Development Core Team R. R: A Language and Environment for Statistical Computing [Internet]. 2011. Available from: http://www.r-project.org.
99.  Storey JD, Tibshirani R. Statistical significance for genomewide studies. Proc Natl Acad Sci U S A. 2003;100:9440–5.
100. Eren AM, Morrison HG, Lescault PJ, Reveillaud J, Vineis JH, Sogin ML. Minimum entropy decomposition: unsupervised oligotyping for sensitive partitioning of high-throughput marker gene sequences. ISME J. 2014;9:968–79.
101. Huse SM, Dethlefsen L, Huber JA, Mark Welch D, Welch DM, Relman DA, et al. Exploring microbial diversity and taxonomy using SSU rRNA hypervariable tag sequencing. PLoS Genet. 2008;4:e1000255. Josephine Bay Paul Center for Comparative Molecular Biology and Evolution, Marine Biological Laboratory, Woods Hole, Massachusetts, United States of America.
102. Ginestet C. ggplot2: elegant graphics for data analysis. J R Stat Soc Ser A Stat Soc. 2011;174:245–6.
103. Seabold S, Perktold J. Statsmodels: econometric and statistical modeling with Python. In Proceedings of the 9th Python in Science Conference. 2010. p. 57–61.
104. Shaiber A, Mark Welch JL, Eren AM. Sequence data for amplicons, long-reads, shotgun metagenomes, and metagenome-assembled genomes from human oral samples. National Center for Biotechnology Information. 2020. https://www.ncbi.nlm.nih.gov/bioproject/?term=PRJNA625082.
105. Shaiber A, Eren AM. Anvi'o reusable data currencies for genome-resolved metagenomics of the human oral cavity. FigShare. 2020. https://figshare.com/projects/Data_for_Oral_TM7s_by_Shaiber_et_al_2020_/74301.

## Publisher's Note