

Research



Cite this article: Achtman M, Zhou Z. 2020 Metagenomics of the modern and historical human oral microbiome with phylogenetic studies on *Streptococcus mutans* and *Streptococcus sobrinus*. *Phil. Trans. R. Soc. B* **375**: 20190573.
<http://dx.doi.org/10.1098/rstb.2019.0573>

Accepted: 9 July 2020

One contribution of 14 to a theme issue 'Insights into health and disease from ancient biomolecules'.

Subject Areas:

bioinformatics, microbiology, genomics, health and disease and epidemiology, evolution, genetics

Keywords:

ancient DNA, dental plaque, dental calculus, saliva, genomic reconstruction, metagenomes

Author for correspondence:

Mark Achtman
e-mail: m.achtman@warwick.ac.uk

Electronic supplementary material is available online at <https://doi.org/10.6084/m9.figshare.c.5095997>.

Metagenomics of the modern and historical human oral microbiome with phylogenetic studies on *Streptococcus mutans* and *Streptococcus sobrinus*

Mark Achtman and Zheming Zhou

Warwick Medical School, University of Warwick, Coventry CV4 7AL, UK

MA, 0000-0001-6815-0070; ZZ, 0000-0001-9783-0366

We have recently developed bioinformatic tools to accurately assign metagenomic sequence reads to microbial taxa: SPARSE for probabilistic, taxonomic classification of sequence reads; EToKi for assembling and polishing genomes from short-read sequences; and GrapeTree, a graphic visualizer of genetic distances between large numbers of genomes. Together, these methods support comparative analyses of genomes from ancient skeletons and modern humans. Here, we illustrate these capabilities with 784 samples from historical dental calculus, modern saliva and modern dental plaque. The analyses revealed 1591 microbial species within the oral microbiome. We anticipated that the oral complexes of Socransky *et al.*, which were defined in 1998, would predominate among taxa whose frequencies differed by source. However, although some species discriminated between sources, we could not confirm the existence of the complexes. The results also illustrate further functionality of our pipelines with two species that are associated with dental caries, *Streptococcus mutans* and *Streptococcus sobrinus*. They were rare in historical dental calculus but common in modern plaque, and even more common in saliva. Reconstructed draft genomes of these two species from metagenomic samples in which they were abundant were combined with modern public genomes to provide a detailed overview of their core genomic diversity.

This article is part of the theme issue 'Insights into health and disease from ancient biomolecules'.

1. Introduction

Multiple research areas have undergone revolutionary changes in the last 10 years due to broad accessibility to high-throughput DNA sequencing at reduced costs. These include the evolutionary biology of microbial pathogens based on metagenomic sequencing. Studies on *Mycobacterium tuberculosis* [1,2], *Mycobacterium leprae* [3,4], *Yersinia pestis* [5–10] and *Salmonella enterica* [11–13] have yielded important insights into the history of infectious diseases by combining modern and historical genomes. In principle, the same approach might also help to elucidate the evolutionary history of both commensal and pathogenic taxa within the human oral microbiome. Periodontitis and dental caries have likely afflicted humans since their origins [14–17]. They may now be amenable to population genetic analyses because a landmark publication by Adler *et al.* in 2013 [18] demonstrated that dental calculus (calcified dental plaque) from the teeth of skeletons that were up to 7500 years old could contain relatively well preserved ancient bacterial DNA. That publication was based on 16S rRNA sequences, which are not informative about intra-species genetic diversity. However, subsequent shotgun sequencing from modern and ancient dental calculus [19–21] has demonstrated that it should be possible to reconstruct genomic sequences

that span millennia of human history from multiple individual species within the oral microbiome.

Reconstructing evolutionary history from the oral microbiome faces numerous technical challenges. Our understanding of the historical evolutionary biology of bacterial pathogens benefitted greatly from existing frameworks for the modern population genomic structure of those bacteria [22–24]. However, extensive bacterial population genetic analyses are largely lacking for the modern oral microbiome. The existing literature largely focuses on taxonomic binning into a traditional subset of 40 cultivatable species from periodontitis [25], whose sub-species population structure has not yet been adequately addressed at the genomic level. Instead, most analyses have focused on the ‘oral complexes’, which consist of groups of multiple species whose co-occurrence is statistically associated with periodontitis [26].

A second barrier to reconstructing evolutionary history is the limits of the currently existing bioinformatic tools. The genetic diversity of metagenomic sequences is usually classified by binning the microbial sequence reads into taxonomic units. Taxonomic assignments can be performed by the *de novo* assembly of metagenomic reads into MAGs (metagenomic assembled genomes) [27,28], or by assigning individual sequence reads to existing reference genomes. However, most current metagenomic classifiers rely on the public genomes in the NCBI database, whose composition is subject to extreme sample bias and which represents a preponderance of genomes from pathogenic bacteria [29]. Furthermore, shotgun metagenomes often include DNA from environmental sources, which include multiple microorganisms that have never been cultivated and may belong to unknown or poorly classified microbial taxa whose abundance is not reflected by existing databases. Recent evaluations have also demonstrated that current taxonomic classifiers either lack sufficient sensitivity for species-level assignments or suffer from false positives, and that they overestimate the number of species in the metagenome [29–31]. Both tendencies are especially problematic for the identification of microbial species which are only present at low-abundance, e.g. detecting pathogens in ancient metagenomic samples.

Over the last few years, we have developed a series of tools which can facilitate comparative metagenomics of modern and ancient samples. SPARSE, a novel taxonomic classifier for short-read sequences in the metagenome, was designed to provide accurate taxonomic assignments of metagenomic reads [32]. SPARSE accounts for the existing bias in reference databases [29,33] by sorting all complete genomes of Bacteria, Archaea, Viruses and Protozoa in RefSeq into sequence similarity-based hierarchical clusters with a cutoff of 99% average nucleotide identity (ANI99%). It subsequently extracts a representative subset from those clusters, consisting of one genome per ANI95% cluster because ANI95% is a common cutoff for individual bacterial species [34,35]. SPARSE then assigns metagenomic sequence reads to these clusters using Minimap2 [36]. However, such alignments are likely to be inaccurate when they are widely dispersed across multiple ANI95% clusters because such wide dispersion reflects either ultra-conserved elements of uncertain specificity or a high probability of homoplasies due to horizontal gene transfer. SPARSE therefore reduces such unreliable alignments by a negative weighting of widely dispersed sequences reads. The remaining metagenomic reads are then assigned to unique species-level clusters on the basis of a probabilistic model and labelled

according to the taxonomic labels and pathogenic potential of the genomes within those clusters. Our methodological comparisons demonstrated that SPARSE has greater precision and sensitivity with simulated metagenomic data than 10 other taxonomic classifiers and yielded more correct identifications of pathogen reads within metagenomes of ancient DNA than five other methods [32]. SPARSE is also suitable for classifying reads from metagenomes from modern samples and can extract reads from any ANI95% taxon of interest.

SPARSE assigns sequence reads to taxa, but does not create genomic assemblies from the selected metagenomic reads. That task is performed by EToKi, a stand-alone package of useful pipelines that are used by Enterobase [5] for manipulations of 100 000s of microbial genomes. EToKi is used to merge overlapping paired-end reads, remove low-quality bases and trim adapter sequences. It then excludes sequence reads with greater sequence similarities to genomes from a related but distinct out-group than to an in-group of genomes from the target taxon of interest. EToKi then masks all nucleotides in an appropriate reference genome and creates a pseudo-MAG by unmasking nucleotides with sufficient coverage among the reads that have passed the in-group/out-group comparisons. Finally, EToKi can create a SNP matrix from pseudo-MAGs plus additional draft genomes and generate a maximum-likelihood phylogeny (RAxML 8.2 [37]), which can be visualized together with its metadata in GrapeTree [38].

Here, we demonstrate the power of this combination of pipelines by examination of the metagenomic diversity of the human oral microbiome from a large number of historical and modern samples from diverse geographic sources. We address the question of which microbial taxa are uniformly present in human saliva, dental plaque and dental calculus, and which are specific to individual niches. We test the associations of oral taxa within the traditional oral complexes and conclude that their very existence needs re-examination. Finally, we examine the population genomic structures of *Streptococcus mutans* and *Streptococcus sobrinus*, which are associated with dental caries in some human populations [39–41].

2. Results

(a) SPARSE analysis of oral metagenomes

We identified 17 public archives containing 1016 sets of metagenomic sequences (table 1) from 791 oral samples from a variety of global sources which had been obtained from modern human saliva, modern human dental plaque or historical dental calculus (electronic supplementary material, table S1). Individual sequence reads from those metagenomes were assigned to taxa with SPARSE. The assignments were made according to an upgraded database of 20 054 genomes of bacteria, archaea or viruses, one genome per ANI95% cluster among 101 680 genomes in the NCBI RefSeq databases in May 2018. Seven metagenomes (ancient dental calculus: 5; modern saliva: 2) lacked bacterial reads from the oral microbiome (electronic supplementary material, table S2). These seven metagenomes were ignored for further analyses, leaving assignments to 1591 microbial taxa from 1009 metagenomes (784 samples) (table 2). Electronic supplementary material, table S3, reports the percentage assignment of the reads in each sample to each of the 1591 taxa, except for assignments with a sequence read frequency of less than 0.0001%, which are reported as 0%. Electronic supplementary material, table

Table 1. Sources of metagenomic reads. Note: ancient calculus refers to ancient dental calculus from historical samples. Plaque and saliva refer to modern dental plaque and saliva. Sets of short reads were downloaded from GenBank except for Archive 17, which was downloaded from the Online Ancient Genome Repository. Seven metagenomes (electronic supplementary material, table S2; Archive 11:2; Archive 17:5) were excluded from further analyses because they contained too few reads from common microbial taxa in the oral microbiome.

archive	accession	sets of short reads	no. samples	source	institute	citation
1	PRJNA445215	62	48	ancient calculus	Max Planck Institute for the Science of Human History	[42]
2	PRJEB30331, PRJNA454196	45	44	ancient calculus	University of Oxford	[21]
3	PRJNA216965	9	2	ancient calculus	University of Oklahoma	[19]
4	PRJNA383868	87	87	plaque	J. Craig Venter Institute	[43]
5	PRJNA255922	48	48	plaque	University of California, Los Angeles	[44]
6	PRJNA78025	7	4	plaque	University of Maryland	[45]
7	PRJNA289925	1	1	plaque	University of Washington	[46]
8	PRJEB6997	298	298	plaque & saliva	BGI	[47]
9	PRJNA230363	12	12	plaque & saliva	Chinese Academy of Sciences	[48]
10	PRJEB24090	61	61	saliva	University of California San Diego	[49]
11	PRJNA380727	56	55	saliva	Peking University School of Stomatology	
12	PRJNA396840	30	30	saliva	University of Copenhagen	[50]
13	PRJEB14383	28	28	saliva	University College London	[51]
14	PRJDB4115	26	26	saliva	University of Tokyo	[52]
15	PRJNA217052	217	18	saliva	Broad Institute	[53]
16	PRJNA188481	8	8	saliva	Broad Institute	[54]
17	http://dx.doi.org/10.4225/55/584775546a409	21	21	ancient calculus	OAGR, University of Adelaide	[55]

S3, includes a column identifying assignments to the oral microbial complexes defined by Socransky *et al.* [26]. SPARSE also identified 152 samples containing Archaea from four species, 214 samples containing at least one of four human viruses and 146 samples containing at least one of 12 bacteriophages (table 3). This dataset may represent the currently broadest sample of the oral microbiome from global sources and over time.

(b) Comparisons of microbiomes from saliva, plaque and historical dental calculus

We tested whether individual oral taxa were particularly enriched or depleted according to source with multiple quantitative approaches, including UMAP (Uniform Manifold Approximation and Projection), principal component analysis (PCA) and hierarchical clustering.

UMAP is a recently described, high-performance algorithm for dimensional reduction of diversity within large amounts of data by nonlinear multidimensional clustering [56]. A UMAP plot of the taxon abundances in each sample showed three clusters (figure 1a). The three clusters are totally discrete (electronic supplementary material, figure S1A) according to a machine learning approach, optimal k-mean clustering of the first three components from the UMAP analysis. With minor exceptions, the three UMAP clusters were also

predominantly associated with the source, with one cluster for taxa from modern saliva, a second one for taxa from modern dental calculus and the third for taxa from ancient dental calculus (figure 1a). Similar results were obtained with a classical PCA, except that the clusters were not as clearly distinguished as with UMAP, and the proportion of exceptions was greater (electronic supplementary material, figure S1B). The assignments of source affiliations to cluster were also largely consistent between UMAP and PCA, with occasional exceptions (electronic supplementary material, figure S1C).

For the third approach, we calculated the Euclidean p-distances between each pair of samples and subjected them to hierarchical clustering by the neighbour-joining algorithm with the results shown in figure 1b. Hierarchical clustering also largely separated the samples by source with only a few exceptions. Samples from modern saliva formed one large cluster. Samples from modern dental plaque formed two related but discrete sub-clusters, one of which included a sub-sub cluster of samples from historical dental calculus. These clusters also largely corresponded to the clusters found by k-mean clustering of UMAP data.

Thus, three primary and distinct clusters were consistently identified by three independent methods from the quantitative numbers of reads in individual microbial taxa. The three clusters were largely source specific for modern saliva, modern plaque and historical dental calculus. This finding

Table 2. Sources of genomes from cultivated bacteria and metagenomic samples. Additional details can be found in electronic supplementary material, table S1.

category	sub-category	number
bacterial genomes		262
	<i>S. mutans</i>	195
	<i>S. sobrinus</i>	50
	others	17
metagenome source		784
	ancient dental calculus	110
	modern plaque	287
	modern saliva	387
metagenome size (nucleotides)	0–2GB	343
	2–4GB	129
	4–6GB	162
	6–8GB	93
	8–10GB	45
	>10GB	12
country	Asia	442
	China	375
	Japan	32
	Philippines	28
	others	7
	North America	159
	USA	157
	Guadeloupe	2
	Europe	166
	UK	75
	Ireland	36
	Denmark	31
	others	24
	Oceania	111
	Australia	92
	Fiji	18
	Papua New Guinea	1
	Africa	9
	South Africa	6
	Sudan	2
Sierra Leone	1	

predicts that the microbiomes from these three sources contain source-specific taxa.

(c) Source-specific taxa

We attempted to identify the most important bacterial taxa for the observed clustering by sample source with a second,

powerful machine learning approach. A supervised support vector machine (SVM) [58] classification was used to identify the most optimal of 300 SVM model variants, and the 40 most discriminating ANI95% taxa according to that model are shown in figure 2, together with mini-histograms that summarize the relative abundance of sequences by source. As predicted from the discrete clustering described above, multiple taxa were dramatically more prominent in samples from one source than from either of the two other sources. The results also show that the most prominent sample source varied with the taxon (figure 2).

Eleven of the 40 most discriminatory taxa belonged to the oral complexes that are associated with periodontitis according to Socransky *et al.* [26]. Seven species from oral complexes (*Veillonella parvula*, *Fusobacterium nucleatum*, *Capnocytophaga gingivalis*, *Streptococcus gordonii*, *Actinomyces naeslundii*, *Actinomyces viscosus* and *Capnocytophaga sputigena*) were most abundant in modern plaque and two other species (*Streptococcus sanguinis*, *Tannerella forsythia*) were most abundant in historical dental calculus. The yellow complex includes *Streptococcus mitis*, which encompasses over 50 distinct ANI95% clusters [59]. Two of these ANI95% clusters, designated *S. mitis* s8897 (ANI95% cluster in electronic supplementary material, table S3; MG_43 in [59]) and *S. mitis* s126097 (MG_56) were included among the 40 most discriminatory taxa, and each of them was more frequent in saliva than in dental plaque or dental calculus.

Seventeen other taxa that were assigned to an oral complex by Socransky *et al.* [26] are not included in figure 2 because they were not among the 40 most discriminatory taxa. We therefore examined the relative abundances of all 28 taxa from oral complexes in greater detail (figure 3). Three of the four taxa in the blue and purple complexes are very abundant in oral metagenomes, and all four are preferentially found in modern plaque. However, the other oral complexes are heterogeneous in their patterns of relative abundances. For example, within the red complex, both *T. forsythia* and *Treponema denticola* were most frequently found in historical dental calculus but *Porphyromonas gingivalis* is most frequent in modern plaque and is generally much less abundant. Similar intra-complex discrepancies were found for the orange, yellow and green complexes. These inconsistent frequencies by source raise questions about the consistency of the compositions of those complexes in individual samples.

(d) Existence of 'oral complexes'?

Socransky *et al.* [26] initially treated oral complexes as a hypothesis. However, they have now attained the status of accepted wisdom and even play a prominent role in routine laboratory investigations and treatment of periodontitis. The oral complexes included 28 cultivated bacterial species, whose presence or absence was determined by DNA hybridization against a small number of probes. This technology is now outdated; the number of known oral taxa has increased dramatically and the data presented here are for relative abundance rather than presence or absence. However even after weighting for genome size, we do not find a close correspondence between the frequencies of cells in sub-gingival dental plaque measured by Socransky *et al.* [25] and the results presented here (electronic supplementary material, text). We therefore re-examined the strengths of association with the oral complexes from the data presented here according to

Table 3. Detailed summary of archaea and Viruses in all 784 samples. No. refers to the numbers of samples after combining metagenomes from a common sample. Percentage of reads refers to the percentage of all reads attributed to a taxon in all the metagenomes from that sample.

taxonomy	no. ancient samples (110)	per cent of reads	no. plaque (287)	per cent of reads	no. saliva (387)	per cent of reads
host (human)	110	0.32	243	9.12	335	7.05
archaea (4)	81	1.78	26	2×10^{-4}	45	1×10^{-4}
<i>Methanobrevibacter oralis</i>	79	1.76	26	2×10^{-4}	43	1×10^{-4}
<i>Methanobrevibacter smithii</i>	1	3×10^{-5}			2	2×10^{-6}
<i>Candidatus Nitrosoarchaeum koreensis</i>	1	1×10^{-5}			0	
<i>Thermoplasmatales archaeon BRNA1</i>	1	7×10^{-6}			0	
human viruses (4)	0		25	3×10^{-4}	189	4×10^{-3}
<i>Human betaherpesvirus 7</i>	0		8	9×10^{-6}	150	6×10^{-4}
<i>Human gammaherpesvirus 4</i>	0		16	3×10^{-4}	86	3×10^{-3}
<i>Human alphaherpesvirus 1</i>	0		1	5×10^{-6}	9	8×10^{-5}
<i>Human betaherpesvirus 6B</i>	0		0		7	2×10^{-5}
bacteriophages (12)	3	1×10^{-5}	26	3×10^{-5}	117	2×10^{-4}
<i>Streptococcus EJ-1</i>	0		14	1×10^{-5}	56	8×10^{-5}
<i>Streptococcus SM1</i>	2	5×10^{-6}	11	1×10^{-5}	41	3×10^{-5}
<i>Streptococcus SpSL1</i>	0		0		9	2×10^{-5}
<i>Streptococcus Dp-1</i>	0		0		7	2×10^{-5}
<i>Streptococcus DT1</i>	0		0		7	2×10^{-5}
<i>Streptococcus PH10</i>	1	6×10^{-6}	2	3×10^{-6}	7	6×10^{-6}
<i>Klebsiella KP15</i>	0		0		6	6×10^{-6}
<i>Lactococcus r1t</i>	0		0		6	4×10^{-6}
<i>Streptococcus YMC-2011</i>	0		0		4	1×10^{-5}
<i>Propionibacterium PHL060L00</i>	0		0		2	2×10^{-6}
<i>Propionibacterium PHL179</i>	0		0		1	2×10^{-6}
<i>Propionibacterium PAD20</i>	0		0		1	2×10^{-6}

similar criteria and similar methods as those used in the Socransky *et al.* 1998 publication [26].

The original assignments to the oral complexes depended strongly on results from hierarchical clustering of the pairwise concordance between species for presence or absence in individual samples. The tree in figure 4 shows neighbour-joining clustering of the common microbial taxa in our dataset by the similarities of their abundances over all samples in our dataset according to SPARSE. This tree contradicts the original composition of the oral complexes: the four areas of the tree where oral complex taxa are clustered each contain representatives from multiple complexes, and none of those four clusters corresponds to the original compositions proposed by Socransky *et al.* [26].

It seemed possible that the discrepancies between figure 4 and the original compositions of the oral complexes might reflect the fact that this study identified many additional taxa, some of which were as common as those used to define the oral complexes (electronic supplementary material, text). We therefore performed cluster analyses of our current data for the original set of 31 cultivatable bacterial species examined by Socransky *et al.* [26]. We compared the neighbour-joining algorithm used here with the less powerful, agglomerative clustering method (UPGMA, unweighted pair group method

with arithmetic mean) that had been used by Socransky *et al.* We also compared the abundances across all samples with abundances in plaque, which was the primary source for bacteria tested by Socransky *et al.* The results (electronic supplementary material, figure S3) show dramatic inconsistencies between independent trees in regard to the clustering of the oral complex bacteria. For example, *T. forsythia*, *T. denticola* and *P. gingivalis* of the red complex cluster together (and also with *C. rectus*) in electronic supplementary material, figure S3A,C,F,G. However, *T. denticola* and *T. forsythia* are separated from *P. gingivalis* in the four other graphs in electronic supplementary material, figure S3. And none of the three cluster together with each other in electronic supplementary material, figure S3E. Similar, or even greater, discrepancies are visible for the other oral complexes in electronic supplementary material, figure S3. Inconsistencies in clustering patterns across minor differences in sampling and clustering algorithms raise severe doubts about the very existence of the oral complexes as defined by Socransky *et al.* [26].

(e) Numbers of taxa per source

The rarefaction curves in figure 5a provide a breakdown of taxa by sample source as additional samples are tested.

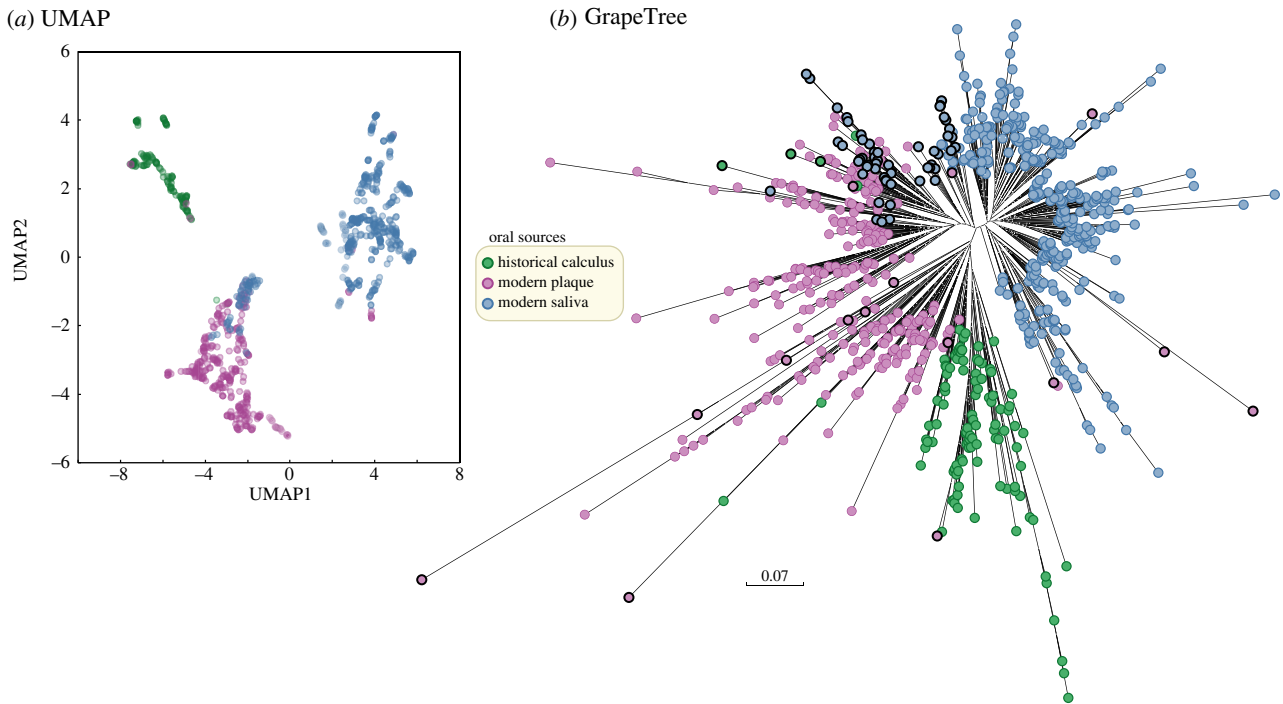


Figure 1. Source specificity of the percentage of species composition in 784 oral metagenomes according to SPARSE. (a) X-Y plot of the first three components from a UMAP (Uniform Manifold Approximation and Projection) [56] dimensional reduction of taxon abundances. (b) Neighbour-joining (FastMe2; [57]) hierarchical clustering based on the Euclidean distances between pairs of metagenomes. Euclidean p-distances were calculated between each pair as the square root of the sum of the squared pairwise differences in the percentage of reads assigned by SPARSE to each microbial taxon. Nodes whose cluster location was inconsistent with the UMAP clustering in (b) are highlighted with black perimeters. Tree visualization: GrapeTree [38].

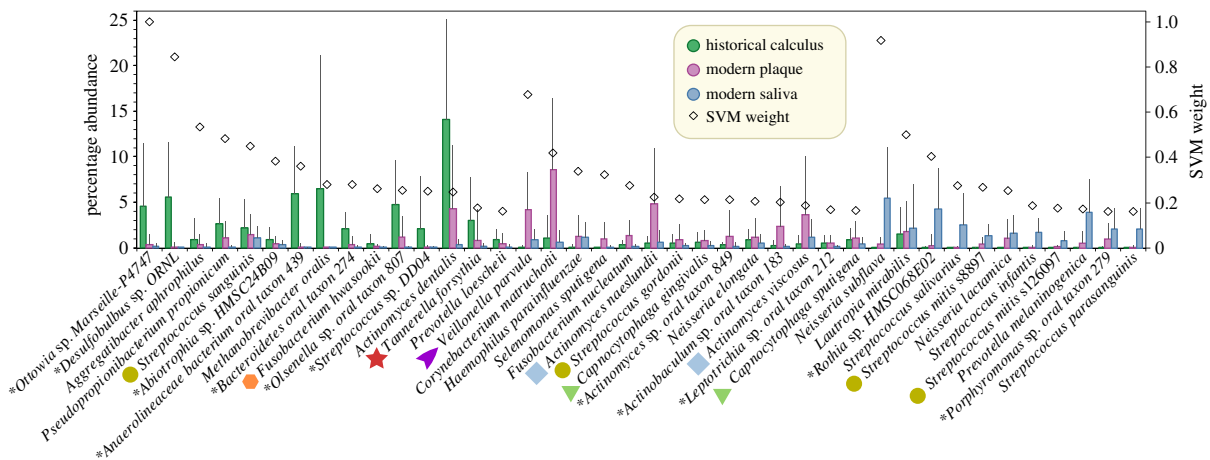


Figure 2. Average percentage abundance (left axis) of bacterial species by source for the 40 most discriminating species according to Support Vector Machine analysis. The relative abundances for each of the three sources are indicated by mini-histograms for each species; error bars indicate standard deviations. Species are sorted in descending order by predominant source and then by SVM weight (squared coefficient) in the optimal model. Species belonging to oral complexes are indicated by oral complex-specific shapes and colours. Key legend: source colours used in the mini-histograms and symbol for SVM weight. Asterisk, species designations assigned by RefSeq to single genomes which have not (yet) been confirmed by taxonomists. *S. mitis* is separated into multiple ANI95% clusters, two of which (s8897; s126097 (electronic supplementary material, table S3)) are among the predominant taxa associated with saliva.

SPARSE detected 1591 microbial taxa over all 784 metagenomic samples: 1389 from modern saliva, 842 from modern plaque and 696 from historical calculus. These estimates will increase as additional samples are added, but at increasingly slower rates because the rarefaction curves seem to be reaching a plateau, except for historical dental calculus where the fewest samples have been evaluated until now.

The median numbers of taxa per sample range from 177 (historical dental calculus) to 288 (modern saliva) and were much smaller than the total numbers. These median values

reflect a bimodal distribution for numbers of taxa per sample (figure 5b), wherein a few samples had jackpots of large numbers of taxa but all other samples had only a few.

The analyses described above focused on differences in taxon composition by source. However, the Venn diagram in figure 5c shows that 447 taxa were common to all three sources, even if their relative abundances varied. Modern plaque yielded only 34 taxa which were not found in either historical dental calculus or modern saliva. More source-specific taxa were found in historical dental calculus, which may possibly

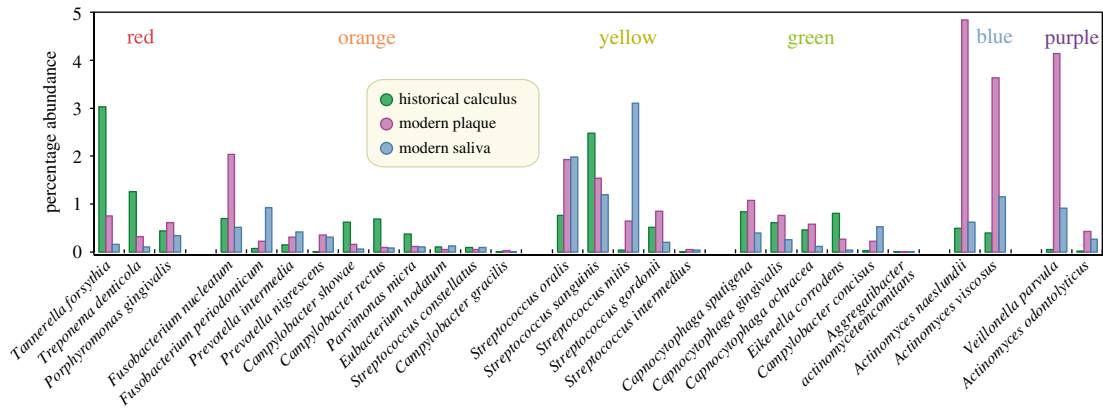


Figure 3. Average percentage abundances in 784 metagenomes by oral source (key legend) of 28 species from six oral complexes described by Socransky *et al.* [26]. The oral sources are indicated by three mini-histogram bars for each species. Species are ordered from left to right by oral complex, whose colour designation is indicated at the top. Within each oral complex, the species order is by decreasing total abundance.

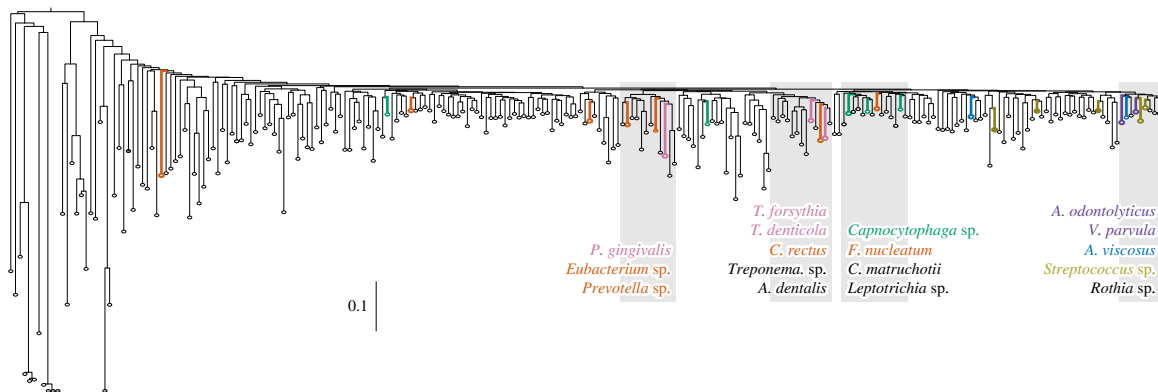


Figure 4. Neighbour-joining (FastMe2; [57]) hierarchical clustering based on the Euclidean distances between pairs of 245 microbial species whose percentage abundance was greater than 2% in at least one metagenome. Members of the six oral complexes [26] are highlighted by coloured species names, whose colours indicate their oral complex membership. These species do not cluster by oral complex, but by other unnamed groupings, four of which are highlighted in grey. An expanded version of the same tree including all species labels is available in electronic supplementary material, figure S2. Branch length distance scale bar is next to the distance of 0.1.

reflect some contamination with environmental material. Alternatively, some taxa may be absent in modern dental plaque because historical lineages have become extinct [11]. Saliva yielded 504 unique taxa, some of which might be transient and do not persist long enough to be incorporated into plaque.

(f) Population genomics of organisms associated with dental caries

The microbiome associated with early stages of dental caries is an unresolved topic that remains under active investigation [40,60–62]. However, it is generally accepted that *S. mutans* and *S. sobrinus* are routinely associated with caries [63]. Our data confirm that reads belonging to these two taxa are abundant in modern dental plaque and also show that they are even more abundant in modern saliva (figure 6*a,c*). However, there was no significant correlation between the relative frequencies of these species across multiple metagenomes (electronic supplementary material, figure S9). Prior analyses based on 16S RNA operational taxonomic units (OTUs) indicated that *S. mutans* was extremely rare in historical dental calculus and argued that this increase was caused by the introduction of high levels of sugar to human diets in industrialized societies in the last 200 years [18]. Our data show that *S. sobrinus* was

undetectable in historical samples (frequency of less than 0.0001% of reads or less than 10 reads per metagenome; figure 6*c*). *S. mutans* was also undetectable in most of these samples, but up to 0.04% of all reads in 10 historical samples spanning the last 1500 years were assigned to *S. mutans* (figure 6*a*), in accordance with archaeological findings that dental caries has been common in multiple eras over the last 10 000 years [14]. The few reads from historical samples that were assigned to *S. mutans* showed increased deamination at their 5'-ends when tested by MapDamage2 [64] (electronic supplementary material, figure S4), confirming that they were truly from ancient DNA.

We exploited the high frequency of sequence reads from these two *Streptococcus* species in modern dental plaque and saliva to illustrate how SPARSE and ETOKi can be used to extract pseudo-MAGs from metagenomic sequence reads and combine them with genomes sequenced from cultivated bacteria (see Methods). These procedures resulted in a total of 31 pseudo-MAGs for *S. mutans* and 15 pseudo-MAGs for *S. sobrinus* in which over 70% of the reference genome had been unmasked (figure 6*e,f*; electronic supplementary material, table S6). Most of these pseudo-MAGs were from Chinese samples [47]. The pseudo-MAGs were combined with genomes from cultivated bacteria of the same species from Brazil, the US and the UK as well as other countries

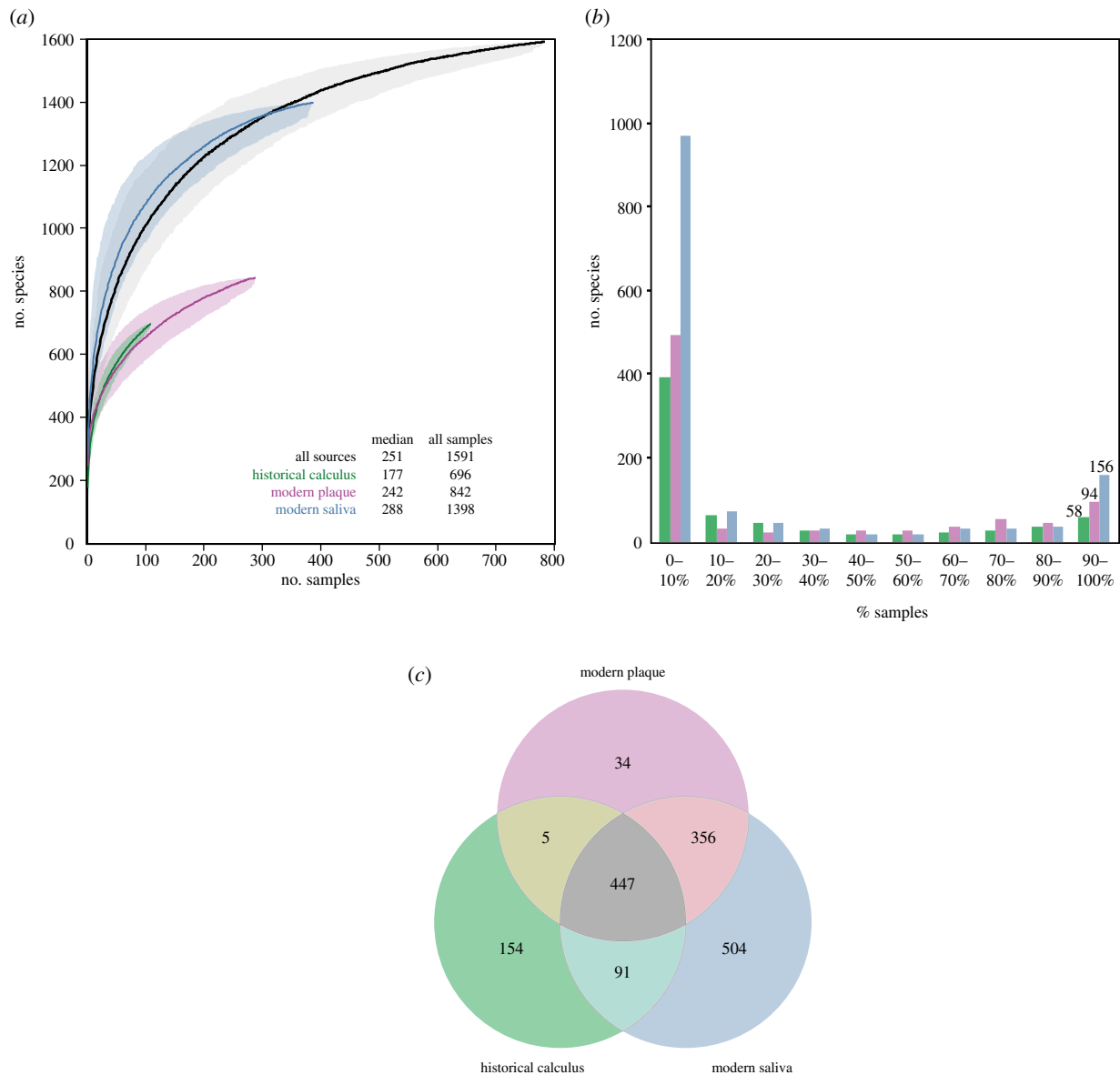


Figure 5. Numbers of microbial taxa by source. (a) Rarefaction curves of numbers of species by source, with 95% confidence estimates (shadow). Inset data indicate median numbers of species per sample by source, as well as the total numbers for all sources. Rarefactions were performed with the program script called SPARSE_curve.py, using 1000 randomized permutations of the order of samples. (b) Binned histograms of number of species by percentage of samples. The data for this plot was also calculated with SPARSE_curve.py. (c) Venn diagram of overlapping presence of taxa ($\geq 0.0001\%$ abundance) for the three oral sources.

(table 2) and maximum-likelihood (ML) phylogenies of non-repetitive SNPs (figure 7) were created with EToKi (see Methods).

The ML phylogenies of the two species showed interesting differences. All 13 Chinese pseudo-MAGs clustered together within the *S. sobrinus* ML tree (figure 7b), whereas almost all the other 44 bacterial genomes from Brazil and elsewhere clustered distantly. By contrast, in the *S. mutans* tree (figure 7a), 20 Chinese pseudo-MAGs did not show any obvious phylogeographic specificities and were inter-dispersed among 196 bacterial genomes from multiple geographic locations. Similar conclusions about a lack of phylogeographic specificity were previously reached by Cornejo *et al.* [65] on a subset of 57 of these *S. mutans* genomes.

3. Discussion

Several years ago, we accidentally became interested in comparing historical and modern genomes reconstructed from metagenomic short-read sequences with draft genomes

assembled from high-throughput sequencing of cultivated bacteria. Our initial efforts involved the deployment of individual bioinformatic tools, comparisons of multiple publicly available algorithms and compilation of draft genomes from publicly available sequence read archives of short-read sequences [2]. In parallel, we were also involved in developing EnterBase, a compendium of 100 000s of draft genome assemblies from multiple genera that can cause enteric diseases in humans, including *Salmonella* [5,24]. These two projects were synergistic for elucidating the evolutionary history of *Salmonella enterica* based on metagenomic sequences from 800-year-old bones, teeth and dental calculus [11]. In that case, sequence reads from *S. enterica* were found in teeth and bone, but not in dental calculus. Our attempts to examine further samples of dental calculus quickly demonstrated that optimized pipelines were needed because manual analyses were too time intensive. However, none of the existing tools were both reliable and sufficiently sensitive for assigning sequence reads from historical metagenomes to the tree of microbial life. We therefore took a step back

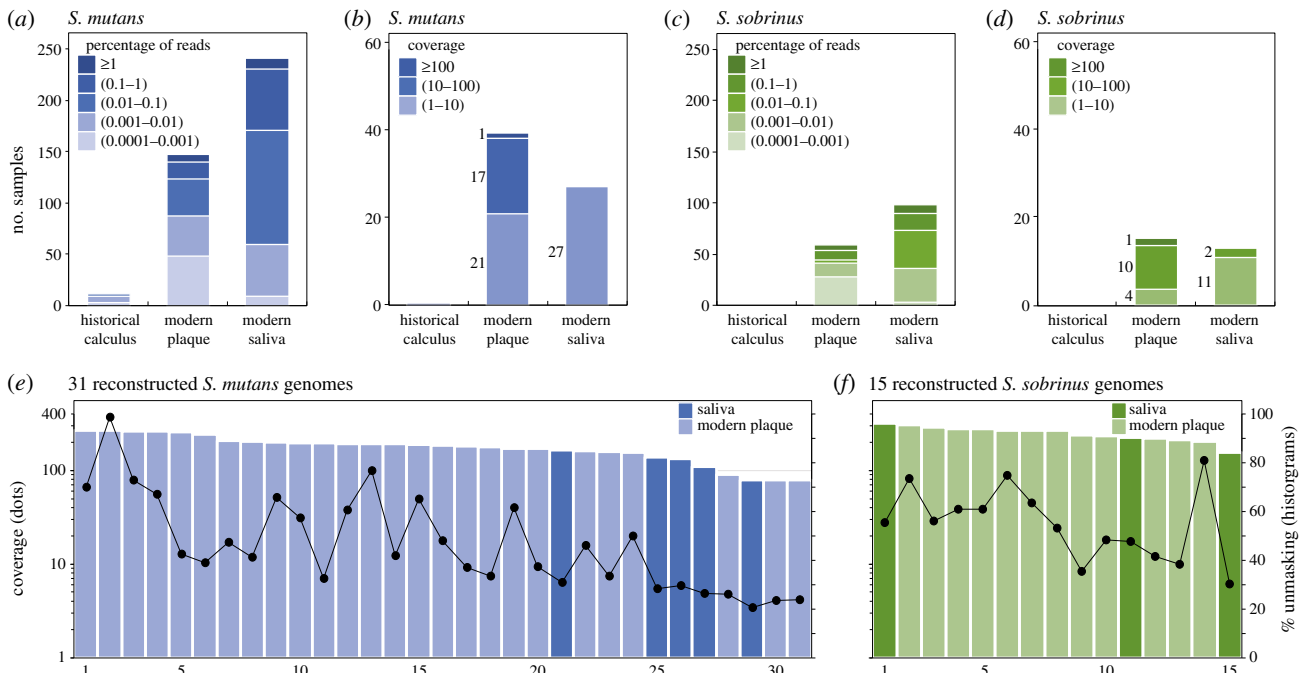


Figure 6. Reconstruction of pseudo-MAGs (metagenomic assembled genomes) of *S. mutans* and *S. sobrinus* from oral metagenomes. (a and c) Numbers of oral samples by source binned by the percentage of reads specific to *S. mutans* (a) and *S. sobrinus* (c). (b and d) Numbers of oral samples by source with an average coverage of at least 1x. The data are binned by the predicted read coverage against a reference genome of *S. mutans* (UA159; (b)) and *S. sobrinus* (NCTC12279; (d)). (e and f) Read coverage (dots; left) and percentage of the reference genome that was unmasked (≥ 3 reads; $\geq 70\%$ consistency) (histogram; right) in *S. mutans* (e) and *S. sobrinus* (f). Ordered by decreasing percentage unmasking.

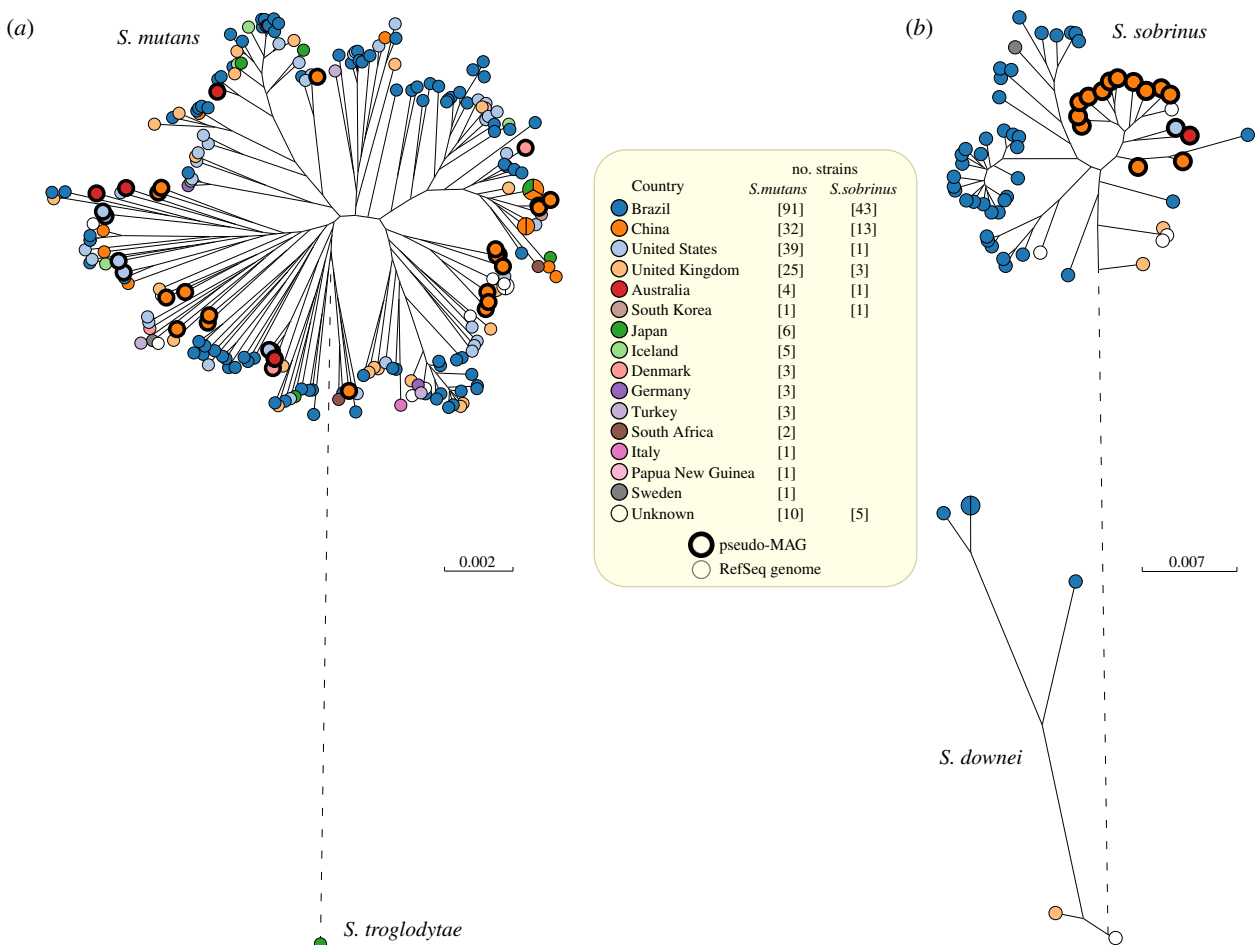


Figure 7. Maximum-likelihood phylogenies of *S. mutans* and *S. sobrinus* genomes. (a) A RaxML [37] tree of 226 genomes of *S. mutans* (RefSeq: 195; pseudo-MAGs: 31) plus one genome of *S. troglodytae* as an out-group. The tree was based on 181 321 non-repetitive SNPs in 1.73 Mb. (b) A RaxML tree of 61 genomes of *S. sobrinus* (RefSeq: 46; pseudo-MAGs: 15) plus six *S. downei* genomes as an out-group. The tree was based on 160 863 non-repetitive SNPs in 1.13 Mb. Pseudo-MAGs are highlighted by thick black perimeters. Visualization with GrapeTree [38]. Branches with a genetic distance of greater than 0.1 were shortened for clarity and are shown as dashed lines. Legend: numbers of strains by country of origin for both trees.

and developed SPARSE [32] to satisfy our requirements. SPARSE replaces the current reference databases, which are strongly biased to multiple, closely related genomes from bacterial pathogens, by a representative subset consisting of one genome per ANI95% hierarchical cluster within RefSeq, and assigns sequence reads to these clusters using a probabilistic model. That model penalizes non-specific mappings of reads and hence reduces false-positive assignments. SPARSE was more reliable than multiple other taxonomic classifiers, and both more sensitive and more reliable for identifying low numbers of reads from ancient metagenomes than multiple other pipelines [32]. In parallel, we expanded the capacities of EToKi [5], an efficient backend pipeline for genomic manipulations, such that it can accurately identify individual sequence reads sieved through SPARSE that are more similar to an in-group of reference genomes from the target species than to an out-group of genomes from a closely related, but distinct taxon. Those reads are then used to unmask nucleotides in a reference genome and generate a pseudo-MAG for SNP-based maximum-likelihood phylogenies. Finally, we developed GrapeTree [38], which facilitates the graphic visualization and manipulation of phylogenetic trees based on large numbers of genomes. Here, we demonstrate how to combine all three tools in order to obtain an overview of the microbial flora in samples from human oral saliva, modern dental plaque and historical dental calculus. We also reconstructed genomes of two taxa present at moderate concentrations within the oral microbiome and compare them with conventional draft genomes. The experimental procedures for processing 1016 metagenomes consisted of running SPARSE in the background for 2 months (approx. 100 000 CPU h). The pipelines described here permitted all other procedures and evaluations described here to be completed in less than two weeks.

Our traditional understanding of oral ecology is largely based on taxonomic assignments of cultivatable bacteria, often performed by checkerboard DNA–DNA hybridization [25]. Currently, 756 species have been cultivated from the human oral cavity and respiratory tract [66]. A subset of 40 are used for checkerboard DNA–DNA hybridization [25], of which 28 were used to define the oral complexes that were thought to be of importance for periodontitis [26]. Our comparisons of those data with the results from the metagenomic analyses presented here show that the frequencies of individual taxa determined by the checkerboard assay were inconsistent with the frequencies determined by our metagenomic analyses (electronic supplementary material, figures S5 and S6). The checkerboard assays also lacked 17 common taxa from dental plaque and dental calculus that were found by metagenomic analyses. These results are not unexpected because our metagenomic analyses included saliva samples as well as ancient dental calculus and identified 1591 taxa, many of which have not been cultivated. Furthermore, it is now well established that the frequencies of certain supposed members of the oral complexes differ very dramatically with geographical source [67]. However, we had anticipated that we might be able to expand the compositions of the oral complexes to include previously uncultivated organisms. Instead, we were unable to reliably identify their very existence (figure 4) because clustering of taxa was affected by minor changes in choice of samples and the choice of clustering algorithm (electronic supplementary material, figure S3). We therefore conclude that the existence and composition of the

oral complexes needs independent verification by modern techniques and new samples.

The data presented here provide an unprecedented comparative overview of the relative proportions of the predominant taxa in public available metagenomes from the modern and historical oral microbiome. Figure 2 identifies 15 taxa, which are particularly common in historical calculus, 14 others that are preferentially found in modern dental plaque and 11 that seem to be specific for saliva. These associations with a particular source in the oral cavity might be used to identify currently undefined ecological complexes of oral taxa that share a common niche. However, species-level OTUs are likely to be conglomerates of multiple microbial populations, each of which may inhabit a somewhat different ecology. For some organisms such as *Salmonella* or *Escherichia*, efforts are currently underway to develop hierarchical clustering of such populations in order to categorize their ecological and pathogenic differentiation [5]. A step in this direction for the oral microbiome is the recognition of ANI95% clusters s8897 and s126097, both of which were preferentially found in saliva. A large study of all streptococci [59] identified multiple other ANI95% clusters within *S. mitis* but their preferential location in the oral cavity has not yet been addressed. Indeed, little is yet known about the sub-species population structure of almost all of the taxa identified here.

Our more detailed investigation of *S. mutans* and *S. sobrinus* may represent a forerunner of future studies on sub-species ecological differences within the oral microbiome. *S. mutans* and *S. sobrinus* are commonly associated with dental caries and may play a causal role in that disease [63]. However, once again these taxa were more common in saliva than in dental plaque (figure 6). We chose *S. mutans* and *S. sobrinus* for more detailed analysis because sufficient reads were found in multiple metagenomes from modern samples to allow the partial reconstruction of multiple genome sequences (pseudo-MAGs). In addition, multiple draft genomes from cultivated bacteria existed in the public domain, which were available for genomic comparisons. We were also intrigued by the claim that *S. mutans* was rare in historical plaque [18]. Our data support that claim, and we found only a few historical samples of dental calculus that contained any reads of *S. mutans*, and none with *S. sobrinus*. Our data also support prior conclusions of a lack of phylogeographic differentiation within *S. mutans* [65]. However, although the data are still somewhat limited, *S. sobrinus* from China tend to cluster distinctly from genomes from Brazil (figure 7). Distinct clustering might reflect phylogeographical signals but other causes of clustering cannot currently be excluded because the Chinese genomes were pseudo-MAGs reconstructed from metagenomes from dental plaque and saliva while the Brazil genomes were from bacteria cultivated from dental plaque. Additional genomes of *S. sobrinus* from other geographical areas would be needed to determine whether the apparent phylogeographical trends are robust. Such analyses could also be facilitated by creating an EnteroBase for *Streptococcus*, which could be done relatively easily [59] if there were interested curators and sufficient interest in the *Streptococcus* community.

In summary, we illustrate the use of a variety of reliable, high-throughput tools for determining microbial diversity within metagenomic data, and for extracting microbial genomes from metagenomes. We illustrate these tools with metagenomes from both modern and historical samples,

and release all the data and methods for further use by others.

4. Methods

(a) SPARSE database update

In its original incarnation in August 2017 [32], SPARSE used MASH [68] to assign 101 680 genomes from the NCBI RefSeq database to 28 732 ANI99% clusters of genomes. By May 2018, 21 540 additional genomes had been added to NCBI RefSeq. These were merged into the existing database in the same manner as previously, by merging that each genome into an existing ANI99% cluster or by creating a new cluster containing one genome if the ANI to all existing clusters was less than 99%. An ANI99% representative microbial database was generated which contained one representative genome for each of the 32 378 ANI99% clusters containing bacteria, archaea or viruses plus a human reference genome (Genome Reference Consortium Human Build 38) such that reads from human DNA could also be called. All the representative genomes were assigned to a superset of 20 054 ANI95% clusters, and this was used for species assignments and genomic extractions as described [32].

(b) SPARSE analyses

'EToKi prepare' was used to collapse paired-end reads and trim all sequence reads. Subsequent SPARSE analyses were performed on all the metagenomes in table 1 and additional metagenomes in electronic supplementary material, figure S7, as described in the SPARSE manual (<https://sparse.readthedocs.io/en/latest/>). The first step was 'SPARSE predict', which identifies ANI95% groups containing greater than or equal to 10 specific reads. Subsequently, 'SPARSE report -low 0.0001' was used to assign taxon designations to the ANI95% groups and produce a table of all metagenome results (electronic supplementary material, table S3) which lists distinct taxa for each metagenome that accounted for $\geq 0.0001\%$ of all its reads. Electronic supplementary material, table S3, also includes the designations of oral complexes and other known pathogens according to a manually curated dictionary. Sequence reads were extracted from the metagenomes for assembling pseudo-MAGs with 'SPARSE extract'.

For electronic supplementary material, figures S5–S8, the taxonomic assignments were inversely weighted by genome size in order to render them comparable to DNA–DNA Checkerboard data and output from Metaphlan2, which calculate cell counts. To this end, the number of metagenomic reads assigned to each species within a metagenome was divided by the genome size of the SPARSE reference genome for that species. These data were then expressed as a proportion of the summed data for all microbial species within that metagenome.

(c) Metagenomes lacking reads from the oral microbiome

We tested all metagenomes to identify any that might be grossly contaminated by collating the 50 most abundant microbial species over all metagenomes (electronic supplementary material, table S4A). The percentage of reads in these 50 taxa was summed for each metagenome and expressed as a percentage of all microbial reads. Seven metagenomes (ancient dental calculus: 5; modern saliva: 2; electronic supplementary material, table S2) were excluded because the percentages of those top oral microbes constituted less than 15% of their total microbial reads.

(d) Dimensional reduction of frequencies of reads

Two forms of dimensional reduction of diversity were used to detect source-specific clustering within the SPARSE results.

UMAP analysis was performed with its Python implementation [56], using the parameters `min_neighbours = 5` and `min_dist = 0.0`. PCA was performed using the `decomposition.PCA` module of the `scikit-learn` Python library [69]. Optimal k-mean clusters of the first three components from the UMAP analysis were calculated with the `sklearn.cluster` module of the `scikit-learn` Python library.

(e) Ranking of microbial species by their associations with source

Microbial species were ranked by their weighting according to an SVM classification [58]. A supervised SVM classification of samples was performed using the SVM module of the `scikit-learn` Python library on the raw SPARSE results (electronic supplementary material, table S3). The SVM classification was performed 300 times on a randomly chosen training set consisting of 60% of all samples with varying penalty hyper-parameter `C` and scored using fivefold cross-validation. The model was then tested with the optimal hyper-parameter from all runs on the remaining 40% of samples and correctly inferred the oral source for greater than 96% of the test samples. The optimal SVM coefficients for each individual species were estimated by training that model once again on all the oral samples. The order of the species in figure 2 consists of the SVM weights (squares of the coefficients; [70]) in descending order. The Python scripts described in Methods (d,e), as well as their outputs are freely accessible online as Dataset S3 in <https://github.com/zheminzhou/OralMicrobiome>.

(f) Genome reconstructions for *Streptococcus mutans* and *Streptococcus sobrinus*

SPARSE identified samples in which the metagenomic sequence reads covered at least 2 MB of the reference genome for *S. mutans* (ANI95% cluster s5; 66 samples) or *S. sobrinus* (s3465; 28 samples) (figure 4*b,d*). The cleaned, species-specific reads generated from these samples as in Methods (b) were processed with the stand-alone version of EToKi as described in electronic supplementary material, fig. S6 of Zhou *et al.* 2020 [5] and in greater detail in the online manual (<https://github.com/zheminzhou/EToKi>). EToKi assemble was then used to identify genome-specific reads after specifying a reference genome, an in-group of related genomes, and a related but distinct out-group of other genomes. For *S. mutans*, the reference genome was UA159 (accession code GCF_000007465), the in-group was 194 other *S. mutans* genomes in RefSeq (electronic supplementary material, table S5), and the out-group was 62 genomes from other species in the Mutans *Streptococcus* group according to Zhou *et al.* 2020 [59]. For *S. sobrinus*, the reference genome was NCTC12279 (accession code GCF_900475395), the in-group was 45 other *S. sobrinus* genomes and the out-group was 211 genomes from other Mutans streptococci (electronic supplementary material, table S5). The assemble module replaces nucleotides in the reference genome by their calculated SNVs after checking that they are supported by at least 70% of at least three metagenomic reads, and that the supporting read frequencies are at least one-third of the average read depth. The resulting pseudo-MAGs are listed in electronic supplementary material, table S6 and are freely accessible online as Datasets S1 and S2 in <https://github.com/zheminzhou/OralMicrobiome>.

'EToKi align' was used to create an alignment of non-repetitive SNPs from 31 *S. mutans* pseudo-MAGs plus all 195 *S. mutans* genomes plus the sole *S. troglodytae* genome in RefSeq (electronic supplementary material, table S5). The alignments spanned 1.73 MB that were shared by at least 95% of the genomes and covered 181 321 core SNPs. Similarly, an alignment of 15 *S. sobrinus* MAGs, 46 draft or complete *S. sobrinus* genomes plus 6 genomes

of *Streptococcus downei* from RefSeq spanned 1.16 MB and contained 160 863 core SNPs. These alignments were subjected to maximum-likelihood phylogeny reconstruction by EToKi phylo. Both ML trees were then visualized with GrapeTree [38].

(g) DNA damage patterns for ancient *Streptococcus mutans* reads

SPARSE assigned low numbers of sequence reads to *S. mutans* in 10 metagenomes from ancient dental calculus (figure 6; electronic supplementary material, table S3). In order to assess their authenticity, these reads were assessed with MapDamage2 [64] for patterns of cytosine deamination that are characteristic of authentic ancient DNA. To this end, all *S. mutans*-specific reads were extracted with SPARSE. They were aligned to the *S. mutans* reference genome UA159 with Minimap2 [36] and reads which were at least 95% identical with the reference genome were used to create BAM alignments. SouthAfr2 contained 11 specific reads according to SPARSE, but only eight survived this filtering step. SouthAfr2 was therefore excluded from further analyses because these were too few reads to provide reliable analyses. The BAM alignments from the remaining nine metagenomes consist of both fully aligned reads (46–72%) and others which were ‘soft-clipped’, i.e. terminal bases were not aligned to the reference genome. In order to ensure that these soft-clipped reads were also specific, we compared the alignment scores for all reads

against UA159 with the alignment scores against the 62 out-group genomes in Mutans streptococci (electronic supplementary material, table S5) and found that the scores with UA159 were highest. We also tested the alignment scores against two other *S. mutans* genomes (SA38, [GCF_000339615]; 4VF1 [GCF_000339215]; electronic supplementary material, table S5), but neither yielded higher alignment scores than UA159. The outputs from MapDamage2 show the soft-clipping ends by a yellow line (electronic supplementary material, figure S4A–D).

Data accessibility. The pseudo-MAGs reconstructed from metagenomes for *S. mutans* and *S. sobrinus* are freely accessible in tar.gz files containing Datasets_S1 and Dataset_S2 at <https://github.com/zheminzhou/OralMicrobiome>, respectively. Python scripts that were used to prepare data for figures 1–5 and S1–S3 are available as Dataset_S3 in the same repository. The taxonomic profiling by SPARSE of all 784 metagenomes is available in electronic supplementary material, table S3. Interactive versions of figure 7 are available at <http://enterobase.warwick.ac.uk/a/42277> (figure 7a) and <http://enterobase.warwick.ac.uk/a/42279> (figure 7b).

Authors' contributions. Z.Z. analysed data and prepared the figures. M.A. and Z.Z. interpreted the results and wrote the manuscript.

Competing interests. We have no competing interests.

Funding. This project was supported by the Wellcome Trust (grant no. 202792/Z/16/Z) and EnteroBase development was funded by the BBSRC (grant no. BB/L020319/1).

References

- Bos KI *et al.* 2014 Pre-Columbian mycobacterial genomes reveal seals as a source of New World human tuberculosis. *Nature* **514**, 494–497. (doi:10.1038/nature13591)
- Kay GL *et al.* 2015 Eighteenth-century genomes show that mixed infections were common at time of peak tuberculosis in Europe. *Nat. Commun.* **6**, 6717. (doi:10.1038/ncomms7717)
- Schilling AK *et al.* 2019 British red squirrels remain the only known wild rodent host for leprosy bacilli. *Front. Vet. Sci.* **6**, 8. (doi:10.3389/fvets.2019.00008)
- Schuenemann VJ *et al.* 2018 Ancient genomes reveal a high diversity of *Mycobacterium leprae* in medieval Europe. *PLoS Pathog.* **14**, e1006997. (doi:10.1371/journal.ppat.1006997)
- Zhou Z, Alikhan N-F, Mohamed K, Fan Y, Group AS, Achtman M. 2020 The EnteroBase user's guide, with case studies on *Salmonella* transmissions, *Yersinia pestis* phylogeny, and *Escherichia* core genomic diversity. *Genome Res.* **30**, 138–152. (doi:10.1101/gr.251678.119)
- Bos KI *et al.* 2011 A draft genome of *Yersinia pestis* from victims of the Black Death. *Nature* **478**, 506–510. (doi:10.1038/nature10549)
- Rasmussen S *et al.* 2015 Early divergent strains of *Yersinia pestis* in Eurasia 5000 years ago. *Cell* **163**, 571–582. (doi:10.1016/j.cell.2015.10.009)
- Damgaard PB *et al.* 2018 137 Ancient human genomes from across the Eurasian steppes. *Nature* **557**, 369–374. (doi:10.1038/s41586-018-0094-2)
- Keller M *et al.* 2019 Ancient *Yersinia pestis* genomes from across Western Europe reveal early diversification during the First Pandemic (541–750). *Proc. Natl Acad. Sci. USA* **116**, 12 363–12 372. (doi:10.1073/pnas.1820447116)
- Spyrou MA *et al.* 2019 Phylogeography of the second plague pandemic revealed through analysis of historical *Yersinia pestis* genomes. *Nat. Commun.* **10**, 4470. (doi:10.1038/s41467-019-12154-0)
- Zhou Z *et al.* 2018 Pan-genome analysis of ancient and modern *Salmonella enterica* demonstrates genomic stability of the invasive Para C Lineage for millennia. *Curr. Biol.* **28**, 2420–2428. (doi:10.1016/j.cub.2018.05.058)
- Vågene ÅJ *et al.* 2018 *Salmonella enterica* genomes from victims of a major sixteenth-century epidemic in Mexico. *Nat. Ecol. Evol.* **2**, 520–528. (doi:10.1038/s41559-017-0446-6)
- Key FM *et al.* 2020 Emergence of human-specific *Salmonella enterica* is linked to the Neolithization process. *Nat. Ecol. Evol.* **4**, 324–333. (doi:10.1038/s41559-020-1106-9)
- Lacy SA. 2014 Oral health and its implications in late Pleistocene Western Eurasian humans. PhD thesis, St. Louis, MO: Washington University.
- Dewitte SN, Bekvalac J. 2010 Oral health and frailty in the medieval English cemetery of St Mary Graces. *Am. J. Phys. Anthropol.* **142**, 341–354. (doi:10.1002/ajpa.21228)
- Carter F, Irish JD. 2019 A sub-continent of caries: prevalence and severity in early Holocene through recent Africans. *Dental Anthropol.* **32**, 22–29. (doi:10.26575/daj.v32i2.285)
- Towle I, Irish JD, De Groot I, Fernée C. 2019 Dental caries in human evolution: frequency of carious lesions in South African fossil hominins. *BioRxiv* 597385. (doi:10.1101/597385)
- Adler CJ *et al.* 2013 Sequencing ancient calcified dental plaque shows changes in oral microbiota with dietary shifts of the Neolithic and Industrial revolutions. *Nat. Genet.* **45**, 450–455. (doi:10.1038/ng.2536)
- Warinner C *et al.* 2014 Pathogens and host immunity in the ancient human oral cavity. *Nat. Genet.* **46**, 336–344. (doi:10.1038/ng.2906)
- Warinner C, Speller C, Collins MJ. 2015 A new era in palaeomicrobiology: prospects for ancient dental calculus as a long-term record of the human oral microbiome. *Phil. Trans. R. Soc. B* **370**, 20130376. (doi:10.1098/rstb.2013.0376)
- Velsko IM *et al.* 2019 Microbial differences between dental plaque and historic dental calculus are related to oral biofilm maturation stage. *Microbiome* **7**, 102. (doi:10.1186/s40168-019-0717-3)
- Coll F *et al.* 2014 A robust SNP barcode for typing *Mycobacterium tuberculosis* complex strains. *Nat. Commun.* **5**, 4812. (doi:10.1038/ncomms5812)
- Achtman M. 2016 How old are bacterial pathogens? *Proc. Biol. Sci.* **283**, 1836. (doi:10.1098/rspb.2016.0990)
- Alikhan N-F, Zhou Z, Sergeant MJ, Achtman M. 2018 A genomic overview of the population structure of *Salmonella*. *PLoS Genet.* **14**, e1007261. (doi:10.1371/journal.pgen.1007261)
- Socransky SS, Haffajee AD. 2005 Periodontal microbial ecology. *Periodontology* **38**, 135–187. (doi:10.1111/j.1600-0757.2005.00107.x)
- Socransky SS, Haffajee AD, Cugini MA, Smith C, Kent Jr RL. 1998 Microbial complexes in subgingival plaque. *J. Clin. Periodontol.* **25**, 134–144. (doi:10.1111/j.1600-051X.1998.tb02419.x)

27. Pasolli E *et al.* 2019 Extensive unexplored human microbiome diversity revealed by over 150,000 genomes from metagenomes spanning age, geography, and lifestyle. *Cell* **176**, 649–662. (doi:10.1016/j.cell.2019.01.001)
28. Nayfach S, Shi ZJ, Seshadri R, Pollard KS, Kyrpides NC. 2019 New insights from uncultivated genomes of the global human gut microbiome. *Nature* **568**, 505–510. (doi:10.1038/s41586-019-1058-x)
29. Velsko IM, Frantz LAF, Herbig A, Larson G, Warinner C. 2018 Selection of appropriate metagenome taxonomic classifiers for ancient microbiome research. *mSystems* **3**, e00080-18. (doi:10.1128/mSystems.00080-18)
30. McIntyre ABR *et al.* 2017 Comprehensive benchmarking and ensemble approaches for metagenomic classifiers. *Genome Biol.* **18**, 182. (doi:10.1186/s13059-017-1299-7)
31. Szyrba A *et al.* 2017 Critical assessment of metagenome interpretation—a benchmark of computational metagenomics software. *Nat. Methods* **14**, 1063–1071. (doi:10.1038/nmeth.4458)
32. Zhou Z, Luhmann N, Alikhan N-F, Quince C, Achtman M. 2018 Accurate reconstruction of microbial strains from metagenomic sequencing using representative reference genomes. In *RECOMB 2018*, pp. 225–240. Cham, Switzerland: Springer.
33. Cribdon B, Ware R, Smith O, Gaffney V, Allaby RG. 2020 PIA: more accurate taxonomic assignment of metagenomic data demonstrated on sedaDNA from the North Sea. *Front. Ecol. Evol.* **8**, 84. (doi:10.3389/fevo.2020.00084)
34. Konstantinidis KT, Tiedje JM. 2005 Genomic insights that advance the species definition for prokaryotes. *Proc. Natl Acad. Sci. USA* **102**, 2567–2572. (doi:10.1073/pnas.0409727102)
35. Jain C, Rodriguez R, Phillippy AM, Konstantinidis KT, Aluru S. 2018 High throughput ANI analysis of 90 K prokaryotic genomes reveals clear species boundaries. *Nat. Commun.* **9**, 5114. (doi:10.1038/s41467-018-07641-9)
36. Li H. 2018 Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**, 3094–3100. (doi:10.1093/bioinformatics/bty191)
37. Stamatakis A. 2014 RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics* **30**, 1312–1313. (doi:10.1093/bioinformatics/btu033)
38. Zhou Z, Alikhan N-F, Sergeant MJ, Luhmann N, Vaz C, Francisco AP, Carrico JA, Achtman M. 2018 GrapeTree: visualization of core genomic relationships among 100,000 bacterial pathogens. *Genome Res.* **28**, 1395–1404. (doi:10.1101/gr.232397.117)
39. Abranches J, Zeng L, Kajfasz JK, Palmer SR, Chakraborty B, Wen ZT, Richards VP, Brady LJ, Lemos JA. 2019 Biology of oral streptococci. In *Gram positive pathogens* (eds VA Fischetti, RP Novick, JJ Ferretti, DA Portnoy, M Braunstein, JI Rood), pp. 426–434. Washington, DC: ASM Press.
40. Johansson I, Witkowska E, Kaveh B, Lif HP, Tanner AC. 2016 The microbiome in populations with a low and high prevalence of caries. *J. Dent. Res.* **95**, 80–86. (doi:10.1177/0022034515609554)
41. Oda Y, Hayashi F, Okada M. 2015 Longitudinal study of dental caries incidence associated with *Streptococcus mutans* and *Streptococcus sobrinus* in patients with intellectual disabilities. *BMC Oral Health* **15**, 102. (doi:10.1186/s12903-015-0087-6)
42. Mann AE *et al.* 2018 Differential preservation of endogenous human and microbial DNA in dental calculus and dentin. *Sci. Rep.* **8**, 9822. (doi:10.1038/s41598-018-28091-9)
43. Espinoza JL *et al.* 2018 Supragingival plaque microbiome ecology and functional potential in the context of health and disease. *MBio* **9**, e01631-18. (doi:10.1128/mBio.01631-18)
44. Shi B *et al.* 2015 Dynamic changes in the subgingival microbiome and their potential for diagnosis and prognosis of periodontitis. *MBio* **6**, e01926-14. (doi:10.1128/mbio.01926-14)
45. Liu B *et al.* 2012 Deep sequencing of the oral microbiome reveals signatures of periodontal disease. *PLoS ONE* **7**, e37919. (doi:10.1371/journal.pone.0037919)
46. McClean JS, Liu Q, Thompson J, Edlund A, Kelley S. 2015 Draft genome sequence of ‘*Candidatus Bacteroides pericalifornicus*,’ a new member of the *Bacteroidetes* phylum found within the oral microbiome of periodontitis patients. *Genome Announc.* **3**, e01485-15. (doi:10.1128/genomeA.01485-15)
47. Zhang X *et al.* 2015 The oral and gut microbiomes are perturbed in rheumatoid arthritis and partly normalized after treatment. *Nat. Med.* **21**, 895–905. (doi:10.1038/nm.3914)
48. Wang J, Jia Z, Zhang B, Peng L, Zhao F. 2019 Tracing the accumulation of *in vivo* human oral microbiota elucidates microbial community dynamics at the gateway to the GI tract. *Gut* **69**, 1355–1356. (doi:10.1136/gutjnl-2019-318977)
49. Marotz CA, Sanders JG, Zuniga C, Zaramela LS, Knight R, Zengler K. 2018 Improving saliva shotgun metagenomics by chemical host DNA depletion. *Microbiome* **6**, 42. (doi:10.1186/s40168-018-0426-3)
50. Belstrom D *et al.* 2017 Metagenomic and metatranscriptomic analysis of saliva reveals disease-associated microbiota in patients with periodontitis and dental caries. *NPJ Biofilms Microbiomes* **3**, 23. (doi:10.1038/s41522-017-0031-4)
51. Lassalle F, Spagnoletti M, Fumagalli M, Shaw L, Dyble M, Walker C, Thomas MG, Bamberg MA, Balloux F. 2018 Oral microbiomes from hunter-gatherers and traditional farmers reveal shifts in commensal balance and pathogen load linked to diet. *Mol. Ecol.* **27**, 182–195. (doi:10.1111/mec.14435)
52. Takayasu L *et al.* 2017 Circadian oscillations of microbial and functional composition in the human salivary microbiome. *DNA Res.* **24**, 261–270. (doi:10.1093/dnares/dsx001)
53. Brito IL *et al.* 2019 Transmission of human-associated microbiota along family and social networks. *Nat. Microbiol.* **4**, 964–971. (doi:10.1038/s41564-019-0409-6)
54. Franzosa EA *et al.* 2014 Relating the metatranscriptome and metagenome of the human gut. *Proc. Natl Acad. Sci. USA* **111**, E2329–E2338. (doi:10.1073/pnas.1319284111)
55. Weyrich LS *et al.* 2017 Neanderthal behaviour, diet, and disease inferred from ancient DNA in dental calculus. *Nature* **544**, 357–361. (doi:10.1038/nature21674)
56. Becht E, McInnes L, Healy J, Dutertre CA, Kwok IWH, Ng LG, Ginhoux F, Newell EW. 2019 Dimensionality reduction for visualizing single-cell data using UMAP. *Nat. Biotechnol.* **37**, 38–44. (doi:10.1038/nbt.4314)
57. Lefort V, Desper R, Gascuel O. 2015 FastME 2.0: a comprehensive, accurate, and fast distance-based phylogeny inference program. *Mol. Biol. Evol.* **32**, 2798–2800. (doi:10.1093/molbev/msv150)
58. Platt JC. 2000 Probabilities for SV machines. In *Advances in large margin classifiers* (eds AJ Smola, P Bartlett, B Schölkopf, D Schuurmans). Boston, MA: MIT Press.
59. Zhou Z, Charlesworth J, Achtman M. In press. Accurate reconstruction of bacterial pan- and core genomes with PEPPAN. *Genome Res.*
60. Simón-Soro A, Mira A. 2015 Solving the etiology of dental caries. *Trends Microbiol.* **23**, 76–82. (doi:10.1016/j.tim.2014.10.010)
61. Richards VP, Alvarez AJ, Luce AR, Bedenbaugh M, Mitchell ML, Burne RA, Nascimento MM. 2017 Microbiomes of site-specific dental plaques from children with different caries status. *Infect. Immun.* **85**, e00106–17. (doi:10.1128/IAI.00106-17)
62. Bowen WH, Burne RA, Wu H, Koo H. 2018 Oral biofilms: pathogens, matrix, and polymicrobial interactions in microenvironments. *Trends Microbiol.* **26**, 229–242. (doi:10.1016/j.tim.2017.09.008)
63. Banas JA, Drake DR. 2018 Are the mutans streptococci still considered relevant to understanding the microbial etiology of dental caries? *BMC Oral Health* **18**, 129. (doi:10.1186/s12903-018-0595-2)
64. Jonsson H, Ginolhac A, Schubert M, Johnson PL, Orlando L. 2013 mapDamage2.0: fast approximate Bayesian estimates of ancient DNA damage parameters. *Bioinformatics* **29**, 1682–1684. (doi:10.1093/bioinformatics/btt193)
65. Cornejo OE *et al.* 2013 Evolutionary and population genomics of the cavity causing bacteria *Streptococcus mutans*. *Mol. Biol. Evol.* **30**, 881–893. (doi:10.1093/molbev/mss278)
66. Fonkou MDM, Dufour J-C, Dubourg G, Raoult D. 2018 Repertoire of bacterial species cultured from the human oral cavity and respiratory tract. *Future Microbiol.* **13**, 1611–1624. (doi:10.2217/fmb-2018-0181)
67. Ryley M, Kilian M. 2008 Prevalence and distribution of principal periodontal pathogens worldwide. *J. Clin. Periodontol.* **35**, 346–361. (doi:10.1111/j.1600-051X.2008.01280.x)
68. Ondov BD, Treangen TJ, Melsted P, Mallonee AB, Bergman NH, Koren S, Phillippy AM. 2016 Mash: fast genome and metagenome distance estimation using MinHash. *Genome Biol.* **17**, 132. (doi:10.1186/s13059-016-0997-x)
69. Pedregosa F *et al.* 2011 Scikit-learn: machine learning in Python. *J. Machine Learning Res.* **12**, 2825–2830.
70. Guyon I, Weston J, Barnhill S, Vapnik V. 2002 Gene selection for cancer classification using Support Vector Machines. *Machine Learning* **46**, 389–422. (doi:10.1023/A:1012487302797)