

1

Doctoral Thesis

2

Metagenomic Profiling

3

in Preterm Birth, Periodontitis, and Colorectal Cancer

4

Jaewoong Lee

5

Department of Biomedical Engineering

6

Ulsan National Institute of Science and Technology

7

2025

⁸

Metagenomic Profiling

⁹

in Preterm Birth, Periodontitis, and Colorectal Cancer

¹⁰

Jaewoong Lee

¹¹

Department of Biomedical Engineering

¹²

Ulsan National Institute of Science and Technology

Metagenomic Profiling in Preterm Birth, Periodontitis, and Colorectal Cancer

A thesis/dissertation submitted to
Ulsan National Institute of Science and Technology
in partial fulfillment of the
requirements for the degree of
Doctor of Philosophy

Jaewoong Lee

04.16.2025 of submission

Approved by

Advisor

Semin Lee

Metagenomic Profiling in Preterm Birth, Periodontitis, and Colorectal Cancer

Jaewoong Lee

This certifies that the thesis/dissertation of Jaewoong Lee is approved.

04.16.2025 of submission

Signature

Advisor: Semin Lee

Signature

Taejoon Kwon

Signature

Eunhee Kim

Signature

Kyemyung Park

Signature

Min Hyuk Lim

Abstract

16 The human microbiome plays a critical role in diseases, influencing immune response, metabolism,
17 and disease progression. Recent advances in microbiome sequencing techniques have highlighted its
18 potential as a diagnostic, prognostic, and therapeutic strategies in various diseases, including preterm
19 birth (Section 2), periodontitis (Section 3), and colorectal cancer (Section 4). Dysbiosis, characterized
20 by alterations in microbiome composition, has been linked to pathogenesis, disease progression, and
21 treatment outcome, emphasizing the need for comprehensive metagenomic analyses. By investigating
22 microbiome profiling, researchers can uncover microbial biomarkers and host-microbiome interactions
23 that contribute to underlying mechanisms of disease. Thus, understanding these complex relationships
24 not only enhances early detection and risk stratification but also paves the way for microbiome-based
25 therapeutic interventions and personalized medicine strategies. Ultimately, as microbiome research
26 continues to evolve, its integration with genomics, metabolomics, and immunology suggests promise for
27 transforming disease management and improving treatment outcomes.

28 Section 2 investigated the association between the prenatal salivary microbiome and preterm birth
29 (PTB) using 16S ribosomal RNA (rRNA) gene sequencing and developed a random forest-based prediction
30 model for risk of preterm birth. A total of 59 pregnant women were included as the study participants, with
31 30 in the preterm birth group and 29 in the full-term birth (FTB) group. Salivary microbiome samples were
32 collected via mouthwash within 25 hours before delivery, and 16S rRNA gene sequencing was performed
33 to analyze microbial taxonomic composition. Differentially abundant taxa (DAT) were identified by
34 DESeq2, revealing the 25 significant taxa, including three PTB-enriched taxa and 22 FTB-enriched taxa,
35 suggesting distinct microbial differences between the two groups. A random forest classifier was applied
36 to predict PTB risk based on salivary microbiome composition, achieving the high balanced accuracy
37 (0.765 ± 0.071) using the nine most important taxa. These findings indicate that salivary microbiome
38 profiling may serve as a novel predictive tool for PTB risk assessment, complementing existing clinical
39 predictors.

40 Section 3 characterized salivary microbiome compositions to classify periodontal health and different
41 stages of periodontitis using 16S rRNA gene sequencing. A total of 250 study participants were included,
42 comprising 100 periodontally healthy controls and 150 periodontitis patients equally classified into
43 stage I, stage II, and stage III. Microbial diversity indices were calculated, and ANCOM was used to
44 identify 20 differentially abundant taxa among the multiple periodontitis stages. A random forest machine
45 learning model was developed to classify periodontitis stages based on the proportions of differentially
46 abundant taxa, achieving an area-under-curve of 0.870 ± 0.079 (mean \pm SD). Among the identified dif-
47 ferentially abundant taxa, *Porphyromonas gingivalis* and *Actinomyces* spp. were the most important
48 features in distinguishing periodontitis stages. Random forest classifier also effectively distinguished
49 healthy individuals from stage I periodontitis with an area-under-curve of 0.852 ± 0.103 (mean \pm SD)
50 and detected periodontitis patients from healthy controls with an area-under-curve of 0.953 ± 0.049 .

51 (mean \pm SD). External validation with Spanish and Portuguese datasets showed a slight performance
52 decrease, likely due to ethnic variations in salivary microbiome composition, emphasizing the need for
53 population-specific models. These findings suggest that salivary microbiome composition profiling may
54 serve as a non-invasive diagnostic technique for periodontitis, aiding in early detection and personalized
55 dental care.

56 Section 4 conducted a comprehensive metagenomic analysis of colorectal cancer using PathSeq,
57 focusing on key clinical outcomes, including recurrence history and overall survival duration. Significant
58 differences in alpha-diversity and beta-diversity indices were observed between tumor and its adjacent
59 normal tissues, with further stratification revealing distinct microbial diversity patterns associated with
60 recurrence status and survival outcomes. Differentially abundant taxa were identified, highlighting
61 microbial signatures may influence CRC progression and prognosis. To evaluate the predictive potential
62 of these selected differentially abundant taxa, we developed a random forest-based machine learning model
63 for CRC recurrence risk and survival duration. While the classification model for recurrence prediction
64 achieved moderate accuracy (0.570 ± 0.164 , mean \pm SD), and the regression model of OS duration showed
65 moderated errors (729.302 ± 179.940 , mean \pm SD), these results suggest that gut microbiome composition
66 alone may not be sufficient for personalized clinical predictions. These findings emphasize the need for
67 multi-omics integration, combining host genomic alterations, *e.g.* somatic and germline mutations, with
68 gut microbiome compositions, to improve CRC risk stratification and personalized medicine applications.
69 This study highlights the potential role of gut microbiome for biomarkers in CRC diagnosis and prognosis
70 while underscoring the complexity of host-microbiome interaction in CRC progression.

71 Together, these studies demonstrate the clinical relevance of microbiome profiling in three distinct
72 yet interconnected diseases by analyzing microbial diversity, identifying differentially abundant taxa, and
73 leveraging machine learning for predictive modeling. While each condition exhibited unique microbial
74 signatures, the findings collectively underscore the broader impact of dysbiosis on pathogenesis and
75 disease progression. These results suggest that microbial biomarkers could serve as valuable tools for
76 early detection, risk assessment, and personalized medicine strategies across multiple disease contexts.
77 However, the predictive performance of machine learning models highlights the requirement for multi-
78 omics integration, incorporating host genomic data to improve the accuracy of disease prediction and
79 personalized therapeutic interventions. Moving forward, further large-scale and multi-cohort validation
80 studies will be essential to refine microbiome-based biomarkers and ensure their clinical applicability in
81 therapeutic guidance. By deepening our understanding of host-microbiome interactions, this dissertation
82 contributes to the growing field of microbiome-driven personalized medicine, paving a novel approaches
83 in disease prevention and management.

84
85 **This doctoral dissertation is an addition based on the following papers that the author has already
86 published:**

- 87 • Hong, Y. M., **Lee, Jaewoong**, Cho, D. H., Jeon, J. H., Kang, J., Kim, M. G., ... & Kim, J. K. (2023).
88 Predicting preterm birth using machine learning techniques in oral microbiome. *Scientific Reports*,
89 13(1), 21105.

Contents

91	1	Introduction	1
92	2	Predicting preterm birth using random forest classifier in salivary microbiome	8
93	2.1	Introduction	8
94	2.2	Materials and methods	10
95	2.2.1	Study design and study participants	10
96	2.2.2	Clinical data collection and grouping	10
97	2.2.3	Salivary microbiome sample collection	10
98	2.2.4	16s rRNA gene sequencing	10
99	2.2.5	Bioinformatics analysis	11
100	2.2.6	Data and code availability	11
101	2.3	Results	12
102	2.3.1	Overview of clinical information	12
103	2.3.2	Comparison of salivary microbiomes composition	12
104	2.3.3	Random forest classification to predict PTB risk	12
105	2.4	Discussion	20
106	3	Random forest prediction model for periodontitis stages based on the salivary microbiomes	22
107	3.1	Introduction	22
108	3.2	Materials and methods	24
109	3.2.1	Study participants enrollment	24
110	3.2.2	Periodontal clinical parameter diagnosis	24
111	3.2.3	Saliva sampling and DNA extraction procedure	26
112	3.2.4	Bioinformatics analysis	26
113	3.2.5	Data and code availability	27
114	3.3	Results	29

115	3.3.1	Summary of clinical information and sequencing data	29
116	3.3.2	Diversity indices reveal differences among the periodontitis severities .	29
117	3.3.3	DAT among multiple periodontitis severities and their correlation . .	29
118	3.3.4	Classification of periodontitis severities by random forest models . .	30
119	3.4	Discussion	51
120	4	Metagenomic signature analysis of Korean colorectal cancer	55
121	4.1	Introduction	55
122	4.2	Materials and methods	57
123	4.2.1	Study participants enrollment	57
124	4.2.2	DNA extraction procedure	57
125	4.2.3	Bioinformatics analysis	57
126	4.2.4	Data and code availability	59
127	4.3	Results	60
128	4.3.1	Summary of clinical characteristics	60
129	4.3.2	Gut microbiome compositions	60
130	4.3.3	Diversity indices	61
131	4.3.4	DAT selection	62
132	4.3.5	Random forest prediction	64
133	4.4	Discussion	82
134	5	Conclusion	88
135	References	90
136	Acknowledgments	109

137

List of Figures

138	1	DAT volcano plot for PTB prediction	14
139	2	Salivary microbiome compositions over DAT for PTB prediction	15
140	3	Random forest-based PTB prediction model	16
141	4	Diversity indices about PTB study participants	17
142	5	PROM-related DAT between FTB and PTB	18
143	6	Validation of random forest-based PTB prediction model	19
144	7	Diversity indices for periodontitis	37
145	8	DAT for periodontitis	38
146	9	Correlation heatmap between periodontitis DAT	39
147	10	Random forest classification metrics for periodontitis prediction	40
148	11	Random forest classification metrics from external datasets	41
149	12	Rarefaction curves for alpha-diversity indices	42
150	13	Salivary microbiome compositions in the different periodontal stages	43
151	14	Correlation plots for periodontitis DAT	44
152	15	Clinical measurements by the periodontitis stages	45
153	16	Number of read counts by the periodontitis stages	46
154	17	Proportions of periodontitis DAT	47

155	18	Random forest classification metrics with the full microbiome compositions and ANCOM-selected DAT compositions	48
156			
157	19	Alpha-diversity indices account for evenness	49
158	20	Gradient Boosting classification metrics for periodontitis prediction	50
159	21	Gut microbiome compositions in genus level	72
160	22	Alpha-diversity indices in genus level	73
161	23	Alpha-diversity indices with recurrence in genus level	74
162	24	Alpha-diversity indices with OS in genus level	75
163	25	Beta-diversity indices in genus level	76
164	26	Beta-diversity indices with recurrence in genus level	77
165	27	Beta-diversity indices with recurrence in genus level	78
166	28	DAT with recurrence in species level	79
167	29	DAT with OS in species level	80
168	30	Random forest classification and regression	81

List of Tables

170	1	Confusion matrix	6
171	2	Standard clinical information of PTB study participants	13
172	3	Clinical characteristics of the study participants	32
173	4	Feature combinations and their evaluations	33
174	5	List of DAT among the periodontally healthy and periodontitis stages	34
175	6	Feature the importance of taxa in the classification of different periodontal statuses.	35
176	7	Beta-diversity pairwise comparisons on the periodontitis statuses	36
177	8	Clinical characteristics of CRC study participants	66
178	9	DAT list for CRC recurrence	67
179	10	DAT list for CRC OS	68
180	11	Random forest classification and their evaluations	70
181	12	Random forest regression and their evaluations	71

182

List of Abbreviations

183 **ACC** Accuracy

184 **ACE** Abundance-based coverage estimator

185 **ASV** Amplicon sequence variant

186 **AUC** Area-under-curve

187 **BA** Balanced accuracy

188 **BMI** Body mass index

189 **C-section** Cesarean section

190 **CAL** Clinical attachment level

191 **DAT** Differentially abundant taxa

192 **F1** F1 score

193 **Faith PD** Faith's phylogenetic diversity

194 **FC** Fold change

195 **FN** False negative

196 **FP** False positive

197 **FTB** Full-term birth

198 **GA** Gestational age

199 **MAE** Mean absolute error

200 **MSI** Microsatellite instability

201 **MSI-H** MSI-High

202 **MSI-L** MSI-Low

203 **MSS** Microsatellite stable

204 **MWU test** Mann-Whitney U-test

205 **OS** Overall survival

206 **PD** Probing depth

207 **PRE** Precision

- 208 **PROM** Prelabor rupture of membrane
- 209 **PTB** Preterm birth
- 210 **qPCR** quantitative-PCR
- 211 **RMSE** Root mean squared error
- 212 **ROC curve** Receiver-operating characteristics curve
- 213 **rRNA** Ribosomal RNA
- 214 **SD** Standard deviation
- 215 **SEN** Sensitivity
- 216 **SPE** Specificity
- 217 **t-SNE** t-distributed stochastic neighbor embedding
- 218 **TN** True negative
- 219 **TP** True positive

220 1 Introduction

221 The microbiome refers to the complex community of microorganisms, including bacteria, viruses, fungi,
222 and other microbes, that inhabit various environment within living organisms (Ursell, Metcalf, Parfrey,
223 & Knight, 2012; Gilbert et al., 2018). In humans, the microbiome plays a crucial role in maintaining
224 health (Lloyd-Price, Abu-Ali, & Huttenhower, 2016), influencing biological processes such as digestion
225 (Lim, Park, Tong, & Yu, 2020), immune response (Thaiss, Zmora, Levy, & Elinav, 2016; Kogut, Lee, &
226 Santin, 2020; C. H. Kim, 2018), and even mental health (Mayer, Tillisch, Gupta, et al., 2015; X. Zhu et
227 al., 2017; X. Chen, D'Souza, & Hong, 2013). These microbial communities are not static nor constant,
228 but rather dynamic ecosystem that interacts with their host and respond to environmental changes. Recent
229 studies have revealed that imbalances in the microbiome, known as dysbiosis, can contribute to a wide
230 range of diseases, including obesity (John & Mullin, 2016; Tilg, Kaser, et al., 2011; Castaner et al.,
231 2018), diabetes (Barlow, Yu, & Mathur, 2015; Hartstra, Bouter, Bäckhed, & Nieuwdorp, 2015; Sharma &
232 Tripathi, 2019), infections (Whiteside, Razvi, Dave, Reid, & Burton, 2015; Alverdy, Hyoju, Weigerinck,
233 & Gilbert, 2017), inflammatory conditions (Francescone, Hou, & Grivennikov, 2014; Peirce & Alviña,
234 2019; Honda & Littman, 2012), and cancers (Helmink, Khan, Hermann, Gopalakrishnan, & Wargo, 2019;
235 Cullin, Antunes, Straussman, Stein-Thoeringer, & Elinav, 2021; Sepich-Poore et al., 2021; Schwabe &
236 Jobin, 2013). Thus, understanding the composition of the human microbiomes is essential for developing
237 new therapeutic approaches that target these microbial populations to promote health and prevent diseases.

238 The microbiome participates a crucial role in overall health, influencing not only digestion and immune
239 function but also systemic and neurological processes through the brain-gut axis (Martin, Osadchiy,
240 Kalani, & Mayer, 2018; Aziz & Thompson, 1998; R. Li et al., 2024). The gut microbiota interact with
241 the host through metabolic byproducts, immune signaling, and the production of neurotransmitters, *e.g.*
242 serotonin and dopamine, which are essential for brain function and cognition. Disruptions in microbial
243 composition, known as dysbiosis, have been linked to various diseases, including inflammatory bowel
244 disease (Sultan et al., 2021; Baldelli, Scaldaferrri, Putignani, & Del Chierico, 2021), obesity (Kang et al.,
245 2022; Hamjane, Mechita, Nourouti, & Barakat, 2024; Pezzino et al., 2023), diabetes (Cai et al., 2024;
246 X. Li et al., 2021; Y. Li et al., 2023), and cardiovascular diseases (Manolis, Manolis, Melita, & Manolis,
247 2022; Tian et al., 2021). Furthermore, the brain-gut axis, a bidirectional communication system between
248 the gut microbiome composition and the central nervous system, has been implicated in mental disorders,
249 *e.g.* anxiety disorder, depressive disorder, and neurodegenerative diseases. Emerging evidence suggested
250 that alterations in the host microbiome can influence mood, cognitive function, and even behavior through
251 immune modulation, vagus nerve signaling, and microbial metabolites. These findings highlight the
252 microbiome as a critical factor in maintaining host health and suggest that targeted interventions, namely
253 probiotics, antibiotics, dietary modification, and microbiome-based therapies, may hold promise for
254 improving both physical and mental comfort. Hence, understanding the microbial effects could lead to
255 novel therapeutic strategies for a wide range of health conditions.

256 16S ribosomal RNA (rRNA) gene sequencing is one of the most extensively applied methods for
257 characterizing microbial communities by targeting the conserved 16S rRNA gene, which contains both

258 highly conserved and variable regions in bacteria (Tringe & Hugenholtz, 2008; Janda & Abbott, 2007).
259 The conserved regions enable universal primer binding, while the variable regions provide the specificity
260 needed to differentiate microbial taxa. Among these regions, the V3-V4 region is frequently selected for
261 sequencing due to its balance between phylogenetic resolution and sequencing efficiency (Johnson et al.,
262 2019; López-Aladid et al., 2023). Therefore, the V3-V4 region offers sufficient variability to classify a
263 wide range of bacteria taxa while maintaining compatibility with widely used sequencing platforms.

264 On the other hand, PathSeq is a computational pipeline designed for the identification and analysis
265 of microbial sequences within short-read human sequencing data, such as next-generation sequencing
266 (Kostic et al., 2011; Walker et al., 2018). PathSeq's scalable and effective processing of massive amounts
267 of sequencing data allows large-scale microbial profiling possible. PathSeq workflow consists of two
268 main phases: a subtractive phase and an analytic phase. The subtractive phase is removing human-derived
269 reads by aligning them to a human reference genome; and, the analytic phase is mapping remaining reads
270 to microbial reference databases, not only bacterial reference genome, but also archaeal, fungal, and viral
271 reference genomes. This approach allows for the comprehensive detection of microbiome compositions,
272 without a requirement for targeted amplification. PathSeq presents a more comprehensive and objective
273 evaluation of microbiome compositions than conventional microbiome profiling techniques including 16S
274 rRNA gene sequencing, capturing an assortment of microbial species beyond bacteria. Therefore, PathSeq
275 is an effective instrument for metagenomic research, infectious disease study, and microbiome analysis in
276 environmental and clinical contexts because of its capacity to operate with complex sequencing datasets
277 (Ojesina et al., 2013; Park et al., 2024; Tejeda et al., 2021).

278 The Anna Karenina principle, originally derived from literature of Leo Tolstoy, has been applied
279 to microbiome research to describe the manner that microbial communities in patients with diseases
280 tend to be more variable and unstable compared to those in healthy individuals (Ma, 2020; W. Li &
281 Yang, 2025). This Anna Karenina principle suggests that while healthy microbiomes exhibit relatively
282 stable and uniform compositions, while disease-associated microbiomes become highly dysregulated due
283 to various environmental, genetic, and pathological influences. Dysbiosis-driven mechanisms, such as
284 inflammation, genotoxic metabolic production, and immune modulation, can contribute pathogenesis
285 and progression of diseases, including periodontitis. In the context of cancer, this Anna Karenina
286 principle suggests that gut microbiome dysbiosis does not follow a single uniform pattern in patients
287 with CRC but rather presents as diverse and individualized disruption in microbial composition. This
288 instability may play a role in field cancerization, where microbial alteration extend beyond the tumor
289 site to adjacent normal-appearing tissues (Curtius, Wright, & Graham, 2018; Rubio, Lang-Schwarz,
290 & Vieth, 2022), potentially priming the tumor microenvironment for malignancy. Therefore, the high
291 inter-individual variability in microbiome alteration across these disease supports the Anna Karenina
292 principle, highlighting the complexity of dysbiosis-driven diseases and the necessity for personalized
293 microbiome-based diagnostic and interventions. Investigating the shared and disease-specific microbial
294 disruptions across these conditions may offer novel insights into microbiome-driven pathogenesis and
295 therapeutic strategies.

296 Diversity indices are essential techniques for evaluating the complexity and variety of microbial

communities, in ecological and microbiological research (Tucker et al., 2017; Hill, 1973). Alpha-diversity index attributes to the heterogeneity within a specific community, obtaining the number of different taxa and the distribution of taxa among the individuals, *i.e.*, richness and evenness. On the other hand, beta-diversity index measures the variations in microbiome compositions between the individuals, highlighting differences among the microbiome compositions of the study participants (B.-R. Kim et al., 2017). Altogether, by providing a thorough understanding of microbiome compositions, diversity indices, *e.g.* alpha-diversity and beta-diversity, allow us to investigate factors that affect community variability and structure.

Differentially abundant taxa (DAT) detection is a key analytical approach in microbiome study to identify microbial taxa that significantly differ in abundance between distinct study participant groups. This DAT detection method is particularly valuable for understanding how microbial communities vary across different conditions, such as disease states, environmental factors, and/or experimental treatments. Various statistical and computational techniques, *e.g.* LEfSe (Segata et al., 2011), DESeq2 (Love, Huber, & Anders, 2014), ANCOM (Lin & Peddada, 2020), and ANCOM-BC (Lin, Eggesbø, & Peddada, 2022; Lin & Peddada, 2024), are commonly used to assess differential abundance while accounting for compositional and sparsity-related challenges in microbiome composition data (Swift, Cresswell, Johnson, Stilianoudakis, & Wei, 2023; Cappellato, Baruzzo, & Di Camillo, 2022). Thus, identifying DAT can provide insights into microbial biomarkers associated with specific health conditions or disease statuses, enabling potential applications in diagnostics and therapeutics. However, due to the nature of microbiome composition data and the influence of sequencing depth, appropriate normalization and statistically adjustments are necessary to ensure reliable and stable detection of differentially abundant microbes (Xia, 2023; Pan, 2021). Integrating DAT detection analysis with functional profiling further enhances our understanding of the biological significance of microbial shifts or dysbiosis. As microbiome research advances, improving methodologies for DAT selection remains essential for uncovering meaningful microbial association and their potential roles in human diseases.

Classification is one of the supervised machine learning techniques used to categorized data into predefined classes based on features within the data (Kotsiantis, Zaharakis, & Pintelas, 2006; Sen, Hajra, & Ghosh, 2020). In other words, the method learns the relationship between input features and their corresponding output classes through the process of training a classification model using labeled data. Classification models are essential for advising choices in a wide range of applications, including medical diagnostics (Omondiagbe, Veeramani, & Sidhu, 2019). Thus, researchers could uncover sophisticated connections in input features and corresponding classes and produce reliable prediction by utilizing machine learning classification.

Random forest classification is one of the ensemble machine learning methods that constructs several decision trees during training and aggregates their results to provide classification predictions (Breiman, 2001; Geurts, Ernst, & Wehenkel, 2006). A portion of the features and classes—known as bootstrapping (Jiang & Simon, 2007; Champagne, McNairn, Daneshfar, & Shang, 2014; J.-H. Kim, 2009) and feature bagging (Bryll, Gutierrez-Osuna, & Quek, 2003; Alelyani, 2021; Yaman & Subasi, 2019)—are utilized to construct each tree in the forest. The majority vote from each tree determines the final classification,

336 which lowers the possibility of overfitting in comparison to a single decision tree. Furthermore, random
337 forest classifier offers several advantages, including its robustness to outliers and its ability to calculate
338 the feature importance.

339 Furthermore, k -fold cross-validation is a widely applied resampling technique that enhances the
340 reliability and robustness of machine learning models by iteratively evaluating their performance across
341 multiple data partitions (Wong & Yeh, 2019; Ghojogh & Crowley, 2019). Instead of relying on a single
342 train-test split, k -fold cross-validation divides the dataset into equally sized k folds, where the machine
343 learning model is trained on $k - 1$ folds and tested on the remaining fold in an iterative manner. This
344 process is repeated k times, with each fold serving as the test set once, and the final performance is
345 averaged across all iterations to provide a more generalizable estimate of model metrics. By reducing the
346 risk of overfitting and minimizing variance in performance evaluation, k -fold cross-validation ensures
347 that the machine learning model is not overly dependent on a specific train-test split. By applying k -fold
348 cross-validation, researchers can ensure that their machine learning models are both robust and reliable,
349 leading to more accurate and reproducible results (Fushiki, 2011).

350 Evaluating the performance of a machine learning classification model is essential to ensure its
351 reliability and effectiveness in real-world solutions and applications (Novaković, Veljović, Ilić, Papić, &
352 Tomović, 2017; Hossin & Sulaiman, 2015; Hand, 2012). A confusion matrix is a tabular representation of
353 predictions of classification, showing the counts of true positives (TP), true negatives (TN), false positives
354 (FP), and false negatives (FN) (Table 1). From this matrix, evaluations can be derived: accuracy (ACC;
355 Equation 1), balanced accuracy (BA; Equation 2), F1 score (F1; Equation 3), sensitivity (SEN; Equation
356 4), specificity (SPE; Equation 5), and precision (PRE; Equation 6). These metrics are in $[0, 1]$ range and
357 high metrics are good metrics. The confusion matrix also helps in identifying specific types of errors, such
358 as a tendency to produce false positive or false negatives, offering valuable insights for improving the
359 classification model. By combining the confusion matrix with other evaluation metrics, researchers can
360 comprehensively assess the classification metrics and refine it for real-world solutions and applications.

361 The receiver-operating characteristics (ROC) curve is a graphical representation used to evaluate
362 the performance of a classification model by plotting the sensitivity against (1-specificity) at multiple
363 threshold setting (Gonçalves, Subtil, Oliveira, & de Zea Bermudez, 2014; Obuchowski & Bullen, 2018;
364 Centor, 1991). The ROC curve illustrates the trade-off between detecting true positives while minimizing
365 false positives, suggesting determining the optimal decision threshold for classification. A key metric
366 derived from the ROC curve is the area-under-curve (AUC), which quantifies overall ability of the
367 classification model to discriminate between positive and negative predictions. An AUC value of 0.5
368 indicates a model performing no better than random chance, while value closer to 1.0 suggests high
369 predictive accuracy. Thus, by analyzing the AUC value of the ROC curve, researchers can compare
370 different models and select the better classification model that offers the best balance between sensitivity
371 and specificity for a given application.

372 Regression is a powerful predictive machine learning approach used to analyze complex relationships
373 between variables and make continuous value predictions (Maulud & Abdulazeez, 2020; Yildiz, Bilbao, &
374 Sproul, 2017). Beside classification, which assigns discrete labels, regression models estimate numerical

375 outcomes based on input features, making them particularly useful in biological research and clinical
376 applications for predicting disease risk, patient outcomes, and biomarker selection. By leveraging high-
377 throughput biological techniques and clinical information, regression model enables the discovery of
378 hidden patterns and the development of precision medicine strategies. As computational methods advance,
379 integrating regression models with metagenomic data can improve predictive accuracy and facilitate
380 data-driven therapeutic guide in healthcare.

381 Evaluating the performance of machine learning regression models requires assessing their prediction
382 errors using appropriate metrics. Mean absolute error (MAE; Equation 7) and root mean squared error
383 (RMSE; Equation 8) are commonly used measures for quantifying the accuracy of regression models. By
384 optimizing regression models based on MAE and RMSE, researchers can improve prediction accuracy
385 and enhance the reliability of machine learning regression models.

386 This dissertation present a comprehensive, multi-disease human microbiome analysis, bridging the
387 association between preterm birth (PTB) (Section 2), periodontitis (Section 3), and colorectal cancer
388 (CRC) (Section 4) through a unified metagenomic approach. While previous studies have examined the
389 role and characteristics of human microbiome in these diseases individually, this dissertation uniquely
390 integrates human microbiome-driven insights across these diseases to identify shared and disease-specific
391 microbial signatures. By applying high-throughput metagenomic sequencing, microbial diversity analysis,
392 and advanced bioinformatics techniques, this dissertation aims to uncover novel microbiome-based
393 biomarkers and mechanistic insights into how microbial communities influence these conditions. These
394 findings contribute to a broader understanding of microbiome-mediated disease interactions and pave the
395 way for personalized medicine strategies, including microbiome-targeted diagnostics and therapeutics.

Table 1: Confusion matrix

		Predicted	
		Positive	Negative
Actual	Positive	True positive (TP)	False negative (FN)
	Negative	False positive (FP)	True negative (TN)

396

$$ACC = \frac{TP + TN}{TP + FN + FP + TN} \quad (1)$$

397

$$BA = \frac{1}{2} \times \left(\frac{TP}{TP + FP} + \frac{TN}{TN + FN} \right) \quad (2)$$

398

$$F1 = \frac{2 \times TP}{2 \times TP + FP + FN} \quad (3)$$

399

$$SEN = \frac{TP}{TP + FP} \quad (4)$$

400

$$SPE = \frac{TN}{TN + FN} \quad (5)$$

401

$$PRE = \frac{TP}{TP + FP} \quad (6)$$

402

$$MAE = \sum_{i=1}^n |Prediction_i - Real_i| / n \quad (7)$$

$$RMSE = \sqrt{\sum_{i=1}^n (Prediction_i - Real_i)^2 / n} \quad (8)$$

403 **2 Predicting preterm birth using random forest classifier in salivary mi-**
404 **crobiome**

405 **This section includes the published contents:**

406 Hong, Y. M., **Lee, Jaewoong**, Cho, D. H., Jeon, J. H., Kang, J., Kim, M. G., ... & Kim, J. K. (2023).
407 Predicting preterm birth using machine learning techniques in oral microbiome. *Scientific Reports*, 13(1),
408 21105.

409 **2.1 Introduction**

410 Preterm birth (PTB), characterized by the delivery of neonates prior to 37 weeks of gestation, is one
411 of the major cause to neonatal mortality and morbidity (Blencowe et al., 2012). Multiple pregnancies
412 including twins, short cervical length, and infection on genitourinary tract are known risk factor for
413 PTB (Goldenberg, Culhane, Iams, & Romero, 2008). Nevertheless, the extent to which these aspects
414 affect birth outcomes is still up for debate. Henceforth, strategies to boost gestation and enhance delivery
415 outcomes can be more conveniently implemented when pregnant women at high risk of PTB are identified
416 early (Iams & Berghella, 2010).

417 Prediction models that can be utilized as a foundation for intervention methods still have an unac-
418 ceptable amount of classification evaluations, including accuracy, sensitivity, and specificity, despite a
419 great awareness of the risk factors that trigger PTB (Sotiriadis, Papatheodorou, Kavvadias, & Makrydi-
420 mas, 2010). Several attempts have been made to predict PTB through integrating data such as human
421 microbiome composition, inflammatory markers, and prior clinical data with predictive machine learn-
422 ing methods (Berghella, 2012). Because it is affordable and straightforward to use, fetal fibronectin is
423 commonly used in medical applications. However, with a sensitivity of only 56% that merely similar to
424 random prediction, it has a low classification evaluation (Honest et al., 2009). Due to the difficulty and
425 imprecision of the method in general, as well as the requirement for a qualified specialist cervical length
426 measuring is also restricted (Leitich & Kaider, 2003).

427 Preterm prelabor rupture of membranes (PROM) brought on by gestational inflammation and infection
428 contribute to about 70% of PTB cases (Romero, Dey, & Fisher, 2014). Nevertheless, as antibiotics and
429 anti-inflammatory therapeutic strategies were ineffective to decrease PTB occurrence rates, the pathology
430 of PTB has not been entirely elucidated by inflammatory and infectious pathways (Romero, Hassan, et al.,
431 2014). Recent researches on maternal microbiomes were beginning to examine unidentified connections
432 of PTB as a consequence of developmental processes in molecular biological technology (Fettweis et al.,
433 2019).

434 However, as anti-inflammatory and antibiotic therapies were insufficient to lower PTB occurrence
435 rates, infectious and inflammatory processes are insufficient to exhaustively clarify the pathogenesis and
436 pathophysiology of PTB. It has been hypothesized that the microbiota linked to PTB originate from either
437 a hematogenous pathway or the female genitourinary tract increasing through the vagina and/or cervix
438 (Han & Wang, 2013). Vaginal microbiome compositions have been found in women who eventually

439 acquire PTB, and recent studies have tried to predict PTB risk using cervico-vaginal fluid (Kindinger et
440 al., 2017). Even though previous investigation have confirmed the potential relationships between the
441 vaginal microbiome compositions and PTB, these studies are only able to clarify an upward trajectory.

442 Multiple unfavorable birth outcomes, including PROM and PTB, have been linked to periodontitis
443 as an independence risk factor, according to numerous epidemiological researches (Offenbacher et al.,
444 1996). It is expected that the oral microbiome will be able to explain additional hematogenous pathways
445 in light of these precedents; however, the oral microbiome composition of fetuses is limited understood.

446 Hence, in order to identify the salivary microbiome linked to PTB and to establish a machine learning
447 prediction model of PTB determined by oral microbiome compositions, this study examined the salivary
448 microbiome compositions of PTB study participants with a full-term birth (FTB) study participants.

449 **2.2 Materials and methods**

450 **2.2.1 Study design and study participants**

451 Between 2019 and 2021, singleton pregnant women who received treatment to Jeonbuk National University Hospital for childbirth were the participants of this study. This study was conducted according to the
452 Declaration of Helsinki (Goodyear, Krleza-Jeric, & Lemmens, 2007). The Institutional Review Board
453 authorized this study (IRB file No. 2019-01-024). Participants who were admitted for elective cesarean
454 sections (C-sections) or induction births, as well as those who had written informed consent obtained
455 with premature labor or PROM, were eligible.
456

457 **2.2.2 Clinical data collection and grouping**

458 Questionnaires and electronic medical records were implemented to gather information on both previous
459 and current pregnancy outcomes. The following clinical data were analyzed:

- 460 • maternal age at delivery
- 461 • diabetes mellitus
- 462 • hypertension
- 463 • overweight and obesity
- 464 • C-section
- 465 • history PROM or PTB
- 466 • gestational week on delivery
- 467 • birth weight
- 468 • sex

469 **2.2.3 Salivary microbiome sample collection**

470 Salivary microbiome samples were collected 24 hours before to delivery using mouthwash. The standard
471 methods of sterilizing were performed. Medical experts oversaw each stage of the sample collecting
472 procedure. Participants received instruction not to eat, drink, or brush their teeth for 30 minutes before
473 sampling salivary microbiome. Saliva samples were gathered by washing the mouth for 30 seconds with
474 12 mL of a mouthwash solution (E-zен Gargle, JN Pharm, Pyeongtaek, Gyeonggi, Korea). The samples
475 were tagged with the anonymous ID for each participant and kept in low temperature (4 °C) until they
476 underwent further processing. Genomic DNA was extracted using an ExgeneTM Clinic SV kit (GeneAll
477 Biotechnology, Seoul, Korea) following with the manufacturer instructions and store at -20 °C.

478 **2.2.4 16s rRNA gene sequencing**

479 Salivary microbiome samples were transported to the Department of Biomedical Engineering of the
480 Ulsan National Institute of Science and Technology . 16S rRNA sequencing was then carried out using a
481 commissioned Illumina MiSeq Reagent Kit v3 (Illumina, San Diego, CA, USA). Library methods were
482 utilized to amplify the V3-V4 areas. 300 base-pair paired-end reads were produced by sequencing the

483 pooled library using a v3 \times 600 cycle chemistry after the samples had been diluted to a final concentration
484 of 6 pM with a 20% PhiX control.

485 **2.2.5 Bioinformatics analysis**

486 The independent *t*-test was utilized to evaluate the differences of continuous values between from the
487 PTB participants than the FTB participants; χ^2 -square test was applied to decide statistical differences of
488 categorical values. Clinical measurement comparisons were conducted using SPSS (version 20.0) (Spss
489 et al., 2011). At $p < 0.05$, statistical significance was taken into consideration.

490 QIIME2 (version 2022.2) was implemented to import 16S rRNA gene sequences from salivary
491 microbiome samples of study participants for additional bioinformatics processing (Bolyen et al., 2019).
492 DADA2 was used to verify the qualities of raw sequences (Callahan et al., 2016). The remain sequences
493 were clustered into amplicon sequence variants (ASVs). Diversity indices, namely Faith PD for alpha
494 diversity index (Faith, 1992) and Hamming distance for beta diversity index (Hamming, 1950), were
495 calculated. MWU test (Mann & Whitney, 1947), and PERMANOVA multivariate test were evaluated for
496 measuring statistical significance (Anderson, 2014; Kelly et al., 2015).

497 Taxonomic assignment were implemented with HOMD (version 15.22) (T. Chen et al., 2010).
498 Afterward, DESeq2 was implemented to identify differentially abundant taxa (DAT) that could dis-
499 tinguish between salivary microbiome from PTB and FTB participants (Love et al., 2014). Taxa with
500 $|\log_2 \text{FoldChange}| > 1$ and $p < 0.05$ were considered as statistically significant.

501 The taxa for predicting PTB using salivary microbiome data were determined using a random forest
502 classifier (Breiman, 2001). Through stratified *k*-fold cross-validation (*k* = 5) that preserves the existence
503 rate of PTB and FTB participants, consistency and trustworthy classification were ensured (Wong & Yeh,
504 2019).

505 **2.2.6 Data and code availability**

506 All sequences from the 59 study participants have been published to the Sequence Read Archives
507 (project ID PRJNA985119): <https://dataview.ncbi.nlm.nih.gov/object/PRJNA985119>. Docker
508 image that employed throughout this study is available in the DockerHub: https://hub.docker.com/r/fumire/helixco_premature. Every code used in this study can be found on GitHub: https://github.com/CompbioLabUnist/Helixco_Premature.

511 **2.3 Results**

512 **2.3.1 Overview of clinical information**

513 In the beginning, 69 volunteer mothers were recruited for this study. However, due to insufficient clinical
514 information or twin pregnancies, 10 participants were excluded from the study participants. Demographic
515 and clinical information of the study participants are displayed in Table 2. Because PROM is one of the
516 leading factors of PTB, it was prevalent in the PTB group than the FTB group. Other maternal clinical
517 factors did not significantly differ between the FTB and PTB groups. There were no cases in both groups
518 that had a history of simultaneous periodontal disease or cigarette smoking.

519 **2.3.2 Comparison of salivary microbiomes composition**

520 The salivary microbiome composition was composed of 13953804 sequences from 59 study participants,
521 with 102305.95 ± 19095.60 and 64823.41 ± 15841.65 (mean \pm SD) reads/sample before and following
522 the quality-check stage, accordingly. There was not a significant distinction between the PTB and FTB
523 groups with regard to on alpha diversity nor beta diversity metrics (Figure 4).

524 DESeq2 was used to select 32 DAT that distinguish between the PTB and FTB groups out of the 465
525 species that were examined (Love et al., 2014): 26 FTB-enriched DAT and six PTB-enriched DAT. Seven
526 PROM-related DAT were removed from these 32 PTB-related DAT to lessen the confounding effect of
527 PROM (Figure 5). Therefore, there were a total of 25 PTB-related DAT: 22 FTB-enriched DAT and three
528 PTB-enriched DAT (Figure 1).

529 A significant negative correlation was found using Pearson correlation analysis between GW and
530 differences between PTB-enriched DAT and FTB-enriched DAT (Pearson correlation $r = -0.542$ and
531 $p = 7.8e-6$; Figure 5).

532 **2.3.3 Random forest classification to predict PTB risk**

533 To classify PTB according to DAT, random forest classifiers were constructed. The nine most significant
534 DAT were used to obtain the best BA (0.765 ± 0.071 ; Figure 3a). Moreover, random forest classification
535 model determined each DAT's importance (Figure 3b). We conducted a validation procedure on nine
536 twin pregnancies that were excluded in the initial study design in order to confirm the reliability and
537 dependability of our random forest-based PTB prediction model (Figure 6). Comparable to the PTB
538 prediction model on the 59 initial singleton study participants, the validation classification on PTB risk of
539 these twin participants have an accuracy of 87.5%.

Table 2: Standard clinical information of PTB study participants.

Continuous variable for independent *t*-test. Categorical variable for Pearson's χ^2 -square test. Continuous variable: mean \pm SD. Categorical variable: count (proportion)

	PTB (n=30)	FTB (n=29)	p-value
Maternal age (years)	31.8 \pm 5.2	33.7 \pm 4.5	0.687
C-section	20 (66.7%)	24 (82.7%)	0.233
Previous PTB history	4 (13.3%)	1 (3.4%)	0.353
PROM	12 (40.0%)	1 (3.4%)	0.001
Pre-pregnant overweight	8 (26.7%)	7 (24.1%)	1.000
Gestational weight gain (kg)	9.0 \pm 5.9	11.5 \pm 4.6	0.262
Diabetes	2 (6.7%)	2 (6.9%)	1.000
Hypertension	11 (36.7%)	4 (13.8%)	0.072
Gestational age (weeks)	32.5 \pm 3.4	38.3 \pm 1.1	\leq 0.001
Birth weight (g)	1973.4 \pm 686.6	3283.4 \pm 402.7	\leq 0.001
Male	14 (46.7%)	13 (44.8%)	1.000

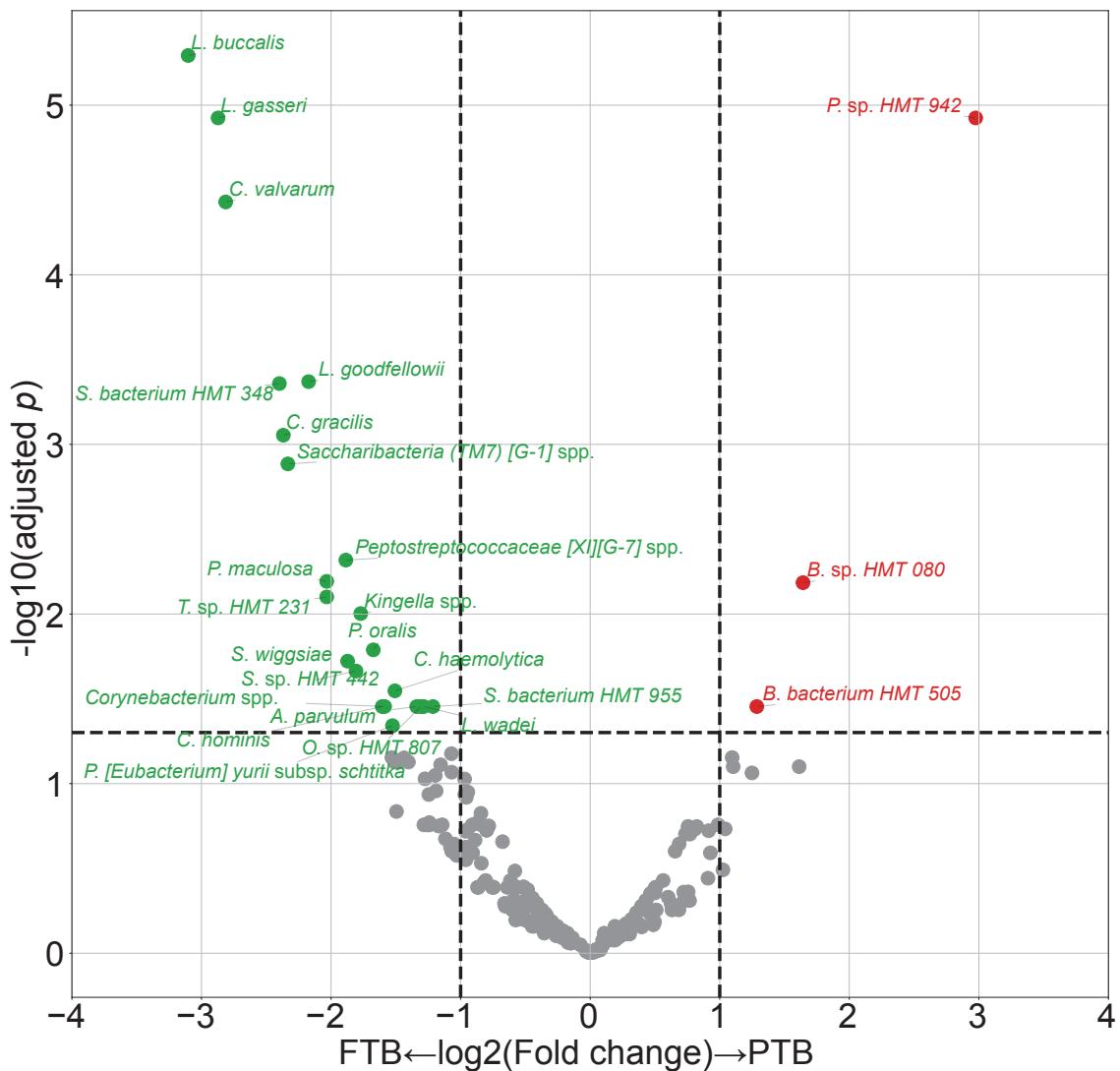


Figure 1: DAT volcano plot for PTB prediction.

Statistical threshold is: adjusted p -value < 0.05 and $|\log_2 \text{Fold Change}| > 1.0$. Red dots represent PTB-enriched DAT, while green dots represent FTB-enriched DAT.

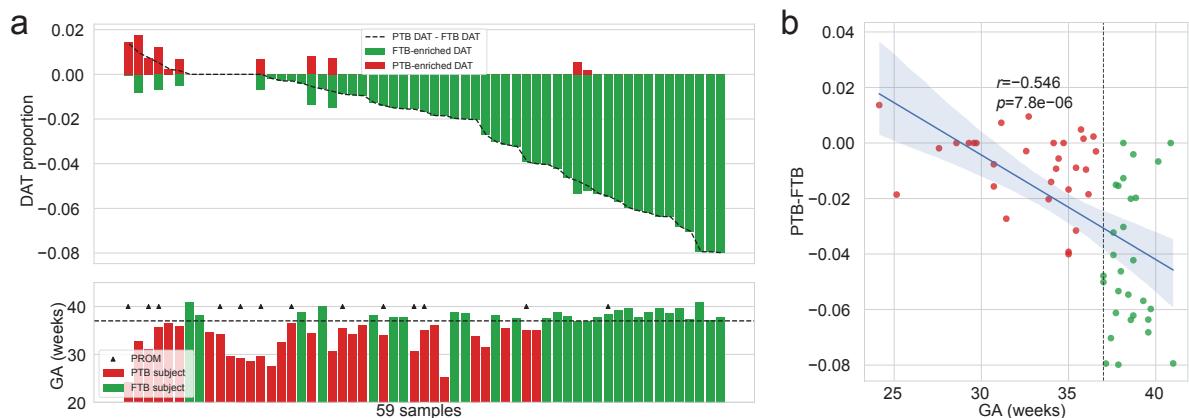


Figure 2: **Salivary microbiome compositions over DAT for PTB prediction.**

(a) Frequencies of DAT of PTB study subjects. The study participants are arranged in respect of (PTB-enriched DAT – FTB-enriched DAT). The study participants' GA is displayed in accordance with the upper panel's order (PTB: red bar, FTB: green bar. PROM: arrow head.) **(b)** Correlation plot with GA and (PTB-enriched DAT – FTB-enriched DAT). Strong negative correlation is found with Pearson correlation ($p = 7.8e - 6$).

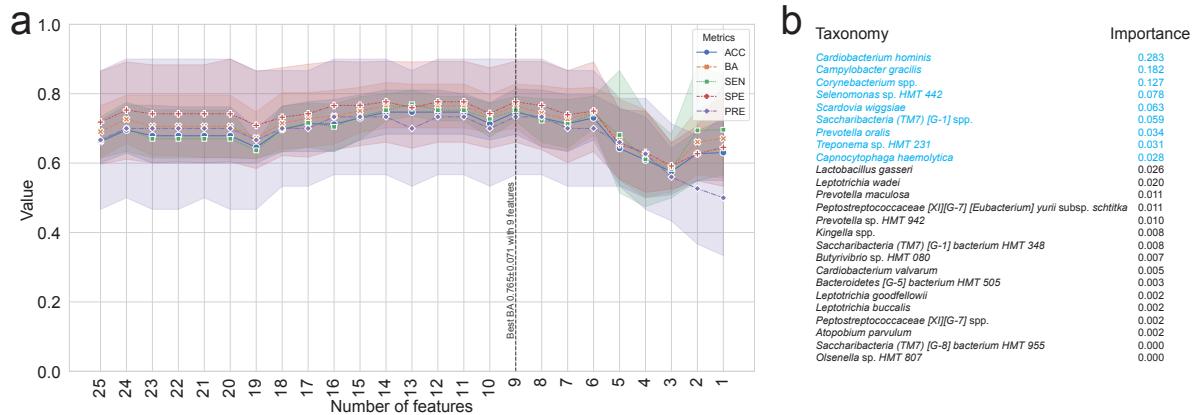


Figure 3: **Random forest-based PTB prediction model.**

(a) Machine learning evaluations upon number of features (DAT). Random Forest classifier has the best BA (0.765 ± 0.071 ; Mean \pm SD) with the nine most important DAT. **(b)** Importance of DAT.

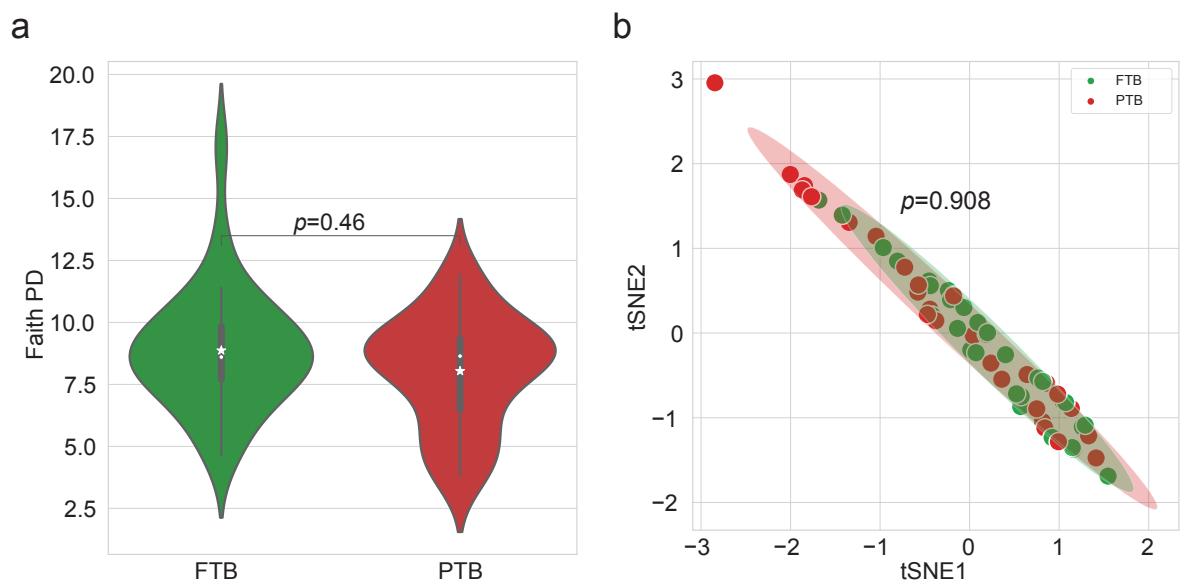


Figure 4: **Diversity indices about PTB study participants.**

(a) Alpha diversity index (Faith PD). There is no statistically significant difference between the PTB and FTB group (MWU test $p = 0.46$). **(b)** t-SNE plot with beta diversity index (Hamming distance). There is no statistically significant difference between the PTB and FTB group (PERMANOVA test $p = 0.908$)

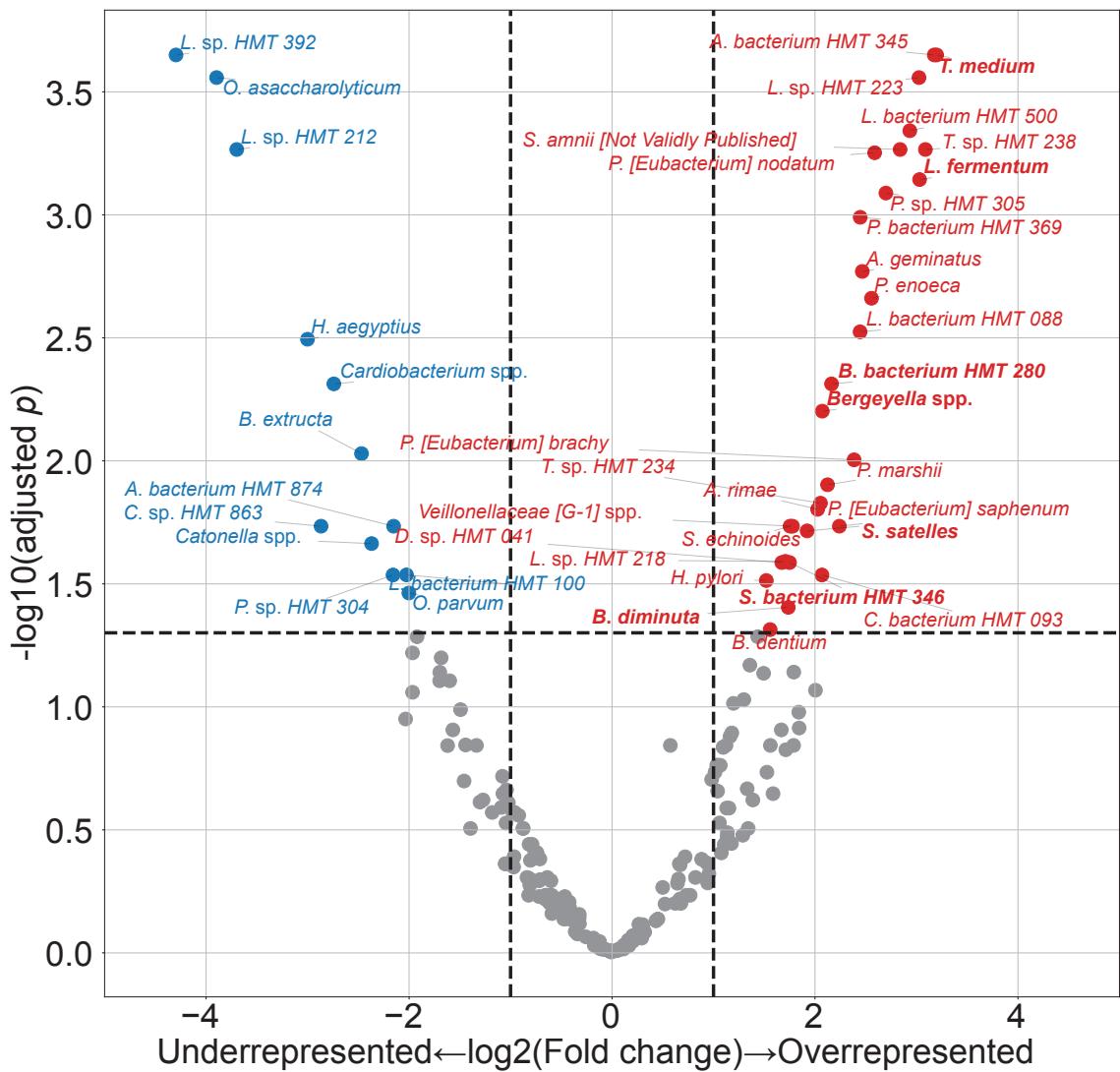


Figure 5: **PROM-related DAT between FTB and PTB.**

Statistical threshold is: adjusted p -value < 0.05 and $|\log_2(\text{Fold Change})| > 1.0$. Only seven of these 42 PROM-related DAT overlapped with PTB-related DAT (bold text). Blue dots represented PROM-underrepresented DAT, while red dots represented PROM-overrepresented DAT.

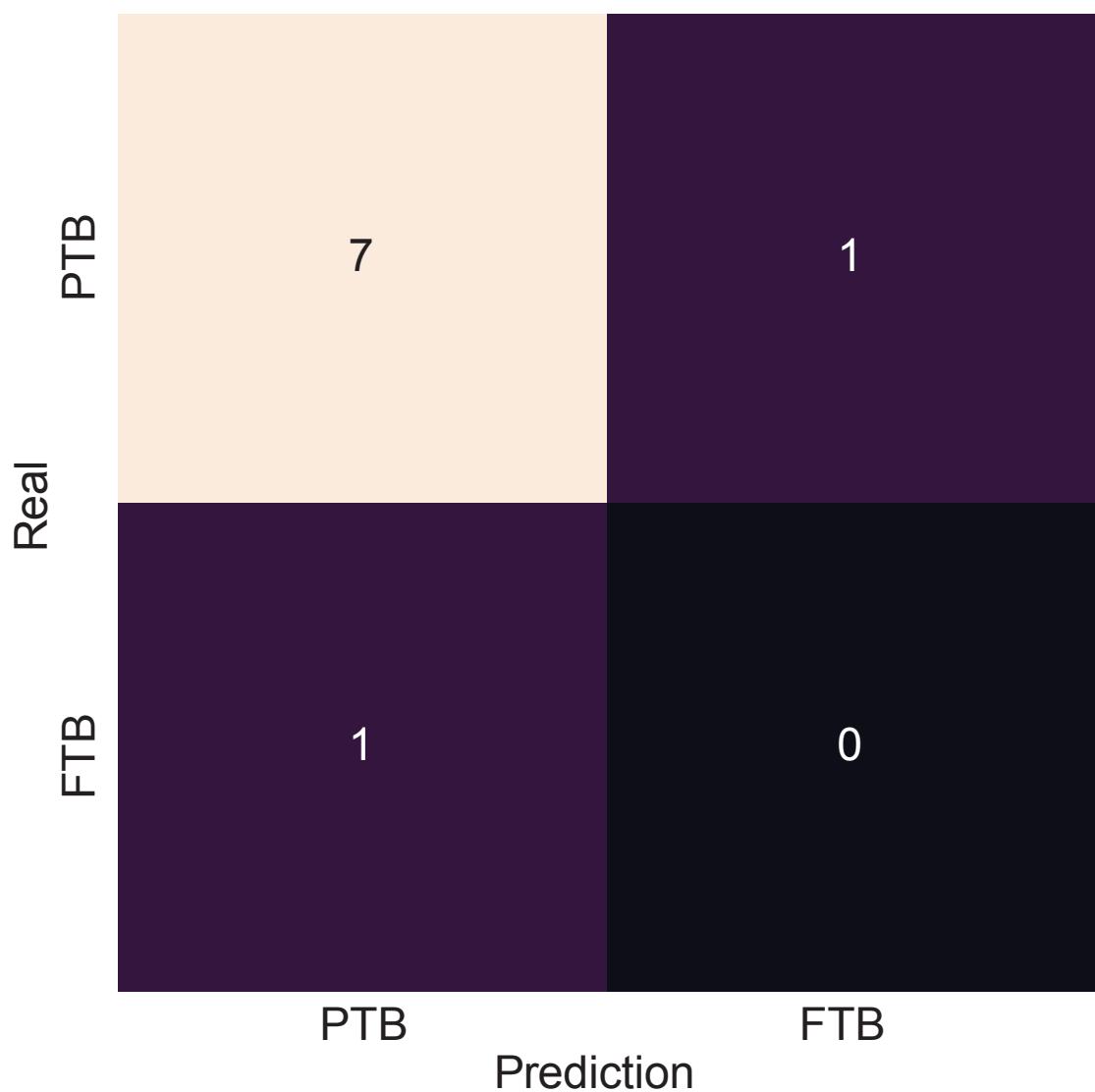


Figure 6: Validation of random forest-based PTB prediction model.

Nine twin pregnancies (eight PTB subjects and a FTB subject) that were excluded in the initial study subjects were subjected to a validation procedure. The random forest-based PTB prediction model shows 87.5% accuracy, comparable to the PTB classification evaluations on the singleton study subjects (0.714 ± 0.061 . Mean \pm SD)

540 **2.4 Discussion**

541 In this study, we employed salivary microbiome compositions to develop the random forest-based PTB
542 prediction models to estimate PTB risks. Previous reports have indicated bidirectional associations
543 between pregnancy outcomes and salivary microbiome compositions (Han & Wang, 2013). Nevertheless,
544 the salivary microbiome composition is not yet elucidated. Salivary microbial dysbiosis, including gingival
545 inflammation and periodontitis, have been connected to unfavorable pregnancy outcomes, such as PTB
546 (Ide & Papapanou, 2013). However, the techniques utilized in recent research that primarily focus on
547 recognized infections have led to inconsistent outcomes.

548 One of the most common salivary taxa that has been examined is *Fusobacterium nucleatum*, that is a
549 Gram-negative, anaerobic, and filamentous bacteria (Han, 2015; Brennan & Garrett, 2019; Bolstad, Jensen,
550 & Bakken, 1996). *Fusobacterium nucleatum* can be separated from not only the salivary microbiome
551 but also the vaginal microbiome (Vander Haar, So, Gyamfi-Bannerman, & Han, 2018; Witkin, 2019). In
552 both animal and human investigation, *Fusobacterium nucleatum* infection has been linked to risk of PTB
553 (Doyle et al., 2014). According to recent researches, the placenta women who give birth prematurely may
554 include additional salivary microbiome dysbiosis, such as *Bergeyella* spp. and *Porphyromonas gingivalis*
555 (León et al., 2007; Katz, Chegini, Shiverick, & Lamont, 2009). Although *Bergeyella* spp. were one of the
556 PROM-overrepresented DAT (Figure 5), it was excluded in the final 25 PTB-related DAT. Furthermore,
557 *Porphyromonas gingivalis* and *Campylobacter gracilis* were pathogens of periodontitis in sub-gingival
558 microbiome (Yang et al., 2022). *Lactobacillus gasseri* was also one of the FTB-enriched DAT (Figure
559 1), and it is well established that early PTB risk can be reduced by *Lactobacillus gasseri* in the vaginal
560 microbiome (Basavaprabhu, Sonu, & Prabha, 2020; Payne et al., 2021).

561 With DAT comprising 22 FTB-enriched DAT and three PTB-enriched DAT (Figure 1), we discovered
562 that the FTB study participants had the majority of the essential DAT that distinguished between the PTB
563 and FTB groups. Thus, we hypothesize that the pathogenesis and pathophysiology of PTB may have been
564 triggered by an absence of species with protective characteristics. The association between unfavorable
565 pregnancy outcomes and a dysfunctional microbiome has been explained through two distinct processes.
566 According to the first hypothesis, periodontal pathogens originating in the gingival biofilm might spread
567 from the infected salivary microbiome over the placenta microbiome, invade the intra-amniotic fluid
568 and fetal circulation, and then have a direct impact on the fetoplacental unit, leading to bacteremia
569 (Hajishengallis, 2015). Based on the second hypothesis, inflammatory mediators and endotoxins that
570 generated by the sub-gingival inflammation and derived from dental plaque of periodontitis may spread
571 throughout the body and reach the fetoplacental unit (Stout et al., 2013; Aagaard et al., 2014). Despite
572 belonging to the same species, some subgroups of the salivary microbiome may influence pregnancy
573 outcomes in both favorable and adverse manners. Following this line of argumentation, the salivary
574 microbiome composition or their dysbiosis are more significant than the existence of particular bacteria.

575 Notably, microbial alteration that take place throughout pregnancy may be expected results of a healthy
576 pregnancy. Those pregnancy-related vulnerabilities to dental problem like periodontitis can be explained
577 by three factors. Because of hormone-driven gingival hyper-reactivity to the salivary microbiome in the

578 oral biofilm including sub-gingival biofilm, these conditions are prevalent in pregnant women. For insight
579 at the relationship between the salivary microbiome compositions and PTB, further studies with pathway
580 analysis are warranted.

581 Our study confirmed that salivary microbiome composition could provide potential biomarkers for
582 predicting pregnancy complications including PTB risks using random forest-based classification models,
583 despite a limited number of study participants and a tiny validation sample size. Another limitation of our
584 study was 16S rRNA gene sequencing. In other words, unlike the shotgun sequencing, 16S rRNA gene
585 sequencing only focused on bacteria, not viruses nor fungi. We did not delve into other variables like
586 nutrition status and socioeconomic statuses of study participants that might affect the salivary microbiome
587 composition.

588 Notwithstanding these limitations, this prospective examination showed the promise of the random
589 forest-based PTB prediction models based on mouthwash-derived salivary microbiome composition.
590 Before applying the methods developed in this study in a clinical context, more multi-center and extensive
591 research is warranted to validate our findings.

592 **3 Random forest prediction model for periodontitis stages based on the**
593 **salivary microbiomes**

594 **3.1 Introduction**

595 Saliva microbial dysbiosis brought on by the accumulation of plaque results in periodontitis, a chronic
596 inflammatory disease of the tissue that surrounds the tooth (Kinane, Stathopoulou, & Papapanou, 2017).
597 Loss of periodontal attachment is a consequence of periodontitis, which may lead to irreversible bone loss
598 and, eventually, permanent tooth loss if left untreated. A new classification criterion of periodontal diseases
599 was created in 2018, about 20 years after the 1999 statements of the previous one (Papapanou et al.,
600 2018). Even with this evolution, radiographic and clinical markers of periodontitis progression remain the
601 primary methods for diagnosing periodontitis (Papapanou et al., 2018). Such tools, nevertheless, frequently
602 demonstrate the prior damage from periodontitis rather than its present condition. Certain individuals have
603 a higher risk of periodontitis, a higher chance of developing severe generalized periodontitis, and a worse
604 response to common salivary bacteria control techniques utilized to prevent and treat periodontitis. As a
605 result, the 2017 framework for diagnosing periodontitis additionally allows for the potential development
606 of biomarkers to enhance diagnosis and treatment of periodontitis (Tonetti, Greenwell, & Kornman, 2018).
607 Instead of only depending on the progression of periodontitis, a new etiological indication based on the
608 current state must be introduced in order to enable appropriate intervention through early detection of
609 periodontitis. Thus, the current clinical diagnostic techniques that rely on periodontal probing can be
610 uncomfortable for patients with periodontitis (Canakci & Canakci, 2007).

611 Due to the development of salivaomics, in this manner, the examination of saliva has emerged as
612 a significant alternative to the conventional ways of identifying periodontitis (Altingöz et al., 2021;
613 Melguizo-Rodríguez, Costela-Ruiz, Manzano-Moreno, Ruiz, & Illescas-Montes, 2020). Given that saliva
614 sampling is non-invasive, painless, and accessible to non-specialists, it may be a valuable instrument
615 for diagnosing periodontitis (C.-Z. Zhang et al., 2016). Furthermore, much research has suggested that
616 periodontitis could be a trigger in the development and exacerbation of metabolic syndrome (Morita et
617 al., 2010; Nesbitt et al., 2010). Consequently, alteration in these levels of salivary microbiome markers
618 may serve as high effective diagnostic, prognostic, and therapeutic indicators for periodontitis and
619 other systemic diseases (Miller, Ding, Dawson III, & Ebersole, 2021; Čižmárová et al., 2022). The
620 pathogenesis of periodontitis typically comprises qualitative as well as quantitative alterations in the
621 salivary microbial community, despite that it is a complex disease impacted by a number of contributing
622 factors including age, smoking status, stress, and nourishment (Abusleme, Hoare, Hong, & Diaz, 2021;
623 Lafaurie et al., 2022). Depending on the severity of periodontitis, the salivary microbial community's
624 diversity and characteristics vary (Abusleme et al., 2021), indicating that a new etiological diagnostic
625 standards might be microbial community profiling based on clinical diagnostic criteria. As a consequence,
626 salivary microbiome compositions have been characterized in numerous research in connection with
627 periodontitis. High-throughput sequencing, including 16S rRNA gene sequencing, has recently used in
628 multiple studies to identify variations in the bacterial composition of sub-gingival plaque collections

629 from periodontal healthy individuals and patients with periodontitis (Altabtbaei et al., 2021; Iniesta
630 et al., 2023; Nemoto et al., 2021). This realization has rendered clear that alterations in the salivary
631 microbial community—especially, shifts to dysbiosis—are significant contributors to the pathogenesis and
632 development of periodontitis (Lamont, Koo, & Hajishengallis, 2018). Yet most of these research either
633 focused only on the microbiome alterations in sub-gingival plaque collection, comprised a limited number
634 of periodontitis study participants, or did not account for the impact of multiple severities of periodontitis.

635 For the objective of diagnosing periodontitis, previous research has developed machine learning-based
636 prediction models based on oral microbiome compositions, such as the sub-gingival microbial dysbiosis
637 index (T. Chen, Marsh, & Al-Hebshi, 2022; Chew, Tan, Chen, Al-Hebshi, & Goh, 2024), which have
638 demonstrated good diagnostic evaluation and could be applied to individual saliva collection. Despite
639 offering valuable details, these indicators are frequently restricted by their limited emphasis on classifying
640 the multiple stages of periodontitis. Furthermore, many of these machine learning models currently in
641 practice are trained solely upon the existence of periodontitis rather than on the multiple severities of
642 periodontitis.

643 Recently, we employed multiplex quantitative-PCR (qPCR) and machine learning-based classification
644 model to predict the stage of periodontitis based on the amount of nine pathogens of periodontitis from
645 saliva collections (E.-H. Kim et al., 2020). On the other hand, the fact that we focused merely at nine
646 pathogens for periodontitis and neglected the variety bacterial species associated to the various severities
647 of periodontitis constrained the breadth of our investigation. By developing a machine learning model
648 that could classify multiple severities of periodontitis based on the salivary microbiome composition,
649 this study aims to fill these knowledge gaps and produce more accurate and therapeutically useful
650 guidance to evaluate progression of periodontitis. Hence, in order to examine the salivary microbiome
651 composition of both healthy controls and patients with periodontitis in multiple stages, we applied
652 16S rRNA gene sequencing. Furthermore, employing the 2018 classification criteria, we sought to find
653 biomarkers (bacterial species) for the precise prediction of periodontitis severities (Papapanou et al.,
654 2018; Chapple et al., 2018).

655 **3.2 Materials and methods**

656 **3.2.1 Study participants enrollment**

657 Between 2018-08 and 2019-03, 250 study participants—100 healthy controls, 50 patients with stage I
658 periodontitis, 50 patients with stage II periodontitis, and 50 patients with stage III periodontitis—visited
659 visited the Department of Periodontics at Pusan National University Dental Hospital. The Institutional
660 Review Board of the Pusan National University Dental Hospital accepted this study protocol and design
661 (IRB No. PNUDH-2016-019). Every study participants provided their written informed authorization after
662 being fully informed about this study's objectives and methodologies. Exclusion criteria for the study
663 participants are followings:

- 664 1. People who, throughout the previous six months, underwent periodontal therapy, including root
665 planing and scaling.
- 666 2. People who struggle with systemic conditions that may affect periodontitis developments, such as
667 diabetes.
- 668 3. People who, throughout the previous three months, were prescribed anti-inflammatory medications
669 or antibiotics.
- 670 4. Women who were pregnant or breastfeeding.
- 671 5. People who have persistent mucosal lesions, *e.g.* pemphigus or pemphigoid, or acute infection, *e.g.*
672 herpetic gingivostomatitis.
- 673 6. Patient with grade C periodontitis or localized periodontitis (< 30% of teeth involved).

674 **3.2.2 Periodontal clinical parameter diagnosis**

675 A skilled periodontist conducted each clinical procedure. Six sites per tooth were used to quantify
676 gingival recession and probing depth: mesiobuccal, midbuccal, distobuccal, mesiolingual, midlingual,
677 and distolingual (Huang et al., 2007). A periodontal probe (Hu-Friedy, IL, USA) was placed parallel to
678 the major axis of the tooth at each tooth location in order to gather measurements. The cementoenamel
679 junction of the tooth was analyzed to determine the clinical attachment level, and the deepest point of
680 probing was taken to determine the periodontal pocket depth from the marginal gingival level of the
681 tooth. Plaque index was measured by probing four surfaces per tooth: mesial, distal, buccal, and palatal
682 or lingual. Plaque index was scored by the following criteria:

- 683 0. No plaque present.
- 684 1. A thin layer of plaque that adheres to the surrounding tissue of the tooth and free gingival margin.
685 Only through the use of a periodontal probe on the tooth surface can the plaque be existed.
- 686 2. Significant development of soft deposits that are visible within the gingival pocket, which is a
687 region between the tooth and gingival margin.

688 3. Considerable amount of soft matter on the tooth, the gingival margin, and the gingival pocket.

689 The arithmetic average of the plaque indices collected from every tooth was determined to calculate
690 plaque index of each study participant. By probing four surfaces per tooth, mesial, distal, buccal, and
691 palatal or lingual, to assess gingival bleeding, the gingival index was scored by the following criteria:

692 0. Normal gingiva: without inflammation nor discoloration.

693 1. Mild inflammation: minimal edema and slight color changes, but no bleeding on probing.

694 2. Moderate inflammation: edema, glazing, redness, and bleeding on probing.

695 3. Severe inflammation: significant edema, ulceration, redness, and spontaneous bleeding.

696 The arithmetic average of the gingival indices collected from every tooth was determined to calculate
697 gingival index of each study participant. The relevant data was not displayed, despite that furcation
698 involvement and bleeding on probing were thoroughly utilized into account during the diagnosis process.

699 Periodontitis was diagnosed in respect to the 2018 classification criteria for periodontitis (Papapanou
700 et al., 2018; Chapple et al., 2018). An experienced periodontist diagnosed the periodontitis stage by con-
701 sidering complexity, depending on clinical examinations including radiographic images and periodontal
702 probing. Periodontitis is categorized into healthy, stage I, stage II, and stage III with the following criteria:

703 • Healthy:

704 1. Bleeding sites < 10%

705 2. Probing depth: \leq 3 mm

706 • Stage I:

707 1. No tooth loss because of periodontitis.

708 2. Inter-dental clinical attachment level at the site of the greatest loss: 1-2 mm

709 3. Radiographic bone loss: < 15%

710 • Stage II:

711 1. No tooth loss because of periodontitis.

712 2. Inter-dental clinical attachment level at the site of the greatest loss: 3-4 mm

713 3. Radiographic bone loss: 15-33%

714 • Stage III:

715 1. Teeth loss because of periodontitis: \leq 3 teeth

716 2. Inter-dental clinical attachment level at the site of the greatest loss: \geq 5 mm

717 3. Radiographic bone loss: > 33%

718 **3.2.3 Saliva sampling and DNA extraction procedure**

719 All study participants received instructions to avoid eating, drinking, brushing, and using mouthwash for
720 at least an hour prior to the saliva sample collection process. These collections were conducted between
721 09:00 and 11:00. Mouth rinse was collected by rinsing the mouth for 30 seconds with 12 mL of a solution
722 (E-zen Gargle, JN Pharm, Korea). All saliva samples were tagged with anonymous ID and stored at -4 °C.

723 Bacteria DNA was extracted from saliva samples using an Exgene™Clinic SV DNA extraction kit
724 (GeneAll, Seoul, Korea), and quality and quantity of bacterial DNA was measured using a NanoDrop
725 spectrophotometer (Thermo Fisher Scientific, Wilmington, DE, USA). Hyper-variable regions (V3-V4)
726 of the 16S rRNA gene were amplified using the following primer:

- 727 • Forward: 5' -TCGTCGGCAGCGTCAGATGTGTATAAGAGACAGCCTACGGGNNGCWGCAG-3'
728 • Reverse: 5' -GTCTCGTGGGCTCGGAGATGTGTATAAGAGACAGGACTACHVGGGTATCTAATCC-3'

729 The standard protocols of the Illumina 16S Metagenomic Sequencing Library Preparation were
730 followed in the preparation of the libraries. The PCR conditions were as follows:

- 731 1. Heat activation for 30 seconds at 95 °C.
732 2. 25 cycles for 30 seconds at 95 °C.
733 3. 30 seconds at 55 °C.
734 4. 30 seconds at 72 °C.

735 NexteraXT Indexed Primer was applied to amplification 10 µL of the purified initial PCR products for
736 the final library creation. The second PCR used the same conditions as the first PCR conditions but with
737 10 cycles. 16S rRNA gene sequencing was performed via 2×300 bp paired-end sequencing at Macrogen
738 Inc. (Macrogen, Seoul, Korea) using Illumina MiSeq platform (Illumina, San Diego, CA, USA).

739 **3.2.4 Bioinformatics analysis**

740 We computed alpha-diversity and beta-diversity indices to quantify the divergence of phylogenetic
741 information. Following alpha-diversity indices were calculated using the scikit-bio Python package
742 (version 0.5.5) (Rideout et al., 2018), and these alpha-diversity indices were compared using the MWU
743 test:

- 744 • Abundance-based Coverage Estimator (ACE) (Chao & Lee, 1992)
745 • Chao1 (Chao, 1984)
746 • Fisher (Fisher, Corbet, & Williams, 1943)
747 • Margalef (Magurran, 2021)
748 • Observed ASVs (DeSantis et al., 2006)
749 • Berger-Parker *d* (Berger & Parker, 1970)
750 • Gini (Gini, 1912)

- 751 • Shannon (Weaver, 1963)
752 • Simpson (Simpson, 1949)

753 Aitchison index for a beta-diversity index was calculated using QIIME2 (version 2020.8) (Aitchison,
754 Barceló-Vidal, Martín-Fernández, & Pawlowsky-Glahn, 2000; Bolyen et al., 2019). We employed the
755 t-SNE algorithm to illustrate multi-dimensional data from the beta-diversity index computation (Van der
756 Maaten & Hinton, 2008). The beta-diversity index was compared using the PERMANOVA test (Anderson,
757 2014; Kelly et al., 2015) and MWU test.

758 DAT between multiple periodontitis stages were identified by ANCOM (Lin & Peddada, 2020).
759 The log-transformed absolute abundances of DAT were analyzed by hierarchical clustering in order to
760 identify sub-groups with similar abundance patterns on periodontitis stages. Additionally, we examined
761 the relative proportions among the 20 DAT in order to reduce the effect of salivary bacteria that differ
762 insignificantly across the multiple severities of periodontitis.

763 Differentially abundant taxa (DAT) among multiple periodontitis severities were selected from the
764 salivary microbiome compositions by ANCOM (Lin & Peddada, 2020). In contrast to conventional
765 techniques that examine raw abundance counts, ANCOM applies log-ratio between taxa to account for
766 the salivary microbiome composition data. The log-transformed abundances of DAT were subjected
767 to hierarchical clustering to discover subgroups of DAT with similar patterns on periodontitis stages.
768 Furthermore, we examined the relative proportion among the DAT in order to reduce the effects of other
769 salivary bacteria that differ non-significantly across the multiple periodontitis severities.

770 As previously stated (E.-H. Kim et al., 2020), we used stratified k -fold cross-validation ($k = 10$)
771 by severity of periodontitis to achieve consistent and trustworthy classification results (Wong & Yeh,
772 2019). Additionally, we utilized various features with confusion matrices and their derivations to evaluate
773 the classification outcomes in order to identify which features optimize classification evaluations and
774 decrease sequencing efforts. Using the DAT discovered by ANCOM, we iteratively removed the least
775 significant taxa from the input features (taxa) of the random forest (Breiman, 2001) and gradient boosting
776 (Friedman, 2002) classification models using the backward elimination method. Random forest classifier
777 builds multiple decision trees independently using bootstrapped samples and aggregates their predictions,
778 enhancing stability and reducing overfitting problems. In contrast, Gradient boosting constructs trees
779 sequentially, where each new tree improves the errors of the previous ones using gradient descent, leading
780 to higher classification evaluations.

781 We investigated external datasets from Spanish individuals (Iniesta et al., 2023) and Portuguese
782 individuals (Relvas et al., 2021) to confirm that our random forest classification was consistent. To
783 ascertain repeatability and dependability, the external datasets were processed using the same pipeline
784 and parameters as those used for our study participants.

785 **3.2.5 Data and code availability**

786 All sequences from the 250 study participants have been published to the Sequence Read Archives (project
787 ID PRJNA976179): <https://www.ncbi.nlm.nih.gov/Traces/study/?acc=PRJNA976179>. Docker

788 image that employed throughout this study is available in the DockerHub: <https://hub.docker.com/>
789 repository/docker/fumire/periodontitis_16s. Every code used in this study can be found on
790 GitHub: https://github.com/CompbioLabUnist/Periodontitis_16S.

791 **3.3 Results**

792 **3.3.1 Summary of clinical information and sequencing data**

793 Among clinical information of the study participants, clinical attachment level, probing depth, plaque
794 index, and gingival index, were significantly increased with periodontitis severity (Kruskal-Wallis test
795 $p < 0.001$), while sex were observed no significant difference (Table 2). Notably, clinical attachment level
796 and probing depth have significant differences among the periodontitis severities (MWU test $p < 0.01$;
797 Figure 15). Additionally, 71461.00 ± 11792.30 and 45909.78 ± 11404.65 (mean \pm SD) reads per sample
798 were obtained before and after filtering low-quality reads and trimming extra-long tails, respectively
799 (Figure 16). In 250 study subjects, we have found a total of 425 bacterial taxa (Figure 13).

800 **3.3.2 Diversity indices reveal differences among the periodontitis severities**

801 Rarefaction curves showed that the sequencing depth was sufficient (Figure 12). Alpha-diversity indices
802 indicated significant differences between the healthy and the periodontitis stages (MWU test $p < 0.01$;
803 Figure 7a-e); however, there were no significant differences between the periodontitis stages. This
804 emphasizes how essential it is to classify the salivary microbiome compositions and distinguish between
805 the stages of periodontitis using machine learning approaches.

806 The confidence ellipses of the tSNE-transformed beta-diversity index (Aitchison index) indicated
807 distinct distributions among the periodontitis severities (PERMANOVA $p \leq 0.001$; Figure 7f). Aitchison
808 index demonstrated significant differences every pairwise of the periodontitis stages (PERMANOVA
809 test $p \leq 0.001$; Table 7). Significant differences in the distances between periodontitis severities further
810 demonstrated the uniqueness of each severity of periodontitis (MWU test $p \leq 0.05$; Figure 7g-j).

811 **3.3.3 DAT among multiple periodontitis severities and their correlation**

812 Of the 425 total taxa that identified in the salivary microbiome composition (Figure 13), 20 DAT were
813 identified (Table 5). Three separate subgroups were formed from the participants-level abundances of the
814 DAT using a hierarchical clustering methodology (Figure 8a):

- 815 • Group 1
- 816 1. *Treponema* spp.
- 817 2. *Prevotella* sp. HMT 304
- 818 3. *Prevotella* sp. HMT 526
- 819 4. *Peptostreptococcaceae [XI][G-5]* saphenum
- 820 5. *Treponema* sp. HMT 260
- 821 6. *Mycoplasma faecium*
- 822 7. *Peptostreptococcaceae [XI][G-9]* brachy
- 823 8. *Lachnospiraceae [G-8]* bacterium HMT 500
- 824 9. *Peptostreptococcaceae [XI][G-6]* nodatum
- 825 10. *Fretibacterium* spp.

- 826 • Group 2
- 827 1. *Porphyromonas gingivalis*
- 828 2. *Campylobacter showae*
- 829 3. *Filifactor alocis*
- 830 4. *Treponema putidum*
- 831 5. *Tannerella forsythia*
- 832 6. *Prevotella intermedia*
- 833 7. *Porphyromonas* sp. HMT 285

- 834 • Group 3
- 835 1. *Actinomyces* spp.
- 836 2. *Corynebacterium durum*
- 837 3. *Actinomyces graevenitzii*

838 Ten DAT that were significant enriched in stage II and stage III, but deficient in healthy formed Group
839 1 (Figure 8). Furthermore, in comparison to the healthy, the seven DAT of Group 2 were significantly
840 enriched in each of the stages of periodontitis. On the other hand, three DAT in Group 3 were deficient in
841 stage II and stage III, but significantly enriched in healthy. The relative proportions of the DAT further
842 supported these findings (Figure 8b), suggesting that the DAT is primarily linked to periodontitis rather
843 than other salivary bacteria.

844 Correlation analysis from the DAT showed that DAT from Group 3 was negatively correlated with
845 Group 1 and Group 2 (Figure 9), and strong correlations were observed the nine pairs of DAT (Figure 14).

846 3.3.4 Classification of periodontitis severities by random forest models

847 To confirm that using selected DAT bacterial profiles could have enhanced sequencing expenses without
848 losing the classification evaluations, we built the random forest classification models based on DAT and
849 full microbiome compositions (Figure 18). DAT based classifier showed non-significant different or better
850 evaluations, by removing confounding taxa.

851 Based on the proportion of DAT, random forest classifier were trained to classify the periodontitis
852 stages (Table 6). We conducted multi-label classification for the multiple periodontitis stages, namely
853 healthy, stage I, stage II, and stage III. In this setting, we classified multiple periodontitis severities with
854 the highest BA of 0.779 ± 0.029 (mean \pm SD) (Table 4). AUC ranged between 0.81 and 0.94 (Figure 10b).

855 Since timely detection in dentistry is demanding (Tonetti et al., 2018), we implemented a random
856 forest classification for both healthy and stage I. Remarkably, the random forest classifier had the highest
857 BA at 0.793 ± 0.123 (mean \pm SD) (Table 4). In this setting, this model showed high AUC value for the
858 classifying of stage I from healthy (AUC=0.85; Figure 10d).

859 Based on the findings that the salivary microbiome composition in stage II is more comparable to
860 those in stage III than to other severities (Figure 7f and Figure 7j), we combined stage II and stage III to
861 perform a multi-label classification.

862 To examine alternative classification algorithms in comparison to random forest classification, we
863 selected gradient boost algorithm because it is another algorithm of the few classification algorithms
864 that can provide feature importances, which is essential for identifying key taxa contributing to the
865 classification of periodontitis stages. Thus, we assessed gradient boosting algorithms (Figure 20). However,
866 the classification evaluations obtained from gradient boosting have non-significant differences compared
867 to random forest classification.

868 Finally, to confirm the reliability and consistency of our random forest classifier, we validated our
869 classification model using openly accessible 16S rRNA gene sequencing from Spanish participants
870 (Iniesta et al., 2023) and Portuguese participants (Relvas et al., 2021) (Figure 11). Although some
871 evaluations, *e.g.* SPE, were low, the other were comparable.

Table 3: Clinical characteristics of the study participants.

Continuous variable: mean \pm SD. Categorical variable: count (proportion). Significant differences were assessed using the Kruskal-Wallis test. NA: Not applicable.

Index	Healthy	Stage I	Stage II	Stage III	p-value
Age (year)	33.83 \pm 13.04	43.30 \pm 14.28	50.26 \pm 11.94	51.08 \pm 11.13	6.18E-17
Gender (Male)	44 (44.0%)	22 (44.0%)	25 (50.0%)	25 (50.0%)	NA
Smoking (Never)	83 (83.0%)	36 (72.0%)	34 (68.0%)	29 (58.0%)	NA
Smoking (Ex)	12 (12.0%)	7 (14.0%)	9 (18.0%)	10 (20.0%)	NA
Smoking (Current)	2 (2.0%)	7 (14.0%)	7 (14.0%)	10 (20.0%)	NA
Number of teeth	28.03 \pm 2.23	27.36 \pm 1.80	26.72 \pm 2.89	25.74 \pm 4.34	8.07E-05
Attachment level (mm)	2.45 \pm 0.29	2.75 \pm 0.38	3.64 \pm 0.83	4.54 \pm 1.14	1.82E-35
Probing depth (mm)	2.42 \pm 0.29	2.61 \pm 0.40	3.27 \pm 0.76	3.95 \pm 0.88	6.43E-28
Plaque index	17.66 \pm 16.21	35.46 \pm 23.75	54.40 \pm 23.79	58.30 \pm 25.25	3.23E-22
Gingival index	0.09 \pm 0.16	0.44 \pm 0.46	0.85 \pm 0.52	1.06 \pm 0.52	2.59E-32

Table 4: Feature combinations and their evaluations.

Classification performance with the most important taxon, the two most important taxa, and taxa with the best balanced accuracy (mean \pm SD). *P. gingivalis* and *Act.* are *Porphyromonas gingivalis* and *Actinomyces* spp., respectively

Classification	Features	ACC	AUC	BA	F1	PRE	SEN	SPE
Healthy vs. Stage I vs. Stage II vs. Stage III	<i>P.gingivalis</i>	0.758 \pm 0.051	0.716 \pm 0.177	0.677 \pm 0.068	0.839 \pm 0.034	0.839 \pm 0.034	0.516 \pm 0.102	
	<i>P.gingivalis+Act.</i>	0.792 \pm 0.043	0.822 \pm 0.105	0.723 \pm 0.057	0.861 \pm 0.029	0.861 \pm 0.029	0.584 \pm 0.086	
Top 5 taxa		0.834 \pm 0.022	0.870 \pm 0.079	0.779 \pm 0.029	0.889 \pm 0.015	0.889 \pm 0.015	0.668 \pm 0.033	
Healthy vs. Stage I	<i>Act.</i>	0.687 \pm 0.116	0.725 \pm 0.145	0.647 \pm 0.159	0.762 \pm 0.092	0.760 \pm 0.128	0.781 \pm 0.116	0.513 \pm 0.224
	<i>Act.+P.gingivalis</i>	0.733 \pm 0.119	0.831 \pm 0.081	0.713 \pm 0.122	0.797 \pm 0.097	0.798 \pm 0.126	0.798 \pm 0.082	0.627 \pm 0.191
Top 9 taxa		0.800 \pm 0.103	0.852 \pm 0.103	0.793 \pm 0.123	0.849 \pm 0.080	0.850 \pm 0.112	0.857 \pm 0.090	0.730 \pm 0.193
Healthy vs. Stage I vs. Stages II/III	<i>P.gingivalis</i>	0.776 \pm 0.042	0.736 \pm 0.196	0.748 \pm 0.047	0.832 \pm 0.031	0.832 \pm 0.031	0.664 \pm 0.062	
	<i>P.gingivalis+Act.</i>	0.843 \pm 0.035	0.876 \pm 0.109	0.823 \pm 0.039	0.882 \pm 0.026	0.882 \pm 0.026	0.764 \pm 0.052	
Top 6 taxa		0.885 \pm 0.036	0.914 \pm 0.027	0.871 \pm 0.038	0.914 \pm 0.025	0.914 \pm 0.025	0.828 \pm 0.051	
Healthy vs. Stages I/II/III	<i>P.gingivalis</i>	0.792 \pm 0.114	0.856 \pm 0.105	0.819 \pm 0.088	0.776 \pm 0.089	0.840 \pm 0.092	0.756 \pm 0.175	0.883 \pm 0.054
	<i>P.gingivalis+Act.</i>	0.828 \pm 0.121	0.926 \pm 0.074	0.847 \pm 0.116	0.797 \pm 0.123	0.800 \pm 0.126	0.830 \pm 0.191	0.864 \pm 0.074
Top 4 taxa		0.860 \pm 0.078	0.953 \pm 0.049	0.885 \pm 0.066	0.832 \pm 0.079	0.840 \pm 0.128	0.864 \pm 0.157	0.905 \pm 0.070

Table 5: **List of DAT among healthy status and periodontitis stages.** Statistical significance was determined by ANCOM W value.

No.	Taxonomy	ANCOM W score
1	<i>Porphyromonas gingivalis</i>	424
2	<i>Actinomyces</i> spp.	424
3	<i>Filifactor alocis</i>	421
4	<i>Prevotella intermedia</i>	419
5	<i>Treponema putidum</i>	418
6	<i>Tannerella forsythia</i>	415
7	<i>Porphyromonas</i> sp. HMT 285	412
8	<i>Peptostreptococcaceae [XI][G-6] nodatum</i>	412
9	<i>Fretibacterium</i> spp.	411
10	<i>Mycoplasma faecium</i>	411
11	<i>Prevotella</i> sp. HMT 304	411
12	<i>Lachnospiraceae [G-8] bacterium</i> HMT 500	409
13	<i>Treponema</i> spp.	408
14	<i>Prevotella</i> sp. HMT 526	401
15	<i>Peptostreptococcaceae [XI][G-9] brachy</i>	400
16	<i>Peptostreptococcaceae [XI][G-5] saphenum</i>	398
17	<i>Campylobacter showae</i>	395
18	<i>Treponema</i> sp. HMT 260	393
19	<i>Corynebacterium durum</i>	393
20	<i>Actinomyces graevenitzii</i>	387

Table 6: Feature the importance of taxa in the classification of different periodontal statuses.

Taxa are ranked in descending order of importance; from most important to least important. Note that $\forall i, 0 \geq \text{importance}_i \geq 1$ and $\sum_i \text{importance}_i = 1$.

Condition	Healthy vs. Stage I vs. Stage II vs. Stage III			Healthy vs. Stage I vs. Stage II/III			Healthy vs. Stage I/II/III		
	Rank	Taxa	Importance	Taxa	Importance	Taxa	Importance	Taxa	Importance
1	<i>Porphyromonas gingivalis</i>	0.297	<i>Actinomyces spp.</i>	0.360	<i>Porphyromonas gingivalis</i>	0.426	<i>Porphyromonas gingivalis</i>	0.461	
2	<i>Actinomyces spp.</i>	0.195	<i>Porphyromonas gingivalis</i>	0.125	<i>Actinomyces spp.</i>	0.244	<i>Actinomyces spp.</i>	0.257	
3	<i>Prevotella intermedia</i>	0.054	<i>Actinomyces graevenitzii</i>	0.095	<i>Actinomyces graevenitzii</i>	0.049	<i>Actinomyces spp.</i>	0.059	
4	<i>Actinomyces graevenitzii</i>	0.052	<i>Porphyromonas sp. HMT 285</i>	0.062	<i>Corynebacterium durum</i>	0.046	<i>Corynebacterium durum</i>	0.035	
5	<i>Filifactor alocis</i>	0.050	<i>Lachnospiraceae (G-8) bacterium HMT 500</i>	0.052	<i>Filifactor alocis</i>	0.036	<i>Filifactor alocis</i>	0.032	
6	<i>Campylobacter showae</i>	0.042	<i>Campylobacter showae</i>	0.050	<i>Prevotella intermedia</i>	0.033	<i>Campylobacter showae</i>	0.023	
7	<i>Porphyromonas sp. HMT 285</i>	0.040	<i>Filifactor alocis</i>	0.039	<i>Tannerella forsythia</i>	0.025	<i>Porphyromonas sp. HMT 285</i>	0.022	
8	<i>Corynebacterium durum</i>	0.032	<i>Corynebacterium durum</i>	0.038	<i>Campylobacter showae</i>	0.023	<i>Prevotella intermedia</i>	0.022	
9	<i>Treponema spp.</i>	0.032	<i>Treponema spp.</i>	0.037	<i>Treponema sp. HMT 285</i>	0.021	<i>Treponema spp.</i>	0.022	
10	<i>Tannerella forsythia</i>	0.026	<i>Tannerella forsythia</i>	0.029	<i>Treponema spp.</i>	0.018	<i>Peptostreptococcaceae (XII/G-9) brachy</i>	0.015	
11	<i>Treponema pritulum</i>	0.025	<i>Prevotella intermedia</i>	0.026	<i>Peptostreptococcaceae (XII/G-9) brachy</i>	0.014	<i>Lachnospiraceae (G-8) bacterium HMT 500</i>	0.010	
12	<i>Freibacterium spp.</i>	0.023	<i>Freibacterium spp.</i>	0.018	<i>Lachnospiraceae (G-8) bacterium HMT 500</i>	0.011	<i>Tannerella forsythia</i>	0.009	
13	<i>Peptostreptococcaceae (XII/G-9) brachy</i>	0.021	<i>Peptostreptococcaceae (XII/G-9) brachy</i>	0.018	<i>Peptostreptococcaceae (XII/G-6) nodatum</i>	0.010	<i>Freibacterium spp.</i>	0.009	
14	<i>Treponema sp. HMT 260</i>	0.019	<i>Treponema pritulum</i>	0.014	<i>Treponema pritulum</i>	0.009	<i>Treponema pritulum</i>	0.006	
15	<i>Prevotella sp. HMT 526</i>	0.018	<i>Prevotella sp. HMT 526</i>	0.011	<i>Prevotella sp. HMT 526</i>	0.008	<i>Peptostreptococcaceae (XII/G-6) nodatum</i>	0.004	
16	<i>Peptostreptococcaceae (XII/G-6) nodatum</i>	0.018	<i>Treponema sp. HMT 260</i>	0.008	<i>Freibacterium spp.</i>	0.008	<i>Treponema sp. HMT 260</i>	0.004	
17	<i>Prevotella sp. HMT 304</i>	0.017	<i>Peptostreptococcaceae (XII/G-6) nodatum</i>	0.008	<i>Treponema sp. HMT 260</i>	0.005	<i>Mycoplasma faecium</i>	0.004	
18	<i>Mycoplasma faecium</i>	0.014	<i>Mycoplasma faecium</i>	0.004	<i>Prevotella sp. HMT 304</i>	0.005	<i>Prevotella sp. HMT 326</i>	0.003	
19	<i>Peptostreptococcaceae (XII/G-5) saphenum</i>	0.014	<i>Prevotella sp. HMT 304</i>	0.003	<i>Mycoplasma faecium</i>	0.005	<i>Peptostreptococcaceae (XII/G-5) saphenum</i>	0.002	
20	<i>Lachnospiraceae (G-8) bacterium HMT 500</i>	0.013	<i>Peptostreptococcaceae (XII/G-5) saphenum</i>	0.003	<i>Peptostreptococcaceae (XII/G-5) saphenum</i>	0.004	<i>Prevotella sp. HMT 304</i>	0.001	

Table 7: Beta-diversity pairwise comparisons on the periodontitis statuses

Statistically significant (p-value) was determined by the PERMANOVA test.

Group 1	Group 2	p-value
Healthy	Stage I	0.001
Healthy	Stage II	0.001
Healthy	Stage III	0.001
Stage I	Stage II	0.001
Stage I	Stage III	0.001
Stage II	Stage III	0.737

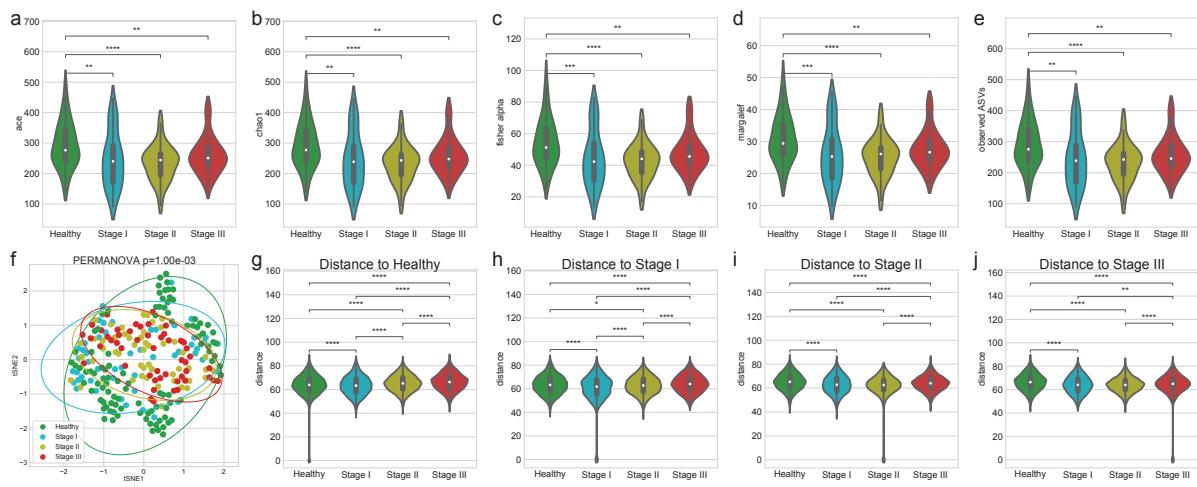


Figure 7: Diversity indices for periodontitis.

Alpha-diversity indices (a-e) indicate that healthy controls have increased heterogeneity than periodontitis stages as measured by: (a) ACE (b) Chao1 (c) Fisher alpha (d) Margalef, and (e) observed ASVs. (f) The beta-diversity index (weighted UniFrac) was visualized using a tSNE-transformed plot. The confidence ellipses are shown to display the distribution of each periodontitis stage. The distance to each stage demonstrated that each periodontitis stage was distinguished from the other periodontitis stages: (g) distance to Healthy (h) distance to Stage I (i) distance to Stage II, and (j) distance to Stage III. Statistical significance determined by the MWU test and the PERMANOVA test: $p < 0.05$ (*), $p < 0.01$ (**), $p < 0.001$ (***) ≤ 0.0001 (****).

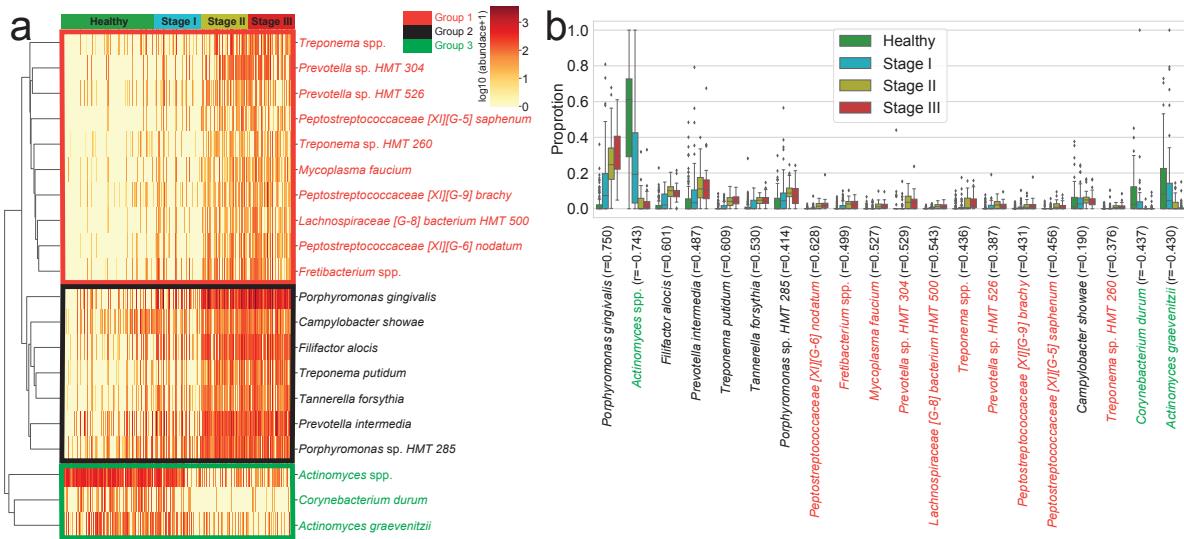


Figure 8: DAT for periodontitis.

DAT that were identified by ANCOM. **(a)** Heatmap of clustered DAT with similar distribution among subjects. Group 1, Group 2, and Group 3 are marked in red, black, and green, respectively. **(b)** Box plots showing the proportions of DAT. Taxa were sorted by their importance according to ANCOM.

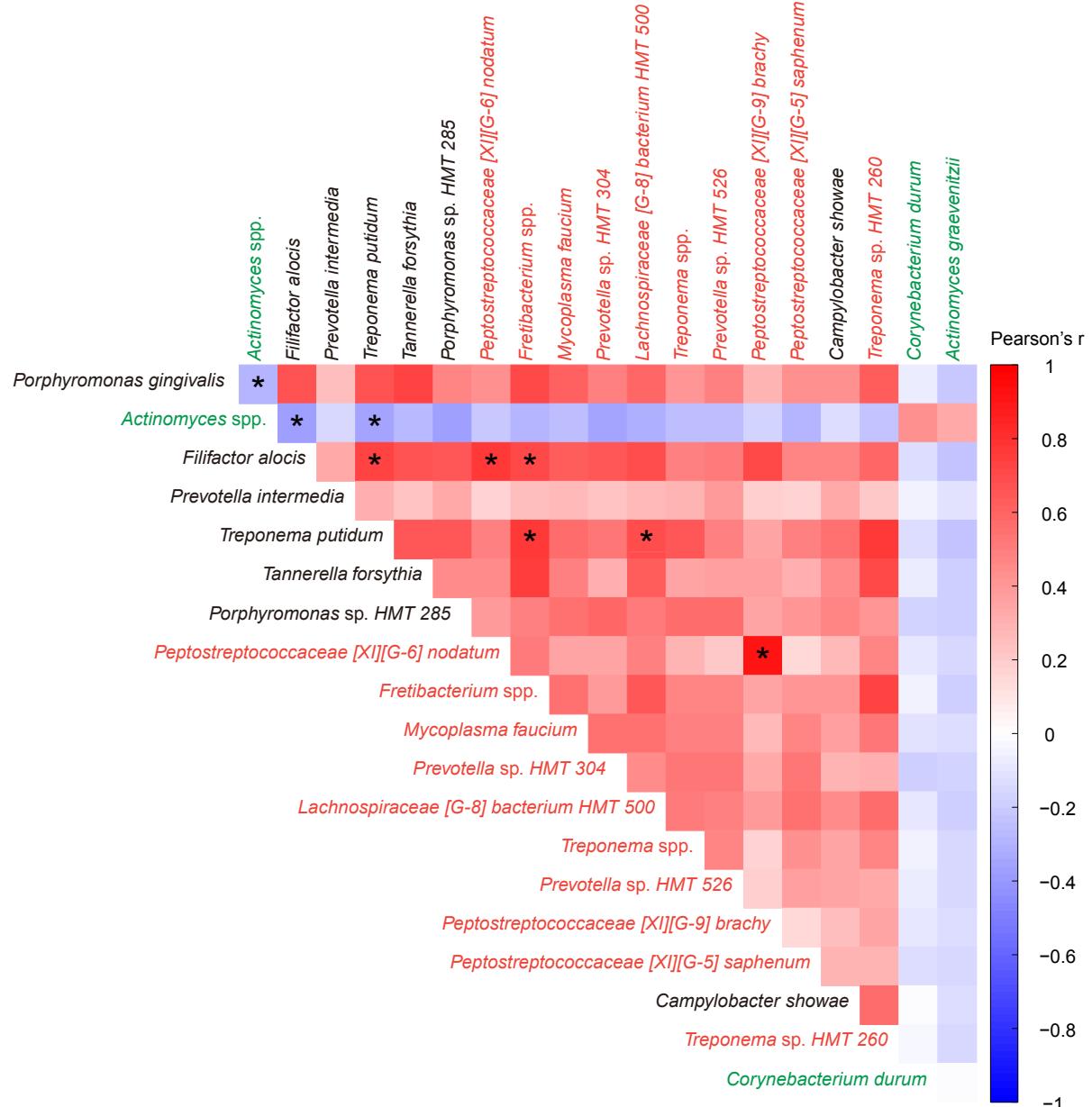


Figure 9: Correlation heatmap between periodontitis DAT.

Pearson's correlations between DAT in healthy status and periodontitis stages. Statistical significance was determined by strong Pearson correlation, i.e., $| \text{coefficient} | \geq 0.5$ (*).

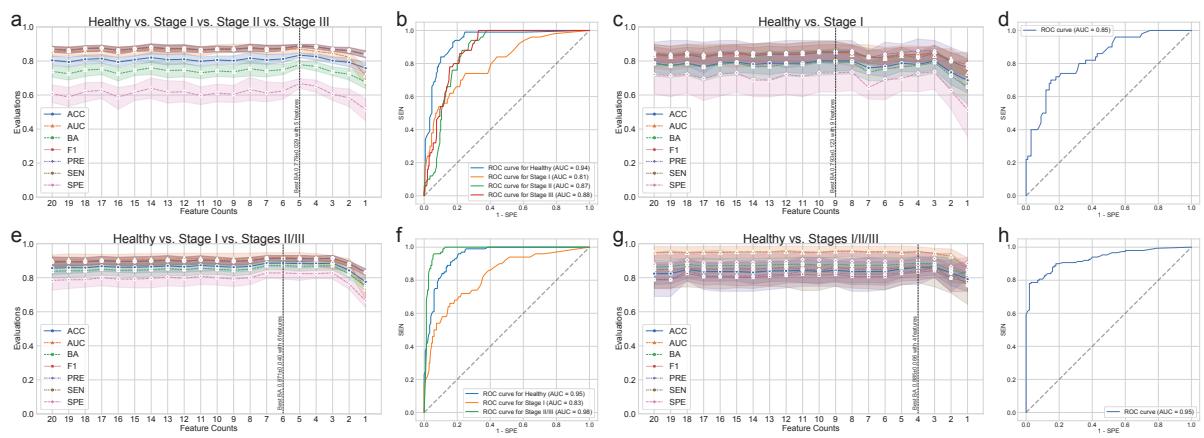


Figure 10: Random forest classification metrics for periodontitis prediction.

The classification metrics in the random forest classifications were as follows: ACC, AUC, BA, F1, PRE, SEN, and SPE. **(a)** Classification performance for healthy vs. stage I vs. stage II vs. stage III. **(b)** ROC curve for the highest BA of (a). **(c)** Classification performance for healthy vs. stage I. **(d)** ROC curve on the highest BA of (c). **(e)** Classification performance for healthy vs. stage I vs. stages II/III. **(f)** ROC curve for the highest BA of (e). **(g)** Classification performance for healthy vs. stages I/II/III. **(h)** ROC curve for the highest BA of (h).

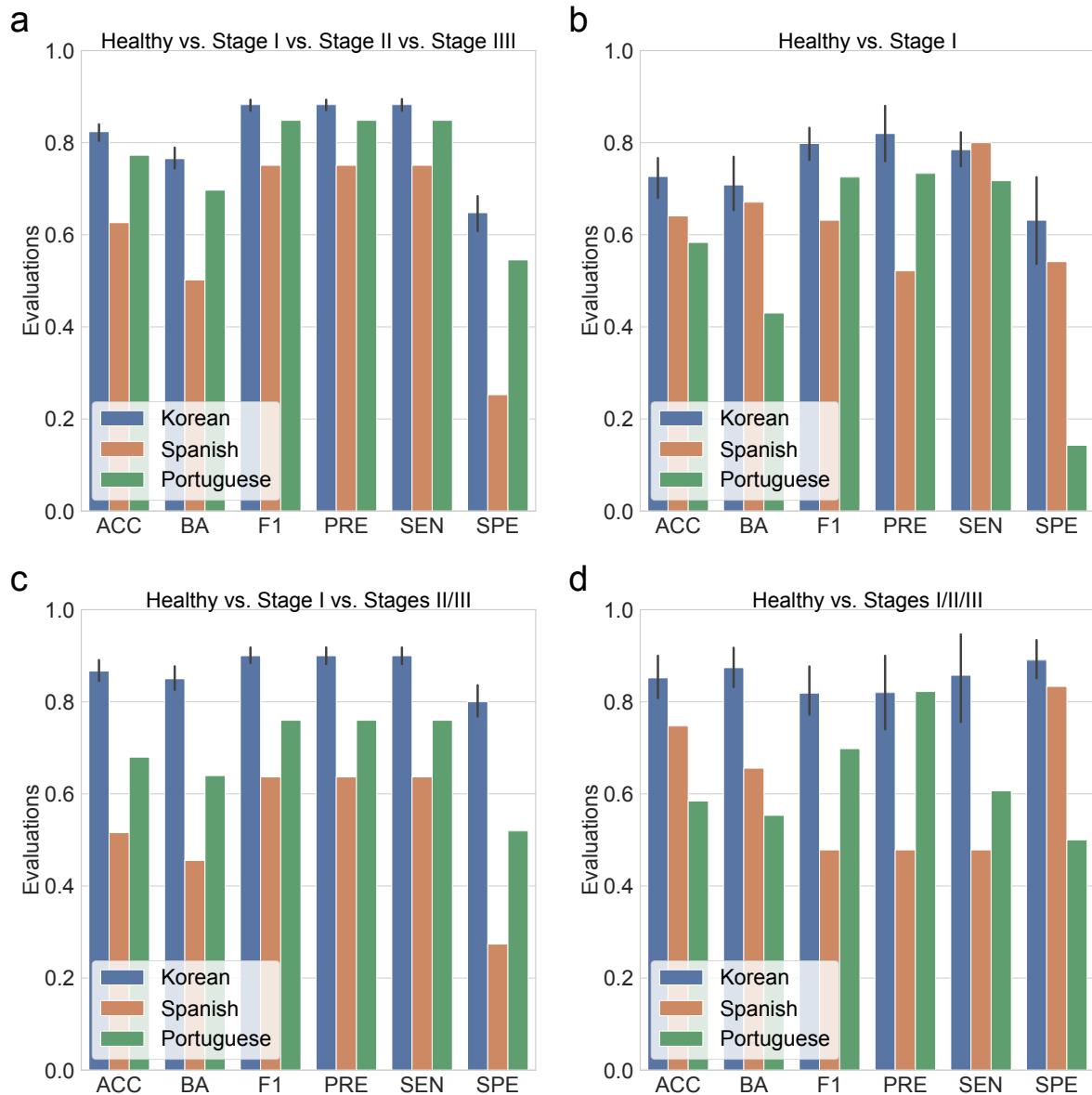


Figure 11: **Random forest classification metrics from external datasets.**

The classification metrics in the random forest classifications were as follows: ACC, AUC, BA, F1, PRE, SEN, and SPE. **(a)** Classification performance for healthy vs. stage I vs. stage II vs. stage III. **(b)** Classification performance for healthy vs. stage I. **(c)** Classification performance for healthy vs. stage I vs. stages II/III. **(d)** Classification performance for healthy vs. stages I/II/III.

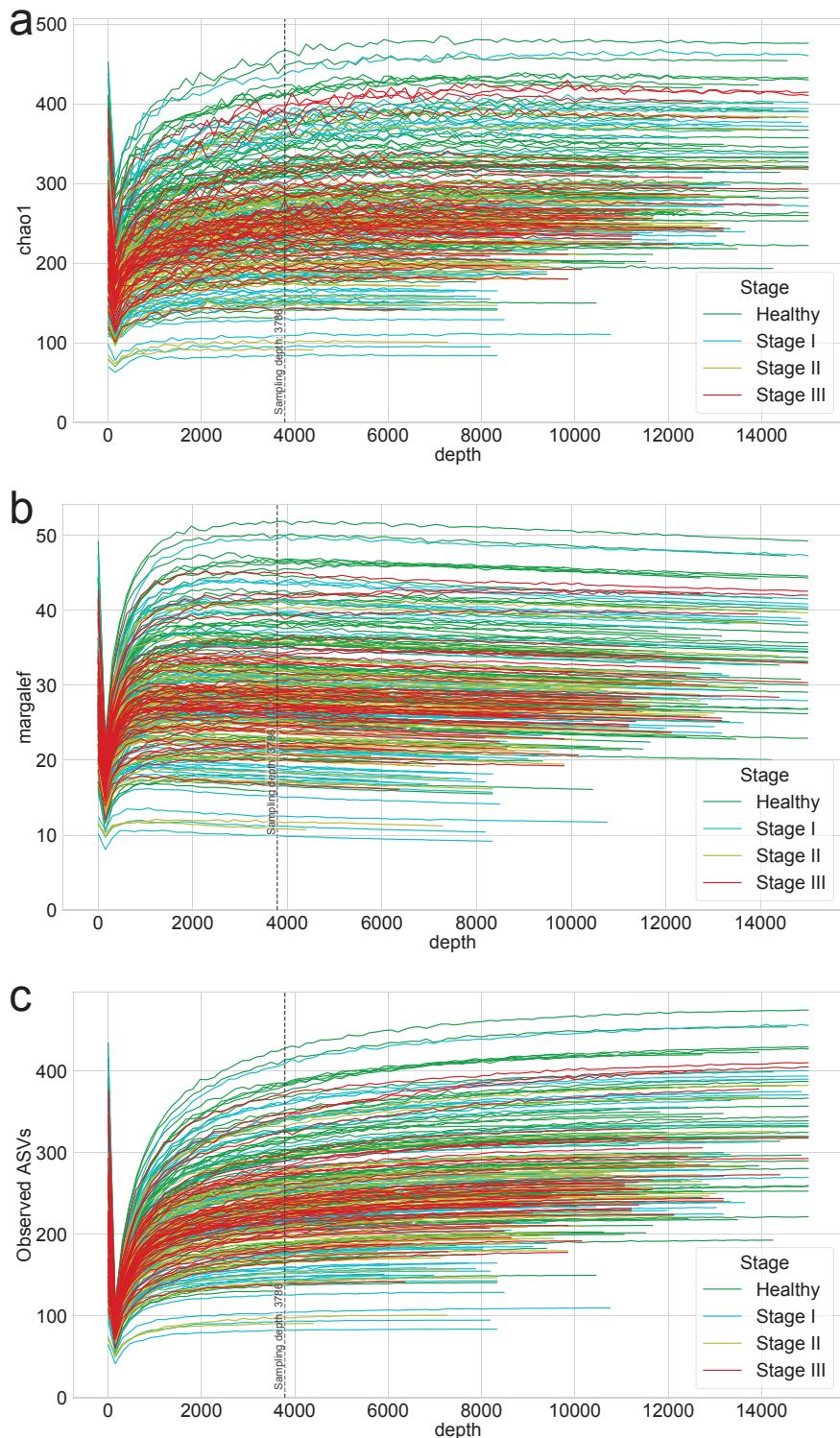


Figure 12: Rarefaction curves for alpha-diversity indices.

Rarefaction of (a) chao1 (b) margalef, and (c) observed ASVs were generated to measure species richness and determine the sampling depth of each sample.

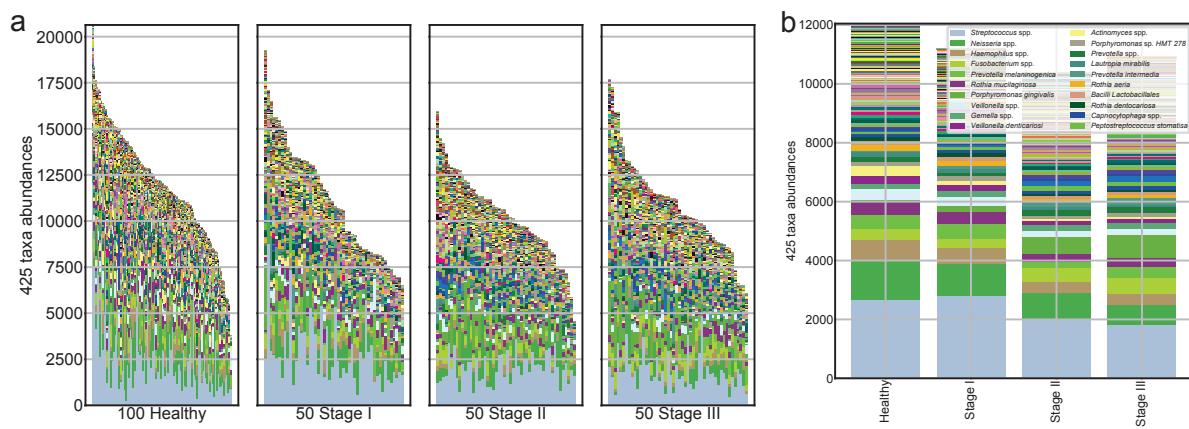


Figure 13: Salivary microbiome compositions in the different periodontal stages.

Stacked bar plot of the absolute abundance of bacterial species for all samples (**a**) and the mean absolute abundance of bacterial species in the healthy, stage I, stage II, and stage III groups (**b**).

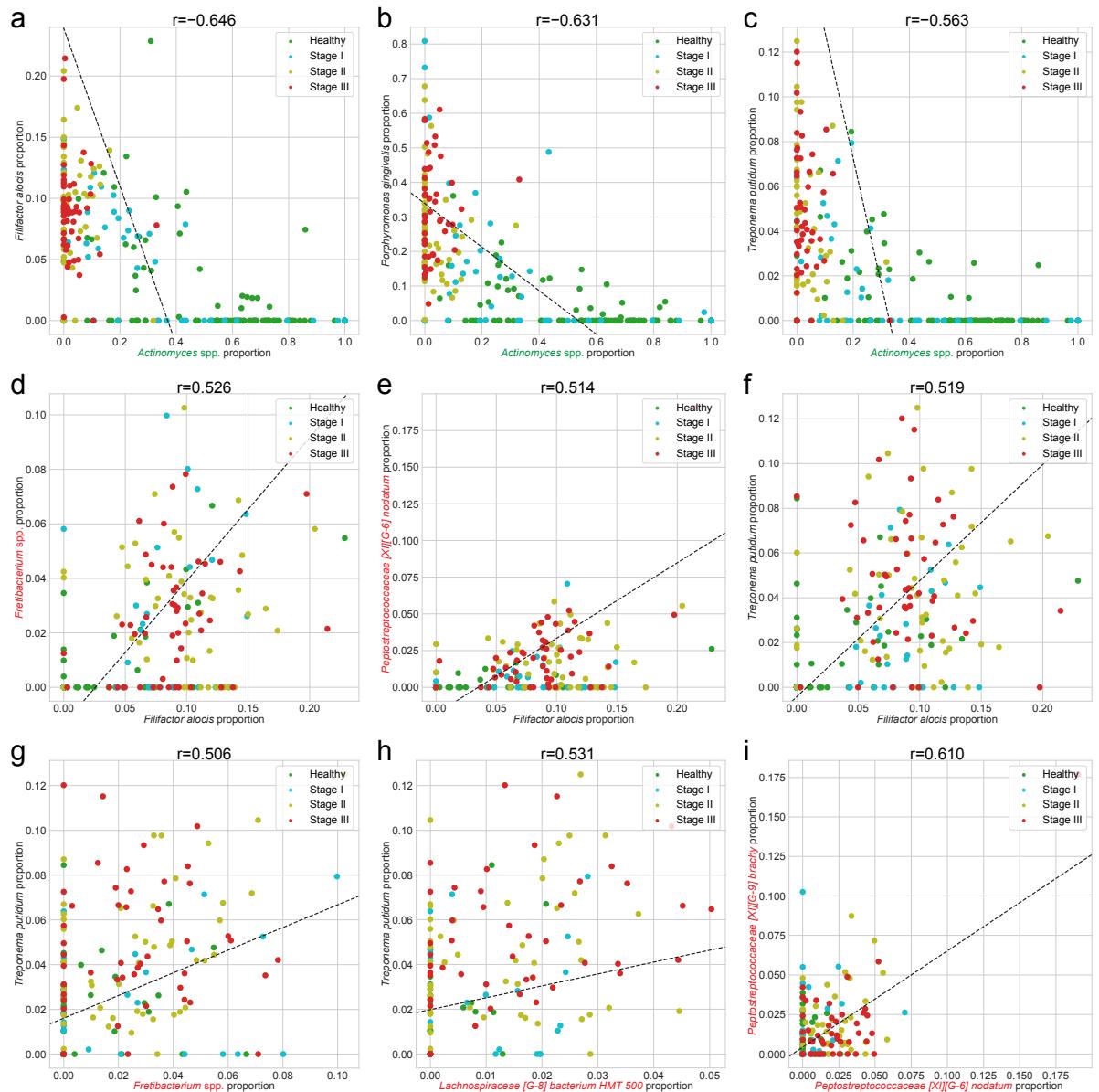


Figure 14: Correlation plots for periodontitis DAT.

We selected the combinations of DAT with absolute Spearman correlation coefficients greater than 0.5. The color represents periodontal healthy periodontal stages (green: healthy, cyan: stage I, yellow: stage II, and red: stage III).

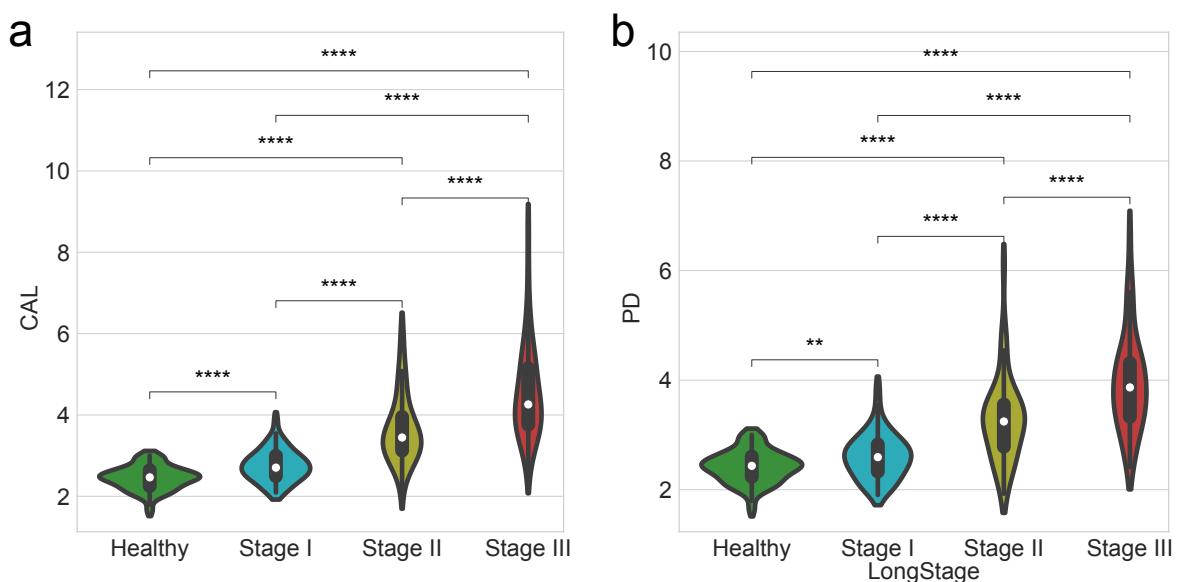


Figure 15: Clinical measurements by the periodontitis stages.

Comparisons of clinical measurement among healthy controls and patients with various periodontitis stages. **(a)** Clinical attachment level (CAL) **(b)** Probing depth (PD). Statistical significance determined by the MWU test: $p < 0.01$ (**) and $p < 0.0001$ (****).

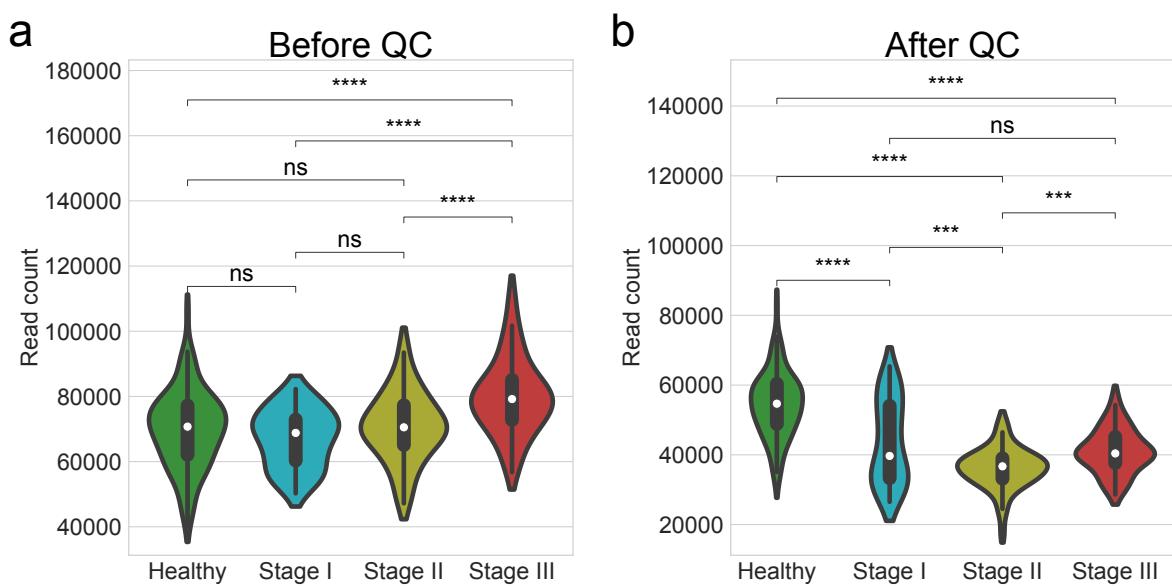


Figure 16: Number of read counts by the periodontitis stages.

Comparisons of the number of read counts among healthy controls and patients with various periodontitis stages. **(a)** Before quality check **(b)** After quality check. Statistical significance determined by the MWU test: $p \geq 0.05$ (ns), $p < 0.05$ (*), $p < 0.01$ (**), $p < 0.001$ (***) $,$ and $p < 0.0001$ (****).

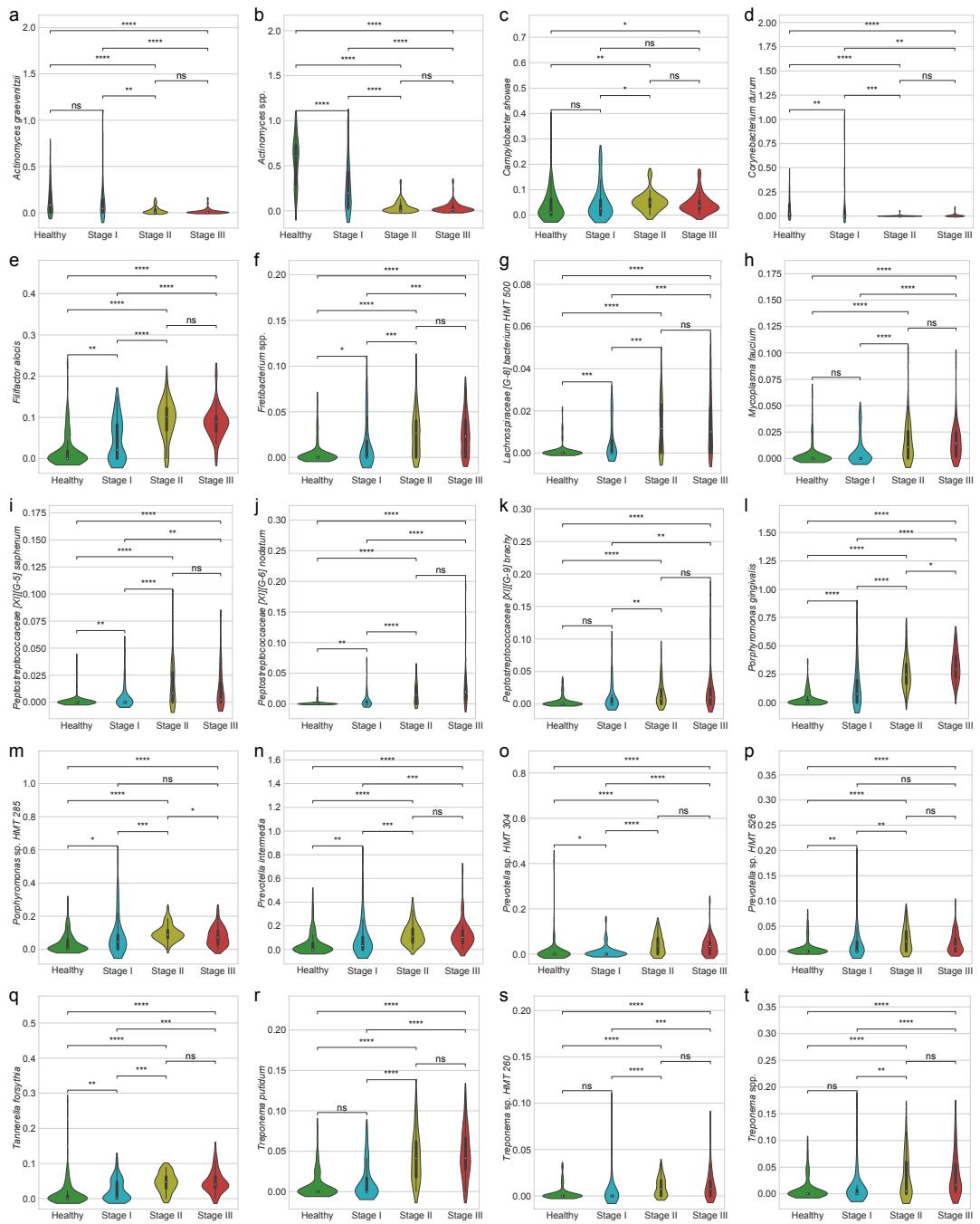


Figure 17: Proportions of periodontitis DAT.

(a) *Actinomyces graevenitzii* **(b)** *Actinomyces* spp. **(c)** *Campylobacter showae* **(d)** *Corynebacterium durum* **(e)** *Filifactor alocis* **(f)** *Fretibacterium* spp. **(g)** *Lachnospiraceae* [G-8] bacterium HMT 500 **(h)** *Mycoplasma faecium* **(i)** *Peptostreptococcaceae* [XI][G-5] saphenum **(j)** *Peptostreptococcaceae* [XI][G-6] nodatum **(k)** *Peptostreptococcaceae* [XI][G-9] brachy **(l)** *Porphyromonas gingivalis* **(m)** *Porphyromonas* sp. HMT 285 **(n)** *Prevotella* intermedia **(o)** *Prevotella* sp. HMT 304 **(p)** *Prevotella* sp. HMT 526 **(q)** *Tannerella forsythia* **(r)** *Treponema putidum* **(s)** *Treponema* sp. HMT 260 **(t)** *Treponema* spp. Statistical significance determined by the MWU test: $p \geq 0.05$ (ns), $p < 0.05$ (*), $p < 0.01$ (**), $p < 0.001$ (***), and $p < 0.0001$ (****).

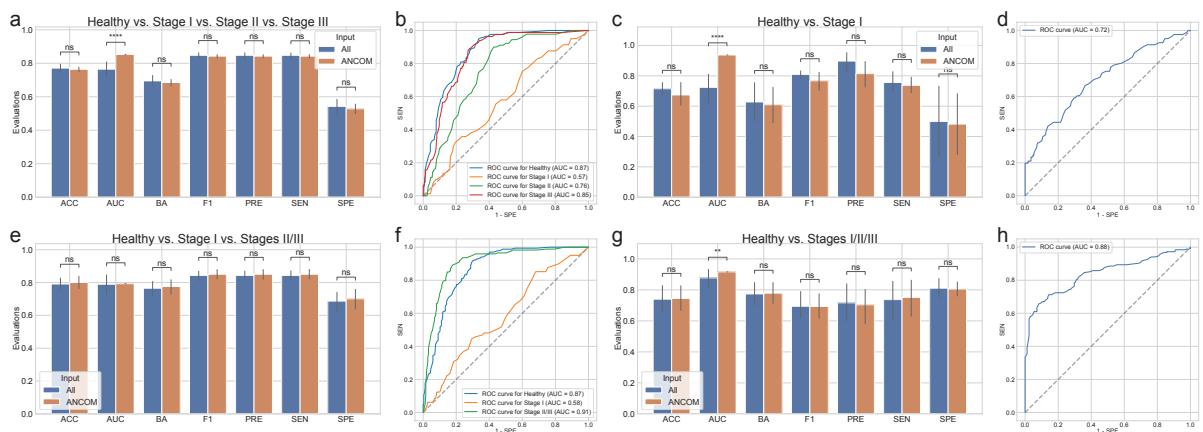


Figure 18: Random forest classification metrics with the full microbiome compositions and ANCOM-selected DAT compositions.

The classification metrics in the random forest classifications were as follows: ACC, AUC, BA, F1, PRE, SEN, and SPE. **(a)** Classification performance for healthy vs. stage I vs. stage II vs. stage III. **(b)** ROC curve for the highest BA of (a). **(c)** Classification performance for healthy vs. stage I. **(d)** ROC curve on the highest BA of (c). **(e)** Classification performance for healthy vs. stage I vs. stages II/III. **(f)** ROC curve for the highest BA of (e). **(g)** Classification performance for healthy vs. stages I/II/III. **(h)** ROC curve for the highest BA of (g). Statistical significance determined by the MWU test: $p \geq 0.05$ (ns), $p < 0.01$ (**), and $p < 0.0001$ (***).

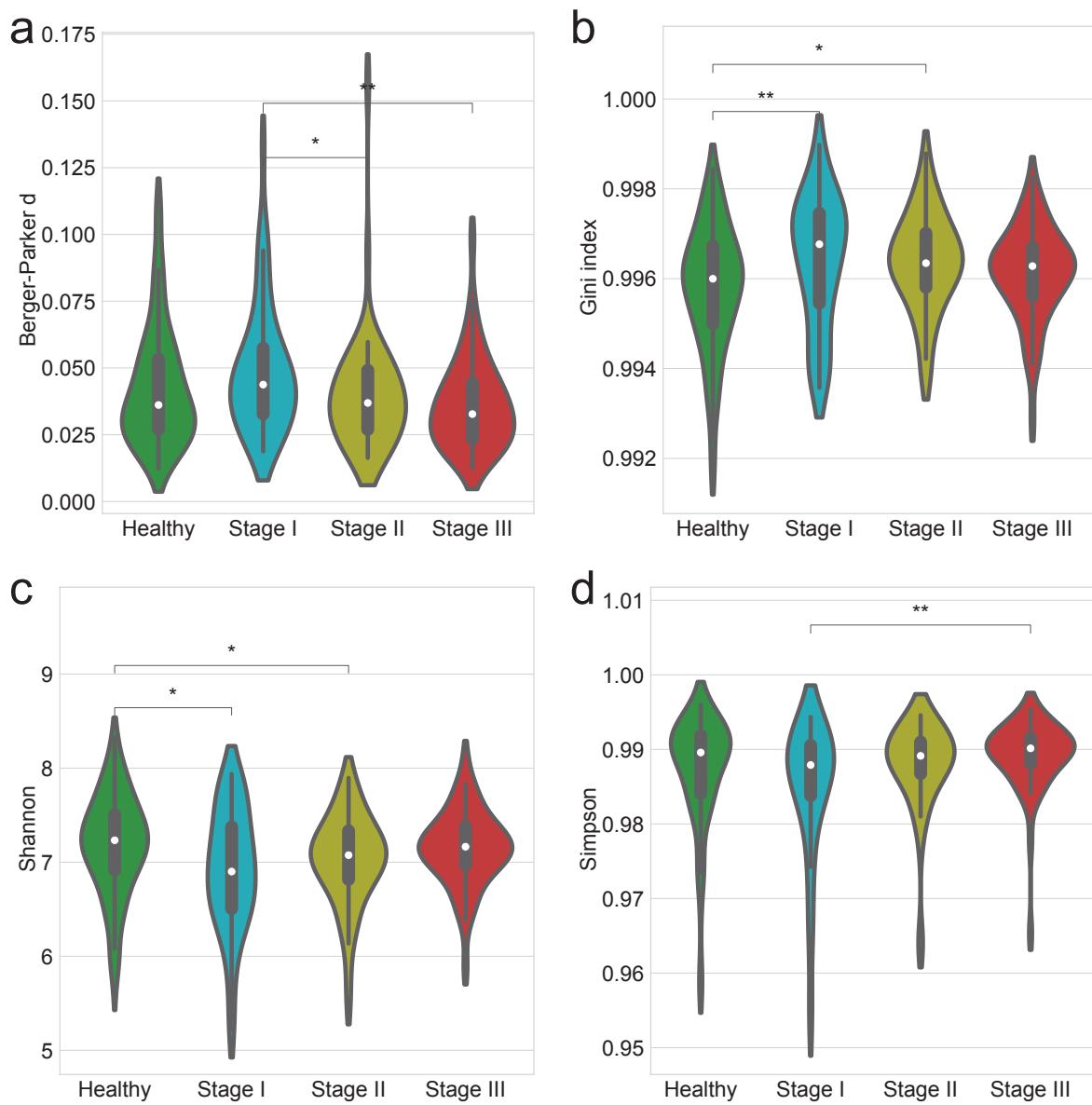


Figure 19: Alpha-diversity indices account for evenness.

Alpha-diversity indices (**a-d**) indicate that the heterogeneity between the periodontitis stages as measured by: **(a)** Berger-Parker *d* **(b)** Gini **(c)** Shannon **(d)** Simpson. Statistical significance determined by the MWU test: $p < 0.05$ (*) and $p < 0.01$ (**)

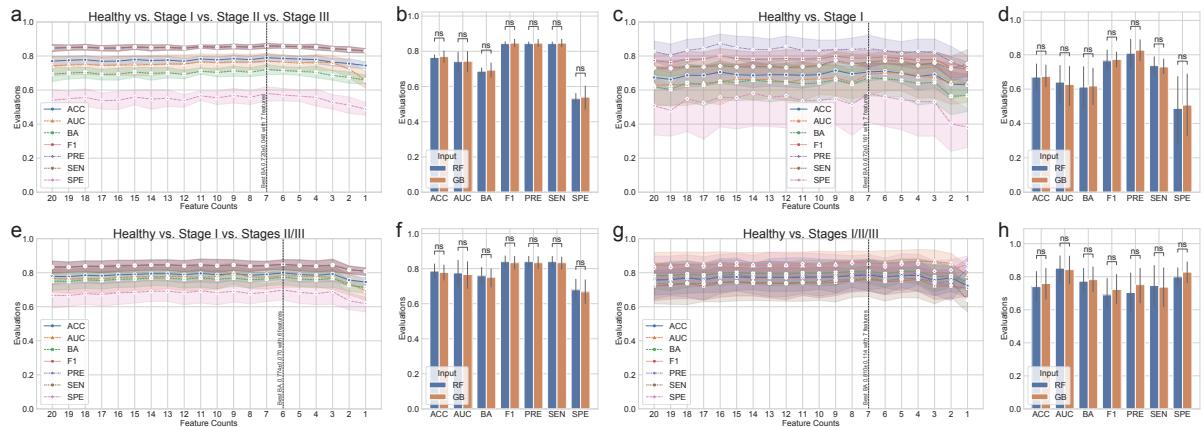


Figure 20: Gradient Boosting classification metrics for periodontitis prediction.

The classification metrics in the random forest classifications were as follows: ACC, AUC, BA, F1, PRE, SEN, and SPE. The feature counts mean that the classification model trained on the most important n features as the Table 5. **(a)** Comparison of Random forest (RF) and Gradient boosting (GB) for healthy vs. stage I vs. stage II vs. stage III. **(b)** Comparison of RF and GB for the highest BA of (a). **(c)** Classification performance for healthy vs. stage I. **(d)** Comparison of RF and GB for healthy vs. stage I vs. stages II/III. **(e)** Comparison of RF and GB for the highest BA of (d). **(f)** Comparison of RF and GB for Healthy vs. Stage I vs. Stages II/III. **(g)** Classification performance for healthy vs. stages I/II/III. **(h)** Comparison of RF and GB for Healthy vs. Stages I/II/III. MWU test: $p \geq 0.05$ (ns)

872 **3.4 Discussion**

873 In order to investigate at potential alterations in the salivary microbiome compositions based on periodontal
874 stages, including healthy, stage I, stage II, and stage III, we employed 16S rRNA gene sequencing to
875 perform a cross-sectional periodontitis analysis. In this study, the 2018 periodontitis classification served
876 as the basis for the classification of periodontitis stages (Papapanou et al., 2018). There were notable
877 variations in the salivary microbiome composition among the multiple stages of periodontitis (Figure
878 13). Furthermore, our random forest classification model based on the proportions of DAT in the salivary
879 microbiome compositions across study participants to predict multiple periodontitis statuses with high
880 AUC of 0.870 ± 0.079 (mean \pm SD) (Table 4).

881 Previous research identified the red complex as the primary pathogens of periodontitis (Listgarten,
882 1986): *Porphyromonas gingivalis*, *Tannerella forsythia*, and *Treponema denticola*. Other studies, however,
883 have shown that periodontal pathogens communicate with other bacteria in the salivary microbiome
884 networks to generate dental plaque prior to the pathogenesis and development of periodontitis (Lamont &
885 Jenkinson, 2000; Rosan & Lamont, 2000; Yoshimura, Murakami, Nishikawa, Hasegawa, & Kawaminami,
886 2009).

887 Using subgingival plaque collections, recent researches have suggested a connection between the
888 periodontitis stage and the salivary microbiome compositions (Altabtbaei et al., 2021; Iniesta et al., 2023;
889 Nemoto et al., 2021). Therefore, we have examined the salivary microbiome compositions of patients
890 with multiple stages of periodontitis and periodontally healthy controls, extending on earlier studies.

891 According to our findings, the salivary microbiome compositions have 425 taxa (Figure 13). We
892 computed the alpha-diversity indices to determine the variability within each salivary microbiome
893 composition, including ace (Chao & Lee, 1992), chao1 (Chao, 1984), fisher alpha (Fisher et al., 1943),
894 margalef (Magurran, 2021), observed ASVs (DeSantis et al., 2006), Berger-Parker *d* (Berger & Parker,
895 1970), Gini (Gini, 1912), Shannon (Weaver, 1963), and Simpson (Simpson, 1949) (Figure 7 and Figure
896 19). Alpha-diversity indices suggested that the microbial richness of periodontally healthy controls was
897 higher than that of patients with periodontitis (Figure 7a-e and Figure 19). These results are in line with
898 findings with that patients with advanced periodontitis, namely stage II and stage III, have less diversified
899 communities than periodontally healthy controls (Jorth et al., 2014). Recognizing that the periodontitis
900 severity increases the amount of *Porphyromonas gingivalis*, the salivary microbiome compositions from
901 periodontally healthy controls conserved microbial networks dominated by *Streptococcus* spp. (Figure
902 13). *Porphyromonas gingivalis* is one of the known periodontal pathogen that could cause dysbiosys
903 in the salivary microbiomes, suggesting in the pathophysiology of periodontitis. Despite this finding,
904 earlier research found that subgingival microbiome of patients with periodontitis had a greater alpha-
905 diversity index (observed ASVs) than that of healthy controls (Iniesta et al., 2023), might due to the
906 different sampling sites between saliva and subgingival plaque. On the other hand, another research
907 has addressed significant discrepancies in alpha-diversity indices from subgingival plaque, saliva, and
908 tongue biofilms from healthy controls and periodontitis patients, resulting the highest alpha-diversity
909 index in saliva collections (Belstrøm et al., 2021). Moreover, early-stage periodontitis, namely stage I,

910 did not determine statistically significant differences in alpha-diversity indices compared to advanced
911 periodontitis, including stage II and stage III (Figure 7a-e). Accordingly, saliva collection of stage I
912 periodontitis may exhibit heterogeneity, indicating a midpoint condition between a healthy state and
913 advanced periodontitis (stage II and stage III). Likewise, gingivitis is often associated with low abundances
914 of the majority of periodontal pathogens, including *Porphyromonas gingivalis*, *Tannerella forsythia*, and
915 *Treponema denticola* (Abusleme et al., 2021). Compared to healthy controls, patients with stage I
916 periodontitis have higher detection rates of *Porphyromonas gingivalis* and *Tannerella forsythia* (Tanner et
917 al., 2006, 2007).

918 Therefore, we calculated beta-diversity indices to analyze the differences between the study partici-
919 pants. The distances for the multiple stages of periodontitis, including stage I, stage II, and stage III, as
920 well as healthy controls (Figure 4g-j and Table 7), suggesting notable differences among the multiple
921 periodontitis stages. In other words, the composition of the salivary microbiome compositions varies
922 depending on the periodontitis stages, so that supporting the findings from a previous study (Iniesta et al.,
923 2023). Taken together that it is nearly impossible to fully restore the attachment level after it has been lost
924 due to the progression and development of periodontitis, the ability to rapidly screen for periodontitis in
925 its early phases using saliva collections would be highly beneficial for effective disease management and
926 treatment.

927 Of the total of 425 taxa in the salivary microbiome composition that have been identified (Figure 13),
928 ANCOM was applied to select 20 taxa as the DAT that indicated notable abundance variation among
929 the periodontitis stages (Figure 8 and Table 5). Three sub-groups were formed from the DAT using
930 hierarchical clustering (Figure 8a). Surprisingly, two of the red complex pathogens (Rôcas, Siqueira Jr,
931 Santos, Coelho, & de Janeiro, 2001), *Porphyromonas gingivalis* and *Tannerella forsythia*, were classified
932 in Group 2 and were more prevalent in stage II and stage III periodontitis compared to healthy controls.
933 *Campylobacter showae* was additionally placed in Group 2 of the orange complex pathogens (Gambin et
934 al., 2021). Furthermore, some of the DAT in Group 2 have reported their crucial roles in pathogenesis
935 and development of periodontitis: *Filifactor alocis* (Aruni et al., 2015), *Treponema putidum* (Wyss et
936 al., 2004), *Tannerella forsythia* (Stafford, Roy, Honma, & Sharma, 2012; W. Zhu & Lee, 2016), and
937 *Prevotella intermedia* (Karched, Bhardwaj, Qudeimat, Al-Khabbaz, & Ellepolo, 2022). Taken together,
938 this indicates that DAT in Group 2 is essential to periodontitis. The portion of some Group 1 DAT,
939 including *Peptostreptococcaceae[XI][G-5] saphenum*, *Peptostreptococcaceae[XI][G-6] nodatum*, and
940 *Peptostreptococcaceae[XI][G-9] brachy*, in healthy controls and patients with periodontitis significantly
941 differed, according to earlier research (Lafaurie et al., 2022). These outcomes support our research,
942 implying that Group 1 DAT are also essential to the etiology and progression of periodontitis. However,
943 in contrast to patients with periodontitis, Group 3 DAT, namely *Corynebacterium durum* and *Actinomyces*
944 *graevenitzii*, were enriched in healthy controls, which is consistent with earlier research (Redanz et al.,
945 2021; Nibali et al., 2020).

946 In our correlation analysis (Figure 9), we have discovered strongly negative correlations (coefficient \leq
947 -0.5) between DAT of Group 3 and these of Group 1 and Group 2; we have also identified nine DAT
948 pairs with strong correlations (coefficient $\leq -0.5 \vee$ coefficient ≥ 0.5) (Figure 14). Interestingly, there

were strongly negative correlations (coefficient ≤ -0.5) between Group 2 DAT and *Actinomyces* spp., taxa which belong to Group 3: *Filifactor alocis* (Figure 14a), *Porphyromonas gingivalis* (Figure 14b), and *Treponema putidum* (Figure 14c). Taken together that pathogens, including *Filifactor alocis* (Aja, Mangar, Fletcher, & Mishra, 2021; Hiranmayi, Sirisha, Rao, & Sudhakar, 2017), *Porphyromonas gingivalis* (Rôças et al., 2001), and *Treponema putidum* (Wyss et al., 2004), become dominant taxa in patients with stage III periodontitis. On the other hand, commensal salivary bacteria, such as *Actinomyces* spp., gradually declined. Additionally, several DAT from Group 1 and Group 2 exhibited strong positive correlations (coefficient ≥ 0.5) (Figure 14d-i). It has been established that all of these DAT from Group 1 and Group 2 are periodontal pathogens: *Filifactor alocis* (Aja et al., 2021; Hiranmayi et al., 2017), *Fretibacterium* spp. (Teles, Wang, Hajishengallis, Hasturk, & Marchesan, 2021), *Lachnospiraceae[G-8] bacterium HMT 500* (Lafaurie et al., 2022), *Peptostreptococcaceae[XI][G-6] nodatum* (Lafaurie et al., 2022; Haffajee, Teles, & Socransky, 2006), *Peptostreptococcaceae[XI][G-9] brachy* (Lafaurie et al., 2022), and *Treponema putidum* (Wyss et al., 2004). Thus, these fundamental roles of identified periodontal pathogens in the pathophysiology and progression of periodontitis are further supported by these strong positive correlations (coefficient ≥ 0.5), suggesting that advanced periodontitis, i.e., stage III, might arise from the additional DAT from Group 1 and Group 2.

Moreover, to predict periodontitis stages from salivary microbiome composition, we have constructed machine-learning classification models based on random forest for four classification settings:

1. healthy vs. stage I vs. stage II vs. stage III
2. healthy vs. stage I
3. healthy vs. stage I vs. stages II/III
4. healthy vs. stages I/II/III

Porphyromonas gingivalis and *Actinomyces* spp. were the two most important taxa (feature) in all classification settings (Table 6). This finding aligns with a recent study that identifies *Actinomyces* spp. as the most prevalent bacteria in both the healthy gingivitis controls, while *Porphyromonas gingivalis* is recognized as the most predominant taxon within the periodontitis patients, based on analyses of subgingival plaque samples (Nemoto et al., 2021). We have previously developed machine learning models for the classification of periodontitis, with the objective of predicting the stages of chronic periodontitis by analyzing the copy numbers of nine known salivary bacteria species. We classified healthy controls and patients with periodontitis utilizing bacterial combinations in conjunction with a random forest model (E.-H. Kim et al., 2020):

- AUC: 94%
- BA: 84%
- SEN: 95%
- SPE: 72%

Another study established a machine-learning model for the classification of periodontitis, employing 266 species derived from the buccal microbiome (Na et al., 2020):

- AUC: 92%
- BA: 84%

- 988 • SEN: 94%
 989 • SPE: 74%
 990 By separating patients with periodontitis from healthy controls using only four DAT, *e.g.* *Actinomyces*
 991 *graevenitzii*, *Actinomyces* spp., *Corynebacterium durum*, and *Porphyromonas gingivalis*, our machine
 992 learning model performed better than previously published models (mean±SD) (Figure 10, Table 4, and
 993 Table 6):
 994 • AUC: 95.3%±4.9%
 995 • BA: 88.5%±6.6%
 996 • SEN: 86.4%±15.7%
 997 • SPE: 90.5%±7.0%
- 998 This result showed that by detecting Group 3 bacteria that were substantially abundant in health
 999 controls than patients with periodontitis, our study increased BA by at least 5% and SPE by at least 17%.
 1000 Furthermore, we have validated our machine-learning prediction model using openly accessible 16S
 1001 rRNA gene sequencing data from Portuguese (Iniesta et al., 2023) and Spanish participants (Relvas et
 1002 al., 2021) in order to ensure the consistency of our random forest classification model (Figure 11). Our
 1003 classification models employed in this study were primarily developed and assessed on Korean study par-
 1004 ticipants, which may limit their generalizability to other ethnic groups with different salivary microbiome
 1005 compositions (Premaraj et al., 2020; Renson et al., 2019). Therefore, the evaluations of this periodonti-
 1006 tis classification models can be affected by ethnic-specific variances and differences, highlighting the
 1007 necessity for additional validation and adjustment across a spectrum of ethnic backgrounds.
- 1008 Regarding the clinical characteristics and potential confounders influencing the analysis of salivary
 1009 microbiome compositions connected with periodontitis severity, this study had a number of limitations
 1010 that were pointed out. We did not offer clinical information, such as the percentage of teeth, the percentage
 1011 of bleeding on probing, nor dental furcation involvement, even though we did gather information on
 1012 attachment level, probing depth, plaque index, and gingival index (Renvert & Persson, 2002); this might
 1013 have it challenging to present thorough and in-depth data about periodontal health. Moreover, the broad age
 1014 range may make it tougher to evaluate the relationship between age and periodontitis statuses, providing
 1015 the necessity for future studies to consider into account more comprehensive clinical characteristics
 1016 associated with periodontitis. Additionally, potential confounders—*e.g.* body mass index (Bombin, Yan,
 1017 Bombin, Mosley, & Ferguson, 2022) and e-cigarette use (Suzuki, Nakano, Yoneda, Hirofumi, & Hanioka,
 1018 2022)—which might have affected dental health and salivary microbiome composition were disregarding
 1019 consideration in addition to smoking status and systemic diseases. Thus, future research incorporating
 1020 these components would offer a more thorough knowledge of how lifestyle factors interact and affect the
 1021 salivary microbiome composition and periodontal health. Throughout, resolving these limitations will
 1022 advance our understanding in pathogenesis and development of periodontitis, offering significant novel
 1023 insights on the causal connection between systemic diseases and the salivary microbiome compositions.

1024 **4 Metagenomic signature analysis of Korean colorectal cancer**

1025 **4.1 Introduction**

1026 Colorectal cancer (CRC) is one of the most prevalent and life-threatening malignancies worldwide
1027 (Kuipers et al., 2015; Center, Jemal, Smith, & Ward, 2009; N. Li et al., 2021), with its incidence
1028 influenced by a combination of genetic (Zhuang et al., 2021; Peltomaki, 2003), environmental (O'Sullivan
1029 et al., 2022; Raut et al., 2021), and lifestyle factors (X. Chen et al., 2021; Bai et al., 2022; Zhou et
1030 al., 2022; X. Chen, Li, Guo, Hoffmeister, & Brenner, 2022). Established risk factors include a often
1031 diet in red and processed meats (Kennedy, Alexander, Taillie, & Jaacks, 2024; Abu-Ghazaleh, Chua, &
1032 Gopalan, 2021), obesity (Mandic, Safizadeh, Niedermaier, Hoffmeister, & Brenner, 2023; Bardou et
1033 al., 2022), cigarette smoking (X. Chen et al., 2021; Bai et al., 2022), alcohol consumption (Zhou et al.,
1034 2022; X. Chen et al., 2022), and a sedentary lifestyle (S. An & Park, 2022), all of which contribute to
1035 chronic inflammation, mutagenesis, and metabolic regulation. Additionally, underlying conditions, e.g.
1036 Lynch syndrome (Vasen, Mecklin, Khan, & Lynch, 1991; Hampel et al., 2008) and familial adenomatous
1037 polyposis (Inra et al., 2015; Burt et al., 2004), significantly increase risk of CRC due to persistent mucosal
1038 inflammation and somatic mutations that promote tumorigenesis.

1039 The gut microbiome plays a fundamental role in maintaining host health by helping digestion
1040 (Joscelyn & Kasper, 2014; Cerqueira, Photenhauer, Pollet, Brown, & Koropatkin, 2020), regulating
1041 metabolism (Dabke, Hendrick, Devkota, et al., 2019; Utzschneider, Kratz, Damman, & Hullarg, 2016;
1042 Magnúsdóttir & Thiele, 2018), adjusting immune function (Kau, Ahern, Griffin, Goodman, & Gordon,
1043 2011; Shi, Li, Duan, & Niu, 2017; Broom & Kogut, 2018), and even coordinating neurological processes
1044 by the brain-gut axis (Martin et al., 2018; Aziz & Thompson, 1998; R. Li et al., 2024). Comprising
1045 these gut microbiota, including archaea, bacteria, fungi, and viruses, the gut microbiome contributes
1046 to the synthesis of essential vitamins, and production of fatty acids, which influence intestinal integrity
1047 and immune responses. Thus, well-balanced gut microbiome composition modulates systemic immune
1048 function by interacting with gut-associated lymphoid tissue, shaping immune tolerance and response
1049 to infections. Hence, emerging evidence suggests that dysbiosis in the gut microbiome composition are
1050 associated not only a narrow range of diseases, e.g. diarrhea and enteritis (Paganini & Zimmermann,
1051 2017; J. Gao, Yin, Xu, Li, & Yin, 2019) but also a wide range of diseases, e.g. obesity, diabetes, and
1052 cancers (Barlow et al., 2015; Hartstra et al., 2015; Helmink et al., 2019; Cullin et al., 2021).

1053 Recent studies have highlighted the crucial role of the gut microbiome in tumorigenesis and progres-
1054 sion of CRC (Song, Chan, & Sun, 2020; Rebersek, 2021), with dysbiosis emerging as a potential risk
1055 factor. Dysbiosis in gut microbiome compositions can promote tumorigenesis of many cancers, including
1056 CRC, through several signaling cascades, including inflammation, mutagenesis, and altered metabolism
1057 in host. Certain bacteria species, such as *Fusobacterium* genus (Hashemi Goradel et al., 2019; Bullman et
1058 al., 2017; Flanagan et al., 2014), *Bacteroides* genus (Ulger Toprak et al., 2006; Boleij et al., 2015), and
1059 *Escherichia coli* (Swidsinski et al., 1998; Bonnet et al., 2014), have been associated with development
1060 and progression of CRC by producing pro-inflammatory signals, generating toxins including mutagens,

1061 and disrupting the intestinal barriers including mucous surface. In contrast, beneficial bacteria, such as
1062 *Lactobacillus* genus (Ghorbani et al., 2022; Ghanavati et al., 2020) and *Bifidobacterium* genus (Le Leu,
1063 Hu, Brown, Woodman, & Young, 2010; Fahmy et al., 2019), are regarded to apply protective roles by
1064 maintaining homeostasis of gut microbiome compositions and regulating immune responses including
1065 inflammation.

1066 Furthermore, identifying metagenome biomarkers in Korean CRC patients is essential, as the gut
1067 microbiome compositions significantly vary by ethnicity due to genetic, dietary, and environmental
1068 factor (Fortenberry, 2013; Merrill & Mangano, 2023; Parizadeh & Arrieta, 2023). Additionally, ethnicity-
1069 specific microbiome composition signatures may affect the reliability of previously established biomarkers
1070 derived from predominantly Western CRC cohorts (Network et al., 2012), necessitating population-
1071 specific investigations. By identifying metagenomic biomarkers tailored to Korean CRC patients, we
1072 can improve early detection rate of early-stage CRC, develop more accurate risk of CRC, and explore
1073 microbiome-targeted therapies that consider host-microbiome interactions within the Korean population.

1074 Accordingly, this study aims to identify microbiome-based biomarkers specific to CRC within
1075 the Korean population, addressing the critical demand for ethnicity-specific microbiome research. By
1076 leveraging metagenomic sequencing and advanced computational biology analysis, this study seeks to
1077 uncover novel microbial signatures associated with Korean CRC patients. As part of the larger "Multi-
1078 genomic analysis for biomarker development in colon cancer" project (NTIS No. 1711055951), this study
1079 investigates microbial signatures within next-generation sequencing data to enhance precision medicine
1080 approaches for CRC and to develop robust microbiome-based biomarkers for early detection, prognosis,
1081 and therapeutic stratification, complementing genomic and epigenomic markers. Hence, this research
1082 represents a crucial step toward personalized cancer diagnostic and therapeutic strategies tailored to the
1083 Korean population.

1084 **4.2 Materials and methods**

1085 **4.2.1 Study participants enrollment**

1086 To achieve metagenomic observations of CRC, a total of 211 Korean CRC patients were enrolled (Table
1087 8). The tissue samples were collected from both the tumor lesion and its corresponding adjacent normal
1088 lesion to enable comparative metagenomic analyses. Tumor tissue samples were obtained from confirmed
1089 CRC lesions, ensuring adequate representation of CRC-associated microbial alterations. Adjacent normal
1090 tissues were collected from non-cancerous regions away from the tumor margin to serve as a control
1091 for baseline molecular and microbial composition. Moreover, clinical information was collected for all
1092 study participants included in this study to investigate potential associations between gut microbiome
1093 compositions and clinical outcomes. Key clinical characteristics recorded included overall survival (OS)
1094 and recurrence. These clinical parameters were integrated with metagenomic data to explore potential
1095 microbiome-based biomarkers for CRC prognosis and progression. Ethical approval was obtained for
1096 clinical data collection, and all patient information was anonymized to ensure confidentiality in accordance
1097 with institutional guidelines.

1098 **4.2.2 DNA extraction procedure**

1099 Tissue samples were immediately processed under sterile conditions to prevent contamination and
1100 preserved in low temperature (-80°C) storage for downstream DNA extraction and whole-genome
1101 sequencing. Furthermore, produced sequencing data were provided by the "Multi-genomic analysis
1102 for biomarker development in colon cancer" project (NTIS No. 1711055951) in mapped BAM format,
1103 aligned to the hg38 human reference genome. The preprocessing pipeline utilized by the main project
1104 included high-throughput whole-genome sequencing using standardized alignment algorithm, BWA
1105 (H. Li & Durbin, 2009). In addition to the mapped human sequences, our whole-genome sequencing
1106 data retained unmapped sequences, which contain potential microbial reads that were not aligned to the
1107 human reference genome.

1108 **4.2.3 Bioinformatics analysis**

1109 To identify microbial signatures associated with CRC, we employed PathSeq (version 4.1.8.1) (Kostic
1110 et al., 2011; Walker et al., 2018), a computational pipeline designed for metagenomic analysis of high-
1111 throughput sequencing data including the whole-genome sequences. After processing these sequencing
1112 data through the PathSeq pipeline, a comprehensive bioinformatics analyses were conducted to characterize
1113 microbial signatures associated with CRC.

1114 Prevalent taxa identification was performed by determining microbial taxa present in the majority of
1115 the study participants, filtering out low-abundance and rare taxa to ensure robust downstream analyses.

1116 To assess microbial community structure, diversity indices were calculated, including alpha-diversity
1117 to evaluate single-sample diversity and beta-diversity to compare microbial composition between the
1118 tumor tissues and their corresponding adjacent normal tissues. Following alpha-diversity indices were

1119 calculated using the scikit-bio Python package (version 0.6.3) (Rideout et al., 2018), and these alpha-
1120 diversity indices were compared using the MWU test:

- 1121 1. Berger-Parker d (Berger & Parker, 1970)
- 1122 2. Chao1 (Chao, 1984)
- 1123 3. Dominance
- 1124 4. Doubles
- 1125 5. Fisher (Fisher et al., 1943)
- 1126 6. Good's coverage (Good, 1953)
- 1127 7. Margalef (Magurran, 2021)
- 1128 8. Mcintosh e (Heip, 1974)
- 1129 9. Observed ASVs (DeSantis et al., 2006)
- 1130 10. Simpson d
- 1131 11. Singles
- 1132 12. Strong (Strong, 2002)

1133 Furthermore, these beta-diversity indices were measured and compared using the PERMANOVA
1134 test (Anderson, 2014; Kelly et al., 2015). To demonstrate multi-dimensional data from the beta-diversity
1135 indices, we utilized the t-SNE algorithm (Van der Maaten & Hinton, 2008).

- 1136 1. Bray-Curtis (Sorensen, 1948)
- 1137 2. Canberra
- 1138 3. Cosine (Ochiai, 1957)
- 1139 4. Hamming (Hamming, 1950)
- 1140 5. Jaccard (Jaccard, 1908)
- 1141 6. Sokal-Sneath (Sokal & Sneath, 1963)

1142 Differentially abundant taxa (DAT) were identified using statistical method, ANCOM (Lin & Peddada,
1143 2020), adjusting for sequencing depth and potential confounders to highlight taxa significantly associated
1144 with categorical clinical information in CRC, such as recurrence. Furthermore, to point attention to
1145 taxa that are substantially linked to continuous clinical measurement in CRC, including OS, DAT were
1146 found using the Spearman correlation and slope from linear regression (Equation 9). Note that both the
1147 Spearman correlation and the slope from linear regression were utilized to provide a more comprehensive
1148 assessment of the relationship between DAT proportions and OS. While the correlation coefficient
1149 measures the strength and direction of a linear relationship between these variables, it does not convey
1150 information about the magnitude of change in independent variable relative to dependent variable. The
1151 slope of the linear regression model, on the other hand, quantifies this change by indicating how much
1152 the dependent variable is expected to increase or decrease per unit change in the independent variable. By
1153 incorporating both the correlation coefficient and the slope from the linear regression, we ensured that
1154 the analysis captured not only whether two variables were associated but also the extent to which one
1155 variable influenced the other. This dual approach enhances the interpretability of results, particularly in
1156 biological and clinical studies where both statistical association and biological effect size are crucial for
1157 meaningful suggestions.

$$\text{slope} = \frac{\Delta \text{OS}}{\Delta \text{DAT proportion}} \quad (9)$$

1158 To assess the predictive potential of microbial signatures in CRC prognosis, we employed a random
 1159 forest machine learning model using DAT proportions as input features. Random forest classification was
 1160 utilized to predict CRC recurrence, where the classification model was trained to distinguish between
 1161 CRC patients with or without recurrence based on the gut microbiome compositions. Additionally,
 1162 random forest regression was applied to predict OS by estimating survival time as a continuous clinical
 1163 outcome based on microbiome features. This approach allowed for the identification of microbial taxa
 1164 that contribute significantly to CRC prognosis, offering insights into potential gut microbiome-based
 1165 biomarkers for cancer progression. By integrating these random forest machine learning models, we
 1166 aimed to improve CRC risk stratification and precision medicine strategies.

1167 This multi-layered bioinformatics approach enabled a comprehensive investigation of gut microbiome
 1168 alteration in CRC, facilitating the identification of potential microbial biomarkers for diagnosis and
 1169 prognosis of CRC.

1170 **4.2.4 Data and code availability**

1171 All sequences from the 211 study participants have been published to the Korea Bioinformation Center
 1172 (data ID KGD10008857): <https://kbds.re.kr/KGD10008857>. Docker image that employed through-
 1173 out this study is available in the DockerHub: <https://hub.docker.com/repository/docker/fumire/unist-crc-copm/general>. Every code used in this study can be found on GitHub: <https://github.com/CompbioLabUnist/CoPM-ColonCancer>.

1176 **4.3 Results**

1177 **4.3.1 Summary of clinical characteristics**

1178 Microsatellite instability (MSI) is one of the key molecular features and risk factors in CRC, resulting
1179 from defects in the DNA mismatch repair system (Boland & Goel, 2010). MSI leads to the accumulation
1180 of mutations in short repetitive DNA sequences (microsatellites), contributing to genomic instability and
1181 tumor development (Søreide, Janssen, Söiland, Körner, & Baak, 2006; Vilar & Gruber, 2010). Therefore,
1182 we compared clinical measurements with MSI status, including microsatellite stable (MSS), MSI-low
1183 (MSI-L), and MSI-high (MSI-H). There were no significant differences in the clinical measurements, *e.g.*
1184 recurrence, sex, OS, and age in diagnosis, in the total of 211 study participants (Table 8).

1185 **4.3.2 Gut microbiome compositions**

1186 In the total of 211 CRC study participants, these ten kingdoms were found in the gut microbiome
1187 composition:

- 1188 1. Archaea kingdom: 31 genera
- 1189 2. Bacteria kingdom: 1508 genera
- 1190 3. Bamfordvirae kingdom: 1 genus
- 1191 4. Eukaryota kingdom: 77 genera
- 1192 5. Fungi kingdom: 137 genera
- 1193 6. Loebvirae kingdom: 2 genera
- 1194 7. Orthornavirae kingdom: 1 genus
- 1195 8. Parnavirae kingdom: 3 genera
- 1196 9. Shotokuvirae kingdom: 6 genera
- 1197 10. Viruses kingdom: 76 genera

1198 Among these kingdoms, the proportions of four major kingdoms, which have at least 50 genera, in
1199 the gut microbiome composition were displayed (Figure 21): bacteria kingdom, eukaryota kingdom,
1200 fungi kingdom, and viruses kingdom. In the bacteria kingdom (Figure 21a), *Bacteroides* genus is the
1201 most prevalent genus in the tumor tissue samples, followed by *Fusobacterium* and *Cutibacterium* genera.
1202 *Toxoplasma* and *Malassezia* genera were the dominant genus, which have over 90% of proportions, in the
1203 eukaryota kingdom (Figure 21b) and the fungi kingdom (Figure 21c), respectively. On the other hand,
1204 *Roseolovirus* genus is the most popular genus of the viruses kingdom in the normal tissue samples (Figure
1205 21d); contrarily, *Lymphocryptovirus* and *Cytomegalovirus* genera had been dominant genera in the tumor
1206 tissue samples. Taken together, these results suggest that the Anna Karenina principle (Ma, 2020; W. Li
1207 & Yang, 2025), *i.e.* in human microbiome-associated diseases, every disease-associated microbiome,
1208 including dysbiosis, is unique and patient-specific, whereas all healthy microbiomes are similar, also
1209 applies to CRC.

1210 **4.3.3 Diversity indices**

1211 In alpha-diversity analysis, which measures within-sample microbial community, revealed a significant
1212 increase in tumor samples compared to adjacent normal samples (Figure 22). Alpha-diversity indices,
1213 including Chao1, Fisher α , and observed features, were consistently higher in CRC tumor tissues (MWU
1214 test $p < 0.05$), indicating a more heterogeneous microbial community, *e.g.* the Anna Karenina principle,
1215 potentially influenced by tumor-associated dysbiosis.

1216 To assess the microbial impact on CRC recurrence, alpha-diversity indices compared between normal
1217 and tumor tissue samples in accordance with recurrence information (Figure 23). In the recurrence patients,
1218 most alpha-diversity indices (11/12 indices; 92% indices), except McIntosh index, exhibited increasing in
1219 tumor samples than normal samples (MWU test $p < 0.05$; Figure 23); In the non-recurrence patients, on
1220 the other hand, some alpha-diversity indices (8/12 indices; 67% indices) amplified in tumor samples than
1221 normal samples (MWU test $p < 0.05$; Figure 23). What is interesting about the alpha-diversity analysis
1222 in this figure is that a few indices, namely Fisher α (Figure 28e) and Margalef (Figure 23g), presented
1223 augmentation in normal sample of the recurrence patients than that of the non-recurrence patients (MWU
1224 test $p < 0.05$). Overall, these alpha-diversity results demonstrate that tumor samples have more diverse
1225 microbiome composition than normal samples. Furthermore, although only two indices significantly
1226 increased, the recurrence patients have diversified microbiome compositions than the non-recurrence
1227 patients in normal samples, not in tumor samples, indicating field cancerization by the gut microbiome
1228 leads to unfavorable prognosis such as recurrence (Curtius et al., 2018; Rubio et al., 2022).

1229 To determine the microbial impact on OS of CRC patients, the Spearman correlation compared
1230 between alpha-diversity indices and OS duration (Figure 24). No significant Spearman correlation was
1231 found between every alpha-diversity indices and OS (Spearman correlation $p \geq 0.1$; Figure 24). However,
1232 a few alpha-diversity indices, *e.g.* Chao1 (Figure 24b), Good's coverage (Figure 24f), and observed
1233 features (Figure 24i), showed negative correlations with OS (Spearman correlation $p < 0.05$). Together
1234 these correlation results provide important insights into heterogeneous microbiome leads to shorter OS,
1235 suggesting the Anna Karenina principle and the field cancerization.

1236 In beta-diversity analysis, which calculates inter-sample microbial community, explain significant
1237 disparity between tumor samples and normal samples (Figure 25). Every six beta-diversity indices
1238 presented discrepancy between normal samples and tumor samples (PERMANOVA test $p < 0.001$),
1239 implying that tumor samples have distinct microbiome compositions from normal tissue samples.

1240 Beta-diversity indices were evaluated between normal and tumor tissue samples along with recurrence
1241 history in order to evaluate the microbial influence on CRC recurrence (Figure 26). All six beta-diversity
1242 indices examined significant difference in microbial community structure between the recurrence patients
1243 and the non-recurrence patients (PERMANOVA test $p < 0.001$; Figure 26), indicating that tumor-
1244 associated gut microbiome composition varies resulting on recurrence status. tSNE-transformed plots
1245 further illustrated clear clustering patterns (Figure 26), suggesting again that the recurrence patients
1246 harbor dissimilar microbial communities compared to the non-recurrence patients. These observed
1247 differences in beta-diversity represent that microbial shifts, including dysbiosis, may be associated with

1248 CRC progression and recurrence risk, possibly due to specific taxa contributing to a tumor-promoting
1249 microenvironment.

1250 Moreover, beta-diversity analysis suggested a potential associated with OS duration in CRC patients.
1251 In all six beta-diversity indices, tSNE-transformed projection plots showed clear clustering patterns
1252 along OS duration (Figure 27), implying that possible microbiome composition shifts related to survival
1253 outcomes in CRC. However, since OS is a continuous variable, statistical significance testing could
1254 not be directly performed for these clustering patterns. Despite this limitation, the observed microbial
1255 community variations suggest that alterations in the gut microbiome composition may be associated to
1256 CRC prognosis and survival duration.

1257 Together, diversity indices analyses revealed significant microbial community alterations between
1258 normal and tumor tissue samples, as well as between the recurrence and non-recurrence CRC patients.
1259 Alpha-diversity indices significantly increased in tumor tissue samples than normal tissue samples (MWU
1260 test $p < 0.05$; Figure 22). This increase was more pronounced in the recurrence patients (11/12 indices;
1261 92% indices) compared to non-recurrence patients (8/12 indices; 67% indices) (Figure 23), indicating a
1262 potential link between microbial diversity and CRC recurrence. Additionally, negative correlation between
1263 OS and alpha-diversity indices were observed in normal samples (Spearman correlation $p < 0.05$; Figure
1264 24), suggesting that lower microbial diversity may be associated with longer survival in CRC. On the
1265 other hand, beta-diversity indices analysis, showed significant separation between tumor and tumor tissue
1266 samples across all six beta-diversity indices (PERMANOVA test $p < 0.001$; Figure 25). Furthermore,
1267 the recurrence and non-recurrence patients displayed significantly discrete microbial compositions
1268 (PERMANOVA test $p < 0.001$; Figure 26), implying that microbial community shifts may reflect CRC
1269 progression and recurrence risk. These findings highlight the importance of microbiome diversity and
1270 gut microbiome composition in CRC prognosis and warrant further investigation into their potential as
1271 predictive biomarkers.

1272 4.3.4 DAT selection

1273 The selection of differentially abundant taxa (DAT) aimed to identify microbial taxa that exhibit significant
1274 differences in relative abundance between clinical information, such as recurrence history or OS in CRC
1275 patients. Identifying and selection these microbial discrepancies is crucial for understanding the role of
1276 the gut microbiome composition in CRC progression, prognosis, and potential therapeutic interventions.

1277 We identified 19 DAT associated with recurrence history across the total samples by ANCOM
1278 (Figure 28a), including 18 non-recurrence-enriched DAT and a recurrence-enriched DAT. When stratified
1279 by sample type, one DAT was enriched in normal samples of the non-recurrence patients (Figure
1280 28b), whereas six DAT exhibited significant differential abundance in tumor samples (Figure 28c).
1281 These findings suggest that microbial composition variations in the tumor microenvironment are more
1282 pronounced in relation to recurrence status (Table 9), potentially indicating a microbial signature linked
1283 to CRC progression. These identified DAT may contribute to tumor-associated dysbiosis, influencing
1284 the likelihood of CRC recurrence through mechanisms such as inflammation, metabolic modulation, or

1285 immune system interaction.

1286 The non-recurrence-enriched DAT have decreased proportions both in normal and tumor samples of
1287 the recurrence patients than those in the non-recurrence patients (MWU test $p < 0.001$; Figure 28d-h).
1288 What is interesting about these non-recurrence-enriched DAT is that they belong to the *Micrococcus* genus.
1289 Among them, *Micrococcus aloeverae* was consistently identified in all three settings—total (Figure 28a),
1290 normal (Figure 28b), and tumor samples (Figure 28c)—indicating its stable presence regardless of tissue
1291 type. Variation in relative proportions of *Micrococcus aloeverae* (Figure 28d) suggests potential ecological
1292 adaptability within tumor microenvironment of CRC. The remaining *Micrococcus* genus DAT showed
1293 less variation between the recurrence and non-recurrence patients, reinforcing their limited associations
1294 with CRC recurrence. Moreover, only one taxon, *Pseudomonas* sp. *NBRC 111133*, was identified as
1295 recurrence-enriched DAT (Figure 28a). This suggests a potential association between *Pseudomonas* sp.
1296 *NBRC 111133* and CRC recurrence, indicating that its presence may contribute to a tumor-supportive
1297 microbial environment. *Pseudomonas* sp. *NBRC 111133* had higher relative proportions both in normal
1298 and tumor tissue samples of the recurrence patients than those of the non-recurrence patients (Figure 28i).
1299 Likewise, *Pseudomonas* sp. *NBRC 111133* were prevalent in tumor samples than normal samples of the
1300 non-recurrence patients (MWU test $p < 0.01$; Figure 28i); however, no significant difference between
1301 normal and tumor tissue samples of the recurrence patients.

1302 These findings imply that while certain species belong to *Micrococcus* genus may be prevalent in
1303 CRC tumor tissues, their roles in cancer progression and recurrence risk remain uncertain. Species of
1304 *Pseudomonas* genus are known for their metabolic involvement in biofilm formation, antibiotic resistance,
1305 and immune modulation, which could play an essential role in CRC progression.

1306 Furthermore, correlation analysis between DAT abundance and OS duration identified a total of 16
1307 over-represented DAT in the total samples (Figure 29a). When analyzed separately, 11 OS-correlated DAT,
1308 which consist of four under-represented and seven over-represented DAT showed significant correlations
1309 with OS in normal samples (Figure 29b), while four under-represented and 45 over-represented DAT
1310 were identified in tumor samples (Figure 29c), indicating that microbial composition shifts in tumor
1311 tissues may have a stronger association with survival outcomes. The higher number of survival-associated
1312 DAT in tumor tissue suggests that the tumor microbiome plays a more dynamic role in progression and
1313 prognosis of CRC. These findings highlight the potential of gut microbial composition as a prognostic
1314 indicator in CRC, warranting further investigation into the functional roles of these DAT in influencing
1315 clinical outcomes.

1316 Among a total of 57 OS-correlated DAT (Table 10) with Spearman correlation and the slope (Equation
1317 9). *Agaricus bisporus* (Figure 29d) and *Corynebacterium* sp. *KPL1824* (Figure 29h) are identified as
1318 over-represented DAT both in normal samples and tumor samples (Spearman correlation $p < 0.05$),
1319 whereas *Corynebacterium lowii* (Figure 29g) and *Paracoccus sphaerophysae* (Figure 29i) are selected
1320 as under-represented DAT both in normal samples and tumor samples (Spearman correlation $p < 0.05$).
1321 On the other hand, *Clostridiales bacterium* (Figure 29e) is classified as under-represented DAT only in
1322 normal samples (Spearman correlation $p < 0.01$), while *Corynebacterium kroppenstedtii* (Figure 29f) is
1323 described as over-represented DAT only in tumor samples (Spearman correlation $p < 0.001$).

1324 These findings highlight the potential influence of microbial dysbiosis on cancer progression and
1325 prognosis. The presence of these OS-correlated DAT in tumor and/or adjacent normal tissues suggests
1326 that microbial alterations may contribute to field cancerization, a phenomenon where histopathologically
1327 benign tissues surrounding the tumor undergo molecular, inflammatory, and microbial shifts, creating
1328 a microenvironment conducive to tumor development and progression. Therefore, these discoveries
1329 reinforce the importance of investigating the gut microbiome as a prognostic biomarker and suggest that
1330 targeting microbial dysbiosis could offer new therapeutic strategies for improving clinical outcomes and
1331 treatment responses of CRC.

1332 **4.3.5 Random forest prediction**

1333 We employed the random forest-based machine learning prediction to assess the predictive power of DAT
1334 from gut microbiome composition for CRC prognosis. To achieve this aim, we utilized random forest
1335 classification to predict recurrence status, training the model to differentiate between recurrence and
1336 non-recurrence patients based on microbial abundance patterns. Additionally, we applied random forest
1337 regression to predict OS, aiming to identify microbial taxa associated with survival duration. By leveraging
1338 random forest models, this study aimed to establish a microbiome-based predictive machine learning
1339 models for CRC recurrence risk assessment and survival prognosis, contributing to the development of
1340 prediction medicine strategies based on gut microbial signatures.

1341 To evaluate the predictive power of gut microbiome composition in CRC recurrence, we implemented
1342 a random forest classification model using two different input sets (Figure 30a-f): the entire gut micro-
1343 biome composition and DAT-selected microbiome. Comparing these models allowed us to assess whether
1344 focusing on DAT-selected microbial features enhances classification performance. While the DAT-based
1345 classification models showed slightly improved classification metrics (MWU test $p \geq 0.05$), including
1346 ACC, AUC, and BA, over the entire microbiome-based model in the total sample (Figure 30a and Figure
1347 29d), normal samples (Figure 30b and Figure 30e), and tumor samples (Figure 30c and Figure 30f),
1348 overall classification metrics remained around 60% (0.570 ± 0.164 , mean \pm SD), suggesting moderated
1349 predictive capability. This relatively low metrics highlight the complexity of CRC recurrence, indicating
1350 that while dysbiosis may contribute to CRC progression, it is likely interwinded with host genetic factors
1351 such as germline and somatic mutations. Thus, the interplay between microbial shifts and tumor genomic
1352 alterations warrants further investigation, as integrating microbiome and genomic sequencing data may
1353 improve therapeutic strategies.

1354 To assess the predictive capability of the gut microbiome composition in OS of CRC patients, we
1355 implemented a random forest regression model, comparing two different input sets (Figure 30g-i): the
1356 entire gut microbiome composition and DAT-selected microbiome. This comparison also aimed to
1357 determine whether focusing on key microbial features (DAT) enhances predictive accuracy. While DAT-
1358 based model showed a slight improvement over the entire microbiome-based model in normal samples
1359 (Figure 30h) and tumor samples (Figure 30i), the regression error remained high (729.302 ± 179.940 ,
1360 mean \pm SD), indicating substantial variability in survival outcomes that cannot be fully explained by gut

1361 microbiome composition alone. This result suggest that while gut microbial dysbiosis may influence
1362 CRC progression, survival duration (OS) is likely also driven by host genetic factors, highlighting the
1363 requirement for multi-omics integration, where combining microbiome and genomic sequencing data
1364 may provide a more accurate and comprehensive predictive model for CRC patients survival.

Table 8: Clinical characteristics of CRC study participants.

Continuous variable: mean \pm SD. Categorical variable: count (proportion). Statistical significance were assessed using the χ^2 -squared test for categorical values and the Kruskal-Wallis test for continuous values. OS: overall survival.

	Overall	MSS	MSI-L	MSI-H	p-value
n	211	181	7	18	
Recurrence, n (%)	132 (62.6%)	112 (61.9%)	4 (57.1%)	13 (72.2%)	0.657
True	79 (37.4%)	69 (38.1%)	3 (42.9%)	5 (27.8%)	
Sex, n (%)					
Male	137 (64.9%)	119 (65.7%)	6 (85.7%)	10 (55.6%)	0.357
Female	74 (35.1%)	62 (34.3%)	1 (14.3%)	8 (44.4%)	
OS, mean \pm SD	1248.5 \pm 770.3	1268.1 \pm 769.5	1416.6 \pm 496.3	1097.7 \pm 903.2	0.580
Age, mean \pm SD	61.2 \pm 13.1	61.7 \pm 12.4	60.1 \pm 15.6	60.2 \pm 19.4	0.867

Table 9: DAT list for CRC recurrence.

Statistical significance was determined by ANCOM W. Significance threshold is $|\log_2 \text{FC}| > 1.0| \wedge W > 9600$. Non-significant values remain blank. DAT are sorted in alphabetical order. FC: fold change

Taxonomy name	Entire-log ₂ FC	Entire-W	Normal-log ₂ FC	Normal-W	Tumor-log ₂ FC	Tumor-W
<i>Cutibacterium acnes</i>	-1.878	10570				
<i>Cutibacterium avidum</i>	-1.383	10266				
<i>Cutibacterium granulosum</i>	-1.476	10271				
<i>Micrococcus aloeverae</i>	-2.280	10740	-1.821	10462	-2.481	10591
<i>Micrococcus luteus</i>	-2.216	10744				
<i>Micrococcus</i> sp. <i>CH3</i>	-2.323	10740			-2.493	10527
<i>Micrococcus</i> sp. <i>CH7</i>	-2.321	10740			-2.493	10542
<i>Micrococcus</i> sp. <i>HMSC31B01</i>	-2.282	10739			-2.458	10519
<i>Micrococcus</i> sp. <i>MS-ASIII-49</i>	-2.284	10740			-2.470	10527
<i>Pseudomonas</i> sp. <i>NBRC 111133</i>	1.139	9732				
<i>Pseudonocardia</i> sp. <i>P2</i>	-2.200	10736			-2.394	10253
<i>Staphylococcus</i> sp. <i>HMSC034A07</i>	-1.341	10050				
<i>Staphylococcus</i> sp. <i>HMSC063F03</i>	-1.322	10001				
<i>Staphylococcus</i> sp. <i>HMSC064E11</i>	-1.064	10163				
<i>Staphylococcus</i> sp. <i>HMSC067B04</i>	-1.343	9952				
<i>Staphylococcus</i> sp. <i>HMSC068G12</i>	-1.344	10173				
<i>Staphylococcus</i> sp. <i>HMSC072H01</i>	-1.298	10197				
<i>Staphylococcus</i> sp. <i>HMSC077C03</i>	-1.331	10115				
<i>Treponema endosymbiont of Eucomomymptha</i> sp.	-1.629	10472				

Table 10: DAT list for CRC OS.

Significance threshold is $\log_{10}|\text{slope}| > 2.0 \wedge |r| > 0.2$. Non-significant values remain blank. DAT are sorted in alphabetical order.

Taxonomy name	Entire-slope	Entire-r	Normal-slope	Normal-r	Tumor-slope	Tumor-r
<i>Acinetobacter venetianus</i>					3.087	0.203
<i>Actinotalea ferrariae</i>					2.574	0.200
<i>Agaricus bisporus</i>	2.329	0.287	2.925	0.276	2.258	0.306
<i>Bifidobacterium boum</i>					2.096	-0.216
<i>Brevundimonas</i> sp. <i>DS20</i>			2.180	0.279		
<i>Clostridiales bacterium</i>			2.631	-0.203		
<i>Corynebacterium kroppenstedtii</i>	2.117	0.220			2.117	0.302
<i>Corynebacterium lipophiloflavum</i>			2.137	0.227		
<i>Corynebacterium lowii</i>			2.006	-0.216		
<i>Corynebacterium</i> sp. <i>KPL1818</i>	2.101	0.209	2.487	0.220	2.044	0.215
<i>Corynebacterium</i> sp. <i>KPL1824</i>	2.057	0.207	2.511	0.212	2.003	0.226
<i>Corynebacterium</i> sp. <i>KPL1986</i>					2.205	0.202
<i>Corynebacterium</i> sp. <i>KPL1996</i>					2.205	0.202
<i>Corynebacterium</i> sp. <i>KPL1998</i>					2.205	0.202
<i>Corynebacterium</i> sp. <i>KPL2004</i>					2.205	0.202
<i>Kocuria flava</i>			2.729	0.214		
<i>Kytococcus sedentarius</i>					2.267	0.206
<i>Lachnospiraceae bacterium AD3010</i>			2.609	-0.203		
<i>Lachnospiraceae bacterium NK4A136</i>					2.538	-0.220
<i>Methylorum extorquens</i>					2.068	0.295
<i>Microbacterium barkeri</i>			2.071	0.389		
<i>Paracoccus sphaerophysae</i>					2.012	-0.209
<i>Pontibacillus litoralis</i>					2.580	-0.209
<i>Porphyromonas macacae</i>			2.476	-0.200		
<i>Pseudomonas balearica</i>					2.117	0.203
<i>Pseudomonas monteilii</i>					2.183	0.228
<i>Rodentibacter myodis</i>					2.444	0.245
<i>Roseovarius tolerans</i>					2.295	0.221
<i>Staphylococcus epidermidis</i>					2.243	0.214
<i>Staphylococcus</i> sp. <i>HMSC034A07</i>					2.183	0.209
<i>Staphylococcus</i> sp. <i>HMSC034D07</i>	2.278	0.206			2.252	0.253
<i>Staphylococcus</i> sp. <i>HMSC034G11</i>	2.362	0.208			2.357	0.261
<i>Staphylococcus</i> sp. <i>HMSC036A09</i>					2.308	0.239
<i>Staphylococcus</i> sp. <i>HMSC055A10</i>					2.168	0.222
<i>Staphylococcus</i> sp. <i>HMSC055B03</i>	2.134	0.202			2.134	0.266
<i>Staphylococcus</i> sp. <i>HMSC058E12</i>					2.106	0.216
<i>Staphylococcus</i> sp. <i>HMSC061C10</i>					2.882	0.207
<i>Staphylococcus</i> sp. <i>HMSC062B11</i>	2.391	0.203			2.377	0.253
<i>Staphylococcus</i> sp. <i>HMSC062D04</i>	2.278	0.202			2.274	0.259
<i>Staphylococcus</i> sp. <i>HMSC063F03</i>	2.376	0.201			2.367	0.251
<i>Staphylococcus</i> sp. <i>HMSC063F05</i>	2.387	0.210			2.381	0.266
<i>Staphylococcus</i> sp. <i>HMSC064E11</i>					2.276	0.218
<i>Staphylococcus</i> sp. <i>HMSC065D11</i>					2.329	0.245

Table 10 continued from previous page

Taxonomy name	Entire-slope	Entire-r	Normal-slope	Normal-r	Tumor-slope	Tumor-r
<i>Staphylococcus</i> sp. <i>HMSC066G04</i>					2.181	0.218
<i>Staphylococcus</i> sp. <i>HMSC067B04</i>	2.332	0.205			2.329	0.260
<i>Staphylococcus</i> sp. <i>HMSC068G12</i>					2.294	0.226
<i>Staphylococcus</i> sp. <i>HMSC070A07</i>	2.360	0.216			2.362	0.287
<i>Staphylococcus</i> sp. <i>HMSC073C02</i>	2.352	0.205			2.334	0.246
<i>Staphylococcus</i> sp. <i>HMSC073E10</i>					2.366	0.255
<i>Staphylococcus</i> sp. <i>HMSC074D07</i>	2.330	0.218			2.308	0.270
<i>Staphylococcus</i> sp. <i>HMSC076H12</i>					2.200	0.219
<i>Staphylococcus</i> sp. <i>HMSC077C03</i>					2.258	0.207
<i>Staphylococcus</i> sp. <i>HMSC077D09</i>					2.245	0.230
<i>Staphylococcus</i> sp. <i>HMSC077G12</i>	2.335	0.200			2.345	0.276
<i>Staphylococcus</i> sp. <i>HMSC077H01</i>					2.214	0.241
<i>Streptomyces cinnamoneus</i>					2.787	0.208
<i>Thauera terpenica</i>					2.975	0.226

Table 11: Random forest classification and their evaluations.

Metrics are shown as mean \pm SD.

	Dataset	ACC	AUC	BA	F1	PRE	SEN	SPE
Entire	Total	0.544 \pm 0.139	0.667 \pm 0.141	0.561 \pm 0.141	0.544 \pm 0.139	0.559 \pm 0.152	0.562 \pm 0.192	0.559 \pm 0.152
	Normal	0.464 \pm 0.214	0.571 \pm 0.182	0.484 \pm 0.210	0.464 \pm 0.214	0.515 \pm 0.200	0.454 \pm 0.255	0.515 \pm 0.200
	Tumor	0.481 \pm 0.176	0.615 \pm 0.087	0.497 \pm 0.181	0.481 \pm 0.176	0.464 \pm 0.189	0.530 \pm 0.212	0.464 \pm 0.189
DAT	Total	0.582 \pm 0.112	0.656 \pm 0.109	0.592 \pm 0.120	0.582 \pm 0.112	0.558 \pm 0.114	0.626 \pm 0.167	0.558 \pm 0.114
	Normal	0.530 \pm 0.117	0.567 \pm 0.102	0.553 \pm 0.123	0.530 \pm 0.117	0.501 \pm 0.117	0.604 \pm 0.194	0.501 \pm 0.117
	Tumor	0.478 \pm 0.122	0.570 \pm 0.164	0.504 \pm 0.143	0.478 \pm 0.122	0.527 \pm 0.240	0.480 \pm 0.119	0.527 \pm 0.240

Table 12: Random forest regression and their evaluations.

Metrics are shown as mean \pm SD.

Dataset		MAE	RMSE
Entire	Total	704.909 \pm 249.010	894.943 \pm 246.192
	Normal	803.487 \pm 145.365	979.334 \pm 158.813
	Tumor	811.505 \pm 204.788	1005.182 \pm 197.351
DAT	Total	823.700 \pm 141.448	994.698 \pm 157.983
	Normal	663.414 \pm 147.203	825.461 \pm 151.120
	Tumor	729.302 \pm 179.940	884.863 \pm 181.154

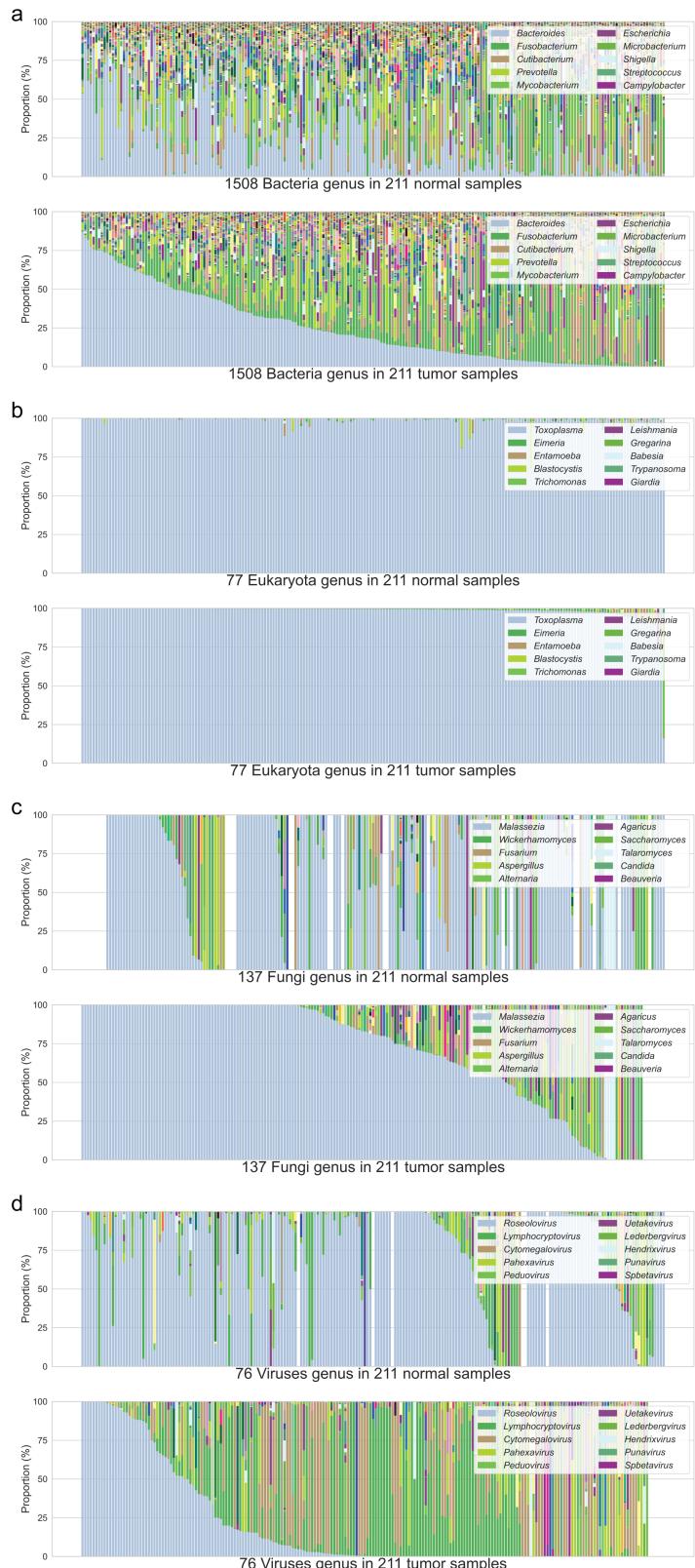


Figure 21: Gut microbiome compositions in genus level.

Taxa were sorted from the most prevalent taxon to the least prevalent taxon. CRC patients were sorted by the most prevalent taxon in descending order. **(a)** Bacteria kingdom **(b)** Eukaryota kingdom **(c)** Fungi kingdom **(d)** Viruses kingdom

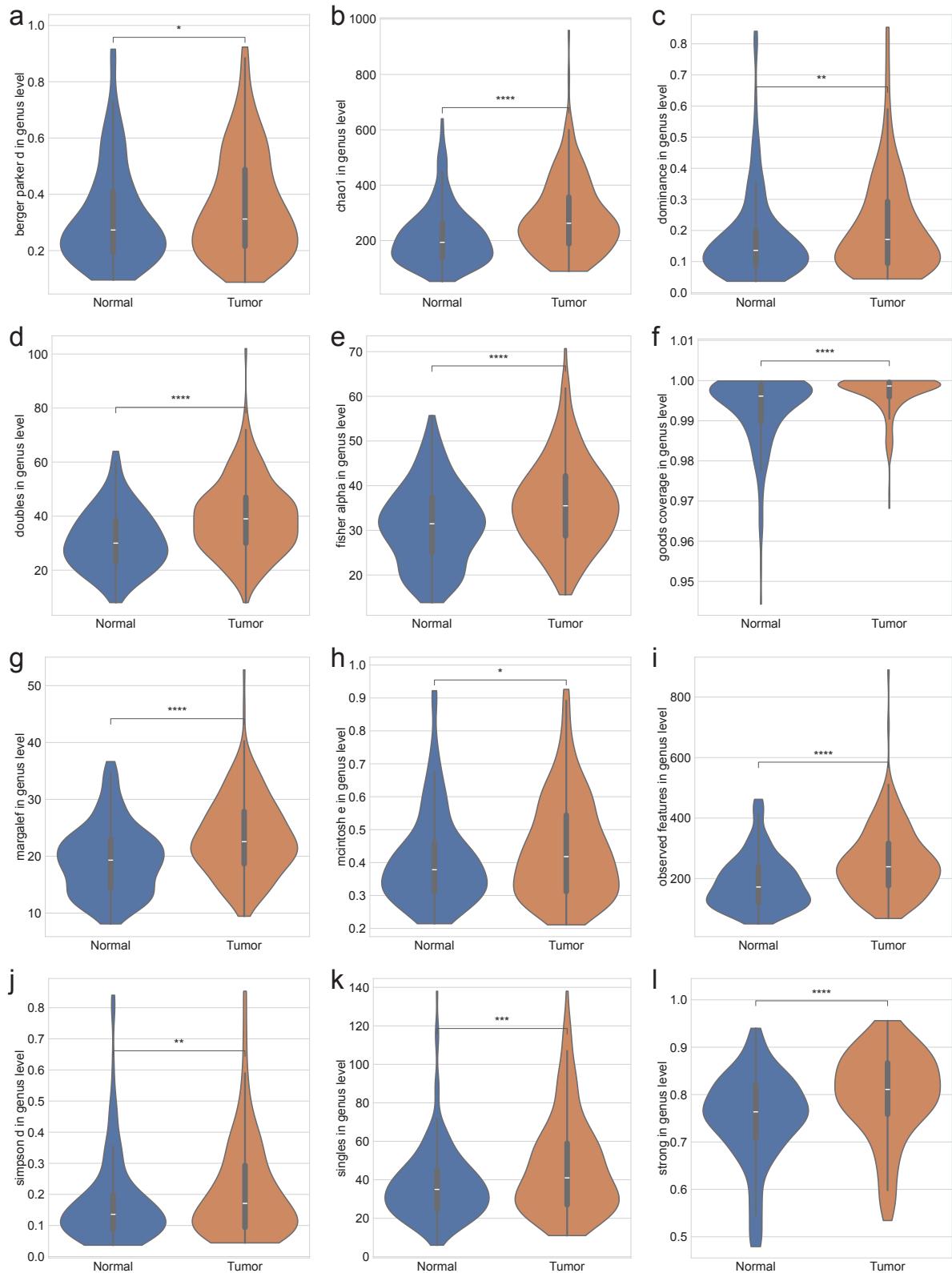


Figure 22: Alpha-diversity indices in genus level.

(a) Berger-Parker d (b) Chao1 (c) Dominance (d) Doubles (e) Fisher α (f) Good's coverage (g) Margalef (h) McIntosh (i) Observed features (j) Simpson d (k) Singles (l) Strong. MWU test: $p < 0.05$ (*), $p < 0.01$ (**), $p < 0.001$ (***), and $p < 0.0001$ (****)

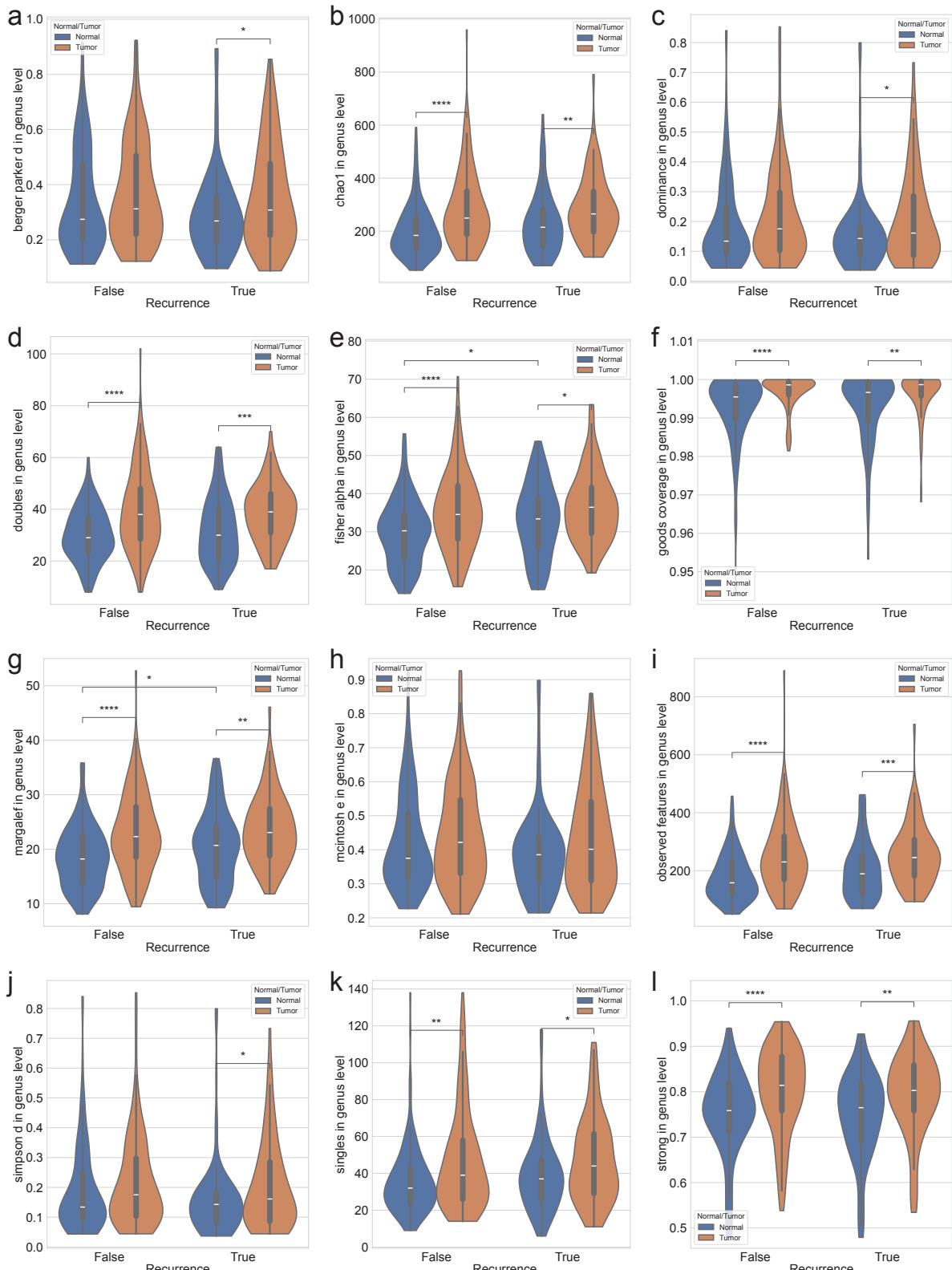


Figure 23: Alpha-diversity indices with recurrence in genus level.

(a) Berger-Parker d (b) Chao1 (c) Dominance (d) Doubles (e) Fisher α (f) Good's coverage (g) Margalef (h) McIntosh (i) Observed features (j) Simpson d (k) Singles (l) Strong. MWU test: $p < 0.05$ (*), $p < 0.01$ (**), $p < 0.001$ (***), and $p < 0.0001$ (****)

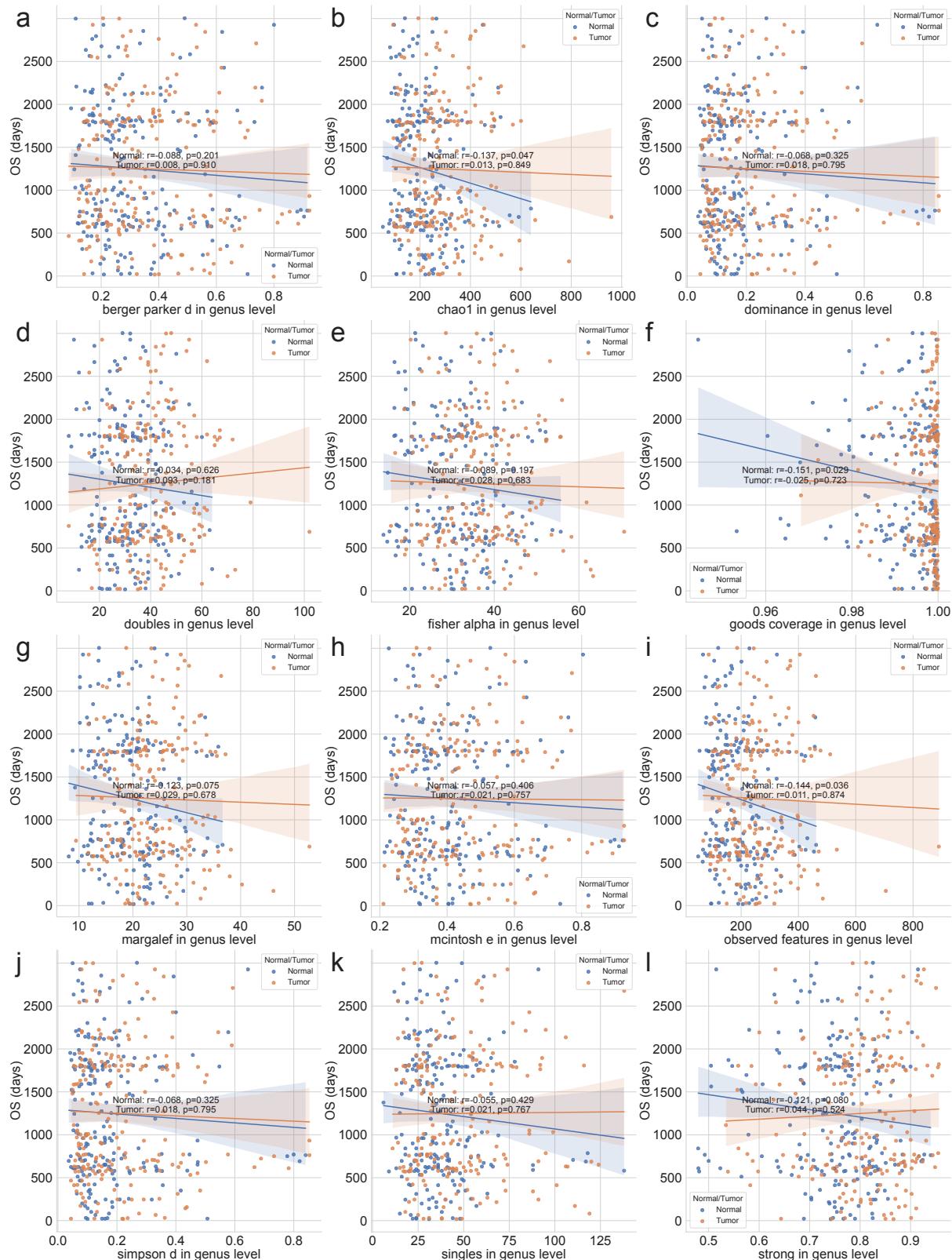


Figure 24: Alpha-diversity indices with OS in genus level.

(a) Berger-Parker d (b) Chao1 (c) Dominance (d) Doubles (e) Fisher α (f) Good's coverage (g) Margalef (h) McIntosh (i) Observed features (j) Simpson d (k) Singles (l) Strong. Statistical significance was calculated by the Spearman correlation.

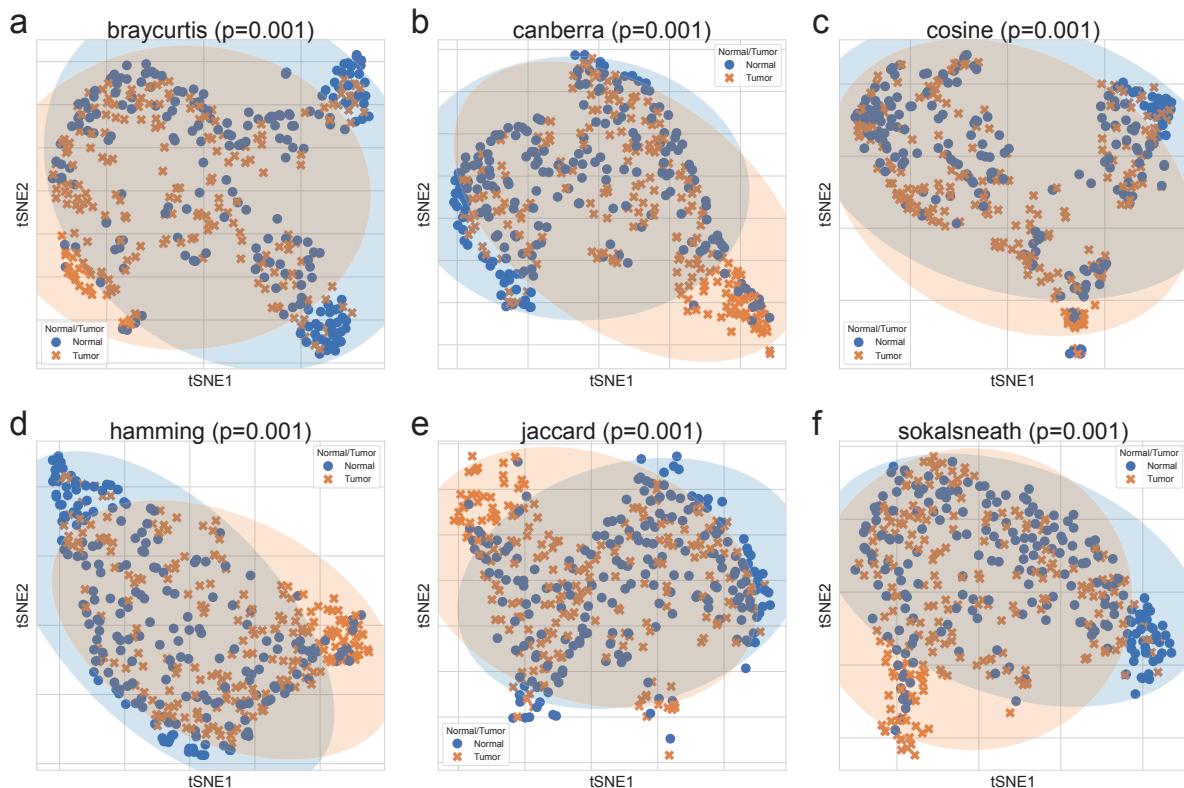


Figure 25: Beta-diversity indices in genus level.

Beta-diversity indices were visualized using a tSNE-transformed plot. The confidence ellipses are shown to display the distribution of each sub-group (Normal or Tumor). **(a)** Bray-Curtis **(b)** Canberra **(c)** Cosine **(d)** Hamming **(e)** Jaccard **(f)** Sokal-Sneath. Statistical significance were determined by PERMANOVA test.

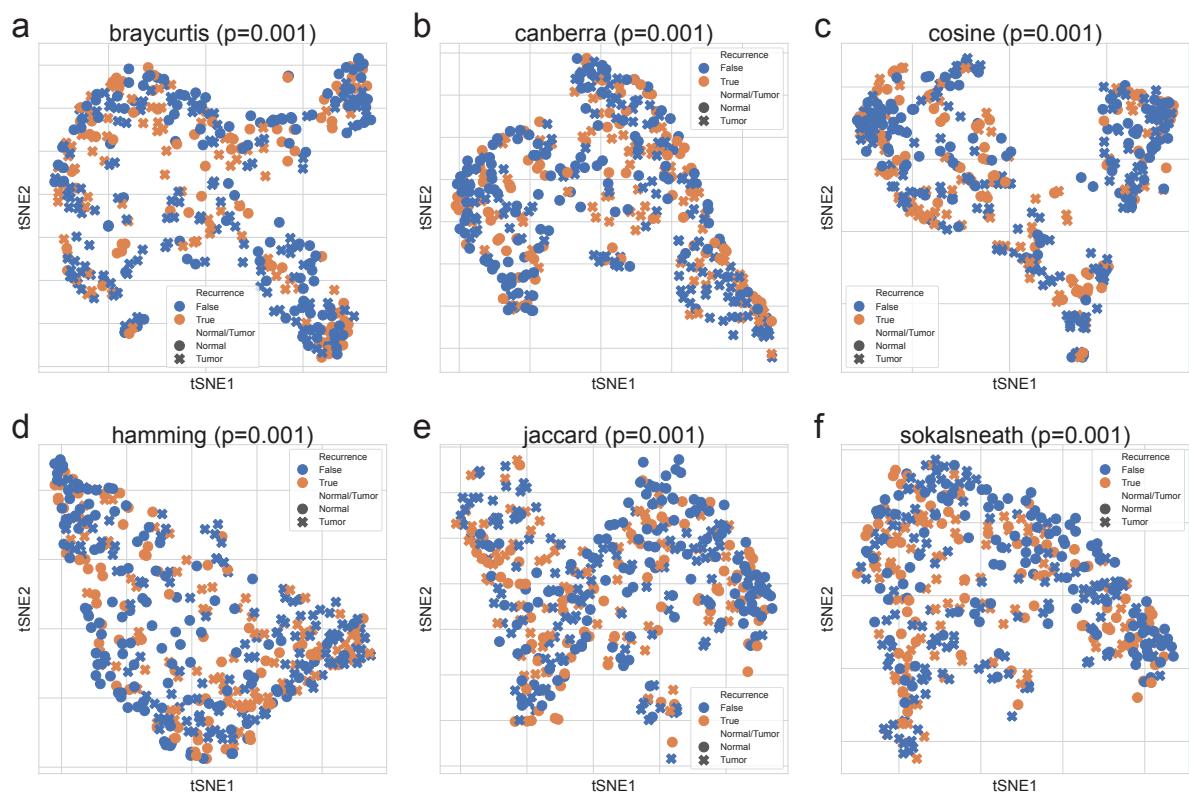


Figure 26: Beta-diversity indices with recurrence in genus level.

Beta-diversity indices were visualized using a tSNE-transformed plot. **(a)** Bray-Curtis **(b)** Canberra **(c)** Cosine **(d)** Hamming **(e)** Jaccard **(f)** Sokal-Sneath. Statistical significance were determined by PERMANOVA test.

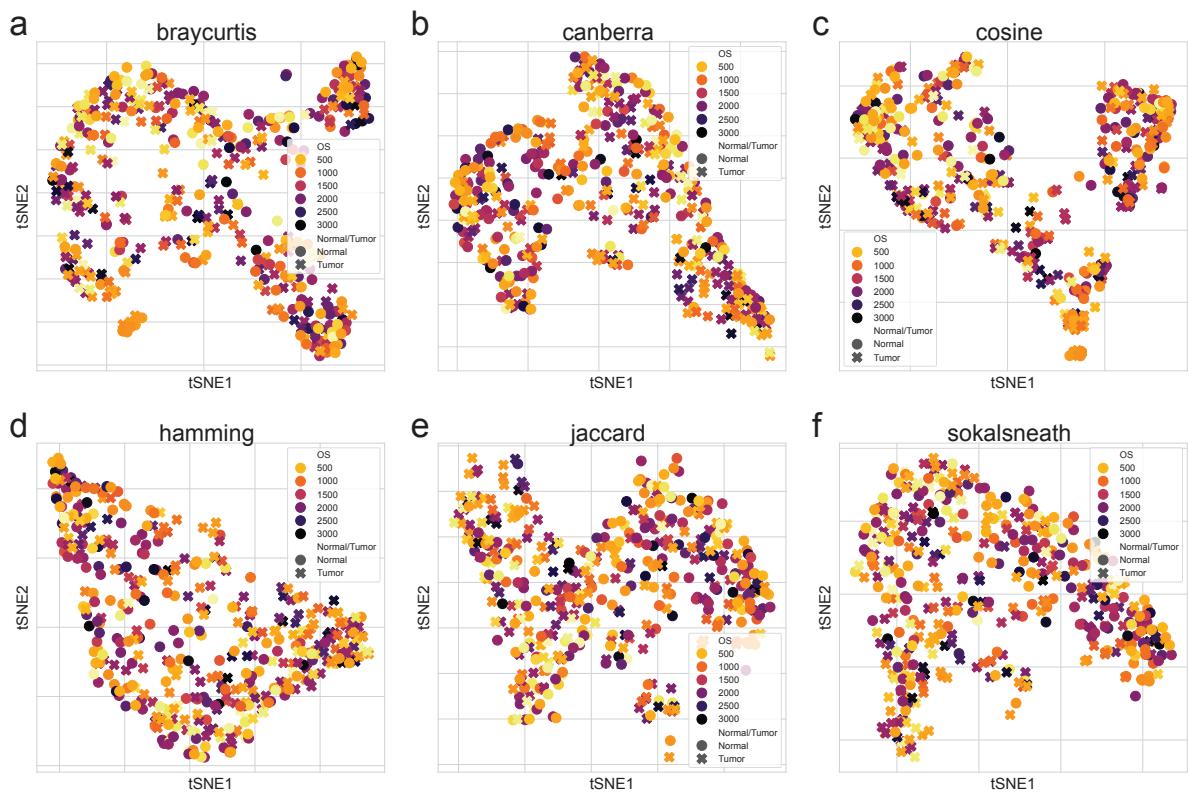


Figure 27: Beta-diversity indices with OS in genus level.

Beta-diversity indices were visualized using a tSNE-transformed plot. **(a)** Bray-Curtis **(b)** Canberra **(c)** Cosine **(d)** Hamming **(e)** Jaccard **(f)** Sokal-Sneath.

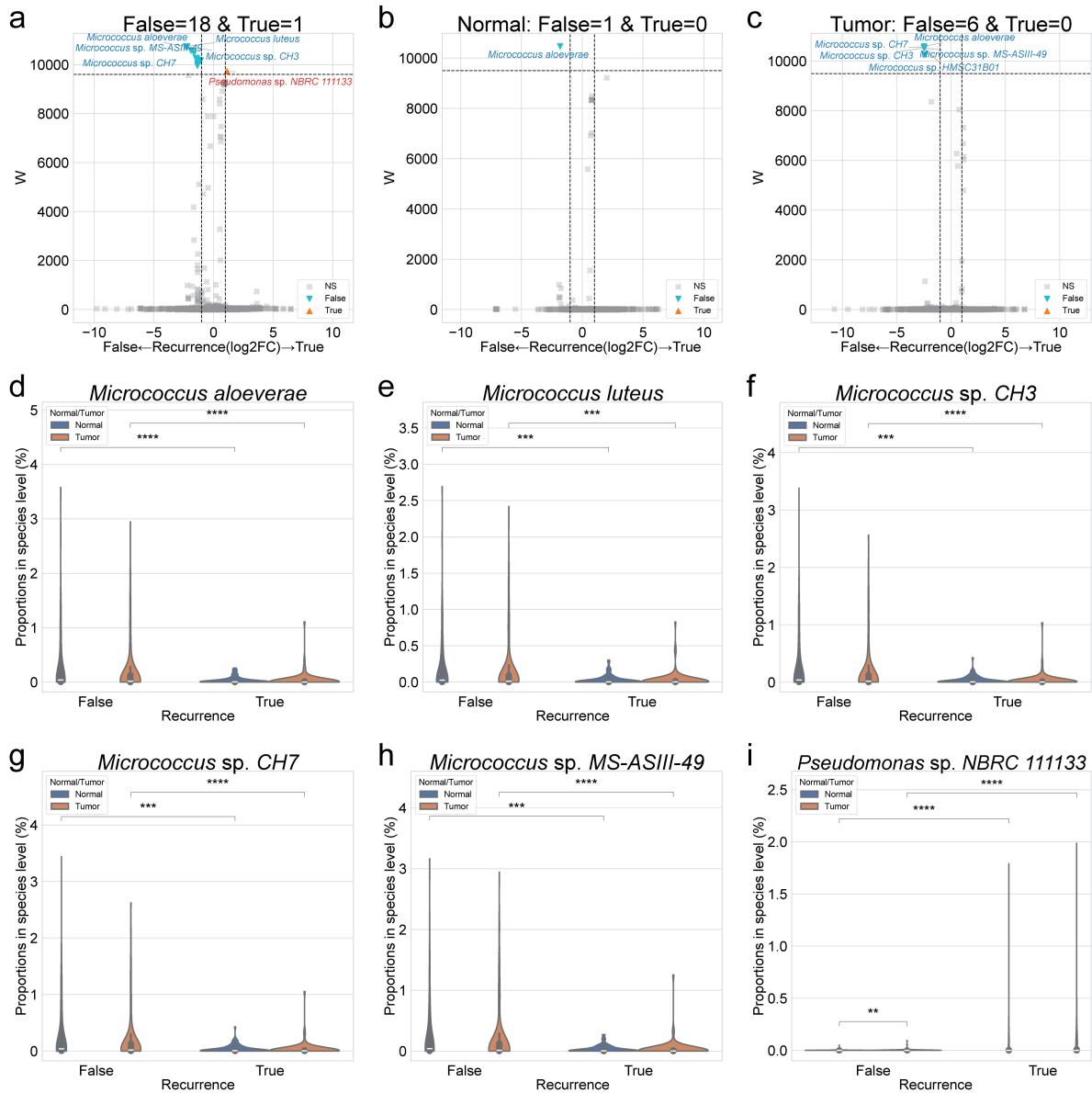


Figure 28: DAT with recurrence in species level.

(a-c) Volcano plots with recurrence. x-axis indicates $\log_2(\text{Fold Change})$ on recurrence, and y-axis indicates ANCOM significance (W). **(a)** Total **(b)** Normal samples **(c)** Tumor samples. **(d-i)** Violin plots of each taxon proportion with recurrence. **(d)** *Micrococcus aloeverae* **(e)** *Micrococcus luteus* **(f)** *Micrococcus* sp. *CH3* **(g)** *Micrococcus* sp. *CH7* **(h)** *Micrococcus* sp. *MS-ASIII-49* **(i)** *Pseudomonas* sp. *NBRC 111133*. Significant threshold is: $|\log_2 \text{Fold Change}| > 1.0$ and $W > 9600$. WU test: $p < 0.05$ (*), $p < 0.01$ (**), $p < 0.001$ (***)¹, and $p < 0.0001$ (****)

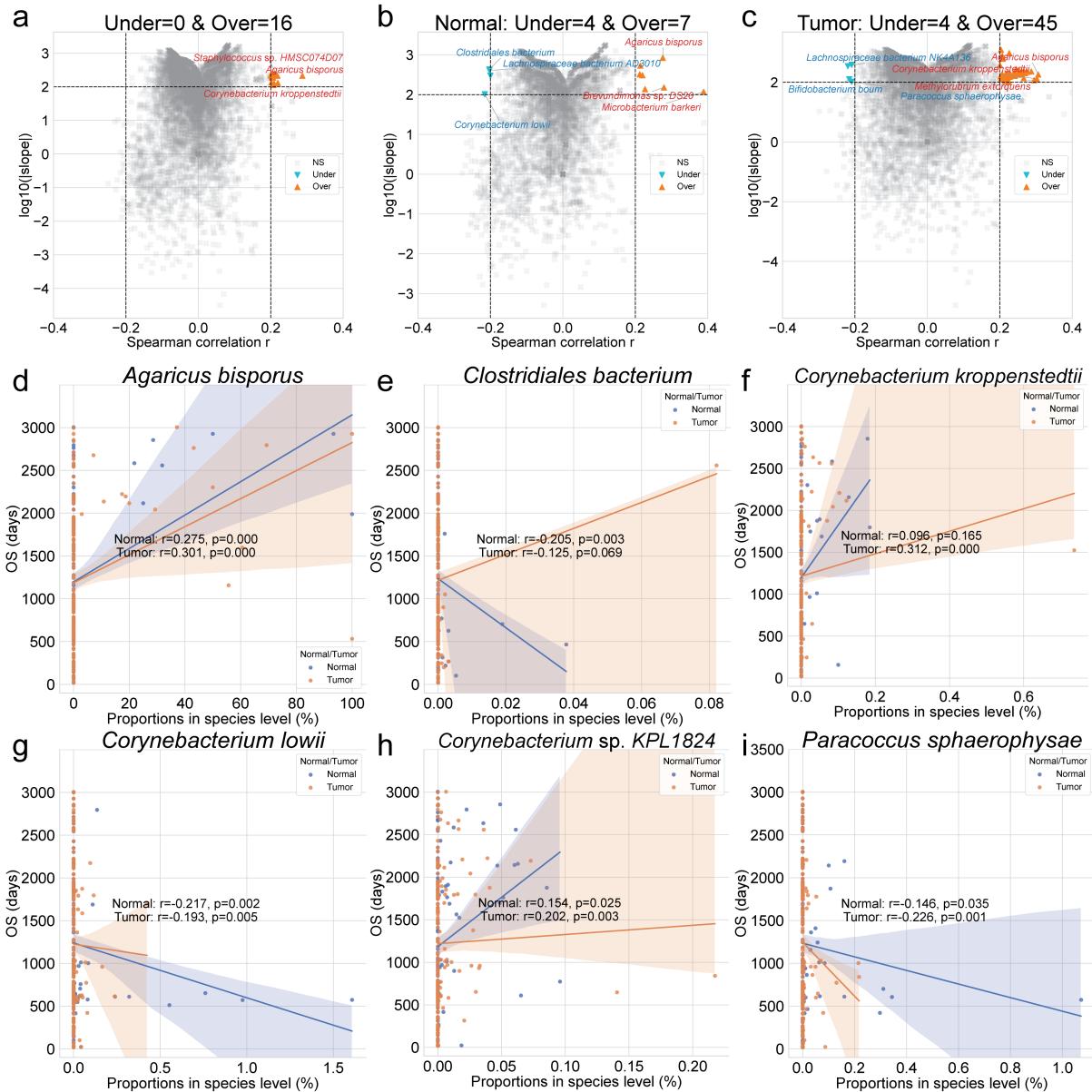


Figure 29: DAT with OS in species level.

(a-c) Volcano plots with OS. x-axis indicates Spearman correlation coefficient (r), and y-axis indicates $\log_{10}(|\text{slope}|)$. **(a)** Total **(b)** Normal samples **(c)** Tumor samples. **(d-li)** Scatter plots of each taxon proportion with OS. **(d)** *Agaricus bisporus* **(e)** *Clostridiales bacterium* **(f)** *Corynebacterium kroppenstedtii* **(g)** *Corynebacterium lowii* **(h)** *Corynebacterium sp. KPL1824* **(i)** *Paracoccus sphaerophysae*. Statistical significance were calculated with Spearman correlation (r and p): $|r| > 0.2$ and \log_{10} slope > 2.0 .

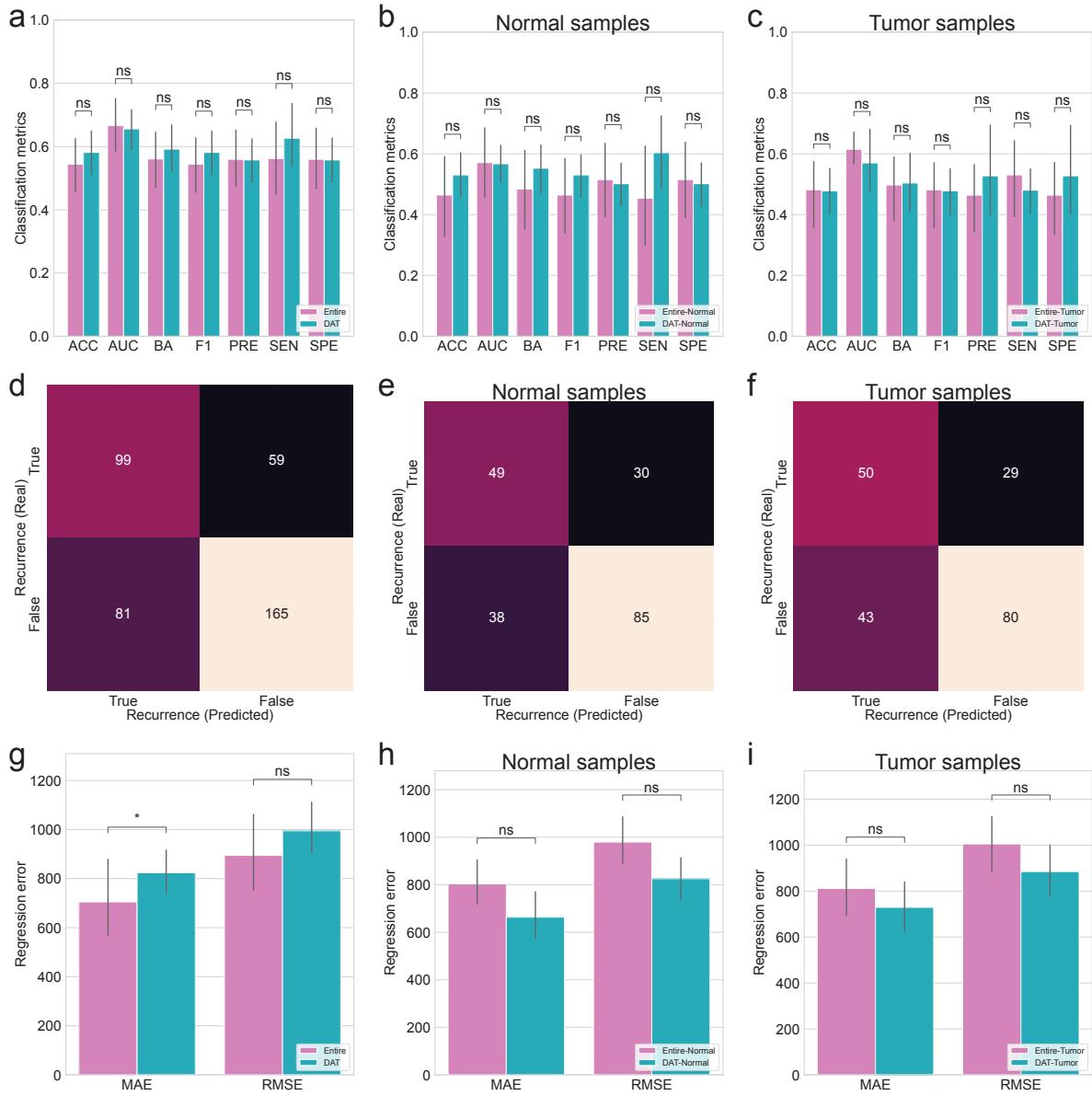


Figure 30: **Random forest classification and regression.**

(a-c) Random forest classification metrics for recurrence. **(a)** Total **(b)** Normal samples **(c)** Tumor samples. **(d-f)** Random forest classification confusion matrices for recurrence. **(d)** Total **(e)** Normal samples **(f)** Tumor samples. **(g-i)** Random forest regression errors for OS. **(g)** Total **(h)** Normal samples **(i)** Tumor samples. MWU test: $p < 0.05$ (*), $p < 0.01$ (**), $p < 0.001$ (***), and $p < 0.0001$ (****)

1365 **4.4 Discussion**

1366 This study provides a comprehensive metagenomic signature analysis of Korean CRC patients by
1367 examining prevalent microbial taxa, diversity indices, DAT selection, and random forest-based predictions
1368 for recurrence and survival outcomes. Our analysis revealed distinct prevalent microbial communities
1369 in CRC patients (Figure 21), with significant difference between tumor tissues and adjacent matched
1370 normal tissues. Alpha-diversity indices analysis showed an overall shift in microbial diversity within
1371 tumor samples (Figure 22, Figure 23, and Figure 24), while beta-diversity analyses indicated significant
1372 changes in microbial composition associated with recurrence history and survival duration (Figure 25,
1373 Figure 26, and Figure 27). Through DAT selection by ANCOM and Spearman correlation, we identified
1374 key microbial taxa link to recurrence history (Table 9 and Figure 28) and OS duration (Table 10 and
1375 Figure 29), highlighting potential microbial biomarkers for CRC prognosis. To evaluate the predictive
1376 capacity of these microbial features, we implemented random forest-based machine learning models,
1377 where random forest classification demonstrated moderate accuracy (0.570 ± 0.164 , mean \pm SD) for CRC
1378 recurrence prediction (Table 11 and Figure 30) and random forest regression showed slightly high
1379 error (729.302 ± 179.940 , mean \pm SD) for OS prediction (Table 12 and Figure 30), suggesting that gut
1380 microbiome alterations alone are insufficient for precise prognosis and may interact with host genetic
1381 factors such as germline and somatic mutations. These findings underscore the potential of microbial
1382 biomarkers in CRC risk stratification, emphasizing the need for multi-omics integration to improve
1383 predictive models and personalized medicine strategies in CRC.

1384 In the bacteria kingdom (Figure 21a), *Bacteroides* genus is the most frequent genus in tumor tis-
1385 sues, then came *Fusobacterium* and *Cutibacterium* genera. These results also accord with previous
1386 studies, which showed that *Bacteroides fragilis* (Scott, Whittle, Jeraldo, & Chia, 2022; Purcell, Permain,
1387 & Keenan, 2022), *Fusobacterium nucleatum* (Wang & Fang, 2023; Zepeda-Rivera et al., 2024), and
1388 *Cutibacterium acnes* (Benej et al., 2024) have significant roles in tumorigenesis and development of CRC.
1389 Further, not only those bacterium genera individually, the association between *Bacteroides* genus and
1390 *Fusobacterium* genus is reported (Viljoen, Dakshinamurthy, Goldberg, & Blackburn, 2015; Joo et al.,
1391 2024; Duy et al., 2024; Conde-Pérez et al., 2024), suggesting possible contribution to CRC pathogenesis
1392 through mechanisms such as biofilm formation, immune evasion, and/or metabolic interactions with
1393 other dysbiotic taxa. Given that *Fusobacterium* genus has been shown to co-aggregate with *Bacteroides*
1394 genus, it is plausible that *Cutibacterium* genus might interact with these genera to influence inflammation,
1395 epithelial barrier integrity, and tumor progression. Thus, further studies integrating functional metage-
1396 nomics, metabolomics, and host-microbiome interactions are warranted to elucidate the precise role of
1397 *Cutibacterium* genus and its relationship with CRC-associated microbial networks.

1398 Analysis of eukaryotic and fungal microbial compositions revealed that the *Toxoplasma* genus was
1399 prevalent in both normal and tumor samples (Figure 21b), while the *Malassezia* genus was more prevalent
1400 in tumor samples (Figure 21c). The consistent presence of *Toxoplasma* genus across both sample types
1401 suggests that this intracellular pathogen may be a stable component of the gut microbiome, although its
1402 role in CRC pathogenesis remains unclear (Yu et al., 2020; Zavareh et al., 2021). In contrast, the increase

1403 prevalence of *Malassezia* genus in tumor tissue aligns with emerging evidence that certain fungal genus
1404 may contribute to CRC-promoting inflammation and metabolic alterations (R. Gao et al., 2017; Yuan et
1405 al., 2025), suggesting a potential role in CRC development and progression. These findings highlight the
1406 need for further investigation into the functional impact of eukaryotic and fungal microbiota in CRC for
1407 shaping the tumor microenvironment.

1408 In normal tissue samples, *Roseolovirus* genus was the most prevalent viral taxon (Figure 21d),
1409 indicating its stable presence in the gut virome of healthy colonic tissues. However, in tumor tissue
1410 samples, *Lymphocryptovirus* and *Cytomegalovirus* genera were more prevalent viral taxa, suggesting
1411 an alteration in viral community structure associated with CRC progression. This viral compositional
1412 shift aligns with the Anna Karenina principle (Ma, 2020; W. Li & Yang, 2025), implying that microbial
1413 communities in diseased states exhibit greater instability and variability compare to their adjacent normal
1414 tissues. The emergence of *Lymphocryptovirus* (Mjelle, Castro, & Aass, 2025; De Flora & Bonanni,
1415 2011) and *Cytomegalovirus* (Harkins et al., 2002; Taher et al., 2014; Bender et al., 2009) genera in
1416 tumor samples raises the possibility that oncogenic viruses may contribute to CRC carcinogenesis by
1417 promoting chronic inflammation, immune modulation, and/or direct viral-host interactions affecting
1418 cellular transformation. Moreover, the detection of tumor-associated viral alterations in adjacent normal
1419 tissues supports the concept of field cancerization (Curtius et al., 2018; Rubio et al., 2022), where viral
1420 dysbiosis may extend beyond the tumor itself, creating a pro-tumorigenic microenvironment even before
1421 malignant transformation occurs. These findings underscore the potential impact of viral communities in
1422 CRC and highlight the requirement for further research into their functional roles in carcinogenesis of
1423 CRC.

1424 Alpha-diversity indices revealed a significant increase in microbial diversity in tumor samples com-
1425 pared to its adjacent normal tissues (MWU test $p < 0.05$; Figure 22), suggesting CRC is associated with
1426 a more heterogeneous gut microbiome (Liu et al., 2021). The increase in alpha-diversity indices within
1427 tumor tissues may support the Anna Karenina principle and/or the concept of field cancerization, where
1428 microbial alterations extend beyond the tumor site and contribute to a pre-malignant microenvironment.
1429 The enrichment of distinct bacterial, eukaryotic, fungal, and viral taxa within tumor samples suggest that
1430 microbial dysbiosis in CRC is not limited to a single pathogenic genus or species but rather involves
1431 complex community-level changes.

1432 Furthermore, alpha-diversity indices in relation to recurrence history revealed distinct microbial
1433 diversity patterns between normal and tumor tissue samples (Figure 23). In recurrence patients, tumor
1434 samples exhibited a greater increase in alpha-diversity indices compared to their adjacent normal tissues
1435 (11/12 indices, 92% indices; Figure 23), suggesting that a more heterogeneous microbial community
1436 may be linked to tumor aggressiveness and recurrence potential (Huo et al., 2022; Vigneswaran &
1437 Shogan, 2020). This trend aligns with a highly diverse but dysregulated microbiome in tumor samples
1438 may contribute to immune evasion, chronic inflammation, and tumor-promoting metabolic changes.
1439 In non-recurrence patients, although tumor samples still exhibited increased alpha-diversity indices
1440 compared to normal tissues, the difference was less pronounced (8/12 indices, 67% indices; Figure 23),
1441 suggesting that a relatively more stable microbiome in tumor tissues may be associated with favorable

1442 survival outcomes (Avuthu & Guda, 2022). These findings reinforce the concept that tumor microbiome
1443 changes are inconsistent across CRC patients, supporting the Anna Karenina principle. Additionally,
1444 the differences in alpha-diversity indices of normal tissues between recurrence and non-recurrence
1445 patients further suggest (Figure 23e and Figure 23g) that specific microbial communities may influence
1446 post-treatment disease progression.

1447 Moreover, alpha-diversity indices and OS duration in CRC patients revealed distinct patterns between
1448 normal and tumor tissues (Figure 24), suggesting that microbial diversity in non-cancerous lesions may
1449 play a role in cancer prognosis (Galeano Niño et al., 2022). While no significant correlation was found
1450 between tumor-associated microbiome and OS duration, three of the 12 alpha-diversity indices exhibited
1451 negative correlations with OS in normal tissues (Figure 29b, Figure 29f, and Figure 29i), indicating
1452 that lower microbial heterogeneity in normal lesions was associated with longer survival. This finding
1453 suggests that a more heterogeneous microbial community in normal colon tissues may contribute to a
1454 microenvironment that fosters tumor progression, aligning with the field cancerization. Therefore, the
1455 negative correlations observed only in normal tissues suggests that pre-onset dysbiosis in non-cancerous
1456 regions could influence prognosis of CRC, potentially serving as an early indicator of cancer progression
1457 risk.

1458 Beta-diversity indices revealed significant differences in gut microbiome compositions between tumor
1459 and normal tissues (Figure 25), aligning with the alpha-diversity indices and further confirming the
1460 presence of dysbiosis in gut microbiome of CRC. The distinct clustering of tumor and normal samples in
1461 beta-diversity indices (PERMANOVA $p < 0.001$) suggests that CRC is associated with a major alteration
1462 in microbial structure. This transformation may be driven by the expansion of tumor-associated taxa and
1463 the shrinkage of protective taxa, resulting in a tumor-supportive microenvironment. This clear separation
1464 in beta-diversity indices between tumor and normal tissues supports again the field cancerization, where
1465 microbial alterations extend beyond tumor lesions and affect surrounding non-cancerous lesions.

1466 Furthermore, beta-diversity indices demonstrated significant microbial composition shifts between
1467 normal and tumor tissues in accordance with recurrence status (Figure 26), suggesting that dysbiosis
1468 in the gut microbiome may play an essential role in CRC progression and post-treatment recurrence.
1469 By the beta-diversity indices, the observed recurrence-associated microbial shifts highlight the potential
1470 of beta-diversity index as predictive markers for recurrence risk of CRC, warranting further studies to
1471 explore their functional significance and potential integration into microbiome-based prognostic models.

1472 Moreover, beta-diversity indices suggested a potential association between the gut microbiome com-
1473 position and OS in CRC patients (Figure 27), as distinct clustering were observed in relation to survival
1474 duration. However, due to the continuous nature of survival duration, direct statistical comparison using
1475 PERMANOVA test could be not performed, limiting the ability to formally quantify these differences.
1476 Despite this limitation, the observed separation of microbial communities along OS suggests that the
1477 gut microbiome composition may play a major role in CRC prognosis, potentially influencing immune
1478 response, tumor progression, and treatment outcomes. This lack of statistical validation highlights the
1479 need for alternative approaches to better assess the relationship between microbiome structure and sur-
1480 vival outcome. Further investigation is required to determine whether specific microbial taxa drive these

1481 compositional shifts and whether gut microbiome profiles could serve as prognostic biomarkers for CRC
1482 survival outcomes.

1483 To identify recurrence-related DAT in CRC, we applied ANCOM to compare the gut microbiome
1484 compositions between recurrence and non-recurrence patients (Table 9 and Figure 28). By applying
1485 ANCOM separately to total samples (Figure 28a), normal samples (Figure 28b), and tumor samples
1486 (Figure 28c), we identified both global and tissue-specific microbial shifts linked to CRC recurrence.
1487 Among these 19 recurrence-related DAT (Table 9), several DAT belonging to the *Micrococcus* and
1488 *Staphylococcus* genera were nominated as non-recurrence-enriched DAT. *Micrococcus* genus has been
1489 reported with anti-bacterial, anti-fungal, and anti-inflammatory activities (Tizabi & Hill, 2023), and
1490 another study has found that the production of carotenoid pigments from *Micrococcus luteus* (Figure 28e)
1491 exhibited promising antibiotics agents (Hegazy, Abu-Hussien, Elsenosy, El-Sayed, & Abo El-Naga, 2024).
1492 In this CRC study participants, *Cutibacterium acnes* was selected one of the non-recurrence-enriched DAT
1493 (Table 9). This finding is consistent with previous studies which have suggested that *Cutibacterium acnes*
1494 inhibits the activities of pathogens, such as *Staphylococcus aureus*, and suppresses tumor growth (Benej
1495 et al., 2024; Ding, Lian, Tam, & Oh, 2024). On the other hand, in this CRC study participants, many
1496 *Staphylococcus* species have chosen as non-recurrence-enriched DAT (Table 9); however, this outcome
1497 is contrary to previous studies which have described that cancer-promoting activity of *Staphylococcus*
1498 *aureus* (Z. Li, Zhuang, Wang, Wang, & Dong, 2021; Cuervo et al., 2010), suggesting opposite behaviors
1499 between *Staphylococcus aureus* and other *Staphylococcus* species. Last but not least, *Pseudomonas* sp.
1500 *NBRC 11113* has been found as the only recurrence-enriched DAT (Figure 28i). This also accords with
1501 earlier studies, which showed that *Pseudomonas aeruginosa* infections in cancer patients (Ohmagari et
1502 al., 2005; Paprocka et al., 2022).

1503 To determine the OS-correlated DAT in CRC, we applied Spearman correlation to measure effects
1504 of the gut microbiome composition with OS (Table 10 and Figure 29). By implementing Spearman
1505 correlation to total samples (Figure 29a), normal samples (Figure 29b), and tumor samples (Figure 29c),
1506 we found that CRC survival is associated with both tissue type-specific and global microbial alterations.
1507 Among these 57 OS-correlated DAT (Table 29), several OS-correlated DAT from the *Corynebacterium*
1508 and *Staphylococcus* genera have significant correlations with survival duration of CRC. *Agaricus bisporus*
1509 has positive correlation with OS both in normal and tumor samples (Figure 29d). In accordance with
1510 this finding, previous studies have demonstrated that a polysaccharide produced from *Agaricus bisporus*
1511 exhibited anti-cancerous activity in colon cancer (Dong, Wang, Tang, Liu, & Gao, 2024; El-Deeb et
1512 al., 2022; N. Zhang, Liu, Tang, Yang, & Wang, 2023). Furthermore, most of *Corynebacterium* genus,
1513 including *Corynebacterium kroppenstedtii* (Figure 29f) and *Corynebacterium* sp. *KPL1824* (Figure
1514 29h), have positive correlations with OS; however, *Corynebacterium lowii* (Figure 29g) has negative
1515 correlation with OS both in normal and tumor samples. Comparison of the findings with those of other
1516 studies confirms a breast cancer risk factor of *Corynebacterium afermentans* (J. An, Kwon, Oh, & Kim,
1517 2025), an increasing of *Corynebacterium appendicis* in CRC (Hasan et al., 2022), an inhibition role of
1518 *Corynebacterium matruchotii* of cancer growth in oral squamous cell carcinoma (Shen et al., 2022), and
1519 a promoting cancer cell apoptosis of *Corynebacterium durum* (S. Kim et al., 2024), warranting future

1520 investigations to selecting pro-tumorigenic and anti-tumorigenic species of *Corynebacterium* genus. *Clostridiales*
1521 *bacterium* has negative correlation with OS in normal samples (Figure 29e). However, this result does not
1522 support previous researches which have demonstrated that anti-cancer activities with immune modulation
1523 of *Clostridiales* genus (Montalban-Arques et al., 2021; Minton, 2003), suggesting different roles from
1524 *Clostridiales* species for immune response against cancer. Many species from *Staphylococcus* genus
1525 have positive correlations with OS (Table 10). Although, these results differ from some published
1526 studies which indicated cancer prevention and treatment via reduction of *Staphylococcus epidermidis*
1527 (Bernardo et al., 2023; Kepp, Zitzgag, & Kroemer, 2023), these results are consistent with other published
1528 researches which suggested that other species of *Staphylococcus* genus exhibited anti-cancer activities
1529 (Hassan, Mustafa, Rahim, & Isa, 2016; M. Zhang et al., 2022). Moreover, *Lachnospiraceae bacterium*
1530 *AD3010* and *Porphyromonas macacae* have negative correlations with OS in normal samples, while
1531 *Lachnospiraceae bacterium NK4A136* and *Paracoccus sphaerophysae* have negative correlations with OS
1532 in tumor samples (Table 10). Previous studies have addressed that high abundance of *Lachnospiraceae*
1533 genus in the gut microbiome showed anti-tumor roles in the CRC (Hexun et al., 2023; X. Zhang et al.,
1534 2023), indicating that more comprehensive investigation of species from *Lachnospiraceae* genus might be
1535 required. *Porphyromonas gingivalis*, a well-known periodontitis pathogens from *Porphyromonas* genus,
1536 was also reported promoting cancer resistance and development on CRC, lung cancer, and oesophageal
1537 cancer (León et al., 2007; Katz et al., 2009; S. Gao et al., 2021), providing a warrant to elucidate cancer-
1538 related roles of not only *Porphyromonas gingivalis* but also other *Porphyromonas* genus. *Paracoccus*
1539 *sphaerophysae* displayed negative correlation with OS in tumor samples and insignificantly negative
1540 correlation with OS (Spearman $|r| \leq 0.2$) in normal samples (Figure 29i), it is consistent with the literature
1541 which have shown *Paracoccus* genus is more prevalent in nasopharyngeal carcinoma group than healthy
1542 individuals (Lu et al., 2024).

1543 To assess the predictive potential of recurrence-related DAT in CRC recurrence risk, we implemented a
1544 random forest classification model (Table 11 and Figure 30). The classification model achieved moderated
1545 classification performance (0.570 ± 0.164 , mean \pm SD), indicating that while gut microbial provides some
1546 predictive value, it is likely insufficient as a standalone biomarker for recurrence risk of CRC. This
1547 limited predictive accuracy may be attributed to the complex and dynamic nature of gut microbiome
1548 network, where epigenetic modifications and immune modulation collectively influence development and
1549 progression of CRC. Additionally, host-microbiome interactions, including metabolic pathways, may
1550 further contribute to recurrence of cancer, warranting a more integrative multi-omics approach. Therefore,
1551 future studies incorporating genomic sequencing data, e.g. somatic mutations and host immune signatures,
1552 could provide a more comprehensive understanding of how microbial dysbiosis interacts with tumor
1553 biology.

1554 To evaluate the predictive potential of OS-related DAT in survival duration of CRC, we employed
1555 a random forest regression model (Table 12 and Figure 30). The regression model exhibited moderate
1556 regression error (729.302 ± 179.940 , mean \pm SD), suggesting that while gut microbiome composition
1557 provides some predictive values for cancer patient survival, it is likely influenced by additional host-
1558 specific and environmental factors. The complex interplay between the gut microbiome and CRC

1559 progression involves sophisticated microbial networks, metabolic interactions, and immune response,
1560 making it difficult to capture survival outcomes solely based on microbiome features. Furthermore,
1561 host-microbiome interactions, including MSI, tumor mutational burden, and epigenetic modifications,
1562 likely play a crucial role in determine favorable or unfavorable survival. These findings highlight the need
1563 for multi-omics integration, combining genomic sequencing data and metagenomic functional analysis,
1564 to gain deeper insights into how microbial dysbiosis interacts with tumor biology and clinical outcomes.
1565 Future studies incorporating machine learning models with multi-layered biological data may improve the
1566 accuracy of survival prediction and contribute to personalized medicine approaches for CRC therapeutics.

1567 **5 Conclusion**

1568 This dissertation underscores the critical character of microbiome research in understanding disease
1569 mechanisms, predicting health outcomes, and advancing personalized medicine. By investigating PTB,
1570 periodontitis, and CRC, this dissertation demonstrated how microbial diversity alters, DAT, and machine
1571 learning-based modeling contribute to disease classification and prognosis. While each condition exhibited
1572 unique microbiome alterations, the findings collectively support the Anna Karenina principle, which
1573 suggests that microbial communities in patients with disease become more variable and dysregulated
1574 compared to their relatively stable and uniform counterparts in healthy individuals. The Anna Karenina
1575 principle was evident in all three diseases examined, where dysbiosis not only disrupted microbial
1576 homeostasis but also contributed to disease progression and development. The ability to identify disease-
1577 specific microbial signature reinforces the importance of microbiome profiling as a new therapeutic
1578 guidance.

1579 In the PTB study (Section 2), salivary microbiome profiling revealed distinct microbial shifts between
1580 PTB and FTB, with a random forest-based model achieving high accuracy in assessing PTB risk.
1581 Similarly, the periodontitis study (Section 3) identified salivary microbial markers that classified between
1582 healthy individuals and multiple stages of periodontitis, suggesting the potential for salivary microbiome-
1583 based diagnostics in management and treatment of periodontitis. The CRC study (Section 4) revealed
1584 significant alpha-diversity and beta-diversity indices differences between tumor and adjacent normal
1585 tissues, with distinct microbial compositions associated with recurrence status and survival duration.
1586 However, while random forest models for predicting recurrence risk and survival duration provided
1587 moderate accuracy, the findings suggest that gut microbiome composition alone may not be sufficient for
1588 precise clinical instruction.

1589 The Anna Karenina principle was particularly evident in the CRC study (Section 4), where gut
1590 microbial alterations were highly individualized among tumor samples, with recurrence status and survival
1591 duration related to divergent microbial community structures. This aligns with the concept of field
1592 cancerization, where dysbiosis extends beyond the tumor lesion, affecting adjacent non-cancerous lesions
1593 and potentially contributing to tumorigenesis and cancer development. These findings reinforce the
1594 complexity of host-microbiome interactions, where microbial imbalances may not only reflect disease
1595 status but actively participate in disease etiology through inflammation, metabolic alterations, and immune
1596 modulation. The variability in microbial community alterations across patients highlights the need for
1597 multi-omics integration, combining host genomic data to enhance personalized treatment and management
1598 strategies.

1599 Despite the promising insights gained from microbiome analyses, this dissertation acknowledges
1600 several limitations. The moderated predictive performance of machine learning models suggests that
1601 microbial features alone may not fully capture disease mechanisms. Future research should integrate
1602 multi-omics datasets, including host genomic mutations, metabolic profiles, and immune signatures,
1603 to improve biomarker discovery and disease prediction models. Additionally, population-specific mi-
1604 crobiome differences must be considered, as external validation in the periodontitis study (Section 3)

1605 showed variations in salivary microbiome composition between different ethnic groups. Large-scale
1606 and multi-center studies are essential to validate microbiome-based biomarkers and ensure their clinical
1607 applicability across diverse populations.

1608 Overall, this dissertation contributes to the growing field of microbiome-driven personalized medicine,
1609 demonstrating the potential of microbiome profiling, diversity analysis, identification DAT, and machine
1610 learning-based modeling in assessing disease risk and progression. By furthering our understanding of
1611 host-microbiome interactions, these findings pave a novel microbiome-targeted therapeutic strategies, ad-
1612 vancing personalized disease prevention and treatment. Moving forward, integrating microbiome research
1613 with genomics, metabolomics, and immunology holds the potential to transform disease management and
1614 personalized medicine, ultimately improving treatment outcomes across a broad spectrum of diseases.

¹⁶¹⁵ References

- ¹⁶¹⁶ Aagaard, K., Ma, J., Antony, K. M., Ganu, R., Petrosino, J., & Versalovic, J. (2014). The placenta harbors
¹⁶¹⁷ a unique microbiome. *Science translational medicine*, 6(237), 237ra65–237ra65.
- ¹⁶¹⁸ Abu-Ghazaleh, N., Chua, W. J., & Gopalan, V. (2021). Intestinal microbiota and its association with
¹⁶¹⁹ colon cancer and red/processed meat consumption. *Journal of gastroenterology and hepatology*,
¹⁶²⁰ 36(1), 75–88.
- ¹⁶²¹ Abusleme, L., Hoare, A., Hong, B.-Y., & Diaz, P. I. (2021). Microbial signatures of health, gingivitis,
¹⁶²² and periodontitis. *Periodontology 2000*, 86(1), 57–78.
- ¹⁶²³ Aitchison, J., Barceló-Vidal, C., Martín-Fernández, J. A., & Pawlowsky-Glahn, V. (2000). Logratio
¹⁶²⁴ analysis and compositional distance. *Mathematical geology*, 32, 271–275.
- ¹⁶²⁵ Aja, E., Mangar, M., Fletcher, H., & Mishra, A. (2021). Filifactor alocis: recent insights and advances.
¹⁶²⁶ *Journal of dental research*, 100(8), 790–797.
- ¹⁶²⁷ Alelyani, S. (2021). Stable bagging feature selection on medical data. *Journal of Big Data*, 8(1), 11.
- ¹⁶²⁸ Altabtbaei, K., Maney, P., Ganesan, S. M., Dabdoub, S. M., Nagaraja, H. N., & Kumar, P. S. (2021). Anna
¹⁶²⁹ karenina and the subgingival microbiome associated with periodontitis. *Microbiome*, 9, 1–15.
- ¹⁶³⁰ Altingöz, S. M., Kurgan, Ş., Önder, C., Serdar, M. A., Ünlütürk, U., Uyanık, M., ... Günhan, M.
¹⁶³¹ (2021). Salivary and serum oxidative stress biomarkers and advanced glycation end products in
¹⁶³² periodontitis patients with or without diabetes: A cross-sectional study. *Journal of periodontology*,
¹⁶³³ 92(9), 1274–1285.
- ¹⁶³⁴ Alverdy, J., Hyoju, S., Weigerinck, M., & Gilbert, J. (2017). The gut microbiome and the mechanism of
¹⁶³⁵ surgical infection. *Journal of British Surgery*, 104(2), e14–e23.
- ¹⁶³⁶ An, J., Kwon, H., Oh, S.-Y., & Kim, Y. J. (2025). Association between breast cancer risk factors and
¹⁶³⁷ blood microbiome in patients with breast cancer. *Scientific Reports*, 15(1), 6115.
- ¹⁶³⁸ An, S., & Park, S. (2022). Association of physical activity and sedentary behavior with the risk of
¹⁶³⁹ colorectal cancer. *Journal of Korean Medical Science*, 37(19).
- ¹⁶⁴⁰ Anderson, M. J. (2014). Permutational multivariate analysis of variance (permanova). *Wiley statsref:
1641 statistics reference online*, 1–15.
- ¹⁶⁴² Aruni, A. W., Mishra, A., Dou, Y., Chioma, O., Hamilton, B. N., & Fletcher, H. M. (2015). Filifactor
¹⁶⁴³ alocis—a new emerging periodontal pathogen. *Microbes and infection*, 17(7), 517–530.
- ¹⁶⁴⁴ Avuthu, N., & Guda, C. (2022). Meta-analysis of altered gut microbiota reveals microbial and metabolic
¹⁶⁴⁵ biomarkers for colorectal cancer. *Microbiology Spectrum*, 10(4), e00013–22.

- 1646 Aziz, Q., & Thompson, D. G. (1998). Brain-gut axis in health and disease. *Gastroenterology*, 114(3),
1647 559–578.
- 1648 Bai, X., Wei, H., Liu, W., Coker, O. O., Gou, H., Liu, C., . . . others (2022). Cigarette smoke promotes
1649 colorectal cancer through modulation of gut microbiota and related metabolites. *Gut*, 71(12),
1650 2439–2450.
- 1651 Baldelli, V., Scaldaferrri, F., Putignani, L., & Del Chierico, F. (2021). The role of enterobacteriaceae in
1652 gut microbiota dysbiosis in inflammatory bowel diseases. *Microorganisms*, 9(4), 697.
- 1653 Bardou, M., Rouland, A., Martel, M., Loffroy, R., Barkun, A. N., & Chapelle, N. (2022). Obesity and
1654 colorectal cancer. *Alimentary Pharmacology & Therapeutics*, 56(3), 407–418.
- 1655 Barlow, G. M., Yu, A., & Mathur, R. (2015). Role of the gut microbiome in obesity and diabetes mellitus.
1656 *Nutrition in clinical practice*, 30(6), 787–797.
- 1657 Basavaprabhu, H., Sonu, K., & Prabha, R. (2020). Mechanistic insights into the action of probiotics
1658 against bacterial vaginosis and its mediated preterm birth: An overview. *Microbial pathogenesis*,
1659 141, 104029.
- 1660 Belstrøm, D., Constancias, F., Drautz-Moses, D. I., Schuster, S. C., Veleba, M., Mahé, F., & Givkov, M.
1661 (2021). Periodontitis associates with species-specific gene expression of the oral microbiota. *npj
1662 Biofilms and Microbiomes*, 7(1), 76.
- 1663 Bender, C., Zipeto, D., Bidoia, C., Costantini, S., Zamò, A., Menestrina, F., & Bertazzoni, U. (2009).
1664 Analysis of colorectal cancers for human cytomegalovirus presence. *Infectious agents and cancer*,
1665 4, 1–6.
- 1666 Benej, M., Hoyd, R., Kreamer, M., Wheeler, C. E., Grencewicz, D. J., Choueiry, F., . . . others (2024). The
1667 tumor microbiome reacts to hypoxia and can influence response to radiation treatment in colorectal
1668 cancer. *Cancer research communications*, 4(7), 1690–1701.
- 1669 Berger, W. H., & Parker, F. L. (1970). Diversity of planktonic foraminifera in deep-sea sediments.
1670 *Science*, 168(3937), 1345–1347.
- 1671 Berghella, V. (2012). Universal cervical length screening for prediction and prevention of preterm birth.
1672 *Obstetrical & gynecological survey*, 67(10), 653–657.
- 1673 Bernardo, G., Le Noci, V., Ottaviano, E., De Cecco, L., Camisaschi, C., Guglielmetti, S., . . . others (2023).
1674 Reduction of staphylococcus epidermidis in the mammary tumor microbiota induces antitumor
1675 immunity and decreases breast cancer aggressiveness. *Cancer Letters*, 555, 216041.
- 1676 Blencowe, H., Cousens, S., Oestergaard, M. Z., Chou, D., Moller, A.-B., Narwal, R., . . . others (2012).
1677 National, regional, and worldwide estimates of preterm birth rates in the year 2010 with time trends
1678 since 1990 for selected countries: a systematic analysis and implications. *The lancet*, 379(9832),
1679 2162–2172.
- 1680 Boland, C. R., & Goel, A. (2010). Microsatellite instability in colorectal cancer. *Gastroenterology*,
1681 138(6), 2073–2087.
- 1682 Boleij, A., Hechenbleikner, E. M., Goodwin, A. C., Badani, R., Stein, E. M., Lazarev, M. G., . . . others
1683 (2015). The bacteroides fragilis toxin gene is prevalent in the colon mucosa of colorectal cancer
1684 patients. *Clinical Infectious Diseases*, 60(2), 208–215.

- 1685 Bolstad, A., Jensen, H. B., & Bakken, V. (1996). Taxonomy, biology, and periodontal aspects of
1686 *fusobacterium nucleatum*. *Clinical microbiology reviews*, 9(1), 55–71.
- 1687 Bolyen, E., Rideout, J. R., Dillon, M. R., Bokulich, N. A., Abnet, C. C., Al-Ghalith, G. A., ... others
1688 (2019). Reproducible, interactive, scalable and extensible microbiome data science using qiime 2.
1689 *Nature biotechnology*, 37(8), 852–857.
- 1690 Bombin, A., Yan, S., Bombin, S., Mosley, J. D., & Ferguson, J. F. (2022). Obesity influences composition
1691 of salivary and fecal microbiota and impacts the interactions between bacterial taxa. *Physiological
1692 reports*, 10(7), e15254.
- 1693 Bonnet, M., Buc, E., Sauvanet, P., Darcha, C., Dubois, D., Pereira, B., ... Darfeuille-Michaud, A. (2014).
1694 Colonization of the human gut by *e. coli* and colorectal cancer risk. *Clinical Cancer Research*,
1695 20(4), 859–867.
- 1696 Breiman, L. (2001). Random forests. *Machine learning*, 45, 5–32.
- 1697 Brennan, C. A., & Garrett, W. S. (2019). *Fusobacterium nucleatum*—symbiont, opportunist and
1698 *oncobacterium*. *Nature Reviews Microbiology*, 17(3), 156–166.
- 1699 Broom, L. J., & Kogut, M. H. (2018). The role of the gut microbiome in shaping the immune system of
1700 chickens. *Veterinary immunology and immunopathology*, 204, 44–51.
- 1701 Bryll, R., Gutierrez-Osuna, R., & Quek, F. (2003). Attribute bagging: improving accuracy of classifier
1702 ensembles by using random feature subsets. *Pattern recognition*, 36(6), 1291–1302.
- 1703 Bullman, S., Pedamallu, C. S., Sicinska, E., Clancy, T. E., Zhang, X., Cai, D., ... others (2017). Analysis
1704 of *fusobacterium* persistence and antibiotic response in colorectal cancer. *Science*, 358(6369),
1705 1443–1448.
- 1706 Burt, R. W., Leppert, M. F., Slattery, M. L., Samowitz, W. S., Spirio, L. N., Kerber, R. A., ... others
1707 (2004). Genetic testing and phenotype in a large kindred with attenuated familial adenomatous
1708 polyposis. *Gastroenterology*, 127(2), 444–451.
- 1709 Cai, Y., Li, Y., Xiong, Y., Geng, X., Kang, Y., & Yang, Y. (2024). Diabetic foot exacerbates gut
1710 mycobiome dysbiosis in adult patients with type 2 diabetes mellitus: revealing diagnostic markers.
1711 *Nutrition & Diabetes*, 14(1), 71.
- 1712 Callahan, B. J., McMurdie, P. J., Rosen, M. J., Han, A. W., Johnson, A. J. A., & Holmes, S. P. (2016).
1713 Dada2: High-resolution sample inference from illumina amplicon data. *Nature methods*, 13(7),
1714 581–583.
- 1715 Canakci, V., & Canakci, C. F. (2007). Pain levels in patients during periodontal probing and mechanical
1716 non-surgical therapy. *Clinical oral investigations*, 11, 377–383.
- 1717 Cappellato, M., Baruzzo, G., & Di Camillo, B. (2022). Investigating differential abundance methods in
1718 microbiome data: A benchmark study. *PLoS computational biology*, 18(9), e1010467.
- 1719 Castaner, O., Goday, A., Park, Y.-M., Lee, S.-H., Magkos, F., Shiow, S.-A. T. E., & Schröder, H. (2018).
1720 The gut microbiome profile in obesity: a systematic review. *International journal of endocrinology*,
1721 2018(1), 4095789.
- 1722 Center, M. M., Jemal, A., Smith, R. A., & Ward, E. (2009). Worldwide variations in colorectal cancer.
1723 *CA: a cancer journal for clinicians*, 59(6), 366–378.

- 1724 Centor, R. M. (1991). Signal detectability: the use of roc curves and their analyses. *Medical decision
making*, 11(2), 102–106.
- 1725
- 1726 Cerqueira, F. M., Photenhauer, A. L., Pollet, R. M., Brown, H. A., & Koropatkin, N. M. (2020). Starch
1727 digestion by gut bacteria: crowdsourcing for carbs. *Trends in Microbiology*, 28(2), 95–108.
- 1728 Champagne, C., McNairn, H., Daneshfar, B., & Shang, J. (2014). A bootstrap method for assessing
1729 classification accuracy and confidence for agricultural land use mapping in canada. *International
1730 Journal of Applied Earth Observation and Geoinformation*, 29, 44–52.
- 1731 Chao, A. (1984). Nonparametric estimation of the number of classes in a population. *Scandinavian
1732 Journal of statistics*, 265–270.
- 1733 Chao, A., & Lee, S.-M. (1992). Estimating the number of classes via sample coverage. *Journal of the
1734 American statistical Association*, 87(417), 210–217.
- 1735 Chapple, I. L., Mealey, B. L., Van Dyke, T. E., Bartold, P. M., Dommisch, H., Eickholz, P., ... others
1736 (2018). Periodontal health and gingival diseases and conditions on an intact and a reduced
1737 periodontium: Consensus report of workgroup 1 of the 2017 world workshop on the classification
1738 of periodontal and peri-implant diseases and conditions. *Journal of periodontology*, 89, S74–S84.
- 1739 Chen, T., Marsh, P., & Al-Hebshi, N. (2022). Smdi: an index for measuring subgingival microbial
1740 dysbiosis. *Journal of dental research*, 101(3), 331–338.
- 1741 Chen, T., Yu, W.-H., Izard, J., Baranova, O. V., Lakshmanan, A., & Dewhirst, F. E. (2010). The human
1742 oral microbiome database: a web accessible resource for investigating oral microbe taxonomic and
1743 genomic information. *Database*, 2010.
- 1744 Chen, X., D’Souza, R., & Hong, S.-T. (2013). The role of gut microbiota in the gut-brain axis: current
1745 challenges and perspectives. *Protein & cell*, 4, 403–414.
- 1746 Chen, X., Jansen, L., Guo, F., Hoffmeister, M., Chang-Claude, J., & Brenner, H. (2021). Smoking,
1747 genetic predisposition, and colorectal cancer risk. *Clinical and translational gastroenterology*,
1748 12(3), e00317.
- 1749 Chen, X., Li, H., Guo, F., Hoffmeister, M., & Brenner, H. (2022). Alcohol consumption, polygenic risk
1750 score, and early-and late-onset colorectal cancer risk. *EClinicalMedicine*, 49.
- 1751 Chew, R. J. J., Tan, K. S., Chen, T., Al-Hebshi, N. N., & Goh, C. E. (2024). Quantifying periodontitis-
1752 associated oral dysbiosis in tongue and saliva microbiomes—an integrated data analysis. *Journal
1753 of Periodontology*.
- 1754 Čižmárová, B., Tomečková, V., Hubková, B., Hurajtová, A., Ohlasová, J., & Birková, A. (2022). Salivary
1755 redox homeostasis in human health and disease. *International Journal of Molecular Sciences*,
1756 23(17), 10076.
- 1757 Conde-Pérez, K., Aja-Macaya, P., Buetas, E., Trigo-Tasende, N., Nasser-Ali, M., Rumbo-Feal, S., ...
1758 others (2024). The multispecies microbial cluster of fusobacterium, parvimonas, bacteroides and
1759 faecalibacterium as a precision biomarker for colorectal cancer diagnosis. *Molecular Oncology*,
1760 18(5), 1093–1122.
- 1761 Cuervo, S. I., Cortés, J. A., Sánchez, R., Rodríguez, J. Y., Silva, E., Tibavizco, D., & Arroyo, P. (2010).
1762 Risk factors for mortality caused by staphylococcus aureus bacteremia in cancer patients. *Enfer-*

- 1763 *medades infecciosas y microbiologia clinica*, 28(6), 349–354.
- 1764 Cullin, N., Antunes, C. A., Straussman, R., Stein-Thoeringer, C. K., & Elinav, E. (2021). Microbiome
1765 and cancer. *Cancer Cell*, 39(10), 1317–1341.
- 1766 Curtius, K., Wright, N. A., & Graham, T. A. (2018). An evolutionary perspective on field cancerization.
1767 *Nature Reviews Cancer*, 18(1), 19–32.
- 1768 Dabke, K., Hendrick, G., Devkota, S., et al. (2019). The gut microbiome and metabolic syndrome. *The
1769 Journal of clinical investigation*, 129(10), 4050–4057.
- 1770 De Flora, S., & Bonanni, P. (2011). The prevention of infection-associated cancers. *Carcinogenesis*,
1771 32(6), 787–795.
- 1772 DeSantis, T. Z., Hugenholtz, P., Larsen, N., Rojas, M., Brodie, E. L., Keller, K., ... Andersen, G. L.
1773 (2006). Greengenes, a chimera-checked 16s rrna gene database and workbench compatible with
1774 arb. *Applied and environmental microbiology*, 72(7), 5069–5072.
- 1775 Ding, R., Lian, S. B., Tam, Y. C., & Oh, C. C. (2024). The cutaneous microbiome in skin cancer—a
1776 systematic review. *JDDG: Journal der Deutschen Dermatologischen Gesellschaft*, 22(2), 177–184.
- 1777 Dong, K., Wang, J., Tang, F., Liu, Y., & Gao, L. (2024). A polysaccharide with a triple helix structure
1778 from agaricus bisporus: Characterization and anti-colon cancer activity. *International Journal of
1779 Biological Macromolecules*, 281, 136521.
- 1780 Doyle, R., Alber, D., Jones, H., Harris, K., Fitzgerald, F., Peebles, D., & Klein, N. (2014). Term and
1781 preterm labour are associated with distinct microbial community structures in placental membranes
1782 which are independent of mode of delivery. *Placenta*, 35(12), 1099–1101.
- 1783 Duy, T. N., Le Huy, H., Thanh, Q. Đ., Thi, H. N., Minh, H. N. T., Dang, M. N., ... Tat, T. N. (2024).
1784 Association between bacteroides fragilis and fusobacterium nucleatum infection and colorectal
1785 cancer in vietnamese patients. *Anaerobe*, 88, 102880.
- 1786 El-Deeb, N. M., Ibrahim, O. M., Mohamed, M. A., Farag, M. M., Farrag, A. A., & El-Aassar, M.
1787 (2022). Alginate/κ-carrageenan oral microcapsules loaded with agaricus bisporus polysaccharides
1788 mh751906 for natural killer cells mediated colon cancer immunotherapy. *International Journal of
1789 Biological Macromolecules*, 205, 385–395.
- 1790 Fahmy, C. A., Gamal-Eldeen, A. M., El-Hussieny, E. A., Raafat, B. M., Mehanna, N. S., Talaat, R. M., &
1791 Shaaban, M. T. (2019). Bifidobacterium longum suppresses murine colorectal cancer through the
1792 modulation of oncomirs and tumor suppressor mirnas. *Nutrition and cancer*, 71(4), 688–700.
- 1793 Faith, D. P. (1992). Conservation evaluation and phylogenetic diversity. *Biological conservation*, 61(1),
1794 1–10.
- 1795 Fettweis, J. M., Serrano, M. G., Brooks, J. P., Edwards, D. J., Girerd, P. H., Parikh, H. I., ... others
1796 (2019). The vaginal microbiome and preterm birth. *Nature medicine*, 25(6), 1012–1021.
- 1797 Fisher, R. A., Corbet, A. S., & Williams, C. B. (1943). The relation between the number of species and
1798 the number of individuals in a random sample of an animal population. *The Journal of Animal
1799 Ecology*, 42–58.
- 1800 Flanagan, L., Schmid, J., Ebert, M., Soucek, P., Kunicka, T., Liska, V., ... others (2014). Fusobacterium
1801 nucleatum associates with stages of colorectal neoplasia development, colorectal cancer and disease

- 1802 outcome. *European journal of clinical microbiology & infectious diseases*, 33, 1381–1390.
- 1803 Fortenberry, J. D. (2013). The uses of race and ethnicity in human microbiome research. *Trends in*
1804 *microbiology*, 21(4), 165–166.
- 1805 Francescone, R., Hou, V., & Grivennikov, S. I. (2014). Microbiome, inflammation, and cancer. *The*
1806 *Cancer Journal*, 20(3), 181–189.
- 1807 Friedman, J. H. (2002). Stochastic gradient boosting. *Computational statistics & data analysis*, 38(4),
1808 367–378.
- 1809 Fushiki, T. (2011). Estimation of prediction error by using k-fold cross-validation. *Statistics and*
1810 *Computing*, 21, 137–146.
- 1811 Galeano Niño, J. L., Wu, H., LaCourse, K. D., Kempchinsky, A. G., Baryiames, A., Barber, B., ... others
1812 (2022). Effect of the intratumoral microbiota on spatial and cellular heterogeneity in cancer. *Nature*,
1813 611(7937), 810–817.
- 1814 Gambin, D. J., Vitali, F. C., De Carli, J. P., Mazzon, R. R., Gomes, B. P., Duque, T. M., & Trentin, M. S.
1815 (2021). Prevalence of red and orange microbial complexes in endodontic-periodontal lesions: a
1816 systematic review and meta-analysis. *Clinical Oral Investigations*, 1–14.
- 1817 Gao, J., Yin, J., Xu, K., Li, T., & Yin, Y. (2019). What is the impact of diet on nutritional diarrhea
1818 associated with gut microbiota in weaning piglets: a system review. *BioMed research international*,
1819 2019(1), 6916189.
- 1820 Gao, R., Kong, C., Li, H., Huang, L., Qu, X., Qin, N., & Qin, H. (2017). Dysbiosis signature of mycobiota
1821 in colon polyp and colorectal cancer. *European Journal of Clinical Microbiology & Infectious*
1822 *Diseases*, 36, 2457–2468.
- 1823 Gao, S., Liu, Y., Duan, X., Liu, K., Mohammed, M., Gu, Z., ... others (2021). *Porphyromonas gingivalis*
1824 infection exacerbates oesophageal cancer and promotes resistance to neoadjuvant chemotherapy.
1825 *British Journal of Cancer*, 125(3), 433–444.
- 1826 Geurts, P., Ernst, D., & Wehenkel, L. (2006). Extremely randomized trees. *Machine learning*, 63, 3–42.
- 1827 Ghanavati, R., Akbari, A., Mohammadi, F., Asadollahi, P., Javadi, A., Talebi, M., & Rohani, M. (2020).
1828 Lactobacillus species inhibitory effect on colorectal cancer progression through modulating the
1829 wnt/β-catenin signaling pathway. *Molecular and Cellular Biochemistry*, 470, 1–13.
- 1830 Ghojogh, B., & Crowley, M. (2019). The theory behind overfitting, cross validation, regularization,
1831 bagging, and boosting: tutorial. *arXiv preprint arXiv:1905.12787*.
- 1832 Ghorbani, E., Avan, A., Ryzhikov, M., Ferns, G., Khazaei, M., & Soleimanpour, S. (2022). Role of
1833 lactobacillus strains in the management of colorectal cancer: An overview of recent advances.
1834 *Nutrition*, 103, 111828.
- 1835 Gilbert, J. A., Blaser, M. J., Caporaso, J. G., Jansson, J. K., Lynch, S. V., & Knight, R. (2018). Current
1836 understanding of the human microbiome. *Nature medicine*, 24(4), 392–400.
- 1837 Gini, C. (1912). Variabilità e mutabilità (variability and mutability). *Tipografia di Paolo Cuppini,*
1838 *Bologna, Italy*, 156.
- 1839 Goldenberg, R. L., Culhane, J. F., Iams, J. D., & Romero, R. (2008). Epidemiology and causes of preterm
1840 birth. *The lancet*, 371(9606), 75–84.

- 1841 Gonçalves, L., Subtil, A., Oliveira, M. R., & de Zea Bermudez, P. (2014). Roc curve estimation: An
1842 overview. *REVSTAT-Statistical journal*, 12(1), 1–20.
- 1843 Good, I. J. (1953). The population frequencies of species and the estimation of population parameters.
1844 *Biometrika*, 40(3-4), 237–264.
- 1845 Goodyear, M. D., Krleza-Jeric, K., & Lemmens, T. (2007). *The declaration of helsinki* (Vol. 335) (No.
1846 7621). British Medical Journal Publishing Group.
- 1847 Haffajee, A., Teles, R., & Socransky, S. (2006). Association of eubacterium nodatum and treponema
1848 denticola with human periodontitis lesions. *Oral microbiology and immunology*, 21(5), 269–282.
- 1849 Hajishengallis, G. (2015). Periodontitis: from microbial immune subversion to systemic inflammation.
1850 *Nature reviews immunology*, 15(1), 30–44.
- 1851 Hamjane, N., Mechita, M. B., Nourouti, N. G., & Barakat, A. (2024). Gut microbiota dysbiosis-associated
1852 obesity and its involvement in cardiovascular diseases and type 2 diabetes. a systematic review.
1853 *Microvascular Research*, 151, 104601.
- 1854 Hamming, R. W. (1950). Error detecting and error correcting codes. *The Bell system technical journal*,
1855 29(2), 147–160.
- 1856 Hampel, H., Frankel, W. L., Martin, E., Arnold, M., Khanduja, K., Kuebler, P., ... others (2008).
1857 Feasibility of screening for lynch syndrome among patients with colorectal cancer. *Journal of
1858 Clinical Oncology*, 26(35), 5783–5788.
- 1859 Han, Y. W. (2015). Fusobacterium nucleatum: a commensal-turned pathogen. *Current opinion in
1860 microbiology*, 23, 141–147.
- 1861 Han, Y. W., & Wang, X. (2013). Mobile microbiome: oral bacteria in extra-oral infections and
1862 inflammation. *Journal of dental research*, 92(6), 485–491.
- 1863 Hand, D. J. (2012). Assessing the performance of classification methods. *International Statistical Review*,
1864 80(3), 400–414.
- 1865 Harkins, L., Volk, A. L., Samanta, M., Mikolaenko, I., Britt, W. J., Bland, K. I., & Cobbs, C. S. (2002).
1866 Specific localisation of human cytomegalovirus nucleic acids and proteins in human colorectal
1867 cancer. *The Lancet*, 360(9345), 1557–1563.
- 1868 Hartstra, A. V., Bouter, K. E., Bäckhed, F., & Nieuwdorp, M. (2015). Insights into the role of the
1869 microbiome in obesity and type 2 diabetes. *Diabetes care*, 38(1), 159–165.
- 1870 Hasan, R., Bose, S., Roy, R., Paul, D., Rawat, S., Nilwe, P., ... Choudhury, S. (2022). Tumor tissue-
1871 specific bacterial biomarker panel for colorectal cancer: Bacteroides massiliensis, alistipes species,
1872 alistipes onderdonkii, bifidobacterium pseudocatenulatum, corynebacterium appendicis. *Archives
1873 of microbiology*, 204(6), 348.
- 1874 Hashemi Goradel, N., Heidarzadeh, S., Jahangiri, S., Farhood, B., Mortezaee, K., Khanlarkhani, N., &
1875 Negahdari, B. (2019). Fusobacterium nucleatum and colorectal cancer: A mechanistic overview.
1876 *Journal of Cellular Physiology*, 234(3), 2337–2344.
- 1877 Hassan, Z., Mustafa, S., Rahim, R. A., & Isa, N. M. (2016). Anti-breast cancer effects of live, heat-killed
1878 and cytoplasmic fractions of enterococcus faecalis and staphylococcus hominis isolated from
1879 human breast milk. *In Vitro Cellular & Developmental Biology-Animal*, 52, 337–348.

- 1880 Hegazy, A. A., Abu-Hussien, S. H., Elsenosy, N. K., El-Sayed, S. M., & Abo El-Naga, M. Y. (2024).
1881 Optimization, characterization and biosafety of carotenoids produced from whey using *micrococcus*
1882 *luteus*. *BMC biotechnology*, 24(1), 74.
- 1883 Heip, C. (1974). A new index measuring evenness. *Journal of the Marine Biological Association of the*
1884 *United Kingdom*, 54(3), 555–557.
- 1885 Helmink, B. A., Khan, M. W., Hermann, A., Gopalakrishnan, V., & Wargo, J. A. (2019). The microbiome,
1886 cancer, and cancer therapy. *Nature medicine*, 25(3), 377–388.
- 1887 Hexun, Z., Miyake, T., Maekawa, T., Mori, H., Yasukawa, D., Ohno, M., ... Tani, M. (2023). High
1888 abundance of lachnospiraceae in the human gut microbiome is related to high immunoscores in
1889 advanced colorectal cancer. *Cancer Immunology, Immunotherapy*, 72(2), 315–326.
- 1890 Hill, M. O. (1973). Diversity and evenness: a unifying notation and its consequences. *Ecology*, 54(2),
1891 427–432.
- 1892 Hiranmayi, K. V., Sirisha, K., Rao, M. R., & Sudhakar, P. (2017). Novel pathogens in periodontal
1893 microbiology. *Journal of Pharmacy and Bioallied Sciences*, 9(3), 155–163.
- 1894 Honda, K., & Littman, D. R. (2012). The microbiome in infectious disease and inflammation. *Annual*
1895 *review of immunology*, 30(1), 759–795.
- 1896 Honest, H., Forbes, C., Durée, K., Norman, G., Duffy, S., Tsourapas, A., ... others (2009). Screening to
1897 prevent spontaneous preterm birth: systematic reviews of accuracy and effectiveness literature with
1898 economic modelling. *Health Technol Assess*, 13(43), 1–627.
- 1899 Hong, Y. M., Lee, J., Cho, D. H., Jeon, J. H., Kang, J., Kim, M.-G., ... J. K. (2023). Predicting preterm
1900 birth using machine learning techniques in oral microbiome. *Scientific Reports*, 13(1), 21105.
- 1901 Hossin, M., & Sulaiman, M. N. (2015). A review on evaluation metrics for data classification evaluations.
1902 *International journal of data mining & knowledge management process*, 5(2), 1.
- 1903 Huang, R.-Y., Lin, C.-D., Lee, M.-S., Yeh, C.-L., Shen, E.-C., Chiang, C.-Y., ... Fu, E. (2007). Mandibular
1904 disto-lingual root: a consideration in periodontal therapy. *Journal of periodontology*, 78(8), 1485–
1905 1490.
- 1906 Huo, R.-X., Wang, Y.-J., Hou, S.-B., Wang, W., Zhang, C.-Z., & Wan, X.-H. (2022). Gut mucosal
1907 microbiota profiles linked to colorectal cancer recurrence. *World journal of gastroenterology*,
1908 28(18), 1946.
- 1909 Iams, J. D., & Berghella, V. (2010). Care for women with prior preterm birth. *American journal of*
1910 *obstetrics and gynecology*, 203(2), 89–100.
- 1911 Ide, M., & Papapanou, P. N. (2013). Epidemiology of association between maternal periodontal
1912 disease and adverse pregnancy outcomes—systematic review. *Journal of clinical periodontology*,
1913 40, S181–S194.
- 1914 Iniesta, M., Chamorro, C., Ambrosio, N., Marín, M. J., Sanz, M., & Herrera, D. (2023). Subgingival
1915 microbiome in periodontal health, gingivitis and different stages of periodontitis. *Journal of*
1916 *Clinical Periodontology*, 50(7), 905–920.
- 1917 Inra, J. A., Steyerberg, E. W., Grover, S., McFarland, A., Syngal, S., & Kastrinos, F. (2015). Racial
1918 variation in frequency and phenotypes of apc and mutyh mutations in 6,169 individuals undergoing

- 1919 genetic testing. *Genetics in Medicine*, 17(10), 815–821.
- 1920 Jaccard, P. (1908). Nouvelles recherches sur la distribution florale. *Bull. Soc. Vaud. Sci. Nat.*, 44, 1921 223–270.
- 1922 Janda, J. M., & Abbott, S. L. (2007). 16s rrna gene sequencing for bacterial identification in the diagnostic 1923 laboratory: pluses, perils, and pitfalls. *Journal of clinical microbiology*, 45(9), 2761–2764.
- 1924 Jiang, W., & Simon, R. (2007). A comparison of bootstrap methods and an adjusted bootstrap approach 1925 for estimating the prediction error in microarray classification. *Statistics in medicine*, 26(29), 1926 5320–5334.
- 1927 John, G. K., & Mullin, G. E. (2016). The gut microbiome and obesity. *Current oncology reports*, 18, 1928 1–7.
- 1929 Johnson, J. S., Spakowicz, D. J., Hong, B.-Y., Petersen, L. M., Demkowicz, P., Chen, L., . . . others (2019). 1930 Evaluation of 16s rrna gene sequencing for species and strain-level microbiome analysis. *Nature 1931 communications*, 10(1), 5029.
- 1932 Joo, J. E., Chu, Y. L., Georgeson, P., Walker, R., Mahmood, K., Clendenning, M., . . . others (2024). 1933 Intratumoral presence of the genotoxic gut bacteria pks+ e. coli, enterotoxigenic bacteroides fragilis, 1934 and fusobacterium nucleatum and their association with clinicopathological and molecular features 1935 of colorectal cancer. *British Journal of Cancer*, 130(5), 728–740.
- 1936 Jorth, P., Turner, K. H., Gumus, P., Nizam, N., Buduneli, N., & Whiteley, M. (2014). Metatranscriptomics 1937 of the human oral microbiome during health and disease. *MBio*, 5(2), 10–1128.
- 1938 Joscelyn, J., & Kasper, L. H. (2014). Digesting the emerging role for the gut microbiome in central 1939 nervous system demyelination. *Multiple Sclerosis Journal*, 20(12), 1553–1559.
- 1940 Kang, Y., Kang, X., Yang, H., Liu, H., Yang, X., Liu, Q., . . . others (2022). Lactobacillus acidophilus ame- 1941 liorates obesity in mice through modulation of gut microbiota dysbiosis and intestinal permeability. 1942 *Pharmacological research*, 175, 106020.
- 1943 Karched, M., Bhardwaj, R. G., Qudeimat, M., Al-Khabbaz, A., & Ellepola, A. (2022). Proteomic analysis 1944 of the periodontal pathogen prevotella intermedia secretomes in biofilm and planktonic lifestyles. 1945 *Scientific Reports*, 12(1), 5636.
- 1946 Katz, J., Chegini, N., Shiverick, K., & Lamont, R. (2009). Localization of p. gingivalis in preterm delivery 1947 placenta. *Journal of dental research*, 88(6), 575–578.
- 1948 Kau, A. L., Ahern, P. P., Griffin, N. W., Goodman, A. L., & Gordon, J. I. (2011). Human nutrition, the 1949 gut microbiome and the immune system. *Nature*, 474(7351), 327–336.
- 1950 Kelly, B. J., Gross, R., Bittinger, K., Sherrill-Mix, S., Lewis, J. D., Collman, R. G., . . . Li, H. (2015). 1951 Power and sample-size estimation for microbiome studies using pairwise distances and permanova. 1952 *Bioinformatics*, 31(15), 2461–2468.
- 1953 Kennedy, J., Alexander, P., Taillie, L. S., & Jaacks, L. M. (2024). Estimated effects of reductions in 1954 processed meat consumption and unprocessed red meat consumption on occurrences of type 2 1955 diabetes, cardiovascular disease, colorectal cancer, and mortality in the usa: a microsimulation 1956 study. *The Lancet Planetary Health*, 8(7), e441–e451.
- 1957 Kepp, O., Zitvogel, L., & Kroemer, G. (2023). *Prevention and treatment of cancers by tumor antigen-*

- 1958 *expressing staphylococcus epidermidis* (Vol. 12) (No. 1). Taylor & Francis.
- 1959 Kim, B.-R., Shin, J., Guevarra, R. B., Lee, J. H., Kim, D. W., Seol, K.-H., ... Isaacson, R. E. (2017).
1960 Deciphering diversity indices for a better understanding of microbial communities. *Journal of*
1961 *Microbiology and Biotechnology*, 27(12), 2089–2093.
- 1962 Kim, C. H. (2018). Immune regulation by microbiome metabolites. *Immunology*, 154(2), 220–229.
- 1963 Kim, E.-H., Kim, S., Kim, H.-J., Jeong, H.-o., Lee, J., Jang, J., ... others (2020). Prediction of chronic
1964 periodontitis severity using machine learning models based on salivary bacterial copy number.
1965 *Frontiers in Cellular and Infection Microbiology*, 10, 571515.
- 1966 Kim, J.-H. (2009). Estimating classification error rate: Repeated cross-validation, repeated hold-out and
1967 bootstrap. *Computational statistics & data analysis*, 53(11), 3735–3745.
- 1968 Kim, S., Lee, M., Kim, N.-Y., Kwon, Y.-S., Nam, G. S., Lee, K., ... Hwang, I. H. (2024). Oxidative
1969 tryptamine dimers from corynebacterium durum directly target survivin to induce aif-mediated
1970 apoptosis in cancer cells. *Biomedicine & Pharmacotherapy*, 173, 116335.
- 1971 Kinane, D. F., Stathopoulou, P. G., & Papapanou, P. N. (2017). Periodontal diseases. *Nature reviews*
1972 *Disease primers*, 3(1), 1–14.
- 1973 Kindinger, L. M., Bennett, P. R., Lee, Y. S., Marchesi, J. R., Smith, A., Caciato, S., ... MacIntyre,
1974 D. A. (2017). The interaction between vaginal microbiota, cervical length, and vaginal progesterone
1975 treatment for preterm birth risk. *Microbiome*, 5, 1–14.
- 1976 Kogut, M. H., Lee, A., & Santin, E. (2020). Microbiome and pathogen interaction with the immune
1977 system. *Poultry science*, 99(4), 1906–1913.
- 1978 Kostic, A. D., Ojesina, A. I., Pedamallu, C. S., Jung, J., Verhaak, R. G., Getz, G., & Meyerson, M. (2011).
1979 Pathseq: software to identify or discover microbes by deep sequencing of human tissue. *Nature*
1980 *biotechnology*, 29(5), 393–396.
- 1981 Kotsiantis, S. B., Zaharakis, I. D., & Pintelas, P. E. (2006). Machine learning: a review of classification
1982 and combining techniques. *Artificial Intelligence Review*, 26, 159–190.
- 1983 Kuipers, E. J., Grady, W. M., Lieberman, D., Seufferlein, T., Sung, J. J., Boelens, P. G., ... Watanabe, T.
1984 (2015). Colorectal cancer. *Nature reviews. Disease primers*, 1, 15065.
- 1985 Lafaurie, G. I., Neuta, Y., Ríos, R., Pacheco-Montealegre, M., Pianeta, R., Castillo, D. M., ... oth-
1986 ers (2022). Differences in the subgingival microbiome according to stage of periodontitis: A
1987 comparison of two geographic regions. *PLoS one*, 17(8), e0273523.
- 1988 Lamont, R. J., & Jenkinson, H. F. (2000). Subgingival colonization by porphyromonas gingivalis. *Oral*
1989 *Microbiology and Immunology: Mini-review*, 15(6), 341–349.
- 1990 Lamont, R. J., Koo, H., & Hajishengallis, G. (2018). The oral microbiota: dynamic communities and
1991 host interactions. *Nature reviews microbiology*, 16(12), 745–759.
- 1992 Leitich, H., & Kaider, A. (2003). Fetal fibronectin—how useful is it in the prediction of preterm birth?
1993 *BJOG: An International Journal of Obstetrics & Gynaecology*, 110, 66–70.
- 1994 Le Leu, R. K., Hu, Y., Brown, I. L., Woodman, R. J., & Young, G. P. (2010). Synbiotic intervention of
1995 bifidobacterium lactis and resistant starch protects against colorectal cancer development in rats.
1996 *Carcinogenesis*, 31(2), 246–251.

- 1997 León, R., Silva, N., Ovalle, A., Chaparro, A., Ahumada, A., Gajardo, M., ... Gamonal, J. (2007).
1998 Detection of porphyromonas gingivalis in the amniotic fluid in pregnant women with a diagnosis
1999 of threatened premature labor. *Journal of periodontology*, 78(7), 1249–1255.
- 2000 Li, H., & Durbin, R. (2009). Fast and accurate short read alignment with burrows–wheeler transform.
2001 *bioinformatics*, 25(14), 1754–1760.
- 2002 Li, N., Lu, B., Luo, C., Cai, J., Lu, M., Zhang, Y., ... Dai, M. (2021). Incidence, mortality, survival,
2003 risk factor and screening of colorectal cancer: A comparison among china, europe, and northern
2004 america. *Cancer letters*, 522, 255–268.
- 2005 Li, R., Miao, Z., Liu, Y., Chen, X., Wang, H., Su, J., & Chen, J. (2024). The brain–gut–bone axis in
2006 neurodegenerative diseases: insights, challenges, and future prospects. *Advanced Science*, 11(38),
2007 2307971.
- 2008 Li, W., & Yang, J. (2025). Investigating the anna karenina principle of the breast microbiome. *BMC*
2009 *microbiology*, 25(1), 1–10.
- 2010 Li, X., Yu, D., Wang, Y., Yuan, H., Ning, X., Rui, B., ... Li, M. (2021). The intestinal dysbiosis of
2011 mothers with gestational diabetes mellitus (gdm) and its impact on the gut microbiota of their
2012 newborns. *Canadian Journal of Infectious Diseases and Medical Microbiology*, 2021(1), 3044534.
- 2013 Li, Y., Qian, F., Cheng, X., Wang, D., Wang, Y., Pan, Y., ... Tian, Y. (2023). Dysbiosis of oral microbiota
2014 and metabolite profiles associated with type 2 diabetes mellitus. *Microbiology spectrum*, 11(1),
2015 e03796–22.
- 2016 Li, Z., Zhuang, H., Wang, G., Wang, H., & Dong, Y. (2021). Prevalence, predictors, and mortality
2017 of bloodstream infections due to methicillin-resistant staphylococcus aureus in patients with
2018 malignancy: systemic review and meta-analysis. *BMC infectious diseases*, 21, 1–10.
- 2019 Lim, J. W., Park, T., Tong, Y. W., & Yu, Z. (2020). The microbiome driving anaerobic digestion and
2020 microbial analysis. In *Advances in bioenergy* (Vol. 5, pp. 1–61). Elsevier.
- 2021 Lin, H., Eggesbø, M., & Peddada, S. D. (2022). Linear and nonlinear correlation estimators unveil
2022 undescribed taxa interactions in microbiome data. *Nature communications*, 13(1), 4946.
- 2023 Lin, H., & Peddada, S. D. (2020). Analysis of compositions of microbiomes with bias correction. *Nature*
2024 *communications*, 11(1), 3514.
- 2025 Lin, H., & Peddada, S. D. (2024). Multigroup analysis of compositions of microbiomes with covariate
2026 adjustments and repeated measures. *Nature Methods*, 21(1), 83–91.
- 2027 Listgarten, M. A. (1986). Pathogenesis of periodontitis. *Journal of clinical periodontology*, 13(5),
2028 418–425.
- 2029 Liu, W., Zhang, X., Xu, H., Li, S., Lau, H. C.-H., Chen, Q., ... others (2021). Microbial commu-
2030 nity heterogeneity within colorectal neoplasia and its correlation with colorectal carcinogenesis.
2031 *Gastroenterology*, 160(7), 2395–2408.
- 2032 Lloyd-Price, J., Abu-Ali, G., & Huttenhower, C. (2016). The healthy human microbiome. *Genome*
2033 *medicine*, 8, 1–11.
- 2034 López-Aladid, R., Fernández-Barat, L., Alcaraz-Serrano, V., Bueno-Freire, L., Vázquez, N., Pastor-
2035 Ibáñez, R., ... Torres, A. (2023). Determining the most accurate 16s rrna hypervariable region for

- 2036 taxonomic identification from respiratory samples. *Scientific reports*, 13(1), 3974.
- 2037 Love, M. I., Huber, W., & Anders, S. (2014). Moderated estimation of fold change and dispersion for
2038 rna-seq data with deseq2. *Genome biology*, 15, 1–21.
- 2039 Lu, Y.-T., Hsin, C.-H., Chuang, C.-Y., Huang, C.-C., Su, M.-C., Wen, W.-S., ... others (2024). Mi-
2040 crobial dysbiosis in nasopharyngeal carcinoma: A pilot study on biomarker potential. *Journal of*
2041 *Otolaryngology-Head & Neck Surgery*, 53, 19160216241304365.
- 2042 Ma, Z. S. (2020). Testing the anna karenina principle in human microbiome-associated diseases. *Iscience*,
2043 23(4).
- 2044 Magnúsdóttir, S., & Thiele, I. (2018). Modeling metabolism of the human gut microbiome. *Current*
2045 *opinion in biotechnology*, 51, 90–96.
- 2046 Magurran, A. E. (2021). Measuring biological diversity. *Current Biology*, 31(19), R1174–R1177.
- 2047 Mandic, M., Safizadeh, F., Niedermaier, T., Hoffmeister, M., & Brenner, H. (2023). Association of
2048 overweight, obesity, and recent weight loss with colorectal cancer risk. *JAMA network Open*, 6(4),
2049 e239556–e239556.
- 2050 Mann, H. B., & Whitney, D. R. (1947). On a test of whether one of two random variables is stochastically
2051 larger than the other. *The annals of mathematical statistics*, 50–60.
- 2052 Manolis, A. A., Manolis, T. A., Melita, H., & Manolis, A. S. (2022). Gut microbiota and cardiovascular
2053 disease: symbiosis versus dysbiosis. *Current Medicinal Chemistry*, 29(23), 4050–4077.
- 2054 Martin, C. R., Osadchiy, V., Kalani, A., & Mayer, E. A. (2018). The brain-gut-microbiome axis. *Cellular*
2055 *and molecular gastroenterology and hepatology*, 6(2), 133–148.
- 2056 Maulud, D., & Abdulazeez, A. M. (2020). A review on linear regression comprehensive in machine
2057 learning. *Journal of Applied Science and Technology Trends*, 1(2), 140–147.
- 2058 Mayer, E. A., Tillisch, K., Gupta, A., et al. (2015). Gut/brain axis and the microbiota. *The Journal of*
2059 *clinical investigation*, 125(3), 926–938.
- 2060 Melguizo-Rodríguez, L., Costela-Ruiz, V. J., Manzano-Moreno, F. J., Ruiz, C., & Illescas-Montes, R.
2061 (2020). Salivary biomarkers and their application in the diagnosis and monitoring of the most
2062 common oral pathologies. *International journal of molecular sciences*, 21(14), 5173.
- 2063 Merrill, L. C., & Mangano, K. M. (2023). Racial and ethnic differences in studies of the gut microbiome
2064 and osteoporosis. *Current Osteoporosis Reports*, 21(5), 578–591.
- 2065 Miller, C. S., Ding, X., Dawson III, D. R., & Ebersole, J. L. (2021). Salivary biomarkers for discriminating
2066 periodontitis in the presence of diabetes. *Journal of clinical periodontology*, 48(2), 216–225.
- 2067 Minton, N. P. (2003). Clostridia in cancer therapy. *Nature Reviews Microbiology*, 1(3), 237–242.
- 2068 Mjelle, R., Castro, Í., & Aass, K. R. (2025). The viral landscape in metastatic solid cancers. *Heliyon*.
- 2069 Montalban-Arques, A., Katkeviciute, E., Busenhart, P., Bircher, A., Wirbel, J., Zeller, G., ... others
2070 (2021). Commensal clostridiales strains mediate effective anti-cancer immune response against
2071 solid tumors. *Cell host & microbe*, 29(10), 1573–1588.
- 2072 Morita, T., Yamazaki, Y., Mita, A., Takada, K., Seto, M., Nishinoue, N., ... Maeno, M. (2010). A cohort
2073 study on the association between periodontal disease and the development of metabolic syndrome.
2074 *Journal of periodontology*, 81(4), 512–519.

- 2075 Na, H. S., Kim, S. Y., Han, H., Kim, H.-J., Lee, J.-Y., Lee, J.-H., & Chung, J. (2020). Identification of
2076 potential oral microbial biomarkers for the diagnosis of periodontitis. *Journal of clinical medicine*,
2077 9(5), 1549.
- 2078 Nemoto, T., Shiba, T., Komatsu, K., Watanabe, T., Shimogishi, M., Shibasaki, M., ... others (2021).
2079 Discrimination of bacterial community structures among healthy, gingivitis, and periodontitis
2080 statuses through integrated metatranscriptomic and network analyses. *Msystems*, 6(6), e00886–21.
- 2081 Nesbitt, M. J., Reynolds, M. A., Shiau, H., Choe, K., Simonsick, E. M., & Ferrucci, L. (2010). Association
2082 of periodontitis and metabolic syndrome in the baltimore longitudinal study of aging. *Aging clinical
2083 and experimental research*, 22, 238–242.
- 2084 Network, C. G. A., et al. (2012). Comprehensive molecular characterization of human colon and rectal
2085 cancer. *Nature*, 487(7407), 330.
- 2086 Nibali, L., Sousa, V., Davrandi, M., Spratt, D., Alyahya, Q., Dopico, J., & Donos, N. (2020). Differences
2087 in the periodontal microbiome of successfully treated and persistent aggressive periodontitis.
2088 *Journal of Clinical Periodontology*, 47(8), 980–990.
- 2089 Novaković, J. D., Veljović, A., Ilić, S. S., Papić, Ž., & Tomović, M. (2017). Evaluation of classification
2090 models in machine learning. *Theory and Applications of Mathematics & Computer Science*, 7(1),
2091 39.
- 2092 Obuchowski, N. A., & Bullen, J. A. (2018). Receiver operating characteristic (roc) curves: review of
2093 methods with applications in diagnostic medicine. *Physics in Medicine & Biology*, 63(7), 07TR01.
- 2094 Ochiai, A. (1957). Zoogeographic studies on the soleoid fishes found in japan and its neighbouring
2095 regions. *Bulletin of Japanese Society of Scientific Fisheries*, 22, 526–530.
- 2096 Offenbacher, S., Katz, V., Fertik, G., Collins, J., Boyd, D., Maynor, G., ... Beck, J. (1996). Periodontal
2097 infection as a possible risk factor for preterm low birth weight. *Journal of periodontology*, 67,
2098 1103–1113.
- 2099 Ohmagari, N., Hanna, H., Graviss, L., Hackett, B., Perego, C., Gonzalez, V., ... others (2005). Risk
2100 factors for infections with multidrug-resistant pseudomonas aeruginosa in patients with cancer.
2101 *Cancer*, 104(1), 205–212.
- 2102 Ojesina, A. I., Pedamallu, C. S., Kostic, A., Jung, J., Auclair, D., Lohr, J., ... Meyerson, M. (2013). High
2103 throughput sequencing-based pathogen discovery in multiple myeloma. *Blood*, 122(21), 5322.
- 2104 Omundiagbe, D. A., Veeramani, S., & Sidhu, A. S. (2019). Machine learning classification techniques
2105 for breast cancer diagnosis. In *Iop conference series: materials science and engineering* (Vol. 495,
2106 p. 012033).
- 2107 O'Sullivan, D. E., Sutherland, R. L., Town, S., Chow, K., Fan, J., Forbes, N., ... Brenner, D. R. (2022).
2108 Risk factors for early-onset colorectal cancer: a systematic review and meta-analysis. *Clinical
2109 gastroenterology and hepatology*, 20(6), 1229–1240.
- 2110 Paganini, D., & Zimmermann, M. B. (2017). The effects of iron fortification and supplementation on the
2111 gut microbiome and diarrhea in infants and children: a review. *The American journal of clinical
2112 nutrition*, 106, 1688S–1693S.
- 2113 Pan, A. Y. (2021). Statistical analysis of microbiome data: the challenge of sparsity. *Current Opinion in*

- 2114 *Endocrine and Metabolic Research*, 19, 35–40.
- 2115 Papapanou, P. N., Sanz, M., Buduneli, N., Dietrich, T., Feres, M., Fine, D. H., ... others (2018).
2116 Periodontitis: Consensus report of workgroup 2 of the 2017 world workshop on the classification of
2117 periodontal and peri-implant diseases and conditions. *Journal of periodontology*, 89, S173–S182.
- 2118 Paprocka, P., Durnaś, B., Mańkowska, A., Król, G., Wollny, T., & Bucki, R. (2022). *Pseudomonas*
2119 *aeruginosa* infections in cancer patients. *Pathogens*, 11(6), 679.
- 2120 Parizadeh, M., & Arrieta, M.-C. (2023). The global human gut microbiome: genes, lifestyles, and diet.
2121 *Trends in Molecular Medicine*.
- 2122 Park, J., Park, S. H., Lee, D., Lee, J. E., Lee, D., Na, K. J., ... Im, H.-J. (2024). Detecting cancer microbiota
2123 using unmapped rna reads on spatial transcriptomics. *Cancer Research*, 84(6_Supplement), 4881–
2124 4881.
- 2125 Payne, M. S., Newnham, J. P., Doherty, D. A., Furfaro, L. L., Pendal, N. L., Loh, D. E., & Keelan, J. A.
2126 (2021). A specific bacterial dna signature in the vagina of australian women in midpregnancy
2127 predicts high risk of spontaneous preterm birth (the predict1000 study). *American journal of*
2128 *obstetrics and gynecology*, 224(2), 206–e1.
- 2129 Peirce, J. M., & Alviña, K. (2019). The role of inflammation and the gut microbiome in depression and
2130 anxiety. *Journal of neuroscience research*, 97(10), 1223–1241.
- 2131 Peltomaki, P. (2003). Role of dna mismatch repair defects in the pathogenesis of human cancer. *Journal*
2132 *of clinical oncology*, 21(6), 1174–1179.
- 2133 Pezzino, S., Sofia, M., Greco, L. P., Litrico, G., Filippello, G., Sarvà, I., ... Latteri, S. (2023). Microbiome
2134 dysbiosis: a pathological mechanism at the intersection of obesity and glaucoma. *International*
2135 *Journal of Molecular Sciences*, 24(2), 1166.
- 2136 Pollard, T. J., Johnson, A. E., Raffa, J. D., & Mark, R. G. (2018). tableone: An open source python
2137 package for producing summary statistics for research papers. *JAMIA open*, 1(1), 26–31.
- 2138 Premaraj, T. S., Vella, R., Chung, J., Lin, Q., Hunter, P., Underwood, K., ... Zhou, Y. (2020). Ethnic
2139 variation of oral microbiota in children. *Scientific reports*, 10(1), 14788.
- 2140 Purcell, R. V., Permain, J., & Keenan, J. I. (2022). Enterotoxigenic bacteroides fragilis activates il-8
2141 expression through stat3 in colorectal cancer cells. *Gut Pathogens*, 14(1), 16.
- 2142 Raut, J. R., Schöttker, B., Holleczeck, B., Guo, F., Bhardwaj, M., Miah, K., ... Brenner, H. (2021).
2143 A microrna panel compared to environmental and polygenic scores for colorectal cancer risk
2144 prediction. *Nature Communications*, 12(1), 4811.
- 2145 Rebersek, M. (2021). Gut microbiome and its role in colorectal cancer. *BMC cancer*, 21(1), 1325.
- 2146 Redanz, U., Redanz, S., Treerat, P., Prakasam, S., Lin, L.-J., Merritt, J., & Kreth, J. (2021). Differential
2147 response of oral mucosal and gingival cells to corynebacterium durum, streptococcus sanguinis, and
2148 porphyromonas gingivalis multispecies biofilms. *Frontiers in cellular and infection microbiology*,
2149 11, 686479.
- 2150 Relvas, M., Regueira-Iglesias, A., Balsa-Castro, C., Salazar, F., Pacheco, J., Cabral, C., ... Tomás, I.
2151 (2021). Relationship between dental and periodontal health status and the salivary microbiome:
2152 bacterial diversity, co-occurrence networks and predictive models. *Scientific reports*, 11(1), 929.

- 2153 Renson, A., Jones, H. E., Beghini, F., Segata, N., Zolnik, C. P., Usyk, M., ... others (2019). Sociodemographic variation in the oral microbiome. *Annals of epidemiology*, 35, 73–80.
- 2154
- 2155 Renvert, S., & Persson, G. (2002). A systematic review on the use of residual probing depth, bleeding on probing and furcation status following initial periodontal therapy to predict further attachment and tooth loss. *Journal of clinical periodontology*, 29, 82–89.
- 2156
- 2157
- 2158 Rideout, J. R., Caporaso, G., Bolyen, E., McDonald, D., Baeza, Y. V., Alastuey, J. C., ... Sharma, K. (2018, December). *biocore/scikit-bio: scikit-bio 0.5.5: More compositional methods added*. Zenodo. Retrieved from <https://doi.org/10.5281/zenodo.2254379> doi: 10.5281/zenodo.2254379
- 2159
- 2160
- 2161 Rôças, I. N., Siqueira Jr, J. F., Santos, K. R., Coelho, A. M., & de Janeiro, R. (2001). “red complex”(bacteroides forsythus, porphyromonas gingivalis, and treponema denticola) in endodontic infections: a molecular approach. *Oral Surgery, Oral Medicine, Oral Pathology, Oral Radiology, and Endodontology*, 91(4), 468–471.
- 2162
- 2163
- 2164
- 2165 Romero, R., Dey, S. K., & Fisher, S. J. (2014). Preterm labor: one syndrome, many causes. *Science*, 345(6198), 760–765.
- 2166
- 2167 Romero, R., Hassan, S. S., Gajer, P., Tarca, A. L., Fadrosh, D. W., Nikita, L., ... others (2014). The composition and stability of the vaginal microbiota of normal pregnant women is different from that of non-pregnant women. *Microbiome*, 2, 1–19.
- 2168
- 2169
- 2170 Rosan, B., & Lamont, R. J. (2000). Dental plaque formation. *Microbes and infection*, 2(13), 1599–1607.
- 2171 Rubio, C. A., Lang-Schwarz, C., & Vieth, M. (2022). Further study on field cancerization in the human colon. *Anticancer Research*, 42(12), 5891–5895.
- 2172
- 2173 Schwabe, R. F., & Jobin, C. (2013). The microbiome and cancer. *Nature Reviews Cancer*, 13(11), 800–812.
- 2174
- 2175 Scott, N., Whittle, E., Jeraldo, P., & Chia, N. (2022). A systemic review of the role of enterotoxic bacteroides fragilis in colorectal cancer. *Neoplasia*, 29, 100797.
- 2176
- 2177 Segata, N., Izard, J., Waldron, L., Gevers, D., Miropolsky, L., Garrett, W. S., & Huttenhower, C. (2011). Metagenomic biomarker discovery and explanation. *Genome biology*, 12, 1–18.
- 2178
- 2179 Sen, P. C., Hajra, M., & Ghosh, M. (2020). Supervised classification algorithms in machine learning: A survey and review. In *Emerging technology in modelling and graphics: Proceedings of iem graph 2018* (pp. 99–111).
- 2180
- 2181
- 2182 Sepich-Poore, G. D., Zitvogel, L., Straussman, R., Hasty, J., Wargo, J. A., & Knight, R. (2021). The microbiome and human cancer. *Science*, 371(6536), eabc4552.
- 2183
- 2184 Sharma, S., & Tripathi, P. (2019). Gut microbiome and type 2 diabetes: where we are and where to go? *The Journal of nutritional biochemistry*, 63, 101–108.
- 2185
- 2186 Shen, X., Zhang, B., Hu, X., Li, J., Wu, M., Yan, C., ... Li, Y. (2022). Neisseria sicca and corynebacterium matruchotii inhibited oral squamous cell carcinomas by regulating genome stability. *Bioengineered*, 13(6), 14094–14106.
- 2187
- 2188
- 2189 Shi, N., Li, N., Duan, X., & Niu, H. (2017). Interaction between the gut microbiome and mucosal immune system. *Military Medical Research*, 4, 1–7.
- 2190
- 2191 Simpson, E. (1949). Measurement of diversity. *Nature*, 163.

- 2192 Sokal, R. R., & Sneath, P. H. (1963). Principles of numerical taxonomy.
- 2193 Song, M., Chan, A. T., & Sun, J. (2020). Influence of the gut microbiome, diet, and environment on risk
2194 of colorectal cancer. *Gastroenterology*, 158(2), 322–340.
- 2195 Söreide, K., Janssen, E., Söiland, H., Körner, H., & Baak, J. (2006). Microsatellite instability in colorectal
2196 cancer. *Journal of British Surgery*, 93(4), 395–406.
- 2197 Sørensen, T. (1948). A method of establishing groups of equal amplitude in plant sociology based on
2198 similarity of species content and its application to analyses of the vegetation on danish commons.
2199 *Biologiske skrifter*, 5, 1–34.
- 2200 Sotiriadis, A., Papatheodorou, S., Kavvadias, A., & Makrydimas, G. (2010). Transvaginal cervical
2201 length measurement for prediction of preterm birth in women with threatened preterm labor: a
2202 meta-analysis. *Ultrasound in Obstetrics and Gynecology: The Official Journal of the International
2203 Society of Ultrasound in Obstetrics and Gynecology*, 35(1), 54–64.
- 2204 Spss, I., et al. (2011). Ibm spss statistics for windows, version 20.0. *New York: IBM Corp*, 440, 394.
- 2205 Stafford, G., Roy, S., Honma, K., & Sharma, A. (2012). Sialic acid, periodontal pathogens and tannerella
2206 forsythia: stick around and enjoy the feast! *Molecular Oral Microbiology*, 27(1), 11–22.
- 2207 Stout, M. J., Conlon, B., Landeau, M., Lee, I., Bower, C., Zhao, Q., ... Mysorekar, I. U. (2013).
2208 Identification of intracellular bacteria in the basal plate of the human placenta in term and preterm
2209 gestations. *American journal of obstetrics and gynecology*, 208(3), 226–e1.
- 2210 Strong, W. (2002). Assessing species abundance unevenness within and between plant communities.
2211 *Community Ecology*, 3(2), 237–246.
- 2212 Sultan, S., El-Mowafy, M., Elgaml, A., Ahmed, T. A., Hassan, H., & Mottawea, W. (2021). Metabolic
2213 influences of gut microbiota dysbiosis on inflammatory bowel disease. *Frontiers in physiology*, 12,
2214 715506.
- 2215 Suzuki, N., Nakano, Y., Yoneda, M., Hirofumi, T., & Hanioka, T. (2022). The effects of cigarette
2216 smoking on the salivary and tongue microbiome. *Clinical and Experimental Dental Research*, 8(1),
2217 449–456.
- 2218 Swidsinski, A., Khilkin, M., Kerjaschki, D., Schreiber, S., Ortner, M., Weber, J., & Lochs, H. (1998).
2219 Association between intraepithelial escherichia coli and colorectal cancer. *Gastroenterology*,
2220 115(2), 281–286.
- 2221 Swift, D., Cresswell, K., Johnson, R., Stilianoudakis, S., & Wei, X. (2023). A review of normalization
2222 and differential abundance methods for microbiome counts data. *Wiley Interdisciplinary Reviews:
2223 Computational Statistics*, 15(1), e1586.
- 2224 Taher, C., Frisk, G., Fuentes, S., Religa, P., Costa, H., Assinger, A., ... others (2014). High prevalence of
2225 human cytomegalovirus in brain metastases of patients with primary breast and colorectal cancers.
2226 *Translational oncology*, 7(6), 732–740.
- 2227 Tanner, A. C., Kent Jr, R., Kanasi, E., Lu, S. C., Paster, B. J., Sonis, S. T., ... Van Dyke, T. E. (2007).
2228 Clinical characteristics and microbiota of progressing slight chronic periodontitis in adults. *Journal
2229 of clinical periodontology*, 34(11), 917–930.
- 2230 Tanner, A. C., Paster, B. J., Lu, S. C., Kanasi, E., Kent Jr, R., Van Dyke, T., & Sonis, S. T. (2006).

- 2231 Subgingival and tongue microbiota during early periodontitis. *Journal of dental research*, 85(4),
2232 318–323.
- 2233 Tejeda, M., Farrell, J., Zhu, C., Haines, J. L., Wang, L.-S., Schellenberg, G. D., ... others (2021). Multiple
2234 viruses detected in human dna are associated with alzheimer disease risk. *Alzheimer's & Dementia*,
2235 17, e054585.
- 2236 Teles, F., Wang, Y., Hajishengallis, G., Hasturk, H., & Marchesan, J. T. (2021). Impact of systemic
2237 factors in shaping the periodontal microbiome. *Periodontology 2000*, 85(1), 126–160.
- 2238 Thaiss, C. A., Zmora, N., Levy, M., & Elinav, E. (2016). The microbiome and innate immunity. *Nature*,
2239 535(7610), 65–74.
- 2240 Tian, R., Liu, H., Feng, S., Wang, H., Wang, Y., Wang, Y., ... Zhang, S. (2021). Gut microbiota dysbiosis
2241 in stable coronary artery disease combined with type 2 diabetes mellitus influences cardiovascular
2242 prognosis. *Nutrition, Metabolism and Cardiovascular Diseases*, 31(5), 1454–1466.
- 2243 Tilg, H., Kaser, A., et al. (2011). Gut microbiome, obesity, and metabolic dysfunction. *The Journal of
2244 clinical investigation*, 121(6), 2126–2132.
- 2245 Tizabi, D., & Hill, R. T. (2023). Micrococcus spp. as a promising source for drug discovery: A review.
2246 *Journal of Industrial Microbiology and Biotechnology*, 50(1), kuad017.
- 2247 Tonetti, M. S., Greenwell, H., & Kornman, K. S. (2018). Staging and grading of periodontitis: Framework
2248 and proposal of a new classification and case definition. *Journal of periodontology*, 89, S159–S172.
- 2249 Tringe, S. G., & Hugenholtz, P. (2008). A renaissance for the pioneering 16s rRNA gene. *Current opinion
2250 in microbiology*, 11(5), 442–446.
- 2251 Tucker, C. M., Cadotte, M. W., Carvalho, S. B., Davies, T. J., Ferrier, S., Fritz, S. A., ... others (2017). A
2252 guide to phylogenetic metrics for conservation, community ecology and macroecology. *Biological
2253 Reviews*, 92(2), 698–715.
- 2254 Ulger Toprak, N., Yagci, A., Gulluoglu, B., Akin, M., Demirkalem, P., Celenk, T., & Soyletir, G. (2006).
2255 A possible role of bacteroides fragilis enterotoxin in the aetiology of colorectal cancer. *Clinical
2256 microbiology and infection*, 12(8), 782–786.
- 2257 Ursell, L. K., Metcalf, J. L., Parfrey, L. W., & Knight, R. (2012). Defining the human microbiome.
2258 *Nutrition reviews*, 70(suppl_1), S38–S44.
- 2259 Utzschneider, K. M., Kratz, M., Damman, C. J., & Hullarg, M. (2016). Mechanisms linking the gut
2260 microbiome and glucose metabolism. *The Journal of Clinical Endocrinology & Metabolism*,
2261 101(4), 1445–1454.
- 2262 Vander Haar, E. L., So, J., Gyamfi-Bannerman, C., & Han, Y. W. (2018). Fusobacterium nucleatum and
2263 adverse pregnancy outcomes: epidemiological and mechanistic evidence. *Anaerobe*, 50, 55–59.
- 2264 Van der Maaten, L., & Hinton, G. (2008). Visualizing data using t-sne. *Journal of machine learning
2265 research*, 9(11).
- 2266 Vasen, H. F., Mecklin, J.-P., Khan, P. M., & Lynch, H. T. (1991). The international collaborative group
2267 on hereditary non-polyposis colorectal cancer (icg-hnpcc). *Diseases of the Colon & Rectum*, 34(5),
2268 424–425.
- 2269 Vigneswaran, J., & Shogan, B. D. (2020). The role of the intestinal microbiome on colorectal cancer

- 2270 pathogenesis and its recurrence following surgery. *Journal of Gastrointestinal Surgery*, 24(10),
2271 2349–2356.
- 2272 Vilar, E., & Gruber, S. B. (2010). Microsatellite instability in colorectal cancer—the stable evidence.
2273 *Nature reviews Clinical oncology*, 7(3), 153–162.
- 2274 Viljoen, K. S., Dakshinamurthy, A., Goldberg, P., & Blackburn, J. M. (2015). Quantitative profiling of
2275 colorectal cancer-associated bacteria reveals associations between *fusobacterium* spp., enterotoxi-
2276 genic *bacteroides fragilis* (etbf) and clinicopathological features of colorectal cancer. *PLoS one*,
2277 10(3), e0119462.
- 2278 Walker, M. A., Pedamallu, C. S., Ojesina, A. I., Bullman, S., Sharpe, T., Whelan, C. W., & Meyerson, M.
2279 (2018). Gatk pathseq: a customizable computational tool for the discovery and identification of
2280 microbial sequences in libraries from eukaryotic hosts. *Bioinformatics*, 34(24), 4287–4289.
- 2281 Wang, N., & Fang, J.-Y. (2023). *Fusobacterium nucleatum*, a key pathogenic factor and microbial
2282 biomarker for colorectal cancer. *Trends in Microbiology*, 31(2), 159–172.
- 2283 Weaver, W. (1963). *The mathematical theory of communication*. University of Illinois Press.
- 2284 Whiteside, S. A., Razvi, H., Dave, S., Reid, G., & Burton, J. P. (2015). The microbiome of the urinary
2285 tract—a role beyond infection. *Nature Reviews Urology*, 12(2), 81–90.
- 2286 Witkin, S. (2019). Vaginal microbiome studies in pregnancy must also analyse host factors. *BJOG: An
2287 International Journal of Obstetrics & Gynaecology*, 126(3), 359–359.
- 2288 Wong, T.-T., & Yeh, P.-Y. (2019). Reliable accuracy estimates from k-fold cross validation. *IEEE
2289 Transactions on Knowledge and Data Engineering*, 32(8), 1586–1594.
- 2290 Wyss, C., Moter, A., Choi, B.-K., Dewhirst, F., Xue, Y., Schüpbach, P., ... Guggenheim, B. (2004).
2291 *Treponema putidum* sp. nov., a medium-sized proteolytic spirochaete isolated from lesions of
2292 human periodontitis and acute necrotizing ulcerative gingivitis. *International journal of systematic
2293 and evolutionary microbiology*, 54(4), 1117–1122.
- 2294 Xia, Y. (2023). Statistical normalization methods in microbiome data with application to microbiome
2295 cancer research. *Gut Microbes*, 15(2), 2244139.
- 2296 Yaman, E., & Subasi, A. (2019). Comparison of bagging and boosting ensemble machine learning methods
2297 for automated emg signal classification. *BioMed research international*, 2019(1), 9152506.
- 2298 Yang, I., Claussen, H., Arthur, R. A., Hertzberg, V. S., Geurs, N., Corwin, E. J., & Dunlop, A. L. (2022).
2299 Subgingival microbiome in pregnancy and a potential relationship to early term birth. *Frontiers in
2300 cellular and infection microbiology*, 12, 873683.
- 2301 Yildiz, B., Bilbao, J. I., & Sproul, A. B. (2017). A review and analysis of regression and machine learning
2302 models on commercial building electricity load forecasting. *Renewable and Sustainable Energy
2303 Reviews*, 73, 1104–1122.
- 2304 Yoshimura, F., Murakami, Y., Nishikawa, K., Hasegawa, Y., & Kawaminami, S. (2009). Surface
2305 components of *porphyromonas gingivalis*. *Journal of periodontal research*, 44(1), 1–12.
- 2306 Yu, Y., Guo, D., Qu, T., Zhao, S., Xu, C., Wang, L., ... Zhou, N. (2020). Increased risk of toxoplasma
2307 gondii infection in patients with colorectal cancer in eastern china: seroprevalence, risk factors,
2308 and a case–control study. *BioMed Research International*, 2020(1), 2539482.

- 2309 Yuan, K., Xu, H., Li, S., Coker, O. O., Liu, W., Wang, L., ... Yu, J. (2025). Intraneoplastic fungal
2310 dysbiosis is associated with colorectal cancer progression and host gene mutation. *EBioMedicine*,
2311 113.
- 2312 Zavareh, F. S. E., Hadiipour, M., Kalantari, R., Mousavi, S., Tavakolifard, N., & Darani, H. Y. (2021).
2313 Effect of toxoplasma gondii on colon cancer growth in mouse model. *Am J Biomed*, 9(2), 168–176.
- 2314 Zepeda-Rivera, M., Minot, S. S., Bouzek, H., Wu, H., Blanco-Míguez, A., Manghi, P., ... others (2024).
2315 A distinct fusobacterium nucleatum clade dominates the colorectal cancer niche. *Nature*, 628(8007),
2316 424–432.
- 2317 Zhang, C.-Z., Cheng, X.-Q., Li, J.-Y., Zhang, P., Yi, P., Xu, X., & Zhou, X.-D. (2016). Saliva in the
2318 diagnosis of diseases. *International journal of oral science*, 8(3), 133–137.
- 2319 Zhang, M., Zhang, Y., Sun, Y., Wang, S., Liang, H., & Han, Y. (2022). Intratumoral microbiota impacts
2320 the first-line treatment efficacy and survival in non-small cell lung cancer patients free of lung
2321 infection. *Journal of Healthcare Engineering*, 2022(1), 5466853.
- 2322 Zhang, N., Liu, Y., Tang, F.-Y., Yang, L.-Y., & Wang, J.-H. (2023). Structural characterization and in vitro
2323 anti-colon cancer activity of a homogeneous polysaccharide from agaricus bisporus. *International
2324 Journal of Biological Macromolecules*, 251, 126410.
- 2325 Zhang, X., Yu, D., Wu, D., Gao, X., Shao, F., Zhao, M., ... others (2023). Tissue-resident lachnospiraceae
2326 family bacteria protect against colorectal carcinogenesis by promoting tumor immune surveillance.
2327 *Cell host & microbe*, 31(3), 418–432.
- 2328 Zhou, X., Wang, L., Xiao, J., Sun, J., Yu, L., Zhang, H., ... others (2022). Alcohol consumption,
2329 dna methylation and colorectal cancer risk: Results from pooled cohort studies and mendelian
2330 randomization analysis. *International journal of cancer*, 151(1), 83–94.
- 2331 Zhu, W., & Lee, S.-W. (2016). Surface interactions between two of the main periodontal pathogens:
2332 *Porphyromonas gingivalis* and *tannerella forsythia*. *Journal of periodontal & implant science*,
2333 46(1), 2–9.
- 2334 Zhu, X., Han, Y., Du, J., Liu, R., Jin, K., & Yi, W. (2017). Microbiota-gut-brain axis and the central
2335 nervous system. *Oncotarget*, 8(32), 53829.
- 2336 Zhuang, Y., Wang, H., Jiang, D., Li, Y., Feng, L., Tian, C., ... others (2021). Multi gene mutation
2337 signatures in colorectal cancer patients: predict for the diagnosis, pathological classification, staging
2338 and prognosis. *BMC cancer*, 21, 1–16.

Acknowledgments

2340 I would like to disclose my earnest appreciation for my advisor, Professor **Semin Lee**, who provided
 2341 solicitous supervision and cherished opportunities throughout the course of my research. His advice and
 2342 consultation encouraged me to become as a researcher and to receive all humility and gentleness. I am also
 2343 grateful to all of my committee members, Professor **Taejoon Kwon**, Professor **Eunhee Kim**, Professor
 2344 **Kyemyung Park**, and Professor **Min Hyuk Lim**, for their meaningful mentions and suggestions.

2345 I extend my deepest gratitude to my Lord, **the Flying Spaghetti Monster**, His Noodly Appendage
 2346 has guided me through the twist and turns of this academic journey. His presence, ever comforting and
 2347 mysterious, has been a source of strength and humor during both highs and lows. In moments of doubt, I
 2348 found solace in the belief that you were there, gently reminding me to keep faith in the process. His Holy
 2349 Noodle has nourished my mind, and for that, I am truly overwhelmed. May His Holy Noodle continue to
 2350 guide me in all my future endeavors. *R’Amen.*

2351 I would like to extend my heartfelt gratitude to Professor **You Mi Hong** for her invaluable guidance
 2352 and insightful advice on PTB study. Her expertise in maternal and fetal health, along with her deep under-
 2353 standing of statistical and clinical interpretations, greatly contributed to refining the analytical framework
 2354 of this study. Her constructive feedback and thoughtful discussions provided critical perspectives that
 2355 enhanced the robustness and relevance of the research findings. I sincerely appreciate her generosity
 2356 in sharing her knowledge and effort, as well as her encouragement throughout my Ph.D. journey. Her
 2357 support has been instrumental in strengthening this work, and I am truly grateful for her contributions.

2358 I also would like to express my sincere gratitude for Professor **Jun Hyeok Lim** for his invaluable
 2359 guidance and insightful advice on lung cancer study. His expertise in cancer genomics and data interpreta-
 2360 tion provided essential perspectives that greatly enriched the analytical approach of my Ph.D. journey. His
 2361 constructive feedback and thoughtful discussion helped refine methodologies and enhance the scientific
 2362 rigor of the research. I deeply appreciate his willingness to share his knowledge and expertise, which has
 2363 been instrumental in shaping key aspects of this work. His support and encouragement have been truly
 2364 inspiring, and I am grateful for the opportunity to have benefited from his mentorship.

2365 I would like to extend my heartfelt gratitude to my colleagues of the **Computational Biology Lab @**
 2366 **UNIST**, whose collaboration, friendship, brotherhood, and support have been an invaluable part of my
 2367 journey. Your willingness to share insights, engage in thoughtful discussions, and offer encouragement
 2368 during the challenging moments of research has significantly shaped my academic experience. The
 2369 camaraderie in Computational Biology Lab made even the most demanding days more enjoyable, and I
 2370 am deeply grateful for the collaborative environment we created together. I appreciate you for standing
 2371 by my side throughout this Ph.D. journey.

2372 I would like to express my heartfelt gratitude to **my family**, whose unwavering support has been the
 2373 foundation of everything I have achieved. Your love, encouragement, and belief in me have sustained me
 2374 through every challenge, and I could not have come this far without you. From your words of wisdom to
 2375 your patience and understanding, each of you has played a vital role in helping me navigate this journey.
 2376 The strength and comfort I have drawn from our family bond have been my greatest source of resilience.

2377 Your presence, both near and far, has filled my life with warmth and motivation. I am deeply grateful for
2378 your unconditional love and for always being there when I needed you the most. Thank you for being my
2379 constant source of strength and inspiration.

2380 I am incredibly pleased to my friends, especially my GSHS alumni (**이망특**), for their unwavering
2381 support and encouragement throughout this journey. The bonds we formed back in our school days have
2382 only grown stronger over the years, and I am fortunate to have had such loyal and understanding friends
2383 by my side. Your constant words of motivation, and even moments of levity during stressful times have
2384 helped keep me grounded. Whether it was a late-night conversations, a shared laugh, or a simple message
2385 of reassurance, you all have played a vital role in keeping me focused and motivated. I am relieved for the
2386 ways you celebrated each small achievement with me and how you patiently listened to my worries. The
2387 memories of our shared past provided me with comfort and a sense of stability when the road ahead felt
2388 uncertain. I could not have reached this point without the love and friendship that you all have generously
2389 given. Each of your, in your unique way, has contributed to this dissertation, even if indirectly, and for
2390 that, I am forever beholden. I look forward to continuing our friendship as we all grow in our individual
2391 paths, knowing that the support we share is something truly special.

2392 I would like to express my deepest recognition to **my girlfriend (expected)** for her unwavering
2393 support, patience, and companionship throughout my Ph.D. journey. Her presence has been a constant
2394 source of comfort and motivation, helping me navigate the challenges of research and writing with
2395 renewed energy. Through moments of frustration and accomplishment alike, her encouragement has
2396 reminded me of the importance of balance and perseverance. Her kindness, understanding, and belief
2397 in me have been invaluable, making even the most difficult days feel lighter. I am truly grateful for her
2398 support and for sharing this journey with me, and I look forward to all the moments we will continue to
2399 experience together.

2400 I would like to express my sincere gratitude to the amazing members of my animal protection groups,
2401 DRDR (**두루두루**) and UNIMALS (**유니멀스**), whose dedication and compassion have been a constant
2402 source of motivation. Your unwavering commitment to improving the lives of animals has inspired me
2403 throughout this journey. I am also thankful for the beautiful cats we have cared for, whose presence
2404 brought both joy and purpose to our allegiance. Their playful spirits and gentle companionship served as
2405 daily reminders of why we continue to fight for animal rights. The bond we share, both with each other
2406 and with the animals we protect, has enriched my life in countless ways. I appreciate you all again for
2407 your support, dedication, and for being part of this meaningful cause.

2408 I would like to express my deepest gratitude to **everyone** I have had the honor of meeting throughout
2409 this journey. Your kindness, encouragement, and support have carried me through both the challenging
2410 and rewarding moments of my life. Whether through a kind word, thoughtful advice, or simply being
2411 there when I needed it most, your presence has made all the difference. I am incredibly fortunate to have
2412 received such generosity and warmth from those around me, and I do not take it for granted. Every act
2413 of kindness, no matter how big or small, has been a source of strength and motivation for me. To all
2414 my friends, colleagues, mentors, and beloved ones, thank you for your unwavering support. I am truly
2415 grateful for each of you, and your kindness has left an indelible mark on my journey.

2416 My Lord, *the Flying Spaghetti Monster*,
2417 give us grace to accept with serenity the things that cannot be changed,
2418 courage to change the things that should be changed,
2419 and the wisdom to distinguish the one from the other.

2420
2421 Glory be to *the Meatball*, to *the Sauce*, and to *the Holy Noodle*.
2422 As it was in the beginning, is now, and ever shall be.

2423 *R'Amen.*



May your progress be evident to all

