



Deposited via The University of Leeds.

White Rose Research Online URL for this paper:

<https://eprints.whiterose.ac.uk/id/eprint/241213/>

Version: Accepted Version

---

**Article:**

Dohlman, A.B., Mjelle, R., Wood, H.M. et al. (2026) Biodiversity and biogeography of the multi-kingdom cancer microbiome. Cell. ISSN: 0092-8674

<https://doi.org/10.1016/j.cell.2026.04.015>

---

**Reuse**

This article is distributed under the terms of the Creative Commons Attribution (CC BY) licence. This licence allows you to distribute, remix, tweak, and build upon the work, even commercially, as long as you credit the authors for the original work. More information and the full terms of the licence here:

<https://creativecommons.org/licenses/>

**Takedown**

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing [eprints@whiterose.ac.uk](mailto:eprints@whiterose.ac.uk) including the URL of the record and the reason for the withdrawal request.

# 1 **Biodiversity and biogeography of the multi-kingdom cancer** 2 **microbiome**

3 Anders B. Dohlman<sup>1,2,3\*</sup>, Robin Mjelle<sup>4</sup>, Henry M. Wood<sup>5</sup>, Kevin Jiang<sup>1,2,3</sup>, Alaina Shumate<sup>1,2,3</sup>, Iris  
4 Lee<sup>1,2,3</sup>, Gianmarco Piccinno<sup>7</sup>, Garazi Serna<sup>8</sup>, Abdul-Rakeem Yakubu<sup>1,2,3\*</sup>, Paolo Nuciforo<sup>8</sup>, Phil  
5 Quirke<sup>5-6</sup>, Curtis Huttenhower<sup>9</sup>, Nicola Segata<sup>7</sup>, Matthew Meyerson<sup>1,2,3,10\*</sup>

6

7 <sup>1</sup> Department of Medical Oncology, Dana-Farber Cancer Institute, Boston, MA, USA

8 <sup>2</sup> Department of Genetics, Harvard Medical School, Boston, MA, USA

9 <sup>3</sup> Cancer Program, Broad Institute of MIT and Harvard, Boston, MA, USA

10 <sup>4</sup> Department of Clinical and Molecular Medicine, Norwegian University of Science and  
11 Technology, Trondheim, Norway

12 <sup>5</sup> Pathology and Data Analytics, Leeds Institute of Medical Research at St. James's, University of  
13 Leeds, Leeds, UK

14 <sup>6</sup> National Institute for Health and Care Research (NIHR) Leeds Biomedical Research Centre,  
15 Leeds, UK

16 <sup>7</sup> Department of Cellular, Computational and Integrative Biology (CiBIO), University of Trento,  
17 Trento, Italy

18 <sup>8</sup> Molecular Oncology Group, Vall d'Hebron Institute of Oncology, Barcelona, Spain

19 <sup>9</sup> Department of Biostatistics, Harvard School of Public Health, Boston, MA, USA

20 <sup>10</sup> Lead contact

21 \* Correspondence:

22 [andersb\\_dohlman@dfci.harvard.edu](mailto:andersb_dohlman@dfci.harvard.edu)

23 [matthew\\_meyerson@dfci.harvard.edu](mailto:matthew_meyerson@dfci.harvard.edu)

24

25 **Summary**

26 Microorganisms represent an important component of the tumor microenvironment, but conflicting  
27 reports have left the extent of microbial prevalence across cancer types unclear, necessitating  
28 more robust methods for characterizing tumor-associated microbiomes. We built and  
29 benchmarked a host-subtraction and classification pipeline to identify microbiota in whole-  
30 genome sequencing data and applied it to 16,369 high-depth tumor whole genomes from the UK  
31 100,000 Genomes Project. After decontamination, microbial signatures were indistinguishable  
32 from background in most cancer types. However, in orodigestive tumors, we detected multi-  
33 kingdom polymicrobial communities including bacteria, fungi, viruses, archaea, and in some  
34 cases, *Trichomonas*, a protozoan parasite. These communities varied by tumor site and subtype,  
35 with increased microbial colonization of microsatellite-unstable and POLE/POLD1-mutated  
36 tumors, supported by a correlation between microbial load and tumor mutation burden observed  
37 across cancer types. This analysis helps to resolve pan-cancer microbial structure and links the  
38 tumor microbiome to host phenotype and tumor genomic context.

39

40 **Keywords**

41 cancer, tumor microbiome, biogeography, decontamination, multi-kingdom, bacteria, fungi,  
42 viruses, archaea, mutation burden

43

44

45

46

47

48

49

50

## 51 Introduction

52 The human microbiome is composed of many trillion cells, including bacteria, fungi,  
53 viruses, and archaea<sup>1-4</sup>. In healthy individuals, these microorganisms are found at anatomical  
54 barrier surfaces, including the skin, upper respiratory tract, urogenital tract, and are most  
55 concentrated in the orodigestive tract<sup>1,2,5</sup>, with substantial variation in community structure across  
56 body sites<sup>4,6-12</sup>. These commensal microbiota play essential roles in human metabolism,  
57 digestion, and immunity<sup>4,13-15</sup>. Perturbations to these microbial communities have been  
58 associated with various human health conditions<sup>16-23</sup>, including cancer<sup>24-27</sup>, where variation in gut  
59 microbial composition has been reported to be linked with response to immunotherapy in both  
60 melanoma and lung cancer<sup>28-33</sup> and to chemotherapy efficacy in pancreatic cancer<sup>34</sup>.

61 Beyond systemic interactions with their hosts, microbiota can directly colonize tumors,  
62 exerting influence on both tumor cells and their surrounding microenvironment<sup>24,35-39</sup>. To date,  
63 some of the most robust evidence for tumor-associated microbiota comes from cancers arising at  
64 body sites known to harbor microbes under normal physiological conditions, including human  
65 papillomavirus (HPV) in oral and cervical cancer<sup>26</sup>; Epstein-Barr virus (EBV) in nasopharyngeal<sup>40</sup>  
66 and gastric<sup>41</sup> cancers; hepatitis B and C viruses in liver cancer<sup>42</sup>; and *Helicobacter pylori* in gastric  
67 cancer<sup>43,44</sup>. In colorectal cancer (CRC), multiple bacteria have been reproducibly linked to the  
68 disease, including *Fusobacterium nucleatum*<sup>35,36</sup>, enterotoxigenic *Bacteroides fragilis*<sup>45,46</sup>, and  
69 *pks+* *Escherichia coli*<sup>47-52</sup>. Collectively, microbial infections are estimated to account for at least  
70 13% of global cancer burden<sup>53</sup>.

71 For cancers occurring at non-barrier sites, there have been inconsistent reports regarding  
72 the scope of microbial colonization. Pan-cancer surveys of tumor tissues have yielded widely  
73 divergent findings regarding the presence, distribution, and cancer-type specificity of microbial  
74 communities<sup>39,54-63</sup>, contributing to significant uncertainty in the field. Many of these prior efforts  
75 have relied on large-scale tumor sequencing datasets such as The Cancer Genome Atlas  
76 (TCGA)<sup>57,59,60,62,64</sup>, the UK 100,000 Genomes Project (100kGP)<sup>61</sup>, and metastatic cancers from  
77 the Hartwig Medical Foundation<sup>65</sup>, each containing thousands of sequenced tumor samples.  
78 While these datasets harbor verifiable tumor-associated microbial sequences<sup>24,57,61,62,64</sup>, their  
79 utility for capturing pan-cancer microbial structure has been hindered by the challenges of  
80 sequence misclassification, contamination, batch effects, and the absence of gold-standard  
81 negative controls<sup>57,58,64</sup>, underscoring the need for more rigorous measures to distinguish true,  
82 tumor-specific microbial signatures from false positives.

83 To address technical false positives arising from misclassification of human sequences,  
84 our group previously developed PathSeq<sup>66,67</sup>, a computational host-subtraction pipeline designed  
85 for detecting microbiota in low-biomass tissues. Applying this approach allowed our discovery of  
86 *Fusobacterium nucleatum* enrichment in colorectal cancer<sup>35</sup> while similar methods were deployed  
87 by others to discover Merkel cell polyomavirus in Merkel cell carcinoma<sup>68</sup>. Such host subtraction  
88 measures are essential not just for mitigating spurious microbial signals but also for preventing  
89 inadvertent exposure of patient-identifiable genomic information<sup>69</sup>. Applying PathSeq to TCGA,  
90 members of our group later used developed and benchmarked models to distinguish tumor-  
91 associated microbial signals from laboratory contamination, enabling pan-cancer characterization  
92 of tumor-associated bacteria<sup>57</sup> and fungi<sup>59</sup>.

93 However, host filtering approaches are inherently limited by the incompleteness of the  
94 host reference genome used to perform subtraction. Moreover, prior microbial analyses have not  
95 been comprehensively verified in independent, comparably large-scale tumor sequencing cohorts  
96 of primary tumor tissue. Today, it is possible to study the cancer microbiome at much greater  
97 resolution, due to the availability of thousands of high-coverage whole-genome sequences (WGS)  
98 of cancer DNA from the UK 100,000 Genomes Project (100kGP)<sup>70-74</sup> as well as the completed  
99 telomere-to-telomere human reference genome T2T-CHM13<sup>75,76</sup> for host subtraction.

100 To address the longstanding challenge of reliably detecting microbiota in tumor  
101 sequencing data and to resolve ongoing uncertainties regarding the distribution of tumor-  
102 associated microbial communities, we implemented and benchmarked PathSeq-T2T, an updated  
103 filtering pipeline that incorporates the complete telomere-to-telomere human genome reference  
104 (T2T-CHM13) for host subtraction and multiple metagenomic classifiers, enabling high-  
105 confidence detection of microbial signals. We then applied this pipeline to the pan-cancer 100kGP  
106 database, which includes 16,369 high-depth tumor whole genomes from 15,237 participants. The  
107 resulting atlas of decontaminated microbial profiles allowed a far more detailed view of the cancer  
108 microbiome than was previously possible.

109

110

111

112

## 113 Results

### 114 ***A host-subtraction pipeline for detecting microbiota in low-biomass human tissues***

115 In tissue samples with low microbial biomass, the incomplete removal of human  
116 sequences can produce false-positive microbial signatures<sup>58</sup> and may expose patient-identifiable  
117 genomic information<sup>69</sup>. To improve the detection of microbial signals in sequenced tumor tissues,  
118 we developed PathSeq-T2T, a host-subtraction and microbial detection pipeline that leverages  
119 the complete T2T-CHM13 human reference genome to isolate and then classify non-human reads  
120 from bulk tumor whole-genome sequencing data (**Figure S1A; STAR Methods**). PathSeq-T2T  
121 performs quality and complexity filtering on reads not mapped to the standard human genome  
122 reference (GRCh38), together with additional subtractive host read filtering against polymorphic  
123 human immune loci (e.g., MHC genes), known breakpoint junctions, vector sequences, and T2T-  
124 CHM13. PathSeq-T2T then screens the remaining putative non-human sequences with three  
125 microbial classification tools: Kraken2<sup>77</sup> (which uses k-mer-counting to classify individual reads),  
126 MetaPhlan4<sup>78</sup> (which detects clade-specific marker genes), and Sylph<sup>79</sup> (which measures k-mer  
127 containment within microbial genomes). This multi-classifier approach is included to allow cross-  
128 validation of microbial signals. After classification, PathSeq-T2T calculates the observed microbial  
129 reads per million starting reads (RPM), allowing quantification of microbial tumor burden.

### 130 ***Benchmarking PathSeq-T2T with in silico mixtures of microbial and human DNA***

131 To assess PathSeq-T2T's ability to detect trace microbial signals in human tissue, we  
132 generated sequences from an *in silico* microbial community containing 17 species of bacteria,  
133 fungi, and archaea, then mixed them with human sequences at varying dilutions. This resulted in  
134 nine, tenfold serial dilution conditions containing pre-defined read counts for each species (**Table**  
135 **S1**). These began at a microbe-to-host ratio of 1:1 and extended to a final dilution condition of  
136 1:10<sup>8</sup> which contained just one microbial read pair. We also included controls containing only  
137 human or only microbial reads.

138 PathSeq-T2T was consistently effective in removing ground-truth human sequences  
139 across the in silico dilution series, leaving zero human reads in every condition (**Figure 1A, Figure**  
140 **S1B**). By comparison, host filtering with just the GRCh38-unaligned read set (an approach used  
141 in studies conducted prior to the release of T2T-CHM13<sup>55,60,64</sup>) resulted in the retention of 5 million  
142 human reads (>17,000 RPM) in the human-only condition. Meanwhile, filtering with the original

143 PathSeq pipeline retained 5,128 human reads (~17 RPM) in the human-only condition. Thus,  
144 PathSeq-T2T substantially outperforms previous methods of host subtraction.

145 After host filtering, PathSeq-T2T correctly identified the absence of microbial sequences  
146 in the human-only control. In conditions containing microbial reads, PathSeq-T2T retained on  
147 average  $79.8\% \pm 7.6\%$  of microbial sequences (**Figure 1B**). Collateral subtraction of true  
148 microbial reads primarily affected microbial sequences with low quality and/or complexity and  
149 duplicated reads (**Figure S1C**). The removal of these sequences resulted in a substantial  
150 reduction in read misclassifications (**Figure S1D**).

151 Overall, the final filtered read counts in the *in silico* mixtures closely matched expected  
152 microbe-to-host ratios. The largest absolute subtraction in human reads occurred during the hg38  
153 alignment step, whereas the largest fractional reduction occurred during T2T-CHM13 subtraction  
154 (**Figure 1C, Table S2**). Among conditions containing microbial reads,  $98.3 \pm 0.4\%$  of PathSeq-  
155 T2T-filtered reads were correctly assigned to on-target microbial taxa, while the remaining reads  
156 could not be uniquely classified (**Figure 1D**). Notably, in the  $1:10^8$  dilution condition, PathSeq-  
157 T2T retained the lone microbial read pair among >300 million starting reads and correctly  
158 classified it as *Veillonella rogosae* after removing all human reads, indicating robust host  
159 subtraction and high sensitivity at low microbial biomass.

### 160 ***Benchmarking PathSeq-T2T with in vitro mixtures of microbial and human DNA***

161 Although PathSeq-T2T performed effectively on our *in silico* microbe-host mixtures, such  
162 data are unable to capture environmental contamination introduced by laboratory manipulations  
163 during sample preparation, DNA extraction, and sequencing. We next generated mixtures of  
164 microbial and human DNA *in vitro*, using the same experimental design as before. Human and  
165 microbial DNA were extracted from a colorectal cancer cell line and microbial community  
166 standard, respectively, then mixed in triplicate across the same serial dilution ladder (1:1 to  $1:10^8$ ).  
167 We then sequenced each mixture to a target of 30X human genome coverage (See STAR  
168 Methods). We also included host-only, microbe-only, and template-free extraction controls.

169 Consistent with the *in silico* mixtures, PathSeq-T2T filtering of *in vitro* mixtures achieved  
170 clear separation of the dilution series, with on-target microbial read counts largely matching  
171 expected ratios (**Figure 1E, Table S3**). Moreover, PathSeq-T2T outperformed host filtering with  
172 GRCh38 alone or with the original PathSeq pipeline (**Figure 1F**), demonstrating effective human

173 read removal and microbial detection on real sequencing data. Similar to our *in silico* mixtures, a  
174 subset of *in vitro* sequencing reads could not be uniquely classified, especially at lower dilutions.

175 In addition to these on-target microbial sequences, we detected evidence of contamination  
176 across nearly all *in vitro* experimental conditions, even after sterile handling during DNA extraction  
177 and mixing (**Figure 1G**). This included off-target, skin-associated taxa<sup>12</sup> such as *Cutibacterium*  
178 *acnes*, *Staphylococcus epidermidis*, and *Malassezia restricta* (**Figure 1H**), which were not  
179 included in the community standard. Each of these off-target species was also detected in  
180 template-free negative controls, confirming their status as contaminants (**Figure 1H**). While  
181 contamination from these species varied across conditions, their read counts were highly  
182 correlated across samples, suggesting they were introduced in unison. The proportion of off-target  
183 microbial sequences generally increased with successive dilutions, despite some variability at  
184 lower dilution conditions (**Figure 1J**). This mirrors prior observations that contamination  
185 disproportionately affects low-biomass samples, even as the absolute level of contamination  
186 remains consistent<sup>80–82</sup>.

187 We also observed low levels of on-target microbial sequences in human-only samples that  
188 did not receive aliquots of microbial DNA (**Figure 1G**). This is a phenomenon that can occur due  
189 to cross-contamination between sample conditions<sup>81,82</sup> or due to barcode swapping on the flow  
190 cell<sup>82–84</sup>. Collectively, these sequences accounted for 82 to 246 reads per sample (0.1 to 0.3  
191 RPM). We therefore determined 0.1 RPM as the threshold for expected background that could be  
192 expected from cross-contamination in otherwise sterile samples.

### 193 **Testing PathSeq-T2T on cancer cell lines**

194 We next sought to evaluate PathSeq-T2T performance on genomically distinct cancer cell  
195 lines, which are not expected to harbor microbiota. We tested microsatellite-stable (MSS;  $n = 12$ )  
196 or microsatellite-unstable (MSI;  $n = 8$ ) colorectal cancer cell lines with whole-genome sequencing  
197 data available from The Cancer Cell Line Encyclopedia<sup>85</sup>. No differences were observed between  
198 MSS and MSI cell lines, indicating that the presence of high mutation burden did not critically  
199 impact host filtering or microbial detection (**Figure S1E**). However, we did observe a substantial  
200 enrichment of microbial sequences in cell lines with DNA libraries prepared using PCR (**Figure**  
201 **S1F**). Examining the microbial composition of these cell lines, we observed variation in taxa  
202 between library types. Both groups included common contaminants such as *Cutibacterium*,  
203 *Pseudomonas*, *Sphingomonas*, *Variovorax*, and *Bradyrhizobium* (**Figure S1G**).

204 **Testing PathSeq-T2T on tumor and germline samples from 100kGP**

205 Having validated and benchmarked PathSeq-T2T with *in silico* and *in vitro* host-microbe  
206 mixtures as cancer cell lines, we began to evaluate the pipeline's performance on real-world,  
207 tumor whole-genome sequencing data from the UK 100,000 Genomes Project (100kGP) (**Table**  
208 **S4**). Patient biospecimens sequenced by this project include both tumor tissues ( $n = 16,369$ ,  
209  $\sim 100X$  coverage) and matched samples collected for the purpose of germline profiling ( $n =$   
210  $15,249$ ,  $\sim 30X$  coverage), including peripheral blood and saliva. As an initial performance analysis  
211 prior to studying the entire 100kGP cancer database, we compared tumor and germline samples  
212 with a range of expected microbial biomass. Among germline samples, we compared saliva ( $n =$   
213  $490$ ), which is typically microbe-rich, and peripheral blood ( $n = 14,127$ ), which is believed to be  
214 microbe-low or sterile under normal conditions<sup>86</sup>. Among primary tumor sequences, we compared  
215 colorectal tumors (CRC,  $n = 2,482$ ) and glioblastoma multiforme brain tumors (GBM,  $n = 295$ ).

216 On these sample types, PathSeq-T2T yielded expected sample-type- and cancer-type-  
217 specific variation in the proportion of sequences retained after filtering (**Figure 2A-B**). Filtering  
218 removed  $>99.99\%$  of reads in blood and GBM, while retaining a median of  $3.87\%$  and  $0.02\%$  of  
219 reads from saliva and CRC samples, respectively. Among germline samples, we detected  
220 microbial reads at much higher levels in saliva (median  $27,500$  RPM) than blood ( $1.7$  RPM), while  
221 among tumor samples, we detected more microbial reads in CRC samples ( $145$  RPM) than GBM  
222 ( $0.7$  RPM). This indicated that PathSeq-T2T was correctly retaining microbial sequences in  
223 sample types known to harbor microbiota, while mostly eliminating them in sample types that are  
224 sterile or microbe-low.

225 To further understand the sequences that remained in these sample types after host  
226 filtering, we examined their microbial composition. The taxa found in saliva samples was highly  
227 distinct from those found in blood (**Figure 2C**), with saliva frequently being dominated by  
228 *Prevotella*, *Streptococcus*, *Rothia*, and *Veillonella*, which are typical of the oral cavity. Similarly,  
229 colorectal cancer (CRC) tissues exhibited a distinct microbial composition compared to  
230 glioblastoma (GBM) (**Figure 2D**), containing high rates of *Bacteroides*, *Phocaeicola*, and  
231 *Fusobacterium*, which are known to be found in CRC tissues<sup>8</sup>. Blood and GBM samples typically  
232 harbored similar genera, including *Achromobacter*, *Sphingomonas*, *Pseudomonas*, *Variovorax*,  
233 *Bradyrhizobium*, and *Cutibacterium*. As many of these taxa are commonly associated with human  
234 skin, the environment, and laboratory reagents<sup>80,81</sup>, we suspected that many of the microbiota  
235 detected in blood and GBM samples were contaminants.

236 To further interrogate this possibility, we compared species prevalence between sample  
237 types. Between saliva and blood samples, most detected species were unique to saliva, while a  
238 smaller subset of species were found at nearly equal rates in both sample types (**Figure 2E**,  
239 **Figure S2A-B**). A similar pattern was observed comparing CRC and GBM (**Figure 2F**, **Figure**  
240 **S2C-D**). However, no species were found to be unique to blood or GBM, suggesting that taxa  
241 found in these microbe-low sample types were likely contaminants shared across cancer types.  
242 Consistent with this, species equiprevalent in saliva/blood ( $n = 125$ ) and in CRC/GBM ( $n = 161$ )  
243 overlapped substantially ( $n = 110$ ; Fisher's  $p = 1.6e-98$ ), suggesting a common pool of  
244 contaminants (**Figure S2E**). Examining these equiprevalent species further supported their status  
245 as contaminants. In blood, the most prevalent such taxa were *Pseudomonas aeruginosa* and  
246 *Variovorax paradoxus* (**Figure 2G**), while among CRC/GBM-equiprevalent taxa, prominent  
247 species included *Cutibacterium acnes* and *Bradyrhizobium sp. BTAi1* (**Figure 2H**). By  
248 comparison, the most frequently detected species in saliva and CRC samples were canonically  
249 oral<sup>11</sup> and colorectal<sup>8</sup> microbiota (**Figure S2F-G**). These results showed that, in addition to  
250 recovering biologically plausible and sample-type-specific microbiota, PathSeq-T2T was also  
251 detecting contaminants that appeared to be present across sample types. This demonstrated the  
252 need for additional decontamination steps beyond host filtering alone.

### 253 **Decontamination of 100kGP microbiomes**

254 We next examined microbial sequences across the entire 100kGP dataset, which  
255 amounted to a total of 16,369 tumor and 15,249 germline samples spanning 28 distinct cancer  
256 types. Since contaminants seemed to appear at similar rates across all sample types, we  
257 designed a pan-cancer equiprevalence (PCE) score to quantify how evenly distributed a species  
258 was across distinct cancer types in the 100kGP cohort. This score was calibrated such that a  
259 species would receive a score of 0 if it was detected in 100% of samples from just one cancer  
260 type, and a score of 1 if it was detected at equal frequency in every cancer type (see STAR  
261 Methods). Applying this score, we found that many of the species that were equiprevalent in  
262 saliva/blood or CRC/GBM also appeared to be equiprevalent across the entire 100kGP cohort  
263 (**Figure 2I**, **Figure S2H-I**, **Table S5**).

264 To assess the ability of pan-cancer equiprevalence to distinguish between environmental  
265 contamination and true tumor-associated taxa, we examined species known to be associated with  
266 cancer (**Figure 2J**). This recapitulated expected cancer-type-specific microbial colonization  
267 patterns. *Helicobacter pylori* (PCE = 0.32), a known driver of gastric cancer, was enriched in  
268 gastric tumors; *Fusobacterium nucleatum* (PCE = 0.55), an oral pathogen associated with

269 colorectal cancer, was enriched in both oral squamous and colorectal cancers; and the human  
270 papillomavirus 16 species group (HPV; PCE = 0.25), implicated in oropharyngeal carcinogenesis,  
271 was specific to that disease. By comparison, skin-associated taxa *Cutibacterium acnes* (PCE =  
272 0.90), *Staphylococcus epidermidis* (PCE = 0.78), and *Malassezia restricta* (PCE = 0.86) were  
273 equiprevalent across cancer types, indicating that 100kGP samples harbored similar  
274 contaminants to those found in our *in vitro* benchmarking experiment. Across the dataset, the  
275 most equiprevalent species were *Pseudomonas aeruginosa* (0.91), *Stutzerimonas stutzeri* (0.91),  
276 *Acidovorax temperans* (0.91), *Pseudomonas oleovorans* (0.91), and *Sphingobium yanoikuyae*  
277 (0.90) (**Figure S2J**), indicating that these species were also likely contaminants.

278 As known contaminants consistently had PCE scores above 0.7 (**Figure 2I**), we selected  
279 this value as a threshold for identifying and removing contaminant taxa. Interestingly, examining  
280 microbial genomes assembled from sequenced tumor tissues, we observed that species  
281 exceeding this threshold shared near-identical, strain-level average nucleotide identity (>99.9%)  
282 despite coming from completely different samples (**Figure S2K, Table S6**), suggesting a  
283 common, exogenous contamination source. Other non-contaminant species exhibited greater  
284 genomic diversity between samples, indicating a more distant common ancestor.

285 Consistent with our analysis of cancer cell lines, putative contaminants were also  
286 frequently enriched in tumor and germline samples whose DNA libraries had undergone PCR  
287 amplification (**Figure S2J-K**), indicating contamination associated with PCR reagents and/or  
288 amplification of contaminating microbial DNA. Contamination was also pronounced in formalin-  
289 fixed, paraffin-embedded (FFPE) samples compared to fresh-frozen tissues (**Figure S2L**). We  
290 therefore restricted downstream analyses to fresh-frozen solid tumors sequenced from PCR-free  
291 libraries.

### 292 ***Tumor microbiota are largely restricted to cancers of the orodigestive tract***

293 We next examined microbial abundance before and after decontamination. Prior to  
294 removing contamination, most cancer sites exhibited microbial signals above background (**Figure**  
295 **S3A-C**), with the abundance of contaminant species generally consistent across cancer sites  
296 (**Figure S3D-F**). However, after decontamination, only a subset of cancer types – specifically  
297 colorectal, oropharyngeal, esophageal, and gastric cancers – consistently retained microbial  
298 signal exceeding background (RPM > 0.1). This finding was reproduced by Kraken2, MetaPhlan4,  
299 and Sylph (**Figure 3A-C**), which were highly concordant in their abundance estimates (**Figure**  
300 **S3G**).

301 Comparing species prevalence across cancer sites, we observed that tumors arising from  
302 oropharyngeal, esophageal, gastric, and colorectal sites also harbored DNA sequences from  
303 many cancer-site-specific taxa, while other cancers contained relatively few (**Figure 3D, Figure**  
304 **S3H**). Broadly, microbial abundance at these sites aligned with prior estimates of microbial density  
305 along the digestive tract<sup>1,2,5</sup>, with colorectal tumors showing the highest abundance (IQR: 18–285  
306 RPM), followed by oropharyngeal tumors (3–189 RPM), and then esophageal (0.3–24 RPM) and  
307 gastric cancers (0.4–19 RPM).

308 These cancer-associated microbial communities were typically polymicrobial (**Figure 3E**),  
309 with more than one species detected in most colorectal (94.7%), oropharyngeal (77.0%),  
310 esophageal (59.8%), and gastric (62.7%) tumors. Their microbial composition mirrored known  
311 patterns of site-specific biogeography and biodiversity observed in healthy individuals<sup>4,6–8,10</sup> as  
312 well as in prior studies of orodigestive cancers<sup>35,36,57,87–92</sup> (**Figure 3F**). Colorectal tumors  
313 commonly featured members of the *Bacteroidaceae* and *Lachnospiraceae* families, with  
314 *Fusobacteriaceae* enriched in both colorectal and oropharyngeal cancers. *Prevotellaceae* and  
315 *Streptococcaceae* were present across all four sites but were especially abundant in  
316 oropharyngeal and esophageal tumors. Gastric cancers were unique in harboring species from  
317 the *Lactobacillaceae* and *Helicobacteraceae* families. Meanwhile, the viral family  
318 *Papillomaviridae* (i.e., HPV-16) was uniquely detected in oropharyngeal tumors (**Figure S3I**).  
319 Cervical cancers, which also typically harbor HPV-16/18, were not included in the 100kGP cohort.

### 320 ***Sparse microbial signals in non-orodigestive tumors***

321 Outside of orodigestive cancer types, microbial signals were highly sparse. Examining the  
322 subset of samples from other cancer types that were microbe-high (>10 RPM), we observed some  
323 signals that were biologically plausible (**Figure S3G**). For example, skin cancers<sup>93,94</sup> occasionally  
324 harbored Merkel cell polyomavirus (2.9% of cases) (**Figure S3J**) and *Staphylococcus aureus*  
325 (6.0%)<sup>95–97</sup>. Some pancreatic tumors<sup>37,38,98,99</sup> harbored gut-typical taxa (*Veillonellaceae*,  
326 *Fusobacteriaceae*, and *Enterobacteriaceae*), potentially reflecting proximity to the digestive tract,  
327 while a subset of bladder tumors<sup>100–102</sup> harbored species previously observed in the genitourinary  
328 tract, including *Peptoniphilus genitalis* and *Actinotignum timonense/shaali*<sup>103,104</sup>.

329 Across other cancer types, microbe-high samples more often resembled acute infection  
330 rather than resident communities. For example, one lung cancer contained DNA from  
331 *Streptococcus pneumoniae* (10 RPM) while another contained *Haemophilus influenzae* (25  
332 RPM), both common causes of pneumonia<sup>105</sup>. The remaining microbe-high samples from non-

333 orodigestive tissues were typically dominated by *Herpesviridae* species, which are mostly  
334 widespread among humans and acquired early in life<sup>106–110</sup>. Among *Herpesviridae*, HSV, EBV,  
335 CMV, and HHV-7 displayed varying degrees of cancer-type specificity (**Figure S3I**), while HHV-  
336 6A and HHV-6B were more broadly distributed (**Figure S3J**).

337 Thus, although microbial signals were generally sparse in non-orodigestive cancers, some  
338 harbored biologically plausible taxa. These signals could represent acute infections or potentially  
339 more prolonged colonization through mechanisms that have not yet been determined. However,  
340 they could also have arisen due to cross-contamination from adjacent organs, barcode swapping,  
341 or sample mislabeling.

### 342 ***Biogeography of oropharyngeal, gastroesophageal, and colorectal cancers***

343 Since tumors of the orodigestive tract consistently harbored abundant, polymicrobial  
344 communities, we next sought to more thoroughly characterize and quantify these communities  
345 across anatomic sites. Concordant with prior biogeography studies of healthy tissues<sup>1,2,5</sup>, we  
346 observed the highest microbial load and diversity in tumors of the lower gastrointestinal tract,  
347 followed by those of the oral cavity (**Figure 4A, Figure S4A**), while microbial abundance and  
348 diversity were generally lower in esophageal and gastric sites.

349 Using microbial read counts to estimate bacterial density (**Figure S4B**), we predicted a  
350 median burden of 20-30 bacteria per 1,000 human cells in oral and sinonasal cancers, while  
351 bacteria-host cell ratios were lower in esophageal and gastric sites (1-3 bacteria per 1,000 cells).  
352 In the lower gastrointestinal tract, microbe-host cell ratios varied by anatomic site, and generally  
353 declined progressively from the cecum (227 per 1,000), through the ascending (251 bacteria per  
354 1,000), transverse (229 per 1,000), and descending colon (119 per 1,000), reaching their lowest  
355 rates in the rectum (77 per 1,000). However, these bacteria-to-host ratios were highly variable  
356 within individual sites and sometimes even exceeded 1:1, consistent with the substantial  
357 heterogeneity observed in spatial analyses<sup>111</sup> and tumor-associated biofilms<sup>112</sup>.

358 We next examined variation in community structure across sites. Principal coordinate  
359 analysis showed that oropharyngeal, gastroesophageal, and colorectal tumors harbored distinct  
360 microbial communities ( $p < 0.001$ ) (**Figure 4B**). Cancers of the oropharynx and gastroesophageal  
361 tract showed some site-specific variation but primarily clustered together, consistent with their  
362 anatomic proximity, while cancers of the lower gastrointestinal tract clustered separately.

363 Consistent with this, we found both distinct and shared microbial taxa at each site (**Figure**  
364 **4C, Figure S4C**). Of the 543 species detected in at least 1% of samples from any of these three  
365 sites, the largest fraction was unique to colorectal cancer (42.0%), while fewer were unique to  
366 oropharyngeal (11.8%) and gastroesophageal (2.6%) cancers. The degree of species overlap  
367 between these sites was greatest between oropharyngeal and gastroesophageal cancers (8.5%)  
368 followed by oropharyngeal and colorectal (7.9%), while very few species (0.9%) were shared  
369 exclusively between gastroesophageal and colorectal cancers. A core set of species (26.3%)  
370 were detected across all three sites.

371 This variation was mirrored in relative abundance profiles (**Figure 4D**), which further  
372 highlighted the biodiversity and polymicrobial nature of these cancer types. *Prevotella* was the  
373 most abundant genus in both oropharyngeal and gastroesophageal tumors, with a mean relative  
374 abundance of 19.4% and 17.7%, respectively, often accompanied by *Haemophilus*,  
375 *Streptococcus*, and *Leptotrichia*. Oropharyngeal cancers also frequently contained high relative  
376 abundances of *Fusobacterium* (9.9%), *Capnocytophaga* (5.3%), and *Porphyromonas* (2.6%),  
377 while a subset were dominated by *Alphapapillomavirus*/HPV (7.2%). Gastroesophageal tumors  
378 exhibited higher relative levels of *Veillonella* (6.6%), *Selenomonas* (6.6%), and *Helicobacter*  
379 (5.3%). Meanwhile, colorectal cancers had high average relative abundances of *Bacteroides*  
380 (27.3%) and *Phocaeicola* (10.6%), while also frequently containing *Fusobacterium* (5.7%),  
381 *Escherichia* (4.1%), and *Segatella* (3.7%).

382 Examining species prevalence along the orodigestive axis (**Figure 4E, Figure S4D**), we  
383 observed enrichment of the periodontal disease associated species<sup>113</sup> *Porphyromonas gingivalis*,  
384 *Treponema socranskii*, and *Capnocytophaga gingivalis* in oropharyngeal cancer, while in gastric  
385 and esophageal cancers, we observed site-specific enrichment for *Helicobacter pylori* and acid-  
386 tolerant species including *Limosilactobacillus fermentum* and *Lactobacillus gasser*<sup>114,115</sup>. Taxa  
387 shared between oropharyngeal and colorectal cancers included several oral-typical pathogens  
388 like *Fusobacterium animalis*, *Parvimonas micra*, *Solobacterium moorei*, and *Prevotella*  
389 *intermedia*, whereas *Bacteroides fragilis*, *Escherichia coli*, and *Akkermansia muciniphila* were  
390 largely colorectal-specific and varied between the proximal and distal colon. *Streptococcus*  
391 *anginosus*, which has been shown to promote gastric cancer in mice<sup>116</sup>, was detected at similar  
392 rates across orodigestive sites.

393 Overall, microbial load and composition along the orodigestive axis broadly recapitulated  
394 known biogeography and biodiversity. However, the frequent presence of oral-typical species in

395 non-oral tumor sites suggests a partial breakdown of normal tissue-site specificity in cancer. Thus,  
396 while anatomic location appears to be the primary determinant of tumor microbiome composition,  
397 cancer-associated changes to the local microenvironment during cancer initiation or progression  
398 may encourage colonization by allochthonous species.

399 ***Orodigestive cancers harbor a multi-kingdom community of bacteria, fungi, archaea,***  
400 ***viruses, and sometimes parasites***

401 Viruses and bacteria have long been known to infect certain tumor tissues<sup>26,35,36,40,41,43,53,68</sup>,  
402 while parasite infections have been linked to higher rates of bladder cancer in endemic areas<sup>117</sup>.  
403 More recently, studies have suggested that other domains, including fungi<sup>59,60</sup>, also colonize  
404 tumor tissues. Therefore, we sought to characterize the multi-kingdom microbial composition of  
405 tumors in the 100kGP cohort, including fungi, archaea, viruses, and parasites.

406 Consistent with our previous findings from TCGA tumor tissues<sup>59</sup>, we observed fungi  
407 across oropharyngeal, gastroesophageal, and colorectal cancers, including *Candida albicans*,  
408 *Nakaseomyces glabratus* (formerly *Candida glabrata*), and *Saccharomyces cerevisiae* (**Figure**  
409 **4F**). Each of these fungi exhibited pronounced site-specific colonization patterns: the most  
410 prevalent site for *C. albicans* was in esophageal (7.8%) cancers and in tumors arising at the  
411 gastroesophageal junction and gastric cardia (6.9%), while *N. glabratus* was enriched in non-  
412 cardia gastric cancers, with greatest prevalence in the antrum/pylorus (5.3%). The yeast *S.*  
413 *cerevisiae* was also detected in tumors of the gastric antrum/pylorus (5.3%) as well as throughout  
414 the lower gastrointestinal tract.

415 Examining other microbial kingdoms, we found enrichment of DNA sequences from  
416 archaea in tumors of the lower gastrointestinal tract. *Methanobrevibacter smithii* (**Figure 4G**), the  
417 predominant archaeon in the human gut<sup>10,118,119</sup>, was most prevalent in tumors of the splenic  
418 flexure (25.9%) and descending colon (25.6%), with the lowest prevalence in the cecum (8.6%).  
419 *M. smithii* is a methane-producing species that has been discussed as a factor in irritable bowel  
420 syndrome (IBS)<sup>120</sup> and colorectal cancer<sup>119</sup>. We also detected the related archaea  
421 *Methanobrevibacter sp. TLL-48-HuF1* and *Candidatus Methanomassiliicoccus intestinalis*  
422 (**Figure S4E**). Less is known about these other archaeal species, but both were originally isolated  
423 from human fecal samples and are reportedly methanogenic<sup>121,122</sup>. We did not detect any archaea  
424 in oropharyngeal or gastroesophageal sites.

425 In addition to fungi and archaea, we detected DNA sequences from the parasite  
426 *Trichomonas* in oropharyngeal cancers (2.6%), rectal cancers (0.8%), and cancers of the  
427 ascending colon (1.8%) (**Figure 4H**). *Trichomonas* are single-celled flagellated protozoan  
428 parasites and include the causative agents of trichomoniasis (*T. vaginalis*), a sexually transmitted  
429 infection that typically affects the urogenital tract<sup>123</sup>, with links to cervical cancer<sup>124</sup>, and *T. tenax*,  
430 which is associated with necrotizing gingivitis. Many *Trichomonas* infections are asymptomatic<sup>123</sup>.  
431 In the 100kGP cohort, *Trichomonas* was most abundant in oropharyngeal cancers, at rates up to  
432 12.75 RPM, corresponding to roughly one parasite per 1,000 human cells. *Trichomonas* presence  
433 was confirmed using both Kraken2 and MetaPhlan4 (no protozoa database was available for  
434 Sylph). The affected oropharyngeal cancer cases occurred almost exclusively in men with HPV-  
435 negative disease.

436 Finally, we detected site-specific viral signals across oropharyngeal, gastroesophageal,  
437 and colorectal cancers (**Figure 4I**). HPV-16 was found exclusively in oropharyngeal cancers  
438 (13%), while herpesviruses HHV-6A and HHV-6B were distributed throughout cancers of the  
439 lower gastrointestinal tract. CMV and EBV were detected in multiple sites, but EBV was most  
440 prevalent in gastric cancers of the fundus/body (13%). Interestingly, we also observed specificity  
441 of HHV-7 to gastric tumors, with highest prevalence in the gastroesophageal junction/cardia  
442 (10%) and the fundus/body (8.7%). Latent HHV-7 has previously been reported in gastric  
443 mucosa<sup>125</sup> and has been associated with EBV reactivation leading to mononucleosis<sup>126</sup>, however  
444 its role in gastric cancer remains largely unexplored.

445 Overall, these findings show that while bacteria are by far the most abundant and diverse  
446 microbial kingdom associated with orodigestive tumors, these cancers may also become  
447 colonized by a broader, multi-kingdom microbial community that includes viruses, fungi, archaea,  
448 and in some cases, parasites.

#### 449 ***Microbial load is elevated in hypermutated genomic subtypes of cancer***

450 In colorectal cancer, the tumor-associated microbiome is a key modulator of the tumor  
451 microenvironment<sup>27</sup>. Prior studies have reported differences in microbial composition and diversity  
452 between anatomical sites within the colon, both in healthy individuals<sup>4,6-10</sup> and in the context of  
453 cancer<sup>127,128</sup>, with the greatest differences observed between proximal (right-sided) and distal  
454 (left-sided) tumor sites.

455 Other studies have indicated that microbial colonization patterns may differ by genomic  
456 subtype, most notably by the observed enrichment of *Fusobacterium nucleatum* in microsatellite-  
457 instable colorectal cancer<sup>129,130</sup>. Genomically, colorectal cancers can become hypermutated due  
458 to deficiencies in DNA mismatch repair or due to mutations in the proofreading domains of DNA  
459 polymerase  $\epsilon$  (*POLE*) or  $\delta$  (*POLD1*)<sup>73,131–133</sup>. These characteristics respectively define  
460 microsatellite-*instable* (MSI; approximately 15% of colorectal cancer cases) and *POLE/POLD1*-  
461 mutated (polymerase-mutant; 1-2% of cases) genomic subtypes. Hypermutated subtypes appear  
462 most frequently in the proximal colon and are typically more immune-infiltrated, with a more  
463 favorable prognosis. In contrast, non-hypermutated cancers, described as microsatellite-*stable*  
464 (MSS; 80-85% of cases), are more common in the distal colon and rectum and are frequently  
465 driven by chromosomal instability<sup>132,133</sup>.

466 Leveraging the large size of the colorectal cancer cohort in 100kGP along with detailed  
467 clinical annotations and tumor genotyping provided by Cornish et al.<sup>73</sup>, we performed an in-depth  
468 characterization of colorectal microbiomes and their phenotypic and genotypic correlates. To  
469 begin, we performed a multivariable analysis (PERMANOVA) which revealed that tumor anatomic  
470 site ( $p < 0.001$ ) and genomic subtype ( $p < 0.001$ ) were the two strongest independent predictors  
471 of tumor microbiome composition in colorectal cancer (**Figure 5A**). These features together  
472 accounted for the majority of explainable variance in composition; however, we also observed  
473 weaker but statistically significant correlations with patient age ( $p < 0.05$ ) and history of  
474 radiotherapy ( $p < 0.05$ ).

475 To investigate how microbial colonization patterns differ by anatomic site and genomic  
476 subtypes, we next examined tumor microbial abundance and diversity. Consistent with Cornish  
477 et al., who reported elevated bacterial load in MSI colorectal cancers<sup>73</sup>, we found that overall  
478 bacterial load was significantly greater in hypermutated colorectal cancer cases compared to non-  
479 hypermutated cases, regardless of genomic mechanism. Compared to MSS cases, we observed  
480 3.9-fold greater microbial density ( $\pm 1.1$ ) in MSI cases ( $p = 4.61e-28$ ) and 6.5-fold greater density  
481 ( $\pm 1.6$ ) in polymerase-mutant cases ( $p = 2.35e-3$ ) (**Figure 5B**). This trend persisted even after  
482 stratifying hypermutated subtypes by tumor site (**Figure S5B**). Proximal and distal MSI tumors  
483 showed 3.2-fold ( $p = 4.26e-13$ ) and 2.9-fold ( $p = 5.28e-3$ ) enrichment in bacterial abundance  
484 compared to MSS, respectively, while proximal polymerase-mutant tumors harbored an 8.8-fold  
485 enrichment compared to proximal MSS ( $p = 2.09e-3$ ). The enrichment of microbiota in MSI and  
486 polymerase-mutant tumors appeared to be driven mostly by bacteria (**Figure S5C**), however we  
487 also observed a modest enrichment of fungal sequences in MSI compared to MSS (FC = 1.5;  $p$

488 = 1.8e-4) (**Figure S5D**). No significant differences in viral or archaeal load were observed (**Figure**  
489 **S5D-E**). Interestingly, despite increased microbial load in hypermutated colorectal cancers, we  
490 observed a decrease in microbial diversity in MSI cases ( $p = 4.64e-3$ ) (**Figure 5D, Figure S5F**),  
491 suggesting that this pattern was driven by a limited number of species rather than by a  
492 proportional expansion across the entire microbial community.

493 To validate these findings, we used decontaminated microbial profiles from The Cancer  
494 Microbiome Atlas<sup>57</sup> (derived from TCGA whole-genome sequencing), which confirmed the  
495 enrichment of bacterial load in MSI colorectal cancers ( $FC = 5.3 \pm 1.78$ ;  $p = 2.43e-2$ ). Examining  
496 gastric cancers in TCGA, we also observed enrichment of microbial load in MSI cases (9.7-fold  
497  $\pm 1.8$ ;  $p = 2.25e-4$ ) (**Figure 5C**), recapitulating a recent finding by Booth et al.<sup>92</sup> and suggesting a  
498 trend that was generalizable beyond colorectal cancer.

#### 499 ***Tumor mutation burden independently predicts tumor microbial load***

500 As MSI and polymerase-mutant tumors are characterized by an elevated mutation rate,  
501 enrichment of microbiota in these subtypes led us to hypothesize a broader relationship between  
502 microbial colonization and tumor mutation burden. Strikingly, microbial load in colorectal cancers  
503 was indeed correlated with tumor mutation burden within both MSS ( $R = 0.132$ ;  $p = 2.36e-7$ ) and  
504 MSI cases ( $R = 0.126$ ;  $p = 1.64e-2$ ) (**Figure 5E**). This association persisted after adjusting for the  
505 effect of genomic subtype on tumor mutation burden and the effect of tumor site on microbial load  
506 ( $R = 0.081$ ;  $p = 2.34e-3$ ) (**Figure 5F**). Quantitatively, each 10-fold increase in tumor mutation  
507 burden corresponded to an average 4.7-fold increase in microbial abundance ( $\pm 1.4$  SE) among  
508 MSS tumors and a 3.6-fold increase ( $\pm 1.2$  SE) in MSI tumors, suggesting a “dose-dependent”  
509 relationship.

510 We next asked whether the previously observed enrichment of microbiota in hypermutated  
511 subtypes could be explained entirely by tumor mutation burden rather than by subtype-specific  
512 biology. To test this, we fit a linear model predicting microbial abundance from genomic subtype,  
513 adjusting for tumor mutation burden and tumor site. When tumor mutation burden was included  
514 in this model, we found that MSI ( $p = 0.40$ ) and polymerase-mutant ( $p = 0.75$ ) subtypes were no  
515 longer predictive of microbial load, even as tumor site remained predictive ( $p = 2.3e-8$ ) (**Figure**  
516 **5G**). Moreover, a model incorporating tumor mutation burden, tumor site, and subtype performed  
517 no better than a model including only tumor mutation burden and tumor site ( $p = 0.67$ ), signifying  
518 that aside from tumor mutation burden, subtype-intrinsic features could not explain additional  
519 variation in microbial abundance.

520 Having observed microbial enrichment in MSI gastric cancers, we wondered if the  
521 relationship between tumor mutation burden and microbial abundance generalized beyond  
522 colorectal cancer. Indeed, oropharyngeal cancers demonstrated a significant positive correlation  
523 between tumor mutation burden and microbial load in both 100kGP ( $R = 0.145$ ;  $p = 2.4e-2$ ) and  
524 TCGA data ( $R = 0.154$ ;  $p = 4.3e-2$ ) (**Figure S5G-H**). The same pattern was not observed in  
525 100kGP gastroesophageal cancers, but microbial load in TCGA gastric cancers strongly  
526 correlated with tumor mutation burden ( $R = 0.251$ ;  $p = 5.2e-4$ ) (**Figure S5I-J**). Together, these  
527 findings implicate tumor mutation burden as a key predictor of microbial colonization across  
528 subtypes, cohorts, and cancer types.

### 529 ***Microbial composition varies independently by anatomic site and tumor mutation burden***

530 Having established that microbial load correlates with both anatomic site and tumor  
531 mutation burden, we next asked whether these features were also associated with differences in  
532 microbial composition. To delineate the effects of these two features, we analyzed associations  
533 between species-level relative abundances and tumor mutation burden (high vs. low) within each  
534 anatomic site. We then performed the complementary analysis, evaluating associations between  
535 species composition and anatomic site (proximal vs. distal), stratifying by tumor mutation burden  
536 (**Figure 5H, Table S8**).

537 This analysis revealed several species that were significantly associated with tumor  
538 mutation burden in colorectal cancer, including *Fusobacterium nucleatum* ( $q = 4.8e-11$ ) and  
539 *Bacteroides fragilis* ( $q = 3.3e-5$ ), which have previously been linked to MSI<sup>129,130,134</sup>. Other  
540 *Fusobacterium* species were also enriched in tumors with high mutation burden, including *F.*  
541 *vincentii* ( $q = 2.7e-8$ ), *F. pseudoperiodonticum* ( $q = 1.5e-5$ ), *F. polymorphum* ( $q = 9.5e-4$ ), and *F.*  
542 *animalis* ( $q = 8.6e-12$ ), a subclade of which has been shown to exhibit tropism for colorectal  
543 tumors<sup>135</sup>. Species not previously tied to hypermutated subtypes were also enriched, including  
544 *Solobacterium moorei* ( $q = 1.4e-8$ ), *Selenomonas sputigena* ( $q = 2.7e-8$ ), *Parvimonas micra* ( $q =$   
545  $4.0e-8$ ), *Hungatella hathewayi* ( $q = 2.0e-6$ ), and *Dialister pneumosintes* ( $q = 5.4e-5$ ). Notably,  
546 many species enriched in hypermutated tumors were oral-typical and/or periodontal disease-  
547 associated, suggesting a relationship between oral-typical microbiota and mutation burden in  
548 colorectal cancer.

549 Conversely, tumors with low mutation burden appeared to be enriched for gut commensals  
550 with reportedly anti-inflammatory properties<sup>136–138</sup>, including *Phocaeicola vulgatus* ( $q = 7.5e-6$ ),  
551 *Faecalibacterium prausnitzii* ( $q = 1.7e-6$ ), and *Bacteroides uniformis* ( $q = 4.9e-6$ ). We additionally

552 observed enrichment of *Escherichia coli* ( $q = 3.3e-4$ ), *Veillonella atypica* ( $q = 2.1e-4$ ), and  
553 *Veillonella parvula* ( $q = 1.4e-3$ ) in tumors with low mutation burden.

554 Meanwhile, another group of species were mutation burden-independent and instead  
555 varied by anatomic site. For example, *Clostridium perfringens* ( $q = 5.8e-29$ ) and *Roseburia*  
556 *intestinalis* ( $q = 1.1e-16$ ) were significantly associated with proximal colorectal tumors, while  
557 *Porphyromonas asaccharolytica* ( $q = 1.4e-15$ ), *Ruthenibacterium lactatiformans* ( $q = 3.4e-16$ ),  
558 and *Akkermansia muciniphila* ( $q = 1.2e-13$ ) were associated with distal tumors.

559 Some species varied both by mutation burden and by anatomic site, suggesting an  
560 interaction between these features to influence microbial composition. *S. sputigena* and *F.*  
561 *animalis* were more abundant in hypermutated tumors of the proximal colon, while *P. micra* was  
562 proportionally greater in hypermutated tumors of the distal colon (**Figure 5I**). Conversely, *V.*  
563 *parvula* and *V. atypica* were enriched in tumors in the proximal colon with low mutation burden,  
564 while *E. coli* was more associated with distal tumors with low mutation burden (**Figure 5J**).

565 Several species associated with mutation burden additionally appeared to accumulate  
566 with advancing tumor stage (**Figure 5K**). While *F. animalis* appeared to be enriched in  
567 hypermutated cases as early as stage I, others like *S. sputigena*, *P. micra*, *P. intermedia*, and *S.*  
568 *moorei* were more enriched at later stages, potentially indicating “early” and “late” colonizers. The  
569 relative increase in these species by stage was accompanied by a corresponding decrease in  
570 commensals such as *P. vulgatus*, *F. prausnitzii*, and *B. uniformis* (**Figure 5L**).

571 Notably, many of the species associated with tumor mutation burden, anatomic location,  
572 or cancer stage had overlap with reproducible fecal biomarkers of colorectal cancer recently  
573 reported by Piccinno et al.<sup>127</sup>. For example, Piccinno et al.’s observation that fecal detection of *V.*  
574 *parvula* and *V. atypica* predicted the presence of proximal colorectal tumors mirrors our  
575 observation that these species are enriched in proximal tumors themselves. This result supports  
576 the idea that these species may be contributing to fecal signatures by shedding from the tumor  
577 into stool, linking tumor microbiome composition to fecal biomarker candidates.

#### 578 ***Akkermansia muciniphila* is depleted in early-onset tumors of the distal colon**

579 Early-onset colorectal cancer is a rising public health concern, with incidence in individuals  
580 under 50 years of age doubling in many countries since the 1990s<sup>139–141</sup>. A recent study reported  
581 enrichment of SBS88 — a mutational signature caused by colibactin, a genotoxin produced by  
582 pks+ *E. coli*<sup>49</sup> — in early-onset colorectal cancer<sup>142</sup>, raising the possibility that gut microbiota may

583 be predictive of colorectal cancer onset in younger adults. Noting that patient age was a predictor  
584 of microbial composition in colorectal tumors (**Figure 5A**), we sought to examine microbial  
585 correlates with age of onset.

586 We compared early- (<50y) and average-onset ( $\geq 50$ ) colorectal cancers, stratified by  
587 anatomic site and restricting to microsatellite-stable (MSS) cases. Clinically, early-onset colorectal  
588 cancer is typically MSS and occurs most frequently in the distal colon<sup>143</sup>. In proximal colorectal  
589 cancers (**Figure S5K**), we observed variation of *Streptococcus sanguinis* ( $p = 8.33e-3$ ),  
590 *Thomasclavelia ramosa* ( $p = 7.01e-3$ ), and *Vescimonas fastidiosa* ( $p = 1.08e-2$ ) by age of onset.  
591 These species are generally considered commensals, but *S. sanguinis* is a common cause of  
592 endocarditis<sup>144</sup>.

593 We next compared tumor microbial composition between early- and average-onset cases  
594 arising in the distal colon, where early-onset cancers are more common (**Figure S5L**).  
595 *Akkermansia muciniphila* was significantly depleted in early-onset cases (FC = 12.5;  $p = 3.4e-3$ )  
596 (**Figure S5M**), as was *Hungatella hathewayi* to a lesser extent (FC = 2.4,  $p = 3.3e-3$ ). *H. hathewayi*  
597 is a widely carried gut commensal noted for its ability to efficiently degrade glycosaminoglycans<sup>149</sup>,  
598 while *A. muciniphila* is a mucin-degrading commensal that has been linked to intestinal barrier  
599 function, metabolic health, and reduced inflammation<sup>150</sup>, and has shown consistent enrichment in  
600 the gut microbiomes of centenarians<sup>151</sup>.

601 To provide independent validation of this association, we examined microbial profiles from  
602 TCMA/TCGA (**Figure S5N**). Consistent with 100kGP, *A. muciniphila* showed no difference in  
603 COAD tumors. However, in READ tumors, *A. muciniphila* was similarly depleted in early-onset  
604 cases (FC = 42.8), mirroring our observation in distal CRCs from the 100kGP cohort, though this  
605 did not reach statistical significance ( $p = 0.158$ ), likely reflecting the low number of early-onset  
606 rectal cancers in the TCGA cohort ( $n = 10$ ).

607 Together, these data indicate that *A. muciniphila* may be consistently depleted in early-  
608 onset CRCs, raising the possibility that loss of mucin-associated commensals might contribute to  
609 and/or reflect patterns of barrier dysfunction and broader mucosal perturbations in younger  
610 patients.

611

612

## 613 **Limitations of the study**

614           Although we observed taxa with a range of cell architectures, our microbial abundance  
615 estimates may be underestimates due to incomplete lysis of microbial cell walls during nucleic  
616 acid extraction. Furthermore, our analysis does not disambiguate “mixed-evidence” species,  
617 whose reads may reflect a mixture of contamination and true signal. This may lead to additional  
618 underestimates of microbial signal in skin cancer, for example, as some of the contaminant taxa  
619 we detected were common skin commensals.

620           Finally, while the tumor genomes analyzed in this study were sequenced to high coverage  
621 (~100X), detection of microbiota in these samples is nevertheless bounded by assay-specific  
622 limits of detection. Normalized to sequencing depth, we determined our species-level detection  
623 threshold to be approximately 0.1 reads per million (RPM). For bacteria, this corresponds to  
624 approximately one microbial cell per ~7,000 tumor cells or per ~50 ng of tumor DNA, which is  
625 comparable to prior thresholds used for estimating tumor microbiome abundance using  
626 quantitative amplification of bacterial 16S rRNA<sup>39</sup>. Nevertheless, the key conclusions of this study  
627 will need to be validated using orthogonal assays to determine the extent to which these patterns  
628 generalize beyond bulk tumor sequencing data.

629

## 630 **Discussion**

631           This study establishes the biodiversity and biogeography of the cancer microbiome by  
632 addressing persistent methodological challenges in tumor microbiome profiling. Robust,  
633 quantitative characterization of tumor-associated microbial signals has historically been elusive  
634 due to incomplete host subtraction, batch effects, and widespread contamination that can mask  
635 real tissue-associated microbial signal in low-biomass samples.

636           To address these challenges, we developed and benchmarked PathSeq-T2T, which uses  
637 the complete human reference genome for host filtering and leverages multiple metagenomic  
638 profilers for quantitative estimation and cross-validation of microbial signals. Experimental  
639 benchmarks using human-microbe DNA mixtures revealed empirical limits of detection and  
640 identified contamination patterns, enabling us to better interpret microbial signals in large-scale  
641 tumor sequencing datasets. We subsequently applied this pipeline to 16,369 high-depth tumor  
642 whole genomes from the UK 100,000 Genomes Project (100kGP), which enabled the largest and  
643 most comprehensive pan-cancer analysis of tumor-associated microbial communities to date.

644           Although contamination was present in this dataset, we were able to identify and remove  
645 microbial contaminants using the principle of equiprevalence, which posits that contaminants will  
646 generally be present at similar rates across distinct sample types. While this approach does not  
647 replace gold-standard negative controls (which are seldom available for large-scale sequencing  
648 initiatives), it has previously proven effective in decontaminating other large-scale sequencing  
649 datasets<sup>57</sup>.

650           Our analysis of the 100kGP cohort affirmed that most solid cancer types lack a tumor-  
651 associated microbiome that is detectable in bulk sequencing data, consistent with and extending  
652 other recent studies of tumor genomes from 100kGP<sup>61</sup> and TCGA<sup>57,62</sup>. A clear exception was  
653 observed in orodigestive cancers — including from oropharyngeal, esophageal, gastric, and  
654 colorectal sites — each of which harbored polymicrobial communities that displayed site-specific  
655 abundance, biogeography, and biodiversity consistent with microbial surveys of healthy  
656 individuals<sup>1–11</sup>. These findings support a model in which microbial colonization of tumors occurs  
657 predominantly at mucosal and epithelial barrier surfaces that are already permissive to  
658 commensal or pathogenic colonization.

659           The conclusions of these collective studies differ from earlier reports of distinct microbial  
660 communities present across all<sup>55,56,60</sup> or many<sup>39</sup> cancer types. Such discrepancies may in some  
661 cases be the result of methodological challenges which have been discussed extensively  
662 elsewhere<sup>54,56,58,61,62</sup>, but in other cases may be due to differences in analytic approach (e.g. bulk  
663 PCR-free sequencing versus PCR-amplified or spatially-resolved assays), thresholds of  
664 detection, and/or strategies for handling contamination (See *Limitations of the study*). These  
665 debates underscore the need for standardized contamination control and transparent  
666 reporting<sup>82,152</sup>.

667           Notwithstanding, our analysis revealed multi-kingdom microbial communities in  
668 orodigestive cancers, including fungi, archaea, viruses, and in rare cases, protozoan parasites.  
669 These results extend prior knowledge of microbial involvement in cancer beyond bacteria and  
670 viruses, reinforce emerging evidence for fungal presence in orodigestive cancers<sup>59,60</sup>, and further  
671 demonstrate that tumor-associated microbial communities can also accommodate archaea and  
672 protozoa. Future work will be needed to determine what role, if any, species from these newly  
673 implicated domains play in cancer development and progression.

674           Extending prior observations of microbial enrichment in hypermutated genomic subtypes  
675 <sup>129,130,134</sup>, we additionally found that tumor mutation burden is a predictor of microbial abundance,  
676 diversity, and composition in orodigestive tumors. Although bacteria (e.g. *pks+* *E. coli*) have been  
677 linked to tumor mutational signatures<sup>47–50</sup>, elevated TMB is unlikely to be the result of increased  
678 microbial colonization, as most hypermutated tumors arise in the setting of genome alterations in  
679 well-defined molecular pathways, notably DNA repair and DNA polymerase mutations. Instead,  
680 the association of high microbial abundance with tumor mutation burden may be mediated by  
681 tumor-immune interactions associated with increased neoantigen load. These include altered  
682 inflammatory pathways and expression of immune checkpoints<sup>153–157</sup>, leading to local  
683 immunosuppression or exhaustion that creates a permissive niche for microbial expansion. In an  
684 alternative but not mutually exclusive scenario, tumor-associated microbiota may themselves  
685 contribute to suppression of anti-tumor immunity, as has been reported for *Fusobacterium*  
686 interactions with TIGIT<sup>158</sup>, providing benefit to neoantigen-rich tumors. Additional work will be  
687 required to further validate the microbial abundance and mutation burden correlation and discern  
688 models for its selective benefit for tumors and/or microbes.

689           We also observed depletion of *Akkermansia muciniphila* in early-onset colorectal cancers  
690 of the distal colon. This mucosa-associated species has attracted interest as a potential probiotic  
691 agent, due to reported links to mucosal barrier function<sup>150</sup>, longevity<sup>151</sup>, and response to  
692 immunotherapy<sup>29,33</sup>. Accordingly, *A. muciniphila* warrants further study in the context of colorectal  
693 cancer and age of onset, including the role of the protective mucin layer in preventing exposure  
694 to luminal genotoxins. Confirmation in larger cohorts of younger colorectal cancer patients is  
695 necessary to assess the potential relevance of *A. muciniphila* to the increasing rate of colorectal  
696 cancer in young adults.

697           In conclusion, our analysis helps to resolve persistent uncertainties stemming from  
698 conflicting reports on the scope and distribution of microbial colonization across cancer sites,  
699 while identifying phenotypic and genotypic determinants of tumor microbiome community  
700 structure. We hope that the methodology, dataset, and insights presented here provide a  
701 foundation for future functional and mechanistic investigations of tumor-associated microbes and  
702 help nominate cancers most likely to benefit from prospective biomarker discovery and  
703 microbiome-focused therapeutic development.

704

## 705 **Resource availability**

### 706 ***Lead contact***

707 Further requests for data should be directed to the lead contact, Dr. Matthew Meyerson  
708 (matthew\_meyerson@dfci.harvard.edu).

### 709 ***Materials availability***

710 This study did not generate new, unique reagents.

### 711 ***Data and code availability***

712 The 100kGP microbiome profiles generated by PathSeq-T2T are accessible from within  
713 the Genomics England research environment, which requires an application to access, found  
714 here: <https://www.genomicsengland.co.uk/join-us>. Microbial profiles and sample metadata  
715 needed to reproduce the analyses described in this study can be found at the following path:  
716 /re\_gecip/shared\_allGeCIPs/adohlman/gel\_manuscript. Due to data privacy considerations,  
717 Genomics England does not allow the export of sample- and patient-level data tables from the  
718 research environment.

719 The code used to perform the analyses and generate the figures associated with this  
720 manuscript are available in Jupyter Notebook format, which are accessible from within the  
721 Genomics England research environment. These can be found in the following directory:  
722 /re\_gecip/shared\_allGeCIPs/adohlman/gel\_manuscript/analysis. The PathSeq-T2T pipeline is  
723 available on GitHub (<https://github.com/abdohlman/pathseq-t2t/>).

## 724 **Acknowledgments**

725 This work was supported by the Damon Runyon Cancer Research Foundation (grant no.  
726 DRG-2504-23 to A.B.D.); Cancer Grand Challenges OPTIMISTIC, funded by Cancer Research  
727 UK (grant no. A27140 to A.B.D., H.W., K.J., G.S., P.N., P.Q., C.H., and M.M.); the NIHR Leeds  
728 Biomedical Research Centre (grant no. NIHR203331 to H.W. and P.Q.); and Yorkshire Cancer  
729 Research (grant no. L386 to P.Q.). P.Q. is an NIHR Senior Investigator. N.S. is supported by the  
730 European Union's Horizon 2020 programme ONCOBIOME (grant no. 825410); the European  
731 Research Council (grant nos. ERC-StG MetaPG-716575 and ERC-CoG microTOUCH-  
732 101045015); the U.S. National Cancer Institute (grant no. 1U01 CA230551); the Premio  
733 Internazionale Lombardia e Ricerca 2019; and MIUR PRIN 2017 (grant no. 2017J3E2W2). N.S.

734 and C.H. are additionally supported through the Cancer Grand Challenges PROSPECT team  
735 (Cancer Research UK grant nos. CGCATF-2023/100036 and CGCATF-2023/100041; National  
736 Cancer Institute grant nos. OT2CA297680 and 1OT2CA297205-01; the Bowelbabe Fund for  
737 Cancer Research UK; and Institut National du Cancer). C.H. is further supported by Prescient  
738 Metabionics. The views expressed are those of the authors and not necessarily those of the NHS,  
739 the NIHR, or the Department of Health and Social Care.

740 This research was made possible through access to data in the National Genomic  
741 Research Library, which is managed by Genomics England Limited (a wholly owned company of  
742 the Department of Health and Social Care). The National Genomic Research Library holds data  
743 provided by patients and collected by the NHS as part of their care. The National Genomic  
744 Research Library is funded by the National Institute for Health Research and NHS England. The  
745 Wellcome Trust, Cancer Research UK and the Medical Research Council have also funded  
746 research infrastructure. We also thank Carrie Cibulskis and Matthew Brown for their contributions,  
747 as well as all of the research participants and their families for making this work possible.

## 748 **Author contributions**

749 A.B.D. conceived of the study, performed the formal data analysis, created the figures, wrote and  
750 edited the manuscript. R.M. performed metagenome assemblies and genome analyses and  
751 contributed to data interpretation. H.W. performed preliminary analyses of colorectal and gastric  
752 cancer microbiomes in the 100kGP cohort and helped curate metadata for these cancer types.  
753 K.J. performed the *in vitro* microbe-host mixture experiments. A.S. generated the *in silico* microbe-  
754 host mixture data. I.L. helped to reanalyze TCGA data. G.P. provided and advised on oral and  
755 fecal biomarker annotations used in the study. P.N. and G.S. performed external validation  
756 analyses and assisted in interpretation of the results. A-R.Y. provided analysis and interpretation  
757 of unclassified sequences. P.Q. advised on preliminary analyses of colorectal and gastric cancer  
758 microbiomes. C.H. and N.S. advised on the analytical approach and data interpretation. M.M.  
759 conceived of the study, supervised the analyses, reviewed and edited the manuscript.

## 760 **Declaration of interests**

761 A.B.D. is an inventor on US patent application no. PCT/US2023/029634, submitted by Cornell  
762 University and Duke University related to the use of fungi for classifying and treating  
763 gastrointestinal and lung tumors. M.M. receives research support from Bayer, holds equity in and  
764 consults for Bayer, Delve Bio, Isabl, and Karyoverse, and is an inventor on patents licensed to

765 Bayer and LabCorp, all outside of the current work. M.M. is also an inventor on an issued patent  
766 and several patent applications regarding *Fusobacterium* in colorectal cancer, which are not  
767 licensed. N.S. is a founder and shareholder of PreBiomics Srl and is on the scientific advisory  
768 board of ZOE Ltd and received consultancy fees from them. C.H. is a member of the Seres  
769 Therapeutics and Empress Therapeutics scientific advisory boards.

## 770 **Declaration of generative AI and AI-assisted technologies**

771 During the preparation of this work, authors used ChatGPT to assist in writing code used  
772 for performing data analyses and visualizations. After the manuscript was written, ChatGPT was  
773 used to assist in editing, such as checking for spelling and grammatical errors. After using this  
774 tool, the authors reviewed and edited the content as needed and take full responsibility for the  
775 content of the publication.

## 776 **Supplemental information**

777 **Table S1.** Pre-defined species-level read counts for *in silico* microbe-human dilution series.

778 **Table S2.** *In silico* read counts after each PathSeq-T2T filtering step.

779 **Table S3.** *In vitro* read counts after each PathSeq-T2T filtering step.

780 **Table S4.** Summary of samples in the Genomics England cohort.

781 **Table S5.** Species equiprevalence scores and contamination classifications.

782 **Table S6.** Summary of genome assembly results from orodigestive cancer types.

783 **Table S7.** Site-specific species associations in oropharyngeal, gastroesophageal, and colorectal  
784 cancers.

785 **Table S8.** Joint analysis of CRC tumor anatomic site and TMB subtype.

786

787

788

## 789 **Figure titles and legends**

### 790 **Figure 1. A host subtraction and classification pipeline for identifying microorganisms in** 791 **low-biomass human tissues**

792 (A) Retention rate of *in silico* human sequences in the human-only condition after each host  
793 filter, expressed as reads per million primary reads.

794 (B) Retention rate of *in silico* microbial reads after each host filter, across microbe-host  
795 mixture conditions ( $n = 9$ ) ranging from a host-microbe ratio of 1:1 (*red*) to 1:10<sup>8</sup> (*blue*).

796 (C) Total read retention rate after each successive PathSeq-T2T filtering step for each  
797 condition in the *in silico* dilution series. Horizontal gray lines indicate the starting microbial  
798 proportion in each condition.

799 (D) Relative abundance of human (*blue*), on-target microbial (*green*), off-target microbial  
800 (*orange*), and not uniquely classified (*gray*) reads after each host filtering step across *in*  
801 *silico* dilution conditions. Reads are “on-target” if assigned to species/genera/families  
802 present in microbial community standard, and “off-target” otherwise.

803 (E) Total read retention rate after each successive PathSeq-T2T filtering step for each  
804 condition in the *in vitro* dilution series ( $n = 3$  replicates each).

805 (F) Relative abundance of human (*blue*), on-target microbial (*green*), off-target microbial  
806 (*orange*), and not uniquely classified (*gray*) reads after each host filtering step across *in*  
807 *vitro* dilution conditions.

808 (G) Abundance in reads per million (RPM) of on-target microbial sequences in each condition  
809 after PathSeq-T2T filtering of *in vitro* microbe-host DNA mixtures.

810 (H) Abundance of off-target (putative contaminants) *Cutibacterium acnes*, *Staphylococcus*  
811 *epidermidis*, and *Malassezia restricta* across *in vitro* dilution conditions.

812 (I) Relative abundance of species detected in negative controls ( $n = 3$ ) with no human or  
813 microbial DNA added, processed in parallel with *in vitro* conditions dilutions.

814 (J) Percentage of off-target (contaminant) reads across successive dilution conditions.

815 *Microbial counts are from Kraken2 unless stated otherwise.*

816

### 817 **Figure S1. A host subtraction pipeline for identifying microorganisms in low-biomass** 818 **human tissues (Related to Figure 1)**

- 819 **(A)** Schematic of the PathSeq-T2T pipeline used to study microbial signals in low-biomass  
820 tissues. Tumor whole-genome sequencing (WGS; delivered as hg38-aligned BAMs) were  
821 (1) filtered to select unmapped reads; (2) processed using quality, complexity, and host k-  
822 mer filtering; and then (3) filtered using T2T-CHM13. Last (4) reads passing these filtering  
823 steps are classified using Kraken2<sup>77</sup>, MetaPhlan4<sup>78</sup>, and Sylph<sup>79</sup>. In parallel, raw hg38-  
824 unmapped reads were assembled for microbial genome reconstruction using MEGAHIT<sup>159</sup>  
825 (See **Table S6**).
- 826 **(B)** Retention rate of *in silico* human sequences in the across microbe-host mixture conditions  
827 (1:1 to 1:10<sup>8</sup>) after each host filter.
- 828 **(C)** Breakdown of *in silico* microbial reads removed at each PathSeq-T2T filtering step  
829 (microbe-only condition), expressed as a percentage of total removed microbial reads.
- 830 **(D)** Misclassification rate of reads removed or retained by PathSeq-T2T, expressed as  
831 misclassifications per million microbial reads.
- 832 **(E)** Total read retention rate after each successive PathSeq-T2T filtering step for colorectal  
833 cancer cell lines annotated as microsatellite-unstable (MSI;  $n = 8$ ) or microsatellite-stable  
834 (MSS;  $n = 12$ ).
- 835 **(F)** Total read retention rate after each successive PathSeq-T2T filtering step for colorectal  
836 cell lines prepared using PCR ( $n = 10$ ) or PCR-free ( $n = 10$ ) methods.
- 837 **(G)** Relative abundance of microbial genera identified in colorectal cancer cell lines by  
838 PathSeq-T2T, representing likely contaminants.

839 \*  $p < 0.05$ , \*\*  $p < 0.01$ , \*\*\*  $p < 0.001$ , \*\*\*\*  $p < 0.0001$  by Wilcoxon's signed-rank test. Microbial  
840 counts are from Kraken2 unless stated otherwise.

841

## 842 **Figure 2. Identification and removal of contamination from the 100kGP cohort**

- 843 **(A)** Median microbial reads per million (RPM) retained in saliva (blue;  $n = 490$ ) and blood (red;  
844  $n = 14,127$ ) after each PathSeq-T2T filtering step. Error bars indicate interquartile range.
- 845 **(B)** Median RPM retained in colorectal cancer (CRC; brown;  $n = 2,482$ ) and glioblastoma  
846 (GBM; green;  $n = 295$ ) tumors after each PathSeq-T2T filtering step. Error bars indicate  
847 interquartile range.
- 848 **(C)** Relative abundance of genera detected in saliva (top) and blood (bottom). Samples are  
849 ordered by germline sequence delivery date.

- 850 (D) Relative abundance of genera detected in CRC (top) and GBM (bottom). Samples are  
851 ordered by germline sequence delivery date.
- 852 (E) Prevalence (i.e. percentage of samples) of species detected in blood samples (x-axis) and  
853 saliva samples (y-axis). Saliva harbors saliva-specific species (blue), while equiprevalent  
854 species (red) are found at similar frequencies in saliva and blood.
- 855 (F) Prevalence of species detected in GBM samples (x-axis) and CRC samples (y-axis). CRC  
856 samples harbor CRC-specific species (blue), while equiprevalent species (red) are found  
857 at similar frequencies in CRC and GBM.
- 858 (G) Bar plot showing ten most blood-prevalent species, among those equiprevalent in  
859 saliva/blood.
- 860 (H) Bar plot showing ten most GBM-prevalent species, among those equiprevalent in  
861 CRC/GBM.
- 862 (I) Histogram showing pan-cancer equiprevalence (PCE) scores of species detected in  
863 100kGP tumor samples. Species equiprevalent in blood/saliva or CRC/GBM are labeled  
864 in red, other species in blue. Species with PCE scores greater than 0.7 (gray shaded area)  
865 were classified as contaminants. Distribution includes species with >5% prevalence in at  
866 least one cancer type.
- 867 (J) Prevalence of selected species of across cancer types, labeled with PCE scores. Known  
868 cancer-associated species (*H. pylori*, *F. nucleatum*, *HPV*) had low PCE scores, while  
869 contaminants (*C. acnes*, *S. epidermidis*, *M. restricta*) had high PCE scores.

870 *Microbial counts are from Kraken2 unless stated otherwise.*

871

872 **Figure S2. Identification and removal of contamination from the 100kGP cohort (Related to**  
873 **Figure 2)**

- 874 (A) Prevalence (i.e. percentage of samples) of species detected in blood (x-axis) and saliva  
875 samples (y-axis) using MetaPhlan4<sup>78</sup>. Saliva samples harbor saliva-specific species  
876 (blue), while equiprevalent species (red) are found at equal frequency in saliva and blood.
- 877 (B) Prevalence of species detected in blood (x-axis) and saliva samples (y-axis) using Sylph<sup>79</sup>.
- 878 (C) Prevalence of species detected in GBM (x-axis) and CRC samples (y-axis) using  
879 MetaPhlan4<sup>78</sup> (left) and Sylph<sup>79</sup> (right). CRC samples harbor CRC-specific species (blue),  
880 while equiprevalent species (red) are found at equal frequency in CRC/GBM.
- 881 (D) Prevalence of species detected in GBM (x-axis) and CRC samples (y-axis) using Sylph<sup>79</sup>.

- 882 (E) Venn diagram depicting the degree of overlap between species that were equiprevalent  
883 in saliva and blood and equiprevalent in CRC and GBM.
- 884 (F) Bar plot showing the ten most saliva-prevalent species.
- 885 (G) Bar plot showing the ten most CRC-prevalent species.
- 886 (H) Histogram of PCE scores of microbial species detected in 100kGP tumor samples by  
887 MetaPhlan4<sup>78</sup> (*left*) and Sylph<sup>79</sup> (*right*). Species that were equiprevalent in blood and saliva  
888 or in CRC and GBM (*red*) typically had high PCE scores, while other species (*blue*) did  
889 not. Species with PCE > 0.7 (*gray shaded area*) were classified as contaminants.  
890 Distribution includes species detected at >5% prevalence in at least one cancer type.
- 891 (I) Histogram of PCE scores of species detected in 100kGP tumor samples by Sylph<sup>79</sup>.
- 892 (J) Prevalence of species with the highest PCE scores in the 100kGP cohort, across cancer  
893 types. Most of these species are environmental and do not typically colonize humans.
- 894 (K) Within-species median average nucleotide identity (ANI) of microbial genomes obtained  
895 through metagenomic assembly, sorted by median ANI. Genomes from species with high  
896 PCE scores (*red*) have strain-level sequence similarity (>99.99%), suggesting a common  
897 contaminant source.
- 898 (L) Volcano plot showing enrichment of equiprevalent species (*red*) in germline samples that  
899 underwent PCR-based ( $\log_2OR > 0$ ) versus PCR-free preparation ( $\log_2OR < 0$ ). Statistical  
900 significance was determined using Fisher's exact test (*y-axis*),
- 901 (M) Volcano plot showing enrichment of equiprevalent species (*red*) in tumor samples that  
902 underwent PCR-based ( $\log_2OR > 0$ ) versus PCR-free preparation ( $\log_2OR < 0$ ). Statistical  
903 significance was determined using Fisher's exact test (*y-axis*),
- 904 (N) Volcano plot showing enrichment of equiprevalent species (*red*) in tumor samples that  
905 were stored as FFPE ( $\log_2OR > 0$ ) versus fresh-frozen ( $\log_2OR < 0$ ).

906 *Microbial counts are from Kraken2 unless stated otherwise.*

907

908 **Figure 3. Tumors of the orodigestive tract harbor a microbiome but most other cancer**  
909 **types do not**

- 910 (A) Cumulative distribution of Kraken2<sup>77</sup> sample-level microbial abundance by cancer site  
911 after decontamination. Few non-orodigestive cancers exceeded background (*shaded gray*  
912 *area*; RPM < 0.1).

- 913 (B) Cumulative distribution of MetaPhlAn4<sup>78</sup> sample-level microbial abundance by cancer site  
914 after decontamination.
- 915 (C) Cumulative distribution of Sylph<sup>79</sup> sample-level microbial abundance by cancer site after  
916 decontamination.
- 917 (D) Species prevalence in each cancer type (*y-axis*) compared to their prevalence in brain (*x-*  
918 *axis*) after decontamination. Few non-oro-digestive cancers exhibit cancer-site specific  
919 microbial signals.
- 920 (E) Letter-value (boxen) plots showing number of species detected per sample (RPM > 0.1)  
921 after decontamination. Few non-oro-digestive cancers consistently harbor more than one  
922 species.
- 923 (F) Average family-level relative abundance of microbe-high samples (>10 RPM) from each  
924 cancer site after decontamination. Only a small fraction of samples from non-oro-digestive  
925 cancers met this criterion and many were dominated by common *Herpesviridae*.

926 *Microbial counts are from Kraken2 unless stated otherwise.*

927

928 **Figure S3. Tumors of the orodigestive tract harbor a microbiome but most other cancer**  
929 **types do not (Related to Figure 3)**

- 930 (A) Cumulative distribution of Kraken2<sup>77</sup> sample-level microbial abundance by cancer site,  
931 prior to decontamination. Nearly all samples and cancer types exceed the background  
932 threshold (shaded gray area, RPM < 0.1).
- 933 (B) Cumulative distribution of MetaPhlAn4<sup>78</sup> sample-level microbial abundance by cancer site,  
934 prior to decontamination.
- 935 (C) Cumulative distribution of Sylph<sup>78</sup> sample-level microbial abundance by cancer site, prior  
936 to decontamination.
- 937 (D) Cumulative distribution of Kraken2<sup>77</sup> sample-level contamination abundance by cancer  
938 site. Cancer types mostly harbor similar rates of contamination.
- 939 (E) Cumulative distribution of MetaPhlAn4<sup>78</sup> sample-level contamination abundance by  
940 cancer site.
- 941 (F) Cumulative distribution of Sylph<sup>78</sup> sample-level contamination abundance by cancer site.
- 942 (G) Spearman correlation (Rho) between microbial abundance estimates from Kraken2,  
943 MetaPhlAn4, and Sylph across cancer types.

- 944 (H) Prevalence of species in each cancer type (*y-axis*) compared to their prevalence in brain  
945 (*x-axis*), prior to decontamination. Species classified as contaminants are shown in red.
- 946 (I) Prevalence of the *Papillomaviridae* species HPV-16 and HPV-18 across cancer sites.
- 947 (J) Prevalence of the *Polyomaviridae* species MCPyV across cancer sites.
- 948 (K) Prevalence of the *Herpesviridae* species HSV-1, EBV, CMV, and HHV-7 across cancer  
949 sites.
- 950 (L) Prevalence of the *Herpesviridae* species HHV-6A and HHV-6B across cancer sites.

951 \*  $p < 0.05$ , \*\*  $p < 0.01$ , \*\*\*  $p < 0.001$ , \*\*\*\*  $p < 0.0001$  by Spearman's rank correlation test. Microbial  
952 counts are from Kraken2 unless stated otherwise.

953

954 **Figure 4. Biogeography and biodiversity of oropharyngeal, gastroesophageal, and**  
955 **colorectal cancers (Related to Figure 4)**

- 956 (A) Microbial abundance (RPM) by across anatomical sites in oropharyngeal (*blue*),  
957 gastroesophageal (*orange*), and colorectal (*green*) cancers.
- 958 (B) Principal coordinates analysis (PCoA) of samples from oropharyngeal, gastroesophageal,  
959 and colorectal cancers.
- 960 (C) Radial plot showing species-level barycentric centroids calculated from cancer site-  
961 specific abundances ( $\log_{10}$  RPM). Abundances in oropharyngeal, gastroesophageal, and  
962 colorectal cancers were used to define a triangle in site-space, the centroid of which was  
963 mapped to polar coordinates (*angle*: site bias; *radius*: magnitude of enrichment).
- 964 (D) Relative abundance of microbial genera in oropharyngeal (*top left*), gastroesophageal (*top*  
965 *right*), and colorectal cancers (*bottom*).
- 966 (E) Prevalence of bacterial species (*P. gingivalis*, *H. pylori*, *F. animalis*, *B. fragilis*, *E. coli*, *S.*  
967 *anginosus*) across anatomical sites of the orodigestive tract.
- 968 (F) Prevalence of fungal species (*C. albicans*, *N. glabratus*, *S. cerevisiae*) across anatomical  
969 sites.
- 970 (G) Prevalence of the archaeon *M. smithii* across anatomical sites of the orodigestive tract.
- 971 (H) Prevalence of the protozoan genus *Trichomonas* across anatomical sites.
- 972 (I) Prevalence of viruses (HPV-16, HHV-7, EBV, HHV-6A, and HHV-6B) across anatomical  
973 sites.

974 Microbial counts are from Kraken2 unless stated otherwise.

975

976 **Figure S4. Biogeography and biodiversity of oropharyngeal, gastroesophageal, and**  
977 **colorectal cancers (Related to Figure 4)**

- 978 (A) Diversity of microbiota (Shannon index) across anatomical sites of the orodigestive tract.  
979 (B) Estimated bacterial cells per thousand human cells across anatomical sites.  
980 (C) Venn diagram showing numbers of overlapping species detected in oropharyngeal,  
981 gastroesophageal, and colorectal cancers, including species detected in >1% of samples  
982 from at least one site.  
983 (D) Prevalence of bacterial species (*T. socransii*, *C. gingivalis*, *L. fermentum*, *L. gasseri*, *P.*  
984 *histicola*, *P. micra*, *S. moorei*, *P. intermedia*, *V. parvula*, *S. sputigena*, and *A. muciniphila*)  
985 across anatomical sites of the orodigestive tract.  
986 (E) Prevalence of the archaea *M. sp. TLL-48-HuF1* and *Candidatus Methanomassiliicoccus*  
987 *intestinalis* across anatomical sites.

988 *Microbial counts are from Kraken2 unless stated otherwise.*

989

990 **Figure 5. The tumor microbiome varies by tumor site, mutation burden, and age of onset.**

- 991 (A) Bar plot showing the percentage of variance explained identifying phenotypic and genomic  
992 features of colorectal cancers associated with microbial composition (*p*-value by  
993 PERMANOVA).  
994 (B) Microbial abundance in reads per million (RPM) in microsatellite-stable (MSS; *gray*)  
995 microsatellite-unstable (MSI; *red*), polymerase-mutated (POL; *blue*) subtypes.  
996 (C) Microbial abundance in MSS (*gray*) and MSI (*red*) colorectal (COAD) and gastric (STAD)  
997 cancers from The Cancer Microbiome Atlas (TCMA)<sup>57</sup>.  
998 (D) Microbial diversity (Shannon index) in MSS (*gray*), MSI (*red*), and POL (*blue*) colorectal  
999 cancer subtypes.  
1000 (E) Pearson correlation between tumor mutation burden (TMB; *x-axis*) and microbial  
1001 abundance (RPM; *y-axis*) in MSS (*gray*), MSI (*red*), and POL (*blue*) subtypes. Microbial  
1002 load and TMB are correlated in MSS and MSI but not POL.  
1003 (F) Correlation between TMB, after regressing out the effect of genomic subtype (*x-axis*) and  
1004 microbial abundance, after regressing out the effect of anatomic site (*y-axis*).

- 1005 (G) Microbial abundance by tumor subtype (MSS/MSI/POL), after regressing out the influence  
1006 of TMB on microbial abundance. Tumor subtype is no longer predictive of microbial load.
- 1007 (H) Scatter plot showing joint associations between relative abundance of CRC-associated  
1008 species and tumor site (proximal vs. distal) and TMB status (TMB high vs. TMB low). Axes  
1009 show the composite Wilcoxon test statistic, calculated for tumor site, blocked by subtype  
1010 (x-axis) and TMB status, blocked by tumor site (y-axis). Species are colored according to  
1011 their association with tumor site (green), TMB status (red), or both (orange), based on a  
1012 combined *p*-values (Fisher's method). Species associated with TMB status are also  
1013 reproducible fecal biomarkers for CRC<sup>127</sup> (triangles) and frequently include oral-typical  
1014 microbes (bold).
- 1015 (I) Relative abundance of TMB-high-associated species (*S. sputigena*, *F. animalis*, and *P.*  
1016 *micra*) by tumor site and genomic subtype.
- 1017 (J) Relative abundance of TMB-low-associated species (*V. parvula*, *P. vulgatus*, and *E. coli*)  
1018 by tumor site and genomic subtype.
- 1019 (K) Relative abundance of TMB-high-associated species (*F. animalis*, *S. sputigena*, *P. micra*,  
1020 *P. intermedia*, and *S. moorei*) by CRC tumor stage and genomic subtype. Data are  
1021 presented as median ± interquartile range.
- 1022 (L) Relative abundance of TMB-low-associated species (*P. vulgatus*, *F. prausnitzii*, and *B.*  
1023 *uniformis*) by CRC tumor stage and genomic subtype.

1024 \* *p* < 0.05, \*\* *p* < 0.01, \*\*\* *p* < 0.001, \*\*\*\* *p* < 0.0001 by Wilcoxon's signed-rank test unless  
1025 specified otherwise. Microbial counts are from Kraken2 unless stated otherwise.

1026

1027 **Figure S5. The tumor microbiome varies by tumor site, mutation burden, and age of onset**  
1028 **(Related to Figure 5)**

- 1029 (A) Microbial abundance (RPM) in MSS (gray), MSI (red), and POL-mutant (blue) colorectal  
1030 cancer subtypes, stratified by anatomic site.
- 1031 (B) Bacterial abundance (RPM) in MSS, MSI, and POL-mutant (blue) colorectal cancer  
1032 subtypes.
- 1033 (C) Fungal abundance (RPM) in MSS, MSI, and POL-mutant (blue) colorectal cancer  
1034 subtypes.
- 1035 (D) Viral abundance (RPM) in MSS, MSI, and POL CRC subtypes.
- 1036 (E) Archaea abundance (RPM) in MSS, MSI, and POL CRC subtypes.

- 1037 (F) Microbial diversity (Shannon index) in MSS, MSI, and POL-mutant (blue) colorectal cancer  
1038 subtypes, stratified by anatomic site.
- 1039 (G) Pearson correlation between tumor mutation burden (TMB; *x-axis*) and microbial  
1040 abundance (RPM; *y-axis*) in oropharyngeal cancers from 100kGP.
- 1041 (H) Pearson correlation between tumor mutation burden (TMB; *x-axis*) and microbial  
1042 abundance (RPM; *y-axis*) in oropharyngeal cancers from TCMA.
- 1043 (I) Pearson correlation between tumor mutation burden (TMB; *x-axis*) and microbial  
1044 abundance (RPM; *y-axis*) in gastroesophageal cancers from 100kGP.
- 1045 (J) Pearson correlation between tumor mutation burden (TMB; *x-axis*) and microbial  
1046 abundance (RPM; *y-axis*) in gastroesophageal cancers from TCMA.
- 1047 (K) Volcano plot showing differences in microbial composition between average-onset (AO;  
1048 *green*) and early-onset (EO; *purple*) in proximal colorectal tumors (MSS only).
- 1049 (L) Volcano plot showing differences in microbial composition between average-onset (AO;  
1050 *green*) and early-onset (EO; *purple*) in distal colorectal tumors (MSS only).
- 1051 (M) Relative abundance of *Akkermansia muciniphila* in EO versus AO colorectal cancer in  
1052 100kGP, stratified by tumor site.
- 1053 (N) Relative abundance of *Akkermansia muciniphila* in EO versus AO colorectal cancer in  
1054 TCMA<sup>57</sup>, stratified by tumor site.

1055 \*  $p < 0.05$ , \*\*  $p < 0.01$ , \*\*\*  $p < 0.001$ , \*\*\*\*  $p < 0.0001$  by Wilcoxon's signed-rank test unless  
1056 specified otherwise. Microbial counts are from Kraken2 unless stated otherwise.

## 1057 STAR Methods

### 1058 In silico microbe-human DNA mixtures

1059 To model sequencing of microbe-human mixtures, we simulated 151bp paired-end reads  
1060 from a defined microbiome community and the human reference, then mixed them at predefined  
1061 ratios. The microbial composition and relative abundance of these mixtures was designed to  
1062 mirror the ZymoBIOMICS Gut Microbiome Standard (Zymo Research, Cat. #D6331), which  
1063 contains 17 species at varying proportions, spanning bacteria (e.g. *F. nucleatum*, *B. fragilis*, *E.*  
1064 *coli*, *A. muciniphila*), fungi (*C. albicans* and *S. cerevisiae*), and archaea (*M. smithii*). Microbial  
1065 genomes were obtained from the vendor-provided reference bundle  
1066 (<https://s3.amazonaws.com/zymo-files/BioPool/D6331.refseq.zip>), while human reads were  
1067 simulated from the CHM13v2.0 assembly (RefSeq accession GCF\_009914755.1).

1068 Illumina sequencing output was simulated using InSilicoSeq<sup>160</sup> v2.0.1 (NovaSeq error  
1069 model), then down-sampled using seqtk v1.2-r94 (<https://github.com/lh3/seqtk>) using a fixed  
1070 random seed (100). For each mixture, simulated microbial and human reads were combined at  
1071 fixed, 10-fold microbe:host ratios ranging from 1:1 to 1:10<sup>8</sup>, keeping total read counts constant at  
1072 300 million per sample (**Table S1**), approximating ~30X human WGS depth. The resulting  
1073 simulated sequencing outputs of microbe-human mixtures were aligned to GRCh38 (NCI  
1074 Genomic Data Commons resource: [https://api.gdc.cancer.gov/data/254f697d-310d-4d7d-a27b-  
1075 27fbf767a834](https://api.gdc.cancer.gov/data/254f697d-310d-4d7d-a27b-27fbf767a834)) using BWA-MEM with GDC parameters (-T 0) to generate GRCh38-aligned  
1076 BAMs. These BAMs were then processed using PathSeq-T2T for host filtering and downstream  
1077 microbial profiling.

### 1078 **In vitro microbe-human DNA mixtures**

1079 Microbial and human DNA were mixed *in vitro* using microbial DNA derived from the  
1080 ZymoBIOMICS Gut Microbiome Standard (Zymo Research, Cat. #D6331) and human DNA  
1081 derived from colorectal cancer cell line HCT116 (ATCC, Manassas, VA, USA), cultured in  
1082 McCoy's 5A medium supplemented with 10% FBS and 1% penicillin-streptomycin. Microbial DNA  
1083 was extracted using the ZymoBIOMICS DNA Miniprep Kit (Zymo Research, Cat. #D4300).  
1084 Human DNA was extracted using the DNeasy Blood & Tissue Kit (QIAGEN, Cat. #69504). All  
1085 DNA extractions were performed following the manufacturers' instructions. DNA concentrations  
1086 were quantified using Qubit™ dsDNA Quantification Assay Kit (Thermo Fisher Scientific, Cat.  
1087 #Q32851).

1088 Purified microbial and human DNA was then combined to generate an eight-point 10-fold  
1089 serial dilution series in triplicate ( $n = 24$ ), spanning microbe:human DNA ratios ranging from 1:1  
1090 to 1:10<sup>8</sup>. To assess technical and biological contamination, we also included human-only ( $n = 3$ )  
1091 and template-free extraction negative controls ( $n = 3$ ). The human-only negative control contained  
1092 purified HCT116 gDNA with no microbial input and was sequenced in parallel with the dilution  
1093 series; this control allowed us to assess the degree of cross-contamination between conditions,  
1094 either due to the inadvertent exchange of DNA between experimental conditions<sup>81,82</sup> or due to  
1095 misassignment of barcodes during the sequencing process<sup>82-84</sup>. The no-template extraction  
1096 control contained aliquots of purified, nuclease-free water which underwent human and microbial  
1097 DNA extraction in parallel then mixed in equal proportions; this control allowed us to assess the  
1098 introduction of contamination during the sample preparation and DNA extraction process. All  
1099 extractions and mixtures were performed under sterile conditions. Mixtures and human-only

1100 controls were prepared at a minimum final concentration of 15 ng/μL in a 100 μL volume, then  
1101 sent for sequencing together with no-template extraction controls.

1102 Library preparation and sequencing were performed by the Broad Institute's Genomics  
1103 Platform. PCR-free whole-genome sequencing (WGS) libraries were prepared from the microbial-  
1104 human mixtures and human-only gDNA and sequenced to ~30X human coverage. Template-free  
1105 negative controls were prepared as non-tagmented libraries and analyzed using mid-output (1.5  
1106 Gb) metagenomic WGS. Samples undergoing 30X WGS were demultiplexed, aggregated, and  
1107 aligned using Illumina DRAGEN and delivered as GRCh38-aligned CRAM files, then subjected  
1108 to host-filtering by PathSeq-T2T. Sequencing output for the template-free negative controls was  
1109 delivered as FASTQ files and screened for microbial sequences using Kraken2<sup>77</sup>.

### 1110 **Analysis of sequenced colorectal cancer cell lines**

1111 Whole genome sequencing data from colorectal cancer cell lines were obtained from The  
1112 Cancer Cell Line Encyclopedia (CCLE)<sup>85</sup> via the Terra workspace. We selected cell lines classified  
1113 as microsatellite-unstable ( $n = 8$ ) or microsatellite stable ( $n = 12$ ) by Cellosaurus<sup>161</sup> and annotated  
1114 by CCLE as being from PCR-Free ( $n = 10$ ) or non-PCR-free ( $n = 10$ ) libraries. Sequencing data  
1115 from these cell lines were analyzed with PathSeq-T2T using default parameters.

### 1116 **Host filtering and microbial classification with PathSeq-T2T**

1117 Whole-genome sequencing data in this study were host-filtered using PathSeq-T2T, a  
1118 pipeline we designed that extends our previous host-filtering tool, PathSeq<sup>66,67</sup>. This pipeline is  
1119 designed to process sequencing data pre-aligned to GRCh38, which is the standard delivery  
1120 format for sequencing data. First, GRCh38-aligned data are pre-filtered using samtools<sup>162</sup>, first  
1121 removing all reads mapped in proper pairs to the reference (parameters: -f 3) to retain  
1122 unmapped/discordant reads. Reads aligned using BWA-MEM are additionally subjected an  
1123 alignment score filter (-e '[AS]>35'). Reads properly paired but aligning to viral decoy sequences  
1124 (often included as alternative contigs in reference genomes to reduce false positives in variant  
1125 calling) are retained and cleared of alignment tags using the "RevertSam" function from Picard<sup>163</sup>.

1126 After prefiltering, both unaligned and viral-mapped sequences are then subjected to  
1127 quality, complexity, and additional host-filtering using GATK PathSeqFilterSpark (parameters: --  
1128 min-clipped-read-length 60 --host-is-aligned true). This step masks and/or removes reads with  
1129 low base quality (Phred < 15) and low complexity (entropy < 0.3), duplicates, as well as sequences  
1130 containing one or more k-mers matching (i) the human reference GRCh38, (ii) sequences from

1131 the highly variable human major histocompatibility complex from the Immuno Polymorphism  
1132 Database, (iii) cloning vector sequences from NCBI UniVec, (iv) known human transcripts from  
1133 Gencode (human v25), and (v) human breakpoint sequences from GenBank (KY503218,  
1134 KY5808060). This results in a set of high-quality, high-complexity, host-filtered paired and  
1135 singleton reads (the mates of some reads are removed at this step, due to the lower quality of  
1136 some reverse reads).

1137         Next, resulting paired and singleton reads are then aligned (independently for pairs and  
1138 singletons) to the complete human reference genome T2T-CHM13<sup>76</sup> with BWA MEM using the  
1139 Genomic Data Commons (GDC) default parameters (-T 0). Properly-paired or singleton reads  
1140 aligning to T2T-CHM13 (i.e. BWA alignment score > 35) are classified as human and removed,  
1141 leaving a final set of high-quality, putatively non-human reads.

1142         Next, reads passing host-filtering and quality-control are processed using Kraken2<sup>77</sup>,  
1143 MetaPhlan4<sup>78</sup>, and Sylph, complementary microbial classification tools commonly used for  
1144 metagenomic analyses. Kraken2<sup>77</sup> is a k-mer-based tool that tends to prioritize sensitivity over  
1145 specificity, allowing detection of microbial signals even when read counts are very low.  
1146 MetaPhlan4<sup>78</sup> relies on microbial marker-genes for classification and prioritizes specificity over  
1147 sensitivity, enabling high-confidence microbial calls provided that sufficient reads are present to  
1148 cover marker genes. Sylph<sup>79</sup> is another k-mer based tool which measures microbial genome  
1149 containment and is demonstrates high specificity.

1150         In this study, Kraken2 classification was performed using the “PlusPF” k-mer database  
1151 (release 2024-06-05; obtained from <https://benlangmead.github.io/aws-indexes/k2>). This  
1152 reference contains k-mers recovered from complete and chromosome-level genome assemblies  
1153 of bacteria, archaea, viruses, fungi, and protozoa as well as the human genome and UniVec  
1154 vector sequences obtained from RefSeq. PathSeq-T2T uses default Kraken2 parameters, with  
1155 the exception of the minimum k-mer fraction, which is increased to 15% (--confidence 0.15) from  
1156 the default (0.0) to increase specificity, which was found to be optimal for accuracy by Wright et  
1157 al.<sup>164</sup>. Because Kraken2 concatenates paired reads prior to classification, host-filtered paired and  
1158 unpaired reads are classified separately; afterwards, counts are merged. Throughout this  
1159 manuscript, given microbial abundance estimates are from Kraken2, unless specified otherwise.

1160         MetaPhlan4 classification in this study was performed using the pre-built CHOCOPHlan  
1161 SGB database (mpa\_vJun23\_CHOCOPHlan\_202403), which contains bacterial and archaeal

1162 marker genes. In PathSeq-T2T and in this study, default parameters are used, except for the  
1163 minimum read length (--read\_min\_len 60) and output format (-t rel\_ab\_w\_read\_stats), which are  
1164 added to enable microbial read count estimates based on marker gene coverage. Because  
1165 MetaPhlan4 does not take read pairing into consideration, paired and unpaired reads are merged  
1166 prior to classification.

1167 Sylph classification was performed using prebuilt reference databases, including  
1168 references for bacteria (gtdb-r226-c200-dbv1.syldb), fungi (fungi-refseq-2025-10-11-c200-  
1169 dbv1.syldb), and viruses (uhgv\_c200\_dbv1.syldb). In PathSeq-T2T and in this study, default  
1170 parameters are used with the exception of "--estimate-read-counts" and "--estimate-unknown",  
1171 which are used for downstream abundance estimation. Like Kraken2, Sylph handles paired and  
1172 unpaired reads differently; therefore these reads were processed separately and counts were  
1173 merged afterwards.

1174 Finally, read counts from Kraken2, MetaPhlan4, and Sylph are normalized to per-sample  
1175 sequencing depth using primary read counts obtained from samtools' flagstat, performed on the  
1176 original, unfiltered bam file. This value is used as a denominator to calculate the microbial reads  
1177 per million primary reads (RPM). This normalization against the total primary reads (human and  
1178 microbial) quantifies the microbial-to-human DNA ratio, providing an approximation of the  
1179 absolute microbial abundance in the sequenced tumor portion. This value be used to calculate  
1180 microbe-host cell ratios after normalizing to relative genome size.

### 1181 **Detection of microbial signals in 100,000 Genomes Project (100kGP) cohort**

1182 The 100kGP cohort includes tumor and normal sample pairs from cancer patients  
1183 recruited to the 100,000 Genomes Project (100kGP v18 release) through 13 Genomic Medicine  
1184 Centres in the United Kingdom. We analyzed the entire 100kGP cohort, consisting of 16,369  
1185 tumor samples and 15,249 matched germline sequencing experiments from a total of 15,237  
1186 participants. Written informed consent was obtained by Genomics England from all participants.  
1187 Tumor samples were sequenced to ~100X coverage and consisted of primary ( $n = 15,041$ ),  
1188 metastatic ( $n = 963$ ), and recurrent ( $n = 365$ ) lesions (**Table S2**). Germline samples were  
1189 sequenced to ~30X coverage and consisted primarily of blood ( $n = 14,127$ ) with additional saliva  
1190 ( $n = 490$ ) and other non-diseased tissues. Most tumor samples were fresh-frozen ( $n = 15,011$ )  
1191 and sequenced from PCR-free libraries ( $n = 14,420$ ), with the remaining samples predominantly  
1192 formalin-fixed, paraffin-embedded (FFPE;  $n = 614$ ) and sequenced using PCR-amplified libraries  
1193 ( $n = 1,949$ ), respectively. Sequencing data for both sample types were delivered as BAM files

1194 aligned to GRCh38 using the DRAGEN pipeline, which we analyzed using PathSeq-T2T  
1195 implemented in “Double Helix”, the high-performance computing (HPC) system within 100kGP  
1196 Research Environment’s virtual desktop. After decontamination, downstream microbiome  
1197 analyses were restricted to primary and recurrent fresh-frozen tumor samples from solid cancer  
1198 sites with a minimum of 50 samples available (i.e. bladder, brain, breast, colorectal, esophageal,  
1199 gastric, kidney, liver, lung, oropharyngeal, ovarian, pancreatic, prostate, skin, testicular, and  
1200 uterus) prepared from PCR-free libraries ( $n = 10,698$ ).

### 1201 **Decontamination of 100kGP tumor microbiomes**

1202 Decontamination was performed on host- and quality-filtered reads produced by PathSeq-  
1203 T2T, using species-level abundance estimates from Kraken2, MetaPhlan4, and Sylph. Our  
1204 approach was based on previously reported observations<sup>57,165</sup> that reagent and sample-handling  
1205 contaminants in cancer sequencing datasets tend to appear broadly and uniformly across distinct  
1206 sample types (e.g. tissue and blood), whereas species that are truly present in tumor samples  
1207 show sample-type-specific biological variation.

1208 To decontaminate 100kGP tumor microbiomes, we computed a pan-cancer  
1209 equiprevalence (PCE) score that quantifies species ubiquity across cancer types. This score was  
1210 obtained by calculating the coefficient of variation (CV) of species prevalence across cancer  
1211 types, rescaled to [0, 1] based on the maximum/minimum theoretical possible CV values a  
1212 species could be given, such that  $PCE = 1$  for species with perfectly equal distribution across  
1213 cancer types and  $PCE = 0$  for detection confined to a single cancer type. Prevalence within each  
1214 cancer type (100kGP “study name”) was defined as the percentage of samples with  $RPM > 0.1$ .  
1215 To reduce variance, PCE was only calculated using 29 clearly defined cancer types with  $\geq 50$   
1216 primary or recurrent tumors.

1217 As expected, common environmental contaminants and species that were equiprevalent  
1218 in blood and saliva or in CRC and GBM also clustered at high PCE (**Figure 2I and S2F**). Based  
1219 on this observation that nearly all of these species had  $PCE > 0.7$ , we selected 0.7 as a cutoff for  
1220 identifying and removing contaminants from the dataset. The only species excepted from this  
1221 threshold were viral taxa known to infect humans (HPV-16, HPV-18, MCPyV, HSV-1, EBV, CMV,  
1222 HHV-6A, HHV-6B, HHV-7), many of which are widely carried in the population.

### 1223 **Assembly of microbial genomes**

1224 To further validate microbial signals and assess contaminant signals, we performed  
1225 metagenome assembly on unmapped tumor sequencing reads from orodigestive cancer types,  
1226 based on the observation that tumors arising at these sites were unique in harboring significant  
1227 microbial signal. Tumor WGS reads not mapped to the human genome reference (GRCh38) were  
1228 extracted from the WGS BAM files using samtools<sup>162</sup> (v.1.16.1) with parameter “-f 4”, then  
1229 converted to FASTQ with the samtools bam2fq function. The unmapped FASTQ files were used  
1230 as input to trim\_galore (v. 0.6.10) with parameters “--illumina --stringency 5 --length 75 --max\_n  
1231 2 --trim-n -j 25”. The trimmed reads were then assembled into contigs using MEGAHIT<sup>159</sup> (v. 1.2.9)  
1232 with default parameters. The contigs were used to create a bowtie2<sup>166</sup> index to which the trimmed  
1233 reads were realigned using bowtie2 with default parameters. The aligned contigs were then  
1234 converted to BAM, sorted, and indexed using samtools (v. 1.16.1). Finally, the BAM files with  
1235 reads mapped to contigs were used as input to MetaBAT2 for genome assembly, with the contig  
1236 FASTA and the sorted BAM file of reads mapped to contigs as inputs. MetaBAT2 was run with  
1237 the default minimum contig length of 2,500 bp. Quality control of the final assembled genomes  
1238 was performed using CheckM2 (parameters: predict -x fa --force --lowmem). ANI analysis was  
1239 performed using pyani (v0.3.0-alpha) with the -m ANIm method. SGBs with at least 10 unique  
1240 genomes and completeness >50% based on CheckM2 results were included in the analysis.

#### 1241 **Analysis of pan-cancer absolute and relative microbial abundance**

1242 To quantify microbial load, we summed species-level RPM values after removing  
1243 contaminant species (PCE > 0.7). To allow efficient data visualization, sites were recoded from  
1244 cancer types according to their anatomic site of origin (e.g. “Glioblastoma multiforme” and “Low  
1245 grade glioma” were assigned to “Brain”). The threshold for “background” microbial signals was  
1246 set to 0.1 RPM, which was determined based on our observation that sequenced human-only  
1247 controls typically contained microbial read counts at approximately this level (**Figure 1G**). This  
1248 threshold was also used for determining species-level presence/absence.

1249 The relative abundance analysis (**Figure 3E**) was calculated by aggregating  
1250 decontaminated species-level read counts to the family level for “microbe-high” samples (RPM >  
1251 10). To visualize the average relative abundances at the level of cancer site, we averaged sample-  
1252 specific relative abundances before re-normalizing to the sum of family-level averages.

#### 1253 **Analysis of orodigestive tumor microbiomes**

1254 For analysis of orodigestive cancer types, samples were categorized as oropharyngeal  
1255 (oral and sinonasal cancers;  $n = 251$ ), gastroesophageal (esophageal and gastric cancers;  $n =$

1256 199), and colorectal (colon and rectal cancers;  $n = 2,422$ ). Curated clinical and anatomic  
1257 annotations of gastroesophageal and colorectal cancers were provided by Booth et al.<sup>92</sup> and  
1258 Cornish et al.<sup>73</sup>, respectively. Total microbial load was calculated by summing species-level RPM  
1259 values after removing contaminant species (PCE > 0.7). Diversity was measured using Shannon's  
1260 entropy.

1261 The site-specific analysis of species-level associations across oropharyngeal,  
1262 gastroesophageal, and colorectal cancers (**Figure 4D**) was performed by calculating a Kruskal-  
1263 Wallis test of site-specific enrichments (**Table S7**) which adjusted for multiple hypothesis tests  
1264 using the Benjamini-Hochberg correction. For radial visualization, the geometric mean abundance  
1265 (average  $\log_{10}$ RPM) was calculated for each species and site, creating a 3-vector defining a  
1266 triangular plane whose vertices represented each site and we calculated the barycentric centroid  
1267 of that triangle. This centroid was then mapped into polar coordinates, with angle encoding site  
1268 bias and radius encoding the magnitude of site selectivity.

#### 1269 **Principal coordinates analysis and PERMANOVA**

1270 Principal coordinates analyses (PCoA) and PERMANOVA analyses were performed in R  
1271 (v. 4.3.3) using the package "vegan" (v. 2.6.10) with genus-level read counts. Distance matrices  
1272 for PCoA and PERMANOVA were calculated using the Bray-Curtis dissimilarity index.  
1273 PERMANOVA was used to identify features predictive of microbial composition and was run using  
1274 vegan's "adonis2" function using the "margin" setting to identify features that independently  
1275 influence microbial composition.

#### 1276 **Validation using tumor microbiome data from The Cancer Microbiome Atlas**

1277 To validate findings regarding 100kGP tumor genomes, we used PathSeq-filtered  
1278 sequences from The Cancer Microbiome Atlas (TCMA)<sup>167</sup>, a decontaminated tumor microbiome  
1279 database based on tumor whole-genome sequencing from The Cancer Genome Atlas (TCGA).  
1280 These sequences were reanalyzed using downstream PathSeq-T2T filtering steps. They were  
1281 then decontaminated using the same methods as the original TCMA paper. Mutation burden for  
1282 correlative analyses were obtained and calculated using the PanCanAtlas<sup>168</sup>a Mutations file  
1283 (mc3.v0.2.8.PUBLIC.maf.gz). Pan-cancer annotations for microsatellite-instability status were  
1284 taken from MANTIS<sup>169</sup>.

#### 1285 **Regression analyses of tumor mutation burden and microbial load**

1286 For our analysis of CRC tumor microbiomes, we compared tumor mutation burden (TMB)  
1287 which we measured using  $\log_{10}$  mutations per megabase (SNVs and indels). This was provided  
1288 by Cornish et al.<sup>73</sup> along with CRC genomic subtyping (MSS/MSI/POL) and anatomical site  
1289 (Proximal/Distal/Rectal) annotations, which were used as covariates. Microbial load was  
1290 measured using  $\log_{10}$ RPM representing the sum of species-level RPM values after removing  
1291 contaminants (PCE > 0.7).

1292 We used ordinary least squares (OLS) with three complementary approaches: (i) within-  
1293 subtype regressions to characterize subtype-specific correlations (**Figure 5E**), and (ii)  
1294 residualization to estimate a tumor subtype- and site-adjusted correlations (**Figure 5F-G**), (iii) a  
1295 nested model comparison to test whether genomic subtype adds predictive value beyond TMB  
1296 and anatomic site. Model designs in R-formula format are given below.

1297 Model 1: Bivariate models, trained separately for each subtype (MSS, MSI, and POL).

1298  $\log\_RPM \sim \log\_TMB\_per\_mb$

1299 Model 2: Computing correlations for residualized TMB and microbial load.

- 1300 1.  $\log\_RPM \sim C(anatomic\_site) \Rightarrow RPM\_resid$   
1301 2.  $\log\_TMB\_per\_mb \sim C(subtype) \Rightarrow TMB\_resid$   
1302 3.  $RPM\_resid \sim TMB\_resid$

1303 Model 3: Nested models predicting microbial load, compared with ANOVA F-test.

- 1304 1.  $\log\_RPM \sim \log\_TMB\_per\_mb + C(anatomic\_site) + C(subtype)$   
1305 2.  $\log\_RPM \sim \log\_TMB\_per\_mb + C(anatomic\_site)$

1306 Regression analyses were performed using python's "statsmodels" module (v. 0.13.5).

### 1307 **Bivariate analysis of CRC tumor site and TMB subtype**

1308 To jointly model the influence of CRC tumor site and TMB subtype (TMB-high versus TMB-  
1309 low), we performed a blocked analysis of species-level relative abundance (**Figure 5H**). For our  
1310 analysis of CRC microbiome variation by tumor site and subtype, we defined CRC samples with  
1311  $\log_{10}TMB > 1.25$  as "High-TMB", and samples with  $\log_{10}TMB \leq 1.25$  as "Low-TMB", and  
1312 categorized CRCs arising in the cecum, ascending colon, transverse colon, or hepatic flexure as

1313 “proximal” and CRCs arising in the splenic flexure, descending colon, sigmoid colon, rectosigmoid  
1314 junction, or rectum as “distal”.

1315 Post-decontamination species-level relative abundances were calculated for each cancer.  
1316 For species detected in >10% of CRC samples, we used a Wilcoxon rank-sum test to evaluate  
1317 site and subtype associations for each species. This test was used to test four pairwise  
1318 comparisons: (1) proximal vs. distal, within TMB-low; (2) proximal vs. distal, within TMB-high; (3)  
1319 TMB-high vs. TMB-low, within proximal; and (4) TMB-high vs. TMB-low, within distal. This  
1320 blocking was implemented to control confounding effects, since TMB-high subtypes occur more  
1321 frequently in the proximal colon while TMB-low subtypes are more common distal CRCs. Test  
1322 statistics for both “proximal vs. distal” comparisons were averaged, as were statistics for both  
1323 “TMB-high vs. TMB-low” comparisons, to obtain a composite Wilcoxon test statistic. The *p*-values  
1324 for matched comparisons were combined using Fisher’s method, then corrected for multiple  
1325 comparisons using Benjamini-Hochberg. Candidate CRC biomarker species and oral-typical  
1326 annotations were provided by Piccinno et al.<sup>127</sup>.

### 1327 **Analysis of microbial composition and early-onset CRC**

1328 To study microbial signals associated with CRC age of onset, we evaluated species-level  
1329 relative abundances between early-onset (EO) and average-onset (AO) CRC. We defined EO-  
1330 CRC as cases for which patient age was <50 years at the time of tumor sampling and AO-CRC  
1331 as cases for which patients were ≥50 years of age at sampling. We restricted our analysis of age  
1332 of onset to MSS cancers. Due to low sample numbers of EO-CRC cases, comparisons were not  
1333 possible for high-TMB subtypes (MSI/POL). We performed relative abundance comparisons for  
1334 both proximal MSS and distal MSS using Wilcoxon’s rank-sum test on relative abundances, as  
1335 described above. Validation analyses were performed using decontaminated CRC tumor  
1336 microbiome profiles from The Cancer Microbiome Atlas<sup>57</sup> (TCMA; available at  
1337 <https://doi.org/10.7924/r4bk1j35s>).

1338

1339

1340

1341

1342 **References**

- 1343 1. Sender, R., Fuchs, S., and Milo, R. (2016). Revised Estimates for the Number of Human and  
1344 Bacteria Cells in the Body. *PLoS Biol* 14, e1002533.  
1345 <https://doi.org/10.1371/journal.pbio.1002533>.
- 1346 2. Berg, R.D. (1996). The indigenous gastrointestinal microflora. *Trends in Microbiology* 4, 430–  
1347 435. [https://doi.org/10.1016/0966-842X\(96\)10057-3](https://doi.org/10.1016/0966-842X(96)10057-3).
- 1348 3. Lloyd-Price, J., Mahurkar, A., Rahnavard, G., Crabtree, J., Orvis, J., Hall, A.B., Brady, A.,  
1349 Creasy, H.H., McCracken, C., Giglio, M.G., et al. (2017). Strains, functions and dynamics in  
1350 the expanded Human Microbiome Project. *Nature* 550, 61–66.  
1351 <https://doi.org/10.1038/nature23889>.
- 1352 4. Huttenhower, C., Gevers, D., Knight, R., Abubucker, S., Badger, J.H., Chinwalla, A.T.,  
1353 Creasy, H.H., Earl, A.M., FitzGerald, M.G., Fulton, R.S., et al. (2012). Structure, function and  
1354 diversity of the healthy human microbiome. *Nature* 486, 207–214.  
1355 <https://doi.org/10.1038/nature11234>.
- 1356 5. Yang, K., Li, G., Li, Q., Wang, W., Zhao, X., Shao, N., Qiu, H., Liu, J., Xu, L., and Zhao, J.  
1357 (2025). Distribution of gut microbiota across intestinal segments and their impact on human  
1358 physiological and pathological processes. *Cell & Bioscience* 15, 47.  
1359 <https://doi.org/10.1186/s13578-025-01385-y>.
- 1360 6. Segata, N., Haake, S.K., Mannon, P., Lemon, K.P., Waldron, L., Gevers, D., Huttenhower,  
1361 C., and Izard, J. (2012). Composition of the adult digestive tract bacterial microbiome based  
1362 on seven mouth surfaces, tonsils, throat and stool samples. *Genome Biology* 13, R42.  
1363 <https://doi.org/10.1186/gb-2012-13-6-r42>.
- 1364 7. Stearns, J.C., Lynch, M.D.J., Senadheera, D.B., Tenenbaum, H.C., Goldberg, M.B.,  
1365 Cvitkovitch, D.G., Croitoru, K., Moreno-Hagelsieb, G., and Neufeld, J.D. (2011). Bacterial  
1366 biogeography of the human digestive tract. *Sci Rep* 1, 170.  
1367 <https://doi.org/10.1038/srep00170>.
- 1368 8. Donaldson, G.P., Lee, S.M., and Mazmanian, S.K. (2016). Gut biogeography of the bacterial  
1369 microbiota. *Nat Rev Microbiol* 14, 20–32. <https://doi.org/10.1038/nrmicro3552>.
- 1370 9. Hillman, E.T., Lu, H., Yao, T., and Nakatsu, C.H. (2017). Microbial Ecology along the  
1371 Gastrointestinal Tract. *Microbes Environ* 32, 300–313.  
1372 <https://doi.org/10.1264/jsme2.ME17017>.
- 1373 10. Eckburg, P.B., Bik, E.M., Bernstein, C.N., Purdom, E., Dethlefsen, L., Sargent, M., Gill, S.R.,  
1374 Nelson, K.E., and Relman, D.A. (2005). Diversity of the Human Intestinal Microbial Flora.  
1375 *Science* 308, 1635–1638. <https://doi.org/10.1126/science.1110591>.
- 1376 11. Baker, J.L., Mark Welch, J.L., Kauffman, K.M., McLean, J.S., and He, X. (2024). The oral  
1377 microbiome: diversity, biogeography and human health. *Nat Rev Microbiol* 22, 89–104.  
1378 <https://doi.org/10.1038/s41579-023-00963-6>.
- 1379 12. Byrd, A.L., Belkaid, Y., and Segre, J.A. (2018). The human skin microbiome. *Nat Rev*  
1380 *Microbiol* 16, 143–155. <https://doi.org/10.1038/nrmicro.2017.157>.

- 1381 13. Oliphant, K., and Allen-Vercoe, E. (2019). Macronutrient metabolism by the human gut  
1382 microbiome: major fermentation by-products and their impact on host health. *Microbiome* 7,  
1383 91. <https://doi.org/10.1186/s40168-019-0704-8>.
- 1384 14. Le Chatelier, E., Nielsen, T., Qin, J., Prifti, E., Hildebrand, F., Falony, G., Almeida, M.,  
1385 Arumugam, M., Batto, J.-M., Kennedy, S., et al. (2013). Richness of human gut microbiome  
1386 correlates with metabolic markers. *Nature* 500, 541–546.  
1387 <https://doi.org/10.1038/nature12506>.
- 1388 15. Donald, K., and Finlay, B.B. (2023). Early-life interactions between the microbiota and  
1389 immune system: impact on immune system development and atopic disease. *Nat Rev*  
1390 *Immunol* 23, 735–748. <https://doi.org/10.1038/s41577-023-00874-w>.
- 1391 16. Hillier, S.L., Nugent, R.P., Eschenbach, D.A., Krohn, M.A., Gibbs, R.S., Martin, D.H., Cotch,  
1392 M.F., Edelman, R., Pastorek, J.G., and Rao, A.V. (1995). Association between bacterial  
1393 vaginosis and preterm delivery of a low-birth-weight infant. The Vaginal Infections and  
1394 Prematurity Study Group. *N Engl J Med* 333, 1737–1742.  
1395 <https://doi.org/10.1056/NEJM199512283332604>.
- 1396 17. Fettweis, J.M., Serrano, M.G., Brooks, J.P., Edwards, D.J., Girerd, P.H., Parikh, H.I., Huang,  
1397 B., Arodz, T.J., Edupuganti, L., Glascock, A.L., et al. (2019). The vaginal microbiome and  
1398 preterm birth. *Nat Med* 25, 1012–1021. <https://doi.org/10.1038/s41591-019-0450-2>.
- 1399 18. Buffie, C.G., Bucci, V., Stein, R.R., McKenney, P.T., Ling, L., Gobourne, A., No, D., Liu, H.,  
1400 Kinnebrew, M., Viale, A., et al. (2015). Precision microbiome reconstitution restores bile acid  
1401 mediated resistance to *Clostridium difficile*. *Nature* 517, 205–208.  
1402 <https://doi.org/10.1038/nature13828>.
- 1403 19. Hsiao, A., Ahmed, A.M.S., Subramanian, S., Griffin, N.W., Drewry, L.L., Petri, W.A., Haque,  
1404 R., Ahmed, T., and Gordon, J.I. (2014). Members of the human gut microbiota involved in  
1405 recovery from *Vibrio cholerae* infection. *Nature* 515, 423–426.  
1406 <https://doi.org/10.1038/nature13738>.
- 1407 20. Garrett, W.S., Lord, G.M., Punit, S., Lugo-Villarino, G., Mazmanian, S.K., Ito, S., Glickman,  
1408 J.N., and Glimcher, L.H. (2007). Communicable ulcerative colitis induced by T-bet deficiency  
1409 in the innate immune system. *Cell* 131, 33–45. <https://doi.org/10.1016/j.cell.2007.08.017>.
- 1410 21. Lloyd-Price, J., Arze, C., Ananthakrishnan, A.N., Schirmer, M., Avila-Pacheco, J., Poon, T.W.,  
1411 Andrews, E., Ajami, N.J., Bonham, K.S., Brislawn, C.J., et al. (2019). Multi-omics of the gut  
1412 microbial ecosystem in inflammatory bowel diseases. *Nature* 569, 655–662.  
1413 <https://doi.org/10.1038/s41586-019-1237-9>.
- 1414 22. Karlsson, F.H., Tremaroli, V., Nookaew, I., Bergström, G., Behre, C.J., Fagerberg, B., Nielsen,  
1415 J., and Bäckhed, F. (2013). Gut metagenome in European women with normal, impaired and  
1416 diabetic glucose control. *Nature* 498, 99–103. <https://doi.org/10.1038/nature12198>.
- 1417 23. Qin, J., Li, Y., Cai, Z., Li, S., Zhu, J., Zhang, F., Liang, S., Zhang, W., Guan, Y., Shen, D., et  
1418 al. (2012). A metagenome-wide association study of gut microbiota in type 2 diabetes. *Nature*  
1419 490, 55–60. <https://doi.org/10.1038/nature11450>.

- 1420 24. Kostic, A.D., Chun, E., Robertson, L., Glickman, J.N., Gallini, C.A., Michaud, M., Clancy, T.E.,  
1421 Chung, D.C., Lochhead, P., Hold, G.L., et al. (2013). *Fusobacterium nucleatum* potentiates  
1422 intestinal tumorigenesis and modulates the tumor-immune microenvironment. *Cell Host*  
1423 *Microbe* 14, 207–215. <https://doi.org/10.1016/j.chom.2013.07.007>.
- 1424 25. Nomura, A., Stemmermann, G.N., Chyou, P.-H., Kato, I., Perez-Perez, G.I., and Blaser, M.J.  
1425 (1991). *Helicobacter pylori* Infection and Gastric Carcinoma among Japanese Americans in  
1426 Hawaii. *New England Journal of Medicine* 325, 1132–1136.  
1427 <https://doi.org/10.1056/NEJM199110173251604>.
- 1428 26. zur Hausen, H. (2002). Papillomaviruses and cancer: from basic studies to clinical application.  
1429 *Nat Rev Cancer* 2, 342–350. <https://doi.org/10.1038/nrc798>.
- 1430 27. El Tekle, G., and Garrett, W.S. (2023). Bacteria in cancer initiation, promotion and  
1431 progression. *Nat Rev Cancer* 23, 600–618. <https://doi.org/10.1038/s41568-023-00594-2>.
- 1432 28. Matson, V., Fessler, J., Bao, R., Chongsuwat, T., Zha, Y., Alegre, M.-L., Luke, J.J., and  
1433 Gajewski, T.F. (2018). The commensal microbiome is associated with anti-PD-1 efficacy in  
1434 metastatic melanoma patients. *Science* 359, 104–108.  
1435 <https://doi.org/10.1126/science.aao3290>.
- 1436 29. Routy, B., Le Chatelier, E., Derosa, L., Duong, C.P.M., Alou, M.T., Daillère, R., Fluckiger, A.,  
1437 Messaoudene, M., Rauber, C., Roberti, M.P., et al. (2018). Gut microbiome influences  
1438 efficacy of PD-1-based immunotherapy against epithelial tumors. *Science* 359, 91–97.  
1439 <https://doi.org/10.1126/science.aan3706>.
- 1440 30. Gopalakrishnan, V., Spencer, C.N., Nezi, L., Reuben, A., Andrews, M.C., Karpnits, T.V.,  
1441 Prieto, P.A., Vicente, D., Hoffman, K., Wei, S.C., et al. (2018). Gut microbiome modulates  
1442 response to anti-PD-1 immunotherapy in melanoma patients. *Science* 359, 97–103.  
1443 <https://doi.org/10.1126/science.aan4236>.
- 1444 31. Park, J.S., Gazzaniga, F.S., Wu, M., Luthens, A.K., Gillis, J., Zheng, W., LaFleur, M.W.,  
1445 Johnson, S.B., Morad, G., Park, E.M., et al. (2023). Targeting PD-L2–RGMB overcomes  
1446 microbiome-related immunotherapy resistance. *Nature* 617, 377–385.  
1447 <https://doi.org/10.1038/s41586-023-06026-3>.
- 1448 32. Gazzaniga, F.S., and Kasper, D.L. (2025). The gut microbiome and cancer response to  
1449 immune checkpoint inhibitors. *J Clin Invest* 135. <https://doi.org/10.1172/JCI184321>.
- 1450 33. Derosa, L., Routy, B., Thomas, A.M., Iebba, V., Zalcman, G., Friard, S., Mazieres, J.,  
1451 Audigier-Valette, C., Moro-Sibilot, D., Goldwasser, F., et al. (2022). Intestinal *Akkermansia*  
1452 *muciniphila* predicts clinical response to PD-1 blockade in patients with advanced non-small-  
1453 cell lung cancer. *Nat Med* 28, 315–324. <https://doi.org/10.1038/s41591-021-01655-5>.
- 1454 34. Tintelnot, J., Xu, Y., Lesker, T.R., Schönlein, M., Konczalla, L., Giannou, A.D., Pelczar, P.,  
1455 Kyliès, D., Puelles, V.G., Bielecka, A.A., et al. (2023). Microbiota-derived 3-IAA influences  
1456 chemotherapy efficacy in pancreatic cancer. *Nature* 615, 168–174.  
1457 <https://doi.org/10.1038/s41586-023-05728-y>.
- 1458 35. Kostic, A.D., Gevers, D., Pedamallu, C.S., Michaud, M., Duke, F., Earl, A.M., Ojesina, A.I.,  
1459 Jung, J., Bass, A.J., Taberero, J., et al. (2012). Genomic analysis identifies association of

- 1460 Fusobacterium with colorectal carcinoma. *Genome Res* 22, 292–298.  
1461 <https://doi.org/10.1101/gr.126573.111>.
- 1462 36. Castellarin, M., Warren, R.L., Freeman, J.D., Dreolini, L., Krzywinski, M., Strauss, J., Barnes,  
1463 R., Watson, P., Allen-Vercoe, E., Moore, R.A., et al. (2012). *Fusobacterium nucleatum*  
1464 infection is prevalent in human colorectal carcinoma. *Genome Res* 22, 299–306.  
1465 <https://doi.org/10.1101/gr.126516.111>.
- 1466 37. Geller, L.T., Barzily-Rokni, M., Danino, T., Jonas, O.H., Shental, N., Nejman, D., Gavert, N.,  
1467 Zwang, Y., Cooper, Z.A., Shee, K., et al. (2017). Potential role of intratumor bacteria in  
1468 mediating tumor resistance to the chemotherapeutic drug gemcitabine. *Science* 357, 1156–  
1469 1160. <https://doi.org/10.1126/science.aah5043>.
- 1470 38. Riquelme, E., Zhang, Y., Zhang, L., Montiel, M., Zoltan, M., Dong, W., Quesada, P., Sahin, I.,  
1471 Chandra, V., San Lucas, A., et al. (2019). Tumor Microbiome Diversity and Composition  
1472 Influence Pancreatic Cancer Outcomes. *Cell* 178, 795-806.e12.  
1473 <https://doi.org/10.1016/j.cell.2019.07.008>.
- 1474 39. Nejman, D., Livyatan, I., Fuks, G., Gavert, N., Zwang, Y., Geller, L.T., Rotter-Maskowitz, A.,  
1475 Weiser, R., Mallel, G., Gigi, E., et al. (2020). The human tumor microbiome is composed of  
1476 tumor type-specific intracellular bacteria. *Science* 368, 973–980.  
1477 <https://doi.org/10.1126/science.aay9189>.
- 1478 40. Raab-Traub, N. (2002). Epstein-Barr virus in the pathogenesis of NPC. *Semin Cancer Biol*  
1479 12, 431–441. <https://doi.org/10.1016/s1044579x0200086x>.
- 1480 41. Shibata, D., and Weiss, L.M. (1992). Epstein-Barr virus-associated gastric adenocarcinoma.  
1481 *Am J Pathol* 140, 769–774.
- 1482 42. Arzumanyan, A., Reis, H.M.G.P.V., and Feitelson, M.A. (2013). Pathogenic mechanisms in  
1483 HBV- and HCV-associated hepatocellular carcinoma. *Nat Rev Cancer* 13, 123–135.  
1484 <https://doi.org/10.1038/nrc3449>.
- 1485 43. Parsonnet, J., Friedman, G.D., Vandersteen, D.P., Chang, Y., Vogelman, J.H., Orentreich,  
1486 N., and Sibley, R.K. (1991). *Helicobacter pylori* infection and the risk of gastric carcinoma. *N*  
1487 *Engl J Med* 325, 1127–1131. <https://doi.org/10.1056/NEJM199110173251603>.
- 1488 44. Uemura, N., Okamoto, S., Yamamoto, S., Matsumura, N., Yamaguchi, S., Yamakido, M.,  
1489 Taniyama, K., Sasaki, N., and Schlemper, R.J. (2001). *Helicobacter pylori* Infection and the  
1490 Development of Gastric Cancer. *New England Journal of Medicine* 345, 784–789.  
1491 <https://doi.org/10.1056/NEJMoa001999>.
- 1492 45. Wu, S., Rhee, K.-J., Albesiano, E., Rabizadeh, S., Wu, X., Yen, H.-R., Huso, D.L., Brancati,  
1493 F.L., Wick, E., McAllister, F., et al. (2009). A human colonic commensal promotes colon  
1494 tumorigenesis via activation of T helper type 17 T cell responses. *Nat Med* 15, 1016–1022.  
1495 <https://doi.org/10.1038/nm.2015>.
- 1496 46. Dejea, C.M., Fathi, P., Craig, J.M., Boleij, A., Taddese, R., Geis, A.L., Wu, X., DeStefano  
1497 Shields, C.E., Hechenbleikner, E.M., Huso, D.L., et al. (2018). Patients with familial  
1498 adenomatous polyposis harbor colonic biofilms containing tumorigenic bacteria. *Science* 359,  
1499 592–597. <https://doi.org/10.1126/science.aah3648>.

- 1500 47. Nougayrède, J.-P., Homburg, S., Taieb, F., Boury, M., Brzuszkiewicz, E., Gottschalk, G.,  
1501 Buchrieser, C., Hacker, J., Dobrindt, U., and Oswald, E. (2006). *Escherichia coli* induces DNA  
1502 double-strand breaks in eukaryotic cells. *Science* 313, 848–851.  
1503 <https://doi.org/10.1126/science.1127059>.
- 1504 48. Wilson, M.R., Jiang, Y., Villalta, P.W., Stornetta, A., Boudreau, P.D., Carrá, A., Brennan, C.A.,  
1505 Chun, E., Ngo, L., Samson, L.D., et al. (2019). The human gut bacterial genotoxin colibactin  
1506 alkylates DNA. *Science* 363, eaar7785. <https://doi.org/10.1126/science.aar7785>.
- 1507 49. Pleguezuelos-Manzano, C., Puschhof, J., Rosendahl Huber, A., van Hoeck, A., Wood, H.M.,  
1508 Nomburg, J., Gurjao, C., Manders, F., Dalmasso, G., Stege, P.B., et al. (2020). Mutational  
1509 signature in colorectal cancer caused by genotoxic pks+ *E. coli*. *Nature* 580, 269–273.  
1510 <https://doi.org/10.1038/s41586-020-2080-8>.
- 1511 50. Arthur, J.C., Perez-Chanona, E., Mühlbauer, M., Tomkovich, S., Uronis, J.M., Fan, T.-J.,  
1512 Campbell, B.J., Abujamel, T., Dogan, B., Rogers, A.B., et al. (2012). Intestinal inflammation  
1513 targets cancer-inducing activity of the microbiota. *Science* 338, 120–123.  
1514 <https://doi.org/10.1126/science.1224820>.
- 1515 51. Xue, M., Kim, C.S., Healy, A.R., Wernke, K.M., Wang, Z., Frischling, M.C., Shine, E.E., Wang,  
1516 W., Herzon, S.B., and Crawford, J.M. (2019). Structure elucidation of colibactin and its DNA  
1517 cross-links. *Science* 365, eaax2685. <https://doi.org/10.1126/science.aax2685>.
- 1518 52. Carlson, E.S., Haslecker, R., Lecchi, C., Aguilar Ramos, M.A., Vennelakanti, V., Honaker, L.,  
1519 Stornetta, A., Millán, E.S., Johnson, B.A., Kulik, H.J., et al. (2025). The specificity and  
1520 structure of DNA cross-linking by the gut bacterial genotoxin colibactin. *Science* 390,  
1521 eady3571. <https://doi.org/10.1126/science.ady3571>.
- 1522 53. de Martel, C., Georges, D., Bray, F., Ferlay, J., and Clifford, G.M. (2020). Global burden of  
1523 cancer attributable to infections in 2018: a worldwide incidence analysis. *Lancet Glob Health*  
1524 8, e180–e190. [https://doi.org/10.1016/S2214-109X\(19\)30488-7](https://doi.org/10.1016/S2214-109X(19)30488-7).
- 1525 54. Gihawi, A., Cooper, C.S., and Brewer, D.S. (2023). Caution regarding the specificities of pan-  
1526 cancer microbial structure. *Microb Genom* 9, mgen001088.  
1527 <https://doi.org/10.1099/mgen.0.001088>.
- 1528 55. Poore, G.D., Kopylova, E., Zhu, Q., Carpenter, C., Fraraccio, S., Wandro, S., Kosciolk, T.,  
1529 Janssen, S., Metcalf, J., Song, S.J., et al. (2020). RETRACTED ARTICLE: Microbiome  
1530 analyses of blood and tissues suggest cancer diagnostic approach. *Nature* 579, 567–574.  
1531 <https://doi.org/10.1038/s41586-020-2095-1>.
- 1532 56. Sepich-Poore, G.D., McDonald, D., Kopylova, E., Guccione, C., Zhu, Q., Austin, G.,  
1533 Carpenter, C., Fraraccio, S., Wandro, S., Kosciolk, T., et al. (2024). Robustness of cancer  
1534 microbiome signals over a broad range of methodological variation. *Oncogene* 43, 1127–  
1535 1148. <https://doi.org/10.1038/s41388-024-02974-w>.
- 1536 57. Dohlman, A.B., Arguijo Mendoza, D., Ding, S., Gao, M., Dressman, H., Iliev, I.D., Lipkin, S.M.,  
1537 and Shen, X. (2021). The cancer microbiome atlas: a pan-cancer comparative analysis to  
1538 distinguish tissue-resident microbiota from contaminants. *Cell Host Microbe* 29, 281–298.e5.  
1539 <https://doi.org/10.1016/j.chom.2020.12.001>.

- 1540 58. Gihawi, A., Ge, Y., Lu, J., Puiu, D., Xu, A., Cooper, C.S., Brewer, D.S., Pertea, M., and  
1541 Salzberg, S.L. (2023). Major data analysis errors invalidate cancer microbiome findings. *mBio*  
1542 *14*, e0160723. <https://doi.org/10.1128/mbio.01607-23>.
- 1543 59. Dohlman, A.B., Klug, J., Mesko, M., Gao, I.H., Lipkin, S.M., Shen, X., and Iliev, I.D. (2022). A  
1544 pan-cancer mycobiome analysis reveals fungal involvement in gastrointestinal and lung  
1545 tumors. *Cell* *185*, 3807-3822.e12. <https://doi.org/10.1016/j.cell.2022.09.015>.
- 1546 60. Narunsky-Haziza, L., Sepich-Poore, G.D., Livyatan, I., Asraf, O., Martino, C., Nejman, D.,  
1547 Gavert, N., Stajich, J.E., Amit, G., González, A., et al. (2022). Pan-cancer analyses reveal  
1548 cancer-type-specific fungal ecologies and bacteriome interactions. *Cell* *185*, 3789-3806.e17.  
1549 <https://doi.org/10.1016/j.cell.2022.09.005>.
- 1550 61. Gihawi, A., Wood, H.M., Clark, J., O'Grady, J., Eeles, R.A., Wedge, D.C., Jakobsdottir, G.M.,  
1551 Magiorkinis, G., Schache, A.G., Masterson, L., et al. (2025). The landscape of microbial  
1552 associations in human cancer. *Science Translational Medicine* *17*, eads6166.  
1553 <https://doi.org/10.1126/scitranslmed.ads6166>.
- 1554 62. Ge, Y., Lu, J., Puiu, D., Revsine, M., and Salzberg, S.L. (2025). Comprehensive analysis of  
1555 microbial content in whole-genome sequencing samples from The Cancer Genome Atlas  
1556 project. *Sci Transl Med* *17*, eads6335. <https://doi.org/10.1126/scitranslmed.ads6335>.
- 1557 63. Ghaddar, B.C., Blaser, M.J., and De, S. (2025). Revisiting the cancer microbiome using  
1558 PRISM. Preprint at bioRxiv, <https://doi.org/10.1101/2025.01.21.634087>  
1559 <https://doi.org/10.1101/2025.01.21.634087>.
- 1560 64. Robinson, K.M., Crabtree, J., Mattick, J.S.A., Anderson, K.E., and Dunning Hotopp, J.C.  
1561 (2017). Distinguishing potential bacteria-tumor associations from contamination in a  
1562 secondary data analysis of public cancer genome sequence data. *Microbiome* *5*, 9.  
1563 <https://doi.org/10.1186/s40168-016-0224-8>.
- 1564 65. Battaglia, T.W., Mimpfen, I.L., Traets, J.J.H., Hoeck, A. van, Zeverijn, L.J., Geurts, B.S., Wit,  
1565 G.F. de, Noë, M., Hofland, I., Vos, J.L., et al. (2024). A pan-cancer analysis of the microbiome  
1566 in metastatic cancer. *Cell* *187*, 2324-2335.e19. <https://doi.org/10.1016/j.cell.2024.03.021>.
- 1567 66. Kostic, A.D., Ojesina, A.I., Peadarallu, C.S., Jung, J., Verhaak, R.G.W., Getz, G., and  
1568 Meyerson, M. (2011). PathSeq: software to identify or discover microbes by deep sequencing  
1569 of human tissue. *Nat Biotechnol* *29*, 393–396. <https://doi.org/10.1038/nbt.1868>.
- 1570 67. Walker, M.A., Peadarallu, C.S., Ojesina, A.I., Bullman, S., Sharpe, T., Whelan, C.W., and  
1571 Meyerson, M. (2018). GATK PathSeq: a customizable computational tool for the discovery  
1572 and identification of microbial sequences in libraries from eukaryotic hosts. *Bioinformatics* *34*,  
1573 4287–4289. <https://doi.org/10.1093/bioinformatics/bty501>.
- 1574 68. Feng, H., Shuda, M., Chang, Y., and Moore, P.S. (2008). Clonal integration of a polyomavirus  
1575 in human Merkel cell carcinoma. *Science* *319*, 1096–1100.  
1576 <https://doi.org/10.1126/science.1152586>.
- 1577 69. Guccione, C., Patel, L., Tomofuji, Y., McDonald, D., Gonzalez, A., Sepich-Poore, G.D.,  
1578 Sonehara, K., Zakeri, M., Chen, Y., Dilmore, A.H., et al. (2025). Incomplete human reference

- 1579 genomes can drive false sex biases and expose patient-identifying information in  
1580 metagenomic data. *Nat Commun* 16, 825. <https://doi.org/10.1038/s41467-025-56077-5>.
- 1581 70. Turnbull, C., Scott, R.H., Thomas, E., Jones, L., Murugaesu, N., Pretty, F.B., Halai, D., Baple,  
1582 E., Craig, C., Hamblin, A., et al. (2018). The 100 000 Genomes Project: bringing whole  
1583 genome sequencing to the NHS. *BMJ* 361, k1687. <https://doi.org/10.1136/bmj.k1687>.
- 1584 71. Turnbull, C. (2018). Introducing whole-genome sequencing into routine cancer care: the  
1585 Genomics England 100 000 Genomes Project. *Annals of Oncology* 29, 784–787.  
1586 <https://doi.org/10.1093/annonc/mdy054>.
- 1587 72. Kinnersley, B., Sud, A., Everall, A., Cornish, A.J., Chubb, D., Culliford, R., Gruber, A.J.,  
1588 Lärkeryd, A., Mitsopoulos, C., Wedge, D., et al. (2024). Analysis of 10,478 cancer genomes  
1589 identifies candidate driver genes and opportunities for precision oncology. *Nat Genet* 56,  
1590 1868–1877. <https://doi.org/10.1038/s41588-024-01785-9>.
- 1591 73. Cornish, A.J., Gruber, A.J., Kinnersley, B., Chubb, D., Frangou, A., Caravagna, G., Noyvert,  
1592 B., Lakatos, E., Wood, H.M., Thorn, S., et al. (2024). The genomic landscape of 2,023  
1593 colorectal cancers. *Nature* 633, 127–136. <https://doi.org/10.1038/s41586-024-07747-9>.
- 1594 74. National Genomic Research Library v5.1, Genomics England (2020).  
1595 [10.6084/m9.figshare.4530893/7](https://doi.org/10.6084/m9.figshare.4530893/7).
- 1596 75. Rhie, A., Nurk, S., Cechova, M., Hoyt, S.J., Taylor, D.J., Altemose, N., Hook, P.W., Koren,  
1597 S., Rautiainen, M., Alexandrov, I.A., et al. (2023). The complete sequence of a human Y  
1598 chromosome. *Nature* 621, 344–354. <https://doi.org/10.1038/s41586-023-06457-y>.
- 1599 76. Nurk, S., Koren, S., Rhie, A., Rautiainen, M., Bzikadze, A.V., Mikheenko, A., Vollger, M.R.,  
1600 Altemose, N., Uralsky, L., Gershman, A., et al. (2022). The complete sequence of a human  
1601 genome. *Science* 376, 44–53. <https://doi.org/10.1126/science.abj6987>.
- 1602 77. Wood, D.E., Lu, J., and Langmead, B. (2019). Improved metagenomic analysis with Kraken  
1603 2. *Genome Biol* 20, 257. <https://doi.org/10.1186/s13059-019-1891-0>.
- 1604 78. Blanco-Míguez, A., Beghini, F., Cumbo, F., McIver, L.J., Thompson, K.N., Zolfo, M., Manghi,  
1605 P., Dubois, L., Huang, K.D., Thomas, A.M., et al. (2023). Extending and improving  
1606 metagenomic taxonomic profiling with uncharacterized species using MetaPhlAn 4. *Nat*  
1607 *Biotechnol* 41, 1633–1644. <https://doi.org/10.1038/s41587-023-01688-w>.
- 1608 79. Shaw, J., and Yu, Y.W. (2025). Rapid species-level metagenome profiling and containment  
1609 estimation with sylph. *Nat Biotechnol* 43, 1348–1359. [https://doi.org/10.1038/s41587-024-](https://doi.org/10.1038/s41587-024-02412-y)  
1610 [02412-y](https://doi.org/10.1038/s41587-024-02412-y).
- 1611 80. Salter, S.J., Cox, M.J., Turek, E.M., Calus, S.T., Cookson, W.O., Moffatt, M.F., Turner, P.,  
1612 Parkhill, J., Loman, N.J., and Walker, A.W. (2014). Reagent and laboratory contamination can  
1613 critically impact sequence-based microbiome analyses. *BMC Biol* 12, 87.  
1614 <https://doi.org/10.1186/s12915-014-0087-z>.
- 1615 81. Eisenhofer, R., Minich, J.J., Marotz, C., Cooper, A., Knight, R., and Weyrich, L.S. (2019).  
1616 Contamination in Low Microbial Biomass Microbiome Studies: Issues and Recommendations.  
1617 *Trends Microbiol* 27, 105–117. <https://doi.org/10.1016/j.tim.2018.11.003>.

- 1618 82. Fierer, N., Leung, P.M., Lappan, R., Eisenhofer, R., Ricci, F., Holland, S.I., Dragone, N.,  
1619 Blackall, L.L., Dong, X., Dorador, C., et al. (2025). Guidelines for preventing and reporting  
1620 contamination in low-biomass microbiome studies. *Nat Microbiol* 10, 1570–1580.  
1621 <https://doi.org/10.1038/s41564-025-02035-2>.
- 1622 83. Effects of Index Misassignment on Multiplexing and Downstream Analysis.
- 1623 84. Larsson, A.J.M., Stanley, G., Sinha, R., Weissman, I.L., and Sandberg, R. (2018).  
1624 Computational correction of index switching in multiplexed sequencing libraries. *Nat Methods*  
1625 15, 305–307. <https://doi.org/10.1038/nmeth.4666>.
- 1626 85. Barretina, J., Caponigro, G., Stransky, N., Venkatesan, K., Margolin, A.A., Kim, S., Wilson,  
1627 C.J., Lehár, J., Kryukov, G.V., Sonkin, D., et al. (2012). The Cancer Cell Line Encyclopedia  
1628 enables predictive modeling of anticancer drug sensitivity. *Nature* 483, 603.  
1629 <https://doi.org/10.1038/nature11003>.
- 1630 86. Tan, C.C.S., Ko, K.K.K., Chen, H., Liu, J., Loh, M., SG10K\_Health Consortium, Chia, M., and  
1631 Nagarajan, N. (2023). No evidence for a common blood microbiome based on a population  
1632 study of 9,770 healthy humans. *Nat Microbiol* 8, 973–985. <https://doi.org/10.1038/s41564-023-01350-w>.
- 1634 87. Zhao, H., Chu, M., Huang, Z., Yang, X., Ran, S., Hu, B., Zhang, C., and Liang, J. (2017).  
1635 Variations in oral microbiota associated with oral cancer. *Sci Rep* 7, 11773.  
1636 <https://doi.org/10.1038/s41598-017-11779-9>.
- 1637 88. Hayes, R.B., Ahn, J., Fan, X., Peters, B.A., Ma, Y., Yang, L., Agalliu, I., Burk, R.D., Ganly, I.,  
1638 Purdue, M.P., et al. (2018). Association of Oral Microbiome With Risk for Incident Head and  
1639 Neck Squamous Cell Cancer. *JAMA Oncol* 4, 358–365.  
1640 <https://doi.org/10.1001/jamaoncol.2017.4777>.
- 1641 89. Böniggen, D., Ren, B., Pickard, R., Li, J., Ozer, E., Hartmann, E.M., Xiao, W., Tickle, T., Rider,  
1642 J., Gevers, D., et al. (2017). Alterations in oral bacterial communities are associated with risk  
1643 factors for oral and oropharyngeal cancer. *Sci Rep* 7, 17686. <https://doi.org/10.1038/s41598-017-17795-z>.
- 1645 90. Yang, W., Chen, C.-H., Jia, M., Xing, X., Gao, L., Tsai, H.-T., Zhang, Z., Liu, Z., Zeng, B.,  
1646 Yeung, S.-C.J., et al. (2021). Tumor-Associated Microbiota in Esophageal Squamous Cell  
1647 Carcinoma. *Front Cell Dev Biol* 9, 641270. <https://doi.org/10.3389/fcell.2021.641270>.
- 1648 91. Ferreira, R.M., Pereira-Marques, J., Pinto-Ribeiro, I., Costa, J.L., Carneiro, F., Machado, J.C.,  
1649 and Figueiredo, C. (2018). Gastric microbial community profiling reveals a dysbiotic cancer-  
1650 associated microbiota. *Gut* 67, 226–236. <https://doi.org/10.1136/gutjnl-2017-314205>.
- 1651 92. Booth, M.E., Wood, H.M., Travis, M.A., Genomics England Research Consortium, Quirke, P.,  
1652 and Grabsch, H.I. (2025). The relationship between the gastric cancer microbiome and  
1653 clinicopathological factors: a metagenomic investigation from the 100,000 genomes project  
1654 and The Cancer Genome Atlas. *Gastric Cancer* 28, 358–371. <https://doi.org/10.1007/s10120-025-01588-9>.
- 1656 93. Bender, M.J., McPherson, A.C., Phelps, C.M., Pandey, S.P., Laughlin, C.R., Shapira, J.H.,  
1657 Medina Sanchez, L., Rana, M., Richie, T.G., Mims, T.S., et al. (2023). Dietary tryptophan

- 1658 metabolite released by intratumoral *Lactobacillus reuteri* facilitates immune checkpoint  
1659 inhibitor treatment. *Cell* 186, 1846-1862.e26. <https://doi.org/10.1016/j.cell.2023.03.011>.
- 1660 94. He, Y., Li, L., Li, Y., Wang, X., Qian, L., Yang, J., and Jiang, M. (2025). Mendelian  
1661 randomization study reveals causal association between skin microbiome and skin cancers.  
1662 *Sci Rep* 15, 21590. <https://doi.org/10.1038/s41598-025-07265-2>.
- 1663 95. Krueger, A., Mohamed, A., Kolka, C.M., Stoll, T., Zaugg, J., Linedale, R., Morrison, M., Soyer,  
1664 H.P., Hugenholtz, P., Frazer, I.H., et al. (2022). Skin Cancer-Associated *S. aureus* Strains  
1665 Can Induce DNA Damage in Human Keratinocytes by Downregulating DNA Repair and  
1666 Promoting Oxidative Stress. *Cancers* 14, 2143. <https://doi.org/10.3390/cancers14092143>.
- 1667 96. Kullander, J., Forslund, O., and Dillner, J. (2009). Staphylococcus aureus and squamous cell  
1668 carcinoma of the skin. *Cancer Epidemiol Biomarkers Prev* 18, 472–478.  
1669 <https://doi.org/10.1158/1055-9965.EPI-08-0905>.
- 1670 97. Krueger, A., Zaugg, J., Chisholm, S., Linedale, R., Lachner, N., Teoh, S.M., Tuong, Z.K.,  
1671 Lukowski, S.W., Morrison, M., Soyer, H.P., et al. (2022). Secreted Toxins From  
1672 Staphylococcus aureus Strains Isolated From Keratinocyte Skin Cancers Mediate Pro-  
1673 tumorigenic Inflammatory Responses in the Skin. *Front. Microbiol.* 12.  
1674 <https://doi.org/10.3389/fmicb.2021.789042>.
- 1675 98. Eckhoff, A.M., Fletcher, A.A., Kelly, M.S., Dohman, A., McIntyre, C.A., Shen, X., Iyer, M.K.,  
1676 Nussbaum, D.P., and Allen, P.J. (2024). Comprehensive Assessment of the Intrinsic  
1677 Pancreatic Microbiome. *Ann Surg.* <https://doi.org/10.1097/SLA.0000000000006299>.
- 1678 99. Pushalkar, S., Hundeyin, M., Daley, D., Zambirinis, C.P., Kurz, E., Mishra, A., Mohan, N.,  
1679 Aykut, B., Usyk, M., Torres, L.E., et al. (2018). The Pancreatic Cancer Microbiome Promotes  
1680 Oncogenesis by Induction of Innate and Adaptive Immune Suppression. *Cancer Discov* 8,  
1681 403–416. <https://doi.org/10.1158/2159-8290.CD-17-1134>.
- 1682 100. Ginwala, R., Bukavina, L., Sindhani, M., Nachman, E., Peri, S., Franklin, J., Drevik, J.,  
1683 Christianson, S., Geynisman, D.M., Kutikov, A., et al. (2025). Bladder cancer microbiome and  
1684 its association with chemoresponse. *Front Oncol* 15, 1506319.  
1685 <https://doi.org/10.3389/fonc.2025.1506319>.
- 1686 101. Bučević Popović, V., Šitum, M., Chow, C.-E.T., Chan, L.S., Roje, B., and Terzić, J. (2018).  
1687 The urinary microbiome associated with bladder cancer. *Sci Rep* 8, 12157.  
1688 <https://doi.org/10.1038/s41598-018-29054-w>.
- 1689 102. Zhang, Y., Lin, H., Liang, L., Jin, S., Lv, J., Zhou, Y., Xu, F., Liu, F., and Feng, N. (2024).  
1690 Intratumoral microbiota as a novel prognostic indicator in bladder cancer. *Sci Rep* 14, 22198.  
1691 <https://doi.org/10.1038/s41598-024-72918-7>.
- 1692 103. Wolfe, A.J., and Brubaker, L. (2019). Urobiome Updates: Advances in Urinary Microbiome  
1693 Research. *Nat Rev Urol* 16, 73–74. <https://doi.org/10.1038/s41585-018-0127-5>.
- 1694 104. Lotte, R., Lotte, L., and Ruimy, R. (2016). *Actinotignum schaalii* (formerly *Actinobaculum*  
1695 *schaalii*): a newly recognized pathogen—review of the literature. *Clinical Microbiology and*  
1696 *Infection* 22, 28–36. <https://doi.org/10.1016/j.cmi.2015.10.038>.

- 1697 105. Torres, A., Cilloniz, C., Niederman, M.S., Menéndez, R., Chalmers, J.D., Wunderink, R.G.,  
1698 and van der Poll, T. (2021). Pneumonia. *Nat Rev Dis Primers* 7, 25.  
1699 <https://doi.org/10.1038/s41572-021-00259-0>.
- 1700 106. Langenberg, A.G.M., Corey, L., Ashley, R.L., Leong, W.P., and Straus, S.E. (1999). A  
1701 Prospective Study of New Infections with Herpes Simplex Virus Type 1 and Type 2. *New*  
1702 *England Journal of Medicine* 341, 1432–1438.  
1703 <https://doi.org/10.1056/NEJM199911043411904>.
- 1704 107. Young, L.S., Yap, L.F., and Murray, P.G. (2016). Epstein–Barr virus: more than 50 years  
1705 old and still providing surprises. *Nat Rev Cancer* 16, 789–802.  
1706 <https://doi.org/10.1038/nrc.2016.92>.
- 1707 108. Griffiths, P., and Reeves, M. (2021). Pathogenesis of human cytomegalovirus in the  
1708 immunocompromised host. *Nat Rev Microbiol* 19, 759–773. [https://doi.org/10.1038/s41579-](https://doi.org/10.1038/s41579-021-00582-z)  
1709 [021-00582-z](https://doi.org/10.1038/s41579-021-00582-z).
- 1710 109. Zerr, D.M., Meier, A.S., Selke, S.S., Frenkel, L.M., Huang, M.-L., Wald, A., Rhoads, M.P.,  
1711 Nguy, L., Bornemann, R., Morrow, R.A., et al. (2005). A Population-Based Study of Primary  
1712 Human Herpesvirus 6 Infection. *New England Journal of Medicine* 352, 768–776.  
1713 <https://doi.org/10.1056/NEJMoa042207>.
- 1714 110. Verbeek, R., Vandekerckhove, L., and Van Cleemput, J. Update on human herpesvirus 7  
1715 pathogenesis and clinical aspects as a roadmap for future research. *J Virol* 98, e00437-24.  
1716 <https://doi.org/10.1128/jvi.00437-24>.
- 1717 111. Galeano Niño, J.L., Wu, H., LaCourse, K.D., Kempchinsky, A.G., Baryames, A., Barber,  
1718 B., Futran, N., Houlton, J., Sather, C., Sicinska, E., et al. (2022). Effect of the intratumoral  
1719 microbiota on spatial and cellular heterogeneity in cancer. *Nature* 611, 810–817.  
1720 <https://doi.org/10.1038/s41586-022-05435-0>.
- 1721 112. Queen, J., Cing, Z., Minsky, H., Nandi, A., Southward, T., Ferri, J., McMann, M., Iyadorai,  
1722 T., Vadivelu, J., Roslani, A., et al. (2025). *Fusobacterium nucleatum* is enriched in invasive  
1723 biofilms in colorectal cancer. *npj Biofilms Microbiomes* 11, 81. [https://doi.org/10.1038/s41522-](https://doi.org/10.1038/s41522-025-00717-7)  
1724 [025-00717-7](https://doi.org/10.1038/s41522-025-00717-7).
- 1725 113. Socransky, S.S., Haffajee, A.D., Cugini, M.A., Smith, C., and Kent, R.L. (1998). Microbial  
1726 complexes in subgingival plaque. *J Clin Periodontol* 25, 134–144.  
1727 <https://doi.org/10.1111/j.1600-051x.1998.tb02419.x>.
- 1728 114. Bao, Y., Zhang, Y., Zhang, Y., Liu, Y., Wang, S., Dong, X., Wang, Y., and Zhang, H.  
1729 (2010). Screening of potential probiotic properties of *Lactobacillus fermentum* isolated from  
1730 traditional dairy products. *Food Control* 21, 695–701.  
1731 <https://doi.org/10.1016/j.foodcont.2009.10.010>.
- 1732 115. Selle, K., and Klaenhammer, T.R. (2013). Genomic and phenotypic evidence for probiotic  
1733 influences of *Lactobacillus gasseri* on human health. *FEMS Microbiol Rev* 37, 915–935.  
1734 <https://doi.org/10.1111/1574-6976.12021>.
- 1735 116. Fu, K., Cheung, A.H.K., Wong, C.C., Liu, W., Zhou, Y., Wang, F., Huang, P., Yuan, K.,  
1736 Coker, O.O., Pan, Y., et al. (2024). *Streptococcus anginosus* promotes gastric inflammation,

- 1737 atrophy, and tumorigenesis in mice. *Cell* 187, 882-896.e17.  
1738 <https://doi.org/10.1016/j.cell.2024.01.004>.
- 1739 117. Fried, B., Reddy, A., and Mayer, D. (2011). Helminths in human carcinogenesis. *Cancer*  
1740 *Letters* 305, 239–249. <https://doi.org/10.1016/j.canlet.2010.07.008>.
- 1741 118. Miller, T.L., and Wolin, M.J. (1982). Enumeration of *Methanobrevibacter smithii* in human  
1742 feces. *Arch Microbiol* 131, 14–18. <https://doi.org/10.1007/BF00451492>.
- 1743 119. Borrel, G., Brugère, J.-F., Gribaldo, S., Schmitz, R.A., and Moissl-Eichinger, C. (2020).  
1744 The host-associated archaeome. *Nat Rev Microbiol* 18, 622–636.  
1745 <https://doi.org/10.1038/s41579-020-0407-y>.
- 1746 120. Triantafyllou, K., Chang, C., and Pimentel, M. (2014). Methanogens, methane and  
1747 gastrointestinal motility. *J Neurogastroenterol Motil* 20, 31–40.  
1748 <https://doi.org/10.5056/jnm.2014.20.1.31>.
- 1749 121. Low, A., Lee, J.K.Y., Gounot, J.-S., Ravikrishnan, A., Ding, Y., Saw, W.-Y., Tan, L.W.L.,  
1750 Moong, D.K.N., Teo, Y.Y., Nagarajan, N., et al. (2022). Mutual Exclusion of  
1751 *Methanobrevibacter* Species in the Human Gut Microbiota Facilitates Directed Cultivation of  
1752 a Candidatus *Methanobrevibacter Intestini* Representative. *Microbiology Spectrum* 10,  
1753 e00849-22. <https://doi.org/10.1128/spectrum.00849-22>.
- 1754 122. Borrel, G., Harris, H.M.B., Parisot, N., Gaci, N., Tottey, W., Mihajlovski, A., Deane, J.,  
1755 Gribaldo, S., Bardot, O., Peyretailade, E., et al. (2013). Genome Sequence of “Candidatus  
1756 *Methanomassiliicoccus intestinalis*” Issoire-Mx1, a Third Thermoplasmatales-Related  
1757 Methanogenic Archaeon from Human Feces. *Genome Announcements* 1,  
1758 10.1128/genomea.00453-13. <https://doi.org/10.1128/genomea.00453-13>.
- 1759 123. Poole, D.N., and McClelland, R.S. (2013). Global epidemiology of *Trichomonas vaginalis*.  
1760 *Sex Transm Infect* 89, 418–422. <https://doi.org/10.1136/sextrans-2013-051075>.
- 1761 124. Hamar, B., Teutsch, B., Hoffmann, E., Hegyi, P., Váradi, A., Nyirády, P., Hunka, Z., Ács,  
1762 N., Lintner, B., Hermáné, R.J., et al. (2023). *Trichomonas vaginalis* infection is associated  
1763 with increased risk of cervical carcinogenesis: A systematic review and meta-analysis of 470  
1764 000 patients. *International Journal of Gynecology & Obstetrics* 163, 31–43.  
1765 <https://doi.org/10.1002/ijgo.14763>.
- 1766 125. Gonelli, A., Boccia, S., Boni, M., Pozzoli, A., Rizzo, C., Querzoli, P., Cassai, E., and Di  
1767 Luca, D. (2001). Human herpesvirus 7 is latent in gastric mucosa. *J Med Virol* 63, 277–283.  
1768 [https://doi.org/10.1002/1096-9071\(200104\)63:4%253C277::aid-jmv1002%253E3.0.co;2-k](https://doi.org/10.1002/1096-9071(200104)63:4%253C277::aid-jmv1002%253E3.0.co;2-k).
- 1769 126. Chiu, H.H., Lee, C.Y., Lee, P.I., Lin, K.H., and Huang, L.M. (1998). Mononucleosis  
1770 syndrome and coincidental human herpesvirus-7 and Epstein-Barr virus infection. *Arch Dis*  
1771 *Child* 78, 479–480. <https://doi.org/10.1136/adc.78.5.479>.
- 1772 127. Piccinno, G., Thompson, K.N., Manghi, P., Ghazi, A.R., Thomas, A.M., Blanco-Míguez,  
1773 A., Asnicar, F., Mladenovic, K., Pinto, F., Armanini, F., et al. (2025). Pooled analysis of 3,741  
1774 stool metagenomes from 18 cohorts for cross-stage and strain-level reproducible microbial  
1775 biomarkers of colorectal cancer. *Nat Med*, 1–14. [https://doi.org/10.1038/s41591-025-03693-](https://doi.org/10.1038/s41591-025-03693-9)  
1776 9.

- 1777 128. Flemer, B., Lynch, D.B., Brown, J.M.R., Jeffery, I.B., Ryan, F.J., Claesson, M.J.,  
1778 O'Riordain, M., Shanahan, F., and O'Toole, P.W. (2017). Tumour-associated and non-  
1779 tumour-associated microbiota in colorectal cancer. [https://doi.org/10.1136/gutjnl-2015-](https://doi.org/10.1136/gutjnl-2015-309595)  
1780 309595.
- 1781 129. Mima, K., Nishihara, R., Qian, Z.R., Cao, Y., Sukawa, Y., Nowak, J.A., Yang, J., Dou, R.,  
1782 Masugi, Y., Song, M., et al. (2016). *Fusobacterium nucleatum* in colorectal carcinoma tissue  
1783 and patient prognosis. *Gut* 65, 1973–1980. <https://doi.org/10.1136/gutjnl-2015-310101>.
- 1784 130. Hamada, T., Zhang, X., Mima, K., Bullman, S., Sukawa, Y., Nowak, J.A., Kosumi, K.,  
1785 Masugi, Y., Twombly, T.S., Cao, Y., et al. (2018). *Fusobacterium nucleatum* in Colorectal  
1786 Cancer Relates to Immune Response Differentially by Tumor Microsatellite Instability Status.  
1787 *Cancer Immunol Res* 6, 1327–1336. <https://doi.org/10.1158/2326-6066.CIR-18-0174>.
- 1788 131. Muzny, D.M., Bainbridge, M.N., Chang, K., Dinh, H.H., Drummond, J.A., Fowler, G.,  
1789 Kovar, C.L., Lewis, L.R., Morgan, M.B., Newsham, I.F., et al. (2012). Comprehensive  
1790 molecular characterization of human colon and rectal cancer. *Nature* 487, 330–337.  
1791 <https://doi.org/10.1038/nature11252>.
- 1792 132. Dienstmann, R., Vermeulen, L., Guinney, J., Kopetz, S., Tejpar, S., and Tabernero, J.  
1793 (2017). Consensus molecular subtypes and the evolution of precision medicine in colorectal  
1794 cancer. *Nat Rev Cancer* 17, 79–92. <https://doi.org/10.1038/nrc.2016.126>.
- 1795 133. Guinney, J., Dienstmann, R., Wang, X., de Reyniès, A., Schlicker, A., Soneson, C.,  
1796 Marisa, L., Roepman, P., Nyamundanda, G., Angelino, P., et al. (2015). The consensus  
1797 molecular subtypes of colorectal cancer. *Nat Med* 21, 1350–1356.  
1798 <https://doi.org/10.1038/nm.3967>.
- 1799 134. Takashima, Y., Kawamura, H., Okadome, K., Ugai, S., Haruki, K., Arima, K., Mima, K.,  
1800 Akimoto, N., Nowak, J.A., Giannakis, M., et al. (2024). Enrichment of *Bacteroides fragilis* and  
1801 enterotoxigenic *Bacteroides fragilis* in CpG island methylator phenotype-high colorectal  
1802 carcinoma. *Clin Microbiol Infect* 30, 630–636. <https://doi.org/10.1016/j.cmi.2024.01.013>.
- 1803 135. Zepeda-Rivera, M., Minot, S.S., Bouzek, H., Wu, H., Blanco-Míguez, A., Manghi, P.,  
1804 Jones, D.S., LaCourse, K.D., Wu, Y., McMahon, E.F., et al. (2024). A distinct *Fusobacterium*  
1805 *nucleatum* clade dominates the colorectal cancer niche. *Nature* 628, 424–432.  
1806 <https://doi.org/10.1038/s41586-024-07182-w>.
- 1807 136. Wang, C., Xiao, Y., Yu, L., Tian, F., Zhao, J., Zhang, H., Chen, W., and Zhai, Q. (2021).  
1808 Protective effects of different *Bacteroides vulgatus* strains against lipopolysaccharide-induced  
1809 acute intestinal injury, and their underlying functional genes. *J Adv Res* 36, 27–37.  
1810 <https://doi.org/10.1016/j.jare.2021.06.012>.
- 1811 137. Sokol, H., Pigneur, B., Watterlot, L., Lakhdari, O., Bermúdez-Humarán, L.G., Gratadoux,  
1812 J.-J., Blugeon, S., Bridonneau, C., Furet, J.-P., Corthier, G., et al. (2008). *Faecalibacterium*  
1813 *prausnitzii* is an anti-inflammatory commensal bacterium identified by gut microbiota analysis  
1814 of Crohn disease patients. *Proc Natl Acad Sci U S A* 105, 16731–16736.  
1815 <https://doi.org/10.1073/pnas.0804812105>.
- 1816 138. Fabersani, E., Portune, K., Campillo, I., López-Almela, I., la Paz, S.M., Romani-Pérez, M.,  
1817 Benítez-Páez, A., and Sanz, Y. (2021). *Bacteroides uniformis* CECT 7771 alleviates

- 1818 inflammation within the gut-adipose tissue axis involving TLR5 signaling in obese mice. *Sci*  
1819 *Rep* 11, 11788. <https://doi.org/10.1038/s41598-021-90888-y>.
- 1820 139. Siegel, R.L., Torre, L.A., Soerjomataram, I., Hayes, R.B., Bray, F., Weber, T.K., and  
1821 Jemal, A. (2019). Global patterns and trends in colorectal cancer incidence in young adults.  
1822 *Gut* 68, 2179–2185. <https://doi.org/10.1136/gutjnl-2019-319511>.
- 1823 140. Vuik, F.E., Nieuwenburg, S.A., Bardou, M., Lansdorp-Vogelaar, I., Dinis-Ribeiro, M.,  
1824 Bento, M.J., Zadnik, V., Pellisé, M., Esteban, L., Kaminski, M.F., et al. (2019). Increasing  
1825 incidence of colorectal cancer in young adults in Europe over the last 25 years.  
1826 <https://doi.org/10.1136/gutjnl-2018-317592>.
- 1827 141. Sung, H., Siegel, R.L., Laversanne, M., Jiang, C., Morgan, E., Zahwe, M., Cao, Y., Bray,  
1828 F., and Jemal, A. (2025). Colorectal cancer incidence trends in younger versus older adults:  
1829 an analysis of population-based cancer registry data. *The Lancet Oncology* 26, 51–63.  
1830 [https://doi.org/10.1016/S1470-2045\(24\)00600-4](https://doi.org/10.1016/S1470-2045(24)00600-4).
- 1831 142. Díaz-Gay, M., dos Santos, W., Moody, S., Kazachkova, M., Abbasi, A., Steele, C.D.,  
1832 Vangara, R., Senkin, S., Wang, J., Fitzgerald, S., et al. (2025). Geographic and age variations  
1833 in mutational processes in colorectal cancer. *Nature*, 1–11. [https://doi.org/10.1038/s41586-](https://doi.org/10.1038/s41586-025-09025-8)  
1834 [025-09025-8](https://doi.org/10.1038/s41586-025-09025-8).
- 1835 143. Sinicrope, F.A. (2022). Increasing Incidence of Early-Onset Colorectal Cancer. *New*  
1836 *England Journal of Medicine* 386, 1547–1558. <https://doi.org/10.1056/NEJMra2200869>.
- 1837 144. Martini, A.M., Moricz, B.S., Ripperger, A.K., Tran, P.M., Sharp, M.E., Forsythe, A.N.,  
1838 Kulhankova, K., Salgado-Pabón, W., and Jones, B.D. (2020). Association of Novel  
1839 *Streptococcus sanguinis* Virulence Factors With Pathogenesis in a Native Valve Infective  
1840 Endocarditis Model. *Front Microbiol* 11, 10. <https://doi.org/10.3389/fmicb.2020.00010>.
- 1841 145. Nijjer, S., and Dubrey, S.W. (2010). *Streptococcus sanguis* endocarditis associated with  
1842 colonic carcinoma. *BMJ Case Rep* 2010, bcr09.2009.2311.  
1843 <https://doi.org/10.1136/bcr.09.2009.2311>.
- 1844 146. Macaluso, A., Simmang, C., and Anthony, T. (1998). *Streptococcus sanguis* bacteremia  
1845 and colorectal cancer. *South Med J* 91, 206–207. [https://doi.org/10.1097/00007611-](https://doi.org/10.1097/00007611-199802000-00016)  
1846 [199802000-00016](https://doi.org/10.1097/00007611-199802000-00016).
- 1847 147. Thomas, R., Gupta, V., and Kwan, B. (2018). Second look at *Streptococcus sanguinis* and  
1848 the colon. *BMJ Case Rep* 2018, bcr2018224799. <https://doi.org/10.1136/bcr-2018-224799>.
- 1849 148. Association of *Streptococcus sanguinis* Infection With Colorectal Carcinoma (2020).  
1850 Consultant360. [https://www.consultant360.com/article/infectious-diseases/bacterial-](https://www.consultant360.com/article/infectious-diseases/bacterial-infections/association-streptococcus-sanguinis-infection)  
1851 [infections/association-streptococcus-sanguinis-infection](https://www.consultant360.com/article/infectious-diseases/bacterial-infections/association-streptococcus-sanguinis-infection).
- 1852 149. Rawat, P.S., Li, Y., Zhang, W., Meng, X., and Liu, W. (2022). *Hungatella hathewayi*, an  
1853 Efficient Glycosaminoglycan-Degrading Firmicutes from Human Gut and Its Chondroitin ABC  
1854 Exolyase with High Activity and Broad Substrate Specificity. *Appl Environ Microbiol* 88,  
1855 e0154622. <https://doi.org/10.1128/aem.01546-22>.

- 1856 150. Ottman, N., Reunanen, J., Meijerink, M., Pietilä, T.E., Kainulainen, V., Klievink, J.,  
1857 Huuskonen, L., Aalvink, S., Skurnik, M., Boeren, S., et al. (2017). Pili-like proteins of  
1858 *Akkermansia muciniphila* modulate host immune responses and gut barrier function. *PLOS*  
1859 *ONE* 12, e0173004. <https://doi.org/10.1371/journal.pone.0173004>.
- 1860 151. Chen, S., Zhang, Z., Liu, S., Chen, T., Lu, Z., Zhao, W., Mou, X., and Liu, S. (2024).  
1861 Consistent signatures in the human gut microbiome of longevous populations. *Gut Microbes*  
1862 16, 2393756. <https://doi.org/10.1080/19490976.2024.2393756>.
- 1863 152. Li, Y., Ma, A., Johnson, E., Eng, C., De, S., Jiang, S., Li, Z., Spakowicz, D., and Ma, Q.  
1864 (2025). The new microbiome on the block: challenges and opportunities of using human tumor  
1865 sequencing data to study microbes. *Nat Methods* 22, 1788–1799.  
1866 <https://doi.org/10.1038/s41592-025-02807-y>.
- 1867 153. Matsushita, H., Vesely, M.D., Koboldt, D.C., Rickert, C.G., Uppaluri, R., Magrini, V.J.,  
1868 Arthur, C.D., White, J.M., Chen, Y.-S., Shea, L.K., et al. (2012). Cancer exome analysis  
1869 reveals a T-cell-dependent mechanism of cancer immunoediting. *Nature* 482, 400–404.  
1870 <https://doi.org/10.1038/nature10755>.
- 1871 154. Gubin, M.M., Zhang, X., Schuster, H., Caron, E., Ward, J.P., Noguchi, T., Ivanova, Y.,  
1872 Hundal, J., Arthur, C.D., Krebber, W.-J., et al. (2014). Checkpoint blockade cancer  
1873 immunotherapy targets tumour-specific mutant antigens. *Nature* 515, 577–581.  
1874 <https://doi.org/10.1038/nature13988>.
- 1875 155. McGranahan, N., Furness, A.J.S., Rosenthal, R., Ramskov, S., Lyngaa, R., Saini, S.K.,  
1876 Jamal-Hanjani, M., Wilson, G.A., Birkbak, N.J., Hiley, C.T., et al. (2016). Clonal neoantigens  
1877 elicit T cell immunoreactivity and sensitivity to immune checkpoint blockade. *Science* 351,  
1878 1463–1469. <https://doi.org/10.1126/science.aaf1490>.
- 1879 156. Le, D.T., Uram, J.N., Wang, H., Bartlett, B.R., Kemberling, H., Eyring, A.D., Skora, A.D.,  
1880 Luber, B.S., Azad, N.S., Laheru, D., et al. (2015). PD-1 Blockade in Tumors with Mismatch-  
1881 Repair Deficiency. *New England Journal of Medicine* 372, 2509–2520.  
1882 <https://doi.org/10.1056/NEJMoa1500596>.
- 1883 157. Rizvi, N.A., Hellmann, M.D., Snyder, A., Kvistborg, P., Makarov, V., Havel, J.J., Lee, W.,  
1884 Yuan, J., Wong, P., Ho, T.S., et al. (2015). Mutational landscape determines sensitivity to PD-  
1885 1 blockade in non-small cell lung cancer. *Science* 348, 124–128.  
1886 <https://doi.org/10.1126/science.aaa1348>.
- 1887 158. Gur, C., Ibrahim, Y., Isaacson, B., Yamin, R., Abed, J., Gamliel, M., Enk, J., Bar-On, Y.,  
1888 Stanietsky-Kaynan, N., Copenhagen-Glazer, S., et al. (2015). Binding of the Fap2 protein of  
1889 *Fusobacterium nucleatum* to human inhibitory receptor TIGIT protects tumors from immune  
1890 cell attack. *Immunity* 42, 344–355. <https://doi.org/10.1016/j.immuni.2015.01.010>.
- 1891 159. MEGAHIT: an ultra-fast single-node solution for large and complex metagenomics  
1892 assembly via succinct de Bruijn graph | *Bioinformatics* | Oxford Academic  
1893 <https://academic.oup.com/bioinformatics/article/31/10/1674/177884>.
- 1894 160. Gourel, H., Karlsson-Lindsjö, O., Hayer, J., and Bongcam-Rudloff, E. (2019). Simulating  
1895 Illumina metagenomic data with InSilicoSeq. *Bioinformatics* 35, 521–522.  
1896 <https://doi.org/10.1093/bioinformatics/bty630>.

1897 161. Bairoch, A. (2018). The Cellosaurus, a Cell-Line Knowledge Resource. *Journal of*  
1898 *Biomolecular Techniques* : JBT 29, 25. <https://doi.org/10.7171/jbt.18-2902-002>.

1899 162. Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis,  
1900 G., Durbin, R., and 1000 Genome Project Data Processing Subgroup (2009). The Sequence  
1901 Alignment/Map format and SAMtools. *Bioinformatics* 25, 2078–2079.  
1902 <https://doi.org/10.1093/bioinformatics/btp352>.

1903 163. Picard toolkit (2019).

1904 164. Wright, R.J., Comeau, A.M., and Langille, M.G.I. (2023). From defaults to databases:  
1905 parameter and database choice dramatically impact the performance of metagenomic  
1906 taxonomic classification tools. *Microbial Genomics* 9, 000949.  
1907 <https://doi.org/10.1099/mgen.0.000949>.

1908 165. Davis, N.M., Proctor, D.M., Holmes, S.P., Relman, D.A., and Callahan, B.J. (2018). Simple  
1909 statistical identification and removal of contaminant sequences in marker-gene and  
1910 metagenomics data. *Microbiome* 6, 226. <https://doi.org/10.1186/s40168-018-0605-2>.

1911 166. Langmead, B., and Salzberg, S.L. (2012). Fast gapped-read alignment with Bowtie 2. *Nat*  
1912 *Methods* 9, 357–359. <https://doi.org/10.1038/nmeth.1923>.

1913 167. Dohlman, A.B., Arguijo Mendoza, D., Ding, S., Gao, M., Dressman, H., Iliev, I.D., Lipkin,  
1914 S.M., and Shen, X. (2021). The cancer microbiome atlas: a pan-cancer comparative analysis  
1915 to distinguish tissue-resident microbiota from contaminants. *Cell Host Microbe* 29, 281-  
1916 298.e5. <https://doi.org/10.1016/j.chom.2020.12.001>.

1917 168. Weinstein, J.N., Collisson, E.A., Mills, G.B., Shaw, K.R.M., Ozenberger, B.A., Ellrott, K.,  
1918 Shmulevich, I., Sander, C., and Stuart, J.M. (2013). The Cancer Genome Atlas Pan-Cancer  
1919 analysis project. *Nat Genet* 45, 1113–1120. <https://doi.org/10.1038/ng.2764>.

1920 169. Kautto, E.A., Bonneville, R., Miya, J., Yu, L., Krook, M.A., Reeser, J.W., and  
1921 Roychowdhury, S. (2016). Performance evaluation for rapid detection of pan-cancer  
1922 microsatellite instability with MANTIS. *Oncotarget* 8, 7452.  
1923 <https://doi.org/10.18632/oncotarget.13918>.

1924

1925

1926

1927

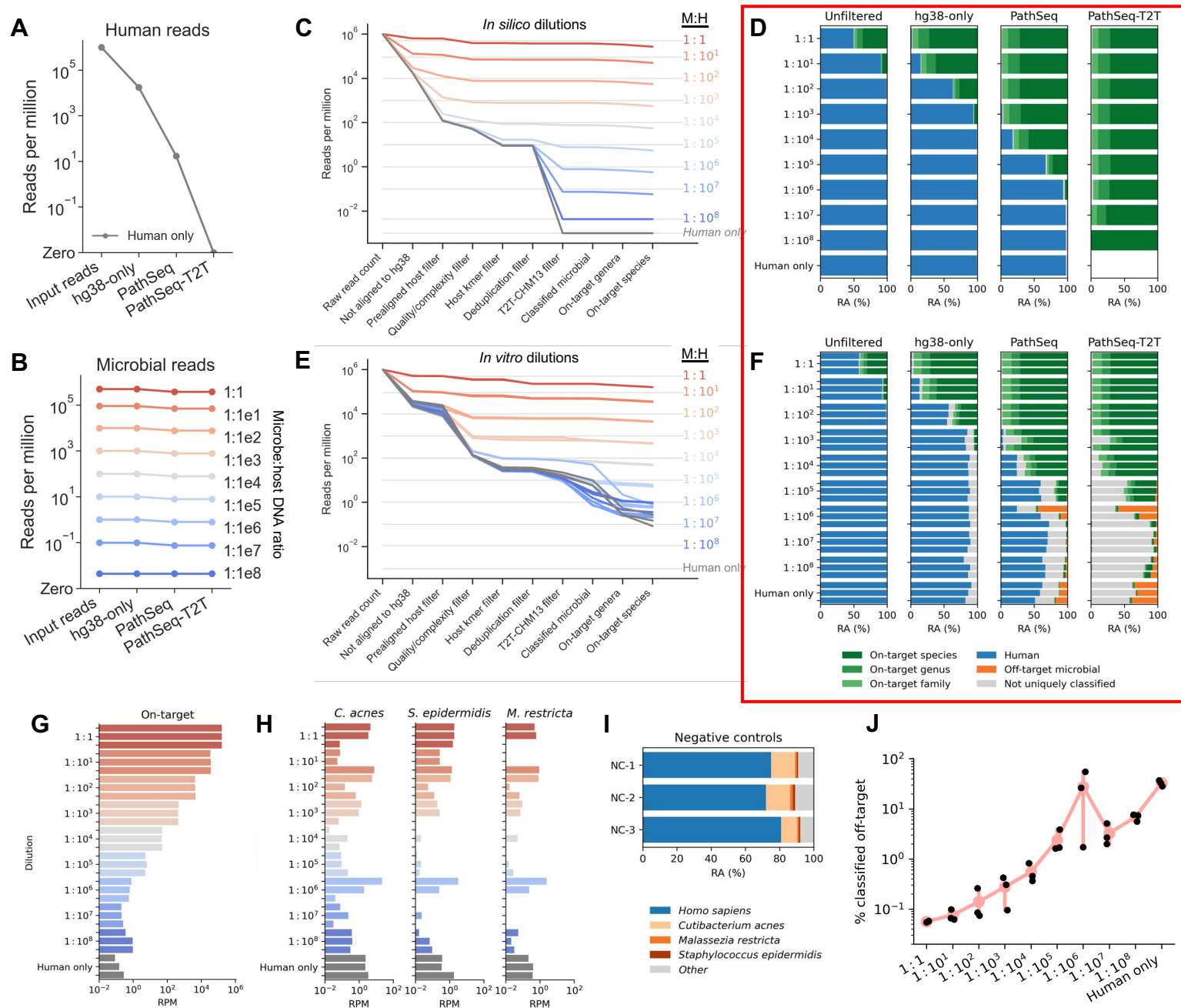
1928

1929

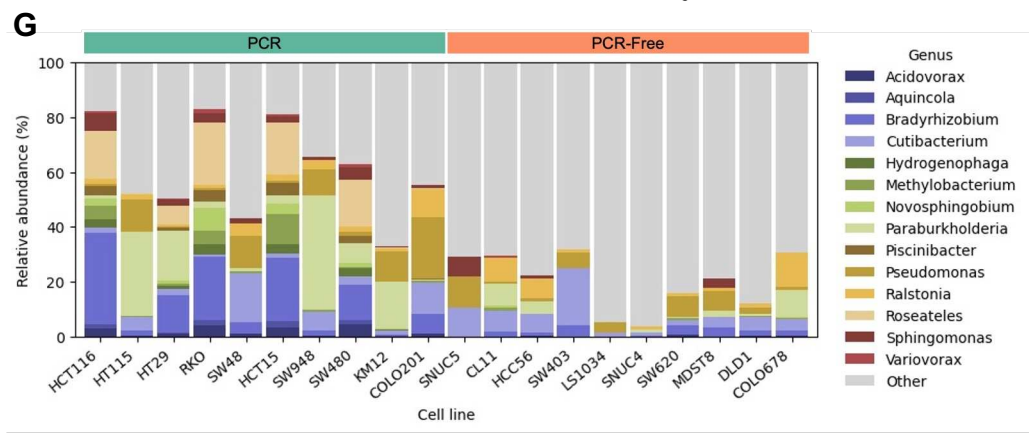
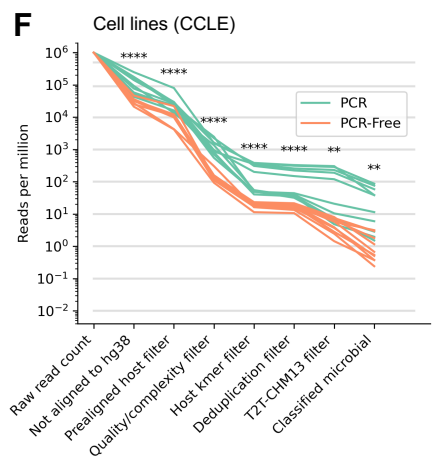
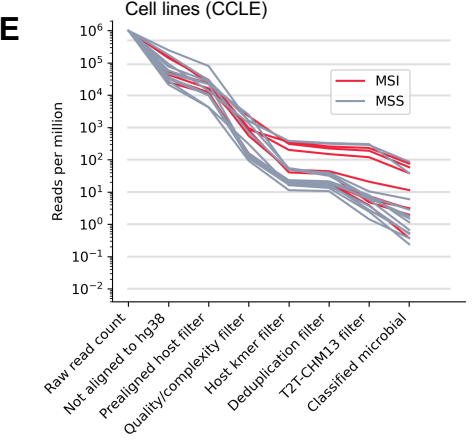
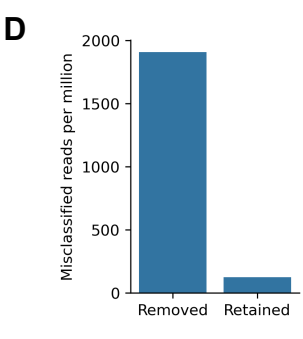
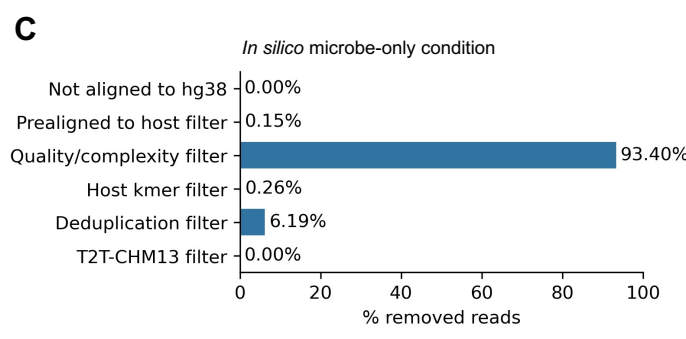
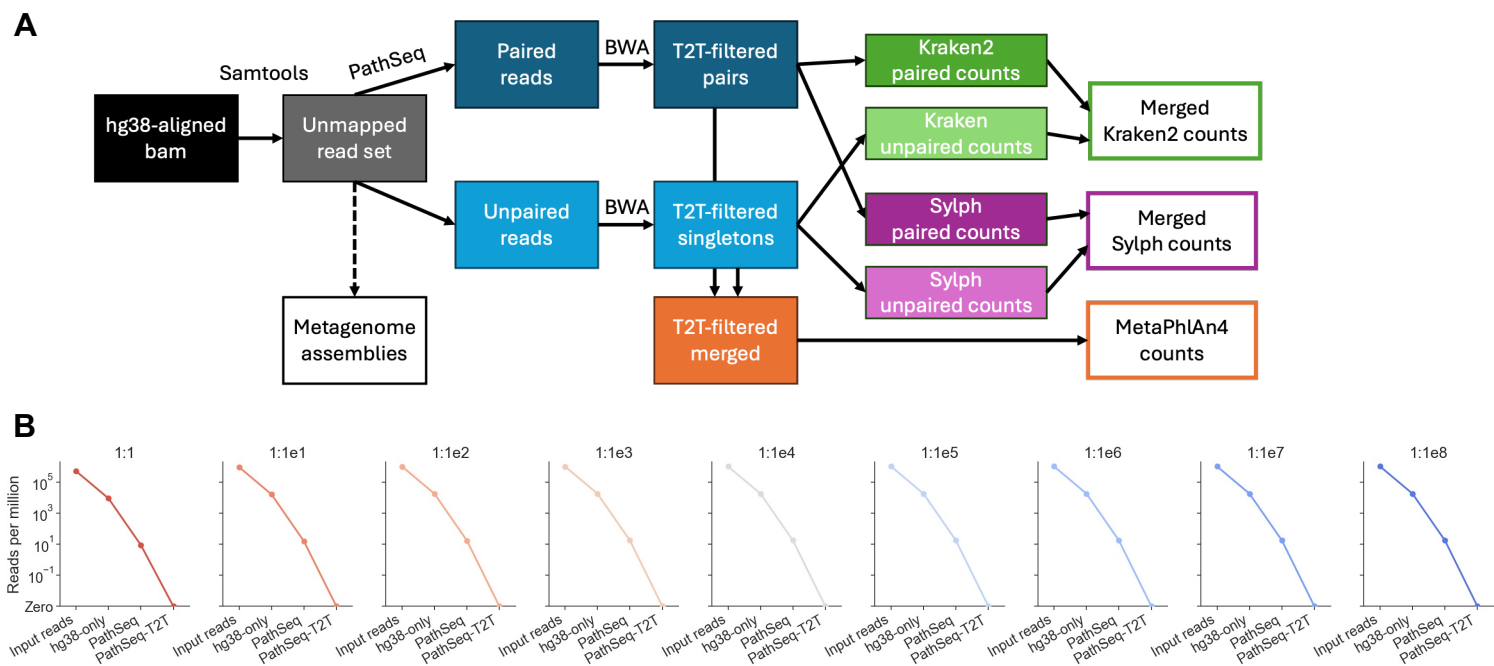
1930

1931

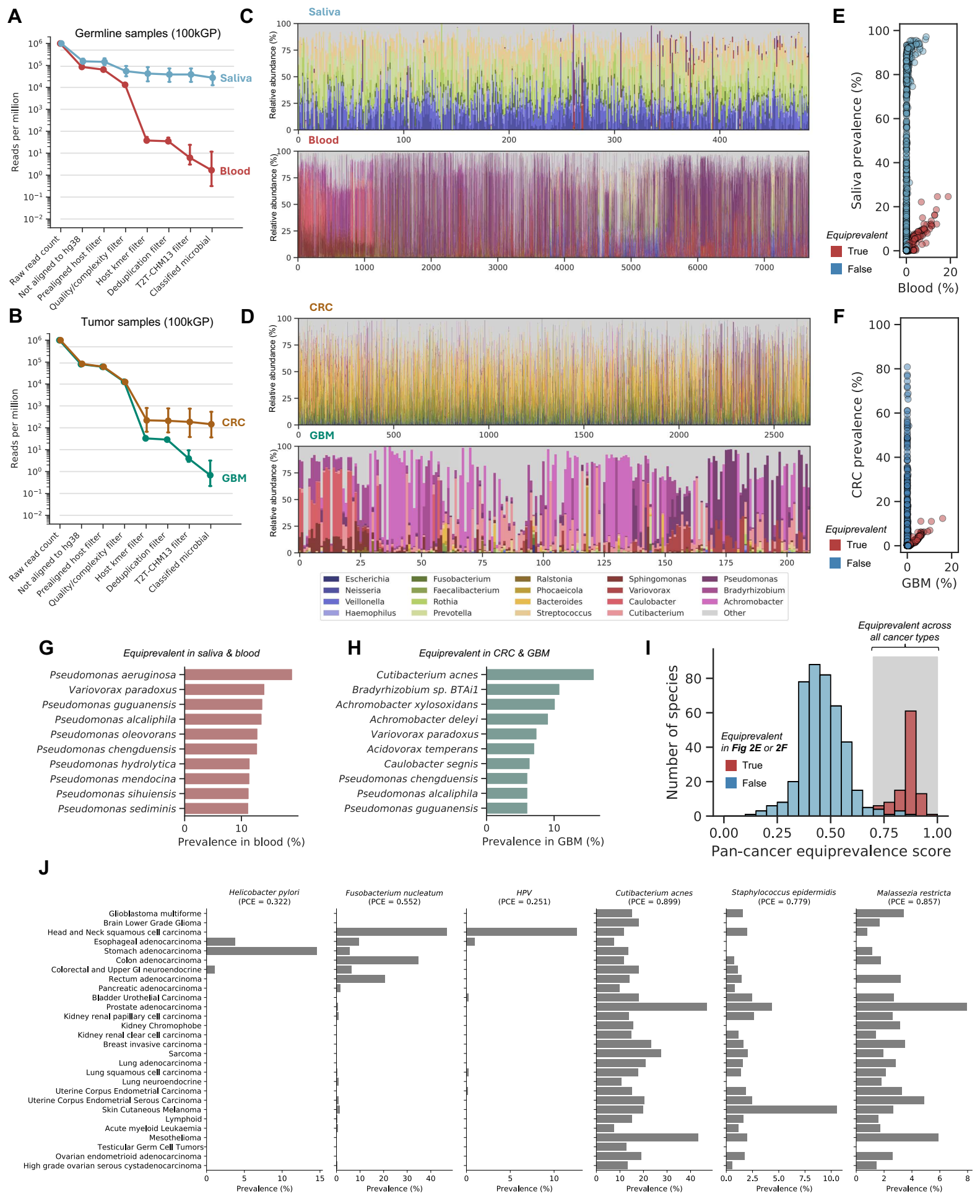
**Figure 1** A host subtraction pipeline for detecting microorganisms in low-biomass human tissues



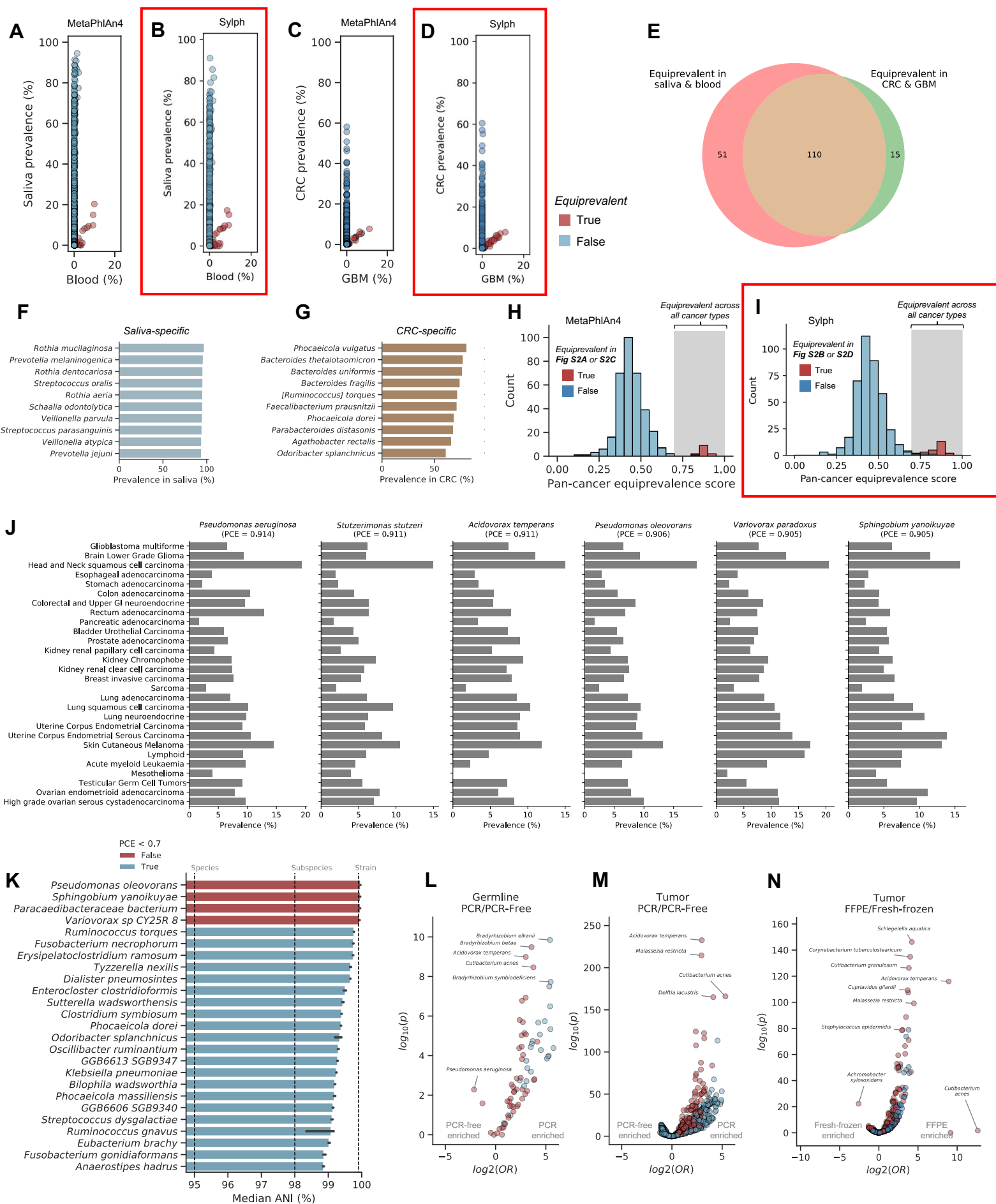
**Figure S1** A host subtraction pipeline for detecting microorganisms in low-biomass human tissues



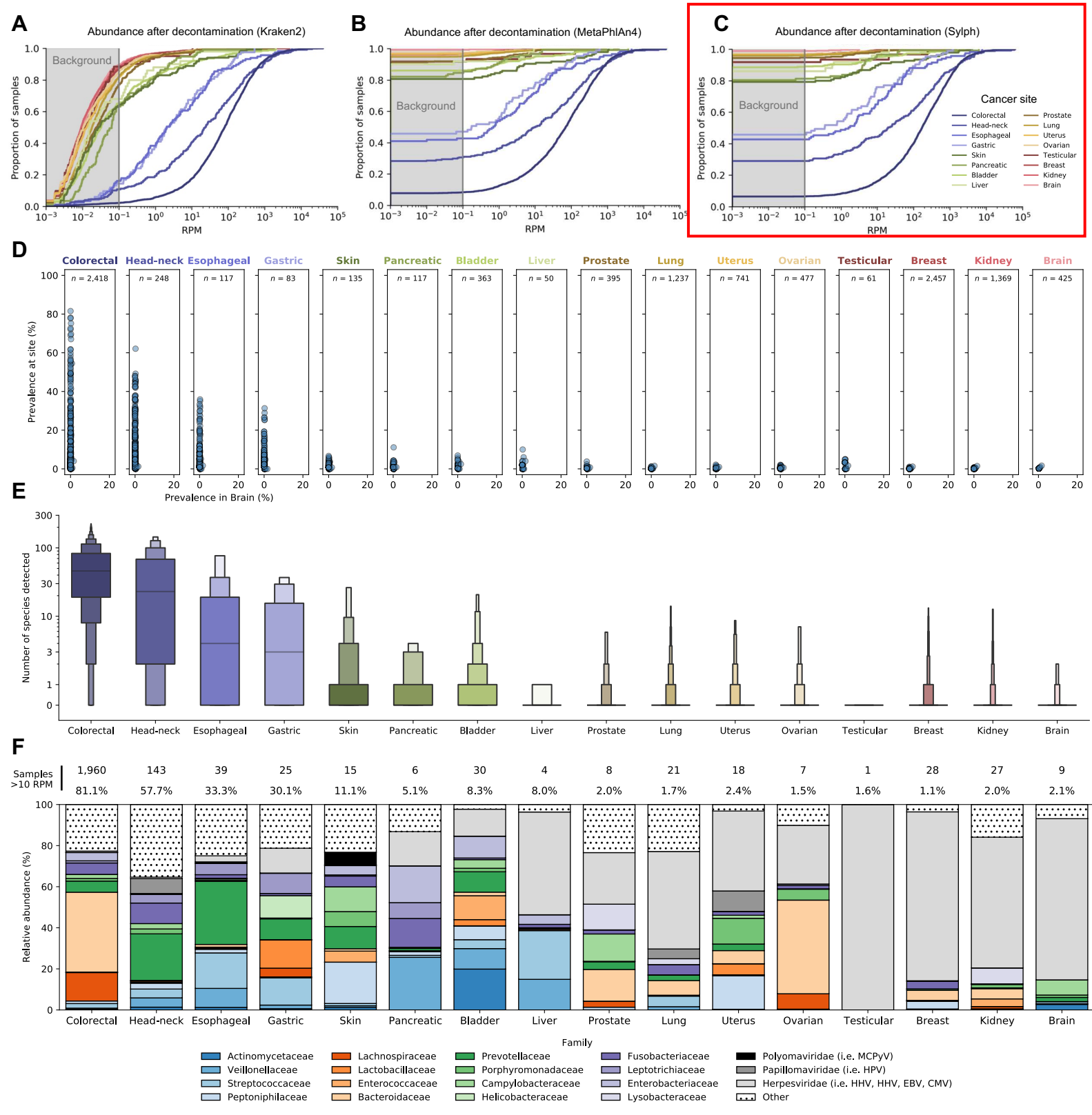
**Figure 2** Identification and removal of contamination from the 100kGP cohort



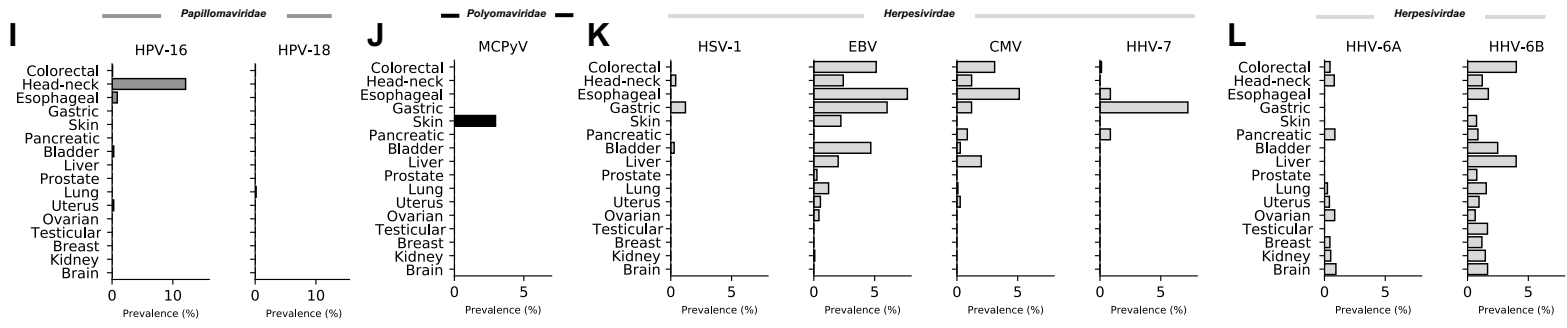
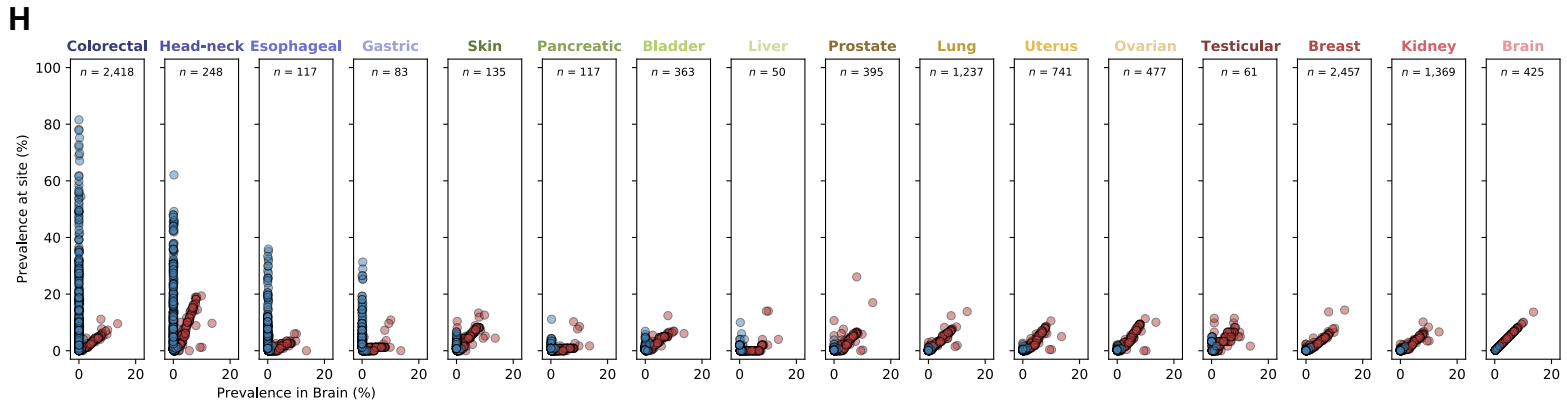
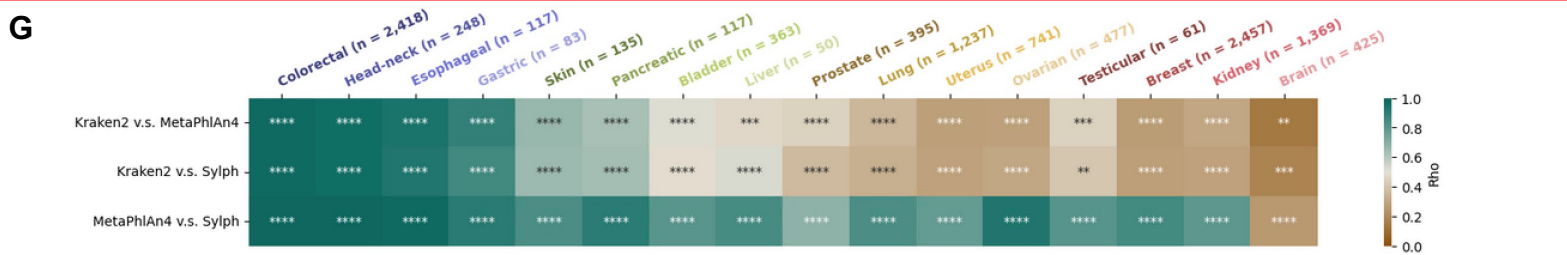
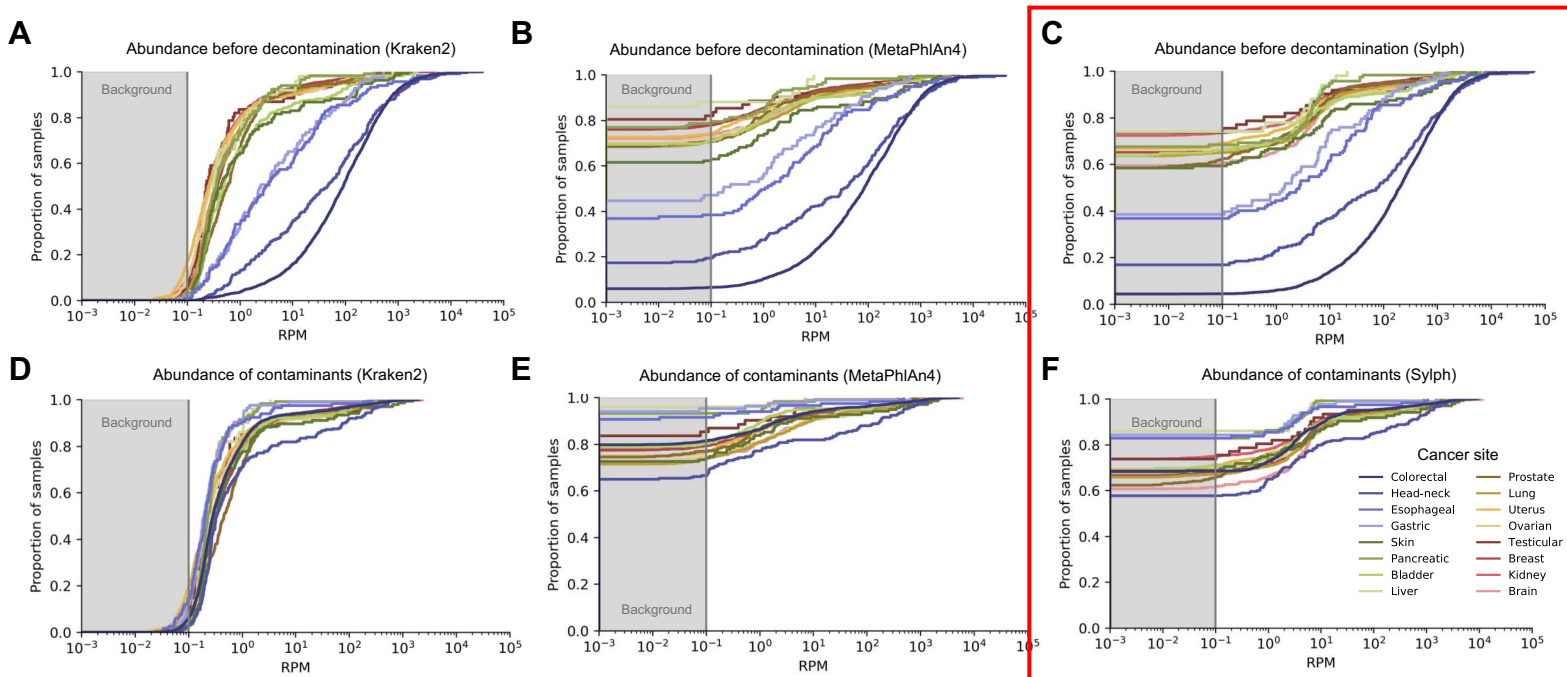
**Figure S2** Identification and removal of contamination from the 100kGP cohort



**Figure 3** Tumors of the oro-digestive tract harbor a microbiome but most other cancer types do not



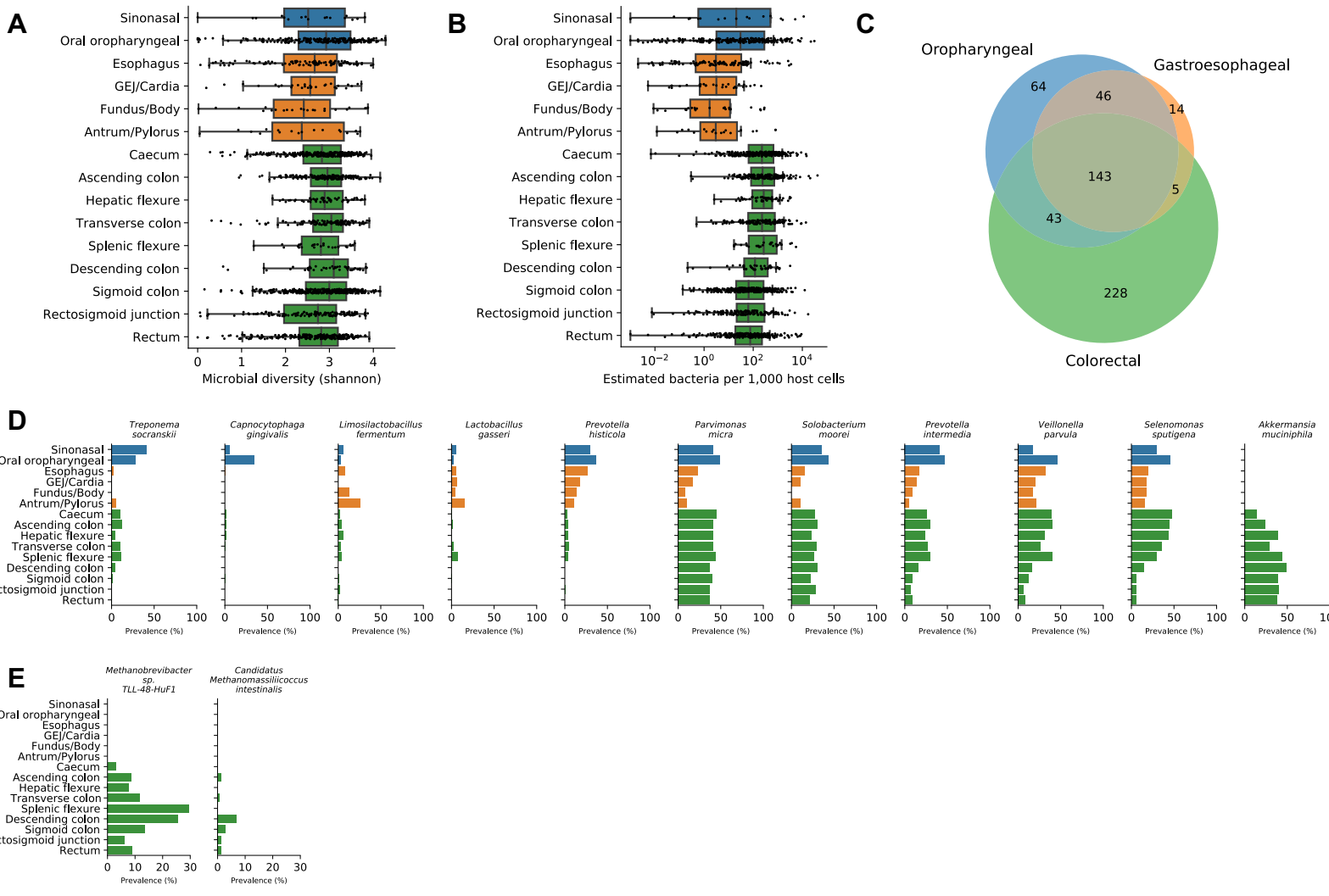
**Figure S3** Tumors of the oro-digestive tract harbor a microbiome but most other cancer types do not



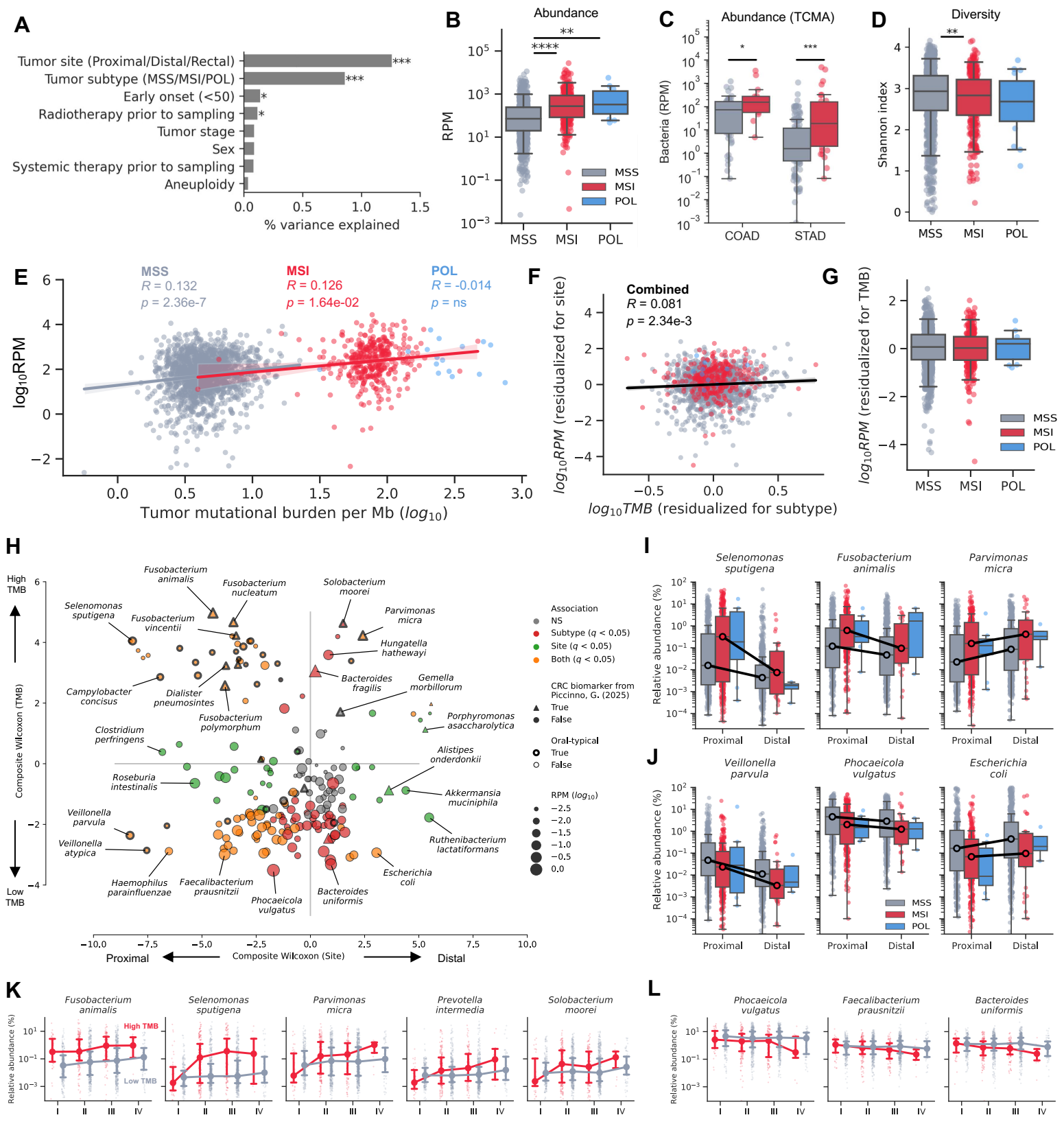
**Figure 4** Biogeography and biodiversity of oropharyngeal, gastroesophageal, and colorectal cancers



**Figure S4** Biogeography and biodiversity of oropharyngeal, gastroesophageal, and colorectal cancers



**Figure 5** The tumor microbiome varies by tumor site, mutation burden, and age of onset



**Figure S5** The tumor microbiome varies by tumor site, mutation burden, and age of onset

