1    **Mucosal microbiomes and *Fusobacterium* genomics in Vietnamese colorectal cancer patients**

2

3    Hoang N. H. Tran[1#], Trang Nguyen Hoang Thu[1#], Phu Huu Nguyen[2], Chi Nguyen Vo[2], Khanh Van Doan[3],

4    Chau Nguyen Ngoc Minh[1], Ngoc Tuan Nguyen[2], Van Ngoc Duc Ta[2], Khuong An Vu[2], Thanh Danh Hua[2],

5    To Nguyen Thi Nguyen[1], Tan Trinh Van[1], Trung Pham Duc[1], Ba Lap Duong[2], Phuc Minh Nguyen[2], Vinh

6    Chuc Hoang[2], Duy Thanh Pham[1,4], Guy E. Thwaites[1,4], Lindsay J. Hall[5,6,7], Daniel J. Slade[8], Stephen

7    Baker[9], Vinh Hung Tran[2], Hao Chung The[1*]

8

9    [1] Oxford University Clinical Research Unit, Ho Chi Minh City, Vietnam

10    [2] Binh Dan Hospital, Ho Chi Minh City, Vietnam

11    [3] Department of Oral Biology, Yonsei University College of Dentistry, Seoul, Korea

12    [4] Centre for Tropical Medicine and Global Health, Nuffield Department of Clinical Medicine, University of Oxford,

13    Oxford, United Kingdom

14    [5] Quadram Institute Biosciences, Norwich Research Park, Norwich, United Kingdom

15    [6] Norwich Medical School, University of East Anglia, Norwich Research Park, Norwich, United Kingdom

16    [7] Intestinal Microbiome, School of Life Sciences, ZIEL - Institute for Food & Health, Technical University of

17    Munich, Freising, Germany

18    [8] Department of Biochemistry, Virginia Tech, Blacksburg, VA 24061, USA

19    [9] Department of Medicine, Cambridge Institute of Therapeutic Immunology and Infectious Diseases (CITIID),

20    University of Cambridge, Cambridge, United Kingdom

21

22    Hoang N.H. Tran and Trang Nguyen Hoang Thu contributed equally to this article.

23    * Corresponding author: Dr. Hao Chung The, Department of Molecular Epidemiology, Oxford University Clinical

24    Research Unit (OUCRU), 764 Vo Van Kiet St., Ward 1, District 5, Ho Chi Minh City, Vietnam.

25    Tel: +84 969937143                          Email: haoct@oucru.org

26

27

28    **Abstract**

29    Perturbations in the gut microbiome have been linked to the promotion and prognosis of colorectal cancer

30    (CRC), with the colonic overabundance of *Fusobacterium nucleatum* shown as the most consistent

31    marker. Despite the increasing health burden inflicted by CRC in low- and middle-income countries like

32    Vietnam, the CRC-specific microbiome in these populations remains underexplored. Here we conducted a

33    study in Vietnam to enrol 43 CRC patients (cases) and 25 patients with non-cancerous colorectal polyps

34    (controls) between December 2018 and January 2020. Our study investigated the mucosal microbiome

35    signature and genomic diversity of *Fusobacterium* in Vietnamese CRC patients, using a combination of

36    16S rRNA gene profiling, anaerobic microbiology, and whole genome sequencing. We found that several

37    oral bacteria, including *F. nucleatum* and *Leptotrichia,* were significantly more abundant in the tumour

38    mucosa, and these two bacteria were also more enriched in tumours of advanced CRC stages (III-IV). We

39    obtained 53 *Fusobacterium* genomes from the saliva, tumour and non-tumour mucosa of six CRC patients.

40    Isolates from the gut mucosa belonged to diverse *F. nucleatum* subspecies (*nucleatum*, *animalis*, *vincentii*,

41    *polymorphum*) and a potential new subspecies of *F. periodonticum*. The *Fusobacterium* population within

42    each individual was distinct and in many cases diverse, with minimal intra-clonal variation. Phylogenetic

43    analyses showed that within each individual, tumour-associated *Fusobacterium* were clonal to those

44    isolated from non-tumour mucosa, but distantly related to those isolated from saliva. Genes encoding

45    major virulence factors (Fap2 and RadD) showed variability in length and evidence of horizontal gene

46    transfer. Our work provides a framework to understand the genomic diversity of *Fusobacterium* within

47    the CRC patients, which can be exploited for the development of CRC diagnostic and therapeutic options

48    targeting this oncobacterium.

49

50    **Keywords:** colorectal cancer microbiome; *Fusobacterium*; *Leptotrichia*; *Fusobacterium* genomic;

51    developing country; *Fusobacterium* diversity; cancer microbiome; mucosal microbiome

**Introduction**

52

53 Colorectal cancer (CRC) is the second leading cause of cancer mortality worldwide, contributing to an

54 estimate of 850,000 deaths and ~1.8 million new cases in 2018 [1,2]. The majority of CRC cases are

55 sporadic (without clear heredity components), with well-established lifestyle risk factors attributed to

56 obesity, alcohol consumption and a diet enriched with red or processed meat [3]. The vast and diverse

57 microbial community inhabiting the colon (termed the gut microbiome) is an integral part of human

58 health, and act as an important interface mediating the interactions between environmental cues, host

59 biology, and CRC [4,5]. Research on CRC gut microbiome has consistently underlined the abundances of

60 certain marker bacteria, among which *Fusobacterium nucleatum* has been most widely reported and

61 intensively studied [6–10].

62

63 The Gram-negative rod-shaped *F. nucleatum* is a common anaerobic member of the human oral

64 microbiome, and it is currently composed of four subspecies (*nucleatum*, *vincentii*, *animalis*, and

65 *polymorphum*) [11]. Mechanistic studies have demonstrated that *F. nucleatum* possesses several virulence

66 factors, most notably FadA and Fap2, which enable the bacteria to potentiate colonic tumourigenesis. The

67 adhesin FadA binds to E-cadherin in CRC cells and activates the β-catenin-dependent oncogenic

68 pathways [12], while the lectin Fap2 further facilitates *F. nucleatum* invasion into CRC cells by

69 specifically binding to the tumour-enriched carbohydrate Gal-GalNAc [13]. Such interaction triggers the

70 secretions of the pro-inflammatory (IL-8) and pro-metastatic (CXCL-1) cytokines, creating a tumour

71 environment conditioned for accelerated growth and migratory tendency [14]. Recent studies have further

72 highlighted that the bacteria could induce DNA damage in oral and colorectal cancerous cells [15,16].

73 Additionally, *F. nucleatum* lipopolysaccharide was shown to induce resistance to chemotherapy via

74 activation of the autophagy machinery in CRC cells, thus complicating effective CRC treatment [17]. As

75 a result, enrichment of *F. nucleatum* in CRC microbiomes has been associated with more severe

76 prognosis and poorer overall survival, particularly in a subset of patients with mesenchymal tumours [18–

77 20]. Preclinical research demonstrated that *F. nucleatum* elimination by antibiotics reduced colorectal

78    tumour proliferation in mice [21]. These evidences strongly support for the utilization of *F. nucleatum* as

79    a target for CRC diagnosis and therapy, but current translational potential is hampered by the lack of

80    insights into *F. nucleatum* diversity and its genomic characteristics in CRC patients.

81

82    The majority of microbiome studies, on either healthy or CRC cohorts, were conducted in high-income

83    countries, and such data are sparse regarding populations in developing settings, where host factors, diet

84    and lifestyle could greatly influence the gut microbiome composition and function. Vietnam has an

85    increasing ageing population adopting a more 'Westernized' diet and sedentary lifestyle [22], where CRC

86    incidence is predicted to climb and rank as among the top three cancers by 2025 [23]. Therefore, CRC

87    microbiome studies in Vietnam are necessary to establish the basis for the implementation of

88    microbiome-oriented strategies for CRC prevention, diagnosis, prognosis and therapy. We set out to

89    investigate the microbiome signatures of Vietnamese CRC patients, by applying 16S-rRNA gene

90    profiling on the saliva and gut mucosa collected from patients with CRC and non-cancerous colorectal

91    polyps. Additionally, different from prior studies, we used anaerobic culturing and whole genome

92    sequencing (WGS) to study the genomic diversity of *Fusobacterium* colonizing these CRC patients,

93    allowing an in-depth and high-resolution examination of these bacterial populations.

94      **Results**

95      **Gut mucosal, but not salivary, microbiomes differ significantly between CRC and controls**

96      We enrolled 43 CRC patients (cases) and 25 patients with colorectal polyps (controls) between December

97      2018 and January 2020. 16S rRNA microbiome profiling was performed for all the saliva and gut mucosa

98      samples collected from the participants, including tissues originating from the diseased (CRC tumour or

99      polyps) and the adjacent normal sites. To limit the scope of this study, we selected participants with

100     tumours/polyps detected in the distal colon or rectum. The patients' demographic and clinical data were

101     summarized in Table 1, which showed that there were no significant differences between the two groups.

102     All polyps showed not more than low-grade dysplasia (i.e. non-cancerous), demonstrating the validity of

103     our control group. Microbiome profiling identified 865 filtered amplicon sequence variants (ASVs – a

104     marker for distinct taxonomic classification) among 66 saliva samples, with a median library size of

105     36,250 paired-end reads [IQR: 31,827 – 50,317]. Due to their lower microbial biomass, the library size of

106     gut mucosal microbiomes was smaller (median: 17,711 [IQR: 9,037 – 30,135]), with 1,073 filtered ASVs

107     detected across 129 mucosa samples (seven removed).

108

109     Ordination by principal coordinate analysis (PCoA), based on phylogenetic-assisted isometric log-ratio

110     (PhILR) transformed value, showed that the salivary microbiomes of CRC and controls completely

111     overlapped (Figure 1A). Only active smoking within the last two years, but not CRC status, was

112     significantly associated with the salivary microbiome structure (RDA, p-value = 0.033). Likewise, only

113     two ASVs belonging to the genera *Leptotrichia* and *Solobacterium* were consistently identified as

114     significantly more abundant in the CRC's salivary microbiome. These point to the high structural

115     similarity in the salivary microbiome between the two groups. By contrast, the gut mucosal microbiomes

116     differ significantly based on CRC status (Figure 1B). CRC and diabetes significantly contributed to the

117     variance in the gut microbiome (RDA, p-value < 0.05). Gut mucosa collected within a participant (tumour

118     and non-tumour for CRC, biopsy and polyp for control) shared more similarity in their microbiomes than

119     those of the same sample type between participants (Figure 1C), resembling findings from previous

120    research [8]. We also conducted these analyses using the weighted Unifrac and Bray-Curtis distances,

121    which produced similar interpretations. Additionally, we performed unsupervised clustering on gut

122    mucosa microbiomes, which showed the presence of two robust community state types (CSTs) supported

123    by an out-of-bag error rate of 10.8% in a random forest classification. This algorithm also identified that

124    several 'balances' (Proteobacteria/Actinobacteria, other bacteria/Lachnospiraceae) contributed

125    significantly in separating the two CSTs (Figure S1). CST1 was generally more enriched in

126    Gammaproteobacteria (mostly *Escherichia*) while CST2 had higher abundance of Actinobacteria (mainly

127    *Collinsella*) and Lachnospiraceae (Figure 1D). The two CSTs were similar in library size (p-value = 0.15,

128    t-test), but different in CRC status (p-value=0.002, Fisher-exact test), with the majority of control samples

129    (72%) belonging to CST1. No other tested covariates were associated with CST grouping. Samples from

130    the same patients mostly shared the same CST membership (90.3%, n=56/62 patients with paired

131    microbiomes), and CRC samples were distributed in both CSTs with different proportion (CST1 = 36,

132    CST2 = 50). These findings suggest that while CRC status mainly explained the dissimilarity observed in

133    the gut mucosal microbiomes, their overall configurations were determined by the dominant presence of

134    Gammaproteobacteria (*Escherichia*), possibly driven by an unknown or stochastic factor.

135

136    **Enrichment of oral bacteria in the tumour gut mucosa**

137    We applied differential abundance analysis to rigorously detect bacteria enriched in the CRC tumours, by

138    comparing results from different approaches, including ANCOMBC, DESeq2 and corncob (see Methods)

139    [24–26]. Our analyses revealed that ASVs classified as bacteria of putative oral origin (*Gemella*,

140    *Peptostreptococcus*, *F. nucleatum*, *Leptotrichia*, *Selenomonas sputigena*, and *Campylobacter rectus*)

141    were overabundant in the tumour mucosa, compared to control biopsies (Figure 2B). This finding

142    corroborates the observation that these oral bacteria had a higher relative abundance in CRC gut

143    microbiomes of both CSTs (Figure 1D). Within the CRC patients, tumours also showed an elevated

144    presence of the aforementioned oral bacteria (alongside *Hungatella*, *Lachnoclostridium*, and *Osillibacter*)

145    when compared to adjacent non-tumour mucosa, albeit with less pronounced fold change (Figure 2A).

146     These increases were coupled with the reduction in abundances of commensal anaerobes in the tumour

147     mucosa, such as *Blautia*, *Parabacteroides*, *Dorea*, and *Collinsella*. When comparing between different

148     cancer stages, the increased abundance of one taxon (*Leptotrichia,* ASV-13) was consistently associated

149     with tumours of advanced stages (III-IV), compared to stage II (Figure 2C). Results from DESeq2 alone

150     additionally showed that *F. nucleatum* was also enriched in advanced CRC stages (adjusted p-value

151     <0.05). ASVs confidently assigned as *F. nucleatum* (n=14) and *Leptotrichia* spp. (n=16) were present

152     with at least 0.1% relative abundance in ~54% and 40% of tumour microbiomes, respectively. On the

153     other hand, tumour samples with low abundances of oral bacteria (<1%) were dominated by other taxa,

154     including *Escherichia/Klebsiella* (n= 4), *Megamonas* (n=4), *Fusobacterium mortiferum* (n=2), and

155     *Helicobacter* (n=1). We performed similar analysis within the control group and showed that only one

156     ASV (*Faecalibacterium*) was consistently depleted in polyps compared to paired biopsies. However,

157     when compared to CRC samples, *F. mortiferum*, *Tyzzerella*, and *Sutterella* were significantly enriched in

158     the control gut microbiomes (Figure 2B).

159

160     To investigate bacterial co-occurrence and their potential interactions, we next used CCLasso and

161     SpiecEasi to construct a correlation network of gut microbiomes from CRC patients (n=86)  (Figure 3)

162     [27,28]. Two oral bacteria clusters emerged from this network, one consisting of several *Streptococcus*

163     and *Veillonella* taxa, and another composed mostly of aforementioned tumour-associated ASVs

164     (*Leptotrichia*, *Selenomonas*, *F. nucleatum*, *Streptococcus*, *Granulicatella*, *Gemella*, *Peptostreptococcus*,

165     and *Parvimonas*). The latter cluster exhibited positive correlation with *E. coli*, and antagonism toward

166     *Blautia*, a member of the gut anaerobic commensal network. Besides, other tumour-associated ASVs such

167     as *Hungatella*, *Lachnoclostridium*, and *C. rectus* were clustered alongside *Negativibacillus* and

168     *Eggerthella*, which showed strong negative correlations with anaerobic gut commensals *Dorea*,

169     *Bacteroides*, and *Faecalibacterium*. These findings highlight the potential competition between tumour-

170     associated taxa and common gut commensal anaerobes (Lachnospiraceae, Oscillospirales). Other

171     *Fusobacterium* species, *F. mortiferum* and *F. varium* were not linked to the oral clusters, showing that

172    they were mainly gut inhabitants. Comparison with the network constructed from salivary microbiomes

173    revealed that the same tumour-associated ASVs (*F. nucleatum*, *Gemella*, *Selenomonas*) formed similar

174    clusters as observed in the CRC gut microbiomes (Figure S2). This indicates that tumour-associated

175    ASVs could have oral origin in this examined cohort, and they likely co-exist in the polymicrobial

176    biofilms (similar to those present in the oral cavity) upon gut colonization. Notably, the tumour-

177    associated *Leptotrichia* (ASV-13) had very low abundance in the salivary microbiome (mean: 0.037%,

178    prevalence: 26%), and it differs from the dominant *Leptotrichia* detected in saliva (ASV-19) in 17

179    nucleotides (pairwise similarity: 93.3%).

180

181    **Diverse *Fusobacterium* colonizes CRC patients**

182    Since *F. nucleatum* was more enriched in the tumour microbiomes and previously demonstrated to

183    promote tumourigenesis, we next studied the population structure of *Fusobacterium* recovered from CRC

184    patients. Six patients with a *Fusobacterium* relative abundance at the tumour site exceeding 10% (except

185    for patient 18) and covered different cancer stages were selected for *Fusobacterium* isolation. In total, we

186    isolated 56 presumptive *Fusobacterium* organisms, as identified by the matrix-assisted laser

187    desorption/ionization time of flight mass spectrometer (MALDI-TOF), from the oral, nontumour and

188    tumour samples of these patients (Table 2). Whole genome short-read sequencing was performed on these

189    isolates, and three were determined contaminated and removed from analyses. Fifty-three recovered

190    genomes belong to *F. nucleatum* (n=38) and *F. periodonticum* (n=15) species complexes, of which

191    phylogenetic reconstruction was performed separately. Core-genome phylogeny of *F. nucleatum* showed

192    that tumour-associated isolates were detected in all four subspecies (*animalis*, *vincentii*, *nucleatum*,

193    *polymorphum*) (Figure 4A). In the *F. periodonticum* phylogeny, tumour-associated isolates (n=8, isolated

194    from P18, P40) formed a distinct cluster that is phylogenetically separated from the available references

195    (Figure 4B). These isolates all showed ~91% average nucleotide identity (ANI) to the closest *F.*

196    *periodonticum* references, suggesting that they constitute a novel subspecies of this species complex,

197    denoted herein as novel *F. periodonticum* (novelFperi). Likewise, two gut isolates (H16-13, H16-14)

198    shared 93% ANI to the closest *F. nucleatum* references and were phylogenetically distant from the

199    remaining *F. nucleatum* isolates, potentially indicative of a novel *F. nucleatum* subspecies. Across the

200    two phylogenies, we identified 14 phylogenetic clusters (PCs; 2 – 6 isolates each) and 12 singletons

201    originating from this study's collection, which were collectively named as PCs herein.

202

203    The *Fusobacterium* population within each individual patient was diverse (2 – 7 PCs). Several

204    *Fusobacterium* species/subspecies were detected in each patient's saliva, sometimes with more than one

205    PCs of the same subspecies (P18, P46) (Table 2). Likewise, we observed similar diversity in gut-

206    associated isolates, with more than one PCs detected in three patients (P10, P16, P18). Most patients did

207    not share the same *Fusobacterium* subspecies recovered from both oral- and gut-associated isolates,

208    except for P16 (*polymorphum*). However, phylogenetic evidence confirmed that the two niches harboured

209    distinct populations, which were ~16,955 SNPs apart (Figure 4A). Particularly, oral *Fusobacterium*

210    isolates from P18 (n=9) belonged to six different PCs (mostly *F. periodonticum* and *F. hwasookii*), while

211    6/7 gut isolates were of a single novelFperi clone. By contrast, *Fusobacterium* from tumour and

212    nontumour sites were frequently clustered in the same PC (n=4; in P10, P16, P18 and P40), indicating

213    that the same bacterial clones have colonized and spread beyond the tumour microenvironment. We used

214    the mapping approach to confidently inspect the intraclonal variations within these PCs, and showed that

215    they shared minimal genetic differences in the core genome (1 – 2 SNVs). These values fall in range with

216    the variation observed in five other gut PCs (with either tumour or nontumour isolates; 0 – 5 SNVs) and

217    five other oral PCs (1 – 10 SNVs).

218

219    **Variation in *Fusobacterium* virulence gene content**

220    We next sought to examine the presence of several *Fusobacterium* virulence factors, of which

221    pathogenicity has been proven in experimental studies, including genes encoding adhesin (*fadA*, *cbpF*),

222    lectin (*fap2*), and bacterial co-aggregation factor (*radD*) [12,13,29,30]. RadD is an autotransporter

223    facilitating *Fusobacterium's* interspecies interaction in polymicrobial biofilms [30], while CbpF inhibits

224    CD4$^+$ T-cell response through CEACAM1 binding and activation [31]. Genomic screening showed that

225    *fap2* was present and intact in the majority of genomes from both species (49/53), with disruptive

226    mutations occurring in some isolates, such as the tumour-associated *F. nucleatum animalis* in P46 (Figure

227    4A). We also detected *fadA* in all isolates (except S18-65), with all *F. periodonticum* variants one amino

228    acid shorter (codon A22) than the canonical FadA found in *F. nucleatum* (129 aa). The other elements

229    showed variable presence among the examined genomes. For example, *cbpF* was present in all *F.*

230    *nucleatum nucleatum*, *F. nucleatum vincentii*, and novelFperi, while *radD* was co-localised with

231    *fadA2/radA* (a 122 aa *fadA* homolog) in 28 isolates. Another *fadA* homolog (*fadA3*) with unknown

232    function was prevalent in both two *Fusobacterium* species. Phylogenies of FadA and CbpF showed that

233    the two tree topologies were largely in agreement with those inferred from the core genomes, suggesting

234    the absence of horizontal gene transfer (Figure S3). By contrast, the clustering pattern observed in the

235    Fap2 phylogeny was concordant to subspecies classification for *F. nucleatum nucleatum*, *F. nucleatum*

236    *vincentii*, and *F. periodonticum*, but was admixed for *F. nucleatum polymorphum*, *F. hwasookii* and *F.*

237    *nucleatum animalis* (Figure S4A). *fap2* encodes a very large protein of variable length (median of 3938 aa

238    [range: 3436 – 4669]), and the protein length showed some correlation with its phylogenetic clustering,

239    with variants > 4200 aa (n=6) all belonging to a monophyly composed of *F. hwasookii* and *F. nucleatum*

240    *polymorphum*. Similarly, the RadD phylogeny did not concur with those inferred from the core genomes,

241    and its length variation (median 3526 aa [range: 3461 – 3602]) also showed association with the tree

242    topology (Figure S4B). *radD* was ~800 bp downstream of *fadA2*, which is flanked by an IS150

243    transposase on the *F. nucleatum* 23726 reference genome. This could explain the mobilization mechanism

244    of *radD-fadA2* across the *Fusobacterium* phylogeny. These data indicate that the autotransporter

245    encoding genes *fap2* and *radD* may have undergone frequent horizontal gene transfer or recombination in

246    the *F. nucleatum* species complex.

247 **Discussion**

248 Our study revealed the composition of microbiome perturbations at the tumour mucosa of Vietnamese

249 patients with CRC and non-cancerous colorectal polyps. Tumour-enriched taxa include mostly bacteria of

250 putative oral origin, such as *F. nucleatum*, *Leptotrichia*, *Gemella*, *C. rectus*, and *Selenomonas*, which

251 agrees with findings from previous studies profiling either gut mucosal or faecal microbiomes in different

252 CRC populations [8,9,32]. This suggests that the proliferation of oral bacteria at the gut mucosa could be

253 a universal signature of CRC microbiomes. Several of these oral taxa shared identical ASVs between the

254 oral and gut niches, pointing to the oral origin of tumour-associated taxa. Our analysis found that these

255 bacteria also display a co-occurrence pattern, indicating that they likely co-exist in a biofilm-like

256 aggregate upon colonization at the gut mucosa. Indeed, previous research has confirmed the frequent

257 presence of polymicrobial biofilms composed of oral taxa (*F. nucleatum*, *Peptostreptococcus*, *Gemella*)

258 in colorectal tumour tissues [33]. Among the oral bacteria, *F. nucleatum* stands out for its ability to form

259 "bridging" interactions with other bacteria via the presence of several adhesins [11]. *F. nucleatum* was

260 recently reported to secrete FadA with amyloid properties, which confers acid tolerance and provides a

261 scaffold for biofilm formation [34]. In addition, our analyses pointed to the significant presence of

262 *Leptotrichia* in tumour mucosa, especially in advanced tumours. This association, however, has only been

263 noted in few studies [32,35]. This may be due to the differences in sampling location, as tumours excised

264 from the distal colon (as performed for all cases in our study) were reported to harbour a higher

265 abundance of *Leptotrichia*, compared to those originating from the proximal colon [35]. The

266 overabundance of *Leptotrichia* in the salivary microbiome has been implicated in patients with malignant

267 oral leukoplakia and pancreatic cancer [36,37], as well as with CRC as shown in this study. *Leptotrichia*

268 belongs to the same order as *Fusobacterium* (Fusobacteriales) and could carry virulence factors similar to

269 those found in the latter genus. It is noteworthy that the predominant tumour-associated *Leptotrichia*

270 taxon (ASV-13) could be detected from different CRC patients, but was in very low abundance in these

271 patients' salivary microbiomes. This suggests that a distinct *Leptotrichia* species/genotype was associated

272 with CRC, which warrants more in-depth investigations.

273

274    Asides from oral taxa, *Hungatella* overabundance was the most significant signature of CRC microbiome

275    in our dataset. This falls in line with results from a recent metagenomic meta-analysis, showing that

276    *Hungatella hathewayi*'s specific choline trimethylamine-lyase gene (*cutC*) was significantly enriched in

277    the faecal microbiomes of CRC patients [10]. Moreover, colonic *H. hathewayi* could induce

278    hypermethylation in prominent tumour suppressor genes, thus silencing their functions and promoting

279    intestinal epithelial cell proliferation [38]. Combining quantitative detection of microbial markers (*H.*

280    *hathewayi*, *F. nucleatum*, *Lachnoclostridium*, and *Bacteroides clarus*) with faecal immunochemical test

281    greatly increased sensitivity (reaching ~94%) for diagnosing CRC in a Chinese population [39].

282    Consistent with this finding, our analyses suggested that *Hungatella* and *Lachnoclostridium* were

283    overabundant and co-occurring in the tumour mucosa of the Vietnamese cohort, which supports the

284    feasibility of applying such microbial detection test for non-invasive CRC screening in our setting. On the

285    other hand, we found that *F. mortiferum* was the most significantly enriched taxon in the polyp control

286    group. *F. mortiferum* was known as a hallmark for dysbiosis in infectious diarrhoea [40], and recent

287    studies have also reported the abundance of *F. mortiferum* in patients with colorectal polyps [41,42].

288    Furthermore, this species was shown to be present in the gut microbiomes of ~60% of a cohort in

289    Southern China, albeit in very low abundance (~0.5%) [43]. Unlike other *Fusobacterium* species, *F.*

290    *mortiferum* was devoid of distinctive virulence factors such as adhesins FadA and Fap2 [44], but could

291    utilize a wide range of sugars for growth independent of amino acid metabolism [45]. The association

292    between *F. mortiferum* and colorectal polyps will need to be further addressed in future studies.

293

294    Despite the increasing importance of *F. nucleatum* in the pathogenesis of CRC and other invasive

295    diseases [11], genomic characterization of these bacteria from patient populations is currently limited due

296    to technical difficulties in *Fusobacterium* isolation. Here, we applied targeted culturomics approach,

297    which combines anaerobic culturing, high-throughput identification by MALDI-TOF and WGS, to study

298    the *Fusobacterium* population in high resolution and help uncover novel bacteria [46]. Indeed, we

299    discovered novel subspecies of both *F. nucleatum* and *F. periodonticum* from culturing the gut mucosa,

300    showing that the microbiomes in non-Western settings offer untapped diversity. Using metagenomic

301    assemblies from Chinese faecal microbiomes, Yeoh and colleagues have proposed several new

302    *Fusobacterium* species (based on 95% ANI cutoff) [44]. However, our WGS approach provided more

303    accurate and complete realization of the bacterial genomes, which contributes to the global representation

304    of *Fusobacterium* diversity (with 26 non-duplicate assemblies added). Furthermore, our approach allows

305    for delineation of bacteria from tumour and non-tumour sites, which is inaccessible by faecal

306    metagenomes. The populations of *Fusobacterium* colonizing the oral cavity and gut mucosa were

307    heterogeneous within each individual, even at the subspecies level, which mirrors the diversity observed

308    previously for gut commensals such as *Bifidobacterium* species [47]. Though this study did not provide

309    evidence of genetic relatedness between oral and gut *Fusobacterium* isolates, this likely point to the high

310    diversity of *Fusobacterium* in the oral niche [48]. Besides, *Fusobacterium* is abundant in subgingival

311    dental biofilms, which our salivary sampling did not fully capture. Previous research deploying WGS has

312    demonstrated that oral and tumour-originated *F. nucleatum* shared little genetic divergence, supporting

313    the notion that oncogenic *Fusobacterium* arise from the patient's oral microbiome [49]. Chronic

314    infections with *Helicobacter pylori* at the stomach, which increases the risk of gastric cancer, usually

315    result in extensive clonal propagations detected by WGS within each patient, though isolates were

316    collected in a single timepoint [50]. This prolonged colonization scenario contrasts with our observations

317    in CRC, in which multiple *Fusobacterium* clones (with minimal intraclonal variation) were present at

318    each patient's gut mucosa. Given that CRC could take years to develop, we speculate that the

319    *Fusobacterium* population at the tumour site fluctuates in response to the frequent seedings from the

320    highly diverse oral source. Longitudinal study design is needed to address this hypothesis, and to further

321    assess how *Fusobacterium* adapts to the gastrointestinal pressure.

322

323    The two well-described major virulence genes (*fadA* and *fap2*) were identified in the majority of

324    *Fusobacterium* genomes, regardless of niche. This concurs with previous research reporting the high

325     prevalence of *fadA* and *fap2* in *F. nucleatum* and *F. periodonticum* metagenomic assemblies from a

326     cohort in China [44]. These suggest that *Fusobacterium* with high virulence potential are prevalent in the

327     human population, and the genetic presence of *fadA* and *fap2* is not suitable for predicting the risk of

328     *Fusobacterium*-related CRC. All gut-derived novelFperi isolates harboured the examined virulence genes

329     (*fadA*, *fap2*, *radD*, and *cbpF*), which was more similar to *F. nucleatum* compared to *F. periodionticum*.

330     Moreover, *fap2* and *radD* showed variation in gene length and evidence of horizontal gene transfer,

331     underlying the significance of dynamic evolutionary processes in shaping *Fusobacterium*'s virulence

332     landscape. Since Fap2 orchestrates *F. nucleatum* invasion into CRC tumour cells via specific binding to

333     Gal-GalNAc, this ligand-receptor interaction was recently proposed as a target for clinical intervention in

334     *Fusobacterium*-enriched CRC [51]. Interestingly, our genetic analysis predicted that *fap2* was either

335     missing or truncated in some gut-associated *Fusobacterium* isolates, which may indicate the complex

336     lifestyle of *Fusobacterium* once colonizing the gut environment.

337

338     Some limitations were notable in our study design. Due to ethical concerns, patients with colorectal

339     polyps were selected as the control group, instead of healthy age-matched individuals. Our interpretations

340     do not extend to cancer in the proximal colon, though previous reports have noted that proximal CRC

341     tumours had a higher *Fusobacterium* abundance [52]. The sample size of cultured *Fusobacterium* isolates

342     was moderate and did not include longitudinal sampling, so it was not possible to investigate the bacterial

343     evolution in longer timeframe. Notwithstanding these shortcomings, our study reconfirmed the prominent

344     role of oral anaerobic conglomerates in CRC microbiome in an understudied Asian population, and

345     provided new insights into the genomic diversity of the oncobacterium *Fusobacterium*. The observed

346     diversity in this organism should be taken into account when designing future diagnostic or therapeutic

347     tools that target *Fusobacterium*.

348    **Material and Methods**

349    *Study design and sample collection*

350    This prospective case-control study enrolled adult Vietnamese patients (≥18 years old) admitted at Binh

351    Dan Hospital, a large surgical hospital in Ho Chi Minh City Vietnam, from December 2018 to January

352    2020. The study received ethical approval from the Ethics Committee of Binh Dan Hospital. Written

353    informed consent was obtained from all study participants. Cases were defined as patients diagnosed with

354    left-sided colorectal cancer (distal colon and rectum) of stage II onward, who received colectomy

355    treatment and underwent non-antibiotic pre-operative bowel preparation. Controls were patients

356    diagnosed with colorectal polyps (single/scattered non-cancerous polyps at distal colon or rectum), who

357    received polypectomy at the hospital. For both cases and controls, patients were excluded if they (1) had

358    received antimicrobial treatments within two weeks prior to enrolment, (2) had additional gastrointestinal

359    infections or obstructions, or (3) were immunocompromised. Additionally, the study excluded patients

360    who had received chemo- and/or radio-therapy within four weeks prior to enrolment (for cases) and those

361    were diagnosed with familial adenomatous polyposis (controls).

362

363    Demographic and clinical information were collected from study participants at recruitment. The

364    calculated body mass index (BMI) was categorized based on WHO recommendation for Asian

365    populations [53]. Cancer stage classification was based on the TNM Staging system [54]. A saliva sample

366    (~3mL) was collected pre-operation from each study participant (by spitting into a sterile container). For

367    cases, the mucosa epithelia at the tumour and adjacent non-tumour (2-10 cm away from the tumour) sites

368    were collected aseptically from the excised colon. For controls, we collected colorectal polyps and 2-3

369    biopsies of non-polyp mucosal epithelium (~50 mg) during colonoscopy. All clinical samples were stored

370    on ice and transported back to the OUCRU laboratory within 4 hours, then were stored in -80$^{o}$C until

371    further experiments.

372

373    *16S rRNA gene sequencing*

374 We selected 43 cases and 25 controls for microbiome profiling. Total DNA was extracted from the gut

375 mucosa samples (n=136) using the FastDNA spin kit for soil (MP Biomedicals, USA), with bead-beating

376 step on Precellys 24 homogenizer (Bertin Instruments, France). DNA from the saliva samples (n=67, one

377 missing) was extracted using the ReliaPrep Blood gDNA Miniprep (Promega, USA). For microbiome

378 profiling, all samples underwent primary PCR amplification (30 cycles) using the conventional V4

379 primers (515F-806R) and KAPA HiFi Hot Start DNA polymerase (KAPA Biosystems, USA), and

380 secondary PCR was performed to add dual-indexes (IDT, USA) to each sample, following procedures

381 optimized in a published protocol [55,56]. Additionally, we applied the same procedures to a positive

382 control (eight species Zymo mock community, Zymo Research, USA) and six negative controls (two for

383 each DNA extraction kits, and two no-template PCR amplifications), in order to respectively evaluate the

384 experiment's efficacy and detect contamination from reagents and kits (kitome) [57]. 16S rRNA

385 sequencing was performed for all samples on one run of the Illumina MiSeq platform, to generate 250 bp

386 paired-end reads.

387

388 *Microbiome data analysis*

389 All data analyses were conducted in R (v4.1.1) and Rstudio using multiple packages, including 'dada2',

390 'phyloseq', 'DESeq2', 'ANCOMBC', 'corncob', 'philr', 'ggplot2', 'vegan', 'SpiecEasi' and others [24–

391 26,28,58–61]. Generated sequence reads were analysed under the amplicon sequence variant framework

392 (ASV) using DADA2, in which statistically denoised forward and reverse read pairs were merged to

393 create error-corrected ASVs with single-nucleotide resolution [62,63]. We retained ASVs with length

394 ranging from 249 to 256 bp, which matches the desired length of the amplified V4 region, and chimeric

395 sequences were detected and removed independently for each sample. Taxonomic assignment (up to the

396 species level) was performed using the RDP Naïve Bayesian Classifier implemented in 'dada2' package,

397 on the SILVA v138 train dataset [64]. Further filtering removed ASVs matching the following criteria (1)

398 classified as 'Mitochondria' or 'Archaea', (2) unclassified at Kingdom or Phylum level, (3) identified as

399 kitome or contamination from mock community (except *Escherichia* and *Enterococcus* ASVs), or (4)

400    identified as low abundant singletons (abundance ≤10 counts and present in only one sample). This

401    resulted in 2,461 ASVs detected across 203 samples (68 participants), totalling 5,250,754 sequences.

402

403    Saliva and gut mucosal microbiomes were then analysed separately. For saliva microbiomes, we removed

404    singleton ASVs with abundances < 79 sequences and one sample with low sequencing depth (837

405    sequences). The filtered ASVs (n=865) were aligned using PASTA [65], and a maximum likelihood

406    phylogeny was constructed under the GTR+G model using IQ-Tree (with 1,000 rapid bootstrap) [66]. The

407    resulting phylogeny was used to transform the ASV count matrix into isometric log-ratio (ILR) 'balances'

408    (weighted log-ratio between two ASVs), using the "philr" package (with zero counts imputed using the

409    "cmultRepl" function) [60,67]. This transformation allowed statistical analyses to be performed

410    accurately in the Euclidean space [68]. Ordination was performed using principle coordinate analysis

411    (PCoA) on a calculated Euclidean distance matrix, implemented in package 'phyloseq'. To identify

412    covariates which explain the salivary microbiome structures, we performed redundancy analysis on the

413    'balance' value matrix of 62 samples with complete metadata (case-control grouping, age, sex, BMI,

414    presence of oral diseases, and status of active smoking in the last two years). We repeat the same

415    analytical procedures on the gut mucosal microbiome data. Low-abundance singleton ASVs (< 44

416    sequences) and seven samples with low sequencing depth (<1,300 sequences each) were removed,

417    retaining 1,073 ASVs across 129 samples for downstream analyses. We tested the association between

418    covariates and the gut mucosal microbiome structures using redundancy analysis, performed on the ILR-

419    transformed 'balance' values of 120 samples with complete metadata (sample type, age, sex, BMI, history

420    of diarrhea in seven days, diabetes, high blood pressure, active consumption of alcohol in the last two

421    years, and anatomical location of samples). The ILR-transformed values were used to calculate the

422    differences among the samples (beta-diversity), within and between participants. In addition, the gut

423    mucosal microbiomes (n=129) were clustered into community state types (CSTs) using the partition

424    around medoid (pam) algorithm on the calculated ILR-transformed distance matrix, with the optimal

425    number of CSTs (k=2) determined by gap statistic and average silhouette width (asw) [69]. The random

426  forest classification algorithm (10,000 trees) was then used to evaluate this clustering performance (in

427  term of error rate) and to identify 'balances' differentiating the two CSTs, using the 'rfsrc' function

428  implemented in package 'randomforestSRC' [70].

429

430  *Evaluating differential abundances*

431  In order to detect ASVs that showed significantly differential abundance between two examined groups,

432  we utilized the compositional data analysis approach implemented in ANCOMBC [24]. In addition, the

433  same comparisons were performed using DESeq2 and corncob to check for consistent results, as

434  recommended in recent benchmark studies [71,72]. The comparisons include salivary microbiomes in

435  cases (n=43) and controls (n=23); paired tumours (n=43) against adjacent non-tumours (n=43); paired

436  polyps (n=16) against non-polyp biopsies (n=16); tumours (n=43) against non-polyp biopsies (n=24);

437  tumours of cancer stage III-IV (n=24) against stage II (n=18). For paired comparison within cases and

438  controls, the model design was set to "~Patient + sample_type" to increase statistical power [73]. Multiple

439  hypothesis testing was corrected using Holm or Benjamini-Hochberg method, setting false discovery rate

440  as 0.05. ANCOMBC and corncob approaches were carried out using default parameters. For DESeq2,

441  count data were normalized using either the native negative binomial distribution (saliva, tumours vs.

442  biopsies) or the zero-inflated negative binomial (ZINB) distribution implemented in the package

443  'zinbwave' (paired samples in cases and controls, between cancer stages) [74]. Library size corrections

444  were estimated using DESeq2's 'poscounts' method. All comparisons were performed using likelihood

445  ratio test, and ASVs with adjusted p-value < 0.05 (and base mean >20 for DESeq2) were considered

446  significant hits. To minimize the number of false positives, ASVs which showed significant hits in at least

447  two tested methods were considered differentially abundant and included in final interpretation. We

448  manually compared the ASV sequences of interest to the expanded Human Oral Microbiome Database

449  (HOMD; www.homd.org/), and species identification was assigned if the ASV showed >99% nucleotide

450  similarity to that in the database.

451

*Correlation network*

We constructed a correlation network of gut mucosal microbiomes from colorectal cancer patients (n=86), using 117 most representative ASVs, defined as ones with abundance of at least 10 sequences detected in at least 15 samples. This filtering resulted in a median sample retainment rate of 77% [70% - 85%]. Zero counts were imputed using the "CmultRepl" function, and the correlation network was constructed using CCLasso, with 250 bootstrap and three-fold cross validation [27]. Interactions with adjusted p values < 0.01 and absolute correlation strength > 0.37 were considered significant hits. Additionally, a separate correlation network was inferred using SpiecEasi (Meinshausen-Buhlmann's neighbourhood selection, nlambda=20) on the same dataset [28]. Both these methods have been demonstrated to produce robust performance in a recent benchmark study [75]. To avoid spurious hits, only significant interactions detected by both the CCLasso and SpiecEasi approaches were included in the final visualization. We applied the same procedures to construct correlation networks of microbiomes in saliva samples (n=66, 115 ASVs) and controls' gut mucosa (n=43, 90 ASVs).


*Fusobacterium isolation and whole genome sequencing*

*Fusobacterium* isolation was performed on six selected case patients (P10, P16, P18, P28, P40, P46), whose *Fusobacterium* relative abundance in the tumour microbiome exceeded 0.5% as inferred by microbiome profiling. The respective samples (saliva, tumour, and non-tumour mucosa) were subjected to anaerobic culturing in a Whitley A35 anaerobic workstation (Don Whitley Scientific, UK) supplied with 5% $CO_2$, 10% $H_2$, and 85% nitrogen gas, following the *Fusobacterium* isolation procedures established previously [76]. Briefly, mucosal tissues were thawed on ice, and ~100 mg tissues were aseptically excised and anaerobically homogenized in Diluent A to maximize bacterial recovery. The suspension (100 µL) was then plated onto the selective media at two-fold and four-fold dilutions (EG agar supplemented with L-cysteine HCl, 50 ml/liter of defibrinated sheep blood, 7 mg/liter of crystal violet, 5 mg/liter of vancomycin, 30 mg/liter of neomycin, and 25 mg/liter of nalidixic acid; Sigma-Aldrich, Germany). Thawed saliva samples were plated directly on the selective media, either undiluted or at two-

478    fold dilutions. Plates were incubated at 37°C for 48-72h, and colonies (up to 10) resembling that of

479    *Fusobacterium* were picked from each plate and sub-cultured on new EG media to confirm purity and

480    select for single colonies. The isolate's taxonomic identities were queried using MALDI-TOF, and those

481    characterized as *Fusobacterium* species were retained. A total of 56 *Fusobacterium* isolates were

482    recovered and subjected to DNA extraction using the Wizard genomic extraction kit (Promega, USA). For

483    each isolate, 1 ng DNA was used to prepare the sequencing library using the Nextera XT library

484    preparation kit, following the manufacturer's instruction. Normalized libraries were pooled and

485    sequenced on an Illumina MiSeq platform to generate 250 bp paired-end reads.

486

487    ***Pangenome analysis, phylogenetic reconstruction and screening for virulence genes***

488    FASTQC was used to check the sequencing quality of each read pair [77], and Trimmomatic v0.36 was

489    used to trim sequencing adapters and low-quality reads (paired end option, SLIDINGWINDOW:10:22,

490    HEADCROP:15, MINLEN:50) [78]. For each isolate, the trimmed read set was input into Unicycler

491    v0.4.9 to construct the de novo assembly, using default parameters, and contigs of size over 500bp were

492    retained for downstream analyses [79]. The assemblies were checked for traces of contamination using

493    Checkm, and three assemblies were shown contaminated and discarded [80]. The resulting assemblies

494    were of adequate quality, with median size of 2,125,169 bp [IQR: 2,067,843 – 2,168,429], median

495    number of contigs of 133 [IQR: 86 – 173] and the median N50 of 35,535 bp [IQR: 21,685 – 51,953].

496    Prokka v1.13 was used to annotate the assemblies, using the well-annotated *F. nucleatum* 23726

497    (accessed via FusoPortal) as reference [81]. To provide preliminary taxonomic classification up to the

498    subspecies level, FastANI was used to calculate the average nucleotide identity (ANI) between the

499    individual assembly and a set of *Fusobacterium* references, with an ANI value ≥ 95% denoting a shared

500    species/subspecies [82]. The pangenomes of 57 *F. nucleatum/hwasookii* isolates (38 sequenced herein

501    plus 19 references) and 25 *F. periodonticum* isolates (15 sequenced herein plus 10 references) were

502    constructed separately using panX, which clusters orthologous proteins based on individual gene tree

503    construction and adaptive postprocessing instead of relying on fixed nucleotide identity cutoff [83]. The

504    respective core genome from each species complex was aligned, with invariant sites removed, producing

505    SNP alignments of 89,900 bp (*F. nucleatum/hwasookii* complex) and 106,738 bp (*F. periodonticum*

506    complex). These were input into RAxML to construct maximum likelihood phylogenies, under the

507    GTRGAMMA substitution model with 300 rapid bootstraps [84]. Using the pangenome analysis output

508    (clustered orthologous proteins and gene presence/absence matrix), we screened for the presence of

509    several known *Fusobacterium* virulence genes (*fap2, fadA, radD, cbpF*). The intact presence or synteny

510    of each genetic element was checked manually by gene alignment (Seaview) or genome visualization

511    (Artemis) tools [85,86]. Since Fap2 is a large protein (3,000 – 4,600 amino acid) with varying size among

512    the *Fusobacterium* species/subspecies, we checked for its intactness using ARIBA and mapping approach,

513    with clone-specific *fap2* serving as mapping reference [87]. Visualization of phylogenetic tree and

514    associated metadata was performed using package 'ggtree' [88]. Individual protein sets were aligned and

515    inspected in Seaview, and phylogenies were constructed in RAxML, using the PROTGAMMAGTR

516    model and 200 rapid bootstraps.

517

518    ***Intra-clonal variation examination***

519    To investigate intra-clonal variation with high confidence, we examined single nucleotide variants (SNV)

520    among isolates belonging to the same phylogenetic cluster (Figure 4 and Table 2), using the mapping

521    approach recommended previously [89]. For each phylogenetic cluster, trimmed fastq files from the

522    isolates were concatenated and input into Unicycler to construct a pan-assembly, with contigs less than

523    500 bp removed. This pan-assembly was ordered against an appropriate *Fusobacterium* reference using

524    ABACAS, depending on its respective species/subspecies, creating a pseudogenome reference [90].

525    Trimmed paired-end reads from each isolate were mapped against this reference using a custom wrapper

526    script. Briefly, mapping was conducted using BWA MEM algorithm and samtools v1.8 [91,92], with

527    duplicate reads removed using PICARD, followed by indel realignment by GATK [93]. Reads with

528    nonoptimal local alignment were subsequently removed using samclip (--max 50;

529    github.com/tseemann/samclip), to avoid false positives during variant calling. SNVs were detected using

530    the haplotype-based caller Freebayes [94], and low quality SNVs were removed using bcftools if they met

531    any of the following criteria: consensus quality < 30, mapping quality < 30, read depth < 4, ratio of SNVs

532    to reads at a position (AO/DP) < 85%, coverage on the forward or reverse strand < 1 [95]. Mapping

533    coverage at each position was inferred using samtools "depth" command (-a -Q 30). The bcftools

534    'consensus' command was used to generate a pseudosequence (with length identical to the mapping

535    reference), integrating the filtered SNVs and invariant sites, and masking the low mapping region (depth

536    < 4) and low-quality SNVs with 'N'. The presence of high quality SNVs were validated by manual

537    visualization of output bam files in Artemis, and SNV pertaining to recombination, transposons, plasmids,

538    or repetitive elements were excluded from interpretation.

539

540    *Data availability*

541    Raw sequence data are available in the NCBI Sequence Read Archive, including ones for 16S rRNA

542    sequencing (BioProject PRJNA791834) and *Fusobacterium* whole genome sequencing (BioProject

543    PRJNA791829). Source data and R codes used for the microbiome analysis and visualization are

544    deposited in Github (https://github.com/Hao-Chung/Vietnam_CRC_microbiome).

545

546 **Figure legends**

547 **Figure 1** The salivary and gut mucosa microbiomes of colorectal cancer patients. Principal coordinate

548 analyses (PCoA), conducted on phylogenetic-assisted isometric log-ratio (PhILR) transformed data, of (A)

549 66 salivary microbiomes, and (B) 129 gut mucosa microbiomes, with different CRC groups and sample

550 types denoted by different colours (see Keys; biopsies and polyps collected from controls, nontumours

551 and tumours collected from cases). (C) Boxplot showing the distribution of pairwise beta-diversity,

552 calculated on PhILR transformed values, observed in each gut microbiome category. (D) Heatmap

553 displaying the proportional abundances of 24 most abundant genera, with headers showing the samples'

554 community state type (CST): CST1 (light gray), CST2 (dark gray), and the corresponding sample type:

555 biopsy (light blue), polyp (dark blue), nontumour (pink), tumour (dark red). Genera were coloured

556 according to their classifications at Phylum level (see Keys). Genera in black box represent ones with

557 probable origin from the oral cavity.

558

559 **Figure 2** Bacterial taxa significantly abundant among the examined classes. Taxa, or amplicon sequence

560 variants (ASVs), were determined as significant and visualized in (A) and (B) if they were detected in at

561 least two of the three tested approaches (ANCOMBC, DESeq2, corncob; adjusted p-value $\leq 0.05$). (A)

562 Log2 fold change of ASVs that differ between paired tumour and non-tumour mucosal microbiomes from

563 case participants, using the full model 'Patient + sample type' (n=86). (B) Log2 fold change of ASVs that

564 differ between tumour and biopsy (control) mucosal microbiomes (n=67). Log2 fold change was derived

565 from ANCOMBC test output, and taxa of oral origin were coloured in pink. (C) Relative abundance of

566 ASVs assigned as *Fusobacterium mortiferum* (n=14), *Fusobacterium nucleatum* (n=14), *Leptotrichia*

567 (n=16), and *Collinsella* (n=14) in the tumour and nontumour mucosal microbiomes, stratified by cancer

568 stages (III-IV vs. II).

569

570 **Figure 3** Correlation network of colorectal cancer gut mucosal microbiomes. The network was

571 constructed from 117 most representative ASVs sampled from 86 mucosal microbiomes, outlining

572    significant interactions detected by both CCLasso (p value ≤ 0.01 and absolute correlation strength > 0.37)

573    and SpiecEasi. Positive and negative interactions were coloured as red and blue lines respectively, with

574    line weight proportional to correlation strength. The ASVs (nodes) were coloured based on taxonomic

575    family (see Legend), with sizes proportional to their relative abundances. The light green shaded area

576    entails ASVs identified as members of the human oral microbiome (comparison with expanded Human

577    Oral Microbiome Database); the blue shaded area covers ASVs identified as gut anaerobic commensals

578    (Lachnospirales, Bacteroidales, Bifidobacteriales, Oscillospirales); and gray shaded area covers other

579    tumour-associated ASVs (as identified in Figure 2A).

580

581    **Figure 4** Maximum likelihood phylogenies of *Fusobacterium* isolates from this study. (A) *F. nucleatum*

582    species phylogeny constructed from the alignment of 516 core genes (89,900 variant sites; N=57), using *F.*

583    *hwasookii* clade as an outgroup. (B) *F. periodonticum* species phylogeny constructed from alignment of

584    863 core genes (106,738 variant sites, N=26). Red circles at internal nodes denote bootstrap values ≥ 80.

585    The associated metadata on the right describe the patient ID and clinical origin of isolates (where

586    reference genomes were left blank), and the genomic presence of several virulence factors (*fap2*, *fadA*,

587    *fadA2*, *radD*, *fadA3*, *cbpF*). Light blue shaded area covers isolates identified as novel *F. periodonticum*

588    subspecies. The scale bars denote the estimated number of substitutions.

589

590 **Table 1** Baseline characteristics of patients recruited in this study. Overweight/obesity was classified using WHO recommendation for Asian
591 populations. Oral diseases include self-reported gingivitis, periodontitis or halitosis. The number in each cell refers to median with interquartile
592 range in brackets, percentage or count number for each category.
593

|  | **CRC cases (n=42)** | **Controls (n=21)** | **p-value** |
|---|---|---|---|
| **Age** | 64 [54 - 69] | 60 [53 - 66] | 0.359 |
| **Male sex** | 62% | 76% | 0.395 |
| **BMI** | 22.9 [20.85 - 24.95] | 22.2 [21.1 - 23.4] | 0.387 |
| **Overweight/obesity** | 47.60% | 33% | 0.409 |
| **Diabetes** | 19% | 19% | 1 |
| **High blood pressure** | 52% | 47.60% | 0.79 |
| **Active smoking in the last two years** | 21.40% | 19% | 1 |
| **Oral diseases** | 33% | 38% | 0.782 |
| **Family history of cancer** | 19% | 19% | 1 |
| **Location of sampled mucosa** |  |  | 0.533 |
| *Descending colon* | 7 | 3 |  |
| *Sigmoid colon* | 28 | 12 |  |
| *Rectum* | 7 | 6 |  |
| **Size of tumour/polyp (cm)** | 5 [4 - 5.75] | 1 [0.7 - 1.2] |  |
| **TNM stage of cancer** | II (18), III (20), IV (4) |  |  |
| **Polyp dysplasia grade** |  | low (4), none (17) |  |

594
595

596  **Table 2** Summary of *Fusobacterium* isolates recovered from this study. Species names in italic represent *Fusobacterium nucleatum* subspecies,
597  and (*) denotes a potential new subspecies. *Fusobacterium* relative abundance, inferred from 16S rRNA gene profiling results, showed the values
598  for non-tumour/tumour samples for each cancer patient. SNV: single nucleotide variant.
599

| Phylogenetic cluster | Species | Isolation source | Number of isolates | intraclonal SNV | Patient ID | Tumour location | Cancer stage | *Fusobacterium* relative abundance |
|---|---|---|---|---|---|---|---|---|
| T10_Fa1 | *animalis* | tumour | 2 | 5 | | | | |
| TH10_Fa2 | *animalis* | tumour-nontumour | 5 | 2 | P10 | Sigmoid | IIB | 0.026/0.24 |
| S16-12 | *polymorphum* | oral | 1 | | | | | |
| S16-17 | periodonticum | oral | 1 | | | | | |
| T16_Fp | *polymorphum* | tumour | 2 | 3 | | | | |
| H16_Fa | *animalis** | nontumour | 2 | 0 | | | | |
| S16_Fa | *animalis* | oral | 2 | 10 | | | | |
| TH16_Fvi | *vincentii* | tumour-nontumour | 6 | 1 | P16 | Sigmoid | IIB | 0.1/0.34 |
| H18-18 | *animalis* | nontumour | 1 | | | | | |
| S18-79 | *polymorphum* | oral | 1 | | | | | |
| S18-78 | *hwasookii* | oral | 1 | | | | | |
| S18-66 | periodonticum | oral | 1 | | | | | |
| S18_Fperi1 | periodonticum | oral | 2 | 1 | | | | |
| S18_Fperi2 | periodonticum | oral | 2 | 1 | | | | |
| S18_Fh | *hwasookii* | oral | 2 | 4 | | | | |
| TH18_novelFperi | novel periodonticum* | tumour-nontumour | 6 | 2 | P18 | Descending | IIA | 0.003/0.008 |
| S28-2 | periodonticum | oral | 1 | | | | | |
| T28_Fn | *nucleatum* | tumour | 3 | 1 | P28 | Sigmoid | IV | 0.14/0.39 |
| S40-28 | *animalis* | oral | 1 | | | | | |
| TH40_novelFperi | novel periodonticum* | tumour-nontumour | 2 | 2 | | | | |
| S40_Fp | *polymorphum* | oral | 2 | 2 | P40 | Sigmoid | IIIB | 0.27/0.1 |
| S46-13 | *vincentii* | oral | 1 | | P46 | Sigmoid | IIIB | 0.038/0.2 |

| | | | | | |
|---|---|---|---|---|---|
| S46-17 | *vincentii* | oral | 1 | | |
| S46-7 | *polymorphum* | oral | 1 | | |
| S46-16 | *polymorphum* | oral | 1 | | |
| T46_Fa | *animalis* | tumour | 3 | 3 | |

600
601

602

**Disclosure of potential conflicts of interest**

610

611    The authors report no potential conflicts of interest.

612

**Consent for publication**

613

614    This study does not publish any identifiable details related to the participant.

615

**References**

1.  Bray F, Ferlay J, Soerjomataram I, Siegel RL, Torre LA, Jemal A: **Global cancer statistics 2018: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries.** *CA Cancer J Clin* 2018, **68**:394–424.

2.  Arnold M, Sierra MS, Laversanne M, Soerjomataram I, Jemal A, Bray F: **Global patterns and trends in colorectal cancer incidence and mortality**. *Gut* 2017, **66**:683–691.

3.  Keum NN, Giovannucci E: **Global burden of colorectal cancer: emerging trends, risk factors and prevention strategies**. *Nat Rev Gastroenterol Hepatol* 2019, **16**:713–732.

4.  Schirmer M, Smeekens SP, Vlamakis H, Jaeger M, Oosting M, Franzosa EA, Jansen T, Jacobs L, Bonder MJ, Kurilshikov A, et al.: **Linking the Human Gut Microbiome to Inflammatory Cytokine Production Capacity**. *Cell* 2016, **167**:1125-1136.e8.

5.  Song M, Chan AT: **Environmental Factors, Gut Microbiota, and Colorectal Cancer Prevention**. *Clin Gastroenterol Hepatol* 2019, **17**:275–289.

6.  Castellarin M, Warren L, Freeman JD, Dreolini L, Krzywinski M, Strauss J, Barnes R, Watson P, Allen-vercoe E, Moore RA, et al.: **Fusobacterium nucleatum infection is prevalent in human colorectal carcinoma**. *Genome Res* 2012, **22**:299–306.

7.  Kostic AD, Gevers D, Pedamallu CS, Kostic AD, Gevers D, Pedamallu CS, Michaud M, Duke F, Earl AM, Ojesina AI, et al.: **Genomic analysis identifies association of Fusobacterium with colorectal carcinoma**. *Genome Res* 2012, **22**:292–298.

8.  Flemer B, Lynch DB, Brown JMR, Jeffery IB, Ryan FJ, Claesson MJ, O'Riordain M, Shanahan F, O'Toole PW: **Tumour-associated and non-tumour-associated microbiota in colorectal cancer**. *Gut* 2017, **66**:633–643.

9.  Wirbel J, Pyl PT, Kartal E, Zych K, Kashani A, Milanese A, Fleck JS, Voigt AY, Palleja A, Ponnudurai R, et al.: **Meta-analysis of fecal metagenomes reveals global microbial signatures that are specific for colorectal cancer**. *Nat Med* 2019, **25**:679–689.

10. Thomas AM, Manghi P, Asnicar F, Pasolli E, Armanini F, Zolfo M, Beghini F, Manara S, Karcher

N, Pozzi C, et al.: **Metagenomic analysis of colorectal cancer datasets identifies cross-cohort microbial diagnostic signatures and a link with choline degradation**. *Nat Med* 2019, **25**:667–678.

11.  Brennan CA, Garrett WS: **Fusobacterium nucleatum — symbiont, opportunist and oncobacterium**. *Nat Rev Microbiol* 2018, doi:10.1038/s41579-018-0129-6.

12.  Rubinstein MR, Wang X, Liu W, Hao Y, Cai G, Han YW: **Fusobacterium nucleatum Promotes Colorectal Carcinogenesis by Modulating E-Cadherin/β-Catenin Signaling via its FadA Adhesin**. *Cell Host Microbe* 2013, **14**:195–206.

13.  Abed J, Emgård JEM, Zamir G, Faroja M, Almogy G, Grenov A, Sol A, Naor R, Pikarsky E, Atlan KA, et al.: **Fap2 Mediates Fusobacterium nucleatum Colorectal Adenocarcinoma Enrichment by Binding to Tumor-Expressed Gal-GalNAc**. *Cell Host Microbe* 2016, **20**:215–225.

14.  Casasanta MA, Yoo CC, Udayasuryan B, Sanders BE, Umana A, Zhang Y, Peng H, Duncan AJ, Wang Y, Li L, et al.: **Fusobacterium nucleatum host cell binding and invasion induces IL-8 and CXCL1 secretion that drives colorectal cancer cell migration**. *Sci Signal* 2020, **1**:1–13.

15.  Geng F, Zhang Y, Lu Z, Zhang S, Pan Y: **Fusobacterium nucleatum Caused DNA Damage and Promoted Cell Proliferation by the Ku70/p53 Pathway in Oral Cancer Cells**. *DNA Cell Biol* 2020, **39**:144–151.

16.  Guo P, Tian Z, Kong X, Yang L, Shan X, Dong B, Ding X, Jing X, Jiang C, Jiang N, et al.: **FadA promotes DNA damage and progression of Fusobacterium nucleatum-induced colorectal cancer through up-regulation of chk2**. *J Exp Clin Cancer Res* 2020, **39**:1–13.

17.  Yu TC, Guo F, Yu Y, Sun T, Ma D, Han J, Qian Y, Kryczek I, Sun D, Nagarsheth N, et al.: **Fusobacterium nucleatum Promotes Chemoresistance to Colorectal Cancer by Modulating Autophagy**. *Cell* 2017, **170**:548-563.e16.

18.  Mima K, Nishihara R, Qian ZR, Cao Y, Sukawa Y, Nowak JA, Yang J, Dou R, Masugi Y, Song M, et al.: **Fusobacterium nucleatum in colorectal carcinoma tissue and patient prognosis**. *Gut*

672      2016, **65**:1973–1980.

673    19.    Serna G, Ruiz-Pace F, Hernando J, Alonso L, Fasani R, Landolfi S, Comas R, Jimenez J, Elez E,

674        Bullman S, et al.: **Fusobacterium nucleatum persistence and risk of recurrence after**

675        **preoperative treatment in locally advanced rectal cancer**. *Ann Oncol* 2020, **31**:1366–1375.

676    20.    Salvucci M, Crawford N, Stott K, Bullman S, Longley DB, Prehn JHM: **Patients with**

677        **mesenchymal tumours and high Fusobacteriales prevalence have worse prognosis in**

678        **colorectal cancer (CRC)**. *Gut* 2021, doi:10.1136/gutjnl-2021-325193.

679    21.    Bullman S, Pedamallu CS, Sicinska E, Clancy TE, Zhang X, Cai D, Neuberg D, Huang K,

680        Guevara F, Nelson T, et al.: **Analysis of Fusobacterium persistence and antibiotic response in**

681        **colorectal cancer - supplementary**. *Science (80- )* 2017, **358**:1443–1448.

682    22.    Althoff T, Sosič R, Hicks JL, King AC, Delp SL, Leskovec J: **Large-scale physical activity data**

683        **reveal worldwide activity inequality**. *Nature* 2017, **547**:336–339.

684    23.    Nguyen SM, Deppen S, Nguyen GH, Pham DX, Bui TD, Tran T Van: **Projecting Cancer**

685        **Incidence for 2025 in the 2 Largest Populated Cities in Vietnam**. *Cancer Control* 2019, **26**:1–

686        13.

687    24.    Lin H, Peddada S Das: **Analysis of compositions of microbiomes with bias correction**. *Nat*

688        *Commun* 2020, **11**:1–11.

689    25.    Love MI, Huber W, Anders S: **Moderated estimation of fold change and dispersion for RNA-**

690        **seq data with DESeq2**. *Genome Biol* 2014, **15**:550.

691    26.    Martin BD, Witten D, Willis AD: **Modeling microbial abundances and dysbiosis with beta-**

692        **binomial regression**. *Ann Appl Stat* 2020, **14**:94–115.

693    27.    Fang H, Huang C, Zhao H, Deng M: **CCLasso: Correlation inference for compositional data**

694        **through Lasso**. *Bioinformatics* 2015, **31**:3172–3180.

695    28.    Kurtz ZD, Müller CL, Miraldi ER, Littman DR, Blaser MJ, Bonneau RA: **Sparse and**

696        **Compositionally Robust Inference of Microbial Ecological Networks**. *PLOS Comput Biol*

697        2015, **11**:e1004226.

698  29.  Brewer ML, Dymock D, Brady RL, Singer BB, Virji M, Hill DJ: **Fusobacterium spp. target**

699       **human CEACAM1 via the trimeric autotransporter adhesin CbpF**. *J Oral Microbiol* 2019, **11**.

700  30.  Wu C, Chen YW, Scheible M, Chang C, Wittchen M, Lee JH, Luong TT, Tiner BL, Tauch A, Das

701       A, et al.: **Genetic and molecular determinants of polymicrobial interactions in Fusobacterium**

702       **nucleatum**. *Proc Natl Acad Sci U S A* 2021, **118**.

703  31.  Galaski J, Shhadeh A, Umaña A, Yoo CC, Arpinati L, Isaacson B, Berhani O, Singer BB, Slade

704       DJ, Bachrach G, et al.: **Fusobacterium nucleatum CbpF Mediates Inhibition of T Cell**

705       **Function Through CEACAM1 Activation**. *Front Cell Infect Microbiol* 2021, **11**:1–8.

706  32.  Warren RL, Freeman DJ, Pleasance S, Watson P, Moore RA, Cochrane K, Allen-vercoe E, Holt

707       RA: **Co-occurrence of anaerobic bacteria in colorectal carcinomas**. *Microbiome* 2013, **1**:1–12.

708  33.  Drewes JL, White JR, Dejea CM, Fathi P, Iyadorai T, Vadivelu J, Roslani AC, Wick EC,

709       Mongodin EF, Loke MF, et al.: **High-resolution bacterial 16S rRNA gene profile meta-analysis**

710       **and biofilm status reveal common colorectal cancer consortia**. *npj Biofilms Microbiomes* 2017,

711       **3**.

712  34.  Meng Q, Gao Q, Mehrazarin S, Tangwanichgapong K, Wang Y, Huang Y, Pan Y, Robinson S,

713       Liu Z, Zangiabadi A, et al.: **Fusobacterium nucleatum secretes amyloid-like FadA to enhance**

714       **pathogenicity**. *EMBO Rep* 2021, **22**:1–19.

715  35.  Gao Z, Guo B, Gao R, Zhu Q, Qin H: **Microbiota disbiosis is associated with colorectal cancer**.

716       *Front Microbiol* 2015, **6**:1–9.

717  36.  Torres PJ, Fletcher EM, Gibbons SM, Bouvet M, Doran KS, Kelley ST: **Characterization of the**

718       **salivary microbiome in patients with pancreatic cancer**. *PeerJ* 2015, **2015**:1–16.

719  37.  Amer A, Galvin S, Healy CM, Moran GP: **The microbiome of potentially malignant oral**

720       **leukoplakia exhibits enrichment for Fusobacterium, Leptotrichia, Campylobacter, and**

721       **Rothia species**. *Front Microbiol* 2017, **8**:1–9.

722  38.  Xia X, Wu WKK, Wong SH, Liu D, Kwong TNY, Nakatsu G, Yan PS, Chuang YM, Chan MWY,

723       Coker OO, et al.: **Bacteria pathogens drive host colonic epithelial cell promoter**

hypermethylation of tumor suppressor genes in colorectal cancer. *Microbiome* 2020, **8**:1–13.

39. Liang JQ, Li T, Nakatsu G, Chen YX, Yau TO, Chu E, Wong S, Szeto CH, Ng SC, Chan FKL, et al.: **A novel faecal Lachnoclostridium marker for the non-invasive diagnosis of colorectal adenoma and cancer**. *Gut* 2019, **69**:1248–1257.

40. Chung The H, Sessions PF de, Jie S, Thanh DP, Thompson CN, Minh CNN, Chu CW, Tran T-A, Thomson NR, Thwaites GE, et al.: **Assessing gut microbiota perturbations during the early phase of infectious diarrhea in Vietnamese children**. *Gut Microbes* 2018, **9**:38–54.

41. Liang S, Mao Y, Liao M, Xu Y, Chen Y, Huang X, Wei C, Wu C, Wang Q, Pan X, et al.: **Gut microbiome associated with APC gene mutation in patients with intestinal adenomatous polyps**. *Int J Biol Sci* 2020, **16**:135–146.

42. Wei PL, Hung CS, Kao YW, Lin YC, Lee CY, Chang TH, Shia BC, Lin JC: **Classification of changes in the fecal microbiota associated with colonic adenomatous polyps using a long-read sequencing platform**. *Genes (Basel)* 2020, **11**:1–14.

43. He Y, Mujagond P, Tang W, Wu W, Zheng H, Chen X, Chen M, Ma W, Chen G, Zhou H: **Non-nucleatum Fusobacterium species are dominant in the Southern Chinese population with distinctive correlations to host diseases compared with F. nucleatum**. *Gut* 2021, **70**:810–812.

44. Yeoh YK, Chen Z, Wong MCS, Hui M, Yu J, Ng SC, Sung JJY, Chan FKL, Chan PKS: **Southern Chinese populations harbour non-nucleatum Fusobacteria possessing homologues of the colorectal cancer-associated FadA virulence factor**. *Gut* 2020, **69**:1998–2007.

45. Robrish SA, Oliver C, Thompson J: **Sugar metabolism by fusobacteria: regulation of transport, phosphorylation, and polymer formation by Fusobacterium mortiferum ATCC 25557.** *Infect Immun* 1991, **59**:4547–54.

46. Lagier JC, Dubourg G, Million M, Cadoret F, Bilen M, Fenollar F, Levasseur A, Rolain JM, Fournier PE, Raoult D: **Culturing the human microbiota and culturomics**. *Nat Rev Microbiol* 2018, **16**:540–550.

47. Chung The H, Nguyen Ngoc Minh C, Tran Thi Hong C, Nguyen Thi Nguyen T, Pike LJ, Zellmer

750    C, Pham Duc T, Tran T-A, Ha Thanh T, Van MP, et al.: **Exploring the Genomic Diversity and**

751    **Antimicrobial Susceptibility of Bifidobacterium pseudocatenulatum in a Vietnamese**

752    **Population**. *Microbiol Spectr* 2021, doi:10.1128/spectrum.00526-21.

753    48.    Richardson M, Ren J, Rubinstein MR, Taylor JA, Friedman RA, Shen B, Han YW: **Analysis of**

754    **16S rRNA genes reveals reduced Fusobacterial community diversity when translocating**

755    **from saliva to GI sites**. *Gut Microbes* 2020, **12**:1–13.

756    49.    Abed J, Maalouf N, Manson AL, Earl AM, Parhi L, Emgård JEM, Klutstein M, Tayeb S, Almogy

757    G, Atlan KA, et al.: **Colon Cancer-Associated Fusobacterium nucleatum May Originate From**

758    **the Oral Cavity and Reach Colon Tumors via the Circulatory System**. *Front Cell Infect*

759    *Microbiol* 2020, **10**:1–12.

760    50.    Ailloud F, Didelot X, Woltemate S, Pfaffinger G, Overmann J, Bader RC, Schulz C, Malfertheiner

761    P, Suerbaum S: **Within-host evolution of Helicobacter pylori shaped by niche-specific**

762    **adaptation, intragastric migrations and selective sweeps**. *Nat Commun* 2019, **10**:2273.

763    51.    Slade DJ: **New Roles for Fusobacterium nucleatum in Cancer: Target the Bacteria, Host, or**

764    **Both?** *Trends in Cancer* 2021, **7**:185–187.

765    52.    Mima K, Cao Y, Chan AT, Qian ZR, Nowak JA, Masugi Y, Shi Y, Song M, Da Silva A, Gu M, et

766    al.: **Fusobacterium nucleatum in Colorectal Carcinoma Tissue According to Tumor Location**.

767    *Clin Transl Gastroenterol* 2016, **7**.

768    53.    WHO expert consultation: **Appropriate body-mass index for Asian populations and its**

769    **implications for policy and intervention strategies**. *Lancet* 2004, **363**:157–63.

770    54.    American Joint Committee on Cancer: *Chapter 20 - Colon and Rectum*. Springer; 2017.

771    55.    Kozich JJ, Westcott SL, Baxter NT, Highlander SK, Schloss PD: **Development of a dual-index**

772    **sequencing strategy and curation pipeline for analyzing amplicon sequence data on the**

773    **MiSeq Illumina sequencing platform.** *Appl Environ Microbiol* 2013, **79**:5112–20.

774    56.    Gohl D, Gohl DM, MacLean A, Hauge A, Becker A, Walek D, Beckman KB: **An optimized**

775    **protocol for high-throughput amplicon-based microbiome profiling**. *Protoc Exch* 2016,

776         doi:10.1038/protex.2016.030.

777    57.    Eisenhofer R, Minich JJ, Marotz C, Cooper A, Knight R, Weyrich LS: **Contamination in Low**

778         **Microbial Biomass Microbiome Studies: Issues and Recommendations**. *Trends Microbiol*

779         2019, **27**:105–117.

780    58.    R Core Team: **R: A language and environment for statistical computing.** 2016,

781    59.    McMurdie PJ, Holmes S: **phyloseq: an R package for reproducible interactive analysis and**

782         **graphics of microbiome census data.** *PLoS One* 2013, **8**:e61217.

783    60.    Silverman JD, Washburne AD, Mukherjee S, David LA: **A phylogenetic transform enhances**

784         **analysis of compositional microbiota data**. *Elife* 2017, **6**:1–20.

785    61.    Oksanen J, Blanchet G, Kindt R, Legendre P, Minchin P, O'Hara R, Simpson G, Solymos P,

786         Stevens M, Wagner H: **vegan: Community Ecology Package. R package version 2.3-5**. 2016,

787    62.    Callahan BJ, Sankaran K, Fukuyama JA, McMurdie PJ, Holmes SP: **Bioconductor workflow for**

788         **microbiome data analysis: from raw reads to community analyses.** *F1000Research* 2016,

789         **5**:1492.

790    63.    Callahan BJ, McMurdie PJ, Holmes SP: **Exact sequence variants should replace operational**

791         **taxonomic units in marker-gene data analysis**. *ISME J* 2017, doi:10.1038/ismej.2017.119.

792    64.    Wang Q, Garrity GM, Tiedje JM, Cole JR: **Naive Bayesian classifier for rapid assignment of**

793         **rRNA sequences into the new bacterial taxonomy.** *Appl Environ Microbiol* 2007, **73**:5261–7.

794    65.    Mirarab S, Nguyen N, Guo S, Wang L-S, Kim J, Warnow T: **PASTA: Ultra-Large Multiple**

795         **Sequence Alignment for Nucleotide and Amino-Acid Sequences**. *J Comput Biol* 2015, **22**:377–

796         386.

797    66.    Nguyen LT, Schmidt HA, Von Haeseler A, Minh BQ: **IQ-TREE: A fast and effective stochastic**

798         **algorithm for estimating maximum-likelihood phylogenies**. *Mol Biol Evol* 2015, **32**:268–274.

799    67.    Washburne AD, Morton JT, Sanders J, McDonald D, Zhu Q, Oliverio AM, Knight R: **Methods**

800         **for phylogenetic analysis of microbiome data**. *Nat Microbiol* 2018, **3**:652–661.

801    68.    Filzmoser P, Hron K, Reimann C: **The bivariate statistical analysis of environmental**

802    **(compositional) data**. *Sci Total Environ* 2010, **408**:4230–4238.

803    69.    Tibshirani R, Walther G, Hastie T: **Estimating the number of clusters in a data set via the gap**

804    **statistic**. *J R Stat Soc* 2001, **63**:411–423.

805    70.    Ishwaran H, Kogalur UB: **Random Forests for Survival, Regression and Classification (RF-**

806    **SRC), R package version 2.2.0**. 2016,

807    71.    Calgaro M, Romualdi C, Waldron L, Risso D, Vitulo N: **Assessment of statistical methods from**

808    **single cell, bulk RNA-seq, and metagenomics applied to microbiome data**. *Genome Biol* 2020,

809    **21**:1–31.

810    72.    Nearing JT, Douglas GM, Hayes M, Macdonald J, Desai D, Allward N, Jones CMA, Wright R,

811    Dhanani A, Comeau AM, et al.: **Microbiome differential abundance methods produce**

812    **disturbingly different results across 38 datasets**. *bioRxiv* 2021,

813    73.    Stevens JR, Herrick JS, Wolff RK, Slattery ML: **Power in pairs: Assessing the statistical value**

814    **of paired samples in tests for differential expression**. *BMC Genomics* 2018, **19**:1–13.

815    74.    Risso D, Perraudeau F, Gribkova S, Dudoit S, Vert JP: **A general and flexible method for signal**

816    **extraction from single-cell RNA-seq data**. *Nat Commun* 2018, **9**.

817    75.    Hirano H, Takemoto K: **Difficulty in inferring microbial community structure based on co-**

818    **occurrence network approaches**. *BMC Bioinformatics* 2019, **20**:1–14.

819    76.    Komiya Y, Shimomura Y, Higurashi T, Sugi Y, Arimoto J, Umezawa S, Uchiyama S, Matsumoto

820    M, Nakajima A: **Patients with colorectal cancer have identical strains of Fusobacterium**

821    **nucleatum in their colorectal cancer and oral cavity**. *Gut* 2018, **0**:7–9.
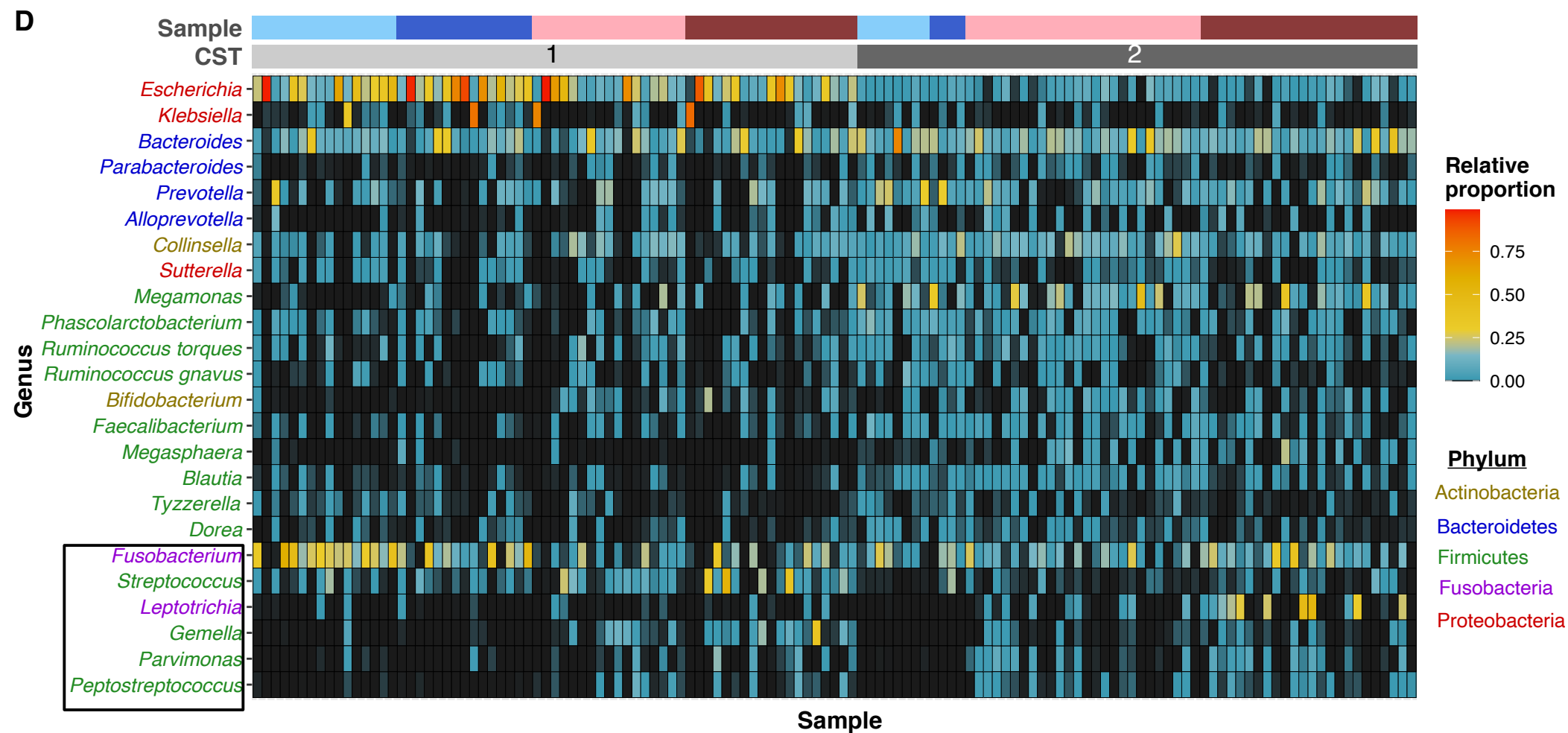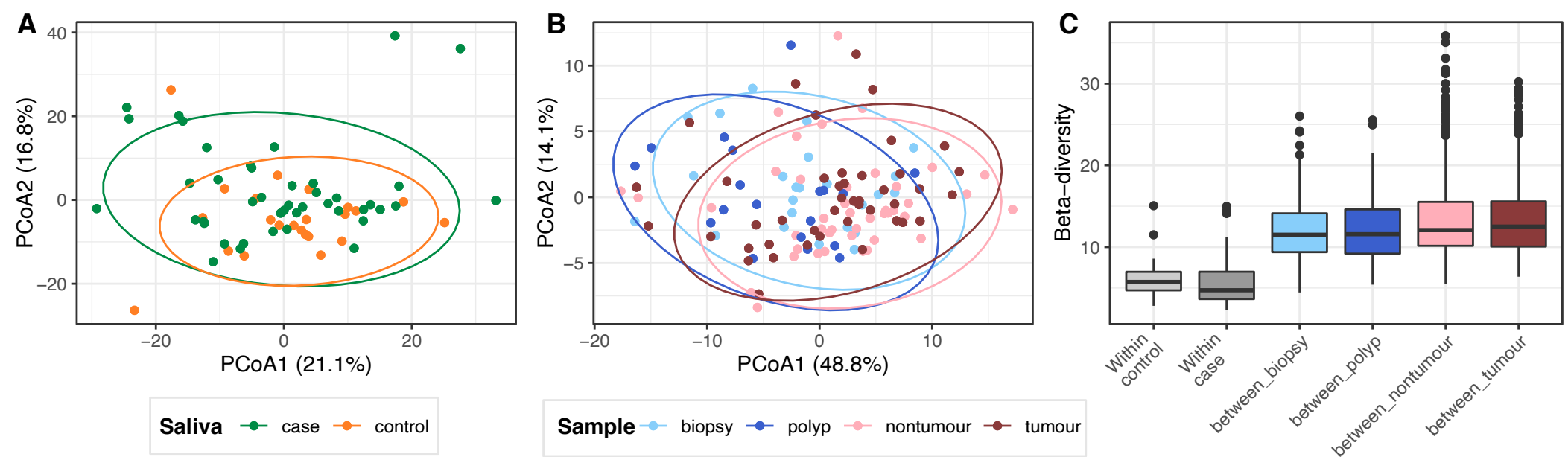
822    77.    Andrews S: **FastQC: A Quality Control Tool for High Throughput Sequence Data**. [date
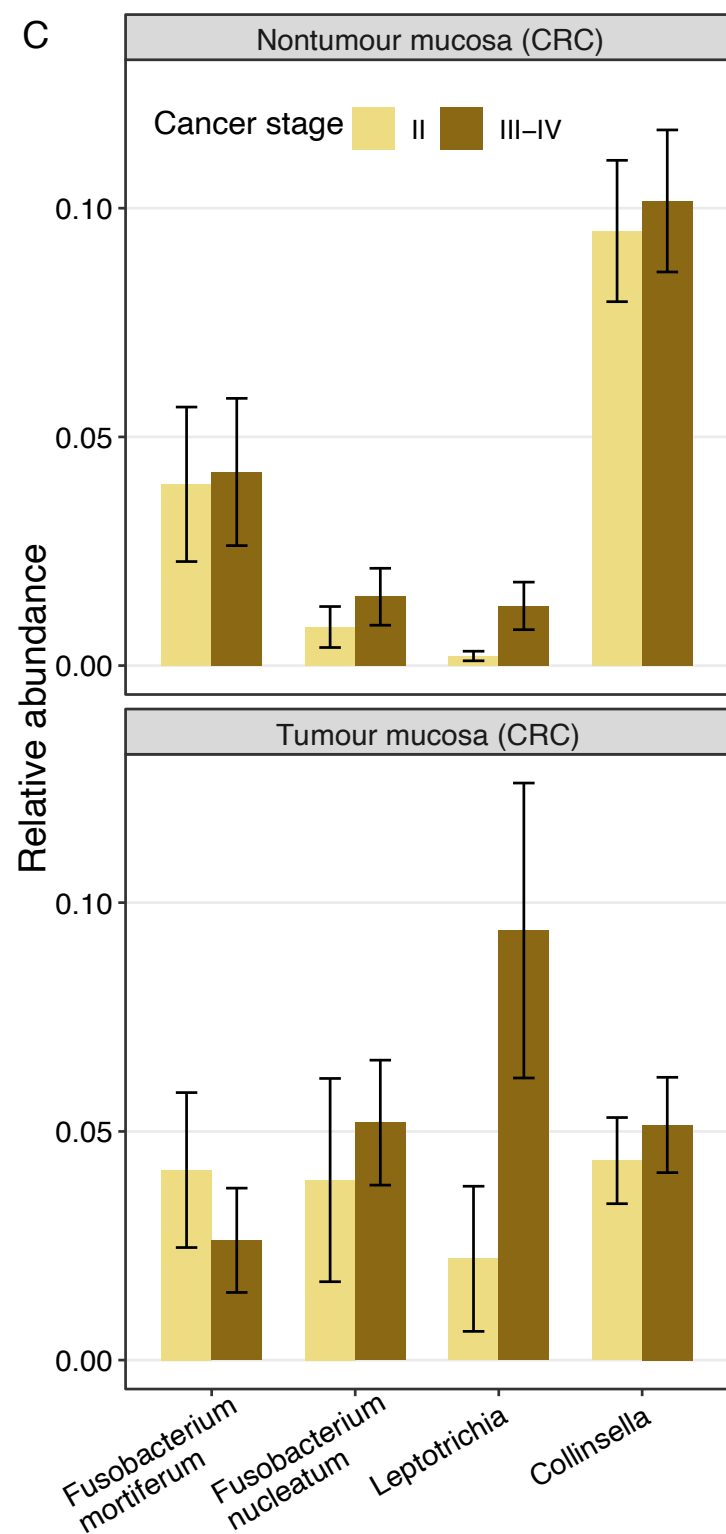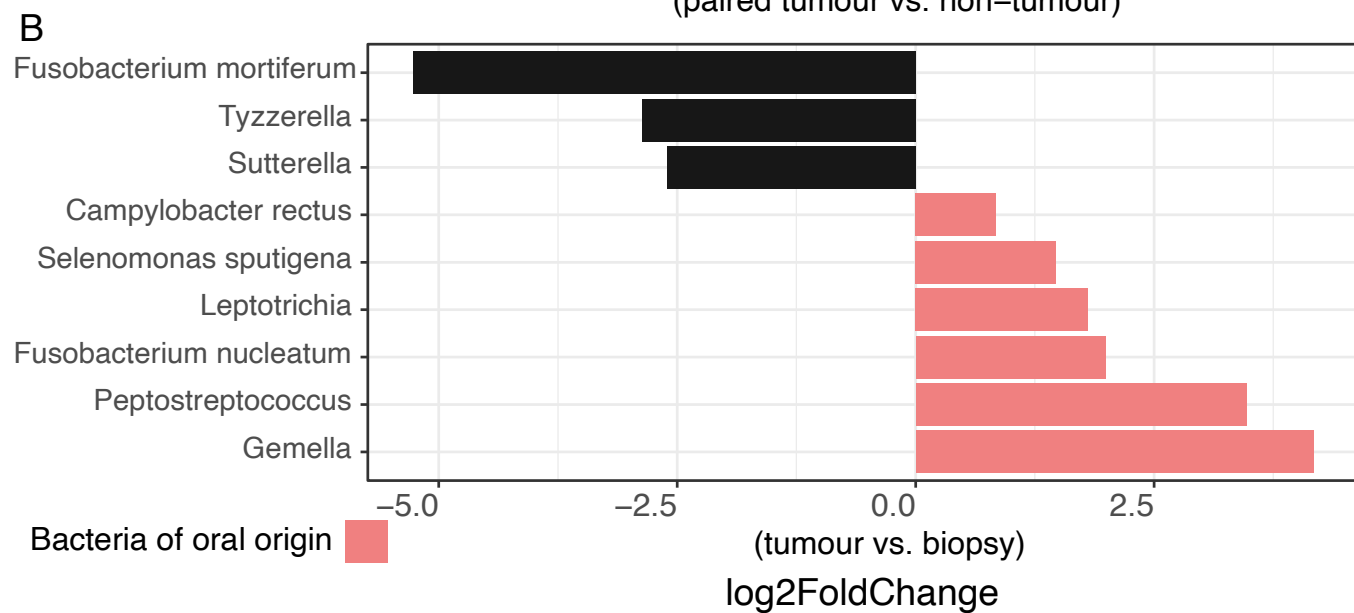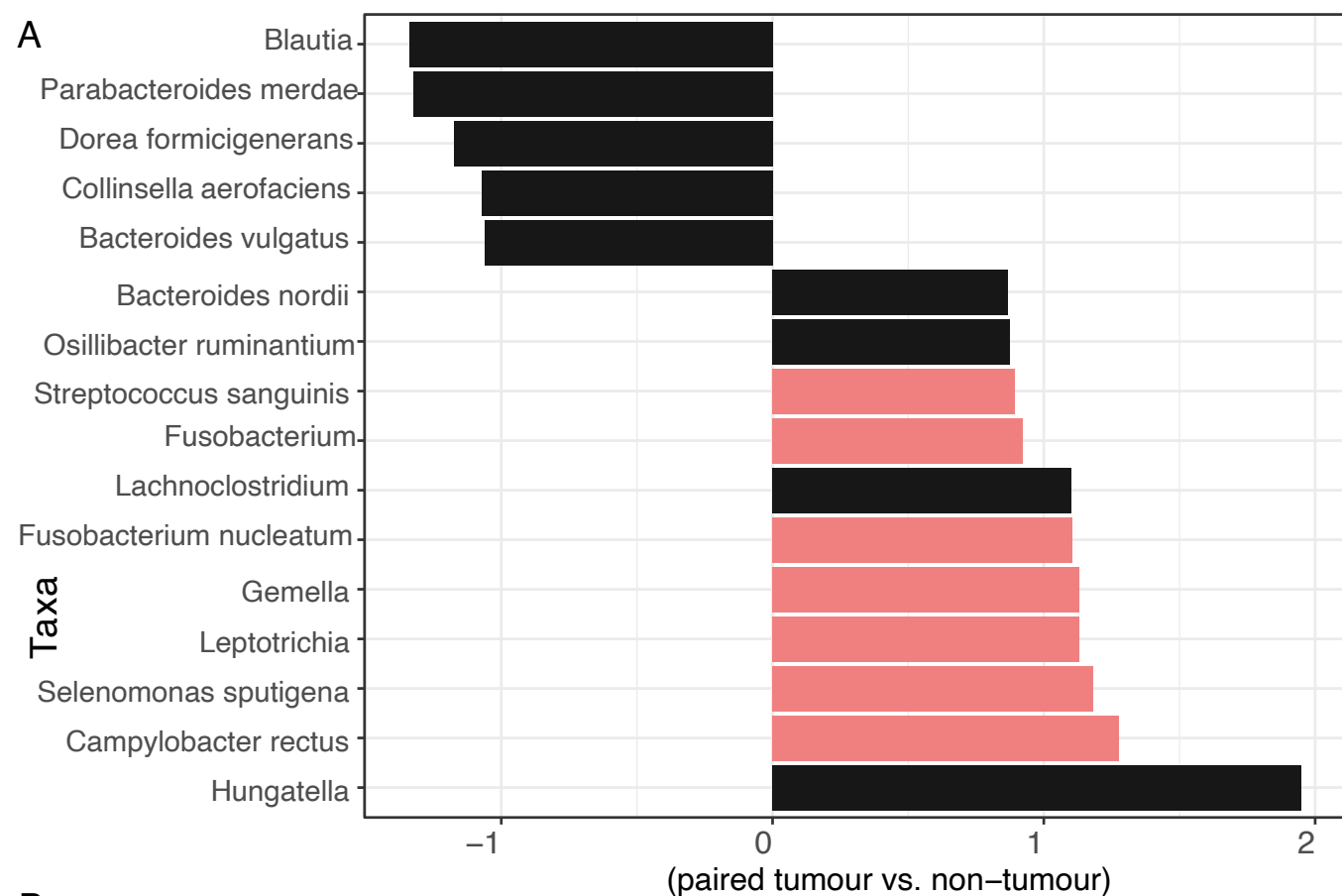
823    unknown],

824    78.    Bolger AM, Lohse M, Usadel B: **Trimmomatic: A flexible trimmer for Illumina sequence data**.
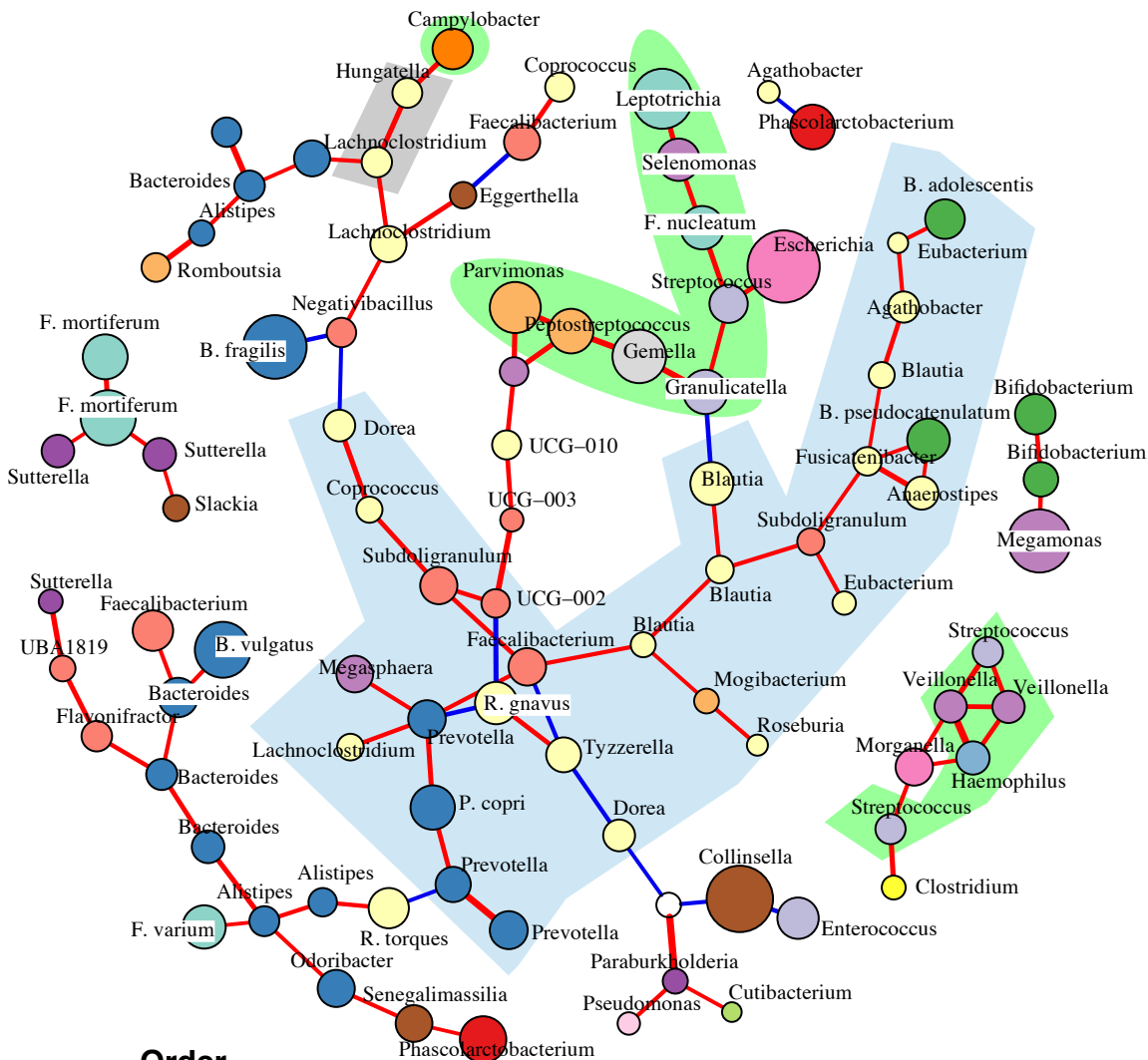
825    *Bioinformatics* 2014, **30**:2114–2120.

826    79.    Wick RR, Judd LM, Gorrie CL, Holt KE: **Unicycler: Resolving bacterial genome assemblies**

827    **from short and long sequencing reads**. *PLoS Comput Biol* 2017, **13**:1–22.

828    80.    Parks DH, Imelfort M, Skennerton CT, Hugenholtz P, Tyson GW: **CheckM: Assessing the**

829           **quality of microbial genomes recovered from isolates, single cells, and metagenomes**. *Genome*

830           *Res* 2015, **25**:1043–1055.

831    81.    Todd SM, Settlage RE, Lahmers KK, Slade DJ: ***Fusobacterium* Genomics Using MinION and**

832           **Illumina Sequencing Enables Genome Completion and Correction**. *mSphere* 2018, **3**:1–9.

833    82.    Jain C, Rodriguez-R LM, Phillippy AM, Konstantinidis KT, Aluru S: **High throughput ANI**

834           **analysis of 90K prokaryotic genomes reveals clear species boundaries**. *Nat Commun* 2018,

835           **9**:1–8.

836    83.    Ding W, Baumdicker F, Neher RA: **panX: pan-genome analysis and exploration**. *Nucleic*

837           *Acids Res* 2018, **46**:1–12.

838    84.    Stamatakis A: **RAxML version 8: a tool for phylogenetic analysis and post-analysis of large**

839           **phylogenies.** *Bioinformatics* 2014, **30**:1312–3.

840    85.    Gouy M, Guindon S, Gascuel O: **SeaView version 4: A multiplatform graphical user interface**

841           **for sequence alignment and phylogenetic tree building.** *Mol Biol Evol* 2010, **27**:221–224.

842    86.    Carver T, Berriman M, Tivey A, Patel C, Böhme U, Barrell BG, Parkhill J, Rajandream MA:

843           **Artemis and ACT: Viewing, annotating and comparing sequences stored in a relational**

844           **database**. *Bioinformatics* 2008, **24**:2672–2676.

845    87.    Hunt M, Mather AE, Sánchez-Busó L, Page AJ, Parkhill J, Keane JA, Harris SR: **ARIBA: Rapid**

846           **antimicrobial resistance genotyping directly from sequencing reads**. *Microb Genomics* 2017,

847           **3**:1–11.

848    88.    Yu G, Smith DK, Zhu H, Guan Y, Lam TTY: **ggtree: An r package for visualization and**

849           **annotation of phylogenetic trees with their covariates and other associated data**. *Methods*

850           *Ecol Evol* 2016, **8**:28–36.

851    89.    Key FM, Khadka VD, Romo-González C, Blake KJ, Deng L, Lynn TC, Lee JC, Chiu IM, García-

852           Romero MT, Lieberman TD: **On-person adaptive evolution of Staphylococcus aureus during**

853           **atopic dermatitis increases disease severity**. *bioRxiv* 2021,

854    90.    Assefa S, Keane TM, Otto TD, Newbold C, Berriman M: **ABACAS: Algorithm-based**

855           **automatic contiguation of assembled sequences**. *Bioinformatics* 2009, **25**:1968–1969.

856    91.    Li H: **Aligning sequence reads , clone sequences and assembly contigs with BWA-MEM**.

857           *bioarxiv* 2013,

858    92.    Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R:

859           **The Sequence Alignment/Map format and SAMtools.** *Bioinformatics* 2009, **25**:2078–9.

860    93.    McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K,

861           Altshuler D, Gabriel S, Daly M, et al.: **The Genome Analysis Toolkit: A MapReduce**

862           **framework for analyzing next-generation DNA sequencing data**. *Genome Res* 2010, **20**:1297–

863           303.

864    94.    Garrison E, Marth G: **Haplotype-based variant detection from short-read sequencing**. *arXiv*

865           *Prepr* 2012,

866    95.    Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO, Whitwham A, Keane T,

867           McCarthy SA, Davies RM, et al.: **Twelve years of SAMtools and BCFtools**. *Gigascience* 2021,
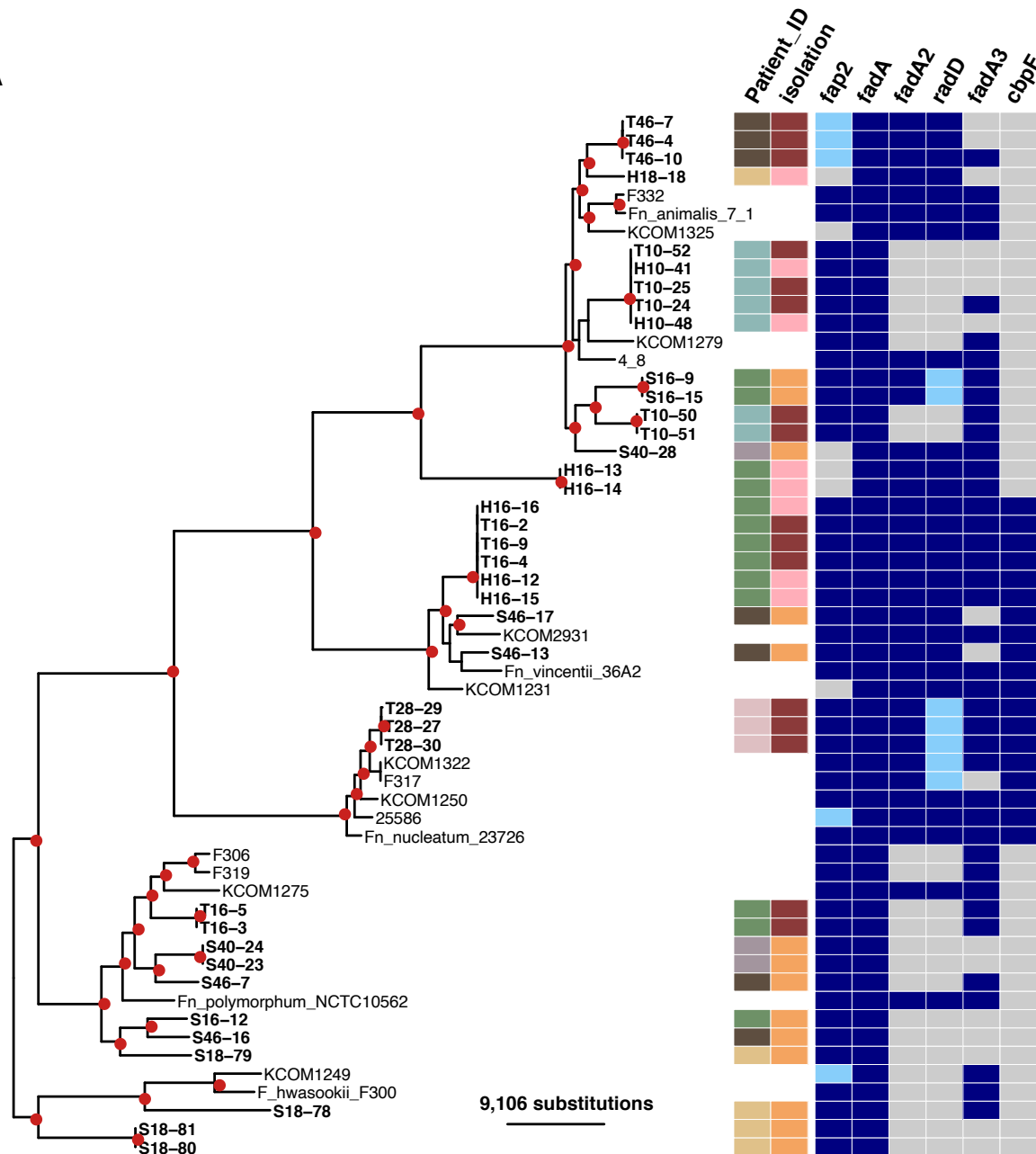
868           **10**:1–4.

869

**Order**

- Acidaminococcales
- Bacteroidales
- Bifidobacteriales
- Burkholderiales
- Campylobacterales
- Clostridiales
- Coriobacteriales
- Enterobacterales
- Fusobacteriales
- Lachnospirales
- Lactobacillales
- Oscillospirales
- Pasteurellales
- Staphylococcales
- Propionibacteriales
- Pseudomonadales
- Peptostreptococcales–Tissierellales
- Veillonellales–Selenomonadales