

Metagenomic analysis of colorectal cancer datasets identifies cross-cohort microbial diagnostic signatures and a link with choline degradation

Authors:

Andrew Maltez Thomas^{1,2,3,#}, Paolo Manghi^{1,#}, Francesco Asnicar¹, Edoardo Pasolli¹, Federica Armanini¹, Moreno Zolfo¹, Francesco Beghini¹, Serena Manara¹, Nicolai Karcher¹, Chiara Pozzi⁴, Sara Gandini⁴, Davide Serrano⁴, Sonia Tarallo⁵, Antonio Francavilla⁵, Gaetano Gallo^{6,7}, Mario Trompetto⁷, Giulio Ferrero⁸, Sayaka Mizutani^{9,10}, Hirotugu Shirota⁹, Satoshi Shiba¹¹, Tatsuhiro Shibata^{11,12}, Shinichi Yachida^{11,13}, Takuji Yamada^{9,14}, Jakob Wirbel¹⁵, Petra Schrotz-King¹⁶, Cornelia M. Ulrich¹⁷, Hermann Brenner^{16,18,19}, Manimozhayan Arumugam^{20,21}, Peer Bork^{15,22,23,24}, Georg Zeller¹⁵, Francesca Cordero⁸, Emmanuel Dias-Neto^{3,25}, João Carlos Setubal^{2,26}, Adrian Tett¹, Barbara Pardini^{5,27}, Maria Rescigno²⁸, Levi Waldron^{29,30,*}, Alessio Naccarati^{5,31,*}, Nicola Segata^{1,*,^}

1 - Department CIBIO, University of Trento, Trento, Italy.

2 - Biochemistry Department, Chemistry Institute, University of São Paulo, São Paulo, Brazil.

3 - Medical Genomics Laboratory, CIPE/A.C. Camargo Cancer Center, São Paulo, Brazil.

4 - IEO, European Institute of Oncology IRCCS, Milan, Italy.

5 - Italian Institute for Genomic Medicine (IIGM), Turin, Italy.

6 - Department of Surgical and Medical Sciences, University of Catanzaro, Catanzaro, Italy.

7 - Department of Colorectal Surgery, Clinica S. Rita, Vercelli, Italy.

8 - Department of Computer Science, University of Turin, Turin, Italy

9 - School of Life Science and Technology, Tokyo Institute of Technology, Tokyo, Japan

10 - Research Fellow of Japan Society for the Promotion of Science

11 - Division of Cancer Genomics, National Cancer Center Research Institute, Tokyo, Japan

12 - Human Genome Center, The Institute of Medical Science, The University of Tokyo, Tokyo, Japan

13 - Department of Cancer Genome Informatics, Osaka University, Osaka, Japan

14 - PRESTO, Japan Science and Technology Agency, Saitama, Japan

15 - Structural and Computational Biology Unit, European Molecular Biology Laboratory, Heidelberg, Germany

16 - Division of Preventive Oncology, National Center for Tumor Diseases (NCT) and German Cancer Research Center (DKFZ), Heidelberg, Germany

17 - Huntsman Cancer Institute and Department of Population Health Sciences, University of Utah, Salt Lake City, Utah, USA

18 - Division of Clinical Epidemiology and Aging Research, German Cancer Research Center (DKFZ), Heidelberg, Germany

19 - German Cancer Consortium (DKTK), German Cancer Research Center (DKFZ), Heidelberg, Germany

20 - Novo Nordisk Foundation for Basic Metabolic Research, Faculty of Health and Medicine, University of Copenhagen, Denmark

21 - Faculty of Healthy Sciences, University of Southern Denmark, Odense, Denmark

22 - Molecular Medicine Partnership Unit (MMPU), Heidelberg, Germany

23 - Max Delbrück Centre for Molecular Medicine, Berlin, Germany

24 - Department of Bioinformatics, Biocenter, University of Würzburg, Würzburg, Germany

25 - Laboratory of Neurosciences (LIM-27), Institute of Psychiatry, University of São Paulo, São Paulo, Brazil.

26 - Biocomplexity Institute of Virginia Tech, Blacksburg VA 24061, USA

27 - Department of Medical Sciences, University of Turin, Turin, Italy.

28 - Mucosal immunology and microbiota Unit, Humanitas Research Hospital, Milan, Italy.

29 - Graduate School of Public Health and Health Policy, City University of New York, New York, USA.

30 - Institute for Implementation Science in Population Health, City University of New York, New York, USA.

31 - Department of Molecular Biology of Cancer, Institute of Experimental Medicine, Prague, Czech Republic.

These authors contributed equally

* Co-senior authors

^ Corresponding author nicola.segata@unitn.it

Abstract

52
53
54 Several studies have investigated links between the gut microbiome and colorectal cancer (CRC), but
55 questions remain about the replicability of biomarkers across cohorts and populations. We
56 performed a meta-analysis of five publicly available datasets and two new cohorts, and validated the
57 findings on two additional cohorts, considering in total 969 fecal metagenomes. Unlike microbiome
58 shifts associated with gastrointestinal syndromes, the gut microbiome in CRC showed reproducibly
59 higher richness than controls ($P < 0.01$), partially due to expansions of species typically from the oral
60 cavity. Meta-analysis of the microbiome functional potential identified gluconeogenesis and the
61 putrefaction and fermentation pathways to be associated with CRC, whereas the stachyose and starch
62 degradation pathways were associated with controls. Predictive microbiome signatures for CRC
63 trained on multiple datasets showed consistently high accuracy in datasets not considered for model
64 training and independent validation cohorts (average AUC 0.84). Pooled analysis of raw
65 metagenomes showed that the choline trimethylamine-lyase gene was over-abundant in CRC ($P =$
66 0.001) identifying a novel relationship between microbiome choline metabolism and CRC. The
67 combined analysis of heterogeneous CRC cohorts thus identified reproducible microbiome
68 biomarkers and accurate disease-predictive models that can form the basis for clinical prognostic
69 tests and hypothesis-driven mechanistic studies.

70
71
72
73
74
75
76
77
78

Introduction

Colorectal cancer (CRC) is the second most common non sex-specific cancer and is responsible for more deaths than any other cancer after lung cancer¹. Because of demographic trends toward an ageing population, the global incidence rate is expected to increase by nearly 80% to 2.2 million cases per year over the next two decades². Sporadic CRCs, as opposed to hereditary CRCs, account for approximately 70%-87% of cases³ and genetics can only explain a small proportion of disease incidence⁴. The missing strong link of CRC with genetics points to the potential role of other variables including lifestyle and environmental factors as disease co-determinants. Reported risk factors associated with CRC include age, tobacco and alcohol consumption, lack of physical activity, increased body weight, and diet^{5,6}. However, many non-genetic risk factors are common to several cancer types and these factors remain largely unsettled for CRC^{7,8}.

The human gut microbiome - defined as the microbial communities that populate our intestinal tract - is emerging as a relevant factor in human diseases^{9,10}. Supported by some evidence of carcinogenic mechanisms induced by bacterial organisms¹¹⁻¹³, the gut microbiome has also been hypothesized to play a crucial role in the development of CRC. Studies using 16S rRNA gene amplicon sequencing have led to the discovery of *Fusobacterium nucleatum*'s association with CRC¹⁴, which was subsequently shown to be causal in animal models of CRC carcinogenesis and progression^{15,16}. Compared to 16S rRNA gene studies, a smaller number of metagenomic sequencing studies have linked other microbial species and potential functional activities of the gut microbiome to CRC¹⁷⁻¹⁹. However, the reproducibility and predictive accuracy of these high-resolution microbial signatures across cohorts and study design choices remain unclear. The potential use of the gut microbiome as a diagnostic tool for CRC has been proposed¹⁷⁻²¹, but not yet validated across multiple independent study populations.

There is thus a need to establish and validate links between the human gut microbiome and CRC carcinogenesis across populations, cohorts, and microbiome tools. Some multi-cohort works have been performed based on 16S rRNA gene studies²², but this technique has important technical limitations²³. The recent availability of whole-metagenome shotgun datasets of CRC cohorts¹⁷⁻²¹ enables a combined multi-population exploration of the CRC-associated microbiome with strain-level resolution^{24,25} and meta-analytic predictive approaches^{10,26}, but the only meta-analysis study performed so far on CRC is affected by overfitting issues²⁷. It is thus crucial to perform large-scale cross-cohort studies to provide an unbiased and well-powered assessment of the link between CRC and the gut microbiome.

In this study, we have sequenced 140 samples from two different cohorts, performed an integrated analysis combining all current metagenomic CRC datasets available, and assessed prediction accuracies of the gut microbiome for CRC detection across populations, datasets, and conditions.

113

Results

114 **A meta-analysis of metagenomic datasets to identify links between the gut microbiome and** 115 **CRC**

116 To identify reproducible relationships between the gut microbiome and CRC, we performed shotgun
117 metagenomic sequencing²⁸ of the stool microbiome of 140 CRC patients and controls recruited in
118 two cohorts, and analyzed these in the context of 624 additional samples from five publicly available
119 and geographically diverse metagenomic studies. We validated the results on two novel datasets of
120 60 CRC and 65 controls²⁹ and 40 CRC and 40 controls (see **Methods**), respectively. In total, we

121 considered 413 samples from CRC patients, 143 from subjects with adenoma and 413 control
122 samples. Participants from all studies underwent colonoscopy to diagnose CRC, adenoma, or to
123 confirm the absence of disease, with samples collected before diagnosis or beginning of treatment
124 (**Suppl. Table 1, Table 1**). All datasets were sequenced at high depth except for the Hannigan *et al.*
125 study³⁰ (**Extended Data 1A, Methods**).

126 **Meta-analysis shows higher species richness in CRC-associated samples**

127 We first tested whether microbial richness and diversity differed between CRC samples and controls,
128 given contrasting current evidence³¹⁻³³. In all but one study, the median species richness was higher
129 in CRC samples compared to controls, and the increase was significant in four of the six deeply
130 sequenced datasets ($P < 0.05$ **Extended Data 1B-C**). Meta-analysis of standardized mean differences
131 by random effects model for the number of microbial species confirmed the higher number of species
132 in CRC compared to controls ($\mu=0.5$, 95% CI [0.16, 0.85], $P = 0.004$), although with significant
133 heterogeneity across datasets ($I^2 = 74.8\%$, $p = 0.0007$, Q-test). This difference was not meaningfully
134 affected when controlling for potential confounding by age, BMI, or sex(**Extended Data 1D-E**).
135 Conversely, we observed no difference in diversity between carcinomas and controls (**Extended**
136 **Data 2A-B**). We thus provide strong evidence that the CRC-associated microbiome has a quantitative
137 species distribution which is consistent with healthy controls, but is significantly enriched in the total
138 number of detected microbes.

139 We further tested whether the CRC-associated microbiome possesses more oral cavity-associated
140 species than controls, as previously hypothesized^{22,34}. Considering the 161 species we identified from
141 multiple existing datasets^{35,36} as being typical colonizers of the oral cavity (see **Methods**), we found
142 increased oral species richness in CRC samples for all but one of the six deeply sequenced datasets
143 compared to controls and the increase was significant in meta-analysis ($\mu = 0.16$, 95% CI [-0.03, 0.35],
144 $P = 0.02$, **Extended Data 2G**). Similarly, the total abundance of oral species in the stool microbiome
145 was also significantly higher in CRC patients compared to controls (meta-analysis $\mu=0.23$, 95% CI
146 [0.07, 0.39], $P = 0.003$). Altogether, greater species richness and abundance may be a sign of an
147 altered gut microbiome in CRC, and it is indicative of an influx of bacterial species originating from
148 the oral cavity.

149 **A panel of microbial biomarkers for CRC is reproducible across cohorts**

150 Individual biomarker discovery efforts can be sensitive to technical artefacts and to heterogeneity of
151 factors implicated in microbial shifts in healthy populations, including biogeography, diet, and host
152 genetics^{25,37}. This is confirmed by the two newly sequenced datasets that have only partially
153 overlapping taxonomic and functional potential biomarkers (**Extended Data 3**). Even so, several CRC
154 biomarker species were identified by univariate statistics³⁸ independently in the majority of the
155 datasets: *F. nucleatum*, *Solobacterium moorei*, *Porphyromonas asaccharolytica*, *Parvimonas micra*,
156 *Peptostreptococcus stomatis*, and *Parvimonas ssp.* Other species were identified in fewer datasets or
157 were dataset-specific (**Figure 1A**, and **Suppl. Table 2**). *F. nucleatum*, whose connection with CRC has
158 been extensively reported^{14,17-19}, had significantly increased abundance in CRC patients in all
159 datasets with adequate sequencing depth, when considering single markers for this species
160 (**Extended Data 4A**). Some of the cross-cohort CRC biomarker species have already been reported
161^{14,22,34} and many of them are commonly found in the oral cavity (8 out of the 39 total biomarkers
162 found in at least 2 datasets), consistent with the increased oral taxa presence in CRC samples
163 mentioned above.

164 We then pooled evidence of differential abundance across datasets by random effects meta-analysis.
165 Among the 26 differentially abundant species at FDR < 0.005, those with the highest effect size were
166 again *F. nucleatum*, *S. moorei*, *P. asaccharolytica*, *P. micra* and *P. stomatis*. The meta-analysis
167 additionally identified *Clostridium symbiosum*, which has been tested as a marker for early CRC
168 detection ³⁹ (**Figure 1B**). Other differentially abundant species at FDR < 0.05 have not been
169 previously reported in CRC microbiome studies, including *Streptococcus tigurinus* and *Streptococcus*
170 *dysgalactiae*, and 3 different *Campylobacter* species. We also confirmed *Gemella morbillorum* and
171 *Streptococcus gallolyticus* to be relevant biomarkers, as previously suggested in smaller cohorts ^{18,40}.
172 In contrast, only 12 species were associated with the control population in the meta-analysis and only
173 four were significantly enriched for the same populations in at least three datasets. Control-
174 associated species with the highest effect sizes were *Gordonibacter pamelae* and *Bifidobacterium*
175 *catenulatum* (**Figure 1B, Suppl. Table 2; Extended Data 4C**), which are generally considered
176 beneficial microbes and have been used as probiotic supplements ⁴¹. Adjustment for potential
177 confounding by host characteristics did not meaningfully affect crude estimates in the meta-analysis
178 (**Figure 1D, Extended Data 4B**). The substantially higher number of species enriched in CRC than in
179 controls (49 vs. 12), even when focusing only on species with putative oral origin (15 vs. 2, **Extended**
180 **Data 5A**), points to the existence of a reproducible taxonomic signature of the CRC-associated
181 microbiome.

182 Functional potential of the microbiome was also significantly associated with CRC samples when
183 compared against healthy controls. We found overall increased richness of UniRef gene families ⁴² in
184 CRC samples in two datasets, with percentages of unmapped reads ranging between 20% and 40%
185 (**Extended Data 5E**). We found 33,840 of the 2,479,274 single gene families detected at least once to
186 be associated with CRC and 30,475 associated with controls (FDR < 0.05, 9,154 and 7,115 differential
187 gene families at FDR < 0.005 respectively). We further observed 136 out of 590 metagenomically
188 reconstructed microbial functional pathways to be CRC-associated, and only 37 associated with
189 controls (**Suppl. Table 3**). Among the most differentially abundant pathways (**Figure 1C**) that are at
190 worst just minimally affected by potential confounding factors (**Figure 1E**), we found starch,
191 stachyose, and galactose degradation to be associated with controls. These associations could
192 indicate how potentially diet-associated changes in the functional repertoire of the microbiome can
193 influence host conditions. The CRC-associated microbiome showed an association with
194 gluconeogenesis and with capacity for uptake and metabolism of amino acids via putrefaction and
195 fermentation pathways (**Suppl. Table 3-4**). These included those pathways responsible for the
196 conversion of different amino acids to tumor-promoting compounds ^{19,43}, such as polyamines (e.g. L-
197 arginine and L-ornithine degradation to putrescine) and ammonia (L-histidine and L-arginine
198 degradation, and L-lysine and L-alanine fermentation to acetate, butyrate and propionate). These
199 pathways (**Figure 1C**) and the set of species described above (**Figure 1A,B**) thus constitute a
200 collection of microbiome biomarkers that is reproducible across cohorts.

201 ***Predicting CRC from single metagenomic datasets in independent cohorts leads to reduced*** 202 ***accuracy***

203 To test the hypothesis that the stool microbiome could be used as a reproducible CRC pre-screening
204 tool, we performed intra-cohort, cross-cohort and combined-cohort prediction validation on the
205 overall set of 621 CRC and controls samples using a Random Forest classifier (**Table 1**). In intra-
206 cohort cross-validation using species-level taxonomic relative abundances, we observed
207 performances ranging from 0.92 to 0.58 AUC score, with an average in the deeply sequenced datasets
208 of 0.81 AUC (**Figure 2A**). When using the functional potential of the gut microbiome by means of

209 pathway abundances, we observed decreased single dataset cross-validation accuracies, with the
210 exception of our Cohort1 (maximum 0.82 AUC, average 0.71 AUC, **Extended Data 6A**). The profiling
211 of the more fine-grained UniRef90 gene family abundances improved the predictions, with AUCs
212 reaching 0.84 AUC for Cohort2 and an average of 0.77 AUC in the deeply sequenced datasets (**Figure**
213 **2B**). These results show that, while cross validation AUCs can be high for predicting CRC in some
214 datasets, they are highly variable and dataset dependent.

215 We then tested whether and how much the microbial signatures of CRC remained predictive across
216 distinct datasets and cohorts. To this end, we trained the classifier on each single “training” dataset
217 and applied the model on each distinct “testing” dataset. For most datasets this led to decreased AUC
218 values when compared to single cross validation AUCs, and AUCs showed a high variability across
219 cohorts (minimum 0.5 and maximum 0.86 cross dataset AUC). These results were consistent when
220 using either pathway or gene family-abundances as predictors (**Extended Data 6** and **Figure 2B**).
221 Overall, we highlight a poor transportability of the microbiome signature from one dataset to the
222 other and experimental choices⁴⁴ and cohort or population characteristics²⁵, may explain the
223 reduced cross-study predictability when considering single datasets to train the model (**Extended**
224 **Data 6C-D**).

225 *Pooling of training cohorts substantially improves prediction across datasets*

226 To overcome the limitations of training on single datasets (**Suppl. Table 5**), we performed a Leave-
227 One-Dataset-Out (LODO) analysis⁴⁵ in which classifiers were trained on six datasets combined, and
228 validated on the left-out dataset, for each dataset in turn. For taxonomic profiles, this approach
229 improved both AUC values and inter-dataset consistency, producing AUCs ≥ 0.80 (average 0.84 s.d.
230 0.03) for all six deeply sequenced datasets (**Figure 2A**). Predictors based on clade-specific markers
231 also produced high, albeit more variable AUC values, outperforming taxonomic profiles in some
232 datasets (**Extended Data 6B**). Gene families achieved slightly reduced performances, whereas
233 pathway abundances produced substantially less accurate predictions (**Figure 2B**). The technical and
234 host population diversity embedded in these training meta-cohorts may be crucial in improving the
235 generalizability of classifiers, as we found this LODO approach to be substantially and consistently
236 more informative than a single-dataset cross-validation, and independent investigations found
237 similarly high LODO performances using different metagenomic profiles and machine learning tools
238²⁹.

239
240 The model trained on taxonomic or functional features was also shown to capture the above whole-
241 microbiome biomarkers because the direct inclusion of alpha-diversity metrics, oral-species
242 abundance, and a measure of metagenome mappability did not provide substantial improvements
243 (mean 0.83, s.d. 0.03 for the deeply sequenced datasets when using the taxonomic model). However,
244 based on the performance and variability of the predictive models across datasets, we recommend
245 using species-level microbial abundance as the main feature set for CRC status prediction in a LODO
246 setting.

247 To assess the relation between population diversity in the training meta-cohort and prediction
248 performance, we considered increasingly larger subsets of the available training cohorts. AUC values
249 sharply increased when moving from one to two training datasets (10% to 13% median AUC
250 improvement depending on the features considered in the model, **Extended Data 7**) with less
251 marked improvements at further dataset additions (**Figure 2C-D**). Large and heterogeneous
252 combined training sets thus generate improved accuracy for identifying CRC cases in independent
253 metagenomic datasets.

254 ***Accurate predictive models using a minimal microbial signature***

255 The predictive CRC-associated microbiome signatures identified above considered all observed
256 species and gene functions and would thus be impractical for clinical application without whole
257 microbiome profiling. We thus sought to identify a minimal set of highly predictive microbial features
258 by exploiting the internal feature ranking of the Random Forest classifier¹⁰. We found that *P. stomatis*
259 was the species with the highest average rank. As expected, other CRC-associated species such as *F.*
260 *nucleatum*, *Parvimonas ssp.*, *P. asaccharolytica*, *G. morbillorum*, *Clostridium symbiosum* and *P. micra*
261 were also crucial to prediction accuracy (**Figure 3A**) with the top seven ranked species for CRC
262 detection amongst those with the largest effect sizes in the meta-analysis. Very few species were
263 ranked high in the learning models, further highlighting that successful discrimination is achieved by
264 CRC-specific rather than control-specific microbial features.

265 To evaluate how many microbial species or gene families are necessary to achieve prediction scores
266 comparable to those obtained using the full set of features, we computed AUC values at increasing
267 numbers of features. Feature ranking was performed internally to each training fold to avoid
268 overfitting. By applying this approach to all datasets (**Figure 3B-C**), we found that using as few as 16
269 species achieved CV AUC >0.8 for the majority of the datasets, with little improvement from using all
270 remaining species (2% improvement in the mean AUC value). We also found that using only 64 gene
271 families achieved prediction values >0.8 for the same datasets, and that using all 8,192 gene families
272 improved AUC only slightly (2% improvement -**Extended Data 8**). Therefore, these results suggest
273 that a stool-based diagnostic test using genetic markers targeting a limited number of microbial
274 species or genes would serve as a promising clinical tool.

275 ***Microbiome signatures for adenomas are only partially predictive***

276 We assessed the ability to discriminate adenomas from controls or carcinomas, using 27 newly
277 sequenced adenoma-associated samples and 116 adenoma-associated samples from available studies
278 (**Table 1**). Adenomas could be distinguished from CRC patients with lower accuracy than controls
279 (mean AUC 0.69 versus 0.79, **Extended Data 6E-F**) and there are only eight species that differentiate
280 adenoma patients from carcinoma patients in the meta-analysis (FDR < 0.05). Seven of these eight
281 biomarkers are in common with the comparison between carcinoma patients and healthy individuals,
282 and the LODO approach did not improve discrimination of adenomas from CRC (average AUC 0.68).
283 Moreover, we found that no dataset could accurately predict adenomas from control samples
284 (maximum AUC 0.58, minimum 0.46), even when using a LODO approach (average AUC 0.54). In the
285 meta-analysis, no species were significantly different when contrasting samples from patients with
286 adenomas and healthy controls. These results reinforce previous findings^{18,19} that the adenoma-
287 associated stool microbiome closely resembles that of the healthy gut.

288 ***Increased abundance of choline TMA-lyase encoding genes in CRC***

289 Microbiome-derived metabolites and specifically polyamines have been implicated in carcinogenesis
290 both in animal models and in humans⁴³. We chose to focus on trimethylamine (TMA), an amine
291 produced by bacteria from choline and carnitine, because it has been shown to play a role in complex
292 diseases such as atherosclerosis and primary sclerosing cholangitis^{9,46}. Since dietary components
293 have been linked with CRC risk^{5,6}, we hypothesized that the TMA-producing potential of the human
294 gut microbiome could also be associated to CRC⁴⁷. To test this hypothesis, we considered the genes
295 belonging to the main TMA-synthesis pathways to reconstruct and quantify the presence of such
296 genes in the CRC-associated metagenomes. The main genes associated with TMA-synthesis are those
297 encoding the choline TMA-lyase (*cutC*), the L-carnitine dioxygenase (*yeaW*) and the L-

298 carnitine/gamma-butyrobetaine antiporter (*caiT*) and we identified them in 923, 5,185 and 5,709
299 available bacterial genomes, respectively.

300 Screening the 7 CRC-associated metagenomic datasets, we found that only one of them had a
301 significant increase of *caiT* in CRC samples compared to controls, whereas no significant differences
302 were detected for *yeaW* (**Extended Data 9A**). However, we found increased abundance of *cutC* in
303 CRC samples compared to controls in all seven datasets ($P < 0.05$ by Wilcoxon Rank Sum test on
304 RPKM abundances for five datasets, **Figure 4A**). Meta-analysis indicated an overall strong association
305 with no evidence of heterogeneity ($P = 0.001$, $\mu = 0.27$, 95% CI [0.1, 0.42], $I^2 = 4.2\%$, Q-test = 0.65,
306 **Figure 4B**). We also analyzed the abundance of the gene encoding the choline TMA-lyase-activating
307 enzyme (*cutD*), finding a significant increase in CRC (meta analysis $P = 0.001$, $\mu = 0.32$, 95% CI [0.16,
308 0.47], $I^2 = 0\%$, Q-test = 0.96, **Extended Data 9B-C**). These results indicate that TMA production might
309 happen preferentially via choline degradation, and not via carnitine, and could substantially affect the
310 amounts of TMA and trimethylamine oxide (TMAO) in an individual⁴⁸. Intermediate levels of *cutC* in
311 adenomas (**Figure 4A**) is further suggestive of a TMA action along the adenoma-carcinoma axis. We
312 validated the increased *cutC* gene abundance in CRC by qPCR⁴⁹ on a subset of samples from Cohort1
313 with enough DNA left after sequencing, and confirmed the metagenomic findings (one-tailed
314 Wilcoxon signed rank test $P = 0.024$, **Figure 4D**). Further quantification of *cutC* transcript abundance
315 from the co-extracted RNA in the same dataset also pointed to an over-expression of this gene in CRC
316 ($P = 0.035$, **Figure 4E**).

317 We further explored the role of *cutC* in the gut microbiome by reconstructing sample-specific
318 sequence variants using a reference-aided targeted assembly approach (see **Methods**). We found a
319 large sequence divergence for the gene encoding this enzyme that is known to occur in single copies
320 in the genomes⁴⁹ and we identified four main sequence variants that are associated with the
321 taxonomic structure (**Figure 4B, Extended Data 9C-D, 10A-B**). Interestingly, the most prevalent
322 (46.5%) *cutC* sequence type belonged (>95% identity over the full length of the gene) to an unknown
323 species that was only recently assembled from metagenomics⁵⁰ and assigned to species-level
324 genome bin (SGB) ID 3957. This candidate species comprises 56 metagenomically-assembled species
325⁵⁰ and is placed within the *Lachnospiraceae* family, but the missing genus assignment confirms that
326 several microbes remain under-characterized in the human microbiome. This *cutC* variant was
327 associated with non-CRC samples (OR 0.38, 95% CI [0.25, 0.57], $P = 0.0001$, Fisher Test), whereas
328 *cutC* sequence types mostly belonging to *Hungatella hathewayi* and *Clostridium asparagiforme*
329 (*Firmicutes*) were significantly CRC-associated (OR 2.14, 95% CI [1.29, 3.56], $P = 0.004$, Fisher test),
330 as were sequence types belonging to *Klebsiella oxytoca* and *Escherichia coli* (OR 1.85, 95% CI [1.13, 3],
331 $P = 0.02$, Fisher Test - **Figure 4B**). Altogether, these novel findings highlight that sequence variants of
332 *cutC* can be strongly associated with disease, potentially because of corresponding differences in the
333 efficacy of choline degradation and TMA production.

334 **Additional independent validation of predictive models**

335 To further validate our meta-analysis results, we considered two additional independent
336 metagenomic cohorts from Germany²⁹ (Validation Cohort1) and Japan (Validation Cohort2)
337 comprising a total of 100 CRC patients and 105 controls (see **Methods**). The metagenomic predictive
338 model was confirmed to be highly accurate on these new cohorts (**Figure 5A**) with an AUC of 0.90
339 and 0.81 for the German and Japanese cohorts respectively, when using the species-level taxonomic
340 abundance model. Species newly associated to the CRC microbiome such as *Streptococcus tigurinus*
341 and *Streptococcus dysgalactiae* were confirmed to have higher prevalence in CRC than in controls In
342 the two validation datasets (blocked Wilcoxon test⁵¹ $P = 0.049$ and $P = 0.011$ for *S. tigurinus* and *S.*

343 *dysgalactiae*, respectively). Enrichment in the CRC-associated microbiome of these two species was
344 confirmed also by the analysis of additional metagenomic datasets of IBD⁵² and type-2 diabetes^{53,54}
345 in which the prevalence of *S. tigurinus* was always below 10% in both cases and controls, whereas *S.*
346 *dysgalactiae* was never detected in these additional datasets. We also confirmed species richness to
347 be significantly higher in CRC ($P = 0.0005$ for both validation datasets after rarefaction at the 10th
348 percentile, **Figure 5B**) as well as richness of oral microbial species in the rarefied samples (blocked
349 Wilcoxon test⁵¹ $P = 0.003$), and the abundance of the gene encoding the choline TMA-lyase enzyme
350 *cutC* in CRC ($P < 1e-6$).

351 **CRC-specificity of microbiome predictive models**

352 We performed additional experiments to validate the discriminative power of the above microbial
353 signatures specifically for CRC and not for other potentially microbiome-linked disease conditions. To
354 this end, we first considered 13 additional fecal samples sequenced from patients that underwent
355 colonoscopy in our Cohort1 that were originally discarded because the final diagnosis pointed at
356 diseases other than adenomas or carcinomas such as ulcerative colitis, Crohn's disease,
357 unclassified colitis, and diverticular diseases. These were distinguishable from CRC samples based
358 on the taxonomic model (0.78 cross-validation AUC, 0.80 AUC using only 16 species), and only
359 slightly decreased the AUC of the model trained on all the other datasets when they were added to
360 the non-disease (i.e. healthy) category (from 0.83 to 0.79 in AUC). We then expanded this analysis to
361 diseases for which at least two distinct large metagenomic datasets are available in the public domain
362 and this includes ulcerative colitis (UC) and Crohn's disease (CD)^{52,55} as well as non-GI diseases such
363 as type-2 diabetes^{53,54}. For this purpose we added samples randomly drawn from each of the case
364 and control conditions of these additional disease cohorts to the control class of the new validation
365 cohort and recorded the variations in AUCs when attempting to predict CRC (see **Methods**). By
366 comparing the AUCs obtained when adding non-CRC external cases and when adding the
367 corresponding external controls, we found for both validation cohorts a small decrease in prediction
368 accuracy for both UC (3% and 4% for Validation Cohort1 and Validation Cohort2, respectively;
369 **Figure 5C**) and CD (5% and 9%, for Validation Cohort1 and Validation Cohort2, **Figure 5C**), pointing
370 to a limited effect on the CRC model of samples from these two diseases. For type-2 diabetes we
371 observed an increase in the predictive power in one dataset⁵³, and a decrease in the other⁵⁴ in both
372 validation datasets, and the CRC model always remained highly predictive ($AUC \geq 0.80$). Altogether,
373 these results point at the existence of a clear microbiome signature of CRC which is distinct from
374 other relevant diseases with a gastrointestinal component.

375 **Relationship to currently available non-invasive clinical screening tests**

376 To assess the potential of microbiome-based prediction models in comparison and in combination
377 with currently used non-invasive clinical screening tests, we considered the Fecal Occult Blood Test
378 (FOBT) and the Wif-1 Methylation test available for 110 samples of the ZellerG_2014 cohort¹⁹. The
379 LODO microbiome model tested on this dataset proved to be slightly superior to the FOBT at multiple
380 combinations of specificity and sensitivity levels (**Figure 5D**) and on par with the Wif-1 Methylation
381 test. Considering the LODO model predictions and the FOBT together in the same test improves the
382 sensitivity/specificity trade-off at high specificity levels when the integration is based on having at
383 least one predictor positive, and at relatively lower specificity levels when requiring both predictors
384 to be positive (**Figure 5D**). Integrating the microbiome model with the Wif-1 Methylation test results
385 in similar performances, and the use of the reduced microbiome model with only 16 species generally
386 improves the results (**Figure 5D**). We thus provide evidence for the potential clinical value of

387 microbiome predictive models especially when considered together with other available non-
388 invasive clinical tests.

389

390

Discussion

391 In the present study, we comprehensively assessed the CRC-associated gut microbiome and its ability
392 to distinguish newly diagnosed CRC patients from tumor-free controls. Our study was performed
393 across multiple datasets and populations, through a combined analysis of fecal CRC metagenomes
394 from four previously unpublished cohorts and five publicly available datasets. Whereas direct specific
395 host-microbe interactions have been shown to cause certain malignancies *in vitro* and *in vivo* animal
396 models^{11-13,56} and genotoxic determinants such as colibactin tend to be over-represented in the
397 analyzed datasets²⁹, indirect metabolite-mediated mechanisms may be more important to the
398 development of carcinomas although causality relations need to be tested. In our analysis, we indeed
399 found a reproducible panel of microbiome species (**Figure 1**), whole microbiome characteristics, and
400 strain-level biomarkers (**Figure 4**) beyond the validated mechanisms of specific variants of
401 *Escherichia coli*^{11,56} and *Bacteroides fragilis*⁵⁶. We found that the gut microbiome in CRC has greater
402 richness than controls, partially due to the presence of oral cavity-associated species rarely found in
403 healthy guts, challenging the widespread assumption that decreased alpha-diversity is generally
404 associated with intestinal dysbiosis^{57,58}.

405 The identification of reproducible microbial biomarkers for CRC may enable the design of non-
406 invasive diagnostic tools. We developed machine learning models able to distinguish between
407 carcinoma patients and controls with an average performance above 0.84 AUC when validated on
408 datasets excluded from the training of the model (**Figure 2A**). Importantly, these performances are
409 quite independent of specific methodological choices given that complementary investigations²⁹
410 using different metagenomic profilers and machine learning approaches achieved very similar
411 results. Further increase in prediction performance can be achieved using larger datasets ($n > 1,000$)
412 rather than different methodologies (**Figure 2C-D, Figure 5C**), and the combination of a microbiome
413 model with other clinical tests and patient risk factors could substantially improve this diagnostic
414 accuracy (**Figure 5D**). Current clinical pre-colonoscopy screening tests (e.g. FOBT, WIF-1) remain
415 cheaper, but the microbiome-based CRC prediction models enable a very high diagnostic potential
416 which increases with the number of microbes or microbial genes used, with single biomarkers being
417 much inferior to multi-featured diagnostic models. However, nearly maximal accuracy was achieved
418 with as few as 15 to 25 microbes (**Figure 3B-C**) or a few hundred genes (**Extended Data 8**),
419 potentially enabling inexpensive clinical microbiological tests to be performed on stool. Prospective
420 studies of these biomarkers are needed to establish whether they can identify individuals at elevated
421 risk of CRC and provide the possibility of disease prevention.

422 The diversity and subject-specificity of the human gut microbiome is not yet fully uncovered, with
423 many microbial genes having unknown function, and with strain-level diversity that is missed by
424 many current analysis pipelines⁵⁰. Large scale shotgun metagenomics can begin to overcome these
425 limitations, as shown here by the novel identification of a link between CRC and the microbial
426 pathway producing trimethylamine from choline⁴⁸. The gene encoding for the key enzyme for this
427 pathway, the CutC choline TMA-lyase, is both more overall abundant and expressed in the gut
428 microbiomes of carcinoma patients, with specific variants of *cutC* characterizing controls, adenomas,
429 and carcinomas (**Figure 4**). TMA-producing choline lyases have been found to be associated with
430 atherosclerosis⁹, and higher plasma trimethylamine oxide and choline levels have been reported to

431 be correlated with CRC risk^{59,60}. We highlighted the importance of strain-level gene resolution in
432 understanding any potential carcinogenic role of *cutC*. CRC-associated variants mostly originated
433 from *Hungatella hathewayi*, *Clostridium asparagiforme*, *Klebsiella oxytoca*, and *Escherichia coli*,
434 whereas no significant enrichment was detected for a *cutC* variant carried by a unexplored recently
435 discovered candidate species in the *Lachnospiraceae* family⁵⁰. Thus, genetic variants in key microbial
436 genes involved in choline-induced TMA production by the gut microbiome are a plausible and novel
437 potential mechanism for colorectal carcinogenesis. Other partially diet-dependent microbiome
438 factors can contribute to promote carcinogenesis, and we found in our parallel work that genes for
439 secondary bile acid conversion are consistently enriched in the CRC-associated microbiomes²⁹.
440 Further work is needed to establish the changes in protein structure and function associated with the
441 genetic variants of the diet-related microbial genes found here to be enriched in the CRC microbiome.

442 Analysis of cancer cohorts that are heterogeneous for geography, ethnicity, and lifestyle, presents a
443 distinct opportunity for studying the cancer-associated microbiome. By combining multiple small
444 cohorts of potentially low generalizability, it is possible to obtain better representation of the
445 spectrum of cancer cases and controls. With appropriate methodology, artifactual findings due to
446 batch effects present in any individual dataset can be avoided. The use of large, diverse training sets
447 enables creation of more accurate diagnostic models, and the availability of independent validation
448 datasets enables more realistic estimation of that accuracy. Future shotgun metagenomic studies of
449 the intestinal mucosa-associated microbiome, which are currently infeasible due to excessive human
450 DNA contamination²⁸, will be important to further refine the list of CRC-associated gut microbes.
451 Nevertheless, this study identifies highly reproducible microbial CRC biomarkers and points to the
452 potential for non-invasive microbial diagnostic tests to supplement existing screening.

453

454 **Acknowledgements**

455 We thank the members of the Segata, Naccarati, and Waldron groups for insightful discussions, all the
456 volunteers enrolled in the study, the NGS facility at University of Trento for performing the
457 metagenomic sequencing, and the HPC facility at University of Trento for supporting the
458 computational experiments. This work was primarily supported by Lega Italiana per La Lotta contro i
459 Tumori to N.S., F.C. and A.N. and by Fundação de Amparo à Pesquisa do Estado de São Paulo (FAPESP -
460 16/23527-2) to A.M.T. This work was also partially supported by the Conselho Nacional de Pesquisa
461 e Desenvolvimento (CNPq, Brazil) to J.C.S. and E.D.-N., FAPESP (14/26897-0), Associação Beneficente
462 Alzira Denise Hertzog Silva (ABADHS, Brazil) and PRONON/SIPAR to E.D.-N., by Coordenação de
463 Aperfeiçoamento de Pessoal de Nível Superior - Brazil (CAPES) - Finance Code 001 to J.C.S., by the
464 Italian Institute for Genomic Medicine (IIGM) and Compagnia di San Paolo Torino to A.N, A.F., B.P. and
465 S.T., by Fondazione Umberto Veronesi “Post-doctoral fellowship Year 2014, 2015, 2016, 2017 and
466 2018” to B.P. and S.T., by the Grant Agency of the Czech Republic (17-16857S) to A.N., by Fondazione
467 Umberto Veronesi (FUV-14-SG-GANDINI) to S.G., by the European Union H2020 Marie-Curie grant
468 (707345) to E.P., by the European Research Council (ERC-STG project MetaPG) to N.S., by MIUR
469 “Futuro in Ricerca” RBFR13EWWI_001 to N.S., by the People Programme (Marie Curie Actions) of the
470 European Union FP7 and H2020 to N.S., and by the National Cancer Institute (U24CA180996) and
471 National Institute of Allergy and Infectious Diseases (1R21AI121784-01) of the National Institutes of
472 Health to L.W. B.P. is recipient of a Fulbright Research Scholarships (year 2018). We acknowledge
473 funding from EMBL, DKFZ, the Huntsman Cancer Foundation, the Intramural Research Program of
474 the National Cancer Institute, ETH Zürich, and the following external sources: the European Research

475 Council (CancerBiome ERC-2010-AdG_20100317 to P.B., Microbios ERC-AdG-669830 to P.B.), the
476 Novo Nordisk Foundation (grant NNF10CC1016515 to M.A.), the Danish Diabetes Academy
477 supported by the Novo Nordisk Foundation and TARGET research initiative (Danish Strategic
478 Research Council [0603-00484B] to M.A.), the Matthias-Lackas Foundation (to C.M.U.), the National
479 Cancer Institute (grants R01 CA189184, R01 CA207371, U01 CA206110, P30 CA042014 II to C.M.U.),
480 the BMBF (the de.NBI network #031A537B to P.B. and the ERA-NET TRANSCAN project 01KT1503 to
481 C.M.U.), and the Helmut Horten Foundation (to S.Sunagawa). For the Validation Cohort2, funding was
482 provided by grants from the National Cancer Center Research and Development Fund (25-A-4,28-A-4,
483 and 29-A-6), Practical Research Project for Rare/Intractable Diseases from the Japan Agency for
484 Medical Research and Development (AMED) (JP18ek0109187), JST (Japan Science and Technology
485 Agency)-PRESTO (JPMJPR1507), JSPS (Japan Society for the Promotion of Science) KAKENHI
486 (16J10135, 142558 and 221S0002), Joint Research Project of the Institute of Medical Science, the
487 University of Tokyo, and the Takeda Science Foundation and Suzuken Memorial Foundation.

488

489 **Author contributions**

490 N.S., A.M.T., L.W., and A.N conceived the study. N.S. supervised the study. C.P., S.G., D.S., S.T., A.F., G.G.,
491 M.T., B.P, M.R., and A.N. organized the clinical study, recruited patients and collected samples.
492 F.Armanini generated metagenomic data. A.M.T., P.M., F.Asnicar, E.P., M.Z., F.B., N.K., and G.F. collected
493 and analyzed the metagenomic data. A.M.T., P.M., F.Asnicar, E.P., M.Z., G.F., J.W., G.Z., and L.W.
494 performed machine learning and statistical analyses. F.Armanini, S.T., S.Manara, A.T., B.P, and A.N.
495 performed validation experiments. S.Mizutani., H.S., S.Shiba, T.S., S.Y., T.Y., J.W., P.S.-K, C.M.U., H.B.,
496 M.A., P.B., and G.Z. provided additional validation data. A.M.T., P.M., L.W., and N.S. designed and
497 produced the figures. A.M.T., P.M., and N.S. wrote the manuscript with contributions from S.Manara,
498 F.C., E.D.-N., J.C.S., M.R., L.W., and A.N. All authors discussed and approved the manuscript.

499

500 **Competing Interests**

501 P. Bork, G. Zeller, A.Y. Voigt, and S. Sunagawa are named inventors on a patent (EP2955232A1:
502 Method for diagnosing colorectal cancer based on analyzing the gut microbiome). All the other
503 authors declare to have no competing interests as defined by Nature Research, or other interests that
504 might be perceived to influence the results and/or discussion reported in this paper.

505

506 **References (for main text only)**

- 507 1. Ferlay, J. *et al.* Cancer incidence and mortality worldwide: sources, methods and major
508 patterns in GLOBOCAN 2012. *Int. J. Cancer* **136**, E359–86 (2015).
- 509 2. Siegel, R., Desantis, C. & Jemal, A. Colorectal cancer statistics, 2014. *CA Cancer J. Clin.* **64**, 104–
510 117 (2014).
- 511 3. Frank, C., Sundquist, J., Yu, H., Hemminki, A. & Hemminki, K. Concordant and discordant
512 familial cancer: Familial risks, proportions and population impact. *Int. J. Cancer* **140**, 1510–
513 1516 (2017).
- 514 4. Foulkes, W. D. Inherited susceptibility to common cancers. *N. Engl. J. Med.* **359**, 2143–2153
515 (2008).

- 516 5. Johnson, C. M. *et al.* Meta-analyses of colorectal cancer risk factors. *Cancer Causes Control* **24**,
517 1207–1222 (2013).
- 518 6. Huxley, R. R. *et al.* The impact of dietary and lifestyle risk factors on risk of colorectal cancer: a
519 quantitative overview of the epidemiological evidence. *Int. J. Cancer* **125**, 171–180 (2009).
- 520 7. Schmidt, T. S. B., Raes, J. & Bork, P. The Human Gut Microbiome: From Association to
521 Modulation. *Cell* **172**, 1198–1215 (2018).
- 522 8. Thomas, R. M. & Jobin, C. The Microbiome and Cancer: Is the ‘Oncobiome’ Mirage Real? *Trends*
523 *Cancer Res.* **1**, 24–35 (2015).
- 524 9. Jie, Z. *et al.* The gut microbiome in atherosclerotic cardiovascular disease. *Nat. Commun.* **8**,
525 845 (2017).
- 526 10. Pasolli, E., Truong, D. T., Malik, F., Waldron, L. & Segata, N. Machine Learning Meta-analysis of
527 Large Metagenomic Datasets: Tools and Biological Insights. *PLoS Comput. Biol.* **12**, e1004977
528 (2016).
- 529 11. Cougnoux, A. *et al.* Bacterial genotoxin colibactin promotes colon tumour growth by inducing
530 a senescence-associated secretory phenotype. *Gut* **63**, 1932–1942 (2014).
- 531 12. Wu, S. *et al.* A human colonic commensal promotes colon tumorigenesis via activation of T
532 helper type 17 T cell responses. *Nat. Med.* **15**, 1016–1022 (2009).
- 533 13. Chung, L. *et al.* Bacteroides fragilis Toxin Coordinates a Pro-carcinogenic Inflammatory
534 Cascade via Targeting of Colonic Epithelial Cells. *Cell Host Microbe* **23**, 203–214.e5 (2018).
- 535 14. Kostic, A. D. *et al.* Genomic analysis identifies association of Fusobacterium with colorectal
536 carcinoma. *Genome Res.* **22**, 292–298 (2012).
- 537 15. Kostic, A. D. *et al.* Fusobacterium nucleatum potentiates intestinal tumorigenesis and
538 modulates the tumor-immune microenvironment. *Cell Host Microbe* **14**, 207–215 (2013).
- 539 16. Rubinstein, M. R. *et al.* Fusobacterium nucleatum promotes colorectal carcinogenesis by
540 modulating E-cadherin/ β -catenin signaling via its FadA adhesin. *Cell Host Microbe* **14**, 195–
541 206 (2013).
- 542 17. Yu, J. *et al.* Metagenomic analysis of faecal microbiome as a tool towards targeted non-
543 invasive biomarkers for colorectal cancer. *Gut* **66**, 70–78 (2017).
- 544 18. Feng, Q. *et al.* Gut microbiome development along the colorectal adenoma–carcinoma
545 sequence. *Nat. Commun.* **6**, 6528 (2015).
- 546 19. Zeller, G. *et al.* Potential of fecal microbiota for early-stage detection of colorectal cancer. *Mol.*
547 *Syst. Biol.* **10**, 766 (2014).
- 548 20. Baxter, N. T., Ruffin, M. T., Rogers, M. A. M. & Schloss, P. D. Microbiota-based model improves
549 the sensitivity of fecal immunochemical test for detecting colonic lesions. *Genome Med.* **8**, 37
550 (2016).
- 551 21. Zackular, J. P., Rogers, M. A. M., Ruffin, M. T., 4th & Schloss, P. D. The human gut microbiome as
552 a screening tool for colorectal cancer. *Cancer Prev. Res.* **7**, 1112–1121 (2014).
- 553 22. Drewes, J. L. *et al.* High-resolution bacterial 16S rRNA gene profile meta-analysis and biofilm
554 status reveal common colorectal cancer consortia. *NPJ Biofilms Microbiomes* **3**, 34 (2017).
- 555 23. Pollock, J., Glendinning, L., Wisedchanwet, T. & Watson, M. The Madness of Microbiome:
556 Attempting To Find Consensus ‘Best Practice’ for 16S Microbiome Studies. *Appl. Environ.*
557 *Microbiol.* **84**, (2018).
- 558 24. Segata, N. On the Road to Strain-Resolved Comparative Metagenomics. *mSystems* **3**, (2018).
- 559 25. Truong, D. T., Tett, A., Pasolli, E., Huttenhower, C. & Segata, N. Microbial strain-level
560 population structure and genetic diversity from metagenomes. *Genome Res.* **27**, 626–638
561 (2017).
- 562 26. Pasolli, E. *et al.* Accessible, curated metagenomic data through ExperimentHub. *Nat. Methods*
563 **14**, 1023 (2017).
- 564 27. Dai, Z. *et al.* Multi-cohort analysis of colorectal cancer metagenome identified altered bacteria
565 across populations and universal bacterial markers. *Microbiome* **6**, 70 (2018).
- 566 28. Quince, C., Walker, A. W., Simpson, J. T., Loman, N. J. & Segata, N. Shotgun metagenomics, from

- 567 sampling to analysis. *Nat. Biotechnol.* **35**, 833–844 (2017).
- 568 29. Wirbel, J. *et al.* Meta-analysis of fecal metagenomes reveals global microbial signatures that
569 are specific for colorectal cancer. *Under Submission* (2018).
- 570 30. Hannigan, G. D., Duhaime, M. B., Ruffin, M. T., 4th, Koumpouras, C. C. & Schloss, P. D. Diagnostic
571 Potential and Interactive Dynamics of the Colorectal Cancer Virome. *MBio* **9**, (2018).
- 572 31. Thomas, A. M. *et al.* Tissue-Associated Bacterial Alterations in Rectal Carcinoma Patients
573 Revealed by 16S rRNA Community Profiling. *Front. Cell. Infect. Microbiol.* **6**, 179 (2016).
- 574 32. Gao, Z., Guo, B., Gao, R., Zhu, Q. & Qin, H. Microbiota disbiosis is associated with colorectal
575 cancer. *Front. Microbiol.* **6**, 20 (2015).
- 576 33. Ahn, J. *et al.* Human gut microbiome and risk for colorectal cancer. *J. Natl. Cancer Inst.* **105**,
577 1907–1911 (2013).
- 578 34. Flemer, B. *et al.* The oral microbiota in colorectal cancer is distinctive and predictive. *Gut*
579 (2017). doi:10.1136/gutjnl-2017-314814
- 580 35. Brito, I. L. *et al.* Mobile genes in the human microbiome are structured from global to
581 individual scales. *Nature* **535**, 435–439 (2016).
- 582 36. Human Microbiome Project Consortium. Structure, function and diversity of the healthy
583 human microbiome. *Nature* **486**, 207–214 (2012).
- 584 37. Bonder, M. J. *et al.* The effect of host genetics on the gut microbiome. *Nat. Genet.* **48**, 1407
585 (2016).
- 586 38. Segata, N. *et al.* Metagenomic biomarker discovery and explanation. *Genome Biol.* **12**, R60
587 (2011).
- 588 39. Xie, Y.-H. *et al.* Fecal Clostridium symbiosum for Noninvasive Detection of Early and Advanced
589 Colorectal Cancer: Test and Validation Studies. *EBioMedicine* **25**, 32–40 (2017).
- 590 40. Boleij, A., van Gelder, M. M. H. J., Swinkels, D. W. & Tjalsma, H. Clinical Importance of
591 Streptococcus gallolyticus infection among colorectal cancer patients: systematic review and
592 meta-analysis. *Clin. Infect. Dis.* **53**, 870–878 (2011).
- 593 41. Fijan, S. Microorganisms with claimed probiotic properties: an overview of recent literature.
594 *Int. J. Environ. Res. Public Health* **11**, 4745–4767 (2014).
- 595 42. Apweiler, R. *et al.* UniProt: the Universal Protein knowledgebase. *Nucleic Acids Res.* **32**, D115–
596 9 (2004).
- 597 43. Gerner, E. W. & Meyskens, F. L., Jr. Polyamines and cancer: old molecules, new understanding.
598 *Nat. Rev. Cancer* **4**, 781–792 (2004).
- 599 44. Costea, P. I. *et al.* Towards standards for human fecal sample processing in metagenomic
600 studies. *Nat. Biotechnol.* **35**, 1069–1076 (2017).
- 601 45. Riester, M. *et al.* Risk prediction for late-stage ovarian cancer by meta-analysis of 1525 patient
602 samples. *J. Natl. Cancer Inst.* **106**, (2014).
- 603 46. Kummen, M. *et al.* Elevated trimethylamine-N-oxide (TMAO) is associated with poor
604 prognosis in primary sclerosing cholangitis patients with normal liver function. *United*
605 *European Gastroenterol J* **5**, 532–541 (2017).
- 606 47. Oellgaard, J., Winther, S. A., Hansen, T. S., Rossing, P. & von Scholten, B. J. Trimethylamine N-
607 oxide (TMAO) as a New Potential Therapeutic Target for Insulin Resistance and Cancer. *Curr.*
608 *Pharm. Des.* **23**, 3699–3712 (2017).
- 609 48. Kalnins, G. *et al.* Structure and Function of CutC Choline Lyase from Human Microbiota
610 Bacterium *Klebsiella pneumoniae*. *J. Biol. Chem.* **290**, 21732–21740 (2015).
- 611 49. Rath, S., Heidrich, B., Pieper, D. H. & Vital, M. Uncovering the trimethylamine-producing
612 bacteria of the human gut microbiota. *Microbiome* **5**, 54 (2017).
- 613 50. Pasolli, E. *et al.* Extensive unexplored human microbiome diversity revealed by over 150,000
614 genomes from metagenomes spanning age, geography, and lifestyle. *Cell* **176**, 1–14 (2019).
- 615 51. Hothorn, T., Hornik, K., van de Wiel, M. A. & Zeileis, A. A Lego System for Conditional
616 Inference. *Am. Stat.* **60**, 257–263 (2006).
- 617 52. Nielsen, H. B. *et al.* Identification and assembly of genomes and genetic elements in complex

- 618 metagenomic samples without using reference genomes. *Nat. Biotechnol.* **32**, 822–828 (2014).
619 53. Karlsson, F. H. *et al.* Gut metagenome in European women with normal, impaired and diabetic
620 glucose control. *Nature* **498**, 99–103 (2013).
621 54. Qin, J. *et al.* A metagenome-wide association study of gut microbiota in type 2 diabetes. *Nature*
622 **490**, 55–60 (2012).
623 55. Integrative HMP (iHMP) Research Network Consortium. The Integrative Human Microbiome
624 Project: dynamic analysis of microbiome-host omics profiles during periods of human health
625 and disease. *Cell Host Microbe* **16**, 276–289 (2014).
626 56. Dejea, C. M. *et al.* Patients with familial adenomatous polyposis harbor colonic biofilms
627 containing tumorigenic bacteria. *Science* **359**, 592–597 (2018).
628 57. Manichanh, C. *et al.* Reduced diversity of faecal microbiota in Crohn’s disease revealed by a
629 metagenomic approach. *Gut* **55**, 205–211 (2006).
630 58. Le Chatelier, E. *et al.* Richness of human gut microbiome correlates with metabolic markers.
631 *Nature* **500**, 541–546 (2013).
632 59. Bae, S. *et al.* Plasma choline metabolites and colorectal cancer risk in the Women’s Health
633 Initiative Observational Study. *Cancer Res.* **74**, 7442–7452 (2014).
634 60. Xu, R., Wang, Q. & Li, L. A genome-wide systems analysis reveals strong link between
635 colorectal cancer and trimethylamine N-oxide (TMAO), a gut microbial metabolite of dietary
636 meat and fat. *BMC Genomics* **16 Suppl 7**, S4 (2015).
637

638 **Figure legends (for main text only)**

639 **Figure 1. Reproducible taxonomic and functional microbial biomarkers across datasets when**
640 **contrasting carcinoma against healthy controls (no adenoma samples considered).** (A) UpSet
641 plot showing the number of taxonomic biomarkers identified using LEfSE on MetaPhlAn2 species
642 profiles shared by combinations of datasets (see **Suppl. Table 3** for all single significant
643 associations). (B) Pooled effect sizes for the 20 significant features with the largest effect size
644 calculated using a meta-analysis of standardized mean differences and a random effects model on
645 MetaPhlAn2 species abundances and on (C) HUMANN2 pathway abundances. Bold lines represent
646 the 95% confidence interval for the random effects model coefficient estimate. (D) Scatter plot of
647 crude and age-, sex-, and BMI-adjusted coefficients obtained from linear models using MetaPhlAn2
648 species abundances. (E) Scatter plot of crude and age-, sex-, and BMI-adjusted coefficients obtained
649 from linear models using HUMANN2 pathway abundances.
650

651 **Figure 2. Assessment of prediction performances of the gut microbiome for CRC detection**
652 **within and across cohorts.** (A) Cross prediction matrix reporting prediction performances as AUC
653 values obtained using a Random Forest (RF) model on species-level relative abundances (see
654 **Methods**). Values on the diagonal refer to 20 times repeated 10-fold stratified cross validations. Off-
655 diagonal values refer to the AUC values obtained by training the classifier on the dataset of the
656 corresponding row and applying it on the dataset of the corresponding column. The Leave-One-
657 Dataset-Out (LODO) row refers to the performances obtained by training the model on the species-
658 level abundances and MetaPhlAn2 markers using all but the dataset of the corresponding column and
659 applying it on the dataset of the corresponding column. See **Extended Data 6** for the marker cross-
660 study validation matrix. (B) Cross prediction matrix of AUC values obtained using HUMANN2
661 UniRef90 gene-family abundances and HUMANN2 pathway relative abundances. See **Extended Data**
662 **6** for the pathway cross-study validation matrix. (C) Prediction performances for the two Italian
663 cohorts at increasing numbers of external datasets considered for training the model. The dark
664 yellow line interpolates the median AUC at each number of training datasets considered. See
665 **Extended Data 7** for the plots referred to the other testing datasets. (D) Prediction performances at

666 increasing number of datasets in the training, using HUMANn2 UniProt90 gene-family abundances.
 667 See **Extended Data 7** for the plots referred to the other testing datasets.
 668

669 **Figure 3. Ranking relevance of each species in the predictive models for each dataset and**
 670 **identification of a minimal microbial signature for CRC detection. (A)** The importance of each
 671 species for the cross-validation prediction performance in each dataset estimated using the internal
 672 RF scores. Only species appearing in the five top ranking features in at least one dataset are reported.
 673 Prediction performances at increasing number of microbial species obtained by re-training the RF
 674 classifier on the N top ranked features identified with a first RF model training in a cross-validation
 675 **(B)** and LODO-setting **(C)**. The rankings are obtained excluding the testing dataset to avoid
 676 overfitting.
 677

678 **Figure 4. Choline TMA-lyase gene *cutC* and its genetic variants are strong biomarkers for CRC-**
 679 **associated stool samples. (A)** Distribution of reads per kilobase million (RPKM) abundances
 680 obtained using ShortBRED for the choline TMA-lyase enzyme gene *cutC*. P-values were computed by
 681 two-tailed Wilcoxon Signed-Rank tests comparing values between controls and carcinomas for each
 682 dataset. **(B)** Forest plot reporting effect sizes calculated using a meta-analysis of standardized mean
 683 differences and a random effects model on *cutC* RPKM abundances between carcinomas and controls.
 684 **(C)** Phylogenetic tree of sample-specific *cutC* sequence variants identified four main sequence
 685 variants. Tips with no circles represent *cutC* sequence variants from genomes absent from the
 686 datasets. Taxonomy was assigned based on mapping against existing *cutC* sequences (criteria of 80%
 687 coverage, >97% identity and minimum 2,000nt alignment length). **(D)** qPCR validation of *cutC* gene
 688 abundance and **(E)** *cutC* transcript abundance (normalized by total 16S rRNA gene/transcript
 689 abundance) on a subset of DNA samples from Cohort1. qPCR validation P-values are obtained by 1-
 690 tail Wilcoxon Signed-Rank test.
 691

692 **Figure 5 - Clinical potential and validation of the predictive biomarkers. (A)** Prediction
 693 performance of the taxonomic models trained on the 7 datasets of **Table 1** and applied on the new
 694 validation cohorts confirmed the strong reproducibility of metagenomic models for CRC across
 695 cohorts when sufficiently large training cohorts are available. Feature ranking of the 16-species
 696 model are obtained the testing cohort to avoid overfitting. **(B)** Species richness, rarefied oral species
 697 richness, and *cutC* gene abundance (RPKM) are confirmed to be strong biomarkers of CRC in the
 698 validation datasets²⁹. P-values underlying the panels refer to one-tailed Wilcoxon Signed Rank test;
 699 P-values overlying the panels refer to the one-sided permutation-based Wilcoxon-Mann-Whitney
 700 tests, blocked for cohort. **(C)** Prediction performances as AUC values on the validation cohorts when
 701 adding external set of case and controls samples from metagenomic cohorts of diseases other than
 702 CRC (Crohn's disease, ulcerative colitis, type-2 diabetes). **(D)** Assessment of the potential of
 703 microbiome-based prediction models in comparison and in combination with current non-invasive
 704 clinical screening tests. Models integrating our LODO machine learning approach with the FOBT or
 705 the Wif-1 Methylation tests are termed OR and AND, depending on whether only one or both need to
 706 be positive for the combined test to be positive.
 707

708 Tables

709 **Table 1.** Size and characteristics of the large-scale CRC metagenomic datasets included in this study.
 710

Dataset	Groups (N)	Age (mean +/- s.d.)	BMI (mean +/- s.d.)	Sex F(%) / M(%)	Country	# of reads (x 10 ⁹)

ZellerG_2014 (Zeller et al. 2014)	Control (61) Adenoma (42) CRC (53)	60.6 +/- 11.4 63 +/- 9.1 66.8 +/- 10.9	24.7 +/- 3.2 25.9 +/- 4.1 25.5 +/- 5.2	54.1/45.9 28.5/71.5 45.2/54.8	France	9.4
YuJ_2015 (Yu et al. 2015)	Control (54) CRC (74)	61.8 +/- 5.7 66 +/- 10.6	23.5 +/- 3 24 +/- 3.2	38.9/61.1 35.1/64.9	China	7.2
FengQ_2015 (Feng et al. 2015)	Control (61)* Adenoma (47) CRC (46)	67 +/- 6.5 66.5 +/- 7.9 67 +/- 10.9	27.6 +/- 3.8 28 +/- 4.7 26.5 +/- 3.5	41/59 51.1/48.9 39.1/60.9	Austria	8.3
VogtmannE_2016 (Vogtmann et al. 2016)	Control (52) CRC (52)	61.2 +/- 11 61.8 +/- 13.6	25.3 +/- 4.2 24.9 +/- 4.2	28.8/71.2 28.8/71.2	USA	6.9
HanniganGD_2018 (Hannigan et al. 2018)	Control (28) Adenoma (27) CRC (27)	NA	NA	NA	USA (54) Canada (28)	0.5
Cohort1 (This study)	Control (24) Adenoma (27) CRC (29)	67.9 +/- 7.1 62.8 +/- 8.6 71.4 +/- 8.2	25.3 +/- 3.5 25.3 +/- 4.1 25.7 +/- 4.1	45.8/54.1 40.7/59.3 20.7/79.3	Italy	8.2
Cohort2 (This study)	Control (28) CRC (32)	57.8 +/- 8.3 58.4 +/- 8.4	24.6 +/- 3.8 26.8 +/- 4.3	42.9/57.1 28.1/71.9	Italy	5.1
Total	Control (308) Adenoma (143) CRC (313)	--	--	--	--	45.6

*Numbers differed from the original sample numbers (N = 61 instead of 63) reported in the article due to metadata and/or sequence processing issues. NA = Not available.

711
712

713

714 **Methods**

715 **Italian cohorts of CRC patients, adenomas and controls**

716 The two clinical studies performed here were approved by the relevant ethics committees (Cohort1:
717 Ethics committee of Azienda Ospedaliera "SS. Antonio e Biagio e C. Arrigo" of Alessandria, Italy,
718 protocol N. Colorectal_miRNA_CEC2014 and Cohort2: Ethics committee of European Institute of
719 Oncology of Milan, Italy, protocol N. R107/14-IEO 118) and informed consent was obtained from all
720 participants.

721 For Cohort1, samples were collected from patients at the Clinica S. Rita in Vercelli, Italy. Patients with
722 hereditary CRC syndromes, with previous history of CRC, and with uncompleted or poorly cleaned
723 colonoscopy, were excluded from the study. Patients were recruited at initial diagnosis and had not
724 received any treatment prior to fecal sample collection. Subjects reporting the use of antibiotics
725 during the 6 months prior to the sample collection were excluded from the study. On the basis of
726 colonoscopy results, recruited subjects were classified into three categories: 1) healthy subjects:
727 individuals with colonoscopy negative for tumor, adenomas and other diseases; 2) adenoma patients:
728 individuals with colorectal adenoma/s; and 3) CRC patients: individuals with newly diagnosed CRC. A
729 total of 93 subjects were initially recruited, and the 80 that passed quality control (see below) are
730 divided into 29 CRC patients, 27 adenomas and 24 controls. An additional 13 subjects that presented
731 inflammatory GI tract diseases (ulcerative and Crohn's colitis, diverticular diseases) were recruited
732 and fecal samples were subsequently used as a part of the final validation. Stool was collected in Stool
733 Nucleic Acid Collection and Transport Tubes with RNA stabilising solution (Norgen Biotek Corp) and
734 returned before performing the colonoscopy. Aliquots of the stool samples were stored at -80°C until
735 use. DNA was extracted from aliquot of fecal samples using the Qiamp DNA stool kit (Qiagen)
736 following manufacturer's instructions. Total RNA from faeces was extracted using the Stool Total
737 RNA Purification Kit (Norgen Biotek Corp) following manufacturer's instructions.

738 For Cohort2, a total of 60 subjects were recruited at the European Oncology Institute in Milan, Italy
739 and were divided into 32 CRC patients and 28 controls. Controls, matched for age (\pm 5 years) and
740 season of blood withdrawn (\pm 2 years), were recruited among subjects who underwent recent
741 colonoscopy and had negative or no other relevant gastrointestinal disorders. Subjects reporting the
742 use of antibiotics in the 6 months prior to the sample collection were excluded. Fecal samples were
743 collected from healthy subjects and patients (before surgery, or any other cancer treatment) and
744 directly frozen at -80°C in resuspension buffer (TES buffer: 50 mM Tris-HCL, 10 mM NaCl, 10 mM
745 EDTA, pH 7.5) and kept in liquid nitrogen until DNA extraction. DNA was extracted from fecal
746 samples with the GNOME DNA isolation kit (MP).

747 Sequencing libraries were prepared using the NexteraXT DNA Library Preparation Kit (Illumina,
748 California, USA), following the manufacturer's guidelines. Sequencing was performed on the
749 HiSeq2500 (Illumina, California, USA) at the internal sequencing facility of the Centre for Integrative
750 Biology, Trento, Italy.

751 **Public metagenomic cohorts of CRC patients, adenomas and controls.**

752 We downloaded 5 public fecal shotgun CRC datasets covering samples from 6 different countries,
753 totaling 313 CRC patients, 143 adenomas and 308 controls (**Table 1**) and now available in
754 curatedMetagenomicData ²⁶. We manually curated metadata tables for the public cohorts according to
755 the curatedMetagenomicData ²⁶ R-package grammatical rules. The metadata table includes ten fields
756 (sampleID, subjectID, body_site, country, sequencing_platform, PMID, number_reads, number_bases,
757 minimum_read_length, median_read_length) that are mandatory for all datasets in addition to other
758 fields that are dataset-specific.

759 **Description of the two validation cohorts**

760 We consider an additional set of samples from two independent cohorts that were not available at the
761 time we performed the meta-analysis on the other seven datasets, and we thus used them as
762 validation cohorts. Validation Cohort1 consists of 60 CRC metagenomes collected in Germany after
763 colonoscopy and 65 sex and age-matched healthy controls and is described in depth in the study
764 accompanying this work ²⁹. Shotgun metagenomic sequencing was performed by Illumina HiSeq 2000
765 / 2500 / 4000 (Illumina, San Diego, USA) platforms at the Genomics Core Facility, European
766 Molecular Biology Laboratory, Heidelberg. Validation Cohort2 consists of 40 CRC samples and 40
767 controls from a Japanese cohort from Tokyo. DNA was extracted for Validation Cohort2 from frozen
768 fecal samples by bead-beating using the GNOME DNA Isolation Kit (MP Biomedicals, Santa Ana, CA)
769 and DNA quality was assessed with an Agilent 4200 TapeStation (Agilent Technologies, Santa Clara
770 CA). Sequencing libraries were generated with a Nextera XT DNA Sample Prep Kit (Illumina, San
771 Diego, CA) and shotgun metagenomics of fecal samples was carried out on the HiSeq2500 platform
772 (Illumina) at a targeted depth of 5.0 Gb (150-bp paired end reads).

773 The samples and clinical information used from both validation cohorts in this study were obtained
774 under conditions of informed consent and with approval of the institutional review boards of each
775 participating institute.

776 **Public metagenomic cohorts of non-CRC patients.**

777 We used the curatedMetagenomicData ²⁶ resource to retrieve taxonomical and functional potential
778 profiles as well as metadata of three public cohorts: NielsenHB_2014 ⁵² comprising 21 Crohn Disease
779 (CD) patients, 127 Ulcerative Colitis (UC) patients and 248 controls; KarlssonFH_2013 ⁵³ comprising
780 53 Type-2 Diabetes (T2D) patients and 43 controls; QinJ_2012 ⁵⁴ comprising 172 T2D patients and

781 174 controls; and we downloaded 1339 metagenomes from the Human Microbiome Consortium
782 phase-2 cohort⁵⁵, comprising 598 Crohn Disease patients, 375 Ulcerative Colitis patients and 365
783 controls.

784 **Sequence pre-processing, taxonomic and functional profiling**

785 Fecal metagenomic shotgun sequences obtained from the Italian cohorts were subjected to a pre-
786 processing pipeline whereby sequences were quality filtered using trim_galore (parameters: --
787 nextera --stringency 5 --length 75 --quality 20 --max_n 2 --trim-n) discarding all reads with quality
788 less than 20 and shorter than 75 nucleotides. Filtered reads were then aligned to the human genome
789 (hg19) and the PhiX genome for human and contaminant DNA removal using bowtie2⁶¹. Thirteen
790 samples, having less than 2Gb of host-decontaminated DNA, were excluded from the study.

791 We used MetaPhlAn2⁶² for quantitative profiling the taxonomic composition of the microbial
792 communities of all metagenomic samples, whereas HUMAnN2⁶³ was used to profile pathway and
793 gene family abundances. The profiles generated for the 6 public cohorts, along with their metadata,
794 and the two newly sequenced cohorts are available through the curatedMetagenomicData R package
795²⁶. Oral species were defined in this work by analyzing the 463 oral samples from the Human
796 Microbiome Project dataset³⁶ and the 140 saliva samples from³⁵. Specifically, all species with > 0.1%
797 abundance and > 5% prevalence were deemed to be of oral origin. For *F. nucleatum* marker analysis,
798 we extracted MetaPhlAn2 clade-specific markers from each sample sam file and considered a marker
799 to be present if the coverage was greater than zero.

800 **The Random Forest based machine learning approach**

801 Our machine learning analyses exploited 4 types of microbiome quantitative profiles: taxonomic
802 species-level relative abundances and marker presence or absence patterns inferred by MetaPhlAn2
803⁶², gene-family and pathway relative abundances estimated by HUMAnN2⁶³.

804 All machine learning experiments used Random Forest⁶⁴, as this algorithm has been shown to
805 outperform, on average, other learning tools for microbiome data¹⁰. The code generating the
806 analyses and the figures is available at
807 https://bitbucket.org/CibioCM/multidataset_machinelearning/src/, and is based on MetAML¹⁰ with
808 the Random Forest implementation taken from Scikit-Learn version 0.19.0,⁶⁵. We used an ensemble
809 of 1000 estimator trees and Shannon entropy to evaluate the quality of a split at each node of a tree.
810 The two hyper-parameters for the minimum number of samples per leaf and for the number of
811 features per tree are set as indicated elsewhere⁶⁶ to 5 and 30% respectively. For the marker
812 presence/absence profiles we used a number of features equal to the square root of the total number
813 of features, and this percentage was further decreased to 1% when using gene-family profiles as they
814 have a substantially higher number of features (> 2M). The experiments ran on reduced sets of input
815 features (**Figure 4, Suppl. Fig. 19**) avoided feature subsampling when less than 128 features were
816 used (**Suppl. Fig. 19**).

817 **Application and evaluation of the learning models**

818 The inside-dataset prediction capability was measured through 10-fold cross-validation, stratified so
819 each fold contained a balanced proportion of positive and negative cases. The procedure of forming
820 the folds and assessing the models was repeated 20 times. The final result is therefore an average
821 over 200 validation folds. In the cross-study validation, datasets are considered two by two: one is
822 used for training the model, the other to validate.

823 The Leave-one-dataset-out (LODO) approach consists of training the model on the pooled samples
824 from all cohorts except the one used for model testing. This mimics the scenario in which all the
825 available samples from multiple cohorts are used to predict CRC-positive samples in a newly
826 established cohort. As a part of the meta-analysis, we iterated along all the cohorts, performing a
827 LODO validation on each set of samples (**Figure 2**).

828 **Additional validation experiments on independent datasets and other diseases**

829 We built a validation LODO model trained on MetaPhlAn2 taxonomic abundances from the previously
830 described set of 7 cohorts and applied it to the independent validation cohorts. To test the
831 performance of the model when challenged with other diseases, we selected 4 metagenomic cohorts
832 ⁵²⁻⁵⁵ covering 3 non-CRC diseases (ulcerative colitis - UC, Crohn's disease - CD, and type-2 diabetes -
833 T2D) and we used them for further experiments. For each disease (UC, CD, T2D) in each dataset, we
834 randomly drawn 60 samples from the control class as well as 60 samples from the cases and added
835 them to each validation dataset in turn, labelled as controls. The random selection was repeated ten
836 times, and the validation AUC computed on the model's prediction accordingly. The rationale is to
837 observe the decrease in AUC when the external cases are added to the controls of the validation
838 cohort with respect the addition of healthy controls.

839 Specificity of the prediction model was also assessed by the addition of 13 IBD samples to Cohort1:
840 we used the 13 samples either as controls for Cohort1 or added to the original controls; we
841 performed a cross-validation and a LODO on Cohort1 (no validation cohorts in the training) using
842 MetaPhlAn2 microbial species.

843 To assess the prediction ability of our Random Forest approach with respect to more traditional non-
844 invasive tests like the FOBt and the Wif-1 Methylation test, we recorded the true positive rate
845 (sensitivity) and the false positive rate (1 - specificity) for a subset of the ZellerG_2014 cohort
846 according to these two tests and one-hundred positive detection thresholds in the case of Random
847 Forest models. We then combined the Random Forest approach with the two tests in turn, first
848 assigning the positive class when both predictors are positive ("AND" model) secondly when just one
849 predictor is ("OR" model).

850 **Statistical analysis**

851 Univariate analyses on a per dataset basis was performed using LEfSe ³⁸ to identify features that were
852 statistically different among groups and estimate their effect size. ANCOM was also applied ⁶⁷ but
853 showed reduced power on our datasets (e.g. it identified *F. nucleatum* as a biomarker in only one
854 dataset) probably due to the low relative abundance of CRC biomarkers that are thus only minimally
855 affected by the problem of compositionality. For these reasons, we chose to use LEfSe for the
856 univariate analysis and focused on the biomarkers with the highest effect size. To overcome the
857 limitations of univariate statistics, we performed multivariate analysis using linear models fitted to
858 the data using the limma R package ⁶⁸ and possible confounders such as age, sex and BMI were
859 included in the models. For the meta-analysis on taxonomic and functional profiles, we converted
860 relative abundances to arcsine-square root transformed proportions and used the *escalc* function
861 from the R metafor package that employed Cohen's standardized mean difference statistic to
862 calculate random effects model estimates. We quantified study heterogeneity using the I^2 estimate
863 (percentage of variation reflecting true heterogeneity) as well as Cochran's Q test to assess
864 statistically significant heterogeneity. P-values obtained from the random effects models were
865 corrected for multiple hypothesis testing correction using the Benjamini-Hochberg procedure and
866 corrected $P < 0.05$ were considered statistically significant. Cluster analysis was conducted by

867 calculating distance matrices from phylogenetic trees using the APE R-package, clustering using
868 partitioning around medoids (PAM) and computing clusters' prediction strength using the cluster R-
869 package. When validating differential species richness, oral-species richness, and increased
870 abundance of the *cutC* gene, we also assessed significance through one-sided permutation-based
871 Wilcoxon-Mann-Whitney tests where we blocked for cohort ⁵¹, as implemented in the 'coin' R-
872 package. The lower and upper hinges of boxplots presented in the figures correspond to the 25th and
873 75th percentiles. The upper and lower whiskers extend from the hinges to the largest (or smallest)
874 value no further than 1.5 * inter-quartile range (IQR) from the hinge, defined as the distance between
875 the 25th and 75th percentiles. Data beyond the end of the whiskers are plotted individually.

876 **Identification and quantification of the genes encoding TMA producing enzymes**

877 In order to obtain a more comprehensive database of choline TMA-lyase enzyme sequences, we
878 downloaded amino acid sequences that matched the keywords "*cutC*" and "*cutD*" from UniProt90 ⁴²,
879 mapped their IDs to EMBL CDS using UniParc and used the resulting DNA sequences to search, using
880 BLASTn ⁶⁹, all 48,902 Prokka ⁷⁰ annotated genomes available in our repository ⁷¹. Matching queries
881 were filtered to include only alignments with >80% identity and length > 1000nt for *cutC* and > 800nt
882 for *cutD*, and an e-value < 1e-15. We used ShortBRED ⁷² to identify short seed sequences that were
883 representative of the filtered queries using UniProt's UniRef100 database and quantified them in the
884 metagenomes, normalizing by the number of reads per kilobase million (RPKM). The pipeline was
885 also applied to identify and quantify the L-carnitine/gamma-butyrobetaine antiporter (*caiT*) and the
886 dioxygenase *yeaW*, responsible for producing TMA preferentially via carnitine degradation. In order
887 to investigate differences in *cutC* sequence types, we clustered *cutC* sequences at 97% sequence
888 identity using UCLUST ⁷³ and aligned raw reads to the clustered *cutC* database using bowtie2 ⁶¹. From
889 the bam files we calculated the breadth and depth of each sequence and generated their
890 corresponding consensus sequence using Samtools ⁷⁴ and VCF utils ⁷⁵. We chose the representative
891 *cutC* sequence for each sample as the one with the highest breadth or the highest depth, if there were
892 multiple *cutC* sequences with the same breadth. We filtered representative *cutC* sequences from each
893 sample to include only those with a breadth > 80%, aligned them using MAFFT ⁷⁶, built a phylogenetic
894 tree using fastTree ⁷⁷ which was refined using RAxML ⁷⁸ and visualized using GraPhlAn ⁷⁹.

895 **Validation of *cutC* gene and transcript abundances by qPCR**

896 Real time qPCR was used to assess differences in *cutC* genes and transcripts between CRC samples
897 and controls. We used a previously described protocol ⁴⁹ which employs 16S rRNA abundances as an
898 internal sample normalization. For first strand cDNA synthesis, 400 ng of RNA templates were
899 retrotranscribed using the High-capacity cDNA Reverse Transcription Kits with Random Primers
900 (Thermofisher Scientific) following the manufacturer's instructions. The *cutC* and 16S rRNA genes
901 (and transcripts from cDNA) were amplified using degenerate primers and cycling conditions as
902 described previously ⁴⁹. Briefly, reactions were performed in triplicate with 10 ng of template DNA or
903 30 ng of cDNA on the Rotor Gene Q (QIAGEN) using HOT FIREPol EvaGreen qPCR mix (SOLIS
904 BIODYNE) with a final primer concentration of 0.5 μ M (16S) or 0.75 μ M (*cutC*). Cycling conditions
905 were as follows: initial denaturation of 95°C for 15 min; followed by 40 cycles of denaturing at 95°C
906 for 45 s, annealing at 57° C (*cutC*) or 55°C for (16S) for 45 s and an extension step of 72°C for 45 s.
907 Melting curves were subsequently performed for all reactions using the following program: 95° for 5
908 s, followed by 65°C for 60s, and a final continuous reading step of seven acquisitions per second
909 between 65 and 97°C.

910 Quantification of the *cutC* gene by means of qPCR protocol was applied to 44 samples belonging to
911 Cohort1 for which enough DNA was available. Samples for which either the *cutC* or the 16S rRNA
912 amplification failed were removed and we retained measurements for a total of 16 CRC and 19
913 control samples. Relative gene fold change was calculated by applying the $\Delta\Delta C_t$ method⁸⁰, with ΔC_t
914 calculated as difference between *cutC* and 16S rRNA C_t values. Significance of the *cutC* vs. 16S rRNA
915 comparison was assessed through the one-tailed Wilcoxon Signed Rank test. The same procedure
916 was applied on the quantification of *cutC* and 16S rRNA transcripts from cDNA, which was computed
917 using 26 CRC and 20 control samples for which we obtained a reliable quantification of both *cutC* and
918 16S rRNA.

919 **Data Availability**

920 Nucleotide sequences for the two new Italian cohorts are available in the Sequence Read Archive
921 (SRA) under the accession number SRP136711. MetaPhlAn2 and HUMAnN2 profiles for the new
922 cohorts were also added to the curatedMetagenomicData R package along with their corresponding
923 metadata. Validation Cohort1 is available in the European Nucleotide Archive (ENA) under the study
924 identifier PRJEB27928, Validation Cohort2 is available in the DDBJ databases under the accession
925 number DRA006684.

926

927 **Methods-only References**

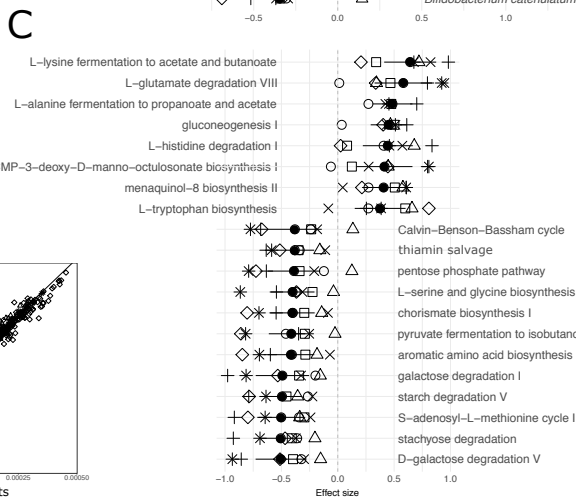
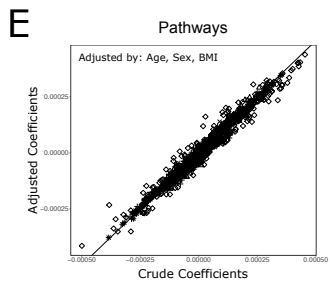
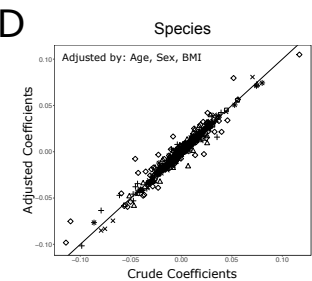
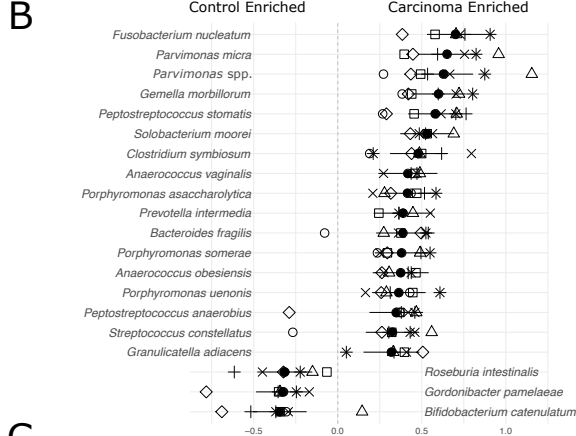
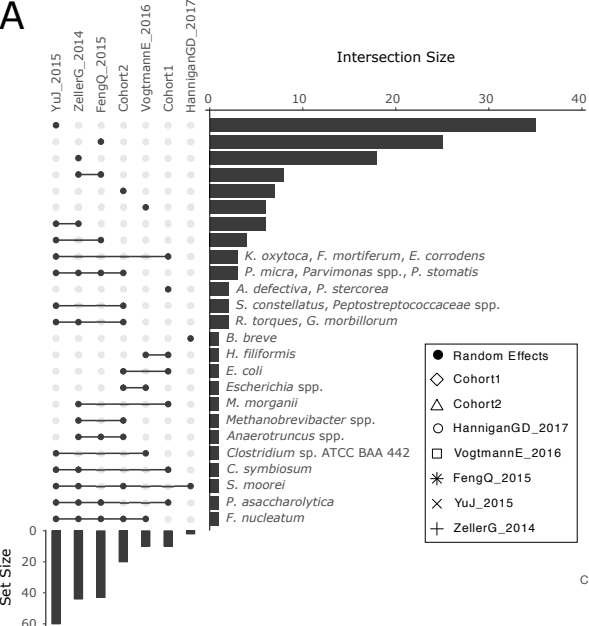
928

- 929 61. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**, 357
930 (2012).
- 931 62. Truong, D. T. *et al.* MetaPhlAn2 for enhanced metagenomic taxonomic profiling. *Nat. Methods*
932 **12**, 902–903 (2015).
- 933 63. Abubucker, S. *et al.* Metabolic reconstruction for metagenomic data and its application to the
934 human microbiome. *PLoS Comput. Biol.* **8**, e1002358 (2012).
- 935 64. Breiman, L. Random Forests. *Mach. Learn.* **45**, 5–32 (2001).
- 936 65. Pedregosa, F. *et al.* Scikit-learn: Machine Learning in Python. *J. Mach. Learn. Res.* **12**, 2825–
937 2830 (2011).
- 938 66. Hastie, T., Tibshirani, R. & Friedman, J. *The Elements of Statistical Learning*. **1**, (Springer-
939 Verlag New York, 2009).
- 940 67. Mandal, S. *et al.* Analysis of composition of microbiomes: a novel method for studying
941 microbial composition. *Microb. Ecol. Health Dis.* **26**, 27663 (2015).
- 942 68. Ritchie, M. E. *et al.* limma powers differential expression analyses for RNA-sequencing and
943 microarray studies. *Nucleic Acids Res.* **43**, e47 (2015).
- 944 69. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment search
945 tool. *J. Mol. Biol.* **215**, 403–410 (1990).
- 946 70. Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068–2069
947 (2014).
- 948 71. Segata, N., Börnigen, D., Morgan, X. C. & Huttenhower, C. PhyloPhlAn is a new method for
949 improved phylogenetic and taxonomic placement of microbes. *Nat. Commun.* **4**, 2304 (2013).
- 950 72. Kaminski, J. *et al.* High-Specificity Targeted Functional Profiling in Microbial Communities
951 with ShortBRED. *PLoS Comput. Biol.* **11**, e1004557 (2015).
- 952 73. Edgar, R. C. Search and clustering orders of magnitude faster than BLAST. *Bioinformatics* **26**,
953 2460–2461 (2010).
- 954 74. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–2079

- 955 (2009).
- 956 75. Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* **27**, 2156–2158 (2011).
- 957 76. Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software version 7:
958 improvements in performance and usability. *Mol. Biol. Evol.* **30**, 772–780 (2013).
- 959 77. Price, M. N., Dehal, P. S. & Arkin, A. P. FastTree 2--approximately maximum-likelihood trees
960 for large alignments. *PLoS One* **5**, e9490 (2010).
- 961 78. Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large
962 phylogenies. *Bioinformatics* **30**, 1312–1313 (2014).
- 963 79. Asnicar, F., Weingart, G., Tickle, T. L., Huttenhower, C. & Segata, N. Compact graphical
964 representation of phylogenetic data and metadata with GraPhlAn. *PeerJ* **3**, e1029 (2015).
- 965 80. Livak, K. J. & Schmittgen, T. D. Analysis of relative gene expression data using real-time
966 quantitative PCR and the 2(-Delta Delta C(T)) Method. *Methods* **25**, 402–408 (2001).

967

968

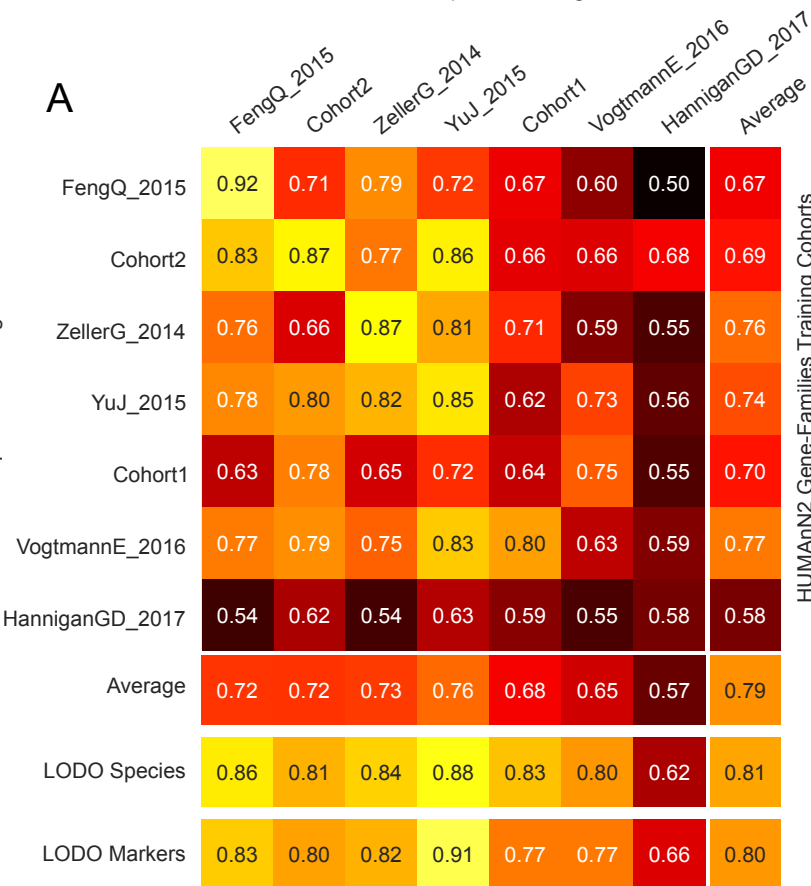


MetaPhlan2 Species Testing Cohorts

HUMANn2 Gene-Families Testing Cohorts

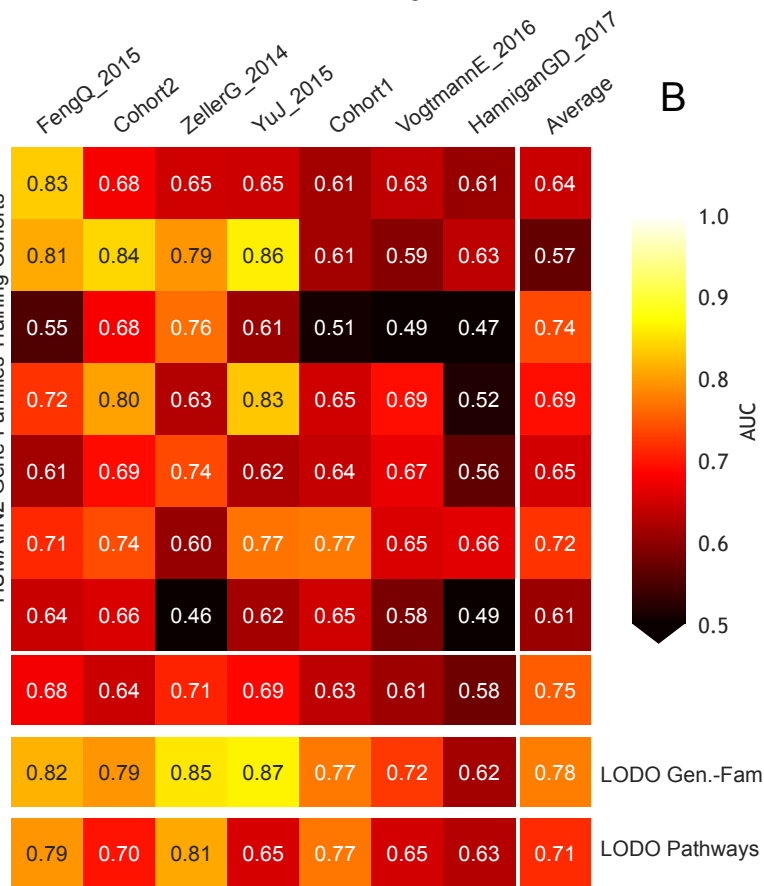
A

MetaPhlan2 Species Training Cohorts

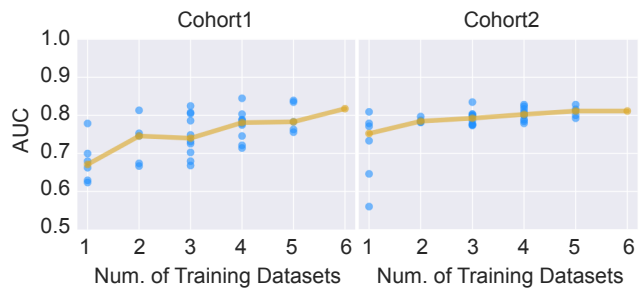


B

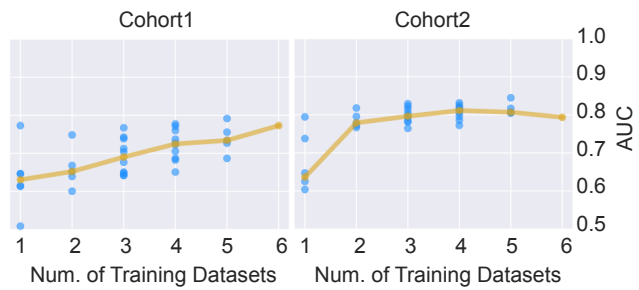
HUMANn2 Gene-Families Training Cohorts



C

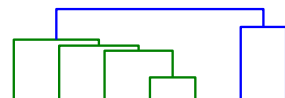


D



A

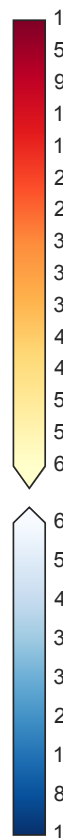
Random Forest Feature Ranking



<i>Peptostreptococcus stomatis</i>	31	18	4	1	1	177	221	1
<i>Fusobacterium nucleatum</i>	1	2	114	4	11	346		2
<i>Parvimonas</i> spp.	26	4	1	3	2	300	207	3
<i>Porphyromonas asaccharolytica</i>	3	5	213	59	70	20	99	4
<i>Gemella morbillorum</i>	121	6	2	27	3	249	162	5
<i>Clostridium symbiosum</i>	6	105	133	22	4	8	205	8
<i>Parvimonas micra</i>	213	13	3	7	7	299		10
<i>Escherichia coli</i>	2	33	20	44	73	2	58	14
<i>Streptococcus parasanguinis</i>	29	21	80	38	100	95	4	23
<i>Clostridium leptum</i>	5	111	40	88	108	67	59	29
<i>Clostridium hathewayi</i>	67	57	122	6	5	24	202	31
<i>Anaerotruncus colihominis</i>	40	148	137	97	90	5	223	32
<i>Prevotella copri</i>	21	1	55	140	52	65	158	50
<i>Lachnospiraceae 3 1 57FAA CT1</i>	4	75	123	110	130	14	69	67
<i>Actinomyces graevenitzii</i>	150	150	5	186	191	185	96	101

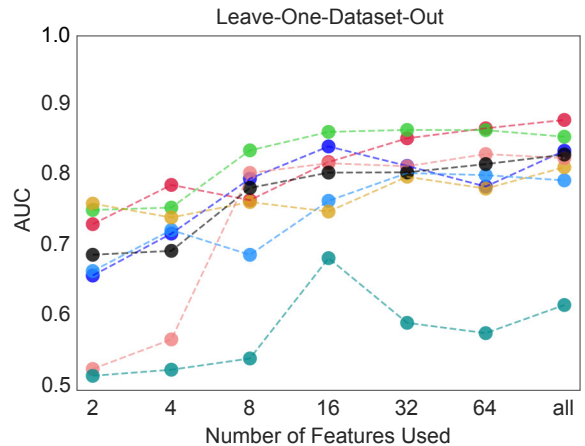
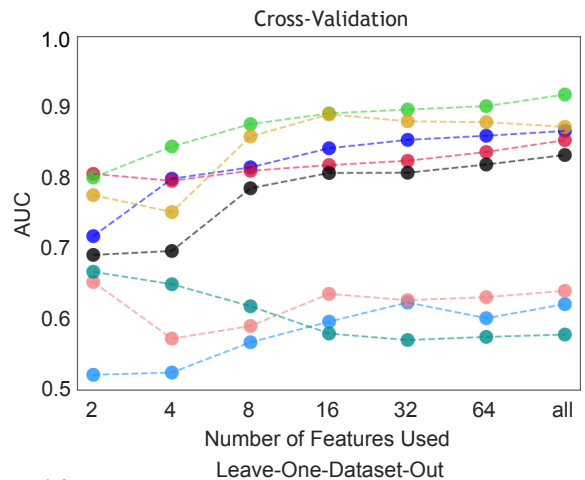
<i>Alistipes</i> spp.	225	207	276	218	238	128	2	166
<i>Lachnospiraceae 8 1 57FAA</i>	193	65	159	195	234	129	3	130
<i>Dialister invisus</i>	146	3	143	149	171	91	46	97
<i>Ruminococcus gnavus</i>	16	71	107	76	84	93	5	64
<i>Bifidobacterium longum</i>	103	96	52	45	35	4	36	42
<i>Subdoligranulum</i> spp.	46	81	36	46	76	3	6	26
<i>Lachnospiraceae 5 1 63FAA</i>	32	109	71	5	16	118	14	15
<i>Eubacterium eligens</i>	63	86	70	11	15	1	70	9
<i>Streptococcus salivarius</i>	25	28	54	2	66	23	1	7

VogtmannE_2016
FengQ_2015
Cohort2
ZellerG_2014
YuJ_2015
Cohort1
HanniganGD_2017
Global Ranking



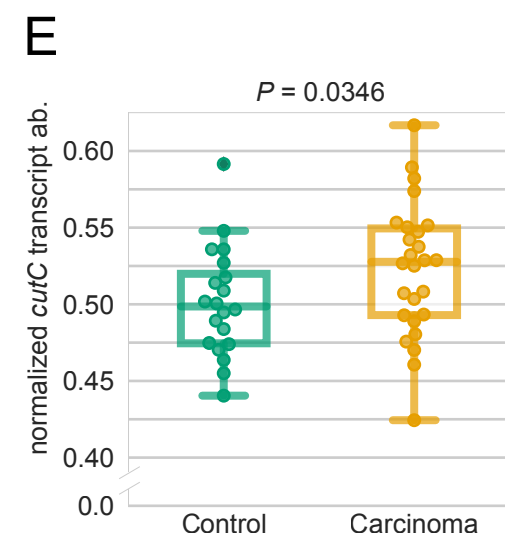
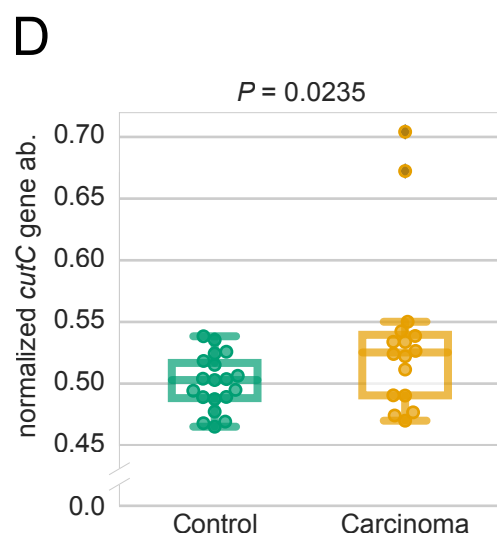
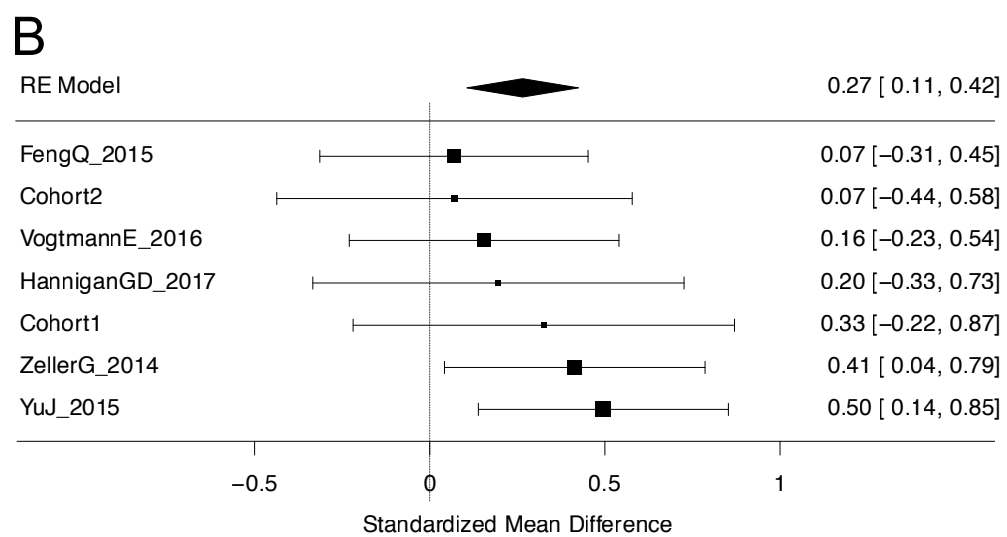
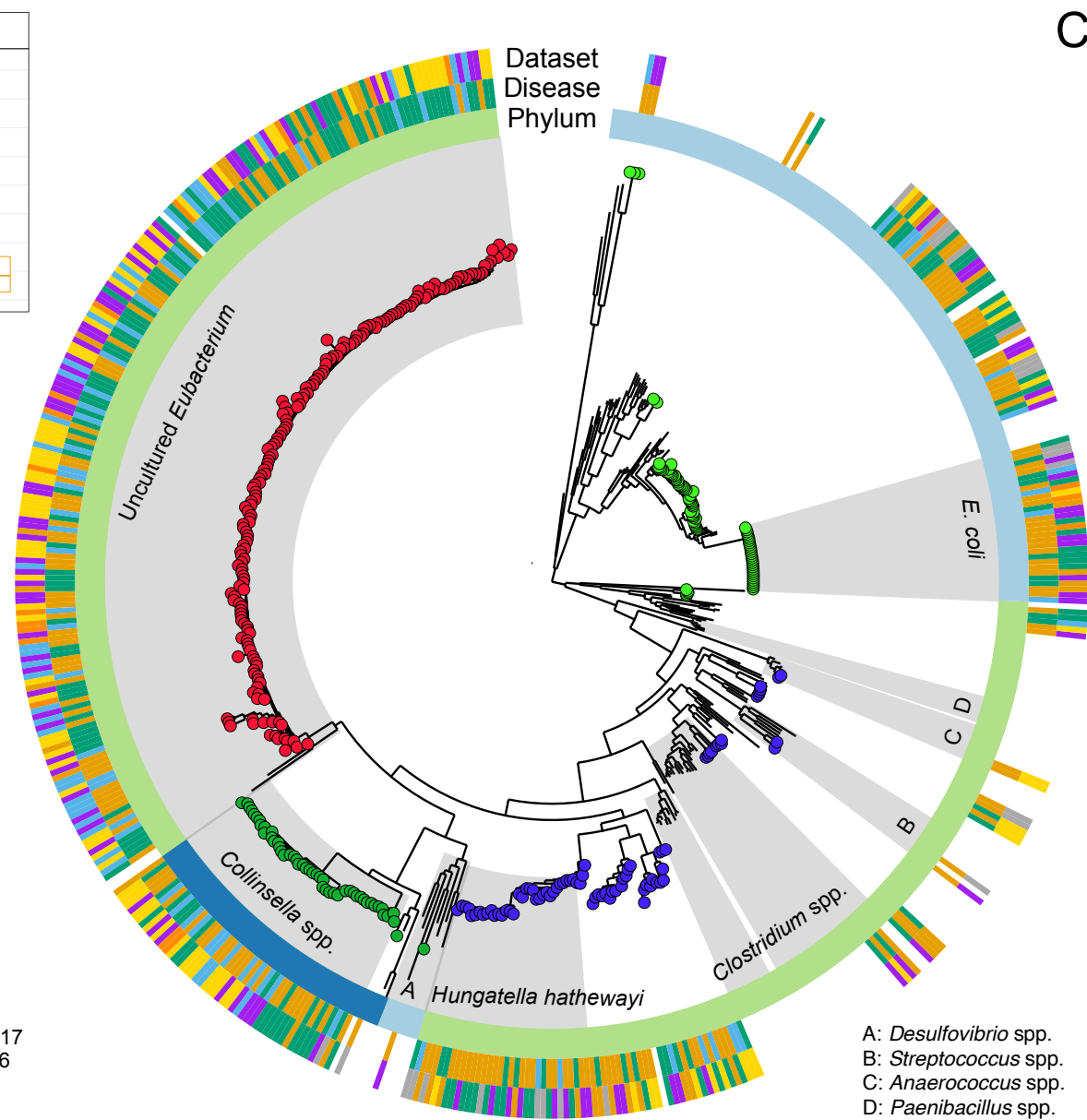
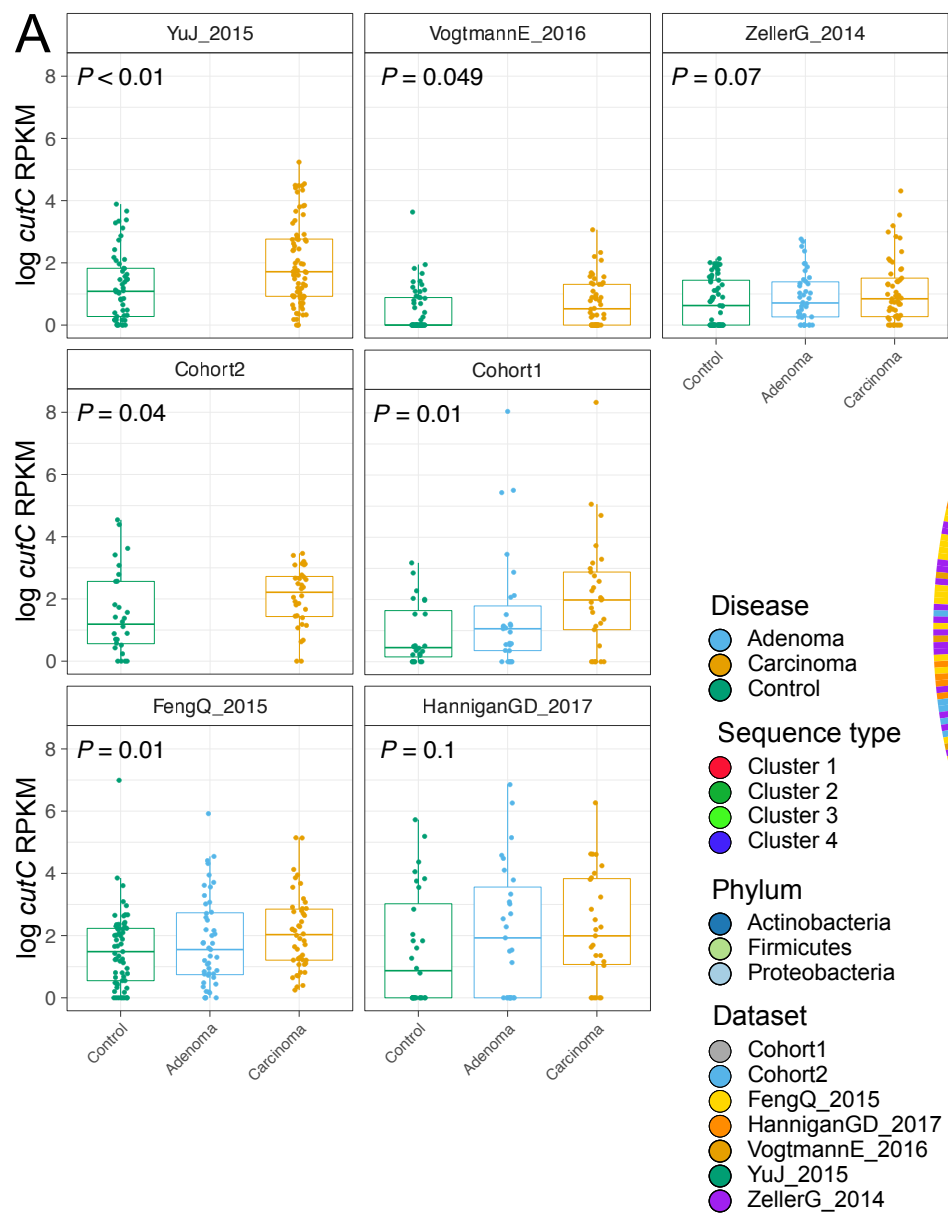
Random Forest Selected Features AUC

FengQ_2015 ZellerG_2014 Cohort2 HanniganGD_2017
YuJ_2015 VogtmannE_2016 Cohort1 Cross-Validation

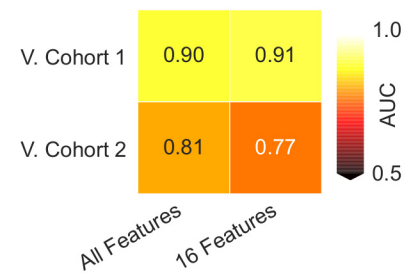


B

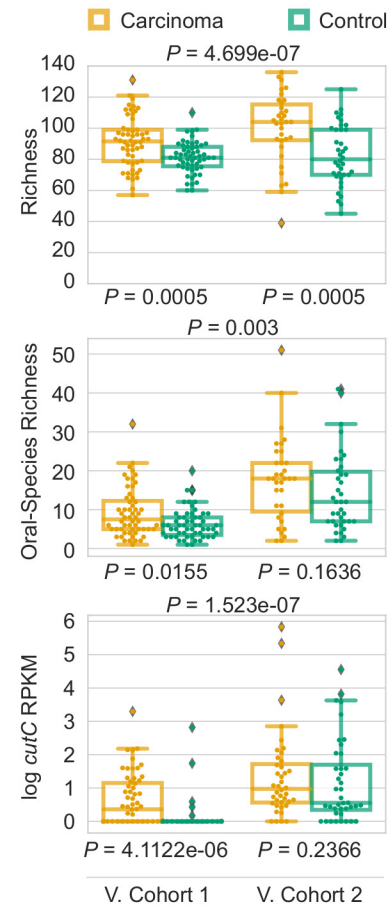
C



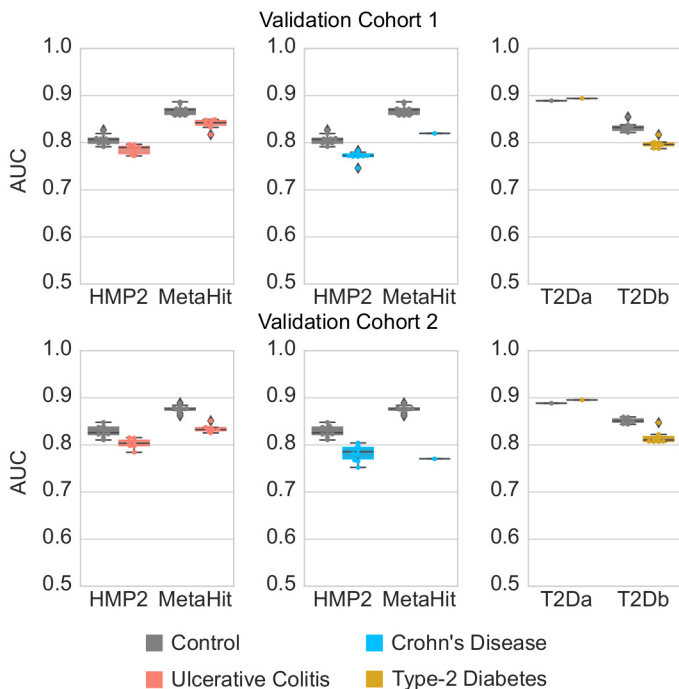
A Random Forest Model



B Validation of predictive biomarkers



CRC-specificity of the predictive models



Relationship to other non-invasive tests

