1  **Meta-analysis of fecal metagenomes reveals global microbial signatures that**
2  **are specific for colorectal cancer**

3

4  **Authors**
5  Jakob Wirbel[1]*, Paul Theodor Pyl[2,3]*, Ece Kartal[1,4], Konrad Zych[1], Alireza Kashani[2], Alessio Milanese[1],
6  Jonas S Fleck[1], Anita Y Voigt[1,5], Albert Palleja[2], Ruby P Ponnudurai[1], Shinichi Sunagawa[1,6], Luis
7  Pedro Coelho[1,‡], Petra Schrotz-King[7], Emily Vogtmann[8], Nina Habermann[9], Emma Niméus[3,10], Andrew
8  M Thomas[11,12], Paolo Manghi[11], Sara Gandini[13], Davide Serrano[13], Sayaka Mizutani[14,15], Hirotsugu
9  Shiroma[14], Satoshi Shiba[16], Tatsuhiro Shibata[16,17], Shinichi Yachida[16,18], Takuji Yamada[14,19], Levi
10  Waldron[20,21], Alessio Naccarati[22,23], Nicola Segata[11], Rashmi Sinha[8], Cornelia M. Ulrich[24], Hermann
11  Brenner[7,25,26], Manimozhiyan Arumugam[2,27]+, Peer Bork[1,4,28,29]+, Georg Zeller[1]+

12

13  **Affiliations**
14      1.  Structural and Computational Biology Unit, European Molecular Biology Laboratory (EMBL),
15          Heidelberg, Germany
16      2.  Novo Nordisk Foundation Center for Basic Metabolic Research, Faculty of Health and
17          Medicine, University of Copenhagen, Copenhagen, Denmark
18      3.  Division of Surgery, Oncology and Pathology, Department of Clinical Sciences Lund, Faculty
19          of Medicine, Lund University, Sweden
20      4.  Molecular Medicine Partnership Unit (MMPU), Heidelberg, Germany
21      5.  The Jackson Laboratory for Genomic Medicine, Farmington, Connecticut, USA
22      6.  Department of Biology, ETH Zürich, Zürich, Switzerland
23      7.  Division of Preventive Oncology, National Center for Tumor Diseases (NCT) and German
24          Cancer Research Center (DKFZ), Heidelberg, Germany
25      8.  Division of Cancer Epidemiology and Genetics, National Cancer Institute, Bethesda,
26          Maryland, USA
27      9.  Genome Biology Unit, European Molecular Biology Laboratory (EMBL), Heidelberg, Germany
28      10. Division of Surgery, Department of Clinical Sciences Lund, Faculty of Medicine, Skane
29          University Hospital, Lund, Sweden
30      11. Department CIBIO, University of Trento, Trento, Italy.
31      12. Biochemistry Department, Chemistry Institute, University of São Paulo, São Paulo, Brazil.
32      13. IEO, European Institute of Oncology IRCCS, Milan, Italy.
33      14. School of Life Science and Technology, Tokyo Institute of Technology, Tokyo, Japan
34      15. Research Fellow of Japan Society for the Promotion of Science
35      16. Division of Cancer Genomics, National Cancer Center Research Institute, Tokyo, Japan
36      17. Laboratory of Molecular Medicine, Human Genome Center, The Institute of Medical Science,
37          The University of Tokyo, Tokyo, Japan
38      18. Department of Cancer Genome Informatics, Graduate School of Medicine/Faculty of
39          Medicine, Osaka University, Osaka, Japan
40      19. PRESTO, Japan Science and Technology Agency, Saitama, Japan

41  20. Graduate School of Public Health and Health Policy, City University of New York, New York,
42      USA.
43  21. Institute for Implementation Science in Population Health, City University of New York, New
44      York, USA.
45  22. Italian Institute for Genomic Medicine (IIGM), Turin, Italy.
46  23. Department of Molecular Biology of Cancer, Institute of Experimental Medicine, Prague,
47      Czech Republic.
48  24. Huntsman Cancer Institute and Department of Population Health Sciences, University of
49      Utah, Salt Lake City, Utah, USA
50  25. Division of Clinical Epidemiology and Aging Research, German Cancer Research Center
51      (DKFZ), Heidelberg, German
52  26. German Cancer Consortium (DKTK), German Cancer Research Center (DKFZ), Heidelberg,
53      German
54  27. Faculty of Healthy Sciences, University of Southern Denmark, Odense, Denmark
55  28. Max Delbrück Centre for Molecular Medicine, Berlin, Germany
56  29. Department of Bioinformatics, Biocenter, University of Würzburg, Würzburg, Germany
57
58  ‡. Present Address: Institute of Science and Technology for Brain-Inspired Intelligence, Fudan
59  University, Shanghai 200433, China
60
61  +  These authors jointly supervised the work. Correspondence should be addressed to
62  zeller@embl.de, bork@embl.de or arumugam@sund.ku.dk
63  *  These authors contributed equally to the work.
64
65
66

67    **ABSTRACT**

68

69    Association studies have linked microbiome alterations with many human diseases, but not always

70    reported consistent results, which necessitates cross-study comparisons. Here, a meta-analysis of

71    eight geographically and technically diverse fecal shotgun metagenomic studies of colorectal cancer

72    (CRC, N = 768), which was controlled for several confounders, identified a core set of 29 species

73    significantly enriched in CRC metagenomes (FDR < 1E-5). CRC signatures derived from single

74    studies maintained accuracy in other studies. By training on multiple studies we improved detection

75    accuracy and disease specificity for CRC. Functional analysis of CRC metagenomes revealed

76    enriched protein and mucin catabolism genes and depleted carbohydrate degradation genes.

77    Moreover we inferred elevated production of secondary bile acids from CRC metagenomes

78    suggesting a metabolic link between cancer-associated gut microbes and a fat- and meat-rich diet.

79    Through extensive validations, this meta-analysis firmly establishes globally generalizable, predictive

80    taxonomic and functional microbiome CRC signatures as a basis for future diagnostics.

81

82

83    **INTRODUCTION**

84

85    Studying microbial communities colonizing the human body in a culture-independent manner has been

86    enabled by metagenomic sequencing technologies [1]. These have yielded glimpses into the complex

87    yet incompletely understood interactions between the gut microbiome – the microbial ecosystem

88    residing primarily in the large intestine – and its host [2]. To explore microbiome-host interactions in a

89    disease context, metagenome-wide association studies (MWAS) have begun to map gut microbiome

90    alterations in diabetes, inflammatory bowel disease, colorectal cancer and many other conditions [3-

91    12]. However, due to the many biological factors possibly influencing gut microbiome composition in

92    addition to the condition studied, a current challenge for MWAS is confounding, which can cause false

93    associations [13, 14]. This issue is further aggravated by a lack of standards in metagenomic data

94    generation and processing, making it difficult to disentangle technical from biological effects [15].

95

96    Robustness of microbiome-disease associations can be assessed through comparisons across

97    multiple metagenomic case-control studies, i.e. meta-analyses. These aim at identifying associations

98    that are consistent across studies and thus less likely attributable to biological or technical

99    confounders. Most informative are meta-analyses of populations from diverse geographic and cultural

100    regions. Previous microbiome meta-analyses based on 16S rRNA gene amplicon data found stark

101    technical differences between studies and the reported taxonomic disease associations were either of

102    low effect size or not well resolved [16-18]. In contrast, shotgun metagenomics enables analyses with

103    higher taxonomic resolution and of gene functions to improve statistical power for fine-mapping

104    disease-associated strains and aid in the interpretation of host-microbial co-metabolism. Thus far

105    however, meta-analyses of shotgun metagenomic data have either reported on features of general

106    dysbiosis in comparisons across multiple diseases [19], or have left it unclear how well microbiome

107 signatures generalize across studies of the same disease when data are rigorously separated to avoid
108 over-optimistic evaluations of their prediction accuracy [20].
109
110 Here, we present a meta-analysis of a total of eight studies of CRC including fecal metagenomic data
111 from 386 cancer cases and 392 tumor-free controls. After consistent data reprocessing, we examined
112 an initial set of five studies for CRC-associated changes in the gut microbiome. Firstly, we investigated
113 potential confounders, followed by identifying (univariate) microbial species associations, and inferring
114 species co-occurrence patterns in CRC. Secondly, we trained multivariable classification models for
115 recognition of CRC status, from both taxonomic and functional microbiome profiles and tested how
116 accurately these models generalized to data from studies not used for training. Moreover, we
117 evaluated performance improvements achieved by pooling data across studies and the disease-
118 specificity of the resulting classification models. Thirdly, targeted investigation of virulence and toxicity
119 genes as candidate functional biomarkers for CRC revealed several of these to be enriched in CRC
120 metagenomes indicative of their prevalence and potential relevance in CRC patients. Three additional,
121 more recent studies were finally used to independently validate these taxonomic and functional CRC
122 signatures.
123
124 **RESULTS**
125
126 **Consistent processing of published and new data for meta-analysis of CRC metagenomes**
127 In this meta-analysis we included four published studies which used fecal shotgun metagenomics to
128 characterize CRC patients compared to healthy controls (referred to by the country codes FR, AT, CN,
129 and US, corresponding to the respective main study population; see **Table 1, Supplementary Table**
130 **S1,** and Methods for inclusion criteria). For an additional fifth study population, we generated new
131 fecal metagenomic data from samples collected in Germany (herein abbreviated as DE); a subset of
132 samples from this patient collective were published previously (**Table 1**, Methods, [8]). These five
133 studies were conducted on three continents and differed in sampling procedures, sample storage, and
134 DNA extraction protocols. Notably, the fecal specimen of the US study were freeze-dried and stored at
135 -80°C for more than 25 years before DNA extraction and sequencing [10]. In all studies, however,
136 samples were collected prior to treatment, thus excluding cancer therapy as a potential confounding
137 effect [14, 21]. Most samples were even taken before bowel preparation for colonoscopy, with some
138 exceptions in the DE, CN and US studies (**Supplementary Table S2**). To ensure consistency in
139 bioinformatic analyses, all raw sequencing data were (re-)processed using mOTUs2 for taxonomic
140 profiling [22] and MOCAT2 for functional profiling [23].
141
142 **Univariate meta-analysis of species associated with CRC**
143 The first aim of the meta-analysis was to determine gut microbial species that are enriched or depleted
144 in CRC metagenomes in a consistent manner across the five study populations. However, as these
145 studies differed from one another in many biological and technical aspects, we first quantified the
146 effect of study-associated heterogeneity on microbiome composition. We contrasted this with other
147 potential confounders ('patient age', 'BMI', 'sex', 'sampling after colonoscopy', and 'library size';

148  additionally, 'smoking status', 'type II diabetes comorbidity', and 'vegetarian diet' where available
149  (**Extended Data 1, Supplementary Table S3**). This analysis revealed the factor 'study' to have a
150  predominant impact on species composition, which is supported by a recent comparison of DNA
151  extraction protocols, as these typically differ between studies [15]. An analysis of microbial alpha and
152  beta diversity showed study heterogeneity to also have a larger effect on overall microbiome
153  composition than CRC in our data (**Extended Data 2**).

155  For the identification of microbial taxa significantly differing in abundance in CRC, parametric effect
156  size measures are not well established, because microbiome data is characterized by non-Gaussian
157  distributions with extreme dispersion; we thus used a generalisation of the fold change (**Extended
158  Data 3)** and non-parametric significance testing. In this permutation test framework [24] (herein
159  referred to as blocked univariate Wilcoxon tests) differential abundance in CRC can be assessed
160  while accounting for 'study' as a nuisance effect that is treated as a blocking factor; additionally,
161  motivated by our confounder analysis, we also blocked for 'colonoscopy' in all analyses (Methods,
162  **Extended Data 1**). To rule out spurious associations due to the compositional nature of microbial
163  relative abundance data, we additionally compared the results of this test with a method [25]
164  employing log-ratio transformation (and found highly correlated results, **Supplementary Fig. 1,
165  Supplementary Table S4**).

167  At a meta-analysis false discovery rate (FDR) of 0.005, we identified 94 microbial species to be
168  differentially abundant in the CRC microbiome, out of 849 species consistently detected across
169  studies (**Supplementary Table S4**, Methods). Among these, we focused on a core set of the 29 most
170  significant markers (FDR < 1E-5, **Fig. 1a**) for further analysis. The latter included members of several
171  genera previously associated with CRC, such as *Fusobacterium*, *Porphyromonas*, *Parvimonas*,
172  *Peptostreptococcus*, *Gemella*, *Prevotella*, and *Solobacterium* (**Fig. 1b**, [8-11]),  and 8 additional
173  species without genomic reference sequences (meta-mOTUs, Methods, [22]) mostly from the
174  *Porphyromonas* and *Dialister* genera and the Clostridiales order (see **Extended Data 4** and
175  **Supplementary Table S4** for genus-level associations). Collectively, these 29 core CRC-associated
176  species show a previously underappreciated diversity of 11 Clostridiales species to be enriched in
177  CRC (**Fig. 1b**). In contrast to the majority of species that are more strongly affected by study
178  heterogeneity than by CRC status, 26 out of the 29 CRC-associated species varied more by disease
179  status (**Fig. 1d**).

181  All of the core CRC-associated species were enriched in patients and were often undetectable in
182  metagenomes from non-neoplastic controls. While previous studies were contradictory in the reported
183  proportion of positive versus negative associations [8, 9, 17, 20], our meta-analysis results are more
184  easily reconciled with a model in which – potentially many – gut microbes contribute to or benefit from
185  tumorigenesis than with the opposing model in which a lack of protective microbes contributes to CRC
186  development (**Fig. 1b**). Although these core taxonomic CRC associations were highly significant and
187  consistent, individual studies showed marked discrepancies in the species identified as significant
188  (**Fig. 1a**). Retrospective examination of the precision and sensitivity with which individual studies

189  detected this core of CRC-associated species showed relatively low sensitivity for the US study

190  (consistent with the original report [10]) and low precision of the AT study due to associations that

191  were not replicated in other studies (**Supplementary Fig. 2**).

192

193  Analyzing patient metagenomes for co-occurrences among the core set of 29 species that are strongly

194  enriched in the CRC microbiome revealed four species clusters with distinct taxonomic composition

195  (**Fig. 2a**, **Extended Data 5**, Methods). Two of them showed strong taxonomic consistency: Cluster 1

196  exclusively comprised *Porphyromonas* spp., and cluster 4 only contained members of the Clostridiales

197  order. In contrast, the other two clusters were taxonomically more heterogeneous with cluster 3

198  grouping together the species with highest prevalence in CRC cases (all among the ten most highly

199  significant markers), consistent with a co-occurrence analysis of one of the data sets included here

200  [11]. Cluster 2 contained species with intermediate prevalence.

201

202  Investigating whether these four clusters were associated with different tumor characteristics, we

203  found the *Porphyromonas* cluster 1 to be significantly enriched in rectal tumors (**Fig. 2b**), consistent

204  with the presence of superoxide dismutase genes in *Porphyromonas* genomes possibly conferring

205  tolerance to a more aerobic milieu in the rectum (**Extended Data 5**). The Clostridiales cluster 4 was

206  significantly more prevalent in female CRC patients. All species clusters showed a slight tendency

207  towards late-stage CRC (i.e. AJCC stages III and IV), but this was only significant for cluster 3.

208  Associations with patient age and BMI were weaker and not significant (**Extended Data 5**). To rule out

209  secondary effects due to differences in patient composition among studies, all of these tests were

210  corrected for study effects (by blocking for 'study' and 'colonoscopy', see Methods). At the level of

211  individual species, significant stage-specific enrichments could not be detected suggesting CRC-

212  associated microbiome changes to be less dynamic during cancer progression than previously

213  postulated [26], although fecal material may be less suitable to address this question than tissue

214  samples.

215

216  **Metagenomic CRC classification models**

217  To establish metagenomic signatures for CRC detection across studies in face of geographic and

218  technical heterogeneity, we developed multivariable statistical modeling workflows with rigorous

219  external validation to avoid prevailing issues of overfitting and over-optimistic reports of model

220  accuracy [19]. As a precaution against over-optimistic evaluation, these workflows are independent of

221  the above-described differential abundance analysis. Instead, LASSO (Least Absolute Shrinkage and

222  Selection Operator) logistic regression classifiers were employed to select predictive microbial

223  features and eliminated uninformative ones (Methods).

224

225  In a first step, we used abundance profiles from five studies including the 849 most abundant microbial

226  species and assessed how well classifiers trained in cross validation (CV) on one study generalize in

227  evaluations on the other four studies (study-to-study transfer of classifiers) (**Fig. 3a**). Within-study

228  cross-validation performance, as quantified by the Area Under the Receiver Operating Characteristics

229  (AUROC) curve, ranged between 0.69 and 0.92 and was generally maintained in study-to-study

230  transfer (AUROC dropping by 0.07±0.12 on average) with two notable exceptions. First, in line with
231  the univariate analysis of species associations, CRC detection accuracy on the US study was lower
232  than for the other studies, both in cross-validation and in study-to-study transfer. This could potentially
233  be explained by the US fecal specimen, unlike in the other studies, being freeze-archived for >25
234  years before metagenomic sequencing [10]. Second, classifiers trained on the AT study did not
235  generalize as well to the other studies, consistent with low study precision seen in univariate meta-
236  analysis (**Supplementary Fig. 2**). Given the microbial co-occurrence clusters described above, we
237  wondered whether species-species interactions would provide additional information relevant for CRC
238  recognition that is not contained in species abundance profiles. However, nonlinear classifiers able to
239  exploit such interactions did not yield significantly better accuracies (**Supplementary Fig. 3**, see also
240  [27]), suggesting that the linear model based on few biomarkers (on average 17 species account for
241  more than 80% of the classifier weight, **Extended Data 6**) is near optimal for CRC prediction.
242
243  We further assessed if including data from all but one study in model training improves prediction on
244  the remaining held-out study (leave-one-study-out validation, LOSO). LOSO performance of species-
245  level models ranged between 0.71 and 0.91, and when disregarding the US study as an outlier was
246  ≥0.83 (**Fig. 3b**). This corresponds to a LOSO accuracy increase of 0.076±0.03 compared to study-to-
247  study transfer. These results suggest that one can expect a CRC detection accuracy ≥0.8 (AUROC)
248  for any new CRC study using similarly generated metagenomic data. We moreover verified that
249  metagenomic CRC classification models trained on species composition were not biased for clinical
250  subgroups. With the exception of slightly more sensitive detection of late stage CRC (P = 0.03, mostly
251  originating from the US study, **Extended Data 7**), we did not observe any classification bias by patient
252  age, sex, BMI, or localization. Together this suggests that these metagenomic classifiers are unlikely
253  to be strongly confounded by the clinical parameters recorded.
254
255  Several previous studies comparing microbiome changes across multiple diseases reported primarily
256  general dysbiotic alterations and highlighted the need to examine the disease specificity of
257  microbiome signatures [17, 19]. Therefore, we assessed false positive (FP) predictions of our
258  metagenomic CRC classifiers on fecal metagenomes of type 2 diabetes [4, 5], Parkinson's disease
259  [12], ulcerative colitis and Crohn's disease [6, 7] patients, reasoning that classifiers relying on
260  biomarkers for general dysbiosis would yield an excess of FPs on these cohorts. However, our LOSO
261  classification models calibrated to have a false-positive rate (FPR) of 0.1 on CRC datasets in fact
262  maintained similarly low FPRs on other disease datasets ranging from 0.09 to 0.13 (**Fig. 3c**).
263  Interestingly, disease specificity of LOSO models was significantly improved over that observed for
264  classifiers trained on a single study, indicating that inclusion of multiple studies in the training set of a
265  classifier can substantially improve its specificity for a given disease.
266
267
268  **Functional metagenomic signatures for CRC**
269  As shotgun metagenomics data, in contrast to 16S rRNA gene amplicon data, allow for a direct
270  analysis of the functional potential of the gut microbiome, we examined how predictive metabolic

271  pathways and orthologous gene families differing in abundance between CRC patients and controls
272  would be of CRC status. When applying the same classification workflow as above to eggNOG
273  orthologous gene family abundances [28], CRC detection accuracy was very similar to that observed
274  for taxonomic models (**Fig. 3de**). AUROC values ranged from 0.70 to 0.81 for study-to-study transfer
275  (per-study averages, **Fig. 3e**) and from 0.78 to 0.89 in LOSO validation with a pattern of generalization
276  across studies resembling that for taxonomic classifiers. The accuracy of functional signatures did not
277  strongly depend on eggNOG as an annotation source, but was similar when based on other
278  comprehensive functional databases, such as KEGG [29] (**Extended Data 8**). When using individual
279  gene abundances from metagenomic gene catalogues as a classifier input [30], we observed higher
280  within-study cross-validation AUROC values of ≥0.96 in all studies, but lower generalization to other
281  studies (AUROC between 0.60 and 0.79) (**Extended Data 8**).

282

283  To explore changes in metabolic capacity of gut microbiomes from CRC patients more broadly, we
284  quantified gut metabolic modules (defined in [31]) and subjected these to the same differential
285  abundance analysis developed for microbial species. Gut metabolic modules with significantly higher
286  abundance (FDR < 0.01, Wilcoxon test blocked for study and colonoscopy) in CRC metagenomes
287  predominantly belonged to pathways for the degradation of amino acids, mucins (glycoproteins) and
288  organic acids. This clear trend was accompanied by a depletion of genes from carbohydrate
289  degradation modules (**Fig. 4ab**). Differences in all four high-level categories were highly significant (P
290  < 1E-6 in all cases, blocked Wilcoxon tests) and consistent across studies (**Fig. 4b**). Overall these
291  results establish a clear shift from dietary carbohydrate utilization in a healthy gut microbiome to amino
292  acid degradation in CRC consistent with an earlier report based on a subset of the data [8].
293  Correlation analysis suggests that increased capacity for amino acid degradation is mostly contributed
294  by CRC-associated Clostridiales (cf. cluster 4 in Fig. 2, **Supplementary Fig. 4**). About one half of
295  these metagenomic pathway enrichments are also in agreement with independent metabolomics data
296  suggesting increased availability of amino acids in epithelial cells or feces of CRC patients
297  (**Supplementary Table S5,** [32-36]). While the observed pathway enrichments could potentially result
298  from many factors, including unmeasured ones [13], they are consistent with established dietary risk
299  factors for CRC, which include red and processed meat consumption [37] and low fiber intake [38].

300

301  The large metagenomic data set analyzed here allowed us to quantify the prevalence of gut microbial
302  virulence and toxicity mechanisms thought to play a role in colorectal carcinogenesis. Prominent
303  examples include the *Fusobacterium nucleatum* adhesion protein A (encoded by the *fadA* gene), the
304  *Bacteroides fragilis* enterotoxin (*bft* gene) and colibactin produced by some *Escherichia coli* strains
305  (*pks* genomic island) [39, 40] . Moreover, intestinal *Clostridium* spp. are known to contribute to the
306  conversion of primary to secondary bile acids using several metabolic pathways including 7α-
307  dehydroxylation, encoded in the *bai* operon [41]. The products of this 7α-dehydroxylation pathway,
308  deoxycholate and lithocholate, are known hepatotoxins associated with liver cancer [42] and
309  hypothesized to also promote CRC [43]. Although intensely studied at a mechanistic level, these
310  factors are not (well) represented in general databases that can be used for metagenome annotation
311  (**Supplementary Fig. 5)**. Thus, we built a targeted metagenome annotation workflow based on

312   Hidden Markov Models to identify and quantify virulence factors and toxicity pathways of interest in
313   CRC. Additionally, we used co-abundance clustering to infer operon completeness for factors encoded
314   by multiple genes (Methods, **Extended Data 9, Supplementary Fig. 5**). While *fadA*, *bft*, the *pks* island
315   and the *bai* operon were clearly detectable in deeply sequenced fecal metagenomes, they varied
316   broadly with respect to abundance, significance and cross-study consistency of enrichment (**Fig. 4c**):
317   *fadA* and *pks* were significantly enriched in CRC metagenomes (P = 5.3E-10 and 4.1E-4 respectively),
318   whereas no significant abundance difference could be detected for *bft* in fecal metagenomes, despite
319   reports on its enrichment in the mucosa of CRC patients [44], its carcinogenic effect in mouse models
320   [45], and synergistic action with *pks* [46]. Our quantification of the *bai* operon showed a highly
321   significant enrichment in CRC metagenomes (P = 1.6E-9) observed across all five studies (**Fig. 4d**) at
322   an average abundance that exceeded *fadA* and *pks* copy numbers (**Fig. 4c**). Metagenome analysis
323   indicated that at least four Clostridiales species (including the well characterized *C. scindens* and *C.*
324   *hylemonae* [47, 48]) have a (near) complete 7α-dehydroxylation pathway contributing to the observed
325   enrichment of *bai* operon copies (**Extended Data 9**). To validate this finding and further explore its
326   value towards diagnostic application, we developed a targeted quantification assay for the *baiF* gene
327   based on quantitative PCR (qPCR, see Methods). Quantification of *baiF* by qPCR using genomic DNA
328   from 47 fecal samples of the DE study population was found to be similar to, yet more sensitive than
329   by metagenomics (**Fig. 4e**). Gut microbial *baiF* copy numbers clearly distinguished CRC patients from
330   controls (P = 0.001) at an AUROC of 0.77, which in this subset of samples is surpassed by only a
331   single species marker for CRC (**Extended Data 9**). Although consistent with increased deoxycholate
332   metabolite levels reported for serum and stool samples of CRC patients [49], this finding does not
333   imply 7α-dehydroxylation pathway activity. We therefore quantified *baiF* expression using RNA
334   extracts from the same set of fecal samples, and found also transcript levels to be elevated in CRC
335   patients (**Fig. 4f**). The observed weak correlation of *baiF* expression with genomic abundance (**Fig. 4f**)
336   might be explained by dynamic transcriptional regulation [47] and *bai* expression in feces might not
337   accurately reflect the tumor microenvironment. Taken together, these data suggest gut microbial
338   metabolic markers to be meaningful and highly predictive of CRC status.
339
340   **Validation of CRC signatures in independent study populations**
341   Even though CRC classification accuracy for both species and functions were evaluated on
342   independent data, we nonetheless sought to confirm it using two additional study populations from
343   Italy (IT1 and IT2, combined N = 61 CRC, N = 62 CTR, [27], see Methods, Table 1) and one from
344   Japan (JP, N = 40 CRC, N = 40 CTR, see Methods, Table 1). The overlap of single species
345   associations detected in the IT2 study and those from the meta-analysis was found to vary within the
346   range seen for the other studies, whereas for IT1 and JP the overlap was slightly lower (cf. study
347   precision in **Supplementary Fig. 2, Extended Data 10**). Nonetheless, the AUROC of LOSO
348   classification models based on species ranged between 0.79 and 0.81 and that for the classifiers
349   based on eggNOG from 0.71 to 0.92 (**Fig. 5ab**). We also validated CRC enrichment of *fadA*, *pks* and
350   *bai* genes in these three study populations (**Fig. 5c**). Altogether these results highlight consistent
351   alterations in the gut microbiome of CRC patients across eight study populations from seven countries
352   in three continents.

9

**353**

**354** **DISCUSSION**

**355**

**356** Through extensive and statistically rigorous validation, in which data from studies used for training is

**357** strictly separated from that for testing, our meta-analysis firmly establishes that gut microbial

**358** signatures are highly predictive of CRC (see also [27]). In particular metagenomic classifiers trained

**359** on species profiles from multiple studies maintained an AUROC of at least 0.8 in seven out of eight

**360** data sets and achieved an accuracy similar to the fecal occult blood test, a standard non-invasive

**361** clinical test for CRC (**Supplementary Fig. 6**, cf. [8]). These results thus suggest that polymicrobial

**362** CRC classifiers are globally applicable and can overcome technical and geographical study

**363** differences, which we found to generally impact observed microbiome composition more than the

**364** disease itself (**Fig. 1c**, **Extended Data 1, 2**). The generalization accuracy of classifiers across studies

**365** seen here is higher than that reported in 16S rRNA gene amplicon sequencing studies, which are

**366** characterized by even larger heterogeneity across studies [16, 18] (**Supplementary Fig. 7**).

**367**

**368** Previous microbiome meta-analyses suggested that the majority of gut microbial taxa differing in any

**369** given case-control study reflect general dysbiosis rather than disease-specific alterations illustrating

**370** the difficulty of establishing disease-specific microbiome signatures [17, 19]. Here, by combining data

**371** across studies for training (LOSO), we were able to develop disease-specific signatures that

**372** maintained false positive control on diabetes and IBD metagenomes at a very similar level as for CRC

**373** (**Fig. 3c**) despite these diseases having shared effects on the gut microbiome [17, 50] and an

**374** increased comorbidity risk [51].

**375**

**376** Although for diagnostic purposes, unresolved causality between microbial and host processes during

**377** CRC development are not a central concern, elucidating the underlying mechanisms would greatly

**378** enhance our understanding of colorectal tumorigenesis. Towards this goal, we developed both broad

**379** and targeted annotation workflows for functional metagenome analysis. First, we found functional

**380** signatures based on the abundances of orthologous groups of microbial genes to yield accuracies as

**381** high as taxonomic signatures (**Fig. 3**), which raises the hope for future improvements in metagenome

**382** annotation to translate into microbiome signature refinements. Second, by investigating potentially

**383** carcinogenic bacterial virulence and toxicity mechanisms taking a targeted metagenome annotation

**384** approach, we confirmed highly significant enrichments of the colibactin-producing *pks* gene cluster

**385** and the *Fusobacterium nucleatum* adhesin *FadA* in CRC metagenomes (**Fig. 4c**). Our results support

**386** the clinical relevance of these factors adding to the experimental evidence for their carcinogenic

**387** potential [46, 52-54]. We further examined the *bai* operon, encoding enzymes that produce secondary

**388** bile acids via 7α-dehydroxylation, as an example of toxic host-microbial co-metabolism (see [27] for

**389** another intriguing example). While α-dehydroxylated bile acids are established liver carcinogens [42],

**390** their contribution to CRC is less clear [43]. Here, we have, for the first time, shown *bai* to be highly

**391** enriched in stool from CRC patients (**Fig. 4cd**) and confirmed this finding at both the genomic and the

**392** transcriptomic level using qPCR (**Fig. 4ef**). As *bai* enrichment (and expression) is likely a

**393** consequence of a diet rich in fat and meat [55], it is intriguing to explore whether *bai* could be used as

394 a surrogate microbiome marker for such difficult-to-measure dietary CRC risk factors. To further
395 unravel the molecular underpinning of these dietary CRC risk factors, molecular pathological
396 epidemiology studies that investigate the mucosal microbiome as part of the tumor microenvironment,
397 hold great potential [56, 57]. However, they will require more comprehensive diet questionnaires,
398 medical records, and molecular tumor characterizations than are available for the study populations
399 analyzed here. In this context, carcinogens possibly contained in the virome also warrant further
400 investigation [58, 59], but for this goal, metagenomic data needs to be generated with protocols
401 optimized for virus enrichment [60].
402
403 Taken together, our results and those by Thomas, Manghi et al. [27], strongly support the promise of
404 microbiome-based CRC diagnostics. Both taxonomic and metabolic gut microbial marker genes
405 established in these meta-analyses could form the basis of future diagnostic assays that are
406 sufficiently robust, sensitive, and cost-effective for clinical application. The targeted qPCR-based
407 quantification of the *baiF* gene is a first step in this direction. Our metagenomic analysis of this and
408 other virulence and toxicity markers bridge to existing mechanistic work in preclinical models and
409 could enable future work aiming to precisely determine the contribution of gut microbiota to CRC
410 development.
411
412

413 **Data and Code Availability**
414 The raw sequencing data for the samples in the DE study that had not been published before (see
415 Methods), are made available in the European Nucleotide Archive (ENA) under the study identifier
416 PRJEB27928. Metadata for these samples are available as **Supplementary Table S6**.
417 For the other studies included here, the raw sequencing data can be found under the following ENA
418 identifiers: PRJEB10878 for [11], PRJEB12449 for [10], ERP008729 for [9], and ERP005534 for [8].
419 The independent validation cohorts can be found in SRA under the identifier SRP136711 for [27] and
420 in the DDBJ database under the ID DRA006684.
421 Filtered taxonomic and functional profiles used as input for the statistical modeling pipeline are
422 available in **Supplementary Data 1**.
423 The code and all analysis results can be found under https://github.com/zellerlab/crc_meta.
424
425

451 **Competing Interest**

454

455 **Author Contributions**

466
467

468 **Figure Captions**

469

470 **Figure 1. Despite study differences, meta-analysis identifies a core set of gut microbes**
471 **strongly associated with CRC.**
472 **(a)** Meta-analysis significance of gut microbial species derived from blocked Wilcoxon tests (n=574
473 independent observations) is given by bar height (false discovery rate, FDR, of 0.05). **(b)** Underneath,
474 species-level significance as computed by two-sided Wilcoxon test (FDR-corrected P-value) and
475 generalized fold change (Methods) within individual studies are displayed as heatmaps in gray and

476 color, respectively (see color bars and Table 1 for details on studies included). Species are ordered by
477 meta-analysis significance and direction of change. **(c)** For a core of highly significant species (meta-
478 analysis FDR 1E-5), association strength is quantified by the area under the Receiver Operating
479 Characteristics curve (AUROC) across individual studies (color coded diamonds) and 95% confidence
480 intervals are indicated by gray lines. Family-level taxonomic information is color-coded above species
481 names (numbers in brackets are mOTU species identifiers, see Methods). **(d)** Variance explained by
482 disease status (CRC vs controls) is plotted against variance explained by study effects for individual
483 microbial species with dot size proportional to abundance (Methods); core microbial markers are
484 highlighted in red. *F. nucleatum – Fusobacterium nucleatum*.

485

486 **Figure 2. Co-occurrence analysis of CRC-associated gut microbial species reveals four**
487 **clusters preferentially linked to specific patient subgroups.**
488 **(a)** The heatmap shows for all CRC patients (n=285 independent samples) if the respective sample is
489 positive for each of the core set of microbial marker species (see Methods for adjustment of positivity
490 threshold). Samples are ordered according to the sum of positive markers and marker species are
491 clustered based on Jaccard similarity of positive samples, resulting in four clusters (Methods). Barplots
492 in **(b)**, **(c)**, and **(d)** show the fraction of CRC samples that are positive for marker species clusters
493 (defined as the union of positive marker species) broken down by patient subgroups based on
494 differences in tumor location, sex, or CRC stage, respectively. Statistically significant associations
495 between CRC subgroups and marker species clusters were identified using the Cochran–Mantel–
496 Haenszel test blocked for study effects and are indicated above bars (P < 0.1).

497

498 **Figure 3. Both taxonomic and functional metagenomic classification models generalize across**
499 **studies in particular when trained on data from multiple studies.**
500 CRC classification accuracy resulting from cross validation within each study (gray boxes along
501 diagonal) and study-to-study model transfer (external validations off diagonal) as measured by
502 AUROC for classifiers trained on **(a)** species and **(d)** eggNOG gene family abundance profiles. The
503 last column depicts the average AUROC across external validations. Classification accuracy, as
504 evaluated by AUROC on a held-out study, improves if taxonomic **(b)** or functional **(e)** data from all
505 other studies are combined for training (leave-one-study-out, LOSO validation) relative to models
506 trained on data from a single study (study-to-study transfer, average and standard deviation shown).
507 Bar height for study-to-study transfer corresponds to the average of four classifiers (error bars indicate
508 standard deviation, n=4). **(c)** Combining training data across studies substantially improves CRC
509 specificity of the (LOSO) classification models relative to models trained on data from a single study
510 (depicted by bar color, as in (c) and (d)) as assessed by the false positive rate (FPR) on fecal samples
511 from patients with other conditions (see legend). Bar height for study-to-study transfer corresponds to
512 the average FPR across classifiers (n=5) with error bars indicating the standard deviation of FPR
513 values observed.

514

515 **Figure 4. Meta-analysis identifies consistent functional changes in CRC metagenomes.**

**(a)** Meta-analysis significance of gut metabolic modules derived from blocked Wilcoxon tests (n=574 independent samples) is indicated by bar height (top panel, FDR of 0.01). Underneath, the generalized fold change (Methods) for gut metabolic modules [31] within individual studies is displayed as heatmap (see color key below (b)). Metabolic modules are ordered by significance and direction of change. A higher-level classification of the modules is color-coded below the heatmap for the four most common categories (colors as in (b), white indicating other classes). **(b)** Normalized log abundances for these selected functional categories is compared between controls (CTR) and colorectal cancer cases (CRC). Abundances are summarized as geometric mean of all modules in the respective category and statistical significance determined using blocked Wilcoxon tests (n=574 independent samples, see Methods). **(c)** Normalized log abundances for virulence factors and toxins compared between metagenomes of controls (CTR) and colorectal cancer cases (CRC) (significant differences P < 0.05 were determined by blocked Wilcoxon test, n=574 independent samples, see Methods for gene identification and quantification in metagenomes; *fadA: gene* encoding *Fusobacterium nucleatum* adhesion protein A, *bft:* gene encoding *Bacteroides fragilis* enterotoxin, *pks*: genomic island in *Escherichia coli* encoding enzymes for the production of genotoxic colibactin, and *bai*: bile acid inducible operon present in some Clostridiales species encoding bile acid converting enzymes). **(d)** Meta-analysis significance (uncorrected P-value) as determined by blocked Wilcoxon tests (n=574 independent samples) and generalized fold change within individual studies are displayed as bars and heatmap, respectively, for the genes contained in the *bai* operon. Due to high sequence similarity to *baiF*, *baiK* was not independently detectable with our approach. **(e)** Metagenomic quantification of *baiF* (metag. ab. – normalized relative abundance) is plotted against qPCR quantification in genomic DNA (gDNA) extracted from a subset of DE samples (n=47), with Pearson correlation (r) indicated (see Methods). **(f)** Expression of *baiF* determined via qPCR on reverse-transcribed RNA from the same samples in contrast to genomic DNA (as in e). The boxplots on the side of (e), (f) show the difference between cancer (CRC) and control (CTR) samples in the respective qPCR quantification (P-values on top were computed using a one-sided Wilcoxon test). All boxplots show interquartile ranges (IQR) as boxes with the median as a black horizontal line and whiskers extending up to the most extreme points within 1.5-fold IQR.

**Figure 5. Meta-analysis results are validated in three independent study populations**

CRC classification accuracy for independent datasets, two from Italy and one from Japan (see **Supplementary Table S2**), is indicated by bar height for single study (white) and leave-one-study-out (grey) models using either **(a)** species or **(b)** eggNOG gene family abundance profiles (cf. Fig. 3). Bar height for single study models corresponds to the average of five classifiers (error bars indicate standard deviation, n=5). **(c)** Normalized log abundances for virulence factors and toxins (cf. Figure 4c) compared between controls (CTR) and colorectal cancer cases (CRC). P-values were determined by blocked, one-sided Wilcoxon tests (n=193 independent samples). Boxes represent interquartile ranges (IQR) with the median as a black horizontal line and whiskers extending up to the most extreme points within 1.5-fold IQR.

557 **Table 1. Fecal metagenomic studies of colorectal cancer included in this meta-analysis.**
558 See Methods for inclusion criteria and **Supplementary Table S2** for extended meta-data. For a
559 detailed description of patient recruitment and data generation for the DE study, see Methods. The
560 data for 38 samples from the DE study had been published previously as part of an independent
561 validation cohort in [8].

| Country Code | Reference | No. of cases | No. of controls |
|---|---|---|---|
| FR | Zeller et al., 2014 [8] | 53 | 61 |
| AT | Feng et al., 2015 [9] | 46 | 63 |
| CN | Yu et al., 2017 [11] | 74 | 54 |
| US | Vogtmann et al., 2016 [10] | 52 | 52 |
| DE | this study | 60 | 60 |
| **External validation cohorts** | | | |
| IT1 | [27] | 29 | 24 |
| IT2 | [27] | 32 | 28 |
| JP | Courtesy of T. Yamada et al. | 40 | 40 |

562

563

564

565 **References**

566 1.    Tringe, S.G. and E.M. Rubin, *Metagenomics: DNA sequencing of environmental samples.*
567       Nat. Rev. Genet., 2005. **6**(11): p. 805-814.
568 2.    Tremaroli, V. and F. Bäckhed, *Functional interactions between the gut microbiota and host*
569       *metabolism.* Nature, 2012. **489**(7415): p. 242-249.
570 3.    Lynch, S.V. and O. Pedersen, *The Human Intestinal Microbiome in Health and Disease.* N.
571       Engl. J. Med., 2016. **375**(24): p. 2369-2379.
572 4.    Qin, J., et al., *A metagenome-wide association study of gut microbiota in type 2 diabetes.*
573       Nature, 2012. **490**(7418): p. 55-60.
574 5.    Karlsson, F.H., et al., *Gut metagenome in European women with normal, impaired and*
575       *diabetic glucose control.* Nature, 2013. **498**(7452): p. 99-103.
576 6.    Qin, J., et al., *A human gut microbial gene catalogue established by metagenomic*
577       *sequencing.* Nature, 2010. **464**(7285): p. 59-65.
578 7.    Schirmer, M., et al., *Dynamics of metatranscription in the inflammatory bowel disease gut*
579       *microbiome.* Nat Microbiol, 2018. **3**(3): p. 337-346.
580 8.    Zeller, G., et al., *Potential of fecal microbiota for early-stage detection of colorectal cancer.*
581       Mol. Syst. Biol., 2014. **10**(11): p. 766.
582 9.    Feng, Q., et al., *Gut microbiome development along the colorectal adenoma-carcinoma*
583       *sequence.* Nat. Commun., 2015. **6**: p. 6528.
584 10.   Vogtmann, E., et al., *Colorectal Cancer and the Human Gut Microbiome: Reproducibility with*
585       *Whole-Genome Shotgun Sequencing.* PLoS One, 2016. **11**(5): p. e0155362.

15

586 11. Yu, J., et al., *Metagenomic analysis of faecal microbiome as a tool towards targeted non-*
587 *invasive biomarkers for colorectal cancer.* Gut, 2017. **66**(1): p. 70-78.
588 12. Bedarf, J.R., et al., *Functional implications of microbial and viral gut metagenome changes in*
589 *early stage L-DOPA-naïve Parkinson's disease patients.* Genome Med., 2017. **9**(1): p. 39.
590 13. Schmidt, T.S.B., J. Raes, and P. Bork, *The Human Gut Microbiome: From Association to*
591 *Modulation.* Cell, 2018. **172**(6): p. 1198-1215.
592 14. Forslund, K., et al., *Disentangling type 2 diabetes and metformin treatment signatures in the*
593 *human gut microbiota.* Nature, 2015. **528**(7581): p. 262-266.
594 15. Costea, P.I., et al., *Towards standards for human fecal sample processing in metagenomic*
595 *studies.* Nat. Biotechnol., 2017. **35**(11): p. 1069-1076.
596 16. Lozupone, C.A., et al., *Meta-analyses of studies of the human microbiota.* Genome Res.,
597 2013. **23**(10): p. 1704-1714.
598 17. Duvallet, C., et al., *Meta Analysis Of Microbiome Studies Identifies Shared And Disease-*
599 *Specific Patterns.* 2017.
600 18. Shah, M.S., et al., *Leveraging sequence-based faecal microbial community survey data to*
601 *identify a composite biomarker for colorectal cancer.* Gut, 2018. **67**(5): p. 882-891.
602 19. Pasolli, E., et al., *Machine Learning Meta-analysis of Large Metagenomic Datasets: Tools and*
603 *Biological Insights.* PLoS Comput. Biol., 2016. **12**(7): p. e1004977.
604 20. Dai, Z., et al., *Multi-cohort analysis of colorectal cancer metagenome identified altered*
605 *bacteria across populations and universal bacterial markers.* Microbiome, 2018. **6**(1): p. 70.
606 21. Maier, L., et al., *Extensive impact of non-antibiotic drugs on human gut bacteria.* Nature, 2018.
607 **555**(7698): p. 623-628.
608 22. Milanese, A., et al., *Microbial abundance, activity, and population genomic profiling with*
609 *mOTUs.* Nature Communications, 2019. **formally accepted for publication**.
610 23. Kultima, J.R., et al., *MOCAT2: a metagenomic assembly, annotation and profiling framework.*
611 Bioinformatics, 2016. **32**(16): p. 2520-2523.
612 24. Hothorn, T., et al., *A Lego System for Conditional Inference.* Am. Stat., 2006. **60**(3): p. 257-
613 263.
614 25. Mandal, S., et al., *Analysis of composition of microbiomes: a novel method for studying*
615 *microbial composition.* Microb Ecol Health Dis, 2015. **26**: p. 27663.
616 26. Tjalsma, H., et al., *A bacterial driver-passenger model for colorectal cancer: beyond the usual*
617 *suspects.* Nat Rev Microbiol, 2012. **10**(8): p. 575-82.
618 27. Thomas, A.M., et al., *Metagenomic analysis of colorectal cancer datasets identifies cross-*
619 *cohort microbial diagnostic signatures and a link with choline degradation.* co-submitted to
620 Nature Medicine, 2018.
621 28. Huerta-Cepas, J., et al., *eggNOG 4.5: a hierarchical orthology framework with improved*
622 *functional annotations for eukaryotic, prokaryotic and viral sequences.* Nucleic Acids Res.,
623 2016. **44**(D1): p. D286-93.
624 29. Kanehisa, M., et al., *Data, information, knowledge and principle: back to metabolism in KEGG.*
625 Nucleic Acids Res., 2014. **42**(Database issue): p. D199-205.
626 30. Li, J., et al., *An integrated catalog of reference genes in the human gut microbiome.* Nat.
627 Biotechnol., 2014. **32**(8): p. 834-841.
628 31. Vieira-Silva, S., et al., *Species-function relationships shape ecological properties of the human*
629 *gut microbiome.* Nat Microbiol, 2016. **1**(8): p. 16088.
630 32. Hirayama, A., et al., *Quantitative metabolome profiling of colon and stomach cancer*
631 *microenvironment by capillary electrophoresis time-of-flight mass spectrometry.* Cancer Res,
632 2009. **69**(11): p. 4918-25.
633 33. Denkert, C., et al., *Metabolite profiling of human colon carcinoma--deregulation of TCA cycle*
634 *and amino acid turnover.* Mol Cancer, 2008. **7**: p. 72.
635 34. Mal, M., et al., *Metabotyping of human colorectal cancer using two-dimensional gas*
636 *chromatography mass spectrometry.* Anal Bioanal Chem, 2012. **403**(2): p. 483-93.
637 35. Weir, T.L., et al., *Stool microbiome and metabolome differences between colorectal cancer*
638 *patients and healthy adults.* PLoS One, 2013. **8**(8): p. e70803.
639 36. Goedert, J.J., et al., *Fecal metabolomics: assay performance and association with colorectal*
640 *cancer.* Carcinogenesis, 2014. **35**(9): p. 2089-2096.
641 37. Aykan, N.F., *Red meat and colorectal cancer.* Oncology Reviews, 2015. **9**(1).
642 38. World Cancer Research Fund / American Institute for Cancer Research, *Diet, Nutrition,*
643 *Physical Activity and Cancer: a Global Perspective*, in *Continuous Update Project Expert*
644 *Report.* 2018.
645 39. Dutilh, B.E., et al., *Screening metatranscriptomes for toxin genes as functional drivers of*
646 *human colorectal cancer.* Best Pract Res Clin Gastroenterol, 2013. **27**(1): p. 85-99.

647    40.    Sears, C.L. and W.S. Garrett, *Microbes, Microbiota, and Colon Cancer.* Cell Host Microbe,
648            2014. **15**(3): p. 317-328.
649    41.    Ridlon, J.M., et al., *Consequences of bile salt biotransformations by intestinal bacteria.* Gut
650            Microbes, 2016. **7**(1): p. 22-39.
651    42.    Yoshimoto, S., et al., *Obesity-induced gut microbial metabolite promotes liver cancer through
652            senescence secretome.* Nature, 2013. **499**(7456): p. 97-101.
653    43.    Ajouz, H., D. Mukherji, and A. Shamseddine, *Secondary bile acids: an underrecognized cause
654            of colon cancer.* World Journal of Surgical Oncology, 2014. **12**(1): p. 164.
655    44.    Boleij, A., et al., *The Bacteroides fragilis toxin gene is prevalent in the colon mucosa of
656            colorectal cancer patients.* Clin. Infect. Dis., 2015. **60**(2): p. 208-215.
657    45.    Wu, S., et al., *A human colonic commensal promotes colon tumorigenesis via activation of T
658            helper type 17 T cell responses.* Nat. Med., 2009. **15**(9): p. 1016-1022.
659    46.    Dejea, C.M., et al., *Patients with familial adenomatous polyposis harbor colonic biofilms
660            containing tumorigenic bacteria.* Science, 2018. **359**(6375): p. 592-597.
661    47.    Ridlon, J.M., D.J. Kang, and P.B. Hylemon, *Isolation and characterization of a bile acid
662            inducible 7alpha-dehydroxylating operon in Clostridium hylemonae TN271.* Anaerobe, 2010.
663            **16**(2): p. 137-46.
664    48.    Mallonee, D.H., W.B. White, and P.B. Hylemon, *Cloning and sequencing of a bile acid-
665            inducible operon from Eubacterium sp. strain VPI 12708.* Journal of Bacteriology, 1990.
666            **172**(12): p. 7011-7019.
667    49.    Ocvirk, S. and S.J.D. O'Keefe, *Influence of Bile Acids on Colorectal Cancer Risk: Potential
668            Mechanisms Mediated by Diet-Gut Microbiota Interactions.* Curr. Nutr. Rep., 2017. **6**(4): p.
669            315-322.
670    50.    Gevers, D., et al., *The treatment-naive microbiome in new-onset Crohn's disease.* Cell Host
671            Microbe, 2014. **15**(3): p. 382-392.
672    51.    Viennot, S., et al., *Colon cancer in inflammatory bowel disease: recent trends, questions and
673            answers.* Gastroenterol. Clin. Biol., 2009. **33 Suppl 3**: p. S190-201.
674    52.    Rubinstein, M.R., et al., *Fusobacterium nucleatum Promotes Colorectal Carcinogenesis by
675            Modulating E-Cadherin/β-Catenin Signaling via its FadA Adhesin.* Cell Host Microbe, 2013.
676            **14**(2): p. 195-206.
677    53.    Kostic, A.D., et al., *Fusobacterium nucleatum potentiates intestinal tumorigenesis and
678            modulates the tumor-immune microenvironment.* Cell Host Microbe, 2013. **14**(2): p. 207-215.
679    54.    Arthur, J.C., et al., *Intestinal inflammation targets cancer-inducing activity of the microbiota.*
680            Science, 2012. **338**(6103): p. 120-123.
681    55.    Reddy, B.S., *Diet and excretion of bile acids.* Cancer Res, 1981. **41**(9 Pt 2): p. 3766-8.
682    56.    Ogino, S., et al., *Integrative analysis of exogenous, endogenous, tumour and immune factors
683            for precision medicine.* Gut, 2018. **67**(6): p. 1168-1180.
684    57.    Ogino, S., et al., *Molecular pathological epidemiology of colorectal neoplasia: an emerging
685            transdisciplinary and interdisciplinary field.* Gut, 2011. **60**(3): p. 397-411.
686    58.    Hannigan, G.D., et al., *Diagnostic Potential and Interactive Dynamics of the Colorectal Cancer
687            Virome.* MBio, 2018. **9**(6).
688    59.    zur Hausen, H., *Red meat consumption and cancer: reasons to suspect involvement of bovine
689            infectious factors in colorectal cancer.* Int J Cancer, 2012. **130**(11): p. 2475-83.
690    60.    Shkoporov, A.N., et al., *Reproducible protocols for metagenomic analysis of human faecal
691            phageomes.* Microbiome, 2018. **6**(1): p. 68.
692

693

694

695 **Methods**
696

697 **Study inclusion and data acquisition**

698 We used PubMed to search for studies that published fecal shotgun metagenomic data of human
699 colorectal cancer patients and healthy controls. The search term, all hits, and the justification for
700 exclusion or inclusion are available in **Supplementary Table S1**. Raw fastq files were downloaded for
701 the four included studies from the European Nucleotide Archive, using the following ENA identifiers:
702 PRJEB10878 for [11], PRJEB12449 for [10], ERP008729 for [9], and ERP005534 for [8].
703

704 **DE study recruitment and sequencing**

705 The German (DE) study population data consist of 60 fecal CRC metagenomes, 38 of which were
706 sequenced and published in [8] under ENA accession ERP005534. The fecal metagenomes from
707 additional 22 CRC patients recruited for the same ColoCare study (DKFZ, Heidelberg, [61, 62]) were
708 sequenced later as part of this work. All fecal samples were collected after colonoscopy. Sixty gender-
709 and age-matched participants of the PRÄVENT study run by the same clinical investigators were
710 included as healthy controls; as these were not subjected to colonoscopy, the presence of
711 undiagnosed colorectal carcinomas cannot be completely ruled out but is expected to be unlikely due
712 to low prevalence of preclinical CRC in the general population [63].
713 Written informed consent was obtained from all additional 22 CRC patients and 60 controls. The study
714 protocol was approved by the institutional review board (EMBL Bioethics Internal Advisory Board) and
715 the ethics committee of the Medical Faculty at the University of Heidelberg. The study is in agreement
716 with the WMA Declaration of Helsinki and the Department of Health and Human Services Belmont
717 Report.
718 Genomic DNA was extracted from the fecal samples (preserved in RNALater) and libraries were
719 prepared as previously described [8]. Whole-genome shotgun sequencing was performed by using
720 Illumina HiSeq 2000 / 2500 / 4000 (Illumina, San Diego, USA) platforms at the Genomics Core
721 Facility, European Molecular Biology Laboratory, Heidelberg.
722

723 **Independent validation cohorts**

724 During the revision of this manuscript, we included three independent study populations for external
725 validation. Two of them were recruited in Italy (IT1 and IT2) with informed consent from all participants
726 and ethical approval by the Ethics committee of Azienda Ospedaliera of Alessandria and that of the
727 European Institute of Oncology of Milan. Shotgun fecal metagenomic data was generated as
728 described in [27].
729 The third study population was recruited in Japan (JP) with informed consent and ethical approval of
730 the institutional review boards of the National Cancer Center Japan - Research Institute and the Tokyo
731 Institute of Technology. DNA was extracted from frozen fecal samples using a GNOME DNA Isolation
732 Kit (MP Biomedicals, Santa Ana, CA) with an additional bead-beating step as previously described
733 [64]. DNA quality was assessed with an Agilent 4200 TapeStation (Agilent Technologies, Santa Clara
734 CA). After final precipitation, the DNA samples were resuspended in TE buffer and stored at -80°C
735 before further analysis. Sequencing libraries were generated with the Nextera XT DNA Sample

736 Preparation Kit (Illumina, San Diego, CA). Library quality was confirmed with an Agilent 4200
737 TapeStation. Whole-genome shotgun sequencing was carried out on the HiSeq2500 platform
738 (Illumina). All samples were paired-end sequenced with a 150-bp read length to a targeted data set
739 size of 5.0 Gb.
740
741 **Taxonomic profiling and data preprocessing**
742 The metagenomic samples were quality controlled using MOCAT2's -rtf procedure, which is based on
743 the 'solexaqa' algorithm [23]. In particular, reads that map with at least 95% sequence identity and
744 alignment length of at least 45 bp to the human genome hg19 were removed. In a second step,
745 taxonomic profiles were generated with the mOTU profiler version 2.0.0 ([22, 65, 66] – see motu-
746 tool.org and GitHub version tag 2.0.0) using the following parameters: -l 75, -g 2 and -c. Briefly, this
747 profiler is based on ten universal single-copy marker-gene families (COG0012, COG0016, COG0018,
748 COG0172, COG0215, COG0495, COG0525, COG0533, COG0541 and COG0552) [66]. These
749 marker-genes were extracted from >25,000 reference genomes and >3,000 metagenomic samples
750 allowing to profile prokaryotic species with a sequenced reference genome (ref-mOTUs) and ones
751 without (meta-mOTUs). The read count for a mOTU was calculated as median of the read count of the
752 genes that belonged to that mOTU.
753 mOTU profiles were first converted to relative abundances to account for library size. Then, profiles
754 were filtered to focus on a set of species that are confidently detectable in multiple studies.
755 Specifically, microbial species that did not exceed a maximum relative abundance of 1E-03 in at least
756 3 of the studies were excluded from further analysis, together with the fraction of unmapped
757 metagenomic reads.
758
759 **Functional metagenome profiling and data preprocessing**
760 High-quality reads (same quality filtering as for taxonomic profiling) were aligned against a combined
761 database (IGChg38 hereafter) consisting of the hg38 release of the human reference genome and the
762 integrated gene catalog (IGC) containing 9.9 million non-redundant microbial genes [30] using BWA
763 mem [67] (Version: 0.7.15-r1140) with default parameters. The purpose of adding the human genome
764 to the reference database was to filter out reads that mapped as well or better to some human
765 sequence than to any bacterial gene. Alignments were computed separately for paired-end and single
766 read libraries (single reads could result from read pairs where one read was filtered out in the quality
767 filtering procedure described above). Alignments were then filtered to only retain those longer than
768 50bp with >95% sequence identity. Then the highest scoring alignment(s) was/were kept for each
769 read. As IGChg38 is a database of predominantly genes and not genomes, there will be a substantial
770 proportion of read-pairs where one end maps within the gene while the other end does not – it either
771 maps to an adjacent gene or remains unmapped due to intergenic regions not contained in the
772 database. Therefore, we counted a whole read-pair aligning to a gene when (i) both ends from a read
773 pair map to the same gene, (ii) only one end from a read-pair maps to the gene, or (iii) a read from the
774 single read library maps to the gene. We then counted only the read-pairs that map uniquely to one
775 gene in the IGC, thus excluding ambiguous read pairs mapping with similarly high scores to multiple
776 genes in the database. For a given metagenomic sample, we further normalized the abundance of

19

each IGC gene by the length of that gene. We then estimated relative abundance of IGC genes by dividing gene abundances by the total abundance of all genes in IGC (excluding the human chromosomes).

Because metagenomes from CRC patients were not included when the IGC was constructed, we analyzed how well CRC-associated species as identified in this meta-analysis were represented in the IGC. Using a phylogenetic marker gene (COG0533), which is also used by the species profiling workflow on which the meta-analysis is based, for 24 out of the 29 core CRC-associated species we found a match in the IGC with at least 90% nucleotide identity, indicating that a sequence from the same species (above 93.1% identity) or a slightly more distant relative is present in the IGC (**Supplementary Fig. 8**).

The relative abundance of eggNOG orthologous groups [28] was estimated by summing relative abundances of genes annotated to belong to the same eggNOG orthologous group as of the most recent annotations provided by MOCAT2 [23]. To obtain KEGG orthologous groups (KO) and pathway abundances, we applied the same procedure, but using KEGG annotations for IGC provided by MOCAT2 [29].

**Overview over statistical analyses**

For univariate association testing between the abundances of microbial taxa or gene functions we used nonparametric tests throughout; all of these were two-sided Wilcoxon tests except were otherwise noted. To account for potential confounding and heterogeneity between data sets we employed a stratified version of the Wilcoxon test [24] (see below for details). ANOVA was conducted on rank-transformed data. Significance of binary co-occurrence patterns was assessed using (stratified) Cochrane-Mantel-Haenszel tests.

Multivariable analysis was done with strict separation between training and test data. This importantly also pertained to feature selection, which was either done via the LASSO [68] or by nested cross-validation procedures to avoid overoptimistic performance assessment [69] (see below for details). All samples included in this meta-analysis came from distinct individuals to ensure that generalization across subjects – rather than across timepoints within a given subject – is assessed.

**Confounder analysis**

To quantify the effect of potential confounding factors relative to that of CRC on single microbial species, we used an ANOVA-type analysis. The total variance within the abundance of a given microbial species was compared to the variance explained by disease status and the variance explained by the confounding factor akin to a linear model including both CRC status and confounding factor as explanatory variables for species abundance. Variance calculations were performed on ranks in order to account for non-Gaussian distribution of microbiome abundance data. Potential confounders with continuous values were transformed into categorical data either as quartiles or for the case of body mass index (BMI) into lean/obese/overweight according to conventional cutoffs (lean: < 25, obese: 25 - 30, overweight: > 30).

**Univariate meta-analysis for the identification of CRC-associated gut microbial species**

20

818    Significance of differential abundance was tested on a per-species basis using a blocked Wilcoxon
819    test implemented in the R coin package [24]. Informed by the results of the preceding confounder
820    analysis, we blocked for `study` and additionally `colonoscopy` in the CN study. Within this framework,
821    significance is tested against a conditional null distribution derived from permutations of the observed
822    data. Notably, permutations are performed within each block in order to control for variations in block
823    size and composition. To adjust for multiple hypothesis testing, P-values were adjusted using the
824    false-discovery rate (FDR) method [70].
825    As nonparametric effect size measures we used the area under the ROC curve (AUROC) with
826    permutation-based confidence intervals computed using the pROC package in R [71]. We further
827    developed a generalization of the (logarithmic) fold change that is widely used for other types of read
828    abundance data. This generalization is designed to have better resolution for sparse microbiome
829    profiles (where 0 entries can render median-based fold change estimates uninformative for the large
830    portion of species with a prevalence below 0.5). The generalized fold change (gFC) is computed as
831    mean difference in a set of pre-defined quantiles of the logarithmic CTR and CRC distributions (see
832    **Extended Data 3** for further details; we used quantiles ranging from 0.1 to 0.9 in increments of 0.1).
833    For the retrospective analysis of study precision and recall for detecting microbial species associations
834    from the meta-analysis, the true set was defined as the species which were associated at a given FDR
835    in the meta-analysis. Then, we checked how well this set of species would be recovered using the
836    single-study significance as determined by the Wilcoxon test. Study precision corresponds to the
837    proportion of meta-analysis significant species among those detected as significant in a single study.
838    Similarly, recall (or sensitivity) corresponds to the proportion of species out of the true set of meta-
839    analysis significant species that were recovered in a given study.
840

841    **Species co-occurrence and cluster analysis in CRC metagenomes**
842    For the analysis of gut bacterial species co-occurring in CRC microbiomes, relative abundances of the
843    core set of associated species (excluding the CRC-depleted *Clostridiales* meta-mOTU [1296]) were
844    discretized into binary values to determine whether a CRC (metagenomic) sample is "positive" or
845    "negative" for a given microbial marker. To normalize for differences in prevalence (and therefore
846    specificity) of these markers we adjusted the threshold value, above which a sample is labeled
847    "positive" based on the abundance in healthy controls. For each microbial species, the 95th percentile
848    in healthy controls was used as threshold, which effectively results in adjusting the per-marker false
849    positive rate to 0.05. Based on the binarized species-by-sample matrix, species were then clustered
850    using the Jaccard dissimilarity as implemented in the vegan package in R [72]. Associations between
851    species clusters and meta-variables were tested as 2-by-n (where n is the number of categories in the
852    meta-variable tested) contingency tables using a Cochrane-Mantel-Haenszel test with study as
853    blocking factor as implemented in the coin package [24].
854

855    **Multivariable statistical modeling workflow and model evaluation**
856    As a main goal of our work is to assess the generalization accuracy of microbiome-based CRC
857    classifiers across technical and geographic differences in patient populations, we extensively validated
858    classification models across studies taking the following two approaches.

859    In *study-to-study transfer* validation, metagenomic classifiers were trained on a single study and their

860    performance externally assessed on all other studies (off-diagonal cells in **Fig. 3ac**). Effectively we

861    implemented a nested cross validation procedure on the training study to compute within-study

862    accuracy (cells on the diagonal in **Fig. 3ac**) and tune the model hyperparameters.

863    In *leave-one-study-out* (LOSO) validation, data from one study was set aside as an external validation

864    set, while the data from the remaining 4 studies was pooled as a training set on which we

865    implemented the same nested cross validation procedure as for study-to-study transfer (see [19] for a

866    more detailed description of LOSO).

867    Data preprocessing, model building, and model evaluation was performed using the SIAMCAT R

868    package (https://bioconductor.org/packages/SIAMCAT, version 1.1.0).

869

870    **Preprocessing of taxonomic abundance profiles for statistical modeling**

871    Relative abundances were first filtered to remove markers with low overall abundance and no variance

872    (an artifact for single-study data arising from the joint data filtering described above), log-transformed

873    (after adding a pseudo-count of 1E-05 to avoid non-finite values resulting from log(0), [73]) and finally

874    standardized as z-scores. Data were split into training and test set for 10 times repeated 10-fold

875    stratified cross validation (balancing class proportions across folds). For each split, a L1-regularized

876    (LASSO) logistic regression model [68] was trained on the training set, which was then used to predict

877    the test set. The lambda parameter, i.e. regularization strength was selected for each model to

878    maximize the area under the precision recall curve under the constraint that the model contained at

879    least 5 non-zero coefficients. Models were then evaluated by calculating the area under the Receiver

880    Operating Characteristics curve (AUROC) based on the posterior probability for the CRC class.

881    In model transfer to a hold-out study, the holdout data were normalized for comparability in the same

882    way as the training dataset by using the frozen normalization function in SIAMCAT, which retains the

883    same features and re-uses the same normalization parameters (e.g. the mean of a feature for z-score

884    standardization). Then, all 100 models derived from the cross validation on the training dataset (10

885    times repeated 10-fold CV) were applied to the holdout dataset and predictions were averaged across

886    all models.

887    In the LOSO setting, data from the four training studies were jointly processed as a single dataset in

888    the same way as described above using 10 times repeated 10-fold stratified cross validation.

889

890    **Preprocessing of functional abundance profiles**

891    Functional profiles, such as eggNOG gene family or KEGG module abundance profiles were

892    preprocessed as described above for species profiles, but using 1E-06 as maximum abundance cutoff

893    and 1E-09 as a pseudo-count during log transformation. Since these abundance tables contained

894    several thousand input features we implemented an additional feature selection step, which was

895    nested properly into the cross-validation procedures as described above. This nested approach is

896    crucial to avoid over-optimistically biased performance estimates ([74], Chapter 7.10). Specifically,

897    features were filtered inside each training fold (without using any information from the test fold) by

898    selecting the 1600 features with highest single-feature AUROC values (for features depleted in CRC,

899    1 - AUROC was used for feature selection).

900

**Preprocessing of gene abundance profiles**

To ascertain the predictive power of a classifiers based on IGC gene abundances [30] we applied a series of filters to the abundance tables to reduce the number of genes that would be the input of the LASSO modelling. These filters where applied once on a per-study level and once in a leave-one-study-out (LOSO) mode, where they were applied jointly to all studies in the training set, with the remaining one being held out for external validation.

The following filters were applied in this order:

1. All genes with 0 abundance in ≥15% of samples (regardless of CRC status) were discarded.
2. The remaining data was discretized using the equal frequencies method implemented in the 'discretize' function of the sideChannelAttack R package (version 1.0-6) as a preparation to the minimal-redundancy-maximal-relevance (mRMR) algorithm [75].
3. As a feature selection procedure, mRMR (code version from 20 April 2009 downloaded from http://home.penglab.com/proj/mRMR/ on 3 Dec 2016) was run on the gene abundance table to retain the 100 top genes as output.

LASSO models were then built on log10-transformed abundances (pseudo-count of 10E-09, centered and scaled) of the sets of 100 top genes returned by mRMR. The whole process was repeated 10 times in a 5-fold stratified cross-validation scheme to allow for an estimation of the confidence of the AUROCs of the resulting models. We used the LiblineaR package (version 2.10-8) to build the LASSO models in R and tested a sequence of 20 cost parameters (equivalent or the lambda parameter controlling regularization strength) evenly spaced from $0.001^2$ to $0.2^2$. The cost parameter was selected to maximize the AUROC within the training set.

**External evaluation of disease-specificity of the metagenomic classifiers**

To assess how disease-specific the predictions of the CRC models are, we applied these to data from case-control studies investigating other human diseases. Fecal metagenomic data of patients with Parkinson's disease [12], type 2 diabetes [4, 5], and inflammatory bowel disease [6, 7] were taxonomically profiled as described above. The parameters for quality control with MOCAT2 and for the mOTU profiler were the same as described above, except for the data from [6], where we used -l 50 (to set the threshold for minimum alignment length to 50) as the read length is shorter (average read length 71) compared to the other more recently generated Illumina shotgun metagenomic data. Relative abundance data were treated exactly as another holdout dataset for each model, i.e. applying the frozen normalization prediction routines as described above. For each CRC model applied to the external datasets, a cutoff on its prediction output was adjusted to yield a false positive rate (FPR) of 0.1 on the controls of its respective (CRC) training set. Subsequently its FPR on metagenomes from patients suffering from the above-mentioned (non-CRC) conditions was assessed to evaluate its disease specificity. The rationale behind this is that a metagenomic classifier recognizing general features of dysbiosis would be expected to predict CRC patients and those suffering from other conditions at a similar rate; such a classifier would thus in the above-described evaluation display a much higher FPR than on the controls of its training set. In contrast maintaining a low FPR in this

940  evaluation indicates that the classification model is based on CRC-specific features rather than
941  hallmarks of general dysbiosis or nonspecific inflammation.
942

943  **Functional profiling of gut metabolic modules (GMMs)**
944  Gut metabolic modules were computed as originally proposed [31], using the KEGG KO profiles based
945  on the IGC (see **Functional metagenome profiling** above) as input. Statistical analysis and
946  generalized fold change calculations were performed analogously to species profiles (see above). Gut
947  metabolic modules were summarized across functional groups (e.g. amino acid degradation) as
948  geometric mean of all modules within the respective group.
949

950  **Targeted functional analysis of virulence and toxicity pathways of potential relevance in CRC**
951  To investigate toxins and virulence mechanisms that have previously been implicated with CRC [40],
952  we constructed for each gene belonging to the respective virulence or toxicity pathway a hidden
953  Markov model (HMM). Each HMM was built from a multiple sequence alignment generated by
954  MUSCLE [76], containing the respective reference sequences and close homologs identified using
955  PSI-Blast [77]. Multiple sequence alignments are available together with the code for this paper
956  (https://github.com/zellerlab/crc_meta). Then, we screened the IGC metagenomic gene catalogue [30]
957  with each HMM using the HMMER software (version 3.1b2) [78]. Genes with an E-value below 1E-10
958  were filtered for uniqueness, since in some cases the HMMs would call different regions in the same
959  gene. For single gene virulence factors (i.e. *fadA* and *bft*), potential IGC hits were aligned against the
960  reference sequence using the Needleman-Wunsch algorithm in the EMBOSS package [79]. Hits were
961  then filtered based on percentage of sequence identity (cutoff: 40%) and sequence similarly to the
962  species relative abundance profiles based on maximum relative abundance (cutoff: 1E-07) in order to
963  exclude genes with limited relevance. Statistical analysis was performed on the sum of all genes.
964  For virulence pathways containing more than one gene, the IGC hits of each functional group within
965  the pathway were aligned against the respective reference sequence and filtered for percentage of
966  sequence identity and maximum abundance. Then, all hits were clustered based on the Pearson
967  correlation of the log-abundances across all samples using the Ward algorithm as implemented in the
968  *hclust* function in R. The gene clusters were filtered based on operon completeness (how many genes
969  of the operon were present in the cluster) and average correlation within the cluster (**Extended Data**
970  **9**). For statistical analysis, the genes in the selected gene clusters were summed up within each group
971  or all together for the overall analysis.
972
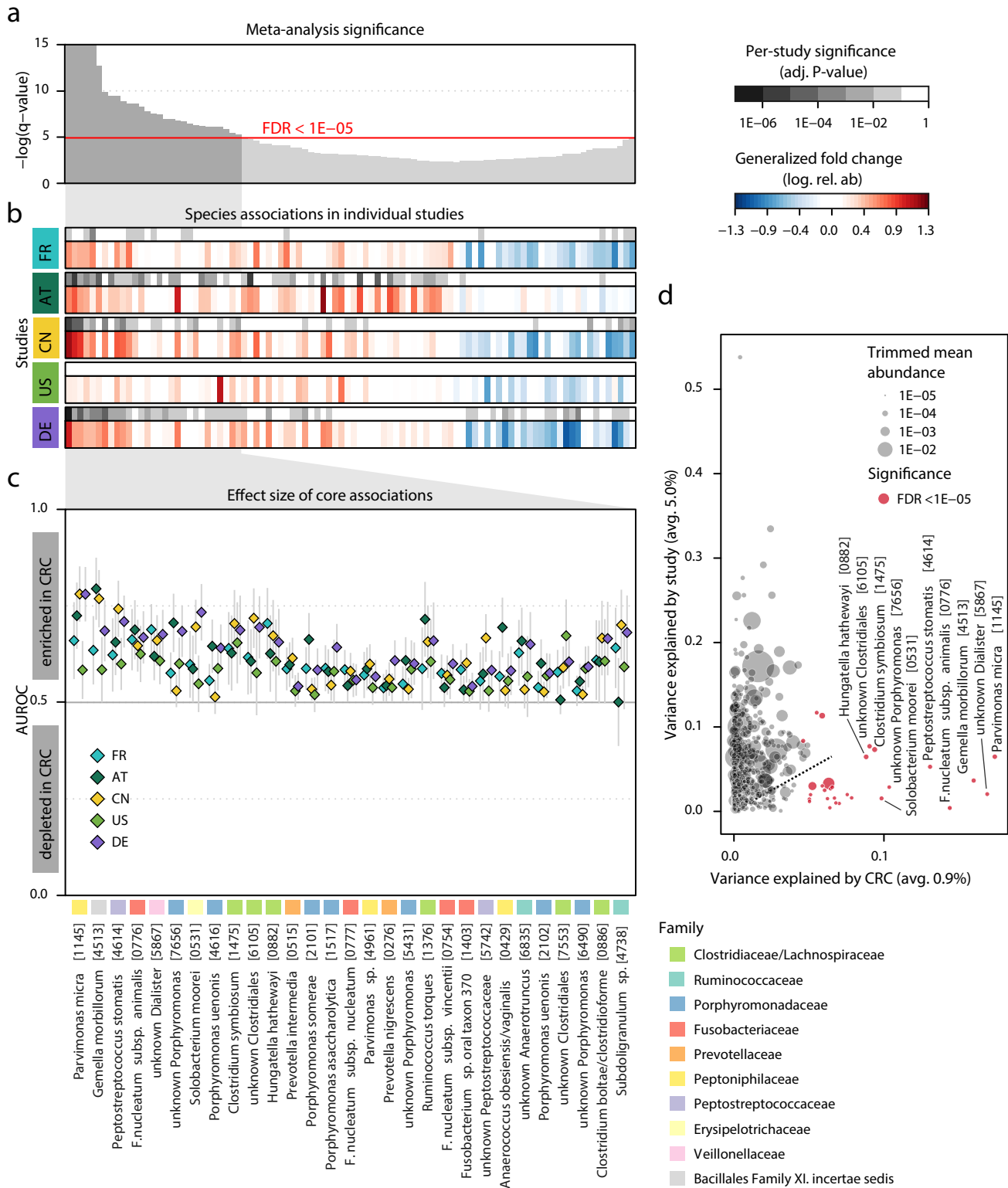
973  **Quantitative PCR for *baiF***
974  Real-time quantitative PCR to quantify the abundance and expression of *baiF* was performed on a
975  subset of samples in the DE cohort (20 control and 24 colorectal cancer samples, see
976  **Supplementary Table S6**). For these samples, DNA and RNA extraction was done with the Allprep
977  PowerFecal DNA/RNA kit (Qiagen, Cat No: 80244) with additional RNAse and DNAse digestion steps,
978  respectively, as described by the manufacturer. DNA and RNA concentrations were determined by
979  Qubit Fluorometer (Invitrogen) and quality control of all RNA samples was done using an Agilent 2100
980  Bioanalyzer in combination with RNA 6000 Nano and Pico LabChip kits.

981 First-strand cDNA was synthesized by SuperScript IV VILO Master Mix with ezDNAse enzyme and
982 random hexamer primers (Invitrogen, catalogue number 11766500) as recommended by the
983 manufacturer. Reaction were performed as described in the protocol with one minor change of
984 temperature (incubation for the reverse transcription step at 55°C).

985 To quantify *baiF* relative to the total bacterial RNA/DNA in a sample, qPCR was performed in
986 triplicates for 16S rRNA and the *baiF* genes, using both cDNA and genomic DNA (gDNA) as template.
987 We used the following primers for *baiF*: TTCAGYTTCTACACCTG (forward),
988 GGTTRTCCATRCCGAACAGCG (reverse), and standard primers F515 and R806 for 16S [80]. RT-
989 PCR reactions were prepared with a final primer concentration of 0.5 $\mu$M, including 5 ng of genomic
990 DNA or 10 ng of cDNA in 20 $\mu$l final reaction volume, and reactions were performed with SYBR Green
991 qPCR mix on StepOne Real-Time PCR system (Thermo Fisler Scientific). Cycling conditions were as
992 follows; initial denaturation of 95°C for 10 min, then 40 cycles of denaturing at 95°C for 15 s, annealing
993 at 60°C for 60 s followed by melt curve analysis.

994 Delta-Ct values were calculated as difference between *baiF* and 16S Ct values. Significance of the
995 comparison between control and colorectal cancer samples was tested on the delta-Ct values using a
996 one-sided Wilcoxon test as a confirmation of metagenomic enrichment.
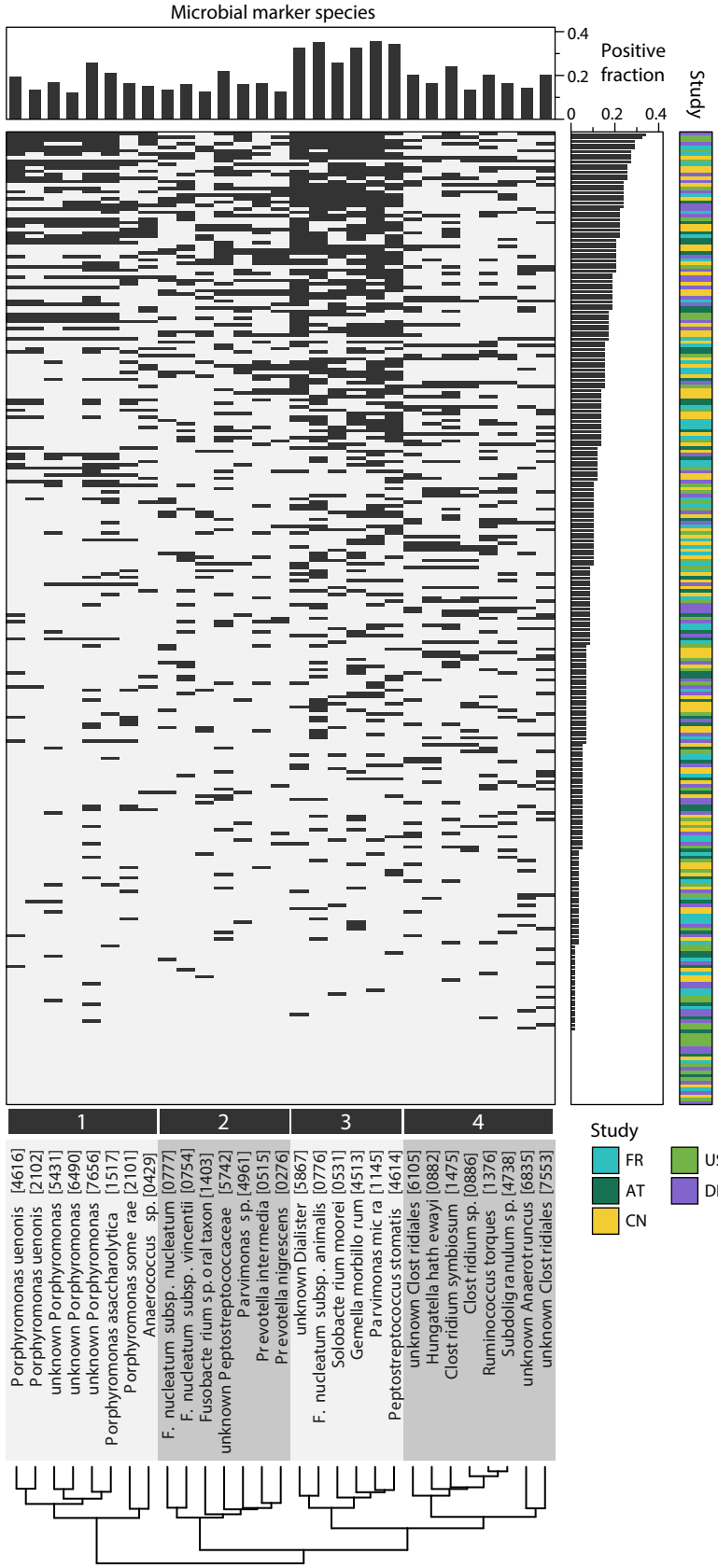
997

998 **Additional References**

999 61.    Bohm, J., et al., *Discovery of novel plasma proteins as biomarkers for the development of*
1000        *incisional hernias after midline incision in patients with colorectal cancer: The ColoCare study.*
1001        Surgery, 2017. **161**(3): p. 808-817.
1002 62.    Liesenfeld, D.B., et al., *Metabolomics and transcriptomics identify pathway differences*
1003        *between visceral and subcutaneous adipose tissue in colorectal cancer patients: the ColoCare*
1004        *study.* Am J Clin Nutr, 2015. **102**(2): p. 433-43.
1005 63.    Pox, C.P., et al., *Efficacy of a nationwide screening colonoscopy program for colorectal*
1006        *cancer.* Gastroenterology, 2012. **142**(7): p. 1460-7 e2.
1007 64.    Furet, J.P., et al., *Comparative assessment of human and farm animal faecal microbiota using*
1008        *real-time quantitative PCR.* FEMS Microbiol Ecol, 2009. **68**(3): p. 351-62.
1009 65.    Mende, D.R., et al., *Accurate and universal delineation of prokaryotic species.* Nat. Methods,
1010        2013. **10**(9): p. 881-884.
1011 66.    Sunagawa, S., et al., *Metagenomic species profiling using universal phylogenetic marker*
1012        *genes.* Nat. Methods, 2013. **10**(12): p. 1196-1199.
1013 67.    Li, H. and R. Durbin, *Fast and accurate short read alignment with Burrows-Wheeler transform.*
1014        Bioinformatics, 2009. **25**(14): p. 1754-60.
1015 68.    Tibshirani, R., *Regression Shrinkage and Selection via the Lasso.* J.R. Stat. Soc. Series B
1016        Stat. Methodol., 1996. **58**(1): p. 267-288.
1017 69.    Smialowski, P., D. Frishman, and S. Kramer, *Pitfalls of supervised feature selection.*
1018        Bioinformatics, 2010. **26**(3): p. 440-3.
1019 70.    Benjamini, Y. and Y. Hochberg, *Controlling the false discovery rate: a practical and powerful*
1020        *approach to multiple testing.* J. R. Stat. Soc. Series B Stat. Methodol., 1995. **57**(1): p. 289–
1021        300.
1022 71.    Robin, X., et al., *pROC: an open-source package for R and S+ to analyze and compare ROC*
1023        *curves.* BMC Bioinformatics, 2011. **12**(1).
1024 72.    Oksanen, J., et al., *vegan: Community Ecology Package.* 2018.
1025 73.    Costea, P.I., et al., *A fair comparison.* Nat. Methods, 2014. **11**(4): p. 359-359.
1026 74.    Hastie, T., R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning: Data Mining,*
1027        *Inference, and Prediction.* 2013: Springer Science & Business Media. 536.
1028 75.    Peng, H., F. Long, and C. Ding, *Feature selection based on mutual information: criteria of*
1029        *max-dependency, max-relevance, and min-redundancy.* IEEE Trans Pattern Anal Mach Intell,
1030        2005. **27**(8): p. 1226-38.

1031   76.    Edgar, R.C., *MUSCLE: multiple sequence alignment with high accuracy and high throughput.*
1032          Nucleic Acids Res., 2004. **32**(5): p. 1792-1797.
1033   77.    Altschul, S.F., et al., *Gapped BLAST and PSI-BLAST: a new generation of protein database*
1034          *search programs.* Nucleic Acids Res., 1997. **25**(17): p. 3389-3402.
1035   78.    Eddy, S.R., *Accelerated Profile HMM Searches.* PLoS Comput. Biol., 2011. **7**(10): p.
1036          e1002195.
1037   79.    Rice, P., I. Longden, and A. Bleasby, *EMBOSS: the European Molecular Biology Open*
1038          *Software Suite.* Trends Genet., 2000. **16**(6): p. 276-277.
1039   80.    Caporaso, J.G., et al., *Global patterns of 16S rRNA diversity at a depth of millions of*
1040          *sequences per sample.* Proc Natl Acad Sci U S A, 2011. **108 Suppl 1**: p. 4516-22.
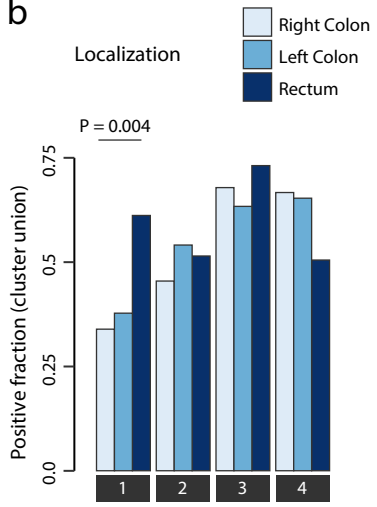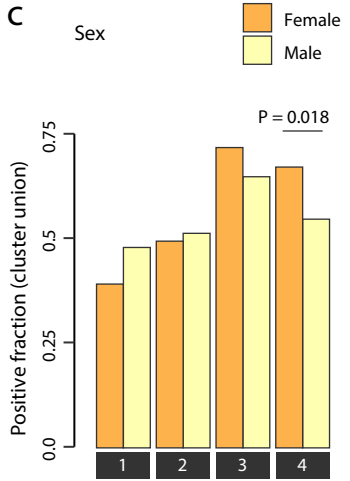1041

**a** Meta-analysis significance

**b** Species associations in individual studies

Per-study significance (adj. P-value)

1E−06   1E−04   1E−02   1

Generalized fold change (log. rel. ab)

−1.3  −0.9  −0.4  0.0  0.4  0.9  1.3

**c** Effect size of core associations

**d**

Trimmed mean abundance

1E−05
1E−04
1E−03
1E−02

Significance
FDR <1E−05

Variance explained by study (avg. 5.0%)

Variance explained by CRC (avg. 0.9%)

Hungatella hathewayi [0882]
unknown Clostridiales [6105]
Clostridium symbiosum [1475]
unknown Porphyromonas [7656]
Peptostreptococcus stomatis [4614]
Solobacterium moorei [0531]
F.nucleatum subsp. animalis [0776]
Gemella morbillorum [4513]
unknown Dialister [5867]
Parvimonas micra [1145]

Family

Clostridiaceae/Lachnospiraceae
Ruminococcaceae
Porphyromonadaceae
Fusobacteriaceae
Prevotellaceae
Peptoniphilaceae
Peptostreptococcaceae
Erysipelotrichaceae
Veillonellaceae
Bacillales Family XI. incertae sedis

**a**

Microbial marker species

Positive fraction

Study

( Sample positive for this marker species )

Colorectal cancer metagenomes

1    2    3    4

Porphyromonas uenonis [4616]
Porphyromonas uenonis [2102]
unknown Porphyromonas [5431]
unknown Porphyromonas [6490]
unknown Porphyromonas [7656]
Porphyromonas asaccharolytica [1517]
Porphyromonas some rae [2101]
Anaerococcus sp. [0429]
F. nucleatum subsp. nucleatum [0777]
F. nucleatum subsp. vincentii [0754]
Fusobacterium s p.s oral taxon [1403]
unknown Peptostreptococcaceae [5742]
Parvimonas s.p. [4961]
Prevotella intermedia [0515]
Prevotella nigrescens [0276]
unknown Dialister [5867]
F. nucleatum subsp. animalis [0776]
Solobacte rium moorei [0531]
Gemella morbillo rum [4513]
Parvimonas mic ra [1145]
Peptostreptococcus stomatis [4614]
unknown Clost ridiales [6105]
Hungatella hath ewayi [0882]
Clost ridium symbiosum [1475]
Clost ridium sp. [0886]
Ruminococcus torques [1376]
Subdolig ranulum s p. [4738]
unknown Anaerot runcus [6835]
unknown Clost ridiales [7553]

Study

FR    US
AT    DE
CN

**b**

Localization

Right Colon
Left Colon
Rectum

P = 0.004

Positive fraction (cluster union)

1    2    3    4

**c**

Sex

Female
Male

P = 0.018

Positive fraction (cluster union)

1    2    3    4

**d**

CRC Stage

Early
Late

P = 0.049

Positive fraction (cluster union)

1    2    3    4

Cluster

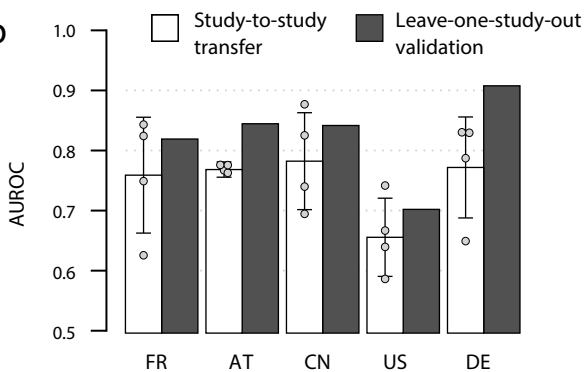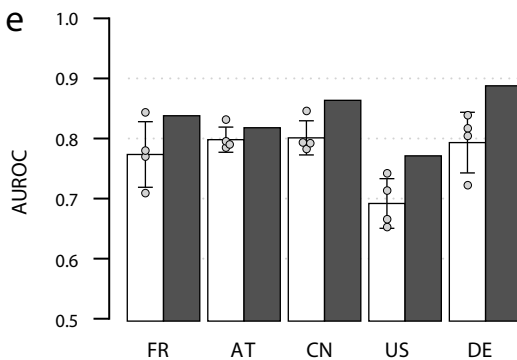**a** Classification models based on species

Test set

| Training set | FR (N = 114) | AT (N = 109) | CN (N = 128) | US (N = 104) | DE (N = 120) | Model average |
|---|---|---|---|---|---|---|
| FR | 0.85 | 0.76 | 0.82 | 0.64 | 0.83 | 0.76 |
| AT | 0.62 | 0.92 | 0.74 | 0.59 | 0.65 | 0.65 |
| CN | 0.82 | 0.76 | 0.81 | 0.67 | 0.83 | 0.77 |
| US | 0.76 | 0.78 | 0.70 | 0.74 | 0.79 | 0.76 |
| DE | 0.84 | 0.79 | 0.88 | 0.74 | 0.79 | 0.81 |

**d** Classification models based on eggNOG

Test set

| Training set | FR (N = 114) | AT (N = 109) | CN (N = 128) | US (N = 104) | DE (N = 120) | Model average |
|---|---|---|---|---|---|---|
| FR | 0.86 | 0.79 | 0.85 | 0.66 | 0.82 | 0.78 |
| AT | 0.71 | 0.80 | 0.79 | 0.67 | 0.72 | 0.72 |
| CN | 0.84 | 0.79 | 0.87 | 0.74 | 0.81 | 0.80 |
| US | 0.78 | 0.83 | 0.80 | 0.75 | 0.84 | 0.81 |
| DE | 0.77 | 0.80 | 0.78 | 0.72 | 0.93 | 0.77 |

**b** Study-to-study transfer / Leave-one-study-out validation

**e**

**c** Other conditions — 10% FPR

| Abbr. | Condition | N |
|---|---|---|
| CTR | Healthy controls from meta-analysis | 290 |
| T2D | Type 2 diabetes | 201 |
| PD | Parkinson's disease | 31 |
| UC | Ulcerative colitis | 98 |
| CD | Crohn's disease | 63 |

a



b



c