# Limitations of using 16S rRNA microbiome sequencing to predict oral squamous cell carcinoma

CHRISTOPHER DELANEY,[1,†] CHANDRA LEKHA RAMALINGAM VEENA,[1,†] MARK C. BUTCHER,[1] WILLIAM MCLEAN,[1] SUROR MOHAMAD AHMAD SHABAN,[1] CHRISTOPHER J. NILE[2] and GORDON RAMAGE[1,*]

[1]Oral Sciences Research Group, Glasgow Dental School, School of Medicine, Dentistry and Nursing, College of Medical, Veterinary and Life Sciences, University of Glasgow, Glasgow; and [2]School of Dental Sciences, Faculty of Medical Sciences, Newcastle University, Newcastle upon Tyne, UK

Delaney C, Veena CLR, Butcher MC, McLean W, Shaban SMA, Nile CJ, Ramage G. Limitations of using 16S rRNA microbiome sequencing to predict oral squamous cell carcinoma. APMIS. 2023; 131: 262–276.

A new era of next-generation sequencing has changed our perception of the oral microbiome in health and disease, and with this there is a growing understanding that the oral microbiome is a contributing factor to oral squamous cell carcinoma (OSCC), a malignancy of the oral cavity. This study aimed to analyse the trends and relevant literature based on the 16S rRNA oral microbiome in head and neck cancer using next-generation sequencing technologies, and to conduct a meta-analysis of the studies with OSCC cases and healthy controls. A literature search using the databases Web of Science and PubMed was conducted in a scoping-like review to collect information based on the study design, and plots were generated using RStudio. We selected case–control studies using 16S rRNA oral microbiome sequencing analysis in OSCC cases versus healthy controls for re-analysis. Statistical analyses were conducted using R. Out of 916 original articles, we filtered and selected 58 studies for review, and 11 studies for meta-analysis. Differences between sampling type, DNA extraction methods, next-generation sequencing technology and region of the 16S rRNA were identified. No significant differences in the α- and β-diversity between health and oral squamous cell carcinoma were observed (p < 0.05). Random Forest classification marginally improved predictability of four studies (training set) when split 80/20. We found an increase in *Selenomonas*, *Leptotrichia* and *Prevotella* species to be indicative of disease. A number of technological advances have been accomplished to study oral microbial dysbiosis in oral squamous cell carcinoma. There is a clear need for standardization of study design and methodology to ensure 16S rRNA outputs are comparable across the discipline in the hope of identifying 'biomarker' organisms for designing screening or diagnostic tools.

Key words: Microbiome; 16S; sequencing; bioinformatics; oral cancer; oral squamous cell carcinoma.

Gordon Ramage, Oral Sciences Research Group, Glasgow Dental School, School of Medicine, Dentistry and Nursing, College of Medical, Veterinary and Life Sciences, University of Glasgow, 378 Sauchiehall Street, Glasgow G2 3JZ, UK. e-mail: gordon.ramage@glasgow.ac.uk

†Authors contributed equally.

Every year, 150 000 lives are lost to head and neck cancer which makes it the sixth most common cancer worldwide (1). Around 90% of head and neck cancers are oral squamous cell carcinoma (OSCC) that can affect oral tissues, lips, tongue, larynx and pharynx (2, 3). Tobacco, alcohol consumption, poor oral hygiene, and human papilloma virus (HPV) are major risk factors for OSCC (4). Due to the high rate of mortality and morbidity, there is an imperative need for early diagnosis through active screening of individuals at-risk (5, 6). While studies have found that early screening of pre-malignant oral disorders (PMODs) can prevent malignant transformation of oral tissues (6), there remains opportunities for improving diagnostics. Notably, chronic inflammation is implicated in disease development, a process which can be mediated by oral microorganisms (7–9), suggesting these could be potential diagnostic biomarkers.

Advances in next-generation sequencing (NGS) have paved the way to decoding the complex relationships of the microbiome, with over 700 species

thought to reside in the oral cavity (5). There has been a marked upward trend in oral microbiome research over the past 10 years, and this is also true in relation to oral cancer (10). Oral microbes can induce carcinogenic changes leading to enrichment of lipopolysaccharide (LPS) biosynthesis and epigenetic modulation causing pro-inflammatory changes in the local tumour microenvironment (5, 11). These changes are often influenced by foreign carcinogenic substances induced through smoking and tobacco by-products, which can be broken down by microbial metabolites (12).

Studies have implicated the role of *Fusobacterium*, *Pseudomonas*, *Porphyromonas*, *Provotella*, *Campylobacter*, *Rothia* and *Leptotrichia* in the progression of OSCC (10, 11, 13). These bacteria are often present in the surrounding tumour environment, however recent evidence has shown that a few species, such as *Fusobacterium nucleatum* and *Porphyromonas. gingivalis* can reside within the tumour itself (14). These 'intratumoural' bacteria play a key role in modulating immune-related changes leading to a more aggressive, enhanced tumour form (15, 16). Several studies have been published utilizing amplicon sequencing of the oral microbiome in OSCC, however, wide-scale variances in study design including sampling technique, nucleic acid extraction methods, sequencing technique and region of 16S rRNA selected for analysis vary dramatically between studies, ultimately hindering the ability to compare findings (13, 14, 17). However, it is now thought that specific members of the oral microbial community promote genetic instability, tumour proliferation and changes to the host metabolism contributing to resistance to therapy (18).

In a systematic review conducted in 2021, Mun et al. (10) concluded that there is evidence for the functional properties of the oral microbiome in OSCC, and analysis of the oral microbiome with meta-transcriptomics could further improve our understanding. The study of the oral microbiome remains a potential resource in diagnostic and therapeutic clinical intervention of OSCC. Although revolutionary, NGS has its own set of limitations, hence we hypothesised that a re-analysis of oral microbiome datasets could resolve gaps within the current research. The aim of the study was to collect and analyse publicly available datasets on 16S rRNA sequencing of the oral microbiome in OSCC using a reproducible, standardized pipeline for downloading, processing, and interpreting the data (19). Additionally, this study aims to profile the functional potential to identify and classify key organisms that may act as predictors of disease within the OSCC microbiome.

## METHODS

### Search criteria

We utilized an analysis pipeline recently developed by our group, as illustrated in Fig. 1. A full methodology has been included as supplementary material and significant deviations have been briefly listed here (19). Studies were collected by using keyword searches on PubMed and Web of Science (Clarivate Analytics, Philadelphia, Pensylvania, USA) to select for microbiome and NGS studies, performed on the oral cavity, and specifically related to oral squamous cell carcinoma (Table S1). These were filtered to exclude any study published before 2012 to coincide with the advancement of microbiome sequencing platforms (20). Remaining studies were exported to the reference manager Endnote X9 (Clarivate Analytics, Philadelphia, Pensylvania, USA).

### Study inclusion and exclusion criteria

Studies were excluded if they were not oral microbiome amplicon studies, had no data accession number, were not mappable to individual samples, had no available metadata for individuals within the study. We also excluded studies with data 'available upon reasonable request', *in vitro* studies, metagenomic studies and transcriptomic studies.

Shortlisting of studies and retention of key information was carried out by CLRV. These then underwent a two-step process for verification and inclusion by another laboratory-based clinician (SS). Shortlisted articles were then assigned a score from 1 to 5 based on data access and cohort metadata inclusion as outlined previously (19).
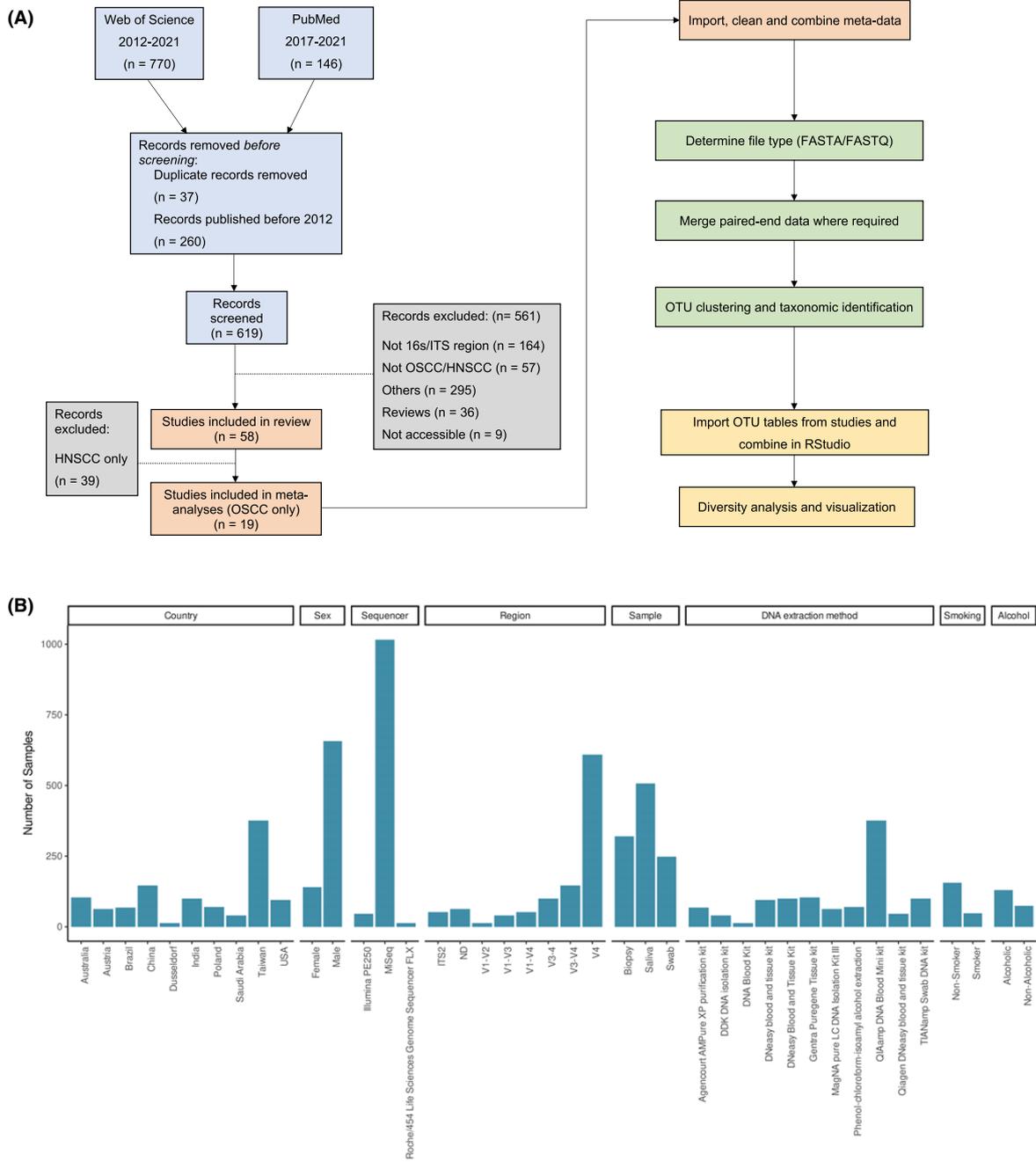
### Data retrieval and processing

Available data were downloaded from the European Nucleotide Archive (ENA) database and processed as described previously unless otherwise stated (19). In brief, Quality Control protocols were carried out in Qiime2 (21). Primers and barcodes were removed and reads trimmed if below 100 bp. Paired end reads were merged and prefiltered using sortMeRNA against the Silva-bac-id90 database (22) before being assigned to operational taxonomic units (OTUs) by mapping to the human oral microbiome database (HOMD) and GreenGenes.

Operational taxonomic units clustering was performed in closed-reference mode using the Vsearch (https://github.com/qiime2/q2-vsearch) package within the Qiime2 (v 2019.10). Phylogenetic trees were constructed for all representative OTUs using the FastTree algorithm within Qiime2 (23). OTU tables were exported and combined with study metadata tables and imported into R for manipulation and visualization.

### Microbiome diversity and composition analyses

Diversity and compositional analyses were carried out as previously described (19). Briefly, α- and β-diversity analyses were calculated using phyloseq. All samples within the case–control studies were normalized to the geometric

**Fig. 1.** Study design, collection of metadata and data processing. (A) Based on our inclusion criteria, a total of 916 studies were screened and 58 studies were included in the review and 19 eligible studies with publicly available data were included in the re-analysis. (B) Overview of study characteristics including countries, sex, sequencing technique, region of 16S rRNA, sample type, DNA extraction method/kit, smoking and alcohol status of the healthy and diseased cohort. Illumina MiSeq and V4 region of the 16S rRNA were most popular for sequencing. Saliva was the most common sampling method, and the number of samples from males, non-smokers and alcoholics was higher.

mean, significance was determined using a Welch's t-test and mean Log fold change was calculated for each of the disease vs health samples.

The centred log-ratio (CLR) Euclidean distance matrix was calculated using the make.CLR function within MicrobeR with replacement of 0 counts using the

zCompositions function and then calculating the distance matrix in base R. The Phylogenetic Isometric Log-Ratio Transform (PhILR) from the phylogenetic trees, created as described above, were visualized using the R package PhiLR (24).

The function ape::PCoA was then used to ordinate the distance matrixes into a 2D plot (25). Ggplot2 was used for visualization and combining of plots within R. Permutational multivariate analysis of variance was performed and assessed for statistical significance using the ADONIS function within the vegan package (https://github.com/vegandevs/vegan) and performed individually on all distance matrix with 999 replications.

### Random forest classifiers

As described previously, four high-impact case-controlled studies were selected as training datasets for random forest classifiers (19, 26). The data was randomly split into a 80:20 ratio to create a training and validation dataset. Centre log-ratio normalized OTUs were used as predictor variables for healthy/disease cases. CLR normalized features from the PICRUSt derived KEGG orthology (KO) feature abundances and PhILR abundances were additionally added.

Random forest models were built using the randomforest package and receiver operator characteristic (ROC) curves were built based on random forest models using the pROC and ggROC packages (27, 28). Prediction and performance metrics were extracted using the predict and performance functions from ROCR (29). The most important features were extracted from random forest models by ranking MeanDecreaseGini scores and plotted using ggplot2. Normalized OTUs from the PICRUSt derived database were used to inform on the functionality of microbiome datasets and the most important features of this model were used to identify enriched metabolic pathways by matching to the KEGG database using the clusterProfiler package in R (30). Random forest models based upon the random assignment 80:20 ratio for 80% training and 20% validation were built iteratively for each individual case-controlled study. The area under the curve was produced using the predict and performance functions for assessment of each study.

## RESULTS

### Data collection and review of publications

The initial search conducted on Web of Science yielded 770 articles and was followed up with an additional search on PubMed produced another 146 studies. The study design, as summarized in Fig. 1A, yielded a total of 58 studies relevant to the oral microbiome in OSCC and subsequently 19 were selected for the meta-analysis of publicly available data. A detailed summary of the study design parameters of all 58 studies is included in Data S1. Out of the 19 studies, 11 case-controlled studies were included in the re-analysis if that had sufficient reads >1000. The study parameters, as well as key findings, of the included 11 studies are summarized in Table 1.

We found nine studies that did not provide accessible sequencing data and were graded 0 for data access. Seventeen studies were graded 1 since the data provided could not be mapped to individual health/disease samples. Eleven studies had publicly available data with the healthy and diseased cohort described within the published article but not individually mappable to the sequencing data and were given a score of 2. A total of 21 studies had publicly available data, out of which five were given a score of 3 based as metadata was both available and mappable to sequencing data. Within the original search criteria we allowed for the inclusion of ITS or fungal amplicon based data and only one dataset remained, which was insufficient for reanalysis.

Sixteen studies were given a score of 4/5 if data were available, mappable, and additional metadata was provided (Data S1).

Overall, the metadata showed wide variances in sampling site, sampling technique, DNA extraction method, sequencing technique, region of 16S rRNA selected for analysis as summarized in Fig. 1B. Out of the total number of samples (n = 1197) included, the most common sampling technique was saliva (48%). Around 55% of the studies had selected the V4 hypervariable region of the 16 rRNA with Illumina MiSeq sequencer used for 95% of studies. DNA extraction methods varied among the individual studies, however, 31% used QIAamp DNA blood mini kit. We collected additional information regarding the alcohol consumption and smoking status among the studies (Fig. S1A,B), and other relevant information such as tumour size, stage and immune status of the enrolled OSCC subjects. Around 4% of the samples were smokers, and 10% were alcohol consumers.

### Bacterial diversity analysis

The individual samples which had <1000 total passed reads were removed and the data were plotted to show the average total reads per sample, per study. Alpha diversity indexes were implemented to represent genus- or species-level diversity within the individual samples which was statistically compared between the cohort groups. The data were analysed using the observed, Shannon, Simpson and Chao1 diversity indexes on all samples. Minor differences were observed between the healthy and OSCC samples which was not statistically significant as seen in Fig. 2. The overall alpha diversity was slightly lower in the OSCC group which was also nonsignificant. However, when looking at the alpha diversity of smokers vs non-smokers, it was slightly lower in the former group, however this was not significant. In alcoholics vs non-alcoholics, it was

**Table 1.** Study characteristics of the 11 studies included in the meta-analysis. Key observed features from all analysed studies. Including: Authorship, Number of Cases, Reference Database, Sample Site, DNA extraction method, Sequencing platform, 16S region and a summary of significant findings from each study

| Authorship | Year | Cases | Controls | Reference database | Sampling technique | DNA extraction method | Sequencer | Region | Significant findings |
|---|---|---|---|---|---|---|---|---|---|
| Al-hebshi | 2017 | 20 | 20 | HOMD | Swab | DDK DNA isolation kit | MiSeq | V1-V3 | *Fusobacterium nucleatum* subsp. *polymorphum* was the most significantly overrepresented species in the tumours followed by *Pseudomonas aeruginosa* and *Campylobacter* sp. Functional prediction showed that genes involved in bacterial mobility, flagellar assembly, bacterial chemotaxis and LPS synthesis were enriched in the tumours |
| Granato | 2021 | 16 | 8 | HOMD | Saliva | Agencourt AMPure XP purification kit | MiSeq | V4 | Relative abundances of *Centipeda*, *Veillonella* and *Gemella* suggested by metagenomics are correlated with tumour size, clinical stage and active lesion. Poor overall patient survival is associated with a higher relative abundance of *Stenophotromonas*, *Staphylococcus*, *Centipeda*, *Selenomonas*, *Alloscordovia* and *Acitenobacter* |
| Kumpitsch | 2020 | 11 | 11 | SILVA | Saliva | MagNA pure LC DNA Isolation Kit III | MiSeq | ND | Increase in *Candida* following chemoradiation, whereas the overall diversity of the microbial and fungal signatures decreased significantly after therapy |
| Lee | 2017 | 125 | 127 | SILVA | Saliva | QIAamp DNA Blood Mini kit | MiSeq | V4 | Compositions of five genera, *Bacillus*, *Enterococcus*, *Parvimonas*, *Peptostreptococcus* and *Slackia*, revealed significant differences between epithelial precursor lesion and cancer patients |
| Perera | 2018 | 27 | 25 | Species-level taxonomy assignment algorithm (BLASTN) | Biopsy | Gentra Puregene Tissue kit | MiSeq | V1-V4 | OSCC tissues had lower species richness and diversity. Genera *Capnocytophaga*, *Pseudomonas* and *Atopobium* were overrepresented in OSCC |
| Sarkar | 2021 | 50 | 50 | GreenGenes | Biopsy | DNeasy Blood and Tissue Kit | MiSeq | V3-V4 | *Prevotella*, *Corynebacterium*, *Pseudomonas*, *Deinococcus* and *Noviherbaspirillum* as significantly enriched genera, whereas genera including *Actinomyces*, *Sutterella*, *Stenotrophomonas*, *Anoxybacillus* and *Serratia* were notably decreased in the OSCC lesions |
| Schmidt | 2014 | 15 | 5 | GreenGenes | Swab | DNeasy blood and tissue kit | MiSeq | V4 | *Firmicutes* and *Actinobacteria* were significantly decreased. Weighted UniFrac principal coordinates analysis based on 12 taxa separated most cancers from other samples |
| Torralba | 2020 | 18 | 18 | SILVA | Biopsy | Phenol-chloroform-isoamyl alcohol extraction | MiSeq | V4 | Higher number of bacteria belonging to the *Fusobacteria*, *Bacteroidetes*, and *Firmicutes* phyla associated with tumour tissue. Saliva metaproteomics revealed a significant increase of *Prevotella* |

**Table 1** *(continued)*

| Authorship | Year | Cases | Controls | Reference database | Sampling technique | DNA extraction method | Sequencer | Region | Significant findings |
|---|---|---|---|---|---|---|---|---|---|
| Wolf | 2017 | 11 | 11 | GreenGenes | Saliva | MagNA pure LC DNA Isolation Kit III | MiSeq | ND | *Firmicutes, Bacteroidetes, Actinobacteria* and *Proteobacteria* were mostly detected |
| Zhang | 2020 | 50 | 50 | RDP | Swab | TIANamp Swab DNA kit | MiSeq | V3-4 | Cancer tissues were enriched in *Prevotellaceae, Fusobacteriaceae, Flavobacteriaceae, Lachnospiraceae, Peptostreptococcaceae* and *Campylobacteraceae* and 13 genera, including *Fusobacterium, Alloprevotella* and *Porphyromonas* |
| Zhou | 2020 | 24 | 24 | GreenGenes | Biopsy | Qiagen DNeasy blood and tissue kit | Illumina PE250 | V3-V4 | *Fusobacterium, Treponema, Streptococcus, Peptostreptococcus, Carnobacterium, Tannerella, Parvimonas* and *Filifactor* were enriched in OSCC |

HOMD, human oral microbiome database; RDP, Ribosomal Database Project; SILVA, Silva ribosomal RNA Gene Database Project.
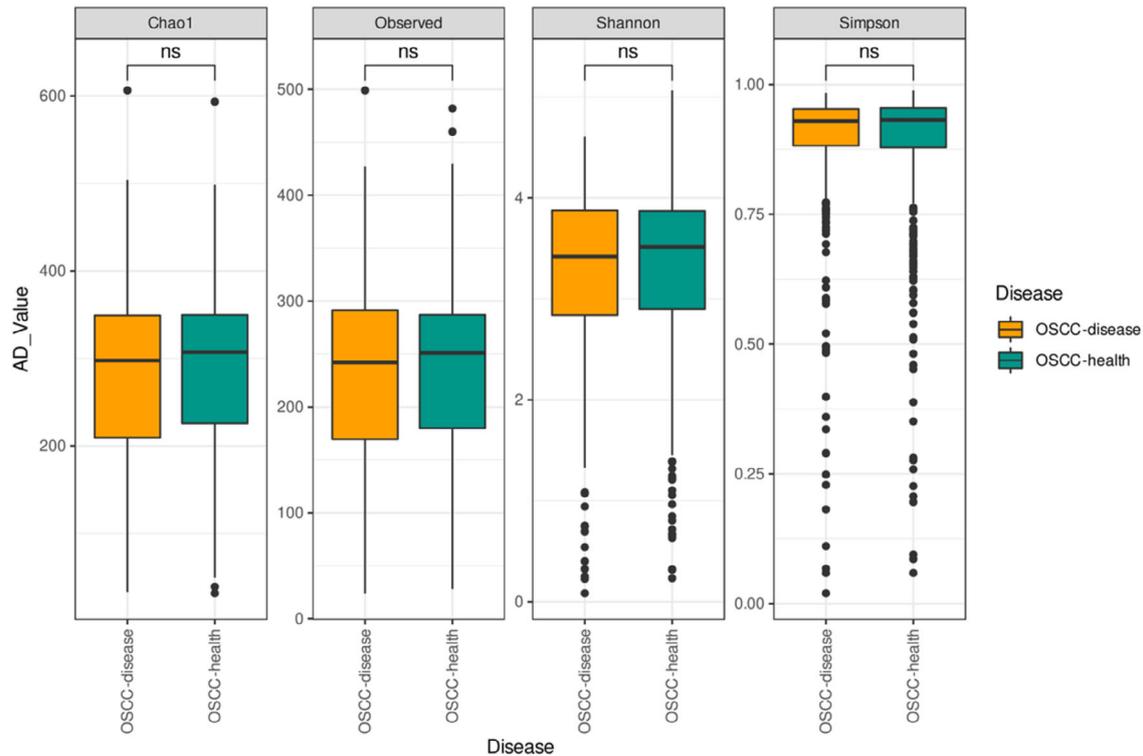
lower in the alcoholic group which was statistically significant as per Simpson index ($p < 0.001$; Fig. S1).

Beta diversity is a metric used to derive differences on a sample-to-sample basis in which diversity can be observed by clustering samples and analysing their level of dissimilarity. The resulting diversity is represented in a distance matrix from which ordination plots are generated to view patterns. Additionally permutational multivariate analysis of variance (PERMANVOVA) was tested *via* ADONIS function (Table S2). Unifrac, Bray proportion, PhiLR Euclidean and Euclidean distance matrices were utilized for determining variation in the microbiome between collated samples. Ordination plots based upon principle coordinate analysis (PCOA) were drawn and colourised to show healthy (blue) vs OSCC (pink) samples in Fig. 3A. There was some clustering observed in the PhiLR distance matrix, while statistically significant differences were observed ($p < 0.001$) across all diversity matrices when comparing health and disease groups *via* ADONIS testing. The PhiLR, weighted UniFrac and Bray-Curtis data types had an $R^2$ value of 0.024, 0.011 and 0.013 respectively. Individual studies clustered together (Fig. 3B), but become less clustered in PhiLR Euclidean ($R^2 = 0.042$) and UniFrac proportion matrices ($R^2 = 0.046$).

Additionally, we analysed the 16S rRNA sequencing region and sampling methods chosen by individual studies and observed clustering within the V4 region and saliva (Fig. S2A,B). ADONIS analysis determined these to be statistically significant with a $p < 0.001$ (Table S2).

## Predictability of OSCC diagnosis among the metadata

The 11 studies were analysed to test the predictability of OSCC samples using the receiver operating characteristic curve (ROC) from random forest classifiers. Four high-impact case-controlled studies, determined by citation number and sample volume, were selected as a training set (Fig. 4). The studies were tested across genus, species, Kegg Orthology (KO) assignment, OTUs and PhiLR. We also applied an 80/20 train-test split and calculated the area under the curve (AUC) for the training test set. The overall predictability slightly improved upon applying the 80/20 split. We found that the predictability of the KO group and the PhiLR group in the training set had an approximate AUC of 0.75. An AUC approaching 1 is a good measure of predictability, and the true positive rate of our health and OSCC samples in the training set was higher than the complete case-controlled study set.

**Fig. 2.** Boxplots of Alpha diversity indexes reflect bacterial abundance and evenness. Chao1, Observed, Shannon and Simpson diversity indexes are shown for each comparison. Higher values in the Chao, Observed and Shannon indexes indicate a higher diversity in the microbiota. The higher the value of the Simpson Index the lower the overall diversity of the microbiota. Boxplots depict the median and upper and lower quartiles of the samples grouped by healthy or OSCC diseased individuals.

In the genus, species and OTU groups, the AUC was around 0.5, and similar in both the complete and training study set.
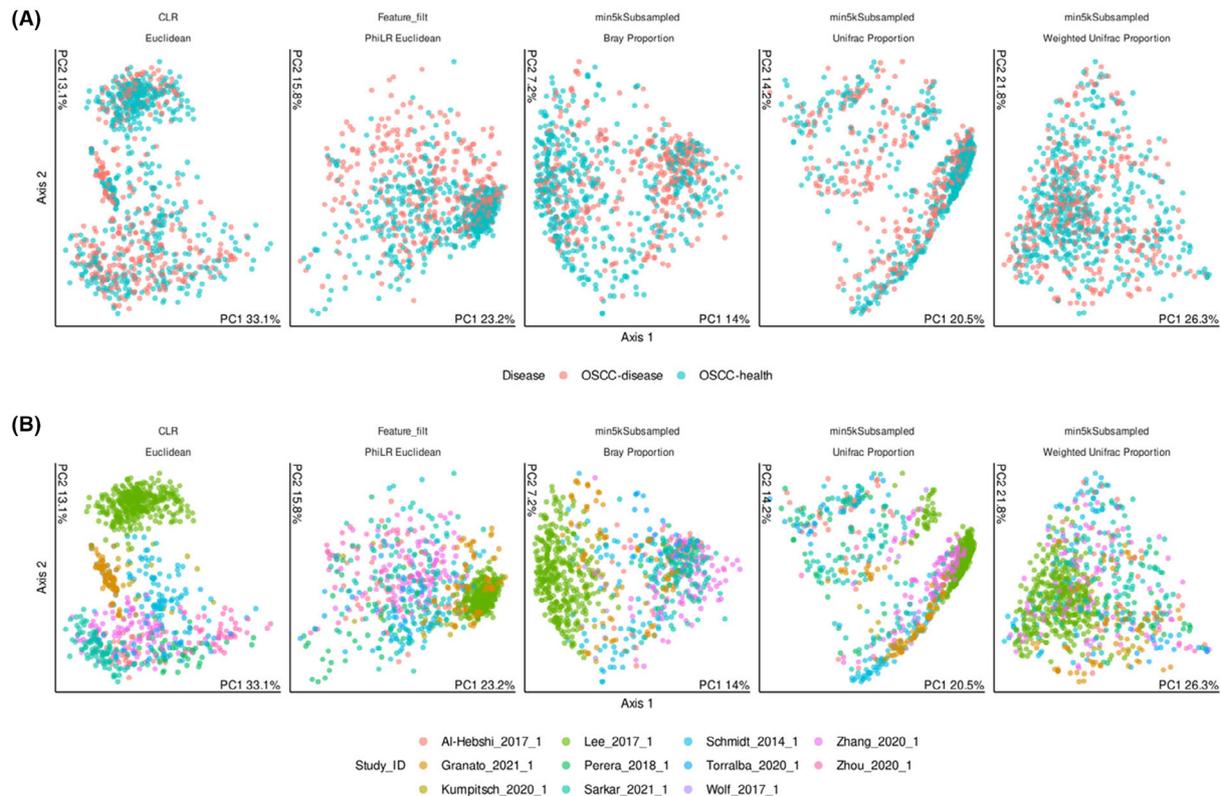
**Predictive capacity of individual studies**

The 11 individual studies were then subjected to random forest classification to test whether they could accurately predict OSCC or health. Overall, the studies Wolf et al. (48), Kumpitsch et al. (47), Zhang et al. (49) had an area under the curve (AUC) over 0.75 (Fig. 5A). Except for Granato et al. (50) which had an overall AUC of 1, Zhang et al. (49) and Kumpitsch et al. (47) had an AUC over 0.75 among genus, species, KO, OTUs and PhiLR. Wolf et al. (48) showed an AUC less than 0.75 for genus, but over 0.75 for the other parameters. These studies are the best predictors of health and disease and furthermore the alpha diversity denoted by the log2fold change was also calculated. Four studies showed a significant increase in the Simpson index, while one study showed an increase across all four indexes (Fig. 5B). Two studies

showed a significant increase in Shannon and observed indexes, while one study was upregulated in Chao1 index. However, no significant difference was observed among the other seven studies. This was calculated using Welch's t-test and a p-value <0.05 was statistically significant. Overall, only 3 of the 11 studies had good predictability, while Granato et al. (50) had an AUC of 1, which could be due to overfitting of data. The overall alpha diversity was also not consistently significantly up or downregulated between the various studies, which may mean that the probability of these studies accurately predicting health and disease is quite low.

**Distribution of organisms across health and disease**

From the random forest classifier based upon our case-controlled studies the variable importance was determined for the individual features from the model built from species-level OTUs and KEGG orthology level features (Fig. 6). The Gini coefficient was utilized to rank each variable into
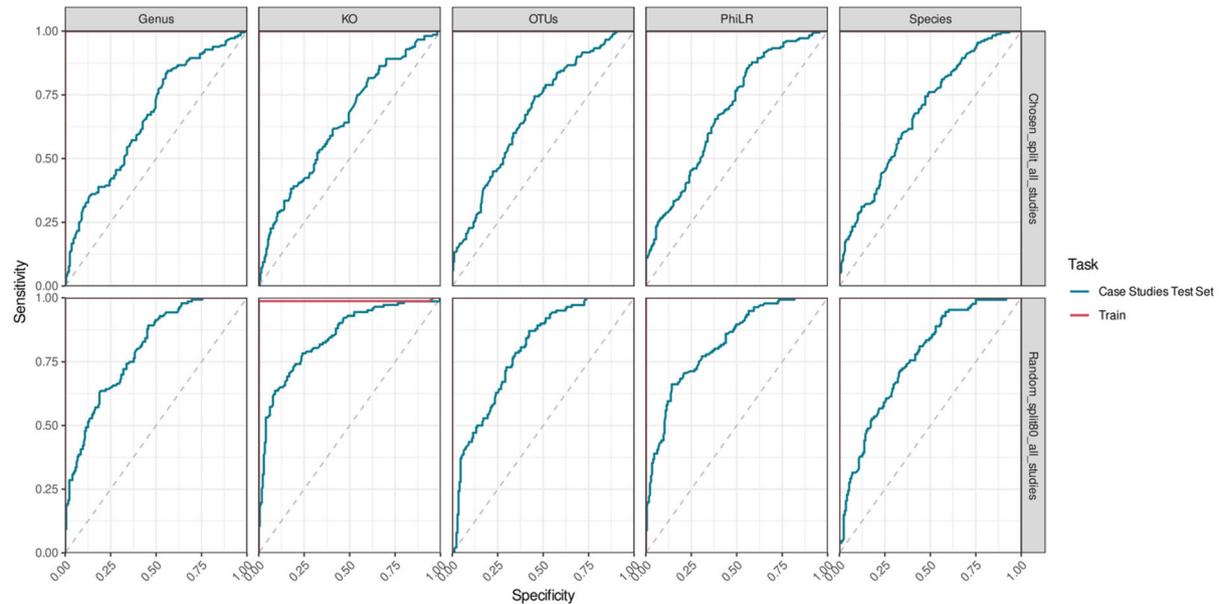
**Fig. 3.** Principal Coordinate Analysis of OSCC Samples by Disease or Study Identifier. Observed differences in beta diversity when comparing sample data from individual microbiome studies (top) or comparing health vs disease (Bottom). Sample data taken from individual studies shows clear clustering together while distance metrics in health vs disease map an overall poor correlation between identifiable features.

importance, this metric indicates the level each variable contributed to creating a strong classification tree. From the HOMD identified species we see that there a varying number of genus within our important features. Our top three species from the HOMD database with a corresponding higher level of abundance in disease according to the mean decrease in Gini were unclassified *Selemonas* sp. HMT 126, unclassified *Leptotrichia* sp. HMT 223 and *Prevotella denticola*. The three species *Acidovorax caeni*, unclassified *Actinomyces* sp. *HMT 175* and unclassified *Stomatabaculum* sp. HMT 373 were the most important features with a higher abundance in health compared to disease. LogFC were low when comparing our two groups with no features exhibiting larger that 1.5 LogFC. Similarly, our models were not able to accurately classify between health and disease using bacterial features between health and disease. This resulted in low values from our mean decrease in Gini as none of our variables were accurate predictors between health and OSCC. A similar finding was observed when using Kegg orthologies derived from our

PICRUSt analysis. Many of the important variables only exhibited a low mean decrease in Gini, and additionally the corresponding fold changes were small between our two groups of health and disease. Important features were grouped into pathways to elucidate any discriminating pathways based upon our highest scoring variables of interest (Fig. 7). We observed some significant clustering of metabolic pathways related to the two component system, nitrogen metabolism and starch and sucrose metabolism, with overrepresentation of metabolites in our disease cohort illustrated by a negative fold change in metabolic potential.

## DISCUSSION

Oral squamous cell carcinoma is an increasingly important area of oral health, so improving our diagnostic capabilities is essential. Focusing microbiologically is one critical line of travel. The purpose of our meta-analysis was to review and analyse the current trends in study design and
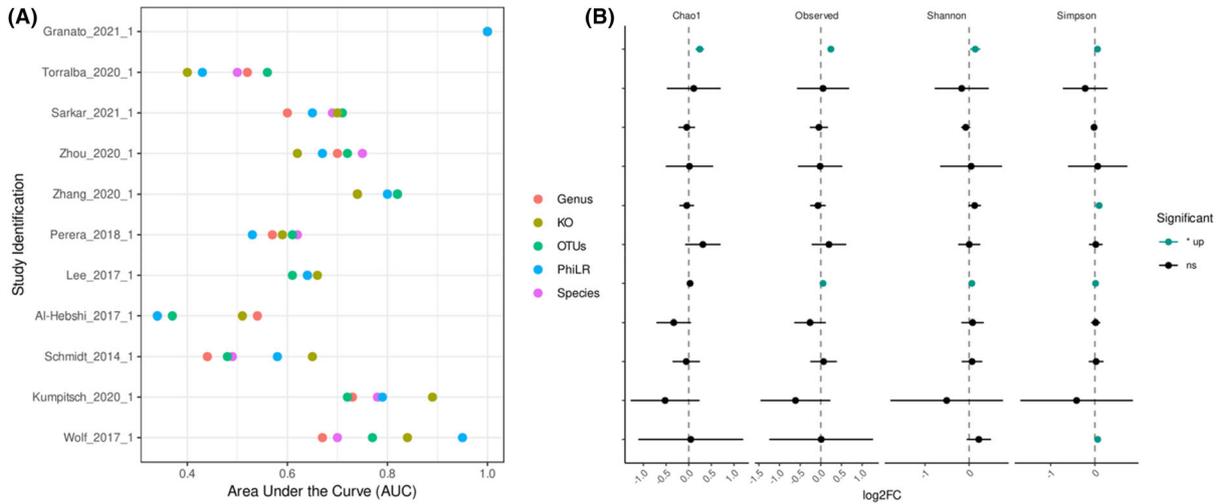
**Fig. 4.** Receiver operator curves from Random Forest based upon feature tables summarized to OTU, Genus, Species, PhiLR and KO (KEGG Orthology). Sample differentiation determined by area over curve compared for each functional group, where a value closer to 1.00 represents a clear ability to predict OSCC status. Data was trained using 4 case-controlled studies compared to all remaining studies (top) and an 80/20 split of all available sample data to training data-set (bottom).

processing of oral microbiome 16S rRNA datasets based on healthy and OSCC patient studies published so far. Overall, the study design seems to be lacking in rationale, with widespread variances in selection of the 16S hypervariable region for sequencing, DNA extraction methods and sampling method for OSCC lesions. Such differences can impact the overall microbiome under study, thus causing a misinterpretation of the oral microbial community, as evidenced by previous systematic reviews (10, 13). A key aim of this study was to identify potential biomarker organisms in OSCC, which could help map early development of the disease as an aid to diagnosis, as evidenced by previous studies conducted with the gut (31).
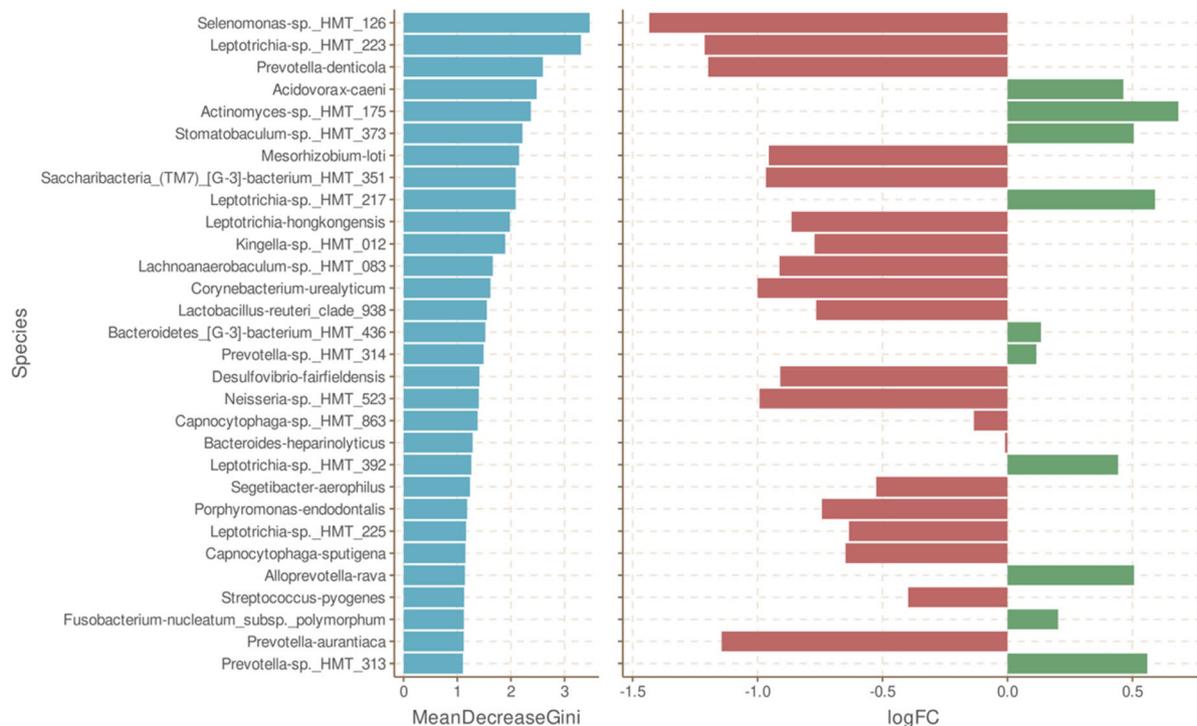
The data for the 19 studies finalized for the re-analyses were successfully downloaded and processed, however we only included 11 case-controlled studies of health vs OSCC and filtered all other samples including those with low median reads. The study characteristics included in our data collection included countries, sex, region of 16S rRNA, sample type, DNA extraction method/kit, sequencing technique, smoking and alcohol status of the healthy and diseased cohort. The selection of DNA extraction kit is crucial for production of good quality genomic DNA (gDNA). When extracting DNA from oral samples, mechanical cell lysis can improve the overall bacterial yield from

saliva (32). One study observed that enzymatic digestion increases the amount of DNA, particularly with phenol-chloroform extraction (33). Our results have shown that QIAGEN® kits have gained popularity, and it was found to yield a higher bacterial diversity, however, it may underestimate the oral microbiome (34). Several studies have found that each kit has its own flaws, and bias can be introduced at any point during processing (35). Therefore, the method of DNA extraction should be considered carefully in conjunction with the hypervariable region selected for analysis and should ideally be standardized for oral microbiome studies. NGS technologies have shorter read lengths, therefore, it is crucial to select the appropriate region of the 16 s rRNA for widespread and diverse bacterial detection (36–38).
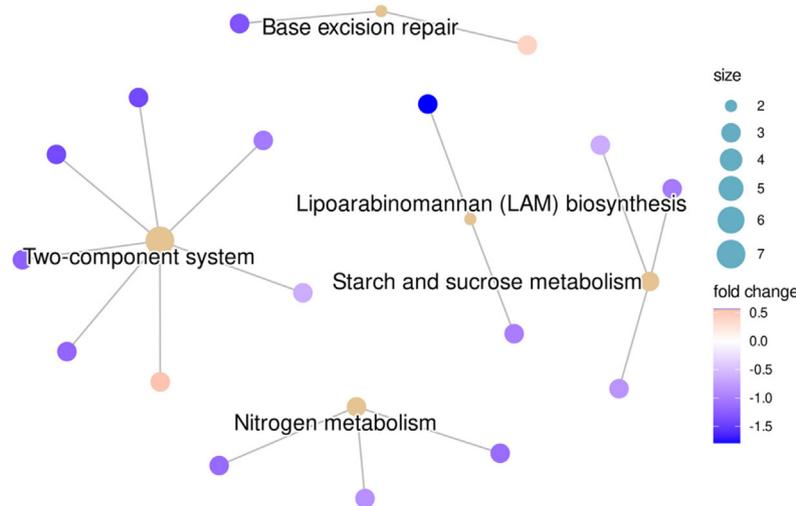
Genetic and molecular research is largely shifting to digital databases, with decreased manual handling of raw data (39). Studies have also discovered that despite discovering new and unknown species, it is not always possible for taxonomists to provide species names or phylogenetic mapping (40, 41). This could be one of the reasons why our meta-analysis showed many unclassified species of microorganisms. Few studies have explored the possibility of direct shotgun sequencing of the whole oral microbial community (metagenome) to reduce some of the bias associated with cloning and PCR (42,

**Fig. 5.** Comparisons of the microbiome in case-controlled studies. (A) Area under the curve was calculated for the performance of the random forest classifier in each of the case-controlled studies between the OSCC and Health groups. Random forest classification was performed on feature tables at OTU, genus, and species level. Additionally, it was performed on Kegg Orthology derived from PICRUST and phylogenetic transformed data (PhILR). (B) For each study, the diversity indexes Chao1, Shannon and Simpson were calculated and are represented as a Log2 fold change between disease and health. Welch's t-test was performed, and points are coloured when the significance = $p < 0.05$.



**Fig. 6.** Species of interest from random forest model. Features were selected by mean decrease in accuracy of the Gini coefficient (MeanDecreaseGini) that distinguish between health (green) and disease (Red). The differential abundance between health and disease is also represented as a log fold change and shown for the top 30 features.

**Fig. 7.** Network of important features from random forest classifier. Kegg orthology term network constructed from PICRUST identified terms classified into networks by over representation analysis. Individual nodes represent upregulated features identified in the model with their relative fold change between health and disease. A positive fold change (red) represents a higher abundance in disease and a negative (blue) fold change represents a higher abundance in disease.

43). Although 95% of the studies in our meta-analysis had used Illumina MiSeq which is a next-generation sequencing (NGS) technology with wide-scale applications, PacBio (Third generation sequencing/TGS) demonstrated much better results. It had longer sequence reads, higher species richness and can identify organisms at a more advanced taxonomic and phylogenetic range (40). However, NGS is dependable, with lower cost when performed in-house; due to its high sensitivity and specificity, the need for additional reference tests or orthogonal validation assays is avoided (44). Therefore, it is overall cost-effective, and a valuable tool, which with a few modifications can significantly improve the future of personalized management of oral cancer.

Although we utilized a standardized pipeline, we failed to find any significant differences in the overall bacterial diversity between health and OSCC. The alpha diversity was lower in the OSCC group, though the difference was small and not statistically significant. The alpha diversity of individual studies was upregulated in 4 studies: Granato et al. (50), Zhang et al. (49), Lee et al. (51), and Wolf et al. (48). Our results showed that the V4 region followed by the V3–V4 region was selected most. Studies have shown through Chao1 and ACE index that the V2-V3 region has higher richness and genetic differences when compared to other regions (36). In our meta-analysis, the V4 region was the most frequently used, yielded the most tightly clustered results, and clustered distinctly from other

regions in the beta diversity analysis. Due to the majority of studies choosing to use the V4 region, future studies should investigate whether the species coverage is high and if this region is a good predictor for microbiome studies. However, we did not observe any specific improvement in classification based upon the V4 compared to other regions. Indeed, work such as that by Johnson et al. (45) has highlighted the potential in producing accurate, high-resolution taxonomic classification of organisms *via* full-length 16S sequencing (45).

The results of our re-analysis of publicly available data returned interesting results. Random Forest analysis to generate ROC curves improved when the data was randomly split 80/20. These randomly selected variables are classified by creating decision trees, comparing the predicted values to the actual values. This can help determine the true positive/false negative values and the AUC is then measured to distinguish between the classes (46). The AUC was higher upon the 80/20 split, which shows a higher true positive rate of the randomly split data. Three of the 11 studies had an AUC over 0.75 (47–49). A significantly increased $\log_2$fold change was also observed in 4 studies (48–51).

The beta diversity analysis showed scarce clustering between health and OSCC, which signifies that there was little to no difference between the microbiome of health and disease, which is consistent with the original results of the studies included in our meta-analysis. We have shown, however, that

ADONIS analysis of PhiLR and UniFrac beta diversity matrices produced statistically significant differences between the health and OSCC groups in individual studies. This may be indicative that pooling of samples from different studies for the meta-analysis might have altered the overall results.

As evidenced by our results, saliva samples have gained popularity over the last few years. The beta diversity analysis also showed clustering of saliva samples, which was distinctly separate from the biopsy and swab samples. Some studies have found that tissue biopsies have a high concentration of *F. nucleatum* localized in both pre-cancerous and cancerous tissues (52, 53). Gopinath et al. performed an extensive study in 2021 on the different sampling types in oral cancer vs healthy tissues and found increased levels of *Prevotella*, *Campylobacter*, *Capnocytophaga*, *Solobacteria*, *Peptostreptococcus* and *Catonella* genera in oral cancer patients. They found significant differences in bacterial composition between tumour biopsies and swabs (14). These differences could be attributed to presence of biofilms or co-aggregation of bacterial pathogens on the diseased oral mucosa. Our meta-analysis showed an increase toward selection of saliva for microbiome sequencing, owing to the ease of saliva collection, storage and non-invasiveness. However, our top three species in the diseased group were *Selemonas* sp. HMT 126, unclassified *Leptotrichia* sp. HMT 223 and *Prevotella denticola*. A few studies in the meta-analysis had found similar results, with *Prevotella* sp. being the most significant overall. Since these studies have found that saliva is a partial indicator of the cancer tissue microbiome, more research is needed to consider it as a conventional method of sampling in OSCC (14).

Studies have also found that tobacco consumption in the form of chewing and smoking could alter the oral microbiome, leading to tumour progression. These parameters are extremely important and need to be included during patient data collection, which seemed to be lacking among most studies (13). Five of 11 studies had included information about smokers and alcoholics, however, it was only 4% and 10% of the overall sample size included in our meta-analysis. Our results showed a slightly higher alpha diversity in smokers in Chao1 and observed index, though this was not statistically significant. Upon further analysis, we found significant differences between alcoholics and non-alcoholics. Alcoholics may have an altered microbiome due to an overproduction of acetaldehyde. High levels of acetaldehyde producing bacteria like *Actinomyces*, *Rothia*, *Streptococcus* and *Prevotella* have been isolated from the oral cavity of chronic alcoholics (13). Indeed, studies, such as those conducted by Mizumoto et al., in 2017 highlight the relationship between acetaldehyde and tumour mutagenesis (54). This would imply that selection for these acetaldehyde producing bacteria could induce further mutations and progression of the tumour. This is further reinforced by our identification of *Prevotella* as a top-three organism present in disease samples during our meta-analysis. Notably, the majority of studies did not consider the yeast *Candida albicans*, an important determinant of oral cancer (55). From our literature screening we also identified one ITS focused paper with available data. We found significantly more 16S rRNA microbiome studies within the literature search, with the majority of studies focusing on the bacterial involvement in OSCC. Therefore, future studies should examine the mycobiome as well as the microbiome in any microbiological investigative studies. Although not included in our search criteria it is noteworthy that there are additional considerations other than the microbiome and mycobiome, including the phageome, virome and meta-transcriptome with the oral environment (56).

Microbiome sequencing utilizing 16S rRNA amplicon also have limitation compared to the more holistic metagenomic shotgun sequencing. Due to limitations of integration methods within our standardized protocol for analysing microbiome data we were unable to include these data types within our study design. However, at the time of writing the majority of studies utilizing the shotgun approach were minimal in comparison to those that used 16S microbiome. Within this study we highlight that although this is the predominant form of profiling the oral microenvironment in OSCC there are other technological developments to be considered. Metagenomics by shotgun has been applied to profile other medically relevant conditions within the human body, including the gut and vaginal microbial communities and other inflammatory oral diseases (57–59). Due to its cost effectiveness, speed of preparation and comparative ease of analysis, amplicon sequencing remains the most popular platform for microbial profiling. Within time, and increase of processing simplicity, it may become preferable to utilize shotgun sequencing. However, for 16S amplicon sequencing to become a useful clinical tool the need for standardization is imminently desirable.

## CONCLUSION

Despite these advancements, the future of oral microbiome research in OSCC is highly dependent on study design characterization. Several systematic reviews conducted over the past few years have concluded that oral microbiome studies must widely

focus on a standard study design, collect essential and adequate metadata, and follow proper pipelines for analyses of the data (10, 13). We have determined that it may be difficult to rely upon microbiome studies for this purpose, due to a lack of specificity, and wide variation between individual study design and outcomes. A consensus in approach is a basic requirement before these studies can be collectively useful. However, in understanding the essential marriage of microbiome and clinical metadata when conducting these analyses, we can also postulate that a targeted multi-omic approach to sample analysis may provide a more promising outcome for early diagnosis of OSCC (60). The oral microbiome had great potential in classification of diseases, including diagnosis and prevention. Future studies could include a more clear and standardized technique for analysis, utilizing the vast number of technological advances in the databases available as a predictive tool in OSCC by applying our knowledge of the microbiome into useful clinical applications.

## CONFLICT OF INTEREST

The authors declare no conflict of interest.

## REFERENCES

1. Bray F, Ferlay J, Soerjomataram I, Siegel RL, Torre LA, Jemal A. Global cancer statistics 2018: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries. CA Cancer J Clin. 2018;68:394–424.
2. Boot A, Ng AWT, Chong FT, Ho S-C, Yu W, Tan DSW, et al. Characterization of colibactin-associated mutational signature in an Asian oral squamous cell carcinoma and in other mucosal tumor types. Genome Res. 2020;30:803–13.
3. Bebek G, Bennett KL, Funchain P, Campbell R, Seth R, Scharpf J, et al. Microbiomic subprofiles and MDR1 promoter methylation in head and neck squamous cell carcinoma. Hum Mol Genet. 2012;21:1557–65.
4. Ram H, Sarkar J, Kumar H, Konwar R, Bhatt ML, Mohammad S. Oral cancer: risk factors and molecular pathogenesis. J Maxillofac Oral Surg. 2011;10:132–7.
5. Li Q, Hu Y, Zhou X, Liu S, Han Q, Cheng L. Role of oral bacteria in the development of oral squamous cell carcinoma. Cancer. 2020;12:2797.
6. Sankaranarayanan R, Ramadas K, Thomas G, Muwonge R, Thara S, Mathew B, et al. Effect of screening on oral cancer mortality in Kerala, India: a cluster-randomised controlled trial. Lancet. 2005;365: 1927–33.
7. Rodriguez-Rabassa M, Lopez P, Rodriguez-Santiago RE, Cases A, Felici M, Sanchez R, et al. Cigarette smoking modulation of saliva microbial composition and cytokine levels. Int J Environ Res Public Health. 2018;15:2479.
8. Al-hebshi NN, Alharbi FA, Mahri M, Chen T. Differences in the Bacteriome of smokeless tobacco products with different oral carcinogenicity: compositional and predicted functional analysis. Genes (Basel). 2017;8: 106.
9. Hsiao J-R, Chang C-C, Lee W-T, Huang C-C, Ou C-Y, Tsai S-T, et al. The interplay between oral microbiome, lifestyle factors and genetic polymorphisms in the risk of oral squamous cell carcinoma. Carcinogenesis. 2018;39:778–87.
10. Su Mun L, Wye Lum S, Kong Yuiin Sze G, Hock Yoong C, Ching Yung K, Kah Lok L, et al. Association of microbiome with oral squamous cell carcinoma: a systematic review of the metagenomic studies. Int J Environ Res Public Health. 2021;18:7224.
11. Arthur RA, Bezerra RS, Bianchi Ximenez JP, Merlin BL, Morraye RA, Neto JV, et al. Microbiome and oral squamous cell carcinoma: a possible interplay on iron metabolism and its impact on tumor microenvironment. Braz J Microbiol. 2021;52:1287–302.
12. Rai AK, Panda M, Das AK, Rahman T, Das R, Das K, et al. Dysbiosis of salivary microbiome and cytokines influence oral squamous cell carcinoma through inflammation. Arch Microbiol. 2021;203:137–52.
13. Ramos RT, Sodré CS, de Sousa Rodrigues P, da Silva AMP, Fuly MS, Dos Santos HF, et al. High-throughput nucleotide sequencing for bacteriome studies in oral squamous cell carcinoma: a systematic review. Oral Maxillofac Surg. 2020;24:387–401.
14. Gopinath D, Menon RK, Wie CC, Banerjee M, Panda S, Mandal D, et al. Differences in the bacteriome of swab, saliva, and tissue biopsies in oral cancer. Sci Rep. 2021;11(1):1181.
15. Chen Q, Shao Z, Liu K, Zhou X, Wang L, Jiang E, et al. Salivary *Porphyromonas gingivalis* predicts outcome in oral squamous cell carcinomas: a cohort study. BMC Oral Health. 2021;21:228.
16. Neuzillet C, Marchais M, Vacher S, Hilmi M, Schnitzler A, Meseure D, et al. Prognostic value of intratumoral *Fusobacterium nucleatum* and association with immune-related gene expression in oral squamous cell carcinoma patients. Sci Rep. 2021;11:7870.
17. Di Bella JM, Bao Y, Gloor GB, Burton JP, Reid G. High throughput sequencing methods and analysis for microbiome research. J Microbiol Methods. 2013;95: 401–14.
18. Huybrechts I, Zouiouich S, Loobuyck A, Vandenbulcke Z, Vogtmann E, Pisanu S, et al. The human microbiome in relation to cancer risk: a systematic review of epidemiologic studies. Cancer Epidemiol Biomarkers Prev. 2020;29:1856–68.
19. Butcher MC, Short B, Veena CLR, Bradshaw D, Pratten JR, McLean W, et al. Meta-analysis of caries microbiome studies can improve upon disease prediction outcomes. APMIS. 2022;130:763–77.
20. Shokralla S, Spall JL, Gibson JF, Hajibabaei M. Next-generation sequencing technologies for environmental DNA research. Mol Ecol. 2012;21:1794–805.

21. Bolyen E, Rideout JR, Dillon MR, Bokulich NA, Abnet CC, Al-Ghalith GA, et al. Reproducible, interactive, scalable and extensible microbiome data science using QIIME 2. Nat Biotechnol. 2019;37:852–7.

22. Kopylova E, Noé L, Touzet H. SortMeRNA: fast and accurate filtering of ribosomal RNAs in metatranscriptomic data. Bioinformatics. 2012;28:3211–7.

23. Price MN, Dehal PS, Arkin AP. FastTree 2 – approximately maximum-likelihood trees for large alignments. PLoS ONE. 2010;5:e9490.

24. Silverman JD, Washburne AD, Mukherjee S, David LA. A phylogenetic transform enhances analysis of compositional microbiota data. Elife. 2017;6:e21887.

25. Paradis E, Schliep K. Ape 5.0: an environment for modern phylogenetics and evolutionary analyses in R. Bioinformatics. 2018;35:526–8.

26. Chen T, Marsh PD, Al-Hebshi NN. SMDI: an index for measuring subgingival microbial dysbiosis. J Dent Res. 2022;101:331–8.

27. Breiman L. Random forests. Mach Learn. 2001;45:5–32.

28. Robin X, Turck N, Hainard A, Tiberti N, Lisacek F, Sanchez JC, et al. pROC: an open-source package for R and S+ to analyze and compare ROC curves. BMC Bioinformatics. 2011;12:77.

29. Sing T, Sander O, Beerenwinkel N, Lengauer T. ROCR: visualizing classifier performance in R. Bioinformatics. 2005;21:3940–1.

30. Wu T, Hu E, Xu S, Chen M, Guo P, Dai Z, et al. clusterProfiler 4.0: a universal enrichment tool for interpreting omics data. The. Innovation. 2021;2:100141.

31. Ren Z, Li A, Jiang J, Zhou L, Yu Z, Lu H, et al. Gut microbiome analysis as a tool towards targeted noninvasive biomarkers for early hepatocellular carcinoma. Gut. 2019;68:1014–23.

32. Vesty A, Biswas K, Taylor MW, Gear K, Douglas RG. Evaluating the impact of DNA extraction method on the representation of human Oral bacterial and fungal communities. PLoS One. 2017;12:e0169877.

33. Rosenbaum J, Usyk M, Chen Z, Zolnik CP, Jones HE, Waldron L, et al. Evaluation of Oral cavity DNA extraction methods on bacterial and fungal microbiota. Sci Rep. 2019;9:1531.

34. Zhou X, Nanayakkara S, Gao J-L, Nguyen K-A, Adler CJ. Storage media and not extraction method has the biggest impact on recovery of bacteria from the oral microbiome. Sci Rep. 2019;9(1):14968.

35. Sergaki C, Anwar S, Fritzsche M, Mate R, Francis RJ, MacLellan-Gibson K, et al. Developing whole cell standards for the microbiome field. Microbiome. 2022;10:123.

36. Bukin YS, Galachyants YP, Morozov IV, Bukin SV, Zakharenko AS, Zemskaya TI. The effect of 16S rRNA region choice on bacterial community metabarcoding results. Sci Data. 2019;6:190007.

37. Claesson MJ, Wang Q, O'Sullivan O, Greene-Diniz R, Cole JR, Ross RP, et al. Comparison of two next-generation sequencing technologies for resolving highly complex microbiota composition using tandem variable 16S rRNA gene regions. Nucleic Acids Res. 2010;38:e200.

38. Hamady M, Knight R. Microbial community profiling for human microbiome projects: tools, techniques, and challenges. Genome Res. 2009;19:1141–52.

39. de la Iglesia D, García-Remesal M, de la Calle G, Kulikowski C, Sanz F, Maojo V. The impact of computer science in molecular medicine: enabling high-throughput research. Curr Top Med Chem. 2013;13(5):526–75.

40. Zhang J, Su L, Wang Y, Deng S. Improved high-throughput sequencing of the human Oral microbiome: from Illumina to PacBio. Can J Infect Dis Med Microbiol. 2020;2020:6678872.

41. Eren AM, Borisy GG, Huse SM, Mark Welch JL. Oligotyping analysis of the human oral microbiome. Proc Natl Acad Sci USA. 2014;111:E2875–E84.

42. Xie G, Chain PS, Lo CC, Liu KL, Gans J, Merritt J, et al. Community and gene composition of a human dental plaque microbiota obtained by metagenomic sequencing. Mol Oral Microbiol. 2010;25:391–405.

43. Belda-Ferre P, Alcaraz LD, Cabrera-Rubio R, Romero H, Simón-Soro A, Pignatelli M, et al. The oral metagenome in health and disease. ISME J. 2012;6:46–56.

44. Tan O, Shrestha R, Cunich M, Schofield DJ. Application of next-generation sequencing to improve cancer management: a review of the clinical effectiveness and cost-effectiveness. Clin Genet. 2018;93(3):533–44.

45. Johnson JS, Spakowicz DJ, Hong B-Y, Petersen LM, Demkowicz P, Chen L, et al. Evaluation of 16S rRNA gene sequencing for species and strain-level microbiome analysis. Nat Commun. 2019;10:5029.

46. Hajian-Tilaki K. Receiver operating characteristic (ROC) curve analysis for medical diagnostic test evaluation. Caspian J Intern Med. 2013;4:627–35.

47. Kumpitsch C, Moissl-Eichinger C, Pock J, Thurnher D, Wolf A. Preliminary insights into the impact of primary radiochemotherapy on the salivary microbiome in head and neck squamous cell carcinoma. Sci Rep. 2020;10:16582.

48. Wolf A, Moissl-Eichinger C, Perras A, Koskinen K, Tomazic PV, Thurnher D. The salivary microbiome as an indicator of carcinogenesis in patients with oropharyngeal squamous cell carcinoma: a pilot study. Sci Rep. 2017;7:5867.

49. Zhang L, Liu Y, Zheng HJ, Zhang CP. The oral microbiota may have influence on oral cancer. Front Cell Infect Microbiol. 2020;9:476.

50. Granato DC, Neves LX, Trino LD, Carnielli CM, Lopes AFB, Yokoo S, et al. Meta-omics analysis indicates the saliva microbiome and its proteins associated with the prognosis of oral cancer patients. Biochim Biophys Acta Proteins Proteom. 2021;1869:140659.

51. Lee WH, Chen HM, Yang SF, Liang C, Peng CY, Lin FM, et al. Bacterial alterations in salivary microbiota and their association in oral cancer. Sci Rep. 2017;7:16540.

52. Healy CM, Moran GP. The microbiome and oral cancer: more questions than answers. Oral Oncol. 2019;89:30–3.

53. Whitmore SE, Lamont RJ. Oral bacteria and cancer. PLoS Pathog. 2014;10:e1003933.

54. Mizumoto A, Ohashi S, Hirohashi K, Amanuma Y, Matsuda T, Muto M. Molecular mechanisms of acetaldehyde-mediated carcinogenesis in squamous epithelium. Int J Mol Sci. 2017;18:1943.

55. Gainza-Cirauqui ML, Nieminen MT, Novak Frazer L, Aguirre-Urizar JM, Moragues MD, Rautemaa R.

Production of carcinogenic acetaldehyde by *Candida albicans* from patients with potentially malignant oral mucosal disorders. J Oral Pathol Med. 2013;42:243–9.

56. Ganly I, Hao Y, Rosenthal M, Wang H, Migliacci J, Huang B, et al. Oral microbiome in nonsmoker patients with oral cavity squamous cell carcinoma, defined by metagenomic shotgun sequencing. Cancers (Basel). 2022;14:6096.

57. Bai X, Narayanan A, Nowak P, Ray S, Neogi U, Sönnerborg A. Whole-genome metagenomic analysis of the gut microbiome in HIV-1-infected individuals on antiretroviral therapy. Front Microbiol. 2021;12:667718.

58. France MT, Fu L, Rutt L, Yang H, Humphrys MS, Narina S, et al. Insight into the ecology of vaginal bacteria through integrative analyses of metagenomic and metatranscriptomic data. Genome Biol. 2022;23:66.

59. Farina R, Severi M, Carrieri A, Miotto E, Sabbioni S, Trombelli L, et al. Whole metagenomic shotgun sequencing of the subgingival microbiome of diabetics and non-diabetics with different periodontal conditions. Arch Oral Biol. 2019;104:13–23.

60. Hasin Y, Seldin M, Lusis A. Multi-omics approaches to disease. Genome Biol. 2017;18:83.

## SUPPORTING INFORMATION

Additional supporting information may be found online in the Supporting Information section at the end of the article.

**Fig. S1.** Additional analysis of alpha diversity examining alcohol and smoking status in health and disease. Chao1, Observed, Shannon and Simpson diversity indexes are shown for each comparison. (A) Significantly lower median values (p < 0.001) for non-alcoholic samples in Chao1, Observed, Shannon and Simpson indexes reveal lower microbial diversity when compared with all other samples. (B) Lower median values (non-significant) for smokers' samples in Observed and Shannon indexes reveal lower microbial diversity when compared with all other samples.

**Fig. S2.** Principal Coordinate Analysis of OSCC samples depicting beta diversity of microbial population. (A) Amplified 16 s rRNA region sequenced and (B) Different sampling methods used by individual studies. Overall, clustering of samples based on sequence region is most prominently identifiable when observed using Euclidian, Bray-Curtis and Unifrac similarity metrics. When comparing health and disease, sample clustering is only marginally distinguishable using Euclidian metrics.

**Table S1.** Search terms used on Web of Science and PubMed for literature search (2012-2021).

**Table S2.** Table displaying results from ADONIS testing for the entire dataset.

**Data S1.** Meta-data table of clinical studies and study specific information.