

1     **EXPERT: Transfer Learning-enabled context-aware microbial**  
2   **source tracking**

3  
4     Hui Chong<sup>1,2,#</sup>, Qingyang Yu<sup>1,#</sup>, Yuguo Zha<sup>1,#</sup>, Guangzhou Xiong<sup>1</sup>, Nan Wang<sup>1</sup>, Xinhe  
5     Huang<sup>1</sup>, Shijuan Huang<sup>1</sup>, Chuqing Sun<sup>1</sup>, Sicheng Wu<sup>1</sup>, Wei-Hua Chen<sup>1</sup>, Luis Pedro Coelho<sup>2</sup>,  
6   Kang Ning<sup>1,\*</sup>

7     <sup>1</sup>Key Laboratory of Molecular Biophysics of the Ministry of Education, Hubei Key  
8     Laboratory of Bioinformatics and Molecular-imaging, Center of AI Biology, Department of  
9     Bioinformatics and Systems Biology, College of Life Science and Technology, Huazhong  
10    University of Science and Technology, Wuhan 430074, Hubei, China

11    <sup>2</sup>Institute of Science and Technology for Brain-Inspired Intelligence, Fudan University,  
12    Shanghai 200433, China

13    <sup>#</sup>These authors contributed equally to this work

14    <sup>\*</sup>Correspondence should be addressed to K.N (Email: ningkang@hust.edu.cn)

15

16    **Abstract**

17    Microbial source tracking quantifies the potential origin of microbial communities, facilitates  
18    better understanding of how the taxonomic structure and community functions were formed  
19    and maintained. However, previous methods involve a tradeoff between speed and accuracy,  
20    and have encountered difficulty in source tracking under many context-dependent settings.  
21    Here, we present EXPERT for context-aware microbial source tracking, in which we adopted  
22    a Transfer Learning approach to profoundly elevate and expand the applicability of source  
23    tracking, enabling biologically informed novel microbial knowledge discovery. We  
24    demonstrate that EXPERT can predict microbial sources with performance superior to other  
25    methods in efficiency and accuracy. More importantly, we demonstrate EXPERT's context-

26 aware ability on several applications, including tracking the progression of infant gut  
27 microbiome development and monitoring the changes of gut microbiome for colorectal  
28 cancer patients. Broadly, transfer learning enables accurate and context-aware microbial  
29 source tracking and has the potential for novel microbial knowledge discovery.

30

## 31 **Introduction**

32 Millions of microbial communities have been accumulated from hundreds of environments  
33 (also known as “biomes”) worldwide<sup>1-3</sup>, which continuously complete the grand picture of  
34 the microbiome world, revolutionizing our understanding of the roles microbes play in  
35 human health and disease<sup>4</sup>, biogeochemical cycling<sup>5</sup> and other processes. The relationships  
36 between microbial community samples and their biomes, on the other hand, are extremely  
37 complicated, owing to the highly dynamic nature of microbial communities and our limited  
38 understanding of how they function. Microbial source tracking (MST) quantifies the potential  
39 origin of microbial communities, thereby could help us to understand how the taxonomic  
40 structure and community functions were formed and maintained<sup>6</sup>. In previous studies, MST  
41 has been widely used to quantify the contamination present in (1) environment<sup>7</sup> and (2) host  
42 driven by the contact of host and environment<sup>8</sup> (e.g., human skin and exposed environments).

43 MST has the potential to go beyond the scope of microbial contamination in a variety of  
44 contexts, including estimating microbial restoration of cesarean-born infants<sup>9</sup>, quantifying the  
45 microbial community differences across diseases, as well as characterizing the gut microbial  
46 communities of patients during cancer progression. In these contexts, MST could reveal  
47 dynamic patterns of microbial communities, and provide insights into the effect of microbial  
48 communities on health care for newborns, chronic diseases, and cancer.

49 Current methods for source tracking, though having made substantial contributions, have  
50 limitations in accuracy and scope of application. SourceTracker<sup>10</sup> and FEAST<sup>9</sup> model the  
51 query (sink) community as a mixture of sources, and estimate source contribution through  
52 Markov Chain Monte Carlo (MCMC) and Expectation-Maximization (EM), which, however,  
53 leads to a tradeoff between running time and accuracy, and source tracking among thousands

54 of sources take hours<sup>9,10</sup>. Also, these two methods are heavily dependent on reference  
55 samples, leading to the needs of manually selecting possible source samples, rendering source  
56 tracking incomplete and error-prone. On the other hand, Random Forest<sup>11</sup> and ONN4MST<sup>10</sup>  
57 utilize supervised learning models for source contribution estimation<sup>6</sup> but are constrained by  
58 adaptivity: Models cannot be directly applied to other MST tasks once they've been built for  
59 specific context.

60 To address these limitations, we developed EXPERT, a method based on an adaptive Neural  
61 Network (NN) framework and Transfer Learning, for solving the MST problem. Previous  
62 studies have shown that Transfer Learning can significantly expand the applicability of  
63 supervised learning models<sup>12</sup>. Here Transfer Learning is used to introduce existing  
64 knowledge learnt from other microbial samples to diverse contexts and to facilitate source  
65 tracking in a context-aware manner. Systematic assessments have shown EXPERT's  
66 capability of quantifying the potential contribution of sources in a fast and accurate manner,  
67 and adapting source tracking in different context-dependent settings. More importantly, we  
68 demonstrate the utility of EXPERT in several representative contexts, including tracking the  
69 development of the infant gut microbiome, as well as tracking the progression of gut  
70 microbiome changes in patients with colorectal cancer (CRC). Transfer Learning-enabled  
71 EXPERT, we reasoned, could make it easier to discover novel microbial knowledge across a  
72 wide range of applications in these contexts.

73

## 74 **Results**

### 75 **Rationale, adaptive modeling, and multi-faceted applications of EXPERT**

76 EXPERT is a context-aware method for MST that employs both the adaptive NN and  
77 Transfer Learning<sup>12</sup> frameworks, enabling knowledge transfer of MST models. The adaptive  
78 NN framework constructs MST models according to a given MST task (**Methods,**  
79 **Supplementary Fig. S1**). Together with Transfer Learning, EXPERT can automatically  
80 construct MST models and utilize the knowledge of fundamental models (i.e., existing MST  
81 models) to aid in the learning of the newly constructed models. In our study, three  
82 fundamental models were introduced for knowledge transfer (**Supplementary Table S1-S5**):  
83 the general model (GM, trained and validated on 118,592 communities from 131  
84 representative biomes), the human model (HM, trained and validated on 52,537 communities  
85 from 27 human-associated biomes), and the disease model (DM, trained and validated on  
86 13,642 fecal communities from patients of 19 diseases and healthy controls). Additionally,  
87 EXPERT utilizes Multi-task Learning<sup>13</sup>, which enables hierarchical MST (**Methods,**  
88 **Supplementary Fig. S1-S2**).

89 The knowledge transfer process of EXPERT is illustrated in **Fig. 1a**. EXPERT adopted the  
90 rationale of Transfer Learning<sup>12</sup>, allowing context-aware MST through three steps, namely  
91 transfer, adaptation, and fine-tuning: In the transfer step, EXPERT adapts the fundamental  
92 model to an MST context; in the adaptation and fine-tuning steps, EXPERT optimizes the  
93 parameters (**Methods, Supplementary Note 1**). The contextualized model can serve a broad-  
94 spectrum of source tracking applications (**Fig. 1b**)

95

### 96 **Efficiency, accuracy, and adaptivity of EXPERT**

97 Benchmark tests have demonstrated EXPERT's superior efficiency, accuracy, and adaptivity  
98 for MST (**Fig. 2**). Specifically, it outperforms Sourcetracker<sup>6</sup> and FEAST<sup>9</sup> in terms of  
99 efficiency, while outperforming the NN approach<sup>10</sup> in terms of accuracy and adaptability. In  
100 this part, we assessed these capabilities using 52,537 communities from 27 human-associated  
101 biomes (**Supplementary Table S1, S3, Fig. S3**).

102 We have compared the performance of EXPERT with FEAST, as SourceTracker was  
103 similarly accurate but slower than FEAST<sup>9</sup>. To compare EXPERT's accuracy and efficiency  
104 with FEAST, we considered community samples from seven biomes (**Supplementary Fig.**  
105 **S3**) as sources, and randomly selected small sets of community samples out of these for  
106 comparison as well (**Fig. 2a, Methods**). As a result, EXPERT could simultaneously reach  
107 high accuracy and efficiency (Maximal F1-measure F-max = 0.923, over 200 queries/second,  
108 **Fig. 2a**). While FEAST faces a severe tradeoff between accuracy and efficiency: FEAST's  
109 accuracy improves as it uses more samples for each biome (F-max = 0.847, 0.884, and 0.911)  
110 while efficiency declines nearly exponentially (0.06, 0.02 and 0.005 queries/second, **Fig. 2a**,  
111 **Supplementary Table S6**).

112 We also compared EXPERT's accuracy with the NN approach, by using different proportions  
113 of source samples (**Fig. 2b**). As the NN approach cannot be directly applied in this context  
114 (i.e., NN has no adaptivity), we have manually implemented a model for the comparison. The  
115 result showed that the EXPERT model outperforms the NN approach on accuracy: while the  
116 MST accuracy steadily increased with the increasing proportion of source samples used, the  
117 EXPERT model only required 10% of source samples to achieve a validation F-max of 0.814,  
118 while the NN approach required three times as many samples to reach a similar validation F-  
119 max of 0.813 (**Supplementary Table S7**). This demonstrated that EXPERT models were  
120 able to “understand” the contextualized microbial community profiles based on only a small  
121 fraction of samples. Notably, as the fine-tuning optimization clearly improved the accuracy  
122 (**Fig. 2b**), the knowledge transfer with fine-tuning was considered the default setting in the  
123 following sections.

124

### 125 **Adaptation to newly introduced microbiome data**

126 In this context, we aim to validate EXPERT's utility in adapting to newly introduced  
127 microbial community samples. Such data could be obtained through new sequencing and  
128 analytical technologies or originate from rarely studied environments. To test EXPERT's  
129 capability in such context, in addition to the 118,592 communities accessed as of Jan. 2020  
130 from MGnify (referred to as “baseline data”, **Supplementary Table S1, S2, Fig. 3a**), we

131 selected 34,209 communities from MGnify between Jan. 2020 and Oct. 2020 (referred to as  
132 “newly introduced data”, **Fig. 3a, Supplementary Table S1, S8, Fig. S4**). Among the newly  
133 introduced data, there are 30,788 communities that originated from biomes included in the  
134 baseline data as well, and 3,421 communities that originated from newly introduced biomes  
135 (**Supplementary Fig. S4**).

136 We first tested the applicability of the general model and EXPERT framework on the newly  
137 introduced data. In this context we only considered the 30,788 communities. We directly  
138 applied the general model (built based on the baseline data, AUROC = 0.982 by cross-  
139 validation) on the data, and obtained a much lower accuracy (AUROC = 0.884,  
140 **Supplementary Note 2**). The reason behind this might be the data heterogeneity and batch  
141 effect between the two datasets (**Supplementary Fig. S5**). However, by using EXPERT, we  
142 could adapt the general model to the newly introduced data, reduce the influence of batch  
143 effect on source tracking analysis, and maintain or even further improve the accuracy  
144 (AUROC = 0.993, **Fig. 3b**).

145 We also tested the applicability of EXPERT on the newly introduced biomes, the results have  
146 shown that based on the EXPERT framework, the general model could also adapt to the  
147 newly introduced biomes (AUROC = 0.988), though the newly introduced biomes were not  
148 included in the baseline data. As demonstrated by these results, EXPERT has the potential for  
149 extending fundamental models into previously unexplored contexts.

150

### 151 **Context-aware microbial source tracking applications**

152 We then demonstrate EXPERT’s utility in context-aware MST, by focusing on patterns of the  
153 human gut microbiome in different contexts: (1) early development of gut microbial  
154 communities for infants, (2) association of gut microbial communities with different types of  
155 diseases, and (3) association of gut microbial communities with the progression of colorectal  
156 cancer. In these contexts, we consider the quantified source contribution generated from  
157 EXPERT as a measure to determine the host status.

158

## 159 **The succession of infant gut microbial communities**

160 We next used EXPERT to characterize small compositional changes among infant gut  
161 microbial communities during the first year of life. Under this circumstance, we could  
162 investigate the dynamic patterns of gut microbial communities from a specific period of life.  
163 For instance, if infant samples from multiple time points and sources are present, EXPERT  
164 can estimate how much of microbial community in the infant's gut originated from birth and  
165 subsequent time points. To confirm this capability, we used longitudinal data from Backhed  
166 et al.<sup>14</sup>, including fecal samples from 98 infants and their mothers, delivered by vaginal  
167 delivery or cesarean section (**Fig. 4a, Supplementary Table S1 and S9**). In this part of the  
168 study, we considered samples from infants at 12 months of age as queries, and samples from  
169 earlier time points or mothers as sources.

170 Based on the hierarchy that divided samples by sampling time first followed by delivery  
171 mode (**Fig. 4a**), we noticed that for infant gut microbial communities at 12 months of age, the  
172 maternal contribution is dominant (**Fig. 4b**). Moreover, there is no significant difference in  
173 the maternal contribution between cesarean-born and vaginal-born infants (Wilcoxon test,  $p =$   
174  $0.929$ , **Fig. 4b**), consistent with Principal Coordinate Analysis (PCoA) using distance metric  
175 either in weighted-UniFrac<sup>15</sup> or Jensen Shannon divergence<sup>16</sup> (**Fig. 4c and Supplementary**  
176 **Fig. S6**). We concluded that the infant gut at 12 months is largely adapted to exposed  
177 environments, resulting in an insignificant difference between samples collected from hosts  
178 of different delivery modes, consistent with previous studies<sup>17,18</sup>.

179 We then assessed the utility of different fundamental models in this context by also  
180 introducing the general model and changing the source biome hierarchy (**Supplementary Fig.**  
181 **S7, S8**). We found that the human model can facilitate MST in this context with significantly  
182 better performance compared with the general model (Transfer (HM) AUROC = 0.773,  
183 Transfer (GM): AUROC = 0.720, Wilcoxon test,  $p = 0.072$ ), suggesting the use of the human  
184 model in this application. Therefore, we suggest that when using EXPERT, it is necessary to  
185 choose a proper fundamental model according to the specific context (**Fig. 4d**).

186

## 187 **EXPERT reveals disease-specific patterns within gut microbial communities**

188 The pattern of gut microbial communities could be disease-specific<sup>19</sup>, reflecting the distinct  
189 inflammation patterns across diseases. In this context, we aimed to demonstrate EXPERT's  
190 utility in characterizing human gut microbial communities associated with different types of  
191 diseases. Using EXPERT, we can measure patterns across multiple diseases. Specifically, we  
192 assembled a large gut microbial community dataset, including 13,462 communities  
193 representing 19 diseases (**Fig. 5b**) and healthy controls, collected from 101 studies and 27  
194 countries (**Fig. 5a, Supplementary Table S1, and S4**). We also introduced the human model  
195 to characterize these diseases. By randomly selecting 10% samples of the dataset as queries,  
196 and considering the remaining samples as microbial sources, we aim to characterize the  
197 pattern across (1) different patients of the same disease, and (2) patients with different  
198 diseases. The results revealed that, except for Crohn's disease, the pattern is shared across  
199 patients with the same disease, but not shared across patients with different diseases (**Fig. 5c**).  
200 This is consistent with a previous study<sup>19</sup>, which found disease-specific patterns within the  
201 human gut microbial communities. Among these diseases examined in this study, we  
202 discovered the disease-specific pattern to Liver Cirrhosis and Irritable Bowel Syndrome,  
203 which had not been reported in a previous cross-disease study<sup>19</sup>.

204 We further validated the disease-specific patterns by utilizing the Independent model, which  
205 was constructed entirely from the same samples. We found that both Independent model and  
206 Transfer (HM) model could distinguish diseases with high AUROC, and confirmed that the  
207 gut microbial communities may be used to discriminate between these diseases (AUROC  
208 over 0.800 for most phenotypes, **Fig. 5d,e**). This demonstrated the utility of EXPERT in  
209 large-scale MST analysis, particularly when comparing a wide variety of microbial  
210 communities from multiple environments.

211

## 212 **EXPERT characterizes gut microbial communities during cancer progression**

213 Gut microbial communities undergo compositional changes as cancer progresses, and this can  
214 be observed in the human gut microbiota, which has been shown to influence the progression  
215 of colorectal cancer (CRC)<sup>20</sup>. In this context, we demonstrate EXPERT's utility in  
216 characterizing the progression of CRC using human gut microbiota. We assessed the



217 applicability of EXPERT by leveraging the disease model with Transfer Learning (**Fig. 6a**).  
218 We considered 635 samples from five stages in the progression of CRC: 0 (Healthy control) I,  
219 II, III, and IV according to the study of Zeller G. et al.<sup>21</sup> (**Fig. 6b, Supplementary Table S1,**  
220 **and S10**). Preliminary analysis using traditional methods<sup>15,16</sup> could not show the  
221 compositional shifts of the human gut within such progression (**Fig. 6c, Supplementary Fig.**  
222 **S9**). However, by randomly selecting 10% of the dataset as queries, and estimating their  
223 resemblant signatures from the remaining samples using the Transfer (DM) model, we found  
224 that for gut microbial communities at each CRC stage, a large proportion of microbes could  
225 also be found in the communities at the same stage: 0.24, 0.49, 0.56, 0.37 and 0.51 for stage 0  
226 to IV (**Fig. 6d**). These results indicated the association between gut microbiota and CRC  
227 progression and suggested the potential of gut microbiota for tracking the progression of  
228 CRC<sup>21,22</sup>.

229 We also assessed the EXPERT's capability in monitoring the progression of CRC, by  
230 comparing the performance of different models: For comparison, we generated a Transfer  
231 (HM) and an Independent model (solely based on the CRC samples) in addition to the  
232 Transfer (DM) model. Results have shown that Transfer (DM) achieved a better performance  
233 (AUROC = 0.977, **Fig. 6e, i**) among these three models, highlighting the EXPERT's utility  
234 on tracking the different stages of CRC progression using gut microbial communities.

235

## 236 **Discussion**

237 Broadly, EXPERT adopted a Transfer Learning approach to profoundly elevate and expand  
238 the applicability of source tracking, enabling biologically informed novel microbial  
239 knowledge discovery. Based on the NN approach and Transfer Learning technique, it could  
240 quickly adapt the supervised model for source tracking tasks in different contexts, thus  
241 providing a fast, accurate, and context-aware computational approach that enables MST  
242 analyses in diverse contexts, for in-depth knowledge discovery.

243 Our analytical results have confirmed that EXPERT has enabled MST with high speed and  
244 fidelity, without the need for pre-defined source samples. Additionally, EXPERT could adapt

245 the fundamental models to newly introduced data, and help reduce the influence of data  
246 heterogeneity and batch effects. More importantly, we have shown that MST solely based on  
247 the fundamental models may be biased by batch effects, whereas EXPERT can significantly  
248 mitigate this influence.

249 We have demonstrated EXPERT's utility in context-aware MST in several applications. First,  
250 EXPERT can characterize the tiny compositional difference associated with environmental  
251 changes. By adapting the human model to microbial communities of infant gut across  
252 delivery modes, we found that due to environmental exposure, cesarean-born infants have a  
253 largely restored gut microbial community compared with infants born vaginally, consistent  
254 with the results of other published analyses<sup>17,18</sup>. Secondly, we demonstrated the utility of  
255 EXPERT beyond traditional MST methods by incorporating a dataset of multi-disease gut  
256 microbial communities. By using EXPERT on the dataset, we discovered that the human gut  
257 microbial community exhibits disease-specific patterns, which is consistent with previous  
258 cross-disease research<sup>19</sup>. Thirdly, we showed EXPERT's utility in characterizing the gut  
259 microbiota for patients at various stages of CRC. By using communities from five stages of  
260 CRC progression, we found hosts sampled at the same stage shared similar gut microbial  
261 communities, enlightening us to realize that the compositional changes within gut microbial  
262 communities could reflect the progression of CRC, supported by Shaoming Z et al.<sup>22</sup>.

263 Several issues need to be looked into further in the future: We noted that in certain contexts  
264 (e.g., characterizing gut microbial communities during cancer progression), the accuracy  
265 could be improved if the fundamental model was properly selected by referring to the  
266 standard ontology<sup>23,24</sup>. EXPERT should provide a collection of fundamental models to enable  
267 effective adaptation in diverse MST contexts (e.g., environmental source tracking<sup>28</sup>), and  
268 provide an approach for intelligently selecting appropriate fundamental models for a given  
269 context. Additionally, the application of EXPERT on the newly introduced data has indicated  
270 its robustness against the batch effect, while the extent to which Transfer Learning could  
271 overcome the batch effect in microbiome context requires further assessment.

272 In conclusion, EXPERT enabled accurate and rapid source tracking, as well as biologically  
273 informed novel microbial knowledge discovery, by utilizing a Transfer Learning approach.  
274 We have demonstrated the applicability of Transfer Learning in the discovery of microbiome

275 knowledge using this method, particularly when dealing with newly introduced data or  
276 context-dependent settings. We believed that EXPERT could facilitate high-fidelity source  
277 tracking in a broad range of applications.

278

279

## 280 **References**

- 281 1. Thompson, L. R. *et al.* A communal catalogue reveals Earth’s multiscale microbial  
282 diversity. *Nature* **551**, 457–463 (2017).
- 283 2. Turnbaugh, P. J. *et al.* The human microbiome project. *Nature* **449**, 804–10 (2007).
- 284 3. Integrative, H. M. P. R. N. C. The Integrative Human Microbiome Project: dynamic  
285 analysis of microbiome-host omics profiles during periods of human health and disease.  
286 *Cell Host Microbe* **16**, 276–89 (2014).
- 287 4. Ding, R. *et al.* Revisit gut microbiota and its impact on human health and disease. *J*  
288 *Food Drug Anal* **27**, 623–631 (2019).
- 289 5. Metcalf, J. L. *et al.* Microbial community assembly and metabolic function during  
290 mammalian corpse decomposition. *Science* **351**, 158–162 (2016).
- 291 6. Knights, D. *et al.* Bayesian community-wide culture-independent microbial source  
292 tracking. *Nat. Methods* **8**, 761–763 (2011).
- 293 7. R, H. *et al.* Into the deep: Evaluation of SourceTracker for assessment of faecal  
294 contamination of coastal waters. *Water Res* **93**, 242–253 (2016).
- 295 8. Lax, S. *et al.* Longitudinal analysis of microbial interaction between humans and the  
296 indoor environment. *Science* **345**, 1048–1052 (2014).
- 297 9. Shenhav, L. *et al.* FEAST: fast expectation-maximization for microbial source tracking.  
298 *Nat. Methods* **16**, 627 (2019).
- 299 10. Zha, Y. *et al.* Ontology-Aware Deep Learning Enables Ultrafast, Accurate and  
300 Interpretable Source Tracking among Sub-Million Microbial Community Samples from  
301 Hundreds of Niches. Preprint at [https://www.biorxiv.org/content/10.1101/20](https://www.biorxiv.org/content/10.1101/2020.11.01.364208v1)  
302 [20.11.01.364208v1](https://www.biorxiv.org/content/10.1101/2020.11.01.364208v1) (2020)
- 303 11. Smith, A., Sterba-Boatwright, B. & Mott, J. Novel application of a statistical technique,  
304 Random Forests, in a bacterial source tracking study. *Water Res* **44**, 4067–4076 (2010).

- 305 12. Pan, S. J. & Yang, Q. A Survey on Transfer Learning. *IEEE Trans Knowl Data Eng* **22**,  
306 15 (2010).
- 307 13. Ruder, S. An Overview of Multi-Task Learning in Deep Neural Networks. Preprint at  
308 <https://arxiv.org/abs/1706.05098> (2017).
- 309 14. Bäckhed, F. *et al.* Dynamics and stabilization of the human gut microbiome during the  
310 first year of life. *Cell Host Microbe* **17**, 690–703 (2015).
- 311 15. Lozupone, C. & Knight, R. UniFrac: a new phylogenetic method for comparing microbial  
312 communities. *J Appl Environ Microbiol* **71**, 8228–8235 (2005).
- 313 16. Lin, J. Divergence measures based on the Shannon entropy. *IEEE Trans Inf Theory* **37**,  
314 145–151 (1991).
- 315 17. Stokholm, J. *et al.* Delivery mode and gut microbial changes correlate with an increased  
316 risk of childhood asthma. *Sci. Transl. Med* **12**, (2020).
- 317 18. Roswall, J. *et al.* Developmental trajectory of the healthy human gut microbiota during  
318 the first 5 years of life. *Cell Host & Microbe* **29**, 765-776.e3 (2021).
- 319 19. Duvallet, C., Gibbons, S. M., Gurry, T., Irizarry, R. A. & Alm, E. J. Meta-analysis of gut  
320 microbiome studies identifies disease-specific and shared responses. *Nat Commun* **8**,  
321 1784 (2017).
- 322 20. Zhu, Q., Gao, R., Wu, W. & Qin, H. The role of gut microbiota in the pathogenesis of  
323 colorectal cancer. *Tumor Biology* **34**, 1285–1300 (2013).
- 324 21. Zeller, G. *et al.* Potential of fecal microbiota for early-stage detection of colorectal cancer.  
325 *Mol Syst Biol* **10**, (2014).
- 326 22. Zou, S., Fang, L. & Lee, M.-H. Dysbiosis of gut microbiota in promoting the  
327 development of colorectal cancer. *Gastroenterol Rep (Oxf)* **6**, 1–12 (2018).
- 328 23. Buttigieg, P. L. *et al.* The environment ontology in 2016: bridging domains with  
329 increased scope, semantic density, and interoperation. *J Biomed Semantics* **7**, 57 (2016).

- 330 24. Schriml, L. M. *et al.* Human Disease Ontology 2018 update: classification, content and  
331 workflow expansion. *Nucleic Acids Res* **47**, D955–D962 (2019).
- 332 25. Mitchell, A. L. *et al.* MGnify: the microbiome analysis resource in 2020. *Nucleic Acids*  
333 *Res* **48**, D570–D578 (2020).
- 334 26. Mukherjee, S. *et al.* Genomes OnLine Database (GOLD) v.8: overview and updates.  
335 *Nucleic Acids Res* **49**, D723–D733 (2021).
- 336 27. Wu, S. *et al.* GMrepo: a database of curated and consistently annotated human gut  
337 metagenomes. *Nucleic Acids Res* **48**, D545–D553 (2020).
- 338 28. Coordinators, N. R. Database resources of the National Center for Biotechnology  
339 Information. *Nucleic Acids Res* **44**, D7-19 (2016).
- 340 29. Kingma, D. P. & Ba, J. Adam: A Method for Stochastic Optimization. in *3rd*  
341 *International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA,*  
342 *May 7-9, 2015, Conference Track Proceedings* (2015).
- 343 30. Rumelhart, D. E., Hinton, G. E. & Williams, R. J. Learning representations by back-  
344 propagating errors. *Nature* **323**, 533–536 (1986).
- 345 31. Abadi, M. *et al.* TensorFlow: A System for Large-Scale Machine Learning. in  
346 *Proceedings of the 12th USENIX conference on Operating Systems Design and*  
347 *Implementation* 265–283 (2016).
- 348 32. Nair, V. & Hinton, G. E. Rectified Linear Units Improve Restricted Boltzmann Machines.  
349 in *Proceedings of the 27th International Conference on International Conference on*  
350 *Machine Learning* 807-814 (2010).
- 351 33. He, K., Zhang, X., Ren, S. & Sun, J. Delving Deep into Rectifiers: Surpassing Human-  
352 Level Performance on ImageNet Classification. Preprint at [https://arxiv.](https://arxiv.org/abs/1502.01852)  
353 [org/abs/1502.01852](https://arxiv.org/abs/1502.01852) (2015).
- 354 34. Xavier Glorot & Yoshua Bengio. Understanding the difficulty of training deep

355 feedforward neural networks. in *Proceedings of the Thirteenth International Conference*  
356 *on Artificial Intelligence and Statistics* (eds. Yee Whye Teh & Mike Titterington) 249–  
357 256 (2010).

358

359

## 360 **Methods**

### 361 **Datasets**

362 We used six datasets to assess the utility of EXPERT (**Supplementary Table S1**). The  
363 hierarchy is essentially a refined subset of an ontology (e.g., Environmental Ontology<sup>23</sup> or the  
364 Human Disease Ontology<sup>24</sup>) or self-defined according to the context of MST. Refer to  
365 **Supplementary Note 3,4** and **Supplementary Table S11** for the unified data processing  
366 pipeline used in the study.

367 For systematic assessment of our general model, the dataset was obtained from MGnify,  
368 which consists of 118,592 communities collected from 131 biomes. Among them, 52,537  
369 samples originated from human biomes, 14,045 samples originated from mammal biomes,  
370 7,189 samples originated from terrestrial biomes, 27,667 samples originated from aquatic  
371 biomes. These samples were analyzed by MGnify<sup>25</sup> before January 2020 (**Supplementary**  
372 **Table S2**). The source environment hierarchy is constructed by referring to the hierarchical  
373 biome classification from MGnify database<sup>25</sup> and the ecosystem classification paths from  
374 GOLD database<sup>26</sup> (**Supplementary Table S12**).

375 For systematic assessment of our human model, the dataset was a part of the first dataset, in  
376 which 52,537 communities from 27 human biomes were selected (**Supplementary Table S3**).  
377 The source environment hierarchy is constructed by referring to the hierarchical biome  
378 classification from MGnify database<sup>25</sup> and the ecosystem classification paths from GOLD  
379 database<sup>26</sup>.

380 We also used the newly introduced data in 2020 from MGnify<sup>25</sup>. Which consists of 34,209  
381 communities collected from 35 biomes. Throughout the dataset, 3,421 samples belonging to 8  
382 biomes were newly added by MGnify<sup>25</sup> after January 2020 (**Supplementary Table S8**). The  
383 source environment hierarchy is constructed by referring to the hierarchical biome  
384 classification from MGnify database<sup>25</sup> and the ecosystem classification paths from GOLD  
385 database<sup>26</sup>.

386 For source tracking the succession of infant gut microbiome, the dataset was obtained from  
387 MGnify<sup>25</sup> which consists of 392 fecal samples collected from 98 infants and their biological  
388 mothers. Among them, 85 infants were born by vaginal delivery and **13** infants were born by



389 cesarean section. The infant samples were collected at three time points including birth, 4  
390 months, and 12 months. The maternal samples were collected during the first week after  
391 delivery (**Supplementary Table S9**).

392 For disease modeling, the dataset was obtained from GMrepo<sup>27</sup>, including 13,642  
393 communities collected from feces of hosts diagnosed with 19 diseases as well as healthy  
394 controls, **Supplementary Table S4**). The source environment hierarchy is constructed by  
395 referring to NCBI MeSH<sup>28</sup> and Human Disease Ontology<sup>24</sup>.

396 For cancer monitoring, the dataset was obtained from GMrepo<sup>27</sup>, which consists of 16, 93,  
397 126, 196, and 204 communities respectively collected at CRC stage 0, I, II, III, and IV, 635 in  
398 total (**Supplementary Table S10**). The source environment hierarchy is constructed by  
399 referring to the five stages of CRC.

400

## 401 **The EXPERT framework**

### 402 **The EXPERT model**

403 Considering a query sample  $q$  represented by its community structure, as well as its potential  
404 sources represented by a hierarchy  $\mathcal{O}$ , to quantify contributions  $\hat{y}_q$  from the sources to  $q$ , we  
405 employed an adaptive and Multi-task NN to learn a mapping  $M$  from a series of source  
406 samples  $s \in \mathcal{D}_s$  to their biome sources,  $y_s = (y_s^2, \dots, y_s^l)$  (where  $y_s^2$  is biome source for source  
407 sample in the second layer of the biome hierarchy), and then apply  $M$  on  $q$  to determine the  
408 contributions for the query community.

$$409 \quad \hat{y}_q = (\hat{y}_q^i)_{0 < i \leq l_{\mathcal{O}}} = M(q)$$

### 410 **Fast inference via forward propagation**

411 We adopt the rationale of Multi-task Learning<sup>13</sup>. EXPERT integrates the representation of  
412 each lower layer (which is calculated by its “inter” modules  $M_{inter}$ ) into its higher layer, by  
413 employing several “integ” modules  $M_{integ}$ . Therefore, together with “output” module  $M_{output}$ ,  
414 the representation of the contributions is given by

$$415 \quad M(q) = (M_{output}^i(R_{integ}^i))_{0 < i \leq I_O}$$

416 Where

$$417 \quad R_{integ}^i(q) = \begin{cases} M_{integ}^i(M_{inter}^i(M_{base}(q)), 0), & \text{if } i = 1 \\ M_{integ}^i(M_{inter}^i(M_{base}(q)), R_{integ}^{i-1}), & \text{otherwise} \end{cases}$$

418 The NN structures of these modules are described in the subsection Adaptive Neural Network.

419

## 420 Robust optimization via backward propagation and Transfer Learning

421 We adopt the rationale of Transfer Learning<sup>12</sup>. Considering  $\tilde{M}_{base}$  of a fundamental model as  
 422 a static mapping, the parameters of the rest modules  $\hat{w}$  could be solved using gradient descent  
 423 as well as backpropagation algorithm<sup>29-31</sup>:

$$424 \quad \hat{w} = \arg \min_{\hat{w}} \sum_{i=0}^{I_O} \left( \alpha(B_O^i) \sum_{s \in S} \beta_s^i L(\hat{y}_s^i(\hat{w}), y_s^i) \right)$$

425 Where

$$426 \quad \beta_s^i = \begin{cases} 1, & \text{if } y_s^i \text{ exists} \\ 0, & \text{otherwise} \end{cases}$$

$$427 \quad L(\hat{y}_s^i, y_s^i) = \sum_{b \in O^i} (\text{CrossEntropy}(\hat{y}_s^i(b), y_s^i(b)))$$

428  $\alpha(B_O^i) = \frac{B_O^i}{B_O}$  stand for the assigned loss weight for  $i$ -th layer (i.e.,  $(i-1)$ -th task in the  
 429 multiple task).  $\beta_s^i$  stand for the sample weight assigned for a sample  $s$  on  $(i-1)$ -th task during  
 430 learning, enabling the learning from partially labeled data.  $B_O^i$  stand for the number of biomes  
 431 contained in the  $i$ -th layer of the biome hierarchy  $O$ .  $O^i$  stand for the  $i$ -th layer of the biome  
 432 hierarchy  $O$ .  $b \in O^i$  is a biome in the biome hierarchy  $i$ -th layer of the biome hierarchy  $O$ .

433 Then, optimizing the parameters of the entire model (including  $\tilde{M}_{base}$ ), the parameters of the  
 434 entire model  $w$  can be solved by using gradient descent as well as backpropagation  
 435 algorithm<sup>29-31</sup>.

436 
$$w = \arg \min_{\hat{w}} \sum_{i=0}^{I_O} \left( \alpha (B_O^i) \sum_{s \in \mathcal{S}} \beta_n^i L(\hat{y}_s^i(\hat{w}), y_s^i) \right)$$

437 For independent optimization (optimization based on completely random initialization),  
438 EXPERT directly optimizes the entire model. See **Supplementary Note 1** for a detailed  
439 description for optimization.

440

#### 441 **Adaptive Neural Network**

442 NN approach has limited capability when there is a series of newly introduced source  
443 environments, as researchers need to modify the NN model at the code level and re-tune its  
444 hyper-parameters. We developed EXPERT's NN model that changes internal NN structure  
445 according to source environments in different contexts, namely the adaptive NN model  
446 (**Supplementary Fig. S1**). The EXPERT framework initializes the model according to the  
447 hierarchy representing source environments. In the model, there are four conceptual modules.

448 To extract low-level representations for input data, the model employs the “base” module  
449 with two Dense NN layers. The NN layers have fixed structures of 1,024 and 512 neurons,  
450 and use ReLU activation<sup>32</sup> and He initializer with Uniform distribution<sup>33</sup>.

451 To extract representations that are specific to different hierarchy layers, the model employs  
452 the “inter” module with three adaptive Dense NN layers. Denoting  $n$  as the number of source  
453 environments in each hierarchy layer, the three NN layers have adaptive structures of  $8 \times n$ ,  
454  $4 \times n$ , and  $2 \times n$  neurons, respectively. The three NN layers use ReLU activation<sup>32</sup> and He  
455 initializer with Uniform distribution<sup>33</sup>.

456 To integrate representation of different hierarchy layers, the model employs the “integ”  
457 module with a Concatenation NN layer and an adaptive Dense NN layer. Denoting the  
458 number of source environments in each hierarchy layer as  $n$ , the NN layer has adaptive  
459 structures of  $1.5 \times n$  neurons, and uses Tanh activation and Xavier initializer with Uniform  
460 distribution<sup>34</sup>.

461 To estimate according to the integrated representations of different hierarchy layers, the  
 462 model employs the “output” module with an adaptive Dense NN layer. Denoting the number  
 463 of source environments in each hierarchy layer as  $n$ , the NN layer has adaptive structures of  $n$   
 464 neurons, and uses Sigmoid activation and Xavier initializer with Uniform distribution<sup>34</sup>.

465

## 466 Performance measures

467 To assess the performance of EXPERT models and other methods, we used these measures:

$$TP_b(t) = \sum_s I(\hat{y}_s(b) > t \wedge b \in y_s)$$

$$TN_b(t) = \sum_s I(\hat{y}_s(b) < t \wedge b \notin y_s)$$

$$FP_b(t) = \sum_s I(\hat{y}_s(b) > t \wedge b \notin y_s)$$

$$FN_b(t) = \sum_s I(\hat{y}_s(b) < t \wedge b \in y_s)$$

$$TPR_b(t) = \frac{TP_b(t)}{TP_b(t) + FN_b(t)}$$

$$FPR_b(t) = \frac{FP_b(t)}{FP_b(t) + TN_b(t)}$$

$$Recall_b(t) = \frac{TP_b(t)}{TP_b(t) + FN_b(t)}$$

$$Precision_b(t) = \frac{TP_b(t)}{TP_b(t) + FP_b(t)}$$

468

469 Where  $TP$  is true positive,  $TN$  is true negative,  $FP$  is false positive,  $FN$  is false negative,  
 470  $y_s(b)$  is the quantified contribution from a biome source  $b$  for a microbial community sample  $s$ ,  
 471 threshold  $t \in [0, 1]$  with a step size of 0.01,  $y_s$  is a set of actual biomes for a sample  $s$ , and  $I$  is  
 472 a logical operation function, the value of  $I$  is 1 when the result of logical operation is TRUE,  
 473 else 0.

474 Then, two evaluation metrics (F-max, AUROC) were introduced. F-max stands for the  
 475 maximal F1-measure and was calculated with the following formula. AUROC stands for the  
 476 area under the ROC (Receiver Operating Characteristics) and was calculated using the  
 477 composite trapezoidal rule.

$$F_{\max_b}(t) = \max_t \frac{2Precision_b(t)Recall_b(t)}{Precision_b(t) + Recall_b(t)}$$

478  
479 Finally, we treated the average performance across all biomes as the performance of the  
480 entire model. Notably, in the subsection “Efficiency, accuracy, and adaptivity of EXPERT”,  
481 we only considered biomes with the number of samples > 100 to compute the average  
482 performance for the general model, the independent model, Transfer (GM) model, and  
483 Transfer (GM0) model.

484

## 485 **Evaluating fundamental models**

486 We assessed each model of the fundamental models through cross-validation, and selected  
487 the best model among all trained models as the final model.

488 We assessed the general model by applying eight-fold cross-validation to the 125,823  
489 microbial community data collected from 132 biomes, and selected the best model among  
490 eight trained models as the general model to be transferred.

491 We assessed the human model by applying repetitive cross-validation (90% as sources to  
492 train a model, the resting 10% as queries to test its performance, repeated for five times) to  
493 the 52,537 microbial community data collected from 25 biomes, and selected the best model  
494 among five trained models as the general model to be transferred.

495 The assessment of the disease model is the same as the assessment of the human model, but  
496 using another dataset consists of 13,462 gut microbial communities associated with 19  
497 diseases.

498

## 499 **Experiment design**

500 We compared EXPERT’s performance with FEAST and the NN approach using the human-  
501 associated dataset (**Supplementary Table S1, S3**). We measured the running time using the  
502 Linux command “time” and considered the real-time usage for comparison. The efficiency

503 was then calculated using the running time we measured. Refer to **Supplementary Note 5** for  
504 detailed comparison procedure for each experiment.

505 We demonstrated EXPERT's utility in context-aware MST in three contexts. In these  
506 contexts, we used standard hyperparameters for training the model (**Supplementary Note 1**).  
507 Detailed descriptions are provided in **Supplementary Note 6**.

508

## 509 **Statistical analysis**

510 Statistical analyses of the contributions have been performed utilizing the Wilcoxon test, at  
511 the significance level of  $\alpha=0.05$ . For all the tests, when the p-value associated is lower than the  
512 significance level, one should reject the null hypothesis  $H_0$ , and accept the alternative  
513 hypothesis  $H_a$ .

## 514 **Visualization of data distribution**

515 Throughout the paper, the box-plot elements are centerline, median; box limits, upper and  
516 lower quartiles; whiskers,  $1.5 \times$  interquartile range (IQR); points, and outliers. The Violin  
517 plot is also used for data distribution analysis, mainly for comparison. The PCoA is also used  
518 for data distribution analysis, with ellipses representing a confidential interval of 0.95. The  
519 Principle Coordination is obtained through applying beta diversity measurement (Scikit-bio  
520 version 0.5.6, **Supplementary Table S5**) on the abundance of all taxa in seven ranks, namely  
521 Superkingdom, Kingdom, Phylum, Class, Order, Family, Genus, and Species. The source  
522 code of the PCoA analysis is hosted on GitHub at <https://github.com/AdeBC/UniPCoA>.

523

## 524 **Data availability**

525 The collected samples from MGnify and GMrepo databases were annotated with their  
526 associated biomes/phenotypes in **Supplementary Table S2-S4, S8-S10**. All the processed  
527 data are uploaded and hosted at <https://github.com/HUST-NingKang-Lab/EXPERT-use-cases>.

528

529 **Code availability**

530 All source codes have been uploaded to the website at [https://github.com/HUST-NingKang-](https://github.com/HUST-NingKang-Lab/EXPERT)  
531 Lab/EXPERT. Detailed software and models used in this study are provided in  
532 **Supplementary Table S5.**

533

534 **Acknowledgments**

535 This work was partially supported by National Natural Science Foundation of China grant  
536 32071465, 31871334, and 31671374, and the China Ministry of Science and Technology's  
537 National Key R&D Program grant (No. 2018YFC0910502).

538

539 **Author contributions**

540 HC and KN designed the study, conceived of and proposed the idea, designed and developed  
541 the framework. HC, QY, GX, NW, SH and XH performed the experiments and analyzed the  
542 data. HC, QY, YZ, GX, and NW visualized the data. CS and SW provided valuable data. HC,  
543 QY, YZ, WC, LPC, and KN contributed to editing and proofreading the manuscript. All  
544 authors read and approved the final manuscript.

545

546 **Competing interests**

547 The authors declare that they have no competing interests.

548

549 **Ethics approval and consent to participate**

550 Not applicable.

551

552

## 553 **Figures Legends**

554 **Figure 1: Illustration of EXPERT's knowledge transfer process. a.** EXPERT can adapt  
555 the knowledge of a fundamental model to an MST context through three steps: transfer (reuse  
556 parameters of a fundamental model and reinitialize contextual layers according to the context,  
557 red dotted arrows), adaptation (quickly optimize only the contextual layers using iterative  
558 forward-backward propagation, green circular arrows), and fine-tuning (further optimize the  
559 entire model using the iterative forward-backward propagation). The fundamental model is a  
560 pre-trained EXPERT model to be adapted, with several NN layers relatively independent to  
561 contexts and a series of contextual NN layers highly specified to a context). Different  
562 background colors of the model indicate the suitability of different modules to the context.  
563 The contextualized model can serve a broad-spectrum of source tracking applications (based  
564 on research purposes, illustrated in **Fig. 1b**). Abbreviations: MST: microbial source tracking;  
565 NN: Neural Network.

566

567 **Figure 2. Efficiency, accuracy, and adaptivity of EXPERT. a.** Comparison of Transfer  
568 (GM) EXPERT model with FEAST on efficiency (number of queries/sinks per second, left  
569 Y-axis) and accuracy (based on cross-validation, right Y-axis). For FEAST, the sources were  
570 randomly selected 70, 140, and 210 samples (10, 20, and 30 samples per biome, respectively).  
571 EXPERT's performance was measured by contextualizing the general model. **b.** The  
572 performances (validation F-max, Y-axis) of three models along with different proportions of  
573 sources used (X-axis). The NN model was trained solely based on contextual data. The  
574 results were obtained by using cross-validation and different proportions (1-10% by a step  
575 size of 1%, and 10-90% by a step size of 10%) of source samples. Loess regression was  
576 applied to these points using the number of source samples used and F-max.

577

578 **Figure 3. Robust adaptation to the newly introduced microbiome data. a.** Partial  
579 representation of the baseline data and the newly introduced data (with sample size  
580 annotations) used to measure the impact of batch effects on MST models and assess the  
581 utility of EXPERT. The baseline data contains 118,592 communities deposited before



582 January 2020. The newly introduced data contains 34,209 communities deposited between  
583 January 2020 and November 2020, including several newly introduced biomes (e.g. fish-  
584 associated biomes). **b.** Performance of EXPERT models on the baseline data and the newly  
585 introduced data (performance for seven representative biomes). Furthermore, we can also  
586 adapt a fundamental model to newly introduced sources to evaluate these potential microbial  
587 sources. Abbreviations: “\*\*\*”: significant difference; “NS”: non-significant difference; GM:  
588 the general model; Transfer (GM): the contextualized model based on the general model.  
589 Representative biomes: biomes in the fourth layer of the MGnify biome hierarchy and with  
590 sample size greater than 100 in both two datasets.

591

592 **Figure 4. EXPERT’s performance in characterizing gut microbial community**  
593 **development over time for infants. a.** The hierarchy representing source environments,  
594 corresponding to infant samples collected from the ENA database. Environments in the  
595 second and third layers were grouped by sampling time and delivery modes. For this part of  
596 the study, sources include the gut microbiome of the mother, infant at birth, and four months,  
597 queries include the gut microbiome of the infant at 12 months. **b.** Estimated contributions by  
598 Transfer (HM) model, separated by two delivery modes. **c.** Distribution of infant gut  
599 microbial communities during their first year, using principal-coordinates analysis (PCoA)  
600 and distance metric of Jensen Shannon divergence. The dotted line refers to samples  
601 delivered vaginally, and the full line refers to samples delivered via cesarean section. The  
602 baby of 4 months is abbreviated to baby 4M, the baby of 12 months is abbreviated to baby  
603 12M. The letters “C” and “V” stand for cesarean section and vaginal delivery, respectively.  
604 Top panel: samples from the infant’s gut are plotted according to their source and collection  
605 date on the Y-axis, and position on the X-axis is plotted according to their first principal  
606 coordinate in the PCoA. **d.** The overall performance of models generated based on different  
607 fundamental models, in which the Independent model was solely based on the samples used  
608 in this context; Transfer (GM) and Transfer (HM) refer to models built based on the general  
609 model and human model with fine-tuning, respectively.

610

611 **Figure 5. EXPERT reveals disease-specific patterns within gut microbial communities. a.**

612 Illustration of knowledge transfer utilized for disease pattern analysis. The knowledge  
613 transfer between MST contexts was illustrated using different colors (white for human-  
614 associated biomes, yellow for gut microbiota-associated disease status). In this analysis, the  
615 knowledge from the human model was contextualized (transferred) to the dataset containing  
616 13,642 samples and 19 diseases as well as healthy control. **b.** The hierarchical organization of  
617 19 diseases and healthy control. The hierarchy was constructed by referring disease names to  
618 Medical Subject Headings and Human Disease Ontology. The hierarchy includes 20 different  
619 health statuses (19 different diseases and infections, plus healthy control) distributed in seven  
620 different layers (X-axis). **c.** Average source contribution among all diseases and healthy  
621 control, obtained by quantifying contribution from 90% samples (randomly selected) of the  
622 dataset to the remaining 10% samples, using Transfer (HM) model based on the human  
623 model. The process of random selection and quantification was repeated five times. The  
624 heatmap was obtained by averaging the contributions to all samples from each one out of 19  
625 diseases and healthy control. There is no sample overlap between source samples and query  
626 samples. **d.** The performance of the EXPERT models on the gut microbial community  
627 associated with each disease or healthy control, evaluated based on the source contribution  
628 (same as in **Fig. 5c**) and biome-specific evaluation (**Methods**). The dashed line indicates an  
629 AUROC of 0.800. **e.** The overall performances of the Transfer (HM) model. Settings of  
630 quantification and assessment were the same as **Fig. 5d**.

631

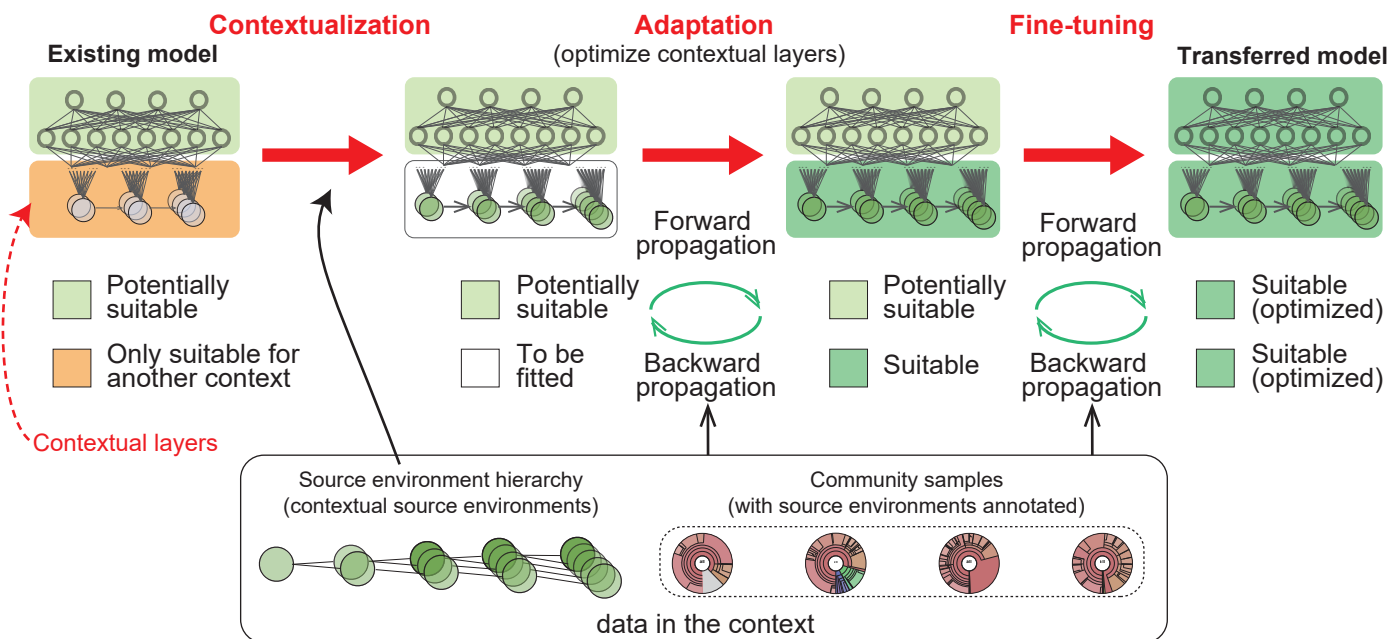
632 **Figure 6. EXPERT characterizes compositional shifts within host gut microbiota during**

633 **CRC progression. a.** Illustration of knowledge transfer utilized for characterizing the  
634 compositional shifts. The knowledge from the human model learned from 52,537 human-  
635 associated communities, as well as the disease model learned from 13,642 human gut  
636 communities associated with 19 diseases and healthy control, were transferred to characterize  
637 the CRC-related compositional shifts. **b.** The five stages of CRC progression, and the number  
638 of samples for each stage. Stage 0: healthy control. **c.** The distribution of gut microbiomes,  
639 visualized by PCoA (using distance metric of weighted-UniFrac). **d.** The average  
640 contribution of different stages of CRC. The source samples were randomly selected 90% out

641 of the entire dataset. The query samples were the remaining 10% samples. This process of  
642 random selection and quantification was repeated five times. There is no sample overlap  
643 between source samples and query samples. **e.** The stage-specific performances (AUROC) of  
644 EXPERT on different CRC stages (see **Methods** for details of stage-specific evaluation).

**a.**

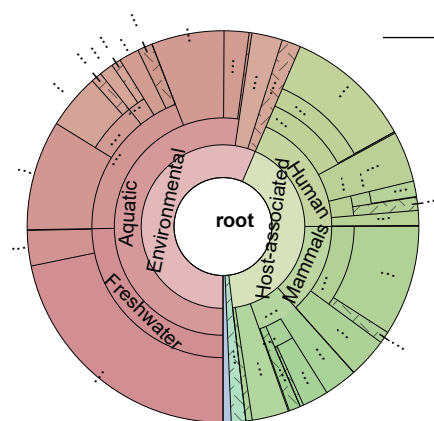
### Transfer process of EXPERT



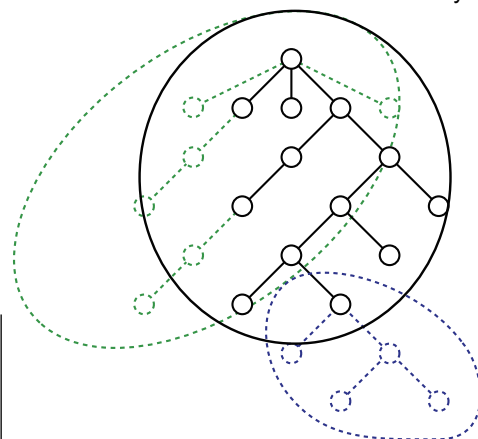
**b.**

### Multi-faceted applications of EXPERT

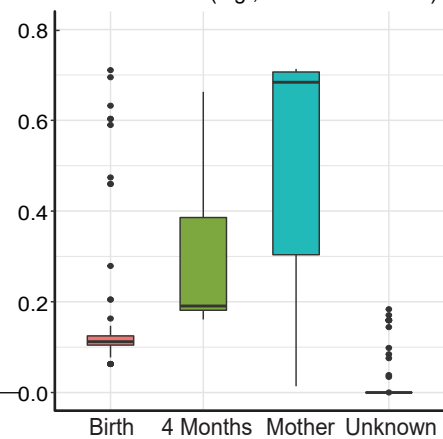
i. Direct estimation of source contribution



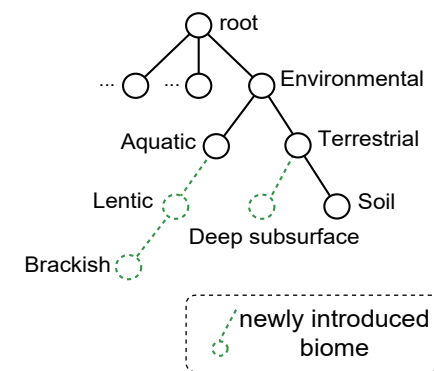
Source environment hierarchy



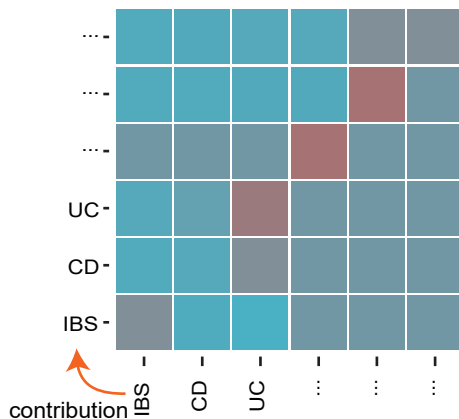
iii. Characterizing small difference among communities (e.g., infant communities)



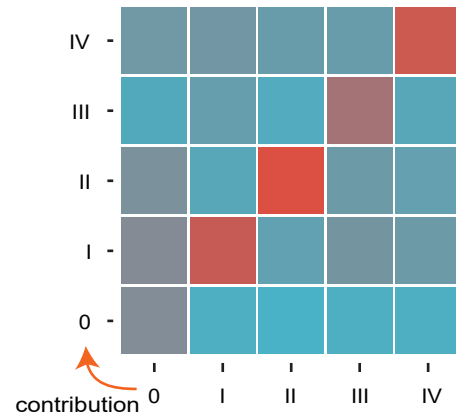
ii. Adapting to newly introduced data against batch effects



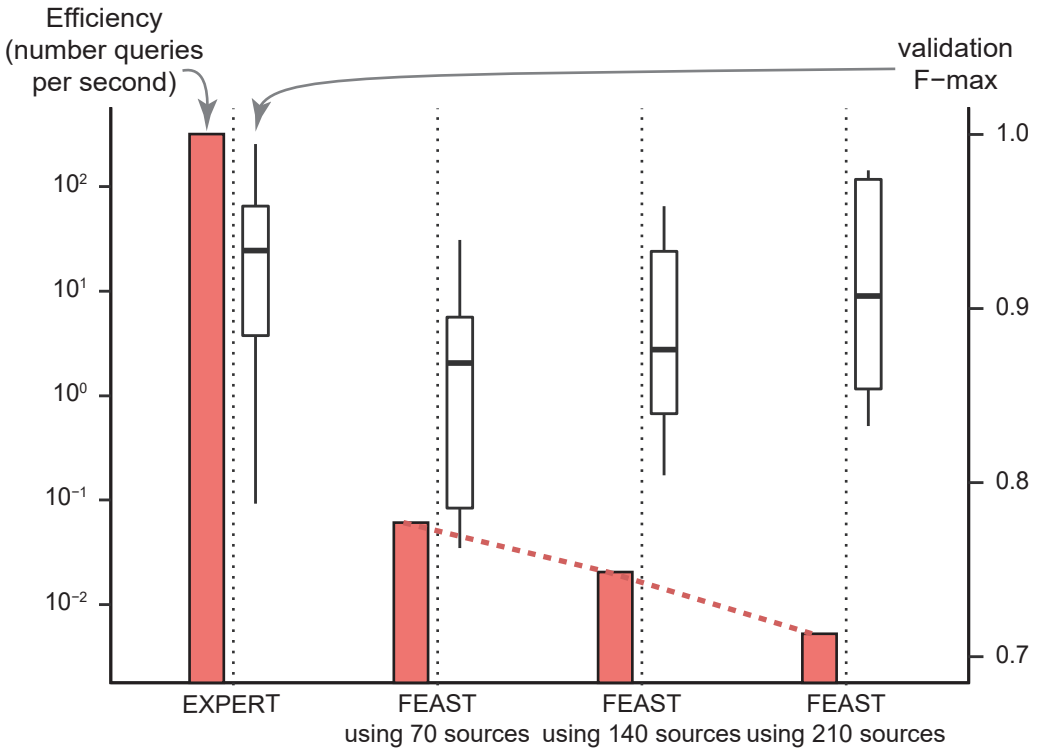
iv. Quantifying similarity of gut microbial communities related with multiple diseases



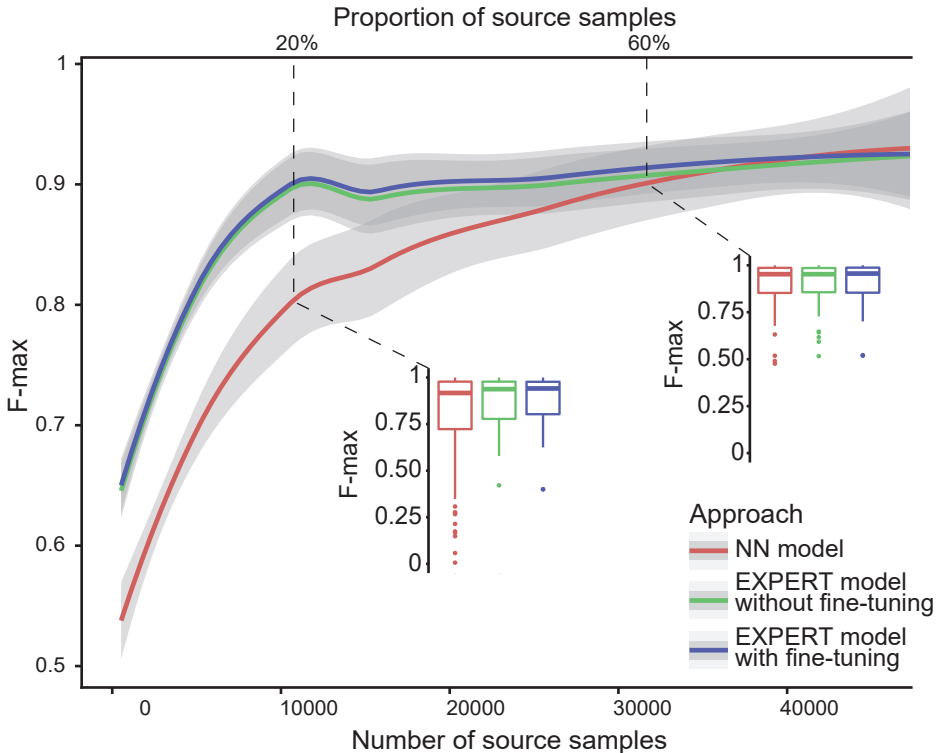
v. Characterizing similarity of gut microbial communities related with CRC progression



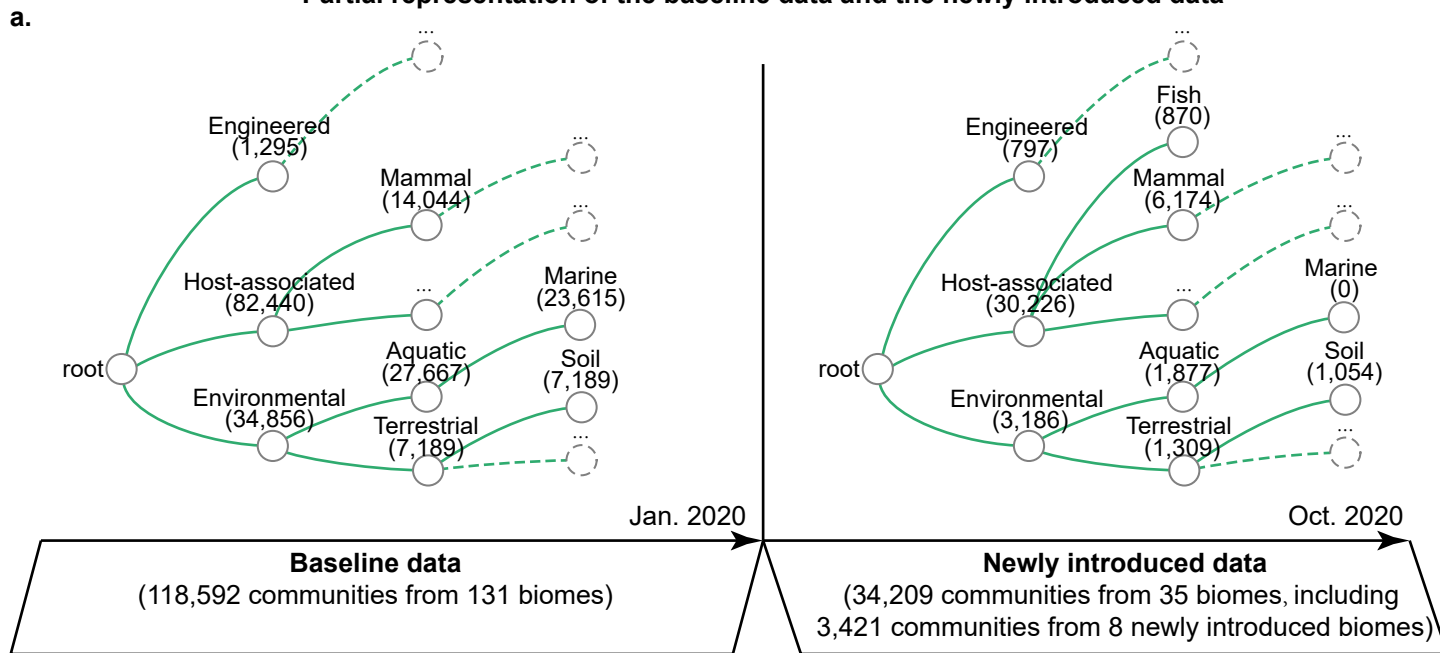
### a. Comparison of EXPERT and FEAST



### b. Prediction results of Transfer and Independent model

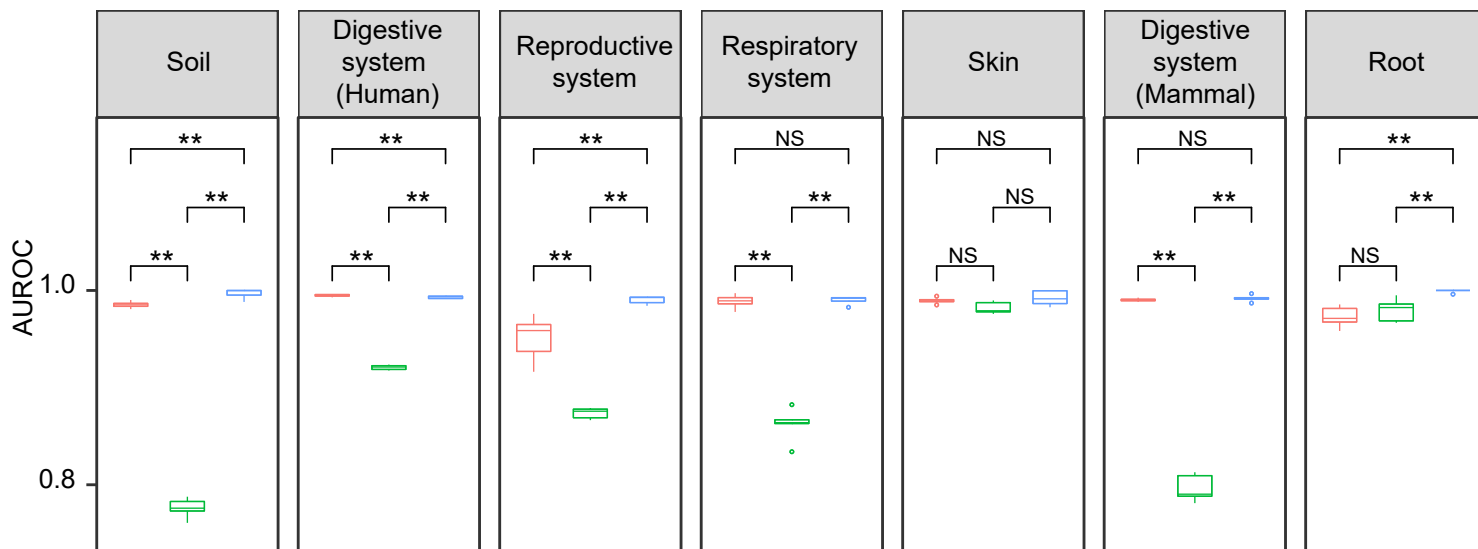


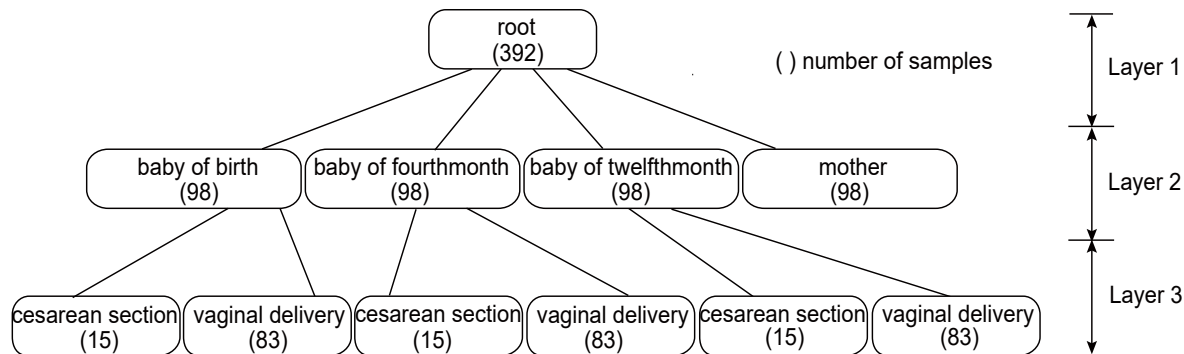
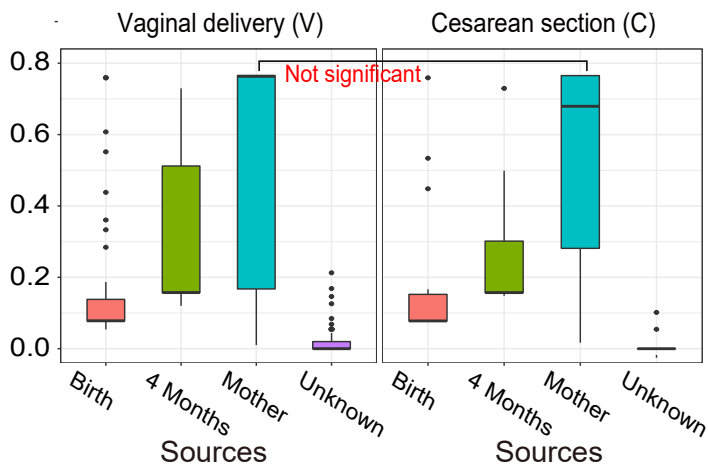
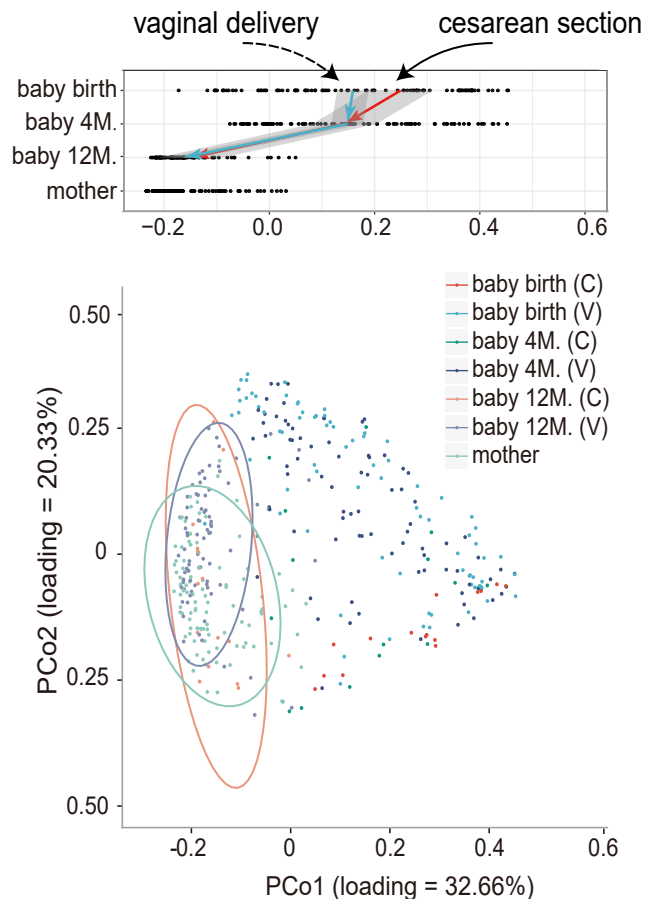
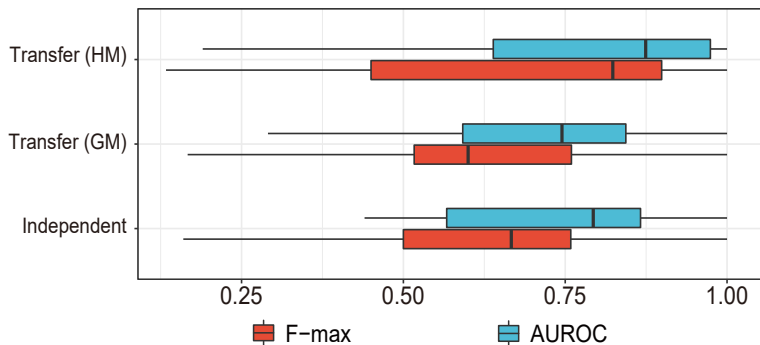
### Partial representation of the baseline data and the newly introduced data



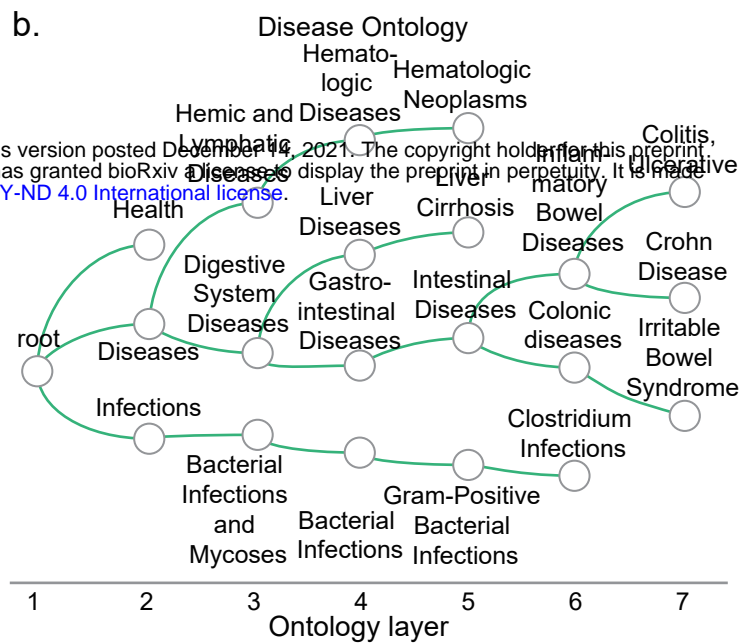
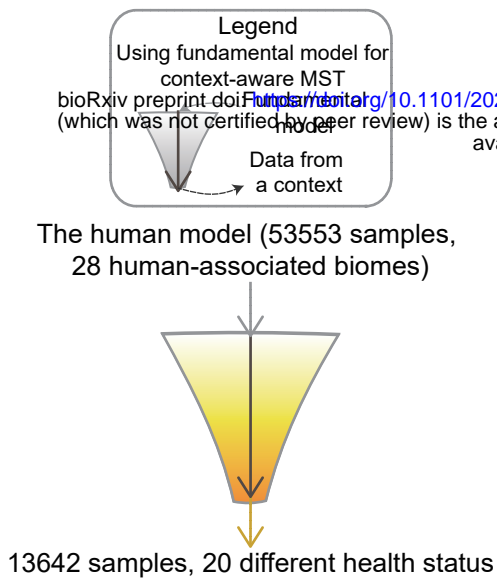
### **b.** Accuracy of the general model and Transfer (GM) for representative biomes

GM on baseline data    GM on emerged data    Transfer (GM) on emerged data

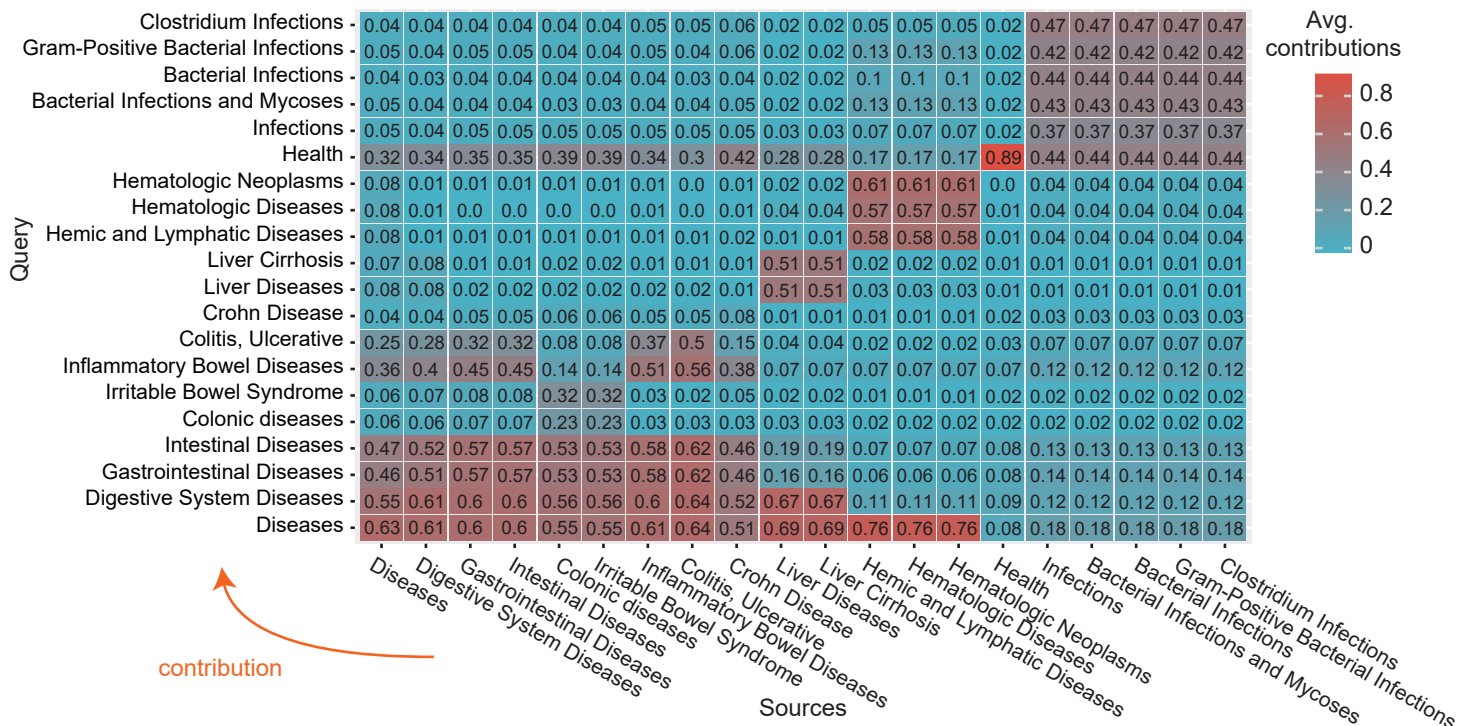


**a. Biome ontology(layer 2: sampling time)****b. Estimated source proportions by different models (queries: samples of 12 months)****c. Distribution of infant gut microbiome during first year of age****d. Overall performances**

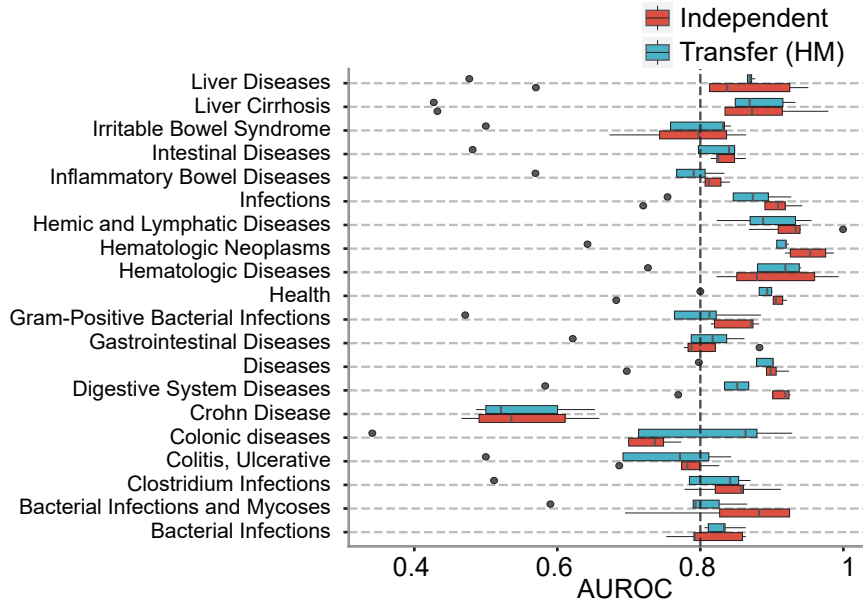
a. Transfer process



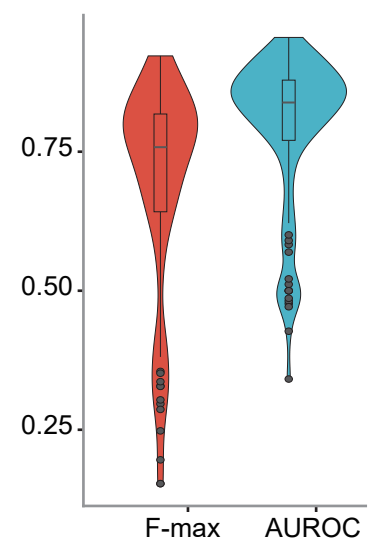
c. Source contribution among host gut microbial communities with different phenotypes



d. Distinguishable patterns of diseases

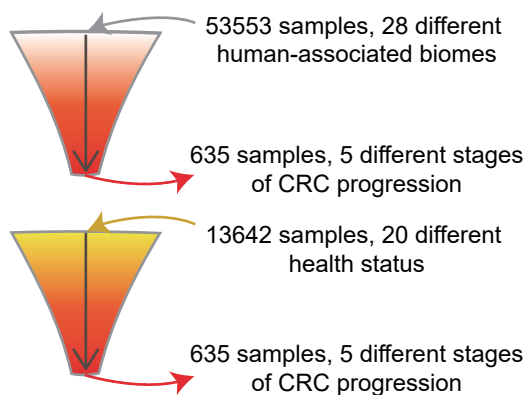


e. Overall performances on differentiating diseases  
Experiment: Transfer (HM)

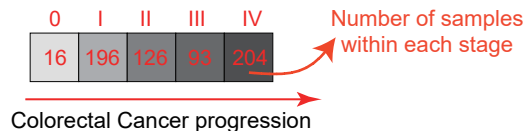




### a. Transfer process

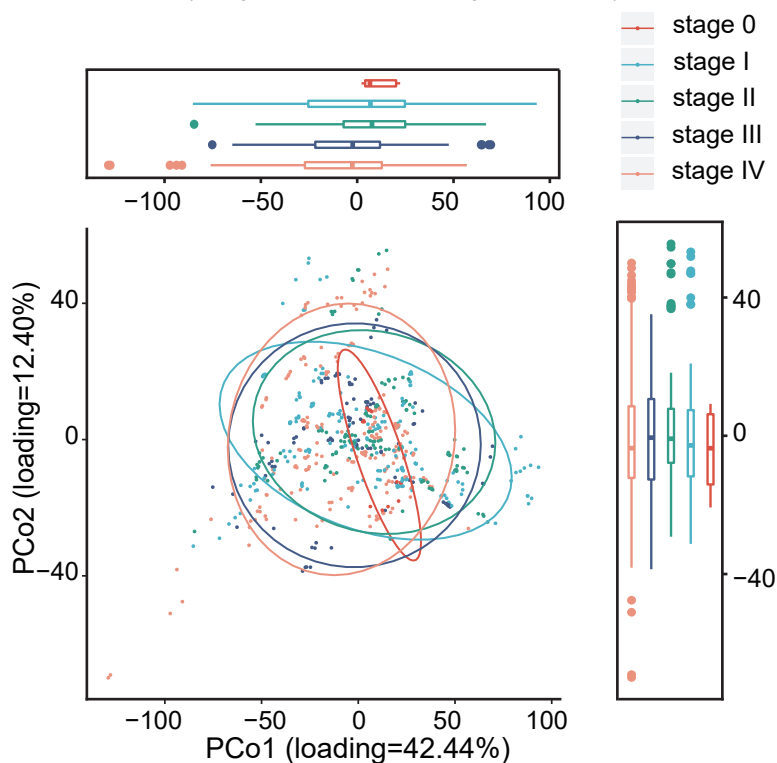


### b. The stage of CRC



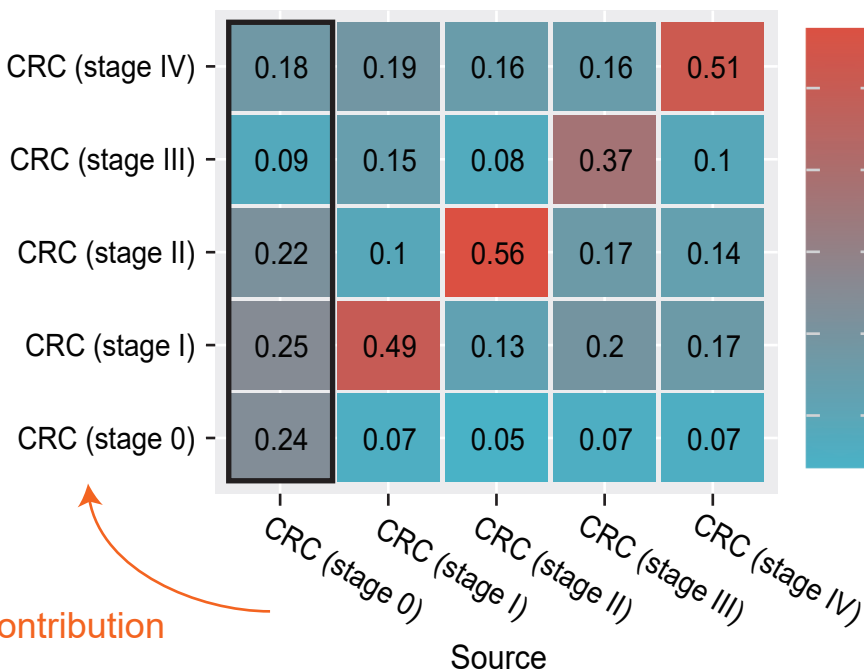
### c. Distribution of gut microbiome across CRC stages

(using distance metric: weighted-Unifrac)



### d. Source contribution within human gut microbiome

Experiment: Transfer (DM)



### e. Differentiating stages for CRC

