1    **An international report on bacterial communities in esophageal squamous cell**

2    **carcinoma**

3

4    Jason Nomburg[1,2,3], Susan Bullman[4], Dariush Nasrollahzadeh[5,6], Eric A. Collisson[7,8],

5    Behnoush Abedi-Ardekani[6], Larry O. Akoko[9], Joshua R. Atkins[6], Geoffrey C. Buckle[7,8],

6    Satish Gopal[10], Nan Hu[11], Bongani Kaimila[12], Masoud Khoshnia[5], Reza Malekzadeh[5],

7    Diana Menya[13], Blandina T. Mmbaga[14,15], Sarah Moody[16], Gift Mulima[17], Beatrice P.

8    Mushi[9], Julius Mwaiselage[18], Ally Mwanga[9], Yulia Newton[19], Dianna L. Ng[7,20], Amie

9    Radenbaugh[19], Deogratias S. Rwakatema[14,15], Msiba Selekwa[9], Joachim Schüz[21],

10    Philip R. Taylor[11], Charles Vaske[19], Alisa Goldstein[11], Michael R. Stratton[16], Valerie

11    McCormack[21], Paul Brennan[6], James A. DeCaprio[1,3,22], Matthew Meyerson[1,2,22,23*], Elia

12    J. Mmbaga[9, 24*], Katherine Van Loon[7,8*]

13

14

15    1 Department of Medical Oncology, Dana-Farber Cancer Institute, Boston, MA

16    2 Broad Institute of MIT and Harvard, Cambridge, MA

17    3 Harvard Program in Virology, Harvard Medical School, Boston, MA

18    4 Fred Hutchinson Cancer Research Center, Seattle, Washington, USA

19    5 Digestive Oncology Research Center, Digestive Disease Research Institute, Tehran

20    University of Medical Sciences, Shariati Hospital. Tehran Iran.

21    6 International Agency for Research on Cancer (IARC/WHO), Genomic Epidemiology

22    Branch, Lyon, France

23   7 University of California, San Francisco (UCSF) Helen Diller Family Comprehensive

24   Cancer Center, San Francisco, CA, USA

25   8 Division of Hematology/Oncology, Department of Medicine, UCSF, San Francisco,

26   California, USA

27   9 Muhimbili University of Health and Allied Sciences, Dar es Salaam, Tanzania

28   10 University of North Carolina, Chapel Hill, North Carolina, USA

29   11 Division of Cancer Epidemiology and Genetics, National Cancer Institute, Bethesda,

30   MD, USA

31   12 UNC Project - Lilongwe, Malawi

32   13 School of Public Health, Moi University, Eldoret, Kenya

33   14 Kilimanjaro Clinical Research Institute, Kilimanjaro Christian Medical Centre, Moshi,

34   Tanzania

35   15 Kilimanjaro Christian Medical University College, Moshi, Tanzania

36   16 Cancer, Ageing and Somatic Mutation, Wellcome Trust Sanger Institute, Wellcome

37   Trust Genome Campus, Hinxton, Cambridgeshire, UK

38   17 Kamuzu Central Hospital, Lilongwe, Malawi

39   18 Ocean Road Cancer Institute, Dar es Salaam, Tanzania

40   19 NantOmics/NantHealth, Inc., El Segundo, California, USA

41   20 Department of Pathology, UCSF, San Francisco, CA, USA

42   21 International Agency for Research on Cancer (IARC/WHO), Environment and Lifestyle

43   Epidemiology Branch, Lyon, France

44   22 Department of Medicine, Brigham and Women's Hospital, Harvard Medical School,

45   Boston, MA

46    23 Department of Genetics, Harvard Medical School, Boston, MA

47    24 Department of Community Medicine and Global Health, University of Oslo, Norway

48

49

50    **\*Correspondence to:**

51    Katherine Van Loon - Katherine.VanLoon@ucsf.edu

52    Elia J. Mmbaga - eliajelia@yahoo.co.uk

53    Matthew Meyerson - matthew_meyerson@dfci.harvard.edu

54

55

56    **ABSTRACT**

57    The incidence of esophageal squamous cell carcinoma (ESCC) is disproportionately

58    high in the eastern corridor of Africa and parts of Asia. Emerging research has identified

59    a potential association between poor oral health and ESCC. One proposed biological

60    pathway linking poor oral health and ESCC involves the alteration of the microbiome.

61    Thus, we performed an integrated analysis of four independent sequencing efforts of

62    ESCC tumors from patients from high- and low-incidence regions of the world. Using

63    whole genome sequencing (WGS) and RNA sequencing (RNAseq) of ESCC tumors

64    and WGS of synchronous collections of saliva specimens from 61 patients in Tanzania,

65    we identified a community of bacteria, including members of the genera *Fusobacterium,*

66    *Selenomonas, Prevotella, Streptococcus, Porphyromonas, Veillonella,* and

67    *Campylobacter*, present at high abundance in ESCC tumors. We then characterized the

68    microbiome of 238 ESCC tumor specimens collected in two additional independent

69    sequencing efforts consisting of patients from other high-ESCC incidence regions

70    (Tanzania, Malawi, Kenya, Iran, China). This analysis revealed a similar tumor

71    enrichment of the ESCC-associated bacterial community in these cancers. Because

72    these genera are traditionally considered members of the oral microbiota, we explored if

73    there is a relationship between the synchronous saliva and tumor microbiomes of ESCC

74    patients in Tanzania. Comparative analyses revealed that paired saliva and tumor

75    microbiomes are significantly similar with a specific enrichment of *Fusobacterium* and

76    *Prevotella* in the tumor microbiome. Together, these data indicate that cancer-

77    associated oral bacteria are associated with ESCC tumors at the time of diagnosis and

78    support a model in which oral bacteria are present in high abundance in both saliva and

4

79     tumors of ESCC patients. Longitudinal studies of the pre-diagnostic oral microbiome are

80     needed to investigate whether these cross-sectional similarities reflect temporal

81     associations.

82

83

84 **INTRODUCTION**

85 Esophageal cancer is the sixth most common cause of cancer-related death worldwide

86 (1).There are two histologic subtypes of esophageal cancer with distinct biological

87 characteristics, geographic distributions, and risk factors (2). Esophageal

88 adenocarcinoma is the most common histologic form of esophageal cancer in high-

89 income countries and is associated with factors including gastroesophageal reflux

90 disease, Barrett's esophagus, and obesity (3, 4). By contrast, esophageal squamous

91 cell carcinoma (ESCC) represents more than 90% of worldwide esophageal cancer

92 cases and is the dominant histology in low-resource settings. In particular, there are two

93 main regions where ESCC is endemic: (1) the Asian esophageal cancer belt, extending

94 from western/northern China to central and southeast Asia; and (2) the eastern corridor

95 of Africa, extending from Ethiopia to South Africa (5, 6).

96

97 Emerging research has identified a possible association between poor oral health and

98 ESCC. Studies from Asia, Europe, Latin America, Kenya, and Iran have reported

99 associations of ESCC with poor oral hygiene, chronic periodontal disease, dental decay,

100 and tooth loss (7-16). Recently, three parallel case-control studies in Kenya and

101 Tanzania, conducted as part of the African Esophageal Cancer Consortium (AfrECC)

102 and ESCCAPE (esccape.iarc.fr) collaborations, reported possible associations of poor

103 or infrequent oral hygiene with increased risk for ESCC in East Africa (17-20).

104

105 Alterations of the oral microbiome due to poor oral health is one proposed biological

106 pathway that could explain the link between oral health and ESCC. Many bacterial

6

107    genera associated with gastrointestinal cancers contain species that are traditionally

108    associated with healthy or diseased oral microbiomes. For example, *Helicobacter pylori*

109    was discovered to be associated with gastric cancers and mucosa-associated lymphoid

110    tissue (MALT) lymphomas, indirectly by promoting gastric inflammation and directly by

111    influencing cellular signaling (21). Similarly, bacteria of the genera *Fusobacterium*,

112    *Selenomonas*, and *Prevotella* are enriched in colorectal cancers (22-24) and can be

113    visualized invasively within tumor tissue (25). *Fusobacterium*, in particular, has been

114    reported to promote carcinogenesis through the selective expansion or inhibition of

115    certain classes of immune cells (26) and may drive cellular proliferation by stimulating

116    Wnt/β-catenin signaling (27, 28). Other bacterial genera such as *Porphyromonas,*

117    *Campylobacter,* and *Streptococcus* have emerging associations with various human

118    gastrointestinal cancers (29-35).

119

120    As part of ongoing investigation into the microbiome's association with ESCC, we

121    performed an integrated analysis of four independent sequencing efforts including

122    ESCC tumors from patients from both high- and low-incidence regions of the world. In

123    addition, we investigated the relationship between the microbiomes of matched ESCC

124    tumors and saliva specimens in a subset of ESCC cases.

125

126    **RESULTS**

127    **Study Population**

128    To evaluate the potential role of the host microbiota in ESCC, we investigated the

129    microbiome of 299 ESCC specimens from patients in five different countries with a high

130    incidence of ESCC. Specimens were collected through four independent sequencing

131    efforts (**Figure 1A**). Specimens consisted of whole genome sequencing (WGS) and

132    RNA sequencing (RNAseq) data from the tumor and saliva of 61 patients from Tanzania

133    (the "MUHAS Tanzania" cohort) (36), RNAseq data from the tumors of 30 ESCC

134    patients in Malawi (the "UNC Project – Malawi" cohort) (37), and WGS from 208

135    additional samples of tumors from patients in high ESCC incidence regions, including

136    specimens from ESCC patients in Tanzania (n=18) and Kenya (n=64) that were

137    collected in the ESCCAPE studies (esccape.iarc.fr) and specimens from ESCC patients

138    in East Golestan, Iran (n=55) and Shanxi, China (n=71) that were sequenced as part of

139    the Cancer Research UK Mutographs project ("Mutographs" cohorts) (38). In addition,

140    we analyzed WGS data of ESCC from The Cancer Genome Atlas (39), which includes a

141    small number of tumors from patients in low-incidence geographic regions including the

142    United States (n=3), Ukraine (n=3), Vietnam (n=22), and Russia (n=8) (the "TCGA"

143    cohort). Patient characteristics are shown in Table 1.

144

145    **Bacterial populations are abundant and diverse in ESCC tumors**

146    We used the metagenomic analysis tool GATK-PathSeq (40) to process the RNAseq

147    and WGS data. GATK-PathSeq uses a sequential mapping strategy to assign reads to

148    human and microbial reference genomes, resulting in detailed information on

149    sequencing reads of human and microbial origin (**Figure S1A**). We likewise used

150    GATK-PathSeq to process WGS data sets from 50 colon adenocarcinoma (COAD)

151    specimens available from TCGA (41) for comparison, as there is strong evidence of

152    microbial associations with COAD (22-25).

153

154    The bacterial burden of ESCC tumors ranged from 10 to 1000 bacterial reads per

155    million human reads, similar to numbers observed in TCGA COAD (**Figure 1B**).

156    Furthermore, the Shannon diversity of bacterial populations at the genus level ranged

157    from 2 to 3 (**Figure 1C**). By comparison, ESCC-associated bacterial communities are as

158    diverse or more diverse than TCGA COAD. At the phylum level, ESCC bacterial

159    populations generally consist of *Firmicutes*, *Bacteroidetes*, *Proteobacteria*,

160    *Actinobacteria*, and *Fusobacteria* (**Figure 1D, Figure S1B**). Of note, the higher than

161    expected abundance of the phylum *Actinobacteria* specifically in the TCGA ESCC

162    samples is attributable, in particular, to a very high abundance of the genus

163    *Tetrasphaera* (**Figure S1C**). This is evidenced by a depressed Shannon diversity of

164    *Actinobacteria* genera in these samples (**Figure S1D**) and may indicate contamination

165    of the TCGA ESCC samples. *Actinobacteria* have been reported as a source of

166    contaminating reads in TCGA gastrointestinal cancer samples (42).

167

168    **Bacterial genera associated with carcinogenesis are observed at high relative**

169    **abundance in ESCC tumors from Tanzania**

170    To determine if bacteria with known associations with cancer are present in ESCC, we

171    first analyzed the sequencing series of the 61 ESCC cases from the MUHAS Tanzania

172    cohort with both WGS and RNAseq data. The paired WGS and RNAseq data from

173    these tumors allowed investigation of bacterial communities at the DNA and RNA levels.

174    Both WGS and RNAseq data revealed high relative abundance of bacterial genera

175    previously associated with carcinogenesis in these ESCC tumors (**Figure 2A, 2B**). The

176    high relative abundance of the *Fusobacterium* genus was particularly notable. Other

177    bacterial genera of interest include *Streptococcus*, *Porphyromonas*, *Campylobacter*,

178    *Prevotella*, *Veillonella*, and *Selenomonas*, many of which have been associated with

179    gastrointestinal malignancies alongside or independently of *Fusobacterium* (25, 29, 32,

180    34, 43). The mean Jaccard similarity index between tumor RNAseq and WGS data from

181    the same tumor is 0.54, greater than the average Jaccard similarity index of random

182    RNAseq-WGS pairs (0.36), indicating that bacterial populations inferred from WGS and

183    RNAseq data are generally consistent (**Figure 2C**).

184

185    Next, we attempted to determine if similar bacterial genera were also present in ESCC

186    from patients in high-incidence countries beyond Tanzania. Investigation of RNA

187    sequencing data from patients in Malawi, WGS data from patients in Kenya, China, and

188    Iran, as well as from the independent ESCCAPE Tanzania patient group revealed

189    pervasive evidence of similar bacterial genera in the tumors of these patients (**Figure

190    2D, Figure S2A**). To investigate if similar microorganisms were found in ESCC tumors

191    from patients in low-incidence regions, we investigated WGS data from ESCC tumors

192    originating from USA, Ukraine, Vietnam, and Russia that were available through TCGA.

193    While the number of samples available from low-incidence regions is low and relies on a

194    single sequencing effort, we found that the tumors of many of these patients contain

195    similar bacterial genera (**Figure 2D, Figure S2A**). Colon cancers from the TCGA COAD

196    cohort revealed evidence of *Fusobacterium,* as expected; however, these COAD

197    samples were notable for much lower relative abundance of the other genera of interest,

198    when compared to ESCC tumors.

10

199

**Evaluation of association between saliva and tumor microbiomes in ESCC**

**patients from Tanzania**

We next investigated the similarity between the saliva and tumor microbiomes of ESCC

patients. Paired saliva samples were only available from patients in the MUHAS

Tanzania cohort (N=45); these paired saliva specimens were analyzed to evaluate

bacterial abundance as a proxy for the oral microbiome.

206

We first assessed the similarity between paired saliva and tumor microbiomes with the

Bray Curtis similarity index (44). To avoid potential confounding due to low bacterial

read counts in some tumor samples, we limited these analyses to the 21 tumor-saliva

pairs that contain appreciable microbial sequencing depth (at least 10,000 bacterial

reads each). We found that the saliva and tumor microbiomes from the same patient in

the Tanzanian samples are significantly more similar than random saliva-tumor pairs

(p=0.0003, Wilcoxon rank sum test) (**Figure 3A**). Next, we asked if there are bacterial

genera whose relative abundance in the saliva correlates with their relative abundance

in the tumor. For this analysis, we included only common-abundant bacterial genera

with at least 1% relative abundance in at least three tumor-oral pairs. The relative

abundance of four bacterial genera (*Fusobacterium*, *Veillonella*, *Streptococcus*, and

*Porphyromonas*) are strongly correlated between tumor and saliva microbiomes, while

other common-abundant bacterial genera were not (**Figure 3B**). To assess if any

bacterial genera are preferentially enriched in the tumor microbiome relative to the

saliva microbiome, we next calculated the difference in the relative abundance of the

11

222 common-abundant bacterial genera between saliva-tumor pairs. Several genera

223 including *Porphyromonas* and *Veillonella* were at higher relative abundance in the

224 saliva, while *Prevotella* and *Fusobacterium* were enriched in the tumor microbiome

225 (**Figure 3C**). Finally, the relative abundance of tumor-associated bacteria including

226 *Fusobacterium*, *Prevotella*, *Selenomonas*, *Veillonella*, *Streptococcus*, and

227 *Campylobacter* are strikingly similar between the microbiomes of tumor and oral pairs

228 (**Figure 3D**). Altogether, these data support the hypothesis that there is an association

229 between the oral and tumor microbiome of ESCC patients in Tanzania.

230

231 **DISCUSSION**

232 This report provides an analysis of bacterial communities present in ESCC tumors from

233 nine countries from different regions of the world, analyzed in four independent

234 sequencing efforts. We found traditionally oral, cancer-associated, bacterial genera in

235 tumors from patients in Tanzania, Malawi, Kenya, China, and Iran. These results

236 provide evidence that these bacterial genera may be associated with ESCC in high-

237 incidence regions. We also identified similar bacterial genera in ESCC tumors from low-

238 incidence regions, although this finding is based on a small sample size and only one

239 sequencing cohort. Finally, in a sub-analysis of tumor and saliva pairs available from

240 Tanzania, we demonstrated that the synchronous collected saliva and tumor

241 microbiomes of ESCC patients are strikingly similar at the time of diagnosis; in

242 particular, we identified a specific correlation between the saliva and tumor relative

243 abundance of the bacterial genera *Fusobacterium*, *Veillonella*, *Streptococcus*, and

12

244    *Porphyromonas*, with *Prevotella* and *Fusobacterium* significantly enriched in the tumor

245    microbiome.

246

247    Many of the bacterial genera identified in this study have been previously implicated in

248    the carcinogenesis of gastrointestinal cancers. For example, studies have found that

249    oral microbiota including *Fusobacterium*, *Prevotella*, *Selenomonas*, *Veillonella*,

250    *Streptococcus*, and *Campylobacter* can be used to distinguish individuals with colorectal

251    cancer from healthy controls (45), and that *Fusobacterium nucleatum* strains that

252    colonize the oral cavity and tumors of patients with colorectal cancer are identical in

253    some patients (46), raising the possibility that the oral cavity is a source of extra-oral

254    cancer microbiota. Our group has previously shown that *Fusobacterium*, *Selenomonas*,

255    and *Prevotella* can be visualized invasively within colorectal tumors and liver

256    metastases (25). Fusobactium nucleatum has been previously identified in esophageal

257    cancers and is associated with shorter survival (47). Members of the genus

258    *Porphyromonas* have been previously observed invasively within ESCC tumors (29)

259    and have been reported to promote oral squamous cell carcinoma through a variety of

260    mechanisms (30, 31). *Campylobacter jejuni* has been reported to promote

261    tumorigenesis in mice (32), and *Streptococcus* species have been identified in human

262    esophageal cancers (33). In addition, the striking association of *Streptococcus bovis*

263    with colorectal cancer has led to the recommendation that colonoscopy be performed

264    upon detection of *Streptococcus bovis* bacteremia or endocarditis (34, 35). Oral

265    commensal bacteria such as *Veillonella* species have been previously implicated in

266    pathogenesis of lung cancer (43). A prospective cohort of American patients (48) and a

13

267    study of Japanese patients (49) likewise found that oral microbiome composition reflects

268    risk of esophageal cancers

269

270    We found that bacterial genera including *Fusobacterium*, *Prevotella*, *Selenomonas*,

271    *Veillonella*, *Streptococcus*, and *Campylobacter* are pervasive in the microbiome of

272    ESCC tumors from patients in high-incidence regions. Moreover, the bacterial

273    composition of ESCC tumors is remarkably similar across countries in those high-

274    incidence regions, raising the possibility that these bacterial genera may be involved in

275    ESCC carcinogenesis or that they may colonize tumors as a result of the common

276    clinical presentation of patients with severe dysphagia. Notably, there are several

277    alternative hypotheses that warrant mention. For example, it is possible that the ESCC-

278    associated bacterial genera simply represent common members of the esophageal

279    microbiome (50) and that the microbial populations we observed in these cancers are

280    not significantly different from those found in normal esophagus tissue. A limitation of

281    our study is a lack of normal esophageal tissue from ESCC cases or healthy controls in

282    these settings, which would allow us to address this possibility. Another possible

283    explanation is that ESCC tumors provide a favorable niche in which these bacteria are

284    sequestered and allowed to colonize due to the propensity of this disease to cause

285    malignant obstruction. Thus, it is plausible that ESCC-associated bacteria are not

286    necessarily promoting ESCC carcinogenesis but rather represent passengers resulting

287    from the sequestration of oral secretions proximal to an obstructing tumor. While the

288    previous association of these bacterial genera with other cancers is consistent with the

289    hypothesis that they influence carcinogenesis of ESCC, future studies are necessary to

14

290    identify which, if any, direct influences these bacterial genera have upon ESCC

291    carcinogenesis. Nevertheless, even if these bacterial genera do not have a role in

292    increasing ESCC risk, but arise at the time of disease onset, they may have an

293    important role to play as part of a non-invasive early-detection biomarker. Finally, a

294    concern of all microbiome analyses is that observed bacteria can be a consequence of

295    contamination at some step between tumor harvest and sequencing. While some TCGA

296    samples may be contaminated by *Actinobacteria* as previously noted, the presence of

297    *Fusobacterium*, *Prevotella*, *Selenomonas*, *Veillonella*, *Streptococcus*, and

298    *Campylobacter* in four independently collected cohorts indicates that these finding are

299    unlikely due to contamination.

300

301    While this study focused on the presence of bacteria with ESCC in high-incidence

302    regions, we found evidence of similar cancer-associated bacteria in tumors in patients

303    from low-incidence regions (USA, Ukraine, Vietnam, and Russia). A limitation of this

304    assessment is the small sample size (n=36) and reliance on a single TCGA cohort that

305    likely contains contaminants (42). Regardless, this finding does not exclude the

306    possibility that the microbiome could be a factor driving patterns of ESCC incidence. For

307    example, it is possible that the prevalence of ESCC-associated bacteria in people could

308    vary across regions, which in turn could drive these differing rates of ESCC incidence.

309    This is an important topic for future study.

310

311    We found that the structure of synchronous paired tumor and oral microbiomes were

312    strikingly similar. It is possible that this similarity is driven by transient contact of saliva

15

313    and its associated microbiome with the tumor (e.g., during swallowing or tumor

314    extraction). However, we found that only four of sixteen common-abundant bacterial

315    genera correlate in abundance between the tumor and oral microbiomes, suggesting

316    tumor-oral microbiome similarity is not driven exclusively by "in-trans" interactions

317    between the saliva and tumor. We also found that genera including *Prevotella* and

318    *Fusobacterium* are often specifically enriched in the tumor microbiome, supporting a

319    model where specific oral bacterial preferentially colonize the tumor. A caveat of this

320    study is that we infer oral bacterial populations from the saliva, despite diverse

321    communities of bacteria throughout the oral cavity (51). However, we do observe

322    *Fusobacterium* in the saliva despite its general association with periodontal plaques

323    (52), suggesting saliva is capable of detecting periodontal pathogens. Additionally,

324    because the samples studied here are from patients with late-stage disease, it is

325    possible that tumor-induced changes to upper-gastrointestinal physiology and

326    dysphagia symptom-induced major dietary changes could themselves alter the oral

327    microbiomes of these patients. The previous findings from the ESCCAPE studies in

328    Kenya and Tanzania (17, 19) which found strong associations with dental staining (ORs

329    > 10) and for which photographic validation studies suggest that most dental staining

330    was not fluorosis, also point to a recent build-up of chromogenic bacteria. Studies of the

331    oral microbiome of patients at earlier stages of ESCC and in prospective studies are

332    necessary to address this possibility. We restricted our analysis to 21 tumor-oral pairs

333    that have a sufficient number of bacterial reads (at least 10,000). It is likely that

334    excluded samples are not molecularly distinct from included samples but that the

335    relatively low bacterial read counts in some tumors is simply reflective of low

336    sequencing depth.

337

338    Our observation of similar tumor and saliva microbiomes in ESCC patients is especially

339    notable considering emerging evidence linking periodontal disease and poor oral health

340    with increased risk of various cancers (17, 53, 54). This raises several important open

341    questions. It will be essential to determine if there is a difference in the oral prevalence

342    of these identified cancer-associated bacteria between ESCC patients and non-patients

343    earlier in the natural history of the disease, for example through comparisons of patients

344    with esophageal squamous dysplasia and healthy controls. Because the prevalence of

345    these bacteria may be associated with factors such as oral health, hygiene, and diet,

346    studies of the impact of these factors on the oral microbiome in the general population

347    would inform whether the oral microbiome is on a pathway linking oral hygiene to ESCC

348    risk and may have a role in prevention.

349

350    In conclusion, we show that cancer-associated, traditionally-oral bacteria including the

351    genera *Fusobacterium, Selenomonas, Prevotella, Streptococcus, Porphyromonas,*

352    *Veillonella,* and *Campylobacter* are highly abundant within ESCC tumors from patients

353    in high-ESCC incidence regions. We also show that there is a correlation between the

354    genus composition of the saliva microbiome and the ESCC tumor microbiome of some

355    ESCC patients. These findings will be foundational for future studies to understand if

356    and how bacteria influence ESCC pathogenesis and to understand the role of the oral

17

357 microbiome in this process. Finally, this study highlights the benefit of collaborative

358 investigation to evaluate the international heterogeneity of this disease.

359

360

361 **MATERIALS AND METHODS**

362 **Sample acquisition and sequencing**

363 The sample acquisition and sequencing methods for the studies from the MUHAS

364 Tanzania cohort (n=61) (36) and UNC Project - Malawi cohort (n=30) (11) have been

365 previously described. Samples sequenced in the Mutographs study (n=210) (38)

366 originated from patients in Golestan, Iran (n=55), ESCCAPE case-control studies in

367 Tanzania (n=18) (19) and Kenya (n=64) (17), and patients in Shanxi, China (n=71).

368 TCGA ESCC (n=36) and COAD samples (n=51) have been previously described (39,

369 41). The TCGA ESCC cohort includes tumors from patients in United States (n=3),

370 Ukraine (n=3), Vietnam (n=22), and Russia (n=8), regions which have lower incidence

371 of ESCC.

372

373 **Metagenomic analysis**

374 GATK-PathSeq (40) was used to conduct computational subtraction of human-mapping

375 reads from input RNAseq and WGS datasets. GATK-PathSeq works by first mapping

376 reads to a host reference database consisting of the human genome grch38 and

377 various supplemental human reference sequences. Next, non-human reads are

378 mapped against a comprehensive microbial database, and microbe read assignments

379 are reported for further study. From the MUHAS Tanzania cohort, a total of 61 tumor

18

380   WGS samples, 45 saliva WGS samples, and 59 RNAseq samples were processed

381   through GATK-PathSeq.

382

383   Bacterial abundance analyses and plotting were conducted in R (v3.5.1). To calculate

384   relative abundance at a phylogenetic level (e.g., phylum or genus), GATK-PathSeq

385   results were filtered for taxa at the level, and relative abundance was calculated for

386   each taxon as follows: (# of taxon reads)/(total # reads at the selected phylogenetic

387   level). The rows of all bacterial abundance heatmaps are arranged according to the

388   mean abundance across all samples. The sample order of relative abundance stacked

389   barplots were determined based on *Fusobacterium* genus relative abundance except

390   where noted. In **Figure 2D**, if any cohort contained more than 50 samples, 50 samples

391   were randomly selected for plotting. The distribution of relative abundances of genera of

392   interest in all samples can be found in **Figure S2**, where width of each violin represents

393   the relative distribution of observed bacterial relative abundance for all patients in each

394   patient cohort.

395

396   Jaccard distance between RNAseq and WGS data from each ESCC tumor was

397   calculated in R based on bacterial genera with at least 1% relative abundance. The

398   qualitative Jaccard index was used in this case because the comparison was between

399   DNA and RNA analytes which would not be expected to be quantitatively identical.

400

401   **Tumor-saliva similarity**

19

402    Only tumor-saliva pairs from the MUHAS Tanzania cohort with at least 10,000 reads

403    mapped to the bacterial superkingdom were available for analysis. This resulted in a

404    total of 21 tumor-oral pairs. Bray-Curtis dissimilarity metrics between tumor-oral pairs

405    were calculated using the R package vegan (55). **Figure 3A** presents the Bray-Curtis

406    *similarity* (1 – Bray-Curtis dissimilarity), for each tumor-oral pair.

407

408    To determine the correlation between the relative abundance of specific genera

409    between tumor and saliva microbiomes, common-abundant genera that are at least 1%

410    abundance in at least 3 tumor-oral pairs were identified. This resulted in the

411    identification of 16 common-abundant genera. Correlations represent a two-sided

412    Pearson correlation coefficient. To determine tumor-oral enrichment of common-

413    abundant genera, the difference in relative abundance of each genus between each

414    tumor-oral pair was plotted (**Figure 3C**). For the relative abundance bar plots of tumor-

415    saliva pairs (**Figure 3D**), bacterial genera that had been highlighted in previous figures

416    are labeled.

417

418    **Code and processed data availability**

419    All GATK-PathSeq output files and reproducible analysis and plotting R Notebooks are

420    available.

421    Zenodo: https://doi.org/10.5281/zenodo.4750577

422    GitHub: https://github.com/jnoms/ESCC_microbiome

423

424 Furthermore, all analysis and figures can be automatically reproduced through a series

425 of Google Colab documents.

426 Figure 1 and Supplementary Figure 1:

427 https://github.com/jnoms/ESCC_microbiome/blob/main/collab/Figure1.ipynb

428 Figure 2 and Supplementary Figure 2:

429 https://github.com/jnoms/ESCC_microbiome/blob/main/collab/Figure2.ipynb

430 Figure 3: https://github.com/jnoms/ESCC_microbiome/blob/main/collab/Figure3.ipynb

431

432 **ABBREVIATIONS**

433 AFRECC – African Esophageal Cancer Consortium

434 COAD – Colon adenocarcinoma

435 ESCA – Esophageal adenocarcinoma

436 ESCC – Esophageal squamous cell carcinoma

437 ESCCAPE – Esophageal Squamous Cell Carcinoma African Prevention Research

438 MUHAS – Muhimbili University of Health and Allied Sciences

439 RNAseq – RNA sequencing

440 TCGA – The Cancer Genome Atlas

441 WGS – Whole genome sequencing

442

443

444

445 **ACKNOWLEDGMENTS**

## CONFLICTS OF INTEREST

### FUNDING SOURCES

480

481

482   **FIGURE LEGENDS**

483

484   **Figure 1**. **Microbiome structure and composition of ESCC tumors**

485   A. Description of ESCC patients, and sample types, assessed in this study. TCGA –

486       The Cancer Genome Atlas; ESCCAPE – Esophageal Squamous Cell Carcinoma

487       African Prevention Research; Mutographs – Cancer Research UK Mutographs

488       Project.

489   B. Bacterial burden of ESCC tumors for each patient cohort. Units are bacterial

490       reads per million human reads as determined by GATK-PathSeq analysis. Each

23

491          dot represents one sample. Analyte type (RNA or DNA) and tumor type (ESCC

492          or COAD) are indicated by color.

493       C.  Shannon diversity of ESCC tumors for each patient cohort. Shannon diversity

494          was determined for each sample at the genus level based on genera that are at

495          least 1% relative abundance. Each dot represents one sample. Analyte type

496          (RNA or DNA) and tumor type (ESCC or COAD) are indicated by color.

497       D.  Heatmap describing the relative abundance of the five top phyla sorted by

498          average phylum relative abundance. Each column represents one sample. Rows

499          represent the indicated phyla. Units are relative abundance. Samples from each

500          cohort are WGS unless noted with "(RNA)", in which case they are RNAseq.

501

502    **Figure 2.  Identification of bacterial genera associated with carcinogenesis**

503       A.  Bacterial genera relative abundance of WGS data from the MUHAS Tanzania

504          cohort. Each column represents a single sample. Samples are ordered by

505          decreasing Fusobacterium relative abundance. Units are relative abundance of

506          bacterial genus-mapping reads. Color indicates the genus, and seven genera are

507          specified. Only patients with GATK-PathSeq analysis from both RNAseq and

508          WGS tumor data are plotted (n=59). Columns are ordered by decreasing relative

509          abundance of *Fusobacterium* genus reads.

510       B.  Bacterial genera relative abundance of RNAseq data from the MUHAS Tanzania

511          cohort. Each column represents a single sample. Here, column order is dictated

512          according to the patient order in Figure 2A. Units are relative abundance of

513          bacterial genus-mapping reads. Color indicates the genus, and seven genera are

514       specified. Only patients with GATK-PathSeq analysis from both RNAseq and

515       WGS tumor data are plotted (n=59). Samples are ordered in the same order as

516       Figure 2A, which is by *Fusobacterium* genus relative abundance in the WGS

517       data.

518    C.  Jaccard index between RNAseq and WGS data of tumors from the MUHAS

519       Tanzania cohort. For the "Paired by Sample" column, Jaccard indices were

520       calculated only between the WGS and RNAseq data from the same tumor (n=59

521       comparisons). For the "Random Pairs" column, Jaccard indices were calculated

522       between all possible WGS-RNAseq pairs independent of patient of origin to

523       represent the expected random distribution of Jaccard indices (n=3,481

524       comparisons). Jaccard index was calculated from relative abundance at the

525       genus level based on genera that are at least 1% relative abundance. The width

526       of the violin represents the relative proportion of comparisons with each Jaccard

527       index, and lines indicate $25^{th}$, $50^{th}$, and $75^{th}$ percentiles.

528    D.  Bacterial genera relative abundance of the remaining patient cohorts, including

529       RNAseq and WGS data as indicated. Each column represents a single sample.

530       Samples are ordered by decreasing Fusobacterium relative abundance within

531       each patient cohort. Units are relative abundance of bacterial genus-mapping

532       reads. Color indicates the genus, and seven genera are specified. Here, if there

533       were more than 50 samples in a patient cohort, 50 samples were randomly

534       selected for visualization. USA – United States, UA – Ukraine, RU – Russia. All

535       cohorts consist of WGS data, with the exception of the tumors from Malawi which

536       are RNAseq. (Number of samples plotted: UNC Project - Malawi 30; ESCCAPE

537    Tanzania 18; ESCCAPE Kenya 50; Shanxi, China 50; Golestan, Iran 50; TCGA

538    ESCC Vietnam 22; TCGA ESCC USA/UA/RU 14).

539

540    **Figure 3.  Association between synchronous saliva and tumor microbiomes in**

541    **Tanzanian ESCC patients**

542    **A.** Bray Curtis Similarity comparing tumor-saliva pairs from patients in the MUHAS

543    Tanzania cohort. Analysis was restricted to the 21 tumor-saliva pairs that

544    contained at least 10,000 bacterial reads. This analysis was conducted at the

545    genus level and using relative abundance. For the "Paired by Patient" column,

546    Bray Curtis Similarity was calculated only between the tumor and saliva WGS

547    data from the same patient. For the "Random Pairs" column, Bray Curtis

548    Similarity was calculated between all possible tumor-saliva pairs independent of

549    patient of origin to represent the expected random distribution of Bray Curtis

550    Similarity. (p=0.0003, Wilcoxon rank sum test).

551    **B.** Correlation between the relative abundance of common-abundant bacterial

552    genera in paired saliva and tumor WGS data. Analysis was restricted to the 21

553    tumor-saliva pairs that contained at least 10,000 bacterial reads. Common-

554    abundant bacterial genera are bacterial genera that are at least 1% abundance in

555    at least 3 tumor-saliva pairs – 16 bacterial genera made this cutoff. Correlation

556    represents a two-sided Pearson correlation. X-axis is the correlation coefficient,

557    and Y axis is the correlation P-Value plotted on a log scale.

558    **C.** Enrichment of genera in the oral or tumor microbiome. Each row details one of

559    the 16 common-abundant bacterial genera. Each row contains one data point per

26

560          patient, for a total of 21 data points. The value of each point represents the

561          difference in the relative abundance of the specified genus in the tumor and oral

562          microbiomes of one patient, with positive values indicating a genus is at higher

563          relative abundance in a patient's tumor. For example, if a genus is at a relative

564          abundance of 0.7 (70%) in the tumor and 0.3 (30%) in the saliva of a patient, the

565          plotted value for that genus and that patient is 0.4. Curves represent the

566          distribution of this relative abundance difference across the tumor-oral pairs, with

567          dots indicating individual tumor-oral pairs. Vertical red lines indicate quartiles.

568   **D.** Relative abundance bar charts of tumor-saliva pairs. Analysis was restricted to

569          the 21 tumor-saliva pairs that contained at least 10,000 bacterial reads. Units are

570          relative abundance of bacterial genus-mapping reads. Color indicates the genus,

571          and seven genera are specified. (abbreviations: T – tumor, S – saliva).

572

573 **Figure S1. GATK-PathSeq statistics and extended phyla and genera information**

574   A. Boxplots indicating the number of GATK-PathSeq Human-mapped reads and

575          GATK-PathSeq microbe-mapped reads for each patient cohort. Samples from

576          each cohort are WGS unless noted with "(RNA)", in which case they are

577          RNAseq.

578   B. Heatmap describing the relative abundance of the 15 top phyla sorted by

579          average phylum relative abundance. Each column represents one sample. Rows

580          represent the indicated phyla. Units are relative abundance. Samples from each

581          cohort are WGS unless noted with "(RNA)", in which case they are RNAseq.

582     C. Heatmap describing the relative abundance of the 15 top genera sorted by

583       average genera relative abundance. Each column represents one sample. Rows

584       represent the indicated genera. Units are relative abundance. Samples from

585       each cohort are WGS unless noted with "(RNA)", in which case they are

586       RNAseq.

587     D. Boxplot representing the Shannon diversity of genera that fall within the phylum

588       *Actinobacteria* for each patient in each cohort. Samples from each cohort are

589       WGS unless noted with "(RNA)", in which case they are RNAseq.

590

591 **Figure S2. Distribution of *Fusobacterium, Selenomonas, Prevotella,***

592 ***Streptococcus, Porphyromonas, Veillonella,* and *Campylobacter* relative**

593 **abundance of genus reads for all samples in each study**

594     A. The distribution of the relative abundance of genus-mapping reads for seven

595       selected genera in all studies. The width of each violin represents the proportion

596       of samples which have the indicated relative abundance of each genus. In

597       contrast to **Figure 2D**, which only plots up to 50 samples per study, this plot

598       includes all patients. Samples from each study are WGS unless noted with

599       "(RNA)", in which case they are RNAseq.

600

601

602

28

## REFERENCES

1.      Sung H, Ferlay J, Siegel RL, Laversanne M, Soerjomataram I, Jemal A, et al. Global cancer statistics 2020: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries. CA: A Cancer Journal for Clinicians.n/a(n/a).

2.      CC MRA, Dawsey S. Oesophageal cancer: A tale of two malignancies. World Cancer Report: Cancer Research for Cancer Prevention Lyon, France: International Agency for Research on Cancer Available from: http://publications.iarc.fr/586. 2020.

3.      Coleman HG, Xie S-H, Lagergren J. The epidemiology of esophageal adenocarcinoma. Gastroenterology. 2018;154(2):390-405.

4.      Rustgi AK, El-Serag HB. Esophageal carcinoma. New England Journal of Medicine. 2014;371(26):2499-509.

5.      Arnold M, Soerjomataram I, Ferlay J, Forman D. Global incidence of oesophageal cancer by histological subtype in 2012. Gut. 2015;64(3):381-7.

6.      Cheng ML, Zhang L, Borok M, Chokunonga E, Dzamamala C, Korir A, et al. The incidence of oesophageal cancer in Eastern Africa: identification of a new geographic hot spot? Cancer epidemiology. 2015;39(2):143-9.

7.      Abnet CC, Kamangar F, Islami F, Nasrollahzadeh D, Brennan P, Aghcheli K, et al. Tooth loss and lack of regular oral hygiene are associated with higher risk of esophageal squamous cell carcinoma. Cancer Epidemiology and Prevention Biomarkers. 2008;17(11):3062-8.

8.      Abnet CC, Qiao Y-L, Mark SD, Dong Z-W, Taylor PR, Dawsey SM. Prospective study of tooth loss and incident esophageal and gastric cancers in China. Cancer Causes & Control. 2001;12(9):847-54.

9.      Dar N, Islami F, Bhat G, Shah I, Makhdoomi M, Iqbal B, et al. Poor oral hygiene and risk of esophageal squamous cell carcinoma in Kashmir. British journal of cancer. 2013;109(5):1367-72.

10.     Chen X, Yuan Z, Lu M, Zhang Y, Jin L, Ye W. Poor oral health is associated with an increased risk of esophageal squamous cell carcinoma-a population-based case-control study in China. International journal of cancer. 2017;140(3):626-35.

11.     Sato F, Oze I, Kawakita D, Yamamoto N, Ito H, Hosono S, et al. Inverse association between toothbrushing and upper aerodigestive tract cancer risk in a Japanese population. Head & neck. 2011;33(11):1628-37.

12.     Liang H, Yang Z, Wang JB, Yu P, Fan JH, Qiao YL, et al. Association between oral leukoplakia and risk of upper gastrointestinal cancer death: a follow-up study of the Linxian general population trial. Thoracic cancer. 2017;8(6):642-8.

13.     Guha N, Boffetta P, Wünsch Filho V, Eluf Neto J, Shangina O, Zaridze D, et al. Oral health and risk of squamous cell carcinoma of the head and neck and esophagus: results of two multicentric case-control studies. American journal of epidemiology. 2007;166(10):1159-73.

14.     Chen Q-L, Zeng X-T, Luo Z-X, Duan X-L, Qin J, Leng W-D. Tooth loss is associated with increased risk of esophageal cancer: evidence from a meta-analysis with dose-response analysis. Scientific reports. 2016;6(1):1-7.

15.     Sheikh M, Poustchi H, Pourshams A, Etemadi A, Islami F, Khoshnia M, et al. Individual and combined effects of environmental risk factors for esophageal cancer

647  based on results from the Golestan Cohort Study. Gastroenterology. 2019;156(5):1416-
648  27.
649  16.    Patel K, Wakhisi J, Mining S, Mwangi A, Patel R. Esophageal cancer, the
650  topmost cancer at MTRH in the Rift Valley, Kenya, and its potential risk factors.
651  International Scholarly Research Notices. 2013;2013.
652  17.    Menya D, Maina SK, Kibosia C, Kigen N, Oduor M, Some F, et al. Dental
653  fluorosis and oral health in the African Esophageal Cancer Corridor: Findings from the
654  Kenya ESCCAPE case–control study and a pan-African perspective. International
655  journal of cancer. 2019;145(1):99-109.
656  18.    Mmbaga EJ, Mushi BP, Deardorff K, Mgisha W, Akoko LO, Paciorek A, et al. A
657  Case–Control Study to Evaluate Environmental and Lifestyle Risk Factors for
658  Esophageal Cancer in Tanzania. Cancer Epidemiology and Prevention Biomarkers.
659  2020.
660  19.    Mmbaga BT, Mwasamwaja A, Mushi G, Mremi A, Nyakunga G, Kiwelu I, et al.
661  Missing and decayed teeth, oral hygiene and dental staining in relation to esophageal
662  cancer risk: ESCCAPE case-control study in Kilimanjaro, Tanzania. International journal
663  of cancer. 2020.
664  20.    Buckle GC, et al. Risk factors associated with early-onset esophageal cancer in
665  Tanzania. (Under Review).
666  21.    Ishaq S, Nunn L. Helicobacter pylori and gastric cancer: a state of the art review.
667  Gastroenterology and hepatology from bed to bench. 2015;8(Suppl1):S6.
668  22.    Kostic AD, Chun E, Robertson L, Glickman JN, Gallini CA, Michaud M, et al.
669  Fusobacterium nucleatum potentiates intestinal tumorigenesis and modulates the
670  tumor-immune microenvironment. Cell host & microbe. 2013;14(2):207-15.
671  23.    Kostic AD, Gevers D, Pedamallu CS, Michaud M, Duke F, Earl AM, et al.
672  Genomic analysis identifies association of Fusobacterium with colorectal carcinoma.
673  Genome research. 2012;22(2):292-8.
674  24.    Castellarin M, Warren RL, Freeman JD, Dreolini L, Krzywinski M, Strauss J, et al.
675  Fusobacterium nucleatum infection is prevalent in human colorectal carcinoma.
676  Genome research. 2012;22(2):299-306.
677  25.    Bullman S, Pedamallu CS, Sicinska E, Clancy TE, Zhang X, Cai D, et al.
678  Analysis of Fusobacterium persistence and antibiotic response in colorectal cancer.
679  Science. 2017;358(6369):1443-8.
680  26.    Gur C, Ibrahim Y, Isaacson B, Yamin R, Abed J, Gamliel M, et al. Binding of the
681  Fap2 protein of Fusobacterium nucleatum to human inhibitory receptor TIGIT protects
682  tumors from immune cell attack. Immunity. 2015;42(2):344-55.
683  27.    Rubinstein MR, Baik JE, Lagana SM, Han RP, Raab WJ, Sahoo D, et al.
684  Fusobacterium nucleatum promotes colorectal cancer by inducing Wnt/β-catenin
685  modulator Annexin A1. EMBO reports. 2019;20(4):e47638.
686  28.    Rubinstein MR, Wang X, Liu W, Hao Y, Cai G, Han YW. Fusobacterium
687  nucleatum promotes colorectal carcinogenesis by modulating E-cadherin/β-catenin
688  signaling via its FadA adhesin. Cell host & microbe. 2013;14(2):195-206.
689  29.    Gao S, Li S, Ma Z, Liang S, Shan T, Zhang M, et al. Presence of Porphyromonas
690  gingivalis in esophagus and its association with the clinicopathological characteristics
691  and survival in patients with esophageal cancer. Infectious agents and cancer.
692  2016;11(1):3.

693    30.    Whitmore SE, Lamont RJ. Oral bacteria and cancer. PLoS pathogens.
694    2014;10(3):e1003933.
695    31.    Inaba H, Sugita H, Kuboniwa M, Iwai S, Hamada M, Noda T, et al. P
696    orphyromonas gingivalis promotes invasion of oral squamous cell carcinoma through
697    induction of pro MMP 9 and its activation. Cellular microbiology. 2014;16(1):131-45.
698    32.    He Z, Gharaibeh RZ, Newsome RC, Pope JL, Dougherty MW, Tomkovich S, et
699    al. Campylobacter jejuni promotes colorectal tumorigenesis through the action of
700    cytolethal distending toxin. Gut. 2019;68(2):289-300.
701    33.    Narikiyo M, Tanabe C, Yamada Y, Igaki H, Tachimori Y, Kato H, et al. Frequent
702    and preferential infection of Treponema denticola, Streptococcus mitis, and
703    Streptococcus anginosus in esophageal cancers. Cancer science. 2004;95(7):569-74.
704    34.    Boleij A, Schaeps RM, Tjalsma H. Association between Streptococcus bovis and
705    colon cancer. Journal of clinical microbiology. 2009;47(2):516-.
706    35.    Ferrari A, Botrugno I, Bombelli E, Dominioni T, Cavazzi E, Dionigi P.
707    Colonoscopy is mandatory after Streptococcus bovis endocarditis: a lesson still not
708    learned. Case report. World journal of surgical oncology. 2008;6(1):49.
709    36.    Van Loon K, et al. A Genomic Analysis of Esophageal Squamous Cell
710    Carcinoma in Eastern Africa. (Under Review).
711    37.    Liu W, Snell JM, Jeck WR, Hoadley KA, Wilkerson MD, Parker JS, et al.
712    Subtyping sub-Saharan esophageal squamous cell carcinoma by comprehensive
713    molecular analysis. JCI insight. 2016;1(16).
714    38.    Moody S, Senkin S, Islam SMA, Wang J, Nasrollahzadeh D, Penha RCC, et al.
715    Mutational signatures in esophageal squamous cell carcinoma from eight countries of
716    varying incidence. medRxiv. 2021:2021.04.29.21255920.
717    39.    Network CGAR. Integrated genomic characterization of oesophageal carcinoma.
718    Nature. 2017;541(7636):169.
719    40.    Walker MA, Pedamallu CS, Ojesina AI, Bullman S, Sharpe T, Whelan CW, et al.
720    GATK PathSeq: a customizable computational tool for the discovery and identification of
721    microbial sequences in libraries from eukaryotic hosts. Bioinformatics.
722    2018;34(24):4287-9.
723    41.    Network CGA. Comprehensive molecular characterization of human colon and
724    rectal cancer. Nature. 2012;487(7407):330.
725    42.    Dohlman AB, Arguijo Mendoza D, Ding S, Gao M, Dressman H, Iliev ID, et al.
726    The cancer microbiome atlas: a pan-cancer comparative analysis to distinguish tissue-
727    resident microbiota from contaminants. Cell Host & Microbe. 2020.
728    43.    Tsay J-CJ, Wu BG, Sulaiman I, Gershner K, Schluger R, Li Y, et al. Lower airway
729    dysbiosis affects lung cancer progression. Cancer Discovery. 2020.
730    44.    Ricotta C, Podani J. On some properties of the Bray-Curtis dissimilarity and their
731    ecological meaning. Ecological Complexity. 2017;31:201-5.
732    45.    Flemer B, Warren RD, Barrett MP, Cisek K, Das A, Jeffery IB, et al. The oral
733    microbiota in colorectal cancer is distinctive and predictive. Gut. 2018;67(8):1454-63.
734    46.    Komiya Y, Shimomura Y, Higurashi T, Sugi Y, Arimoto J, Umezawa S, et al.
735    Patients with colorectal cancer have identical strains of Fusobacterium nucleatum in
736    their colorectal cancer and oral cavity. Gut. 2019;68(7):1335-7.

737    47.    Yamamura K, Baba Y, Nakagawa S, Mima K, Miyake K, Nakamura K, et al.
738    Human microbiome Fusobacterium nucleatum in esophageal cancer tissue is
739    associated with prognosis. Clinical Cancer Research. 2016;22(22):5574-81.
740    48.    Peters BA, Wu J, Pei Z, Yang L, Purdue MP, Freedman ND, et al. Oral
741    microbiome composition reflects prospective risk for esophageal cancers. Cancer
742    research. 2017;77(23):6777-87.
743    49.    Kawasaki M, Ikeda Y, Ikeda E, Takahashi M, Tanaka D, Nakajima Y, et al. Oral
744    infectious bacteria in dental plaque and saliva as risk factors in patients with esophageal
745    cancer. Cancer. 2021;127(4):512-9.
746    50.    Corning B, Copland AP, Frye JW. The esophageal microbiome in health and
747    disease. Current gastroenterology reports. 2018;20(8):1-7.
748    51.    Dewhirst FE, Chen T, Izard J, Paster BJ, Tanner AC, Yu W-H, et al. The human
749    oral microbiome. Journal of bacteriology. 2010;192(19):5002-17.
750    52.    Signat B, Roques C, Poulet P, Duffaut D. Role of Fusobacterium nucleatum in
751    periodontal health and disease. Curr Issues Mol Biol. 2011;13(2):25-36.
752    53.    Michaud DS, Lu J, Peacock-Villada AY, Barber JR, Joshu CE, Prizment AE, et
753    al. Periodontal disease assessed using clinical dental measurements and cancer risk in
754    the ARIC study. JNCI: Journal of the National Cancer Institute. 2018;110(8):843-54.
755    54.    Ahrens W, Pohlabeln H, Foraita R, Nelis M, Lagiou P, Lagiou A, et al. Oral
756    health, dental care and mouthwash associated with upper aerodigestive tract cancer
757    risk in Europe: the ARCAGE study. Oral oncology. 2014;50(6):616-25.
758    55.    Oksanen J, Kindt R, Legendre P, O'Hara B, Stevens MHH, Oksanen MJ, et al.
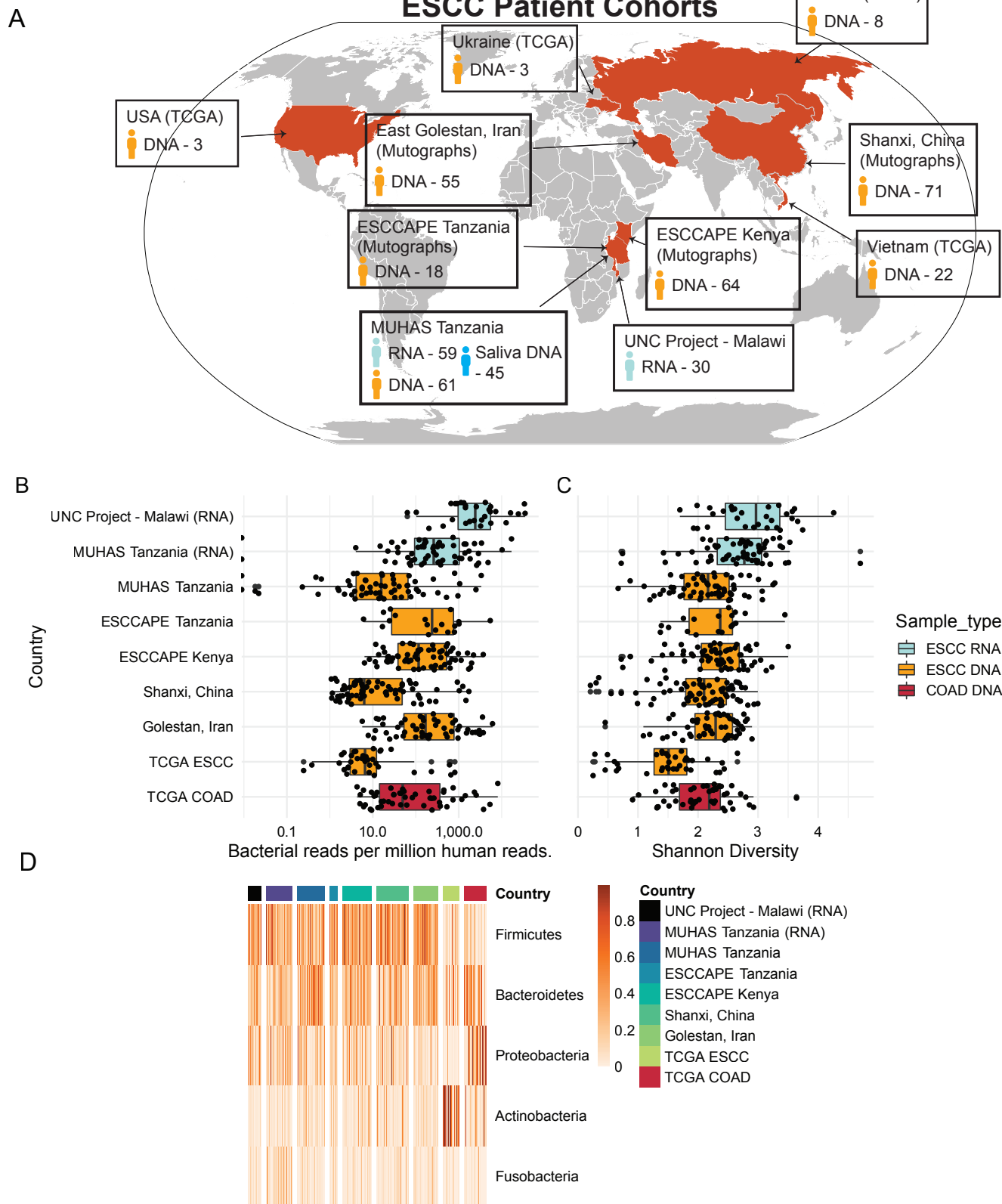759    The vegan package. Community ecology package. 2007;10:631-7.
760

Figure 1



Figure 1. Microbiome structure and composition of ESCC tumors.

A.      Description of ESCC patients, and sample types, assessed in this study. TCGA – The Cancer Genome Atlas; ESCCAPE – Esophageal Squamous Cell Carcinoma African Prevention Research; Mutographs – Cancer Research UK Mutographs Project.

B.      Bacterial burden of ESCC tumors for each patient cohort. Units are bacterial reads per million human reads as determined by GATK-PathSeq analysis. Each dot represents one sample. Analyte type (RNA or DNA) and tumor type (ESCC or COAD) are indicated by color.

C.      Shannon diversity of ESCC tumors for each patient cohort. Shannon diversity was determined for each sample at the genus level based on genera that are at least 1% relative abundance. Each dot represents one sample. Analyte type (RNA or DNA) and tumor type (ESCC or COAD) are indicated by color.

D.      Heatmap describing the relative abundance of the five top phyla sorted by average phylum relative abundance. Each column represents one sample. Rows represent the indicated phyla. Units are relative abundance. Samples from each cohort are WGS unless noted with "(RNA)", in which case they are RNAseq.
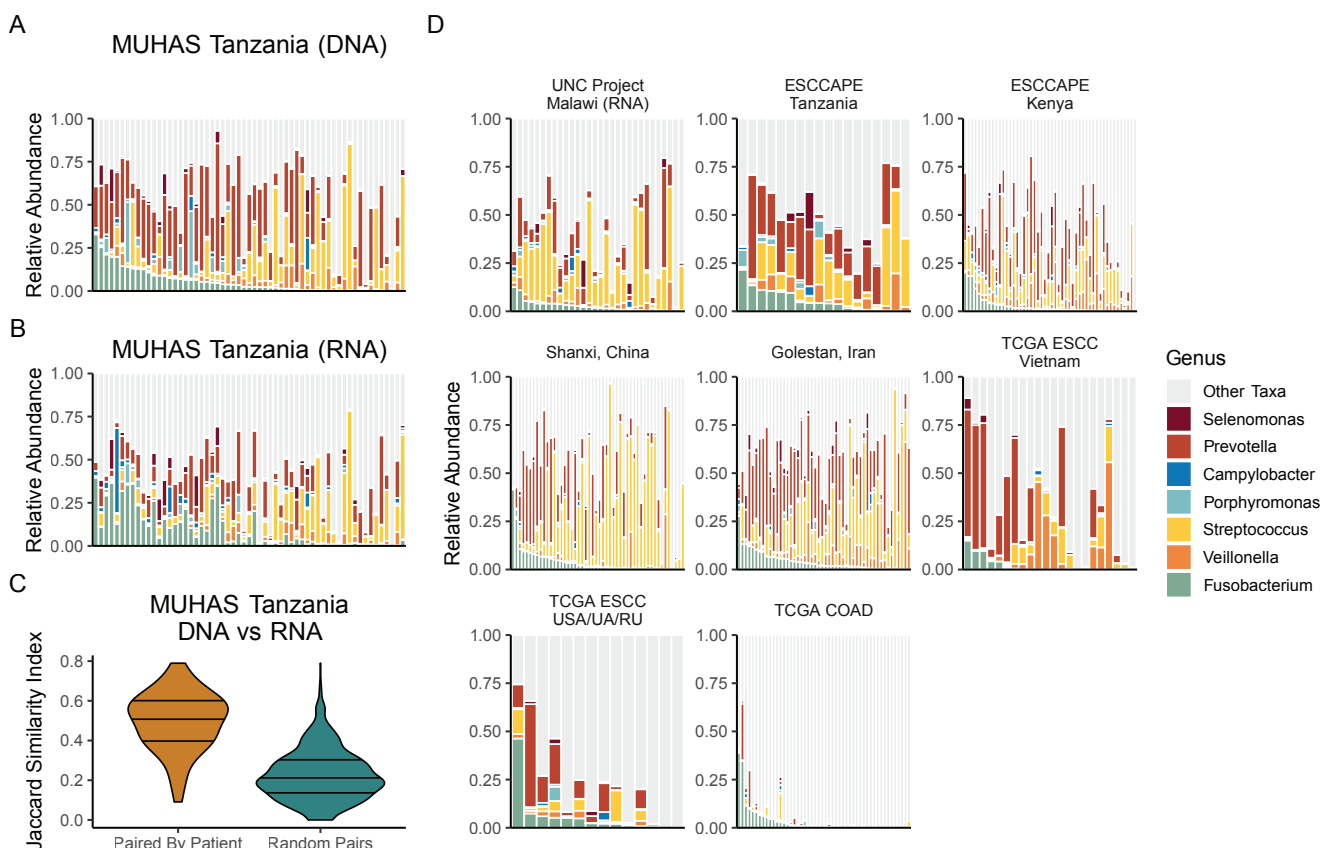
Figure 2



Figure 2. Identification of bacterial genera associated with carcinogenesis.

A.    Bacterial genera relative abundance of WGS data from the MUHAS Tanzania cohort.  Each column represents a single sample. Samples are ordered by decreasing Fusobacterium relative abundance. Units are relative abundance of bacterial genus-mapping reads. Color indicates the genus, and seven genera are specified. Only patients with GATK-PathSeq analysis from both RNAseq and WGS tumor data are plotted (n=59). Columns are ordered by decreasing relative abundance of Fusobacterium genus reads.

B.    Bacterial genera relative abundance of RNAseq data from the MUHAS Tanzania cohort. Each column represents a single sample. Here, column order is dictated according to the patient order in Figure 2A. Units are relative abundance of bacterial genus-mapping reads. Color indicates the genus, and seven genera are specified. Only patients with GATK-PathSeq analysis from both RNAseq and WGS tumor data are plotted (n=59). Samples are ordered in the same order as Figure 2A, which is by Fusobacterium genus relative abundance in the WGS data.

C.    Jaccard index between RNAseq and WGS data of tumors from the MUHAS Tanzania cohort. For the "Paired by Sample" column, Jaccard indices were calculated only between the WGS and RNAseq data from the same tumor (n=59 comparisons). For the "Random Pairs" column, Jaccard indices were calculated between all possible WGS-RNAseq pairs independent of patient of origin to represent the expected random distribution of Jaccard indices (n=3,481 comparisons). Jaccard index was calculated from relative abundance at the genus level based on genera that are at least 1% relative abundance. The width of the violin represents the relative proportion of comparisons with each Jaccard index, and lines indicate 25th, 50th, and 75th percentiles.

D.    Bacterial genera relative abundance of the remaining patient cohorts, including RNAseq and WGS data as indicated. Each column represents a single sample. Samples are ordered by decreasing Fusobacterium relative abundance within each patient cohort. Units are relative abundance of bacterial genus-mapping reads. Color indicates the genus, and seven genera are specified. Here, if there were more than 50 samples in a patient cohort, 50 samples were randomly selected for visualization. USA – United States, UA – Ukraine, RU – Russia. All cohorts consist of WGS data, with the exception of the tumors from Malawi which are RNAseq. (Number of samples plotted: UNC Project - Malawi 30; ESCCAPE Tanzania 18; ESCCAPE Kenya 50; Shanxi, China 50; Golestan, Iran 50; TCGA ESCC Vietnam 22; TCGA ESCC USA/UA/RU 14).
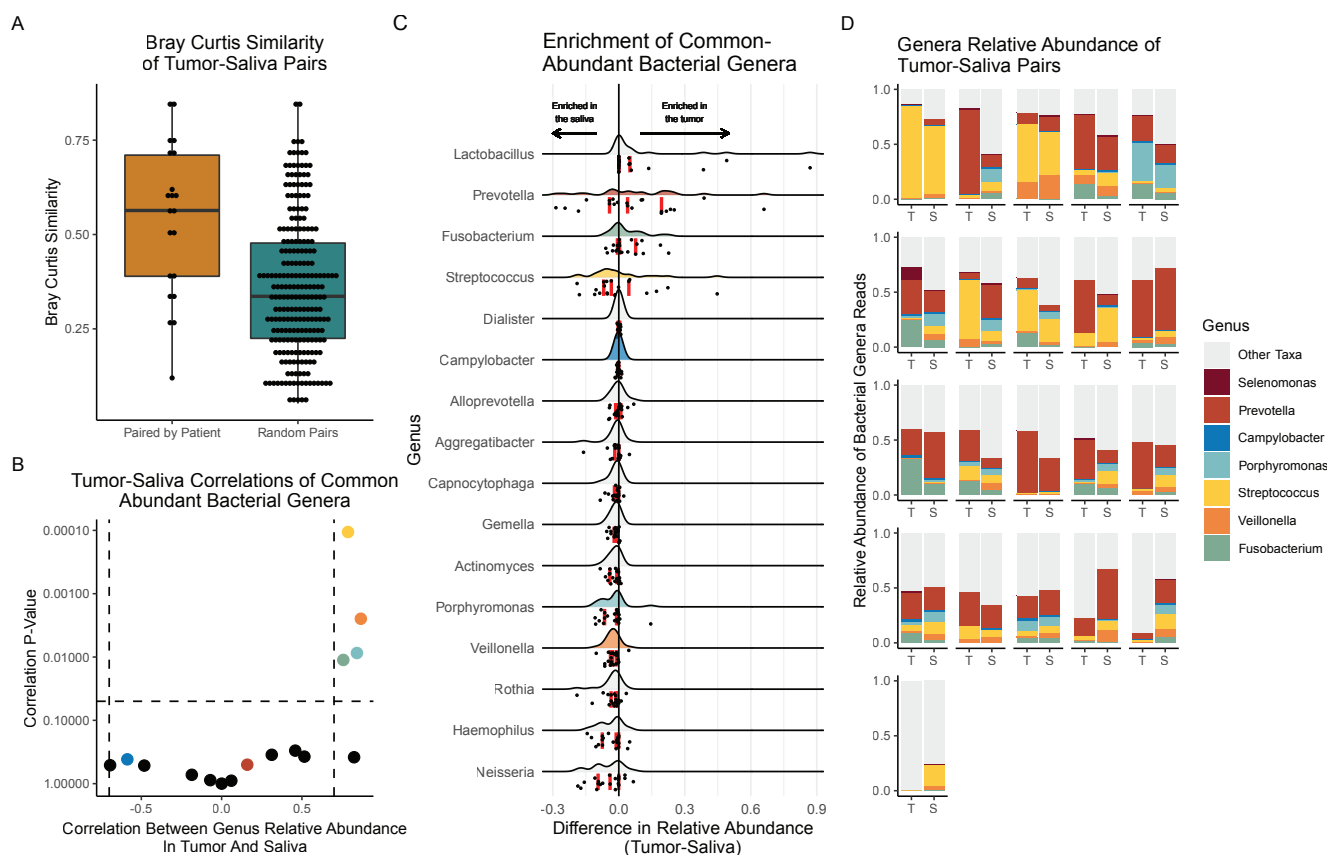
Figure 3



Figure 3. Association between synchronous saliva and tumor microbiomes in Tanzanian ESCC patients.

A. Bray Curtis Similarity comparing tumor-saliva pairs from patients in the MUHAS Tanzania cohort. Analysis was restricted to the 21 tumor-saliva pairs that contained at least 10,000 bacterial reads. This analysis was conducted at the genus level and using relative abundance. For the "Paired by Patient" column, Bray Curtis Similarity was calculated only between the tumor and saliva WGS data from the same patient. For the "Random Pairs" column, Bray Curtis Similarity was calculated between all possible tumor-saliva pairs independent of patient of origin to represent the expected random distribution of Bray Curtis Similarity. (p=0.0003, Wilcoxon rank sum test).

B. Correlation between the relative abundance of common-abundant bacterial genera in paired saliva and tumor WGS data. Analysis was restricted to the 21 tumor-saliva pairs that contained at least 10,000 bacterial reads. Common-abundant bacterial genera are bacterial genera that are at least 1% abundance in at least 3 tumor-saliva pairs – 16 bacterial genera made this cutoff. Correlation represents a two-sided Pearson correlation. X-axis is the correlation coefficient, and Y axis is the correlation P-Value plotted on a log scale.

C. Enrichment of genera in the oral or tumor microbiome. Each row details one of the 16 common-abundant bacterial genera. Each row contains one data point per patient, for a total of 21 data points. The value of each point represents the difference in the relative abundance of the specified genus in the tumor and oral microbiomes of one patient, with positive values indicating a genus is at higher relative abundance in a patient's tumor. For example, if a genus is at a relative abundance of 0.7 (70%) in the tumor and 0.3 (30%) in the saliva of a patient, the plotted value for that genus and that patient is 0.4. Curves represent the distribution of this relative abundance difference across the tumor-oral pairs, with dots indicating individual tumor-oral pairs. Vertical red lines indicate quartiles.

D. Relative abundance barcharts of tumor-saliva pairs. Analysis was restricted to the 21 tumor-saliva pairs that contained at least 10,000 bacterial eads. Units are relative abundance of bacterial genus-mapping reads. Color indicates the genus, and seven genera are specified. (Abbreviations: T – tumor, S – saliva.)
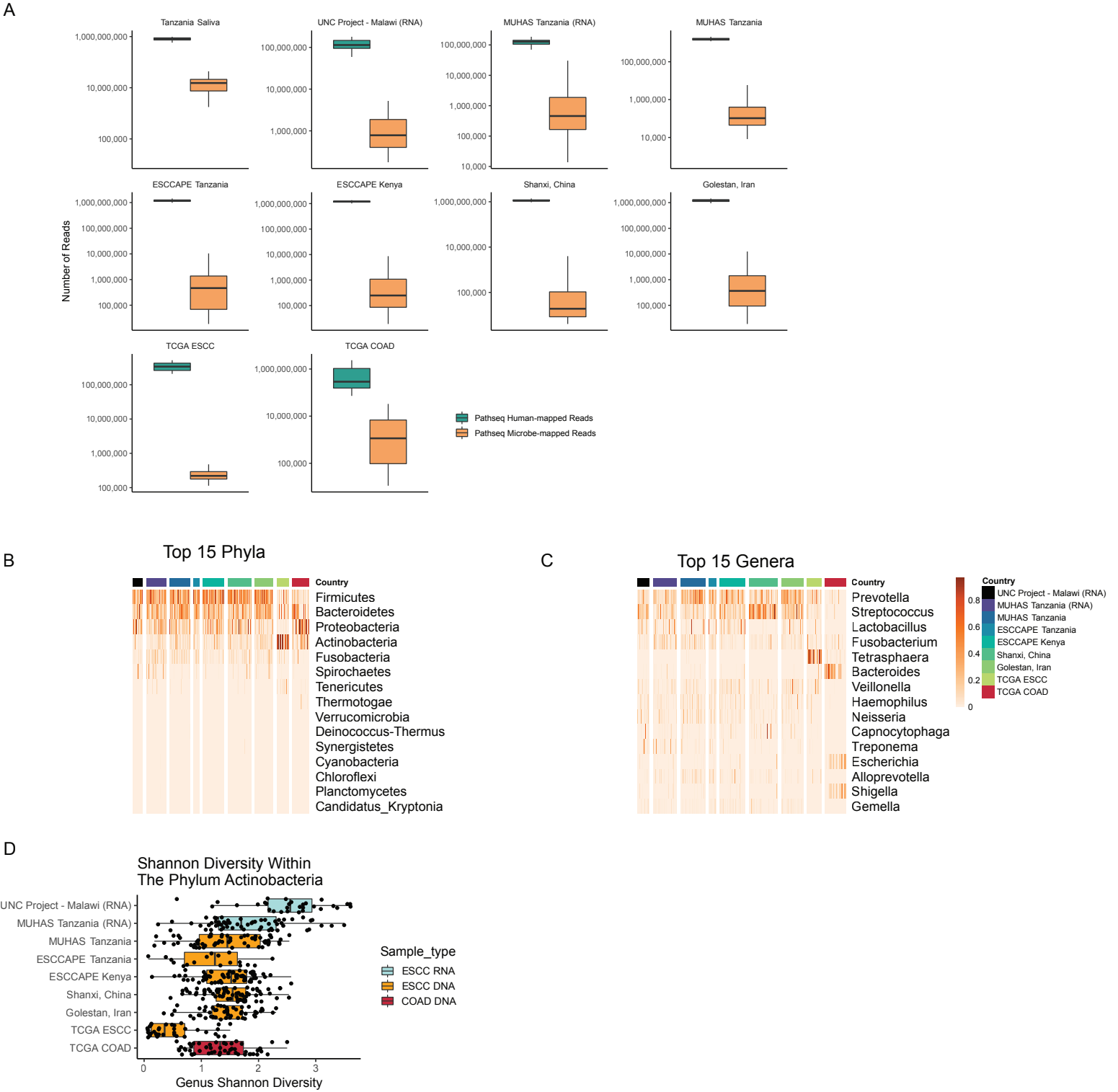
Figure S1



Figure S1. GATK-PathSeq statistics and extended phyla and genera information.

A.   Boxplots indicating the number of GATK-PathSeq Human-mapped reads and GATK-PathSeq microbe-mapped reads for each patient cohort. Samples from each cohort are WGS unless noted with "(RNA)", in which case they are RNAseq.

B.   Heatmap describing the relative abundance of the 15 top phyla sorted by average phylum relative abundance. Each column represents one sample. Rows represent the indicated phyla. Units are relative abundance. Samples from each cohort are WGS unless noted with "(RNA)", in which case they are RNAseq.

C.   Heatmap describing the relative abundance of the 15 top genera sorted by average genera relative abundance. Each column represents one sample. Rows represent the indicated genera. Units are relative abundance. Samples from each cohort are WGS unless noted with "(RNA)", in which case they are RNAseq.

D.   Boxplot representing the Shannon diversity of genera that fall within the phylum Actinobacteria for each patient in each cohort. Samples from each cohort are WGS unless noted with "(RNA)", in which case they are RNAseq.
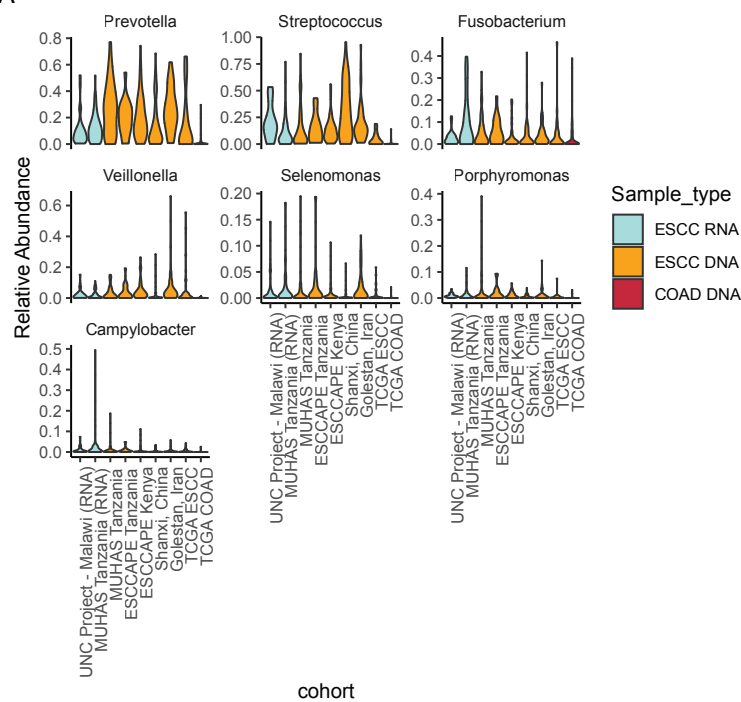
Figure S2

A



Figure S2. Distribution of Fusobacterium, Selenomonas, Prevotella, Streptococcus, Porphyromonas, Veillonella, and Campylobacter relative abundance of genus reads for all samples in each study.

A.    The distribution of the relative abundance of genus-mapping reads for seven selected genera in all studies. The width of each violin represents the proportion of samples which have the indicated relative abundance of each genus. In contrast to Figure 2D, which only plots up to 50 samples per study, this plot includes all patients. Samples from each study are WGS unless noted with "(RNA)", in which case they are RNAseq.

| TABLE 1 | | | | | | |
|---|---|---|---|---|---|---|
| **Study** | **Tanzania** | **Malawi**\*\* | **ESCCAPE Tanzania**\*\*\* | **ESCCAPE Kenya**\*\*\* | **East Golestan, Iran**\*\*\*\* | **Shanxi, China**\*\*\*\* |
| No. cases included | 61 | 30 | 18 | 65 | 55 | 71 |
| *Demographics* | | | | | | |
| Median age (IQR) | 49 (44-62) | 56 | 65 (61-73) | 64 (53, 71) | 62 (54,73) | 56 (50, 64) |
| % male | 67% | 45.8% | 61% | 68% | 55% | 56% |
| *Status at diagnosis* | | | | | | |
| Weight (kg), median (IQR) | | | 44 (40-52) | 52 (46, 60) | | |
| Body mass index (kg/m$^2$) median (IQR) | | | 15.8 (15.4, 19.1) | 19.5 (15.6, 22.0) | | |
| Median months ill before coming to endoscopy (IQR) | | | 2 (1, 6) | 3 (2, 4.5) | | |
| HIV status:       Positive | 2 (3.2%) | 10 (16.9%) | 1 (5%) | 5 (8%) | | |
| Negative | 36 (59.0%) | 44 (74.6%) | 10 (56%) | 48 (74%) | | |
| Not known | 23 (37.7%) | 5 (8.5%) | 7 (39%) | 12 (18%) | | |
| *Key lifestyle habits* | | | | | | |
| N (%) ever tobacco users | | | 11 (61%) | 38 (58%) | 17 (31%) | 35 (49%) |
| N (%) who brush teeth daily:       With toothbrush | | | 12 (67%) | 16 (25%)\* | | |
| With stick | | | 6 (33%) | 10 (15%) | | |
| Median no missing teeth (IQR) | | | 3 (1, 5) | 4 (1, 8) | | |

\*N=22 (34%) brush once per week or never, n=17 (26%) brush 2 to 6 times/week

\*\*Indicates demographics are from the entire patient population, consisting of both included and unincluded patients.

\*\*\*Indicates demographic percentages are from the entire patient population, with discrete counts scaled to the number of cases included.

\*\*\*\*Indicates demographic information is exclusively for included patients.