

1 **Enabling technology for microbial source tracking based on transfer lea**
2 **rning: From ontology-aware general knowledge to context-**
3 **aware expert systems**

4

5 Hui Chong¹, Qingyang Yu¹, Yuguo Zha¹, Guangzhou Xiong¹, Nan Wang¹, Chuqing Sun¹,
6 Sicheng Wu¹, Weihua Chen¹, Kang Ning^{1,*}

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

11 *Correspondence should be addressed to K.N (Email: ningkang@hust.edu.cn)

12

13 **Abstract**

14 Habitat specific patterns reflected by microbial communities, as well as complex
15 interactions between the community and their environments or hosts' characteristics, have
16 created obstacles for microbial source tracking: diverse and context-dependent
17 applications are asking for quantification of the contributions of different niches (biomes),
18 which have already overwhelmed existing methods. Moreover, existing source tracking
19 methods could not extend well for source tracking samples from understudied biomes, as
20 well as samples from longitudinal studies.

21 Here, we introduce EXPERT (<https://github.com/HUST-NingKang-Lab/EXPERT>), an
22 exact and pervasive expert model for source tracking microbial communities based on
23 transfer learning. Built upon the biome ontology information and transfer learning
24 techniques, EXPERT has acquired the context-aware flexibility and could easily expand
25 the supervised model's search scope to include the context-dependent community
26 samples and understudied biomes. While at the same time, it is superior to current

27 approaches in source tracking accuracy and speed. EXPERT's superiority has been
28 demonstrated on multiple source tracking tasks, including source tracking samples
29 collected at different disease stages and longitudinal samples. For example, when dealing
30 with 635 samples from a recent study of colorectal cancer, EXPERT could achieve an
31 AUROC of 0.977 when predicting the host's phenotypical status. In summary, EXPERT
32 has unleashed the potential of model-based source tracking approaches, enabling source
33 tracking in versatile context-dependent settings, accomplishing pervasive and in-depth
34 knowledge discovery from microbiome.

35

36 **Introduction**

37 Advances in sequencing technology and informatics are producing an exponential
38 increase in data acquisition and integration, and revolutionizing our understanding of the
39 roles microbes play in health and disease, biogeochemical cycling, etc.¹⁻⁴. Hundreds of
40 thousands of microbial community samples have been accumulated yearly, corresponding
41 to hundreds of niches (biomes) around the globe⁵⁻⁷, and continuously completing the
42 grand picture about the microbiome world. However, currently annotated biomes only
43 represented the pick of an iceberg of such a grand picture, while a considerable amount of
44 these samples are from understudied and newly discovered biomes, including (1) newly
45 discovered biomes from where microbial community samples were only collected
46 recently, (2) the biomes which were more context-dependent such as those representing
47 the development of gut microbial communities, or those representing different stages of
48 diseases, and (3) the biomes which include the longitudinal samples. These understudied
49 biomes represent a grand pool of unexplored and previously less studied microbiome
50 research sphere, have created huddles for sample comparison and source tracking⁸. The
51 microbial community dark matters, referring to those samples and niches that have been
52 unseen or understudied, have grown exponentially, which urged the context-aware
53 method to source track the biomes. However, there is no context-aware method for such
54 purposes, as such context-dependent investigations requiring both comprehensiveness
55 (serving for characterizing (1)) and scalability (serving for characterizing both (2) and
56 (3)).

57 Previous approaches for community-wide microbial source tracking (MST) have serious
58 tradeoffs regarding efficiency, accuracy and scalability. Markov Chain Monte Carlo
59 (MCMC)⁹ and Expectation-Maximization (EM)¹⁰ methods have a tradeoff between
60 accuracy and efficiency when facing an increasing number of possible sources (referred
61 to **Supplementary Note 1** for a detailed discussion). They can hardly serve for source
62 tracking among thousands of samples, which often requires weeks to years to complete
63 (linearly extrapolated based on results in L. Shenhav et al.¹⁰). Whereas for Neural
64 Network (NN) method¹¹, the tradeoff lies in scalability: It can hardly be expanded to
65 context-dependent applications and longitudinal analyses, limiting its utility in current
66 microbiome studies.

67 We have developed EXPERT, an exact and pervasive expert model for source tracking
68 microbial community samples to address these limitations. EXPERT combines the
69 superior efficiency and accuracy of the Neural Network method, as well as the transfer
70 learning approach's inherent scalability, enabling knowledge transfer into context-
71 dependent settings to better understand the microbial community dark matters. EXPERT
72 is designed to quantitatively assign the contribution of biomes to a specific microbial
73 community, in both ontology-aware and context-aware manners. EXPERT's advantages
74 include its ability to rapidly infer the contributions from multiple sources via ontology-
75 aware forward propagation, and its ability to adapt to emerging research via alterable
76 biome ontology structure (i.e., a hierarchy of biomes involved in the target application).
77 EXPERT has made it possible to look into the dark matter of microbial communities that
78 include millions of samples from hundreds of context-dependent biomes. Thus it has
79 enabled fast, accurate, flexible, and interpretable source tracking at an unprecedented
80 scale towards deeper knowledge discovery from microbial communities.

81 EXPERT has demonstrated superior performance on a diverse set of source tracking tasks:
82 For source tracking among newly discovered biomes, it has been shown to be able to
83 adapt to emerging samples and biomes with unprecedented AUROC higher than 0.990.
84 For the more context-dependent biomes such as those representing the progression of
85 colorectal cancer, it has achieved an AUROC of 0.977 when predicting the host's
86 phenotypical status, exhibiting its potential in disease diagnosis. It has also enabled the

87 longitudinal sample source tracking, revealed the compositional shifts of individuals' gut
88 microbial communities along the timeline. In summary, EXPERT has enabled the
89 flexibility of model-based method for context-aware source tracking, expanded the scope
90 of source tracking to understudied biomes, and boosted in-depth knowledge discovery
91 from the microbial communities.

92

93 **Results**

94 **Rationale, modeling and multi-facet applications of EXPERT**

95 EXPERT is an efficient approach for microbial community source tracking employing
96 both neural network modeling and transfer learning techniques¹², enabling knowledge
97 transfer from ontology-aware general knowledge to context-aware expert systems. The
98 data preprocessing pipeline, the limitation of directly use of ontology-aware general
99 knowledge, as well as knowledge transferring procedure of EXPERT are illustrated in **Fig.**
100 **1a-c**. Firstly, EXPERT is agnostic to the sequencing data type (16S ribosomal RNA or
101 shotgun sequencing) or the analysis pipeline. And before source tracking, EXPERT
102 applies remapping, normalizing, aggregation (**Supplementary Note 2**), and Z-score
103 standardization for the standardization of microbial community data (**Fig. 1a**).

104 Secondly, one of the fundamental models in EXPERT before transfer learning is a neural
105 network model, which is either a general ontology-aware neural network (ONN) model
106 ONN4MST¹¹ (<https://github.com/HUST-NingKang-Lab/ONN4MST>) or any other neural
107 network built for the same purpose. This neural network model is also referred to as the
108 ontology-aware general knowledge. The model-based source tracking follows this
109 rationale: For a community sample as query, by assuming that quantifying contributions
110 from a set of independent biomes will also contribute to quantifying contributions from
111 the biomes serving as the superset of these independent biomes (e.g., "Human" is the
112 superset biome of "Human: Oral" and "Human: Fecal"), EXPERT quantifies the
113 contributions of biomes along the ontology layers via a propagation procedure, where
114 messages of higher layers are integrated into those of lower layers, namely ontology-
115 aware forward propagation. The ontology-aware general knowledge modeling has

116 enabled EXPERT to efficiently source track among up to hundreds of sources biomes
117 (containing sub-millions of source samples). However, the ontology-aware general
118 knowledge comes with poor scalability and can hardly be expanded to context-dependent
119 applications (**Fig. 1b**).

120 Thirdly, EXPERT can adapt the general knowledge of an existing model to context-
121 dependent applications through three steps, namely transfer, adaptation, and finetuning
122 (**Fig. 1c**). It has achieved this by utilizing both ontology-aware general knowledge, and
123 context-dependent biome ontology, which is essentially a hierarchy of biomes involved in
124 the target application. During the transfer process, EXPERT reuses the existing model's
125 parameters to optimize a context-aware model and encode context-dependent biome
126 ontology into the model through reinitializing context-dependent layers (containing only
127 5% parameters of the entire model) according to the knowledge of the application-
128 dependent context. The phylogenetic tree from the existing model was also reused
129 (**Supplementary Table S1**). During the adaptation process, EXPERT quickly optimizes
130 only context-dependent layers to enable the model to become suitable for the context.
131 During the finetuning process, EXPERT further optimizes the entire model to thoroughly
132 adapt it to the application. Note that EXPERT can also learn from partially-labeled data
133 by masking the losses of unlabeled layers (**Supplementary Fig. 1, Supplementary Note**
134 **3**).

135 The multi-facets of EXPERT applications are illustrated in **Fig. 1c**. First, we can quantify
136 the contributions of biomes for any microbial community by using EXPERT with
137 superior speed and accuracy (**Fig. 1c i**). Secondly, EXPERT is able to adapt to source
138 tracking among merging microbiome data (also referred to as the "Grafting", **Fig. 1c ii**).
139 Thirdly, we can introduce more detailed biomes in order to investigate the small
140 differences among these closely related biomes (i.e., the biomes representing the infant of
141 different ages, **Fig. 1c iii**), which is also referred to as the "Scattering". Fourthly, by
142 introducing knowledge about diseases (disease ontology), we can build an EXPERT
143 model for characterizing the health status of hosts. Alternatively, by introducing
144 knowledge about disease progression, we can build a model for monitoring disease
145 progression (**Fig. 1c iv**). Finally, we can leverage time points and additional context-

146 aware knowledge to investigate dynamics of microbial communities along the timeline,
147 such as longitudinal dynamics of communities in accordance with dietary shifts or
148 seasonal shifts (**Fig. 1c v**). All of these context-aware applications were also implemented
149 and assessed in details in the following parts of this work.

150

151 **Transferring general knowledge into context-dependent applications: systematic** 152 **assessment**

153 We verified the utility of our transfer scheme by assessing the performance of transfer
154 learning under two application scenarios: quantification of contributions under general
155 application scenario, as well as quantifying human-associated source contributions under
156 context-dependent application scenario. We introduced two datasets to assess the
157 performances systematically: a combined dataset consists of 118,592 samples collected
158 from 132 biomes (**Supplementary Table S2, Supplementary Table S3**), and a human
159 dataset consists of 52,538 samples collected from 28 human-associated biomes
160 (**Supplementary Table S2, Supplementary Table S4**). The general EXPERT model
161 (also referred to as the general model) is generated based on using samples from the
162 combined dataset, and we have first assessed the performance of the general model under
163 general application scenario, namely on source tracking samples in the combined dataset.
164 Through random cross-validation (see **Methods**), we found that the general model was
165 able to quantify source contributions for communities with high AUROC of 0.971
166 (**Supplementary Fig. 2**).

167 We then evaluated the transfer learning approach in a context-dependent application. We
168 examined the application of knowledge transfer and finetuning (refer to **Methods** for
169 details) on quantifying human-associated source contributions under context-dependent
170 application scenario, in three aspects: efficiency, accuracy, . In this context, the transfer
171 models built based on transferring knowledge from the general model with and without
172 finetune were referred to as Transfer (GM) and Transfer (GM0), respectively
173 (**Supplementary Table S5**). These transfer models were compared with independent
174 model, which was generated by general ontology-aware neural network approach¹¹ based
175 on the samples and biomes in the context-dependent application (**Supplementary Table**

176 **S5**). The results have shown that the knowledge transfer scheme with finetuning enabled
177 more accurate quantification of source contributions for query samples (average AUROC
178 of 0.960, **Fig. 2c**) compared to the independent model. We have also found that 95.7% of
179 the parameters from the original general model has been transferred to Transfer (GM),
180 representing 99.3% parameters of Transfer (GM), confirming that the high source
181 tracking accuracy actually come from the transfer learning process (**Supplementary**
182 **Table S5**). Notably, the finetune optimization process comes up with a cost: Three times
183 as much time was spent on performing this optimization (**Fig. 2d**), primarily due to the
184 low learning rate (1×10^{-5}) utilized by finetuning. Nevertheless, considering the
185 advantages in accuracy (**Fig. 2c**), we determined finetuning as a default setting in the
186 following sections.

187 We further compared the performance of EXPERT with FEAST¹⁰ in a flat setting, where
188 the potential sources (only the bottom layer of the human ontology) are considered
189 independent of each other (**Fig. 2a**), as FEAST cannot recognize the hierarchical
190 relationships among biomes (**Supplementary Note 4**). Results have shown that EXPERT
191 could perform source tracking with excellent performance (F-max of 0.914) and ultrahigh
192 search speed (over 200 samples per second). However, though FEAST has been proven
193 to be much faster than SourceTracker⁹ (**Fig. 2e, Supplementary Note 1**), we have
194 noticed a severe tradeoff in FEAST between accuracy and running time, which is heavily
195 depend on the number of sources: when the numbers of samples as sources for FEAST
196 increase, the accuracies could reach those similar with EXPERT, but at the cost of
197 magnitude more time (**Fig. 2e, Supplementary Table S6**).

198

199 **EXPERT enables adaptation for emerging microbiome data**

200 In this context, we aim to adapt the general model to emerging microbiome data from
201 understudied or newly discovered biomes, which is also referred to as "Grafting". An
202 increasing number of biomes are being studied in metagenomics, accompanied by tens of
203 thousands of communities deposited into public databases. Many of these biomes are
204 understudies, and a profound number of biomes are newly discovered, such as root: Host-
205 associated: Fish and root: Host-associated: Birds, just in year 2020 (**Fig. 3a**). To assess

206 the capability of EXPERT on such emerging microbiome data, we collected 34,209
207 samples, which were collected and analyzed by MGnify⁶ in 2020. Among them, there are
208 3,419 samples originated from 8 newly discovered biomes (**Supplementary Table S2**,
209 **Supplementary Table S7**). We used the hierarchical biome organization of these samples
210 to construct a new biome ontology (Biome ontology (2020)), which included 36 biomes
211 in total (**Fig. 3a**). When comparing the accuracy between the independent model
212 (AUROC = 0.993, F-max = 0.986) and the Transfer (GM) model (AUROC = 0.991, F-
213 max = 0.978), we confirmed that EXPERT is able to source track accurately among
214 biomes in the Biome ontology (2020) (**Fig. 3b**). Again, we confirm that the high source
215 tracking accuracy actually come from the transfer learning process (**Supplementary**
216 **Table S5**). The average search time of the transfer model is also less than that of the
217 independent model (**Fig. 3c**). Two understudied biomes (root: Host-associated: Fish and
218 root: Host-associated: Birds: Digestive System: Ceca) were chosen to illustrate the source
219 tracking accuracy of EXPERT at specific layers of the Biome ontology (2020). We
220 noticed that the contribution of the correct biome (Fish and Birds) has a high value on
221 multiple layers of the Biome ontology (2020) (**Fig. 3d**). In a nutshell, EXPERT could
222 accurately source track samples from understudies or newly discovered biomes, enabling
223 lifelong adaptation for the emerging microbiome data.

224

225 **EXPERT for monitoring the succession of infant gut microbial communities**

226 In this context, we aim to utilize knowledge transfer to source track among more detailed
227 biomes, which is also referred to as "Scattering". Under this circumstance, we could
228 explore the dynamic patterns of gut microbiota from a specific period of life. For instance,
229 if infant samples from multiple time points and sources are offered, EXPERT could
230 estimate how much of microbial community in the infant's gut is originated from birth
231 and subsequent time points. To confirm this capability, we used longitudinal data from
232 Backhed et al.¹³, which consists of fecal samples from 98 infants and their mothers,
233 delivered by vaginal delivery or cesarean section (**Fig. 4a, Supplementary Table S2**,
234 **Supplementary Table S8**). Except for the general EXPERT model, here we introduced
235 the human EXPERT model (also referred to as the human model), which is generated
236 based on using samples from the human dataset consists of 53,553 samples collected

237 from 27 human-associated biomes (**Supplementary Table S2, Supplementary Table**
238 **S4**). The context-dependent ontology is organized by dividing samples by development
239 stage first, followed by birth mode (**Fig. 4a**).

240 In this part of the study, we considered samples from infants at 12 months of age as
241 queries, and samples from early time points or mothers were treated as sources. Firstly,
242 there is no significant difference among the samples from mothers (Wilcoxon test, $p =$
243 0.388 for Transfer (GM), $p = 0.929$ for Transfer (HM), **Fig. 4b**), which is consistent with
244 the results of J. Stokholm et al.¹⁴. This indistinguishable pattern is further supported by
245 Principal Coordination Analysis (PCoA) using distance metric either in weighted-
246 Unifrac¹⁵ or Jensen Shannon divergence¹⁶ (**Fig. 4c** and **Supplementary Fig. 3**), in which
247 samples from mothers and infants at 12 months are indistinguishable.

248 We then assessed the performances of different source tracking models via cross-
249 validation. For infant gut microbial communities at 12 months of age, even samples were
250 collected from hosts of different delivery modes, the maternal contribution is dominant
251 (**Fig. 4b**). When comparing with other methods, we found that the transfer model built
252 from the human model (Transfer (HM)) has the best performance (AUROC = 0.773)
253 compared to the independent model ((AUROC = 0.738) and the Transfer (GM) model
254 (AUROC = 0.720), indicating the outstanding performance of EXPERT for source
255 tracking gut microbial communities for infants at different stages (**Fig. 4d**). Again, we
256 confirm that the high source tracking accuracy actually come from the transfer learning
257 process (**Supplementary Table S5**). When examining the results of source tracking on
258 individual samples, we also observe that for quite a few samples the predicted biome
259 contributions by independent model were not consistent with ground truth, while the
260 predictions by Transfer (HM) were in agreement with ground truth, confirming the
261 advantage of transfer learning in this context, especially when samples are difficult to
262 distinguish (**Fig. 4e**). Furthermore, we made attempt by changed the context-dependent
263 setting: we have used the infant gut microbial community sample at birth as queries
264 (**Supplementary Fig. 4**), with results showing that samples from infant at birth have
265 most contributions from infant of 4 months, consistent with the results of J. Stokholm et
266 al.¹⁴. We have also changed the biome ontology, by means of dividing samples by birth

267 mode first, followed by development stage (**Supplementary Fig. 5**), and the source
268 tracking results of Transfer (HM) is still superior than other models.

269

270 **EXPERT for disease diagnosis and monitoring**

271 It has been reported that the human gut microbiota is potentially strongly associated with
272 a vast array of complex, chronic diseases¹⁷⁻²¹. Considering these potential associations
273 between the human gut microbial community and diseases, we also attempted to
274 characterize disease phenotypes (**Fig. 5a**). We collected 13,642 microbial community
275 samples from the GMrepo database⁷ and retrieved their corresponding disease
276 classification from the NCBI MeSH database²² (**Supplementary Table S2**,
277 **Supplementary Table S9**, and **Fig. 5b**). Except for the general EXPERT model and the
278 human EXPERT model, here we introduced the disease EXPERT model (also referred to
279 as the disease model), which is generated based on using samples from these 13,642
280 samples collected from 20 disease-associated biomes. Built based on the disease-related
281 biome ontology, EXPERT can precisely distinguish 20 phenotypes, including health and
282 19 gut-associated hematologic diseases, liver diseases, intestinal diseases, and bacterial
283 infections (over 0.8 AUROC for most phenotypes, **Fig. 5c**). Notably, a high AUROC
284 (over 0.8) was also observed when the model was rebuilt based on the human model
285 (Transfer (HM), **Fig. 5d**), suggesting the potential of EXPERT in host health status
286 prediction. To the best of our knowledge, this is the most comprehensive study in terms
287 of multiple disease diagnosis (13,642 samples, 20 health-associated phenotypes) with
288 superior accuracy²³. Again, we confirm that the high source tracking accuracy actually
289 come from the transfer learning process (**Supplementary Table S5**). These results also
290 suggests the potential role of the gut microbes in the alteration of host health status^{24,25}.

291 One confirmed association is that the gut microbiota is involved in colorectal cancer
292 (CRC) progression²⁶. It holds great potential to investigate whether CRC progression
293 could be monitored through gut microbiota. We assessed such applicability of EXPERT
294 by introducing CRC knowledge from Zeller, G. et al.²³: we considered five stages in the
295 progression of CRC: 0 (Healthy control) I, II, III, and IV according to the study of Zeller

296 G. et al.²³ (**Fig. 5f, Supplementary Table S2, Supplementary Table S10**). We first
297 found that the compositional shifts of the human gut within such progression are
298 indistinguishable by traditional methods, exemplified by Principle Coordination Analysis
299 (PCoA) using distance metric either in weighted-Unifrac¹⁵ or Jensen Shannon
300 divergence¹⁶ (**Fig. 5g, Supplementary Fig. 6**). Then we have compared the performance
301 of the independent model and transfer models built from the human model (Transfer
302 (DM)) and that from the disease mode (Transfer (DM)). Results have shown that Transfer
303 (DM) achieved a better performance (AUROC over 0.95, **Fig. 5h, i**) among these three
304 models, proving the superior applicability of EXPERT as a method for early detection of
305 the occurrence of colorectal cancers. We also notice that the samples used for this
306 analysis is not from a prospective study, and we expected that for an actual prospective
307 study on the development of cancers, more insights could be gained through such a
308 knowledge transfer approach.

309

310 **EXPERT implicates dynamic patterns of microbial communities in longitudinal** 311 **studies**

312 Using EXPERT for time-series analysis offers a quantitative way to characterize time-
313 related microbial compositional shifts, such as the dynamics of gut microbial
314 communities during international travel. In this context, by leveraging additional
315 metadata such as geographic regions and time points, we can characterize the
316 community's compositional shifts associated with exceptional events. To demonstrate this
317 capability, we used longitudinal samples from Liu H. et al.²⁷, including samples from ten
318 Chinese travelers (MT1-9 normally returned, while MT10 returned with an exceptional
319 early-back, **Fig. 6a,b**) who had a six-month long-stay from China to Trinidad and Tobago,
320 as well as both China and Trinidad and Tobago native persons (**Supplementary Table S2**
321 and **Supplementary Table S11**). Specifically, we performed an individual (people)-level
322 leave-one-out experiment for these ten travelers: each time considering all samples from
323 a specific traveler as queries and the rest samples (except for samples from MT10) as
324 sources (referred to **Methods** for details), ten times. And we generated a transfer model
325 based on transferring knowledge from the disease model (containing samples all from

326 human gut) for source tracking in this context. In these settings, our model can reveal the
327 compositional shifts of Chinese travelers gut microbial communities over time and
328 accurately pinpoint the early return timepoint of MT10 (**Fig. 6c**). These results are
329 consistent with their travel trajectory within ten months (from Dec 2015 to Sep 2016)²⁷.
330 And the significant difference (p-value 1.9×10^{-6} , 2.4×10^{-4} , and 3.9×10^{-5} for three stages,
331 Wilcoxon test) in source contributions revealed by our model also demonstrated
332 predictive power (for such compositional shifts) of EXPERT.

333 We further explored whether we could observe such time-related compositional shifts in
334 another cohort, about the seasonal changes of the Hadza persons' gut microbial
335 communities, by using longitudinal samples from Smith S. A. et al.²⁸ (**Supplementary**
336 **Table S2** and **Supplementary Table S12**). By dividing samples into "Dry" and "Wet"
337 categories (**Fig. 6d**), an average AUROC of 0.82 was achieved for distinguishing
338 seasonal patterns ("Wet" or "Dry") among gut microbial communities from Hadza hunter-
339 gatherers (**Fig. 6e**). These results are consistent with the results by the original study of
340 Smith S. A. et al.²⁸, confirmed the capability of EXPERT in longitudinal microbiome
341 analysis.

342

343 **Discussion**

344 This work presents an enabling technology for microbial source tracking based on
345 transfer learning, EXPERT, for source tracking microbial community samples in different
346 context-dependent applications, including prediction and monitoring the development of
347 diseases and longitudinal studies. Based on transfer learning techniques, it provides a fast,
348 accurate, flexible, and interpretable computational approach that could quickly adapt the
349 supervised model to source tracking samples from understudied and newly discovered
350 biomes. According to the needs of different applications, EXPERT can quantitatively
351 assign the contribution of biomes to a specific microbial community, in both ontology-
352 aware and context-aware manners, providing a method that could potentially illuminate
353 the microbial dark matters.

354 The utility of EXPERT is established in three contexts. First, we used EXPERT as it was

355 originally intended—to quantify the contribution of different source biomes. In this
356 context, we were able to address questions about more detailed biomes, including
357 microbial communities from infants of different stages. Specifically, using EXPERT, we
358 quantitatively reaffirmed the findings of J. Stokholm et al.¹⁴, suggesting that for 12
359 months infant, either delivered by cesarean section or vaginal delivery, there is no
360 significant difference in their gut microbial communities with those of the adults. Second,
361 we used EXPERT for disease diagnosis. In this context, EXPERT can serve for the early
362 detection of host health status and the scrutinizing of progression of cancer. Third,
363 EXPERT could be used for longitudinal analysis to better understand the dynamics of the
364 microbial communities. In this context, EXPERT can identify important events such as
365 “turning point” along the timeline, which might be linked to many aspects of human
366 physiology and health, including obesity¹⁸, inflammatory diseases²⁴, cancer²⁶, metabolic
367 diseases¹⁸, aging, etc. These context-dependent applications have highlighted the
368 necessity of the EXPERT model for microbiome studies, especially when facing the
369 understudied or newly discovered microbiome data.

370 The current approaches, including EXPERT, inevitably come up with multi-facet
371 challenges⁸ of the microbial dark matters . One is about the ontology structure, for which
372 a better representation might be a graph neural network²⁹ rather than a tree-like
373 hierarchical structure, for more precise quantification of contributions from biomes.
374 Another is when facing countless context-dependent applications, a collection of transfer
375 models (like Transfer (GM), Transfer (HM), Transfer (DM), etc.) might also be needed to
376 quickly adapt for applications either for environmental source tracking, large-scale time-
377 series analysis, or global scale pattern discovery.

378 Taken together, EXPERT is an ontology-aware neural network method based on biome
379 ontology information and transfer learning techniques, which may contribute to obtaining
380 biological insights from understudied or emerging microbiome dark matter. Combining
381 the supervised model-based efficiency and accuracy, together with the flexibility of the
382 transfer learning approach¹², EXPERT has enabled a broad spectrum of context-aware
383 applications, including adaptation to emerging data and host health monitoring as well as
384 longitudinal monitoring of microbial communities.

385

386 **References**

- 387 1. Gilbert, J. A., Jansson, J. K. & Knight, R. The Earth Microbiome project:
388 successes and aspirations. *BMC Biol.* **12**, 69 (2014).
- 389 2. Thompson, L. R. *et al.* A communal catalogue reveals Earth’s multiscale
390 microbial diversity. *Nature* vol. 551 457–463 (2017).
- 391 3. Turnbaugh, P. J. *et al.* The human microbiome project. *Nature* vol. 449 804–10
392 (2007).
- 393 4. Integrative, H. M. P. R. N. C. The Integrative Human Microbiome Project:
394 dynamic analysis of microbiome-host omics profiles during periods of human health and
395 disease. *Cell Host Microbe* vol. 16 276–89 (2014).
- 396 5. Gonzalez, A. *et al.* Qiita: rapid, web-enabled microbiome meta-analysis. *Nature*
397 *methods* vol. 15 796–798 (2018).
- 398 6. Mitchell, A. L. *et al.* MGnify: the microbiome analysis resource in 2020. *Nucleic*
399 *acids research* vol. 48 D570–D578 (2020).
- 400 7. Wu, S. *et al.* GMrepo: a database of curated and consistently annotated human gut
401 metagenomes. *Nucleic Acids Res* vol. 48 D545–D553 (2020).
- 402 8. Bernard, G., Pathmanathan, J. S., Lannes, R., Lopez, P. & Bapteste, E. Microbial
403 Dark Matter Investigations: How Microbial Studies Transform Biological Knowledge
404 and Empirically Sketch a Logic of Scientific Discovery. *Genome Biol Evol* **9** (2018).
- 405 9. Knights, D. *et al.* Bayesian community-wide culture-independent microbial
406 source tracking. *Nature methods* vol. 8 761–763 (2011).
- 407 10. Shenhav, L. *et al.* FEAST: fast expectation-maximization for microbial source
408 tracking. *Nature Methods* vol. 16 627 (2019).
- 409 11. Zha, Y. *et al.* Ontology-Aware Deep Learning Enables Ultrafast, Accurate and
410 Interpretable Source Tracking among Sub-Million Microbial Community Samples from
411 Hundreds of Niches. *bioRxiv* 2020.11.01.364208 (2020) doi:10.1101/2020.11.01.364208.
- 412 12. Pan, S. J. & Yang, Q. A Survey on Transfer Learning. *IEEE Trans. Knowl. DATA*
413 *Eng.* **22**, 15 (2010).
- 414 13. Bäckhed, F. *et al.* Dynamics and stabilization of the human gut microbiome
415 during the first year of life. *Cell Host Microbe* vol. 17 690–703 (2015).

- 416 14. Stokholm, J. *et al.* Delivery mode and gut microbial changes correlate with an
417 increased risk of childhood asthma. *Sci. Transl. Med.* **12**, (2020).
- 418 15. Lozupone, C. & Knight, R. UniFrac: a new phylogenetic method for comparing
419 microbial communities. *Applied and environmental microbiology* vol. 71 8228–8235
420 (2005).
- 421 16. Lin, J. Divergence measures based on the Shannon entropy. *IEEE Transactions on*
422 *Information theory* vol. 37 145–151 (1991).
- 423 17. Gupta, V. K. *et al.* A predictive index for health status using species-level gut
424 microbiome profiling. *Nat Commun* vol. 11 4635 (2020).
- 425 18. Chen, X. & Devaraj, S. Gut Microbiome in Obesity, Metabolic Syndrome, and
426 Diabetes. *Curr. Diab. Rep.* **18**, 129 (2018).
- 427 19. Shreiner, A. B., Kao, J. Y. & Young, V. B. The gut microbiome in health and in
428 disease. *Curr. Opin. Gastroenterol.* **31**, 69–75 (2015).
- 429 20. Ni, J., Wu, G. D., Albenberg, L. & Tomov, V. T. Gut microbiota and IBD:
430 causation or correlation? *Nat. Rev. Gastroenterol. Hepatol.* **14**, 573–584 (2017).
- 431 21. Theriot, C. M. & Young, V. B. Interactions Between the Gastrointestinal
432 Microbiome and *Clostridium difficile*. *Annu. Rev. Microbiol.* **69**, 445–461 (2015).
- 433 22. Sayers, E. W. *et al.* Database resources of the National Center for Biotechnology
434 Information. *Nucleic Acids Res.* **48**, D9–D16 (2020).
- 435 23. Zeller, G. *et al.* Potential of fecal microbiota for early-stage detection of colorectal
436 cancer. *Mol Syst Biol* vol. 10 766 (2014).
- 437 24. Balato, A. *et al.* Human Microbiome: Composition and Role in Inflammatory
438 Skin Diseases. *Arch. Immunol. Ther. Exp. (Warsz.)* **67**, 1–18 (2019).
- 439 25. Sirisinha, S. The potential impact of gut microbiota on your health: Current status
440 and future challenges. *Asian Pac. J. Allergy Immunol.* **34**, 249–264 (2016).
- 441 26. Zhu, Q., Gao, R., Wu, W. & Qin, H. The role of gut microbiota in the
442 pathogenesis of colorectal cancer. *Tumor Biol.* **34**, 1285–1300 (2013).
- 443 27. Liu, H. *et al.* Resilience of human gut microbial communities for the long stay
444 with multiple dietary shifts. *Gut* vol. 68 2254–2255 (2019).
- 445 28. Smits, S. A. *et al.* Seasonal cycling in the gut microbiome of the Hadza hunter-
446 gatherers of Tanzania. *Science* vol. 357 802–806 (2017).

- 447 29. Wu, Z. *et al.* A Comprehensive Survey on Graph Neural Networks. *IEEE Trans.*
448 *Neural Netw. Learn. Syst.* **32**, 4–24 (2021).
- 449 30. Leinonen, R. *et al.* The European Nucleotide Archive. *Nucleic Acids Res.* **39**,
450 D28–D31 (2011).
- 451 31. Kuczynski, J. *et al.* Using QIIME to analyze 16S rRNA gene sequences from
452 Microbial Communities. *Curr. Protoc. Bioinforma. Ed. Board Andreas Baxevanis Al*
453 **CHAPTER**, Unit10.7 (2011).
- 454 32. Kebschull, J. M. *et al.* High-Throughput Mapping of Single-Neuron Projections
455 by Sequencing of Barcoded RNA. *Neuron* **91**, 975–987 (2016).
- 456 33. Mukherjee, S. *et al.* Genomes OnLine Database (GOLD) v.8: overview and
457 updates. *Nucleic Acids Res.* **49**, D723–D733 (2021).
- 458 34. Schriml, L. M. *et al.* Human Disease Ontology 2018 update: classification,
459 content and workflow expansion. *Nucleic Acids Res* vol. 47 D955–D962 (2019).
- 460 35. Caruana, R. Multitask learning. *Machine learning* vol. 28 41–75 (1997).
- 461 36. Rumelhart, D. E., Hinton, G. E. & Williams, R. J. Learning representations by
462 back-propagating errors. *nature* **323**, 533–536 (1986).
- 463 37. Kingma, D. P. & Ba, J. Adam: A method for stochastic optimization. *arXiv*
464 *preprint arXiv:1412.6980* (2014).
- 465
- 466

467 **Methods**

468 **Datasets**

469 We used eight datasets to evaluate the performance of EXPERT.

470 For systematic assessment of our general model, the dataset was obtained from MGnify,
471 which consists of 118,592 communities collected from 132 biomes. Among them, 52,537
472 samples originated from human biomes, 14,045 samples originated from mammal biomes,
473 7,189 samples originated from terrestrial biomes, 27,667 samples originated from aquatic
474 biome. These samples were analyzed by MGnify⁶ before February 2020 (**Supplementary**
475 **Table S2 Supplementary Table S3**).

476 For systematic assessment of our human model, the dataset was a part of the first dataset,
477 in which 52,537 communities from 28 human biomes were selected (**Supplementary**
478 **Table S2, Supplementary Table S4**).

479 We also used emerging data in 2020 from MGnify⁶. Which consists of 34,209
480 communities collected from 35 biomes. Throughout the dataset, 3,421 samples belong to
481 8 biomes were newly added by MGnify⁶ after January 2020 (**Supplementary Table S2,**
482 **Supplementary Table S7**).

483 For source tracking the succession of infant gut microbiome, the dataset was obtained
484 from MGnify⁶ which consists of 392 fecal samples collected from 98 infant and their
485 biological mothers. Among them, 85 infants were born by vaginal delivery and 13 infants
486 were born by cesarean section. The infant samples were collected at three time points
487 including birth, fourth-month and twelfth-month. The maternal samples were collected
488 during the first week after delivery (**Supplementary Table S2, Supplementary Table**
489 **S8**).

490 For disease modeling, the dataset was obtained from GMrepo⁷, including 13,642
491 communities collected from feces of hosts diagnosed with 20 different phenotypical
492 status (different kinds of diseases, plus healthy controls, **Supplementary Table S2,**
493 **Supplementary Table S9**).

494 For cancer monitoring, the dataset was obtained from GMrepo⁷, which consists of 16, 93,

495 126, 196, and 204 communities respectively collected at CRC stage 0, I, II, III, and IV
496 (**Supplementary Table S2, Supplementary Table S10**).

497 For longitudinal dietary shifts analysis, the dataset was obtained from ENA³⁰ and analyzed
498 mainly using QIIME (version 1.9)^{27,31}, including 280 samples collected from travelers
499 and other native persons, (**Supplementary Table S2, Supplementary Table S11**).

500 For investigation of seasonal patterns within Hadza hunter-gatherers' gut microbiome,
501 The dataset used was obtained from ENA and analyzed mainly by using MapSEQ³²,
502 including 203 samples collected from the Hadza hunter-gatherers' gut (referred to Wu, S.
503 et al.⁷ for details, **Supplementary Table S2, Supplementary Table S12**).

504

505 **Data preprocess**

506 *Relative abundance calculation according to reference database*

507 In order to bypass the impact from the sequencing depth of the rich-sourced data, we
508 firstly regularized the abundance data by calculating relative abundances according to
509 only the taxa mapped to our phylogenetic tree (**Supplementary Table S1**), which is a
510 part of taxonomical classification tree in NCBI taxonomy database²² and reflects the
511 input feature set for EXPERT model (as explained below).

512 *Universal feature set for heterogeneous transfer learning*

513 In order to realize knowledge transfer among multi-faceted microbial source tracking
514 applications, we utilized a uniform collection of features in all these applications. Such
515 universal feature set was established according to only the variance of relative
516 abundances at genes level: Among the first dataset (for which the most comprehensive),
517 6,006 genera with variance of relative abundances above the threshold average variance \times
518 10⁻³. We constructed a phylogenetic tree based on the classification of these 6,006 genera
519 (three taxon ranks were included: "phylum", "order", "genus"), and served the tree for
520 abundance mapping (**Supplementary Note 2**).

521 *Z-score standardization of relative abundances*

522 In order to speed up the optimization of EXPERT model, we standardized the relative

523 abundances before feeding them into the model, by applying z-score standardization on
524 normalized relative abundances.

525

526 **Establishment of the ontology**

527 Throughout the paper, the reference information utilized for biome ontology
528 constructions are: (1) hierarchical biome classification from MGnify database⁶ and
529 ecosystem classification paths from GOLD database³³ for assessing the generalizing
530 performance of general knowledge (**Supplementary Table S13**), assessing the
531 performance in context-dependent settings (**Fig. 2a**), and assessing the performance on
532 emerging data (**Fig. 3a**); (2) sampling time of infant gut¹³ and delivery modes of infants¹³
533 for source tracking the succession of infant gut microbiome (**Fig. 4a, Supplementary Fig.**
534 **4**); (3) disease classification from NCBI MeSH database²² and Human Disease
535 Ontology³⁴ for disease diagnosis and monitoring (**Fig. 5a**); (4) geographic origins of
536 samples²⁷ and sampling time²⁸ for longitudinal data analyses (**Fig. 6a,b**).

537

538 **The EXPERT model**

539 *The probabilistic model of EXPERT*

540 Considering a query sample q represented by its community structure, as well as its
541 potential sources represented by a biome ontology O , to quantify contributions \hat{y}_q from
542 ontologically organized biome sources to q , we employed a multi-task neural network to
543 learn a mapping M from a series of source samples $s \in D_s$ to their biome sources,
544 $y_s = (y_s^2, \dots, y_s^I)$ (where y_s^2 is biome source for source sample s in the second layer of
545 the biome ontology), and then apply M on q to determine the contributions from biome
546 sources.

$$547 \quad \hat{y}_q = (\hat{y}_q^i)_{0 < i \leq I_o} = M(q)$$

548 *Fast inference via forward propagation*

549 We assume that the learning in the higher ontology layers, such as the distinguishing

550 between "Human" and "Environmental", are helpful to the learning in the lower ontology
 551 layers, such as the distinguishing between "Human: Digestive System", "Environmental:
 552 Aquatic" and "Environmental: Terrestrial"³⁵. EXPERT integrates the representation of the
 553 lower layer (which is calculated by its layer-specific modules M_{inter} , into higher layer),
 554 by employing several integrator module M_{integ} . Therefore, together with layer-specific
 555 output module M_{output} , the representation of the contributions is given by

$$556 \quad M(q) = (M_{output}^i(R_{integ}^i))_{0 < i \leq I_o}$$

557 Where

$$558 \quad R_{integ}^i(q) = \begin{cases} M_{integ}^i(M_{inter}^i(M_{base}(q)), 0) & \text{if } i = 1 \\ M_{integ}^i(M_{inter}^i(M_{base}(q)), R_{integ}^{i-1}) & \text{otherwise} \end{cases}$$

559 ***Robust optimization via backward propagation and transfer learning***

560 We also assume that the quantification of contributions from a biome ontology is helpful
 561 to the quantification of the contributions from a series of associated biome ontologies¹².

562 Considering \tilde{M}_{base} of a source model as a static mapping, the parameters of the rest
 563 modules \hat{w} can be solved using gradient descent as well as backpropagation
 564 algorithm^{36,37}

$$565 \quad \hat{w} = \underset{\hat{w}}{\operatorname{argmin}} \sum_{i=0}^{I_o} \alpha(B_O^i) L(\hat{y}_s^i(\hat{w}), y_s^i)$$

566 Where

$$\alpha(B_O^i) = \frac{B_O^i}{B_O f}$$

$$567 \quad L(\hat{y}_s^i, y_s^i) = \sum_{b \in O^i} (\operatorname{CrossEntropy}(\hat{y}_s^i(b), y_s^i(b)))$$

$$b \in O^i$$

568 B_O^i stand for the number of biomes contained in the i -th layer of the biome ontology

569 O^i, O^i stand for the i -th layer of the biome ontology O^i

570 Then, optimizing the parameters of the entire model (including \tilde{M}_{base}), the parameters of
571 the rest modules w can be solved by using gradient descent as well as backpropagation
572 algorithm^{36,37}

$$573 \quad w = \underset{w}{\operatorname{argmin}} \sum_{i=0}^{I_o} \alpha(B_O^i) L(\hat{y}_s^i(w), y_s^i)$$

574 For independent optimization (optimization based on completely random initialization),
575 EXPERT straightforwardly optimizes the entire model. See **Supplementary Note 5** for
576 detailed description for optimization.

577

578 **Cross-validation**

579 We evaluated performances of EXPERT utilizing multiple cross-validation settings.
580 When assessing the performances of general model (including adapted model on
581 emerging data), human model, and disease model, we repeatedly performed random
582 cross-validation for five times--each time randomly select 10% of the dataset as queries
583 and the rest as sources.

584 When source tracking the succession of infant gut microbiome, we performed
585 proportional sampling cross-validation by randomly select 10% of the samples from
586 mother, infants at birth, 4 months and 12 months of age as query, respectively. This
587 process was also repeated for five times.

588 When analyzing the compositional shifts of ten Chinese individuals during the travel to
589 Trinidad and Tobago, we performed individual-level leave one out cross-validation--each
590 time consider all samples of an individual as queries and all samples from the other
591 individuals as sources. Notably, we didn't consider samples from MT10 as sources as he
592 returned with an exception.

593 When analyzing the compositional shifts of Hadza hunter-gatherers, we performed

594 proportional sampling cross validation by randomly select 50% of the samples from from
595 young adults ($18 \leq \text{age} < 50$) in each season ("2013-Late Dry", "2014-Early Wet",
596 "2014-Late Wet", "2014-Early Dry", and "2014-Late Dry") as queries and the rest as
597 sources. This process was also repeated for five times.

598

599 **Performance measures**

600 To benchmark and compare the performance of EXPERT model based on knowledge
601 transfer (Transfer (GM), Transfer (HM), Transfer (DM),) and independent model, as well as
602 other methods, we used these measures:

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

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

603

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

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

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

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

604

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

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

605 Where TP is true positive, TN is true negative, FP is false positive, FN is false negative,
606 $\hat{y}_s(b)$ is the quantified contribution from a biome source b for a microbial
607 community sample s , threshold $t \in [0,1]$ with a step size of 0.01, y_s is a set of actual
608 biomes for a sample s , and I is a logical operation function, the value of I is 1 when the
609 result of logical operation is TRUE, else 0.

610 Then, three evaluation metrics (Accuracy, F-max, AUROC) was introduced. These

611 evaluation metrics were calculated with the following formulas:

$$Accuracy_b(t) = \frac{TP_b(t) + TN_b(t)}{TP_b(t) + TN_b(t) + FP_b(t) + FN_b(t)}$$

612
$$F-max_b(t) = \max_t \frac{2Precision_b(t)Recall_b(t)}{Precision_b(t) + Recall_b(t)}$$

$$AUROC_b(t) = \sum_{t=0}^1 \frac{(TPR_b(t) + TPR_b(t + 0.01))(FPR_b(t + 0.01) - FPR_b(t))}{2}$$

613 Then, we treated the average performance across all biomes in a specific ontology layer
614 as the performance of the entire model.

615

616 **Statistical analysis**

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

621

622 **Data distribution**

623 Throughout the paper, the box-plot elements are: center line, median; box limits, upper
624 and lower quartiles; whiskers, $1.5 \times$ interquartile range (IQR); points and outliers. The
625 Violin plot is also used for data distribution analysis, mainly for comparison.

626

627 **Data availability**

628 The collected samples from MGnify/GMrepo were annotated with their associated
629 biomes/phenotypes in **Supplementary Table S3-4** and **Supplementary Table S7-12**. All
630 the processed data are uploaded and hosted at [https://github.com/HUST-NingKang-](https://github.com/HUST-NingKang-Lab/EXPERT-use-cases)
631 Lab/EXPERT-use-cases.

632 **Code availability**

633 All source codes have been uploaded to the website at: [https://github.com/HUST-](https://github.com/HUST-NingKang-Lab/EXPERT)
634 NingKang-Lab/EXPERT. Detailed parameters of software and package used in this study
635 are provided in Supplementary Table S2.

636

637 **Acknowledgments**

638 This work was partially supported by National Science Foundation of China grant
639 32071465, 31871334 and 31671374, and the Ministry of Science and Technology's
640 national key research and development program grant (No. 2018YFC0910502).

641

642 **Author contributions**

643 KN and HC conceived of and proposed the idea, and designed the study. HC, QY, GX,
644 NW, CS, and SW performed the experiments and analyzed the data. HC, QY, YZ, WC,
645 and KN contributed to editing and proof-reading the manuscript. All authors read and
646 approved the final manuscript.

647

648 **Competing interests**

649 The authors declare that they have no competing interests.

650

651 **Ethics approval and consent to participate**

652 Not applicable.

653

654 **Figures**

655 **Figure 1**

656 **Figure 1: Illustration of EXPERT's data processing, knowledge transfer process, and multi-**
657 **faceted applications. a.** The data preprocessing workflow of EXPERT. For each sample, we mapped
658 the taxonomical abundances to the phylogenetic tree, which is compatible with NCBI phylogenetic
659 tree, to obtain a regular abundance matrix. And then, the matrix is normalized and standardized in
660 order to obtain a standard abundance matrix. Blue and White boxes indicate entities of data and
661 operations, respectively. The final input data is a matrix, in which each row represents abundances for
662 a sample, and each column represents abundances for a genus (in total, 6,006 genera were used in this
663 study). **b.** The general ontology-aware neural network (ONN) model (built based on a biome ontology
664 with more than a hundreds biomes, referred to <https://github.com/HUST-NingKang-Lab/ONN4MST>
665 for details) has a fixed structure and poor scalability, thus cannot perform well on source tracking for
666 context-dependent applications. **c.** EXPERT can adapt the general knowledge of an existing model to
667 such applications through three steps: transfer (reuse parameters of an existing model and reinitialize
668 context-dependent layers according to prior pieces of knowledge about the context, green arrows),
669 adaptation (quickly optimize only the context-dependent layers using iterative forward-backward
670 propagation, green circular arrows), and finetuning (further optimize the entire model using the
671 iterative forward-backward propagation). The existing model is a general ONN model (either in GM,
672 HM, or DM model used in this work, refer to **Supplementary Table S2** for details) to be adapted,
673 with two fully connected layers (relatively independent to context) and a series of context-dependent
674 layers (highly specified to a context, **Supplementary Fig. 1**). Red arrows indicate the three steps of
675 the transfer process. Different background colors of the model indicate the suitability of the layers to
676 the context-dependent application. This model could be applied to community samples collected from
677 new biome ontology (representing context-dependent knowledge) that cannot be source-tracked by the
678 general ONN model directly. The transferred model can serve a broad spectrum of source tracking
679 applications (based on research purposes, further illustrated in **Fig. 1d**). **d.** The applications of
680 knowledge transfer utilizing EXPERT. **i.** EXPERT quantifies the source biome contributions for query
681 communities. **ii.** EXPERT enables adaptation of emerging microbiome data from understudied biomes
682 (also referred to as Grafting). **iii.** Knowledge transfer can source track among more detailed biomes
683 (also referred to as Scattering), such as quantifying gut microbial communities' contribution from
684 mothers to their developing infants. **iv.** EXPERT enables disease diagnosis, as well as disease
685 progression monitoring. **v.** EXPERT enables investigation of dynamics of microbial communities
686 along the timeline.

687

688 **Figure 2**

689 **Figure 2. Assessment of the transfer learning model on context-dependent application. a.** Biome
690 ontology of human-associated biomes, used for assessing the generalizing performance of EXPERT.
691 The ontology is constructed based on 52,537 samples, with four layers in total (from top down). The
692 biomes with orange color indicates biomes with enough samples (> 100) within each biome. The
693 biomes with orange color and in dashed line box with yellow background indicates biomes for
694 comparison with FEAST, which are selected as they represent most detailed biomes (comply with the
695 culture-independent assumption of SourceTracker and FEAST,), and each of them has enough samples
696 (> 100). **b.** The performance (AUROC and F-max, X-axis) of Transfer (GM) for each biome (Y-axis)
697 by using repetitive cross-validation (5 times, 90% for training, 10% for testing, see **Methods**). **c.** The
698 performances (AUROC and F-max, Y-axis) of three models (X-axis): Transfer (GM) model, Transfer
699 (GM0) model, and Independent model (built based on independent training). **d.** The training time and
700 query time (Y-axis) across different optimization schemes (Independent, Transfer, Transfer (0), X-
701 axis). **e.** Comparison of Transfer (GM) with FEAST in source tracking accuracy (first Y-axis) and
702 efficiency (second Y-axis). FEAST10, FEAST20, FEAST30 represent FEAST based on using 10, 20
703 and 30 samples per biomes for all biomes with enough samples (> 100) in the human dataset as source
704 samples, respectively. Transfer (GM) and Transfer (GM0) refer to models built based on transferring
705 knowledge from the general EXPERT model with and without finetune, respectively.

706

707 **Figure 3**

708 **Figure 3. EXPERT enables adaptation for emerging microbiome data. a.** Biome ontology (2020)
709 constructed based on 34,209 samples collected and analyzed by MGnify after January 2020. These
710 samples are from existing biomes and the newly discovered biomes, hierarchically organized in four
711 layers (indicated by the black interval on the left, from top down). The yellow nodes indicate newly
712 discovered biomes in the MGnify database, which have hardly ever been analyzed before January
713 2020. **b.** Generalizing performance (AUROC and F-max, Y-axis) of models optimized using different
714 training schemes (Independent model based on 34,209 samples used in this application, and Transfer
715 (GM) model built based on transferring knowledge from the general model with finetune, X-axis). **c.**
716 Searching time per sample by different models. **d.** Estimated biome contributions for query samples,
717 exemplified by samples from two newly discovered biomes ("Fish" and "Birds") on different layers.
718 Notice that on the third layer, the contribution of Fish is 88.29% for query samples from Fish, and the
719 contribution of Birds is 73.78% for query samples from Birds; and on the fourth layer, the contribution
720 of Digest system (birds) is 72.22% for query samples from Birds; while on the fifth layer, the
721 contribution of Ceca is 56.35% for query samples from Birds.

722

723 **Figure 4**

724 **Figure 4. EXPERT's performance in characterizing gut microbial community development over**
725 **time for infants. a.** The biome ontology, corresponding to infant samples collected from the ENA
726 database. Layer 2 is based on sampling time, while layer three is based on delivery modes (Backhed et
727 al.). For this part of the study, sources include gut microbiome of the mother, infant at birth, and four
728 months, queries include the gut microbiome of the infant at 12 months. **b.** Estimated contributions of
729 biomes by different models of EXPERT, separated by two delivery modes and transferred from the
730 general model (Transfer (GM), above) and human model (Transfer (HM), below). **c.** Distribution of
731 infant gut microbial communities during their first year, using principal-coordinates analysis (PCoA)
732 and distance metric of Jensen Shannon divergence. Infant gut given by different delivery modes
733 developed in an adult-like pattern over time, and the compositional shifts of infant's gut at the age of
734 12 months are indistinguishable. Dotted line refers to samples delivered by vaginal delivery, and the
735 full line refers to samples delivered by cesarean section. The baby of 4 month is abbreviated to baby
736 4M., the baby of 12 month is abbreviated to baby 12M. "C" represents cesarean section, "V"
737 represents vaginal delivery. Top panel: samples from infant's gut are plotted according to their source
738 and collection date on the Y-axis, and position on the X-axis is plotted according to their first principal
739 coordinate in the PCoA. **d.** The overall performance of models generated based on different models, in
740 which Independent model based on the samples and biomes used in this application, Transfer (GM)
741 and Transfer (HM) refer to models built based on transferring knowledge from the general model and
742 human model with finetune, respectively. Violin plot: Red represents AUROC, and blue represents F-
743 max. **e.** Samples whose predicted source contributions were not consistent with the ground truth by the
744 independent model, but were consistent with the ground truth by the Transfer (HM) model. Different
745 colors refer to different biomes: purple represents unknown sources, blue represents the baby of 12th
746 month, green represents the baby of 4th month, red represents the baby of birth. Baby of 4 month is
747 abbreviated to baby 4M., the baby of 12 month is abbreviated to baby 12M.

748

749 **Figure 5**

750 **Figure 5. EXPERT for disease diagnosis and monitoring. a.** Illustration of knowledge transfer
751 utilized for disease phenotype differentiating. The knowledge transfer between models trained using
752 different datasets (named Source dataset and Target dataset) are illustrated using different colors
753 (white for human, yellow for disease, and red for colorectal cancer). In health status prediction, the
754 knowledge from human model containing 53,553 samples and 28 biomes, were transferred to disease
755 model containing 13,642 samples and 20 biomes. **b.** the disease ontology constructed based on host
756 phenotype and considering NCBI MeSH database (<https://www.ncbi.nlm.nih.gov/mesh>) and Human
757 Disease Ontology (<https://www.ebi.ac.uk/ols/ontologies/doid>) as references of disease classification.
758 The disease ontology includes 20 phenotypes (19 different disease and infections, plus healthy control)
759 distributed in seven different layers (X-axis). **c.** The diagnostic model's performance on each
760 phenotype, evaluated based on repetitive cross-validation (5 times, 90% for training, and 10% for
761 testing) and biome-specific evaluation (see **Methods**). A dashed line indicates AUROC of 0.8. **d.** The
762 Transfer (HM) model's overall performances on differentiation of diseases. **e.** Illustration of
763 knowledge transfer utilized for differentiating disease phenotypes. We generated two transferred
764 models: Transfer (HM) and Transfer (DM) refer to models built based on transferring knowledge from
765 the human model and disease model with finetune, respectively. In CRC progression monitoring, the
766 knowledge from Transfer (HM) containing 53,553 samples and 28 biomes, as well as Transfer
767 (DM) containing 13,642 samples and 20 biomes, were transferred to CRC applications in which there
768 are 5 stages. **f.** The five stages of colorectal cancer progression, and the number of samples for each
769 stage. Among them, stage 0 stand for healthy control. **g.** The distribution of gut microbiome,
770 visualized by PCoA (utilizing distance metric of weighted-Unifrac). **h.** The performances (AUROC
771 and F-max) of different models (Independent model based on independent training, Transfer (HM)
772 based on the human model, and Transfer (DM) based on disease model). **i.** The stage-specific
773 performances (AUROC, see **Methods**) of EXPERT on different CRC stages, evaluated based on
774 repetitive cross-validation (5 times, 90% for training, and 10% for testing).

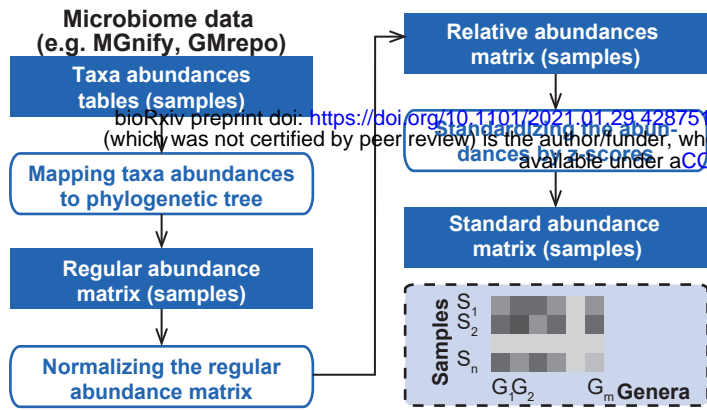
775

776 **Figure 6**

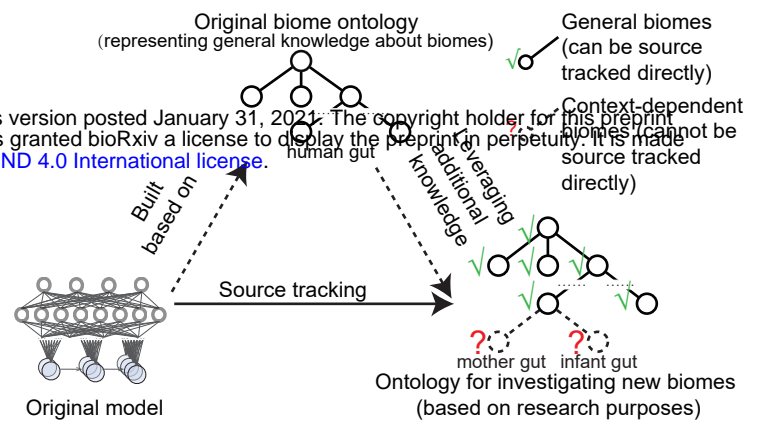
777 **Figure 6. The EXPERT model could reveal compositional shifts within hosts' gut microbiome**
778 **over time. a.** Biome ontology used for analyzing dynamic patterns among gut microbiome during
779 international travel. Two countries (China and Trinidad and Tobago, black nodes) leveraged in the
780 travel are modeled into the biome ontology. Additional information (dotted green nodes) can be
781 inferred based on the model's prediction and additional metadata. **b.** Sampling time points (30 time
782 points in total, X-axis) of metagenomic samples for each individual (Y-axis). Three stages are included:
783 before travel (in China), during travel (in Trinidad and Tobago), and after travel (returned to China),
784 and transitions between these stages (flights) are marked on the top. Different colors indicate samples
785 from different individuals. **c.** Estimated source contributions (Y-axis) by EXPERT over the timeline
786 (also 30 time points in total, X-axis), using individual-level leave one out and considering additional
787 samples from native people as potential sources (see **Methods**). **d.** Biome ontology used for analyzing
788 seasonal patterns among Hadza hunter-gatherers' gut microbial communities. Samples are divided into
789 “Dry” and “Wet” biomes in the biome ontology, Additional information (dotted green nodes) can be
790 inferred based on the model's prediction and additional metadata. **e.** Estimated source contributions
791 (Y-axis) by EXPERT over the timeline (15 time points in total, X-axis), by using samples from the
792 “Dry” season as sources. Result are evaluated based on repetitive proportional sampling cross-
793 validation (see **Methods**).

794

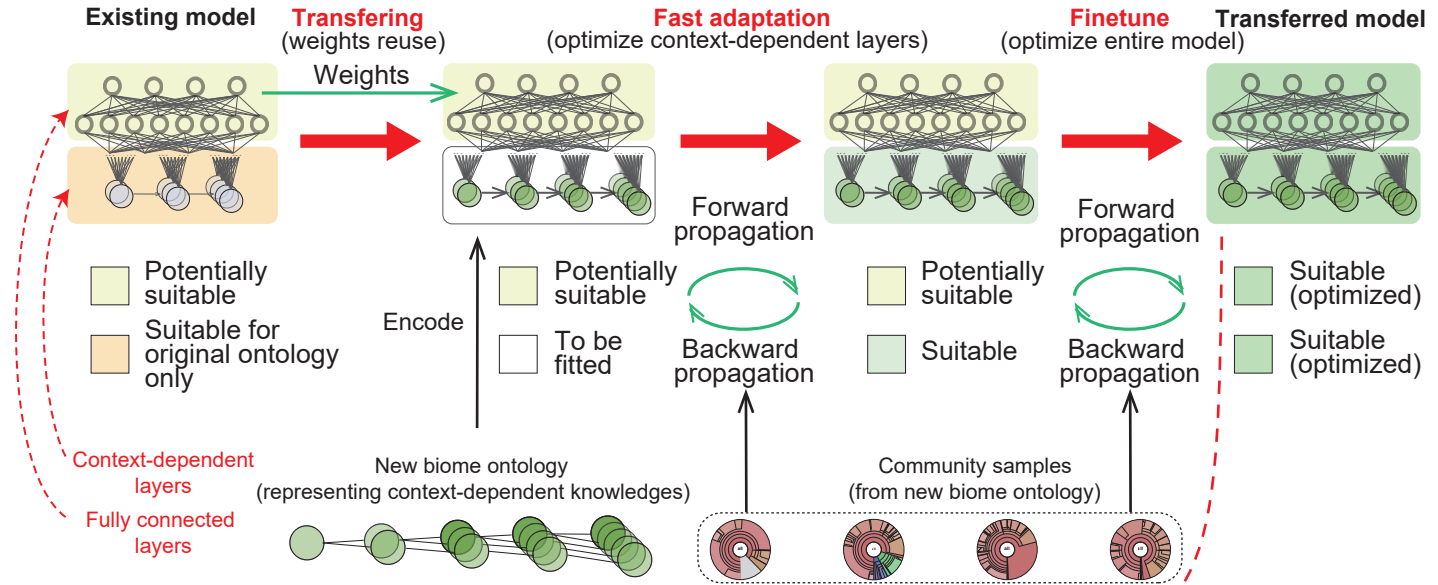
a. Preprocessing microbiome data



b. Source tracking for context-dependent applications

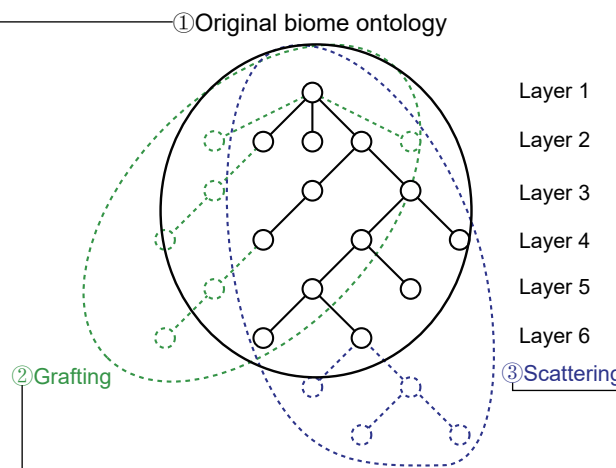
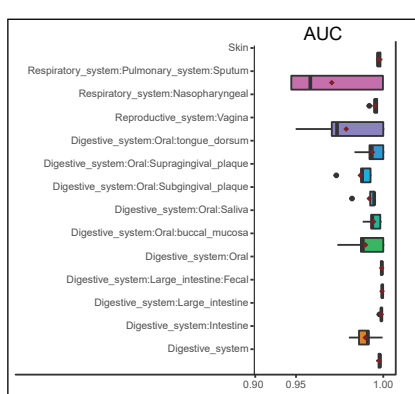


c. Transfer learning-enabled adaptation of original model into context-dependent applications

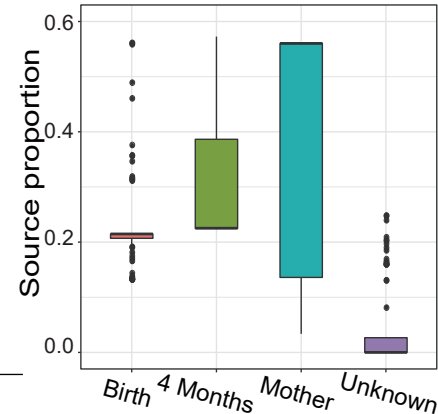


d. EXPERT for microbial source tracking

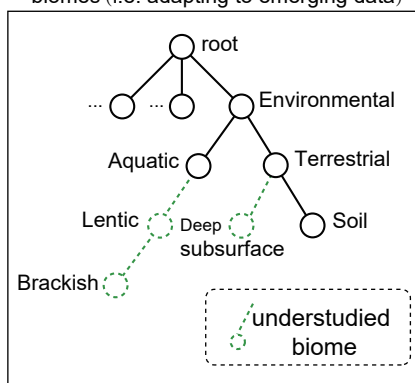
Multi-facets of EXPERT applications



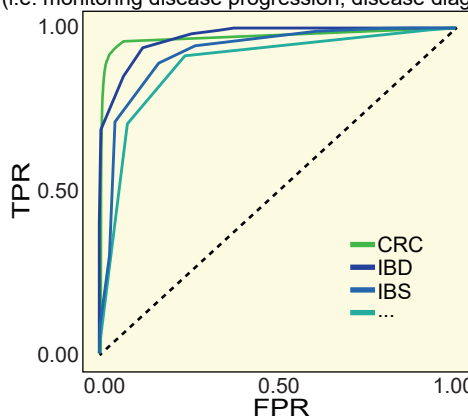
iii. Source tracking among detailed biomes (i.e. mother & infant relationship)



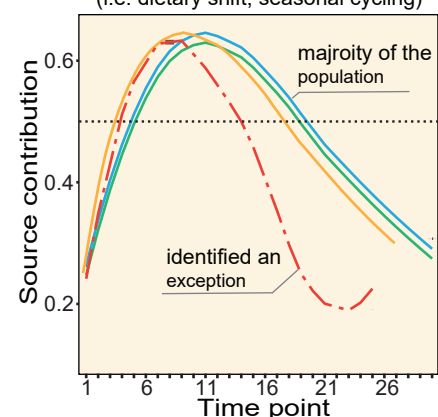
ii. Source tracking among newly discovered biomes (i.e. adapting to emerging data)



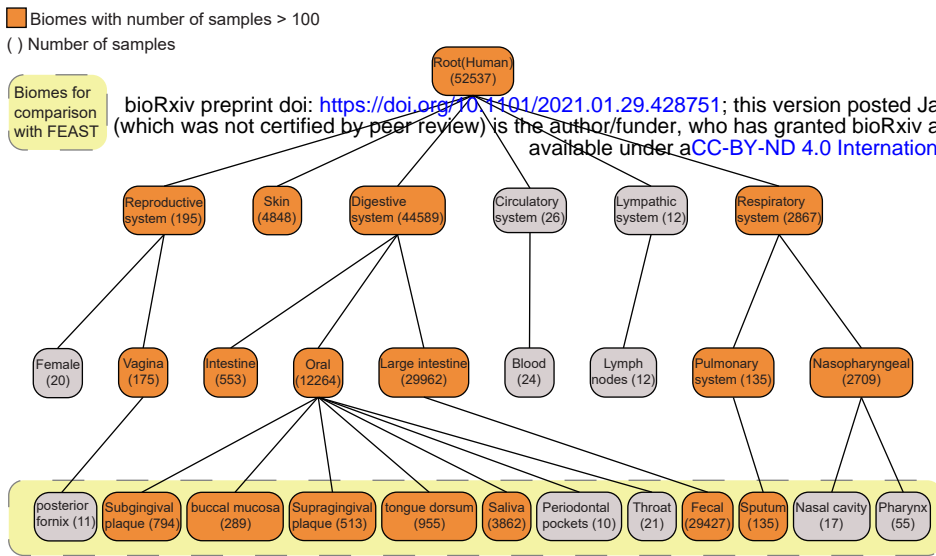
iv. EXPERT for disease diagnosis (i.e. monitoring disease progression, disease diagnosis)



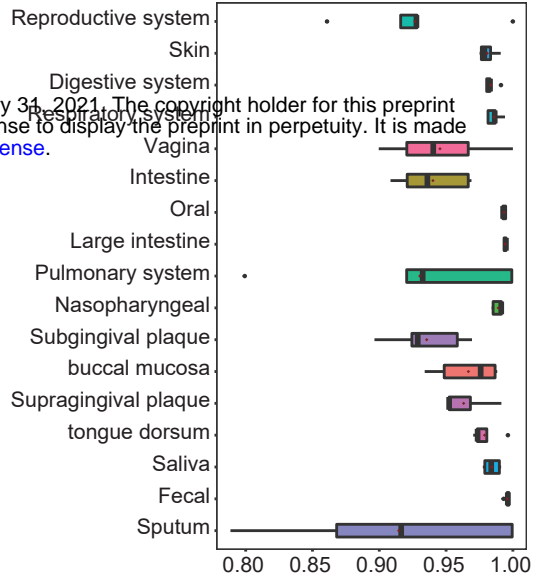
v. EXPERT for longitudinal analysis (i.e. dietary shift, seasonal cycling)



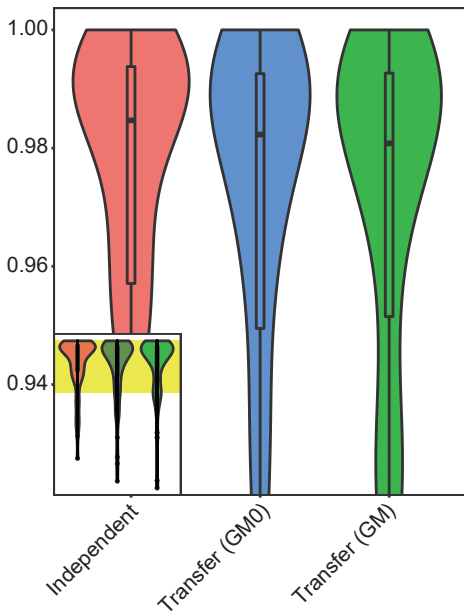
a Human biomes ontology



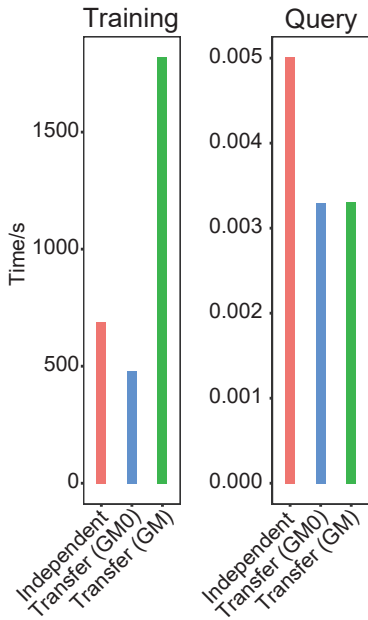
b. Prediction results by Transfer (GM)



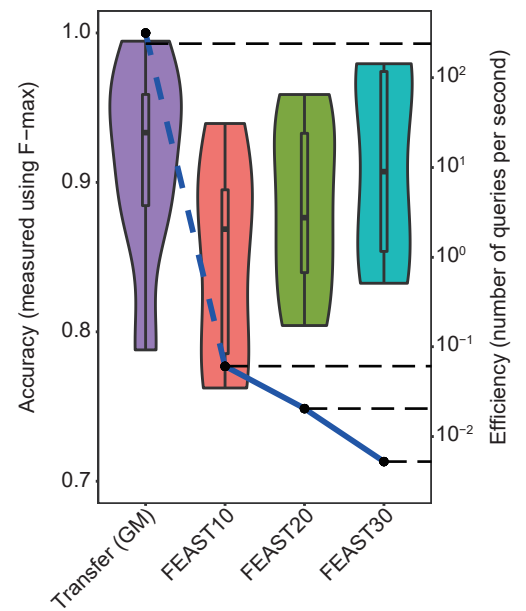
c Transfer and Independent model prediction results



d Training & sample query time



e Comparison of Transfer (GM) and FEAST



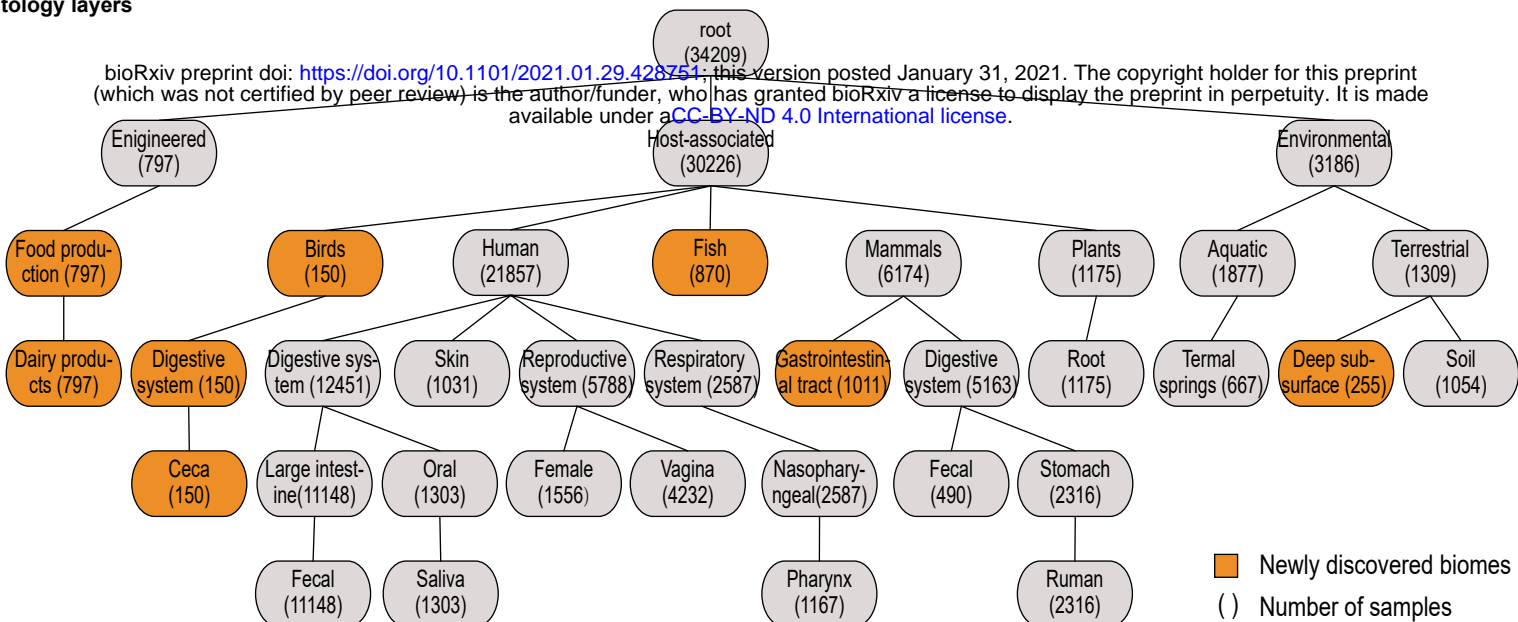
Biome ontology(2020)

a.

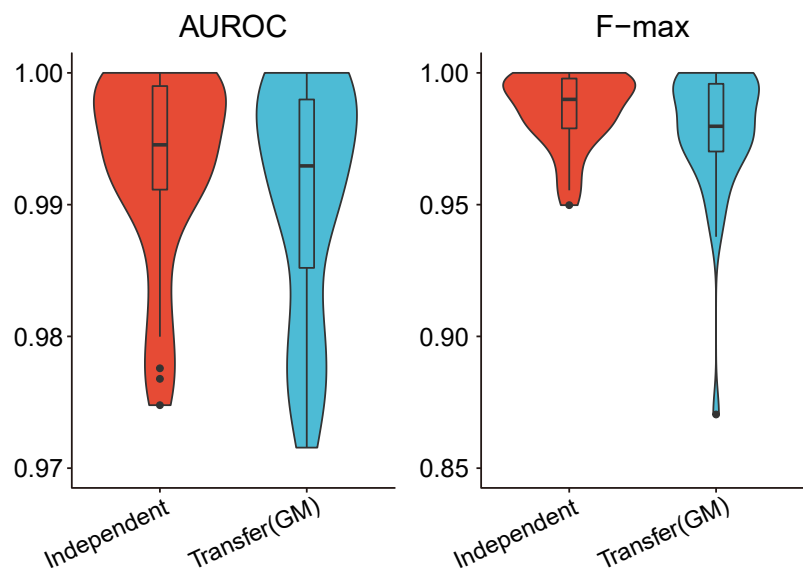
Ontology layers

1
2
3
4
5
6

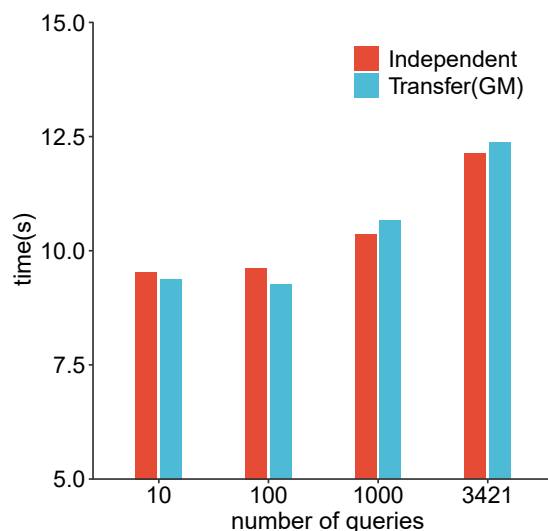
bioRxiv preprint doi: <https://doi.org/10.1101/2021.01.29.428751>; this version posted January 31, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-ND 4.0 International license](https://creativecommons.org/licenses/by-nd/4.0/).



b. Performance of different models after adding new biomes



c. Query time usage



d.

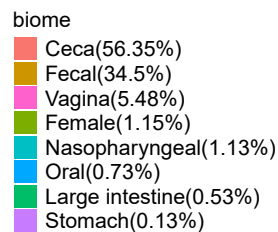
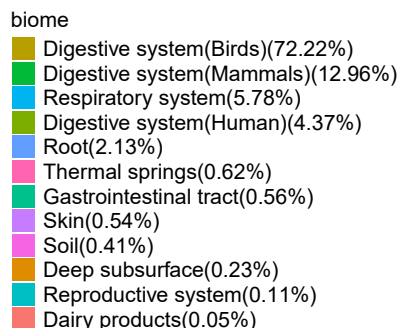
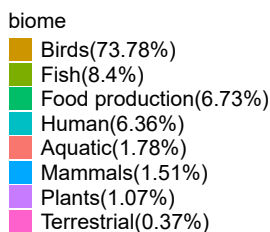
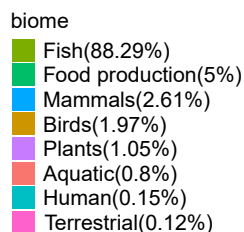
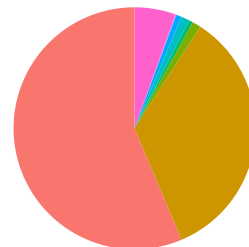
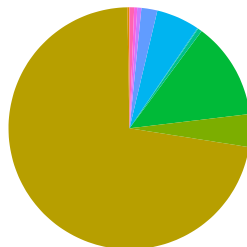
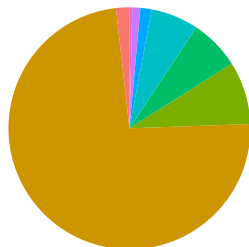
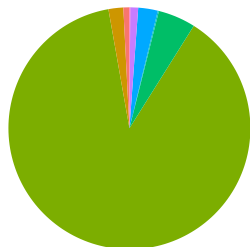
Contributions from different biomes

layer3 contribution of fish

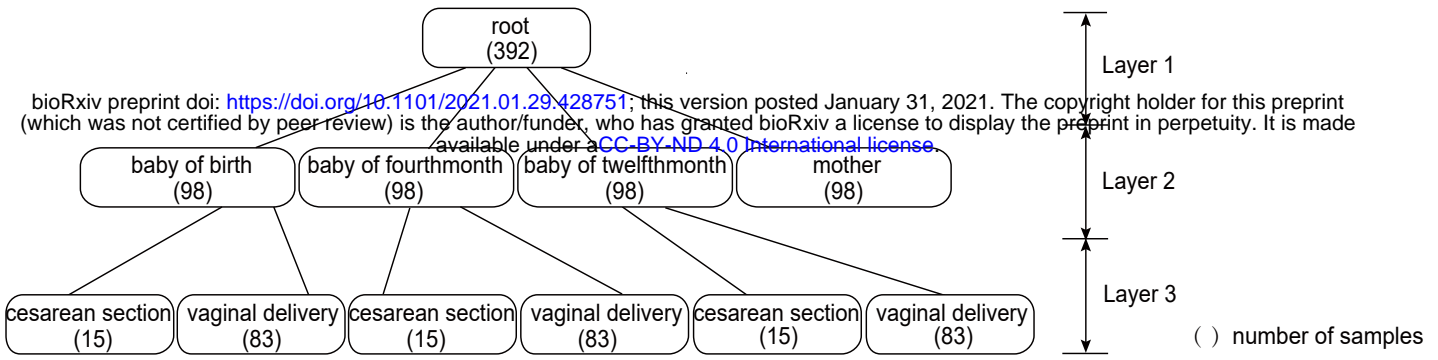
layer3 contribution of birds

layer4 contribution of birds

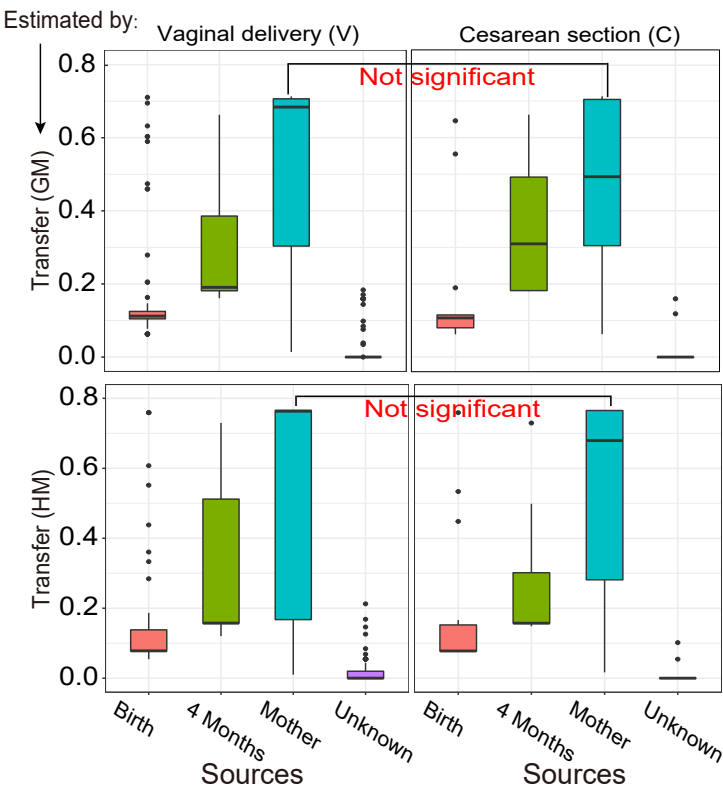
layer5 contribution of birds



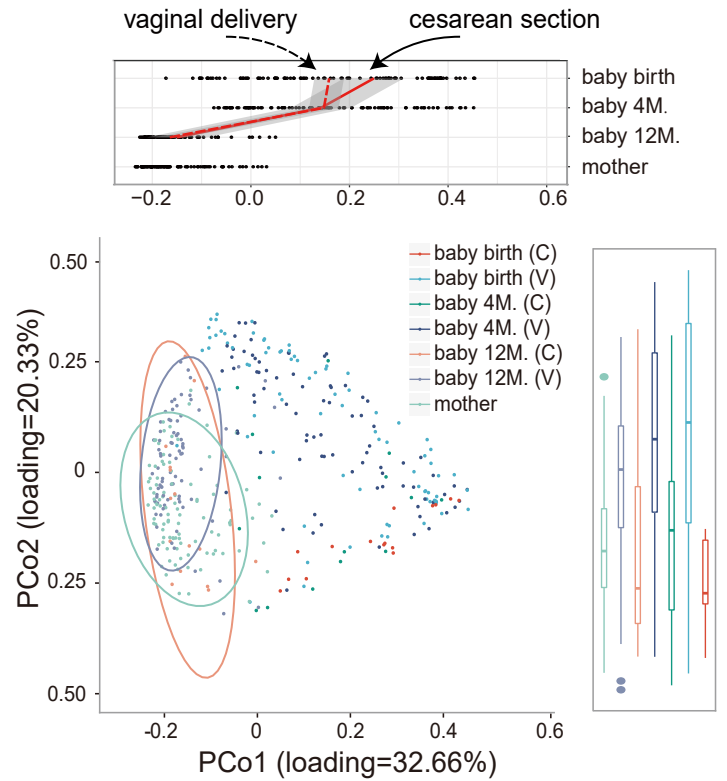
a. Biome ontology(layer 2: sampling time)



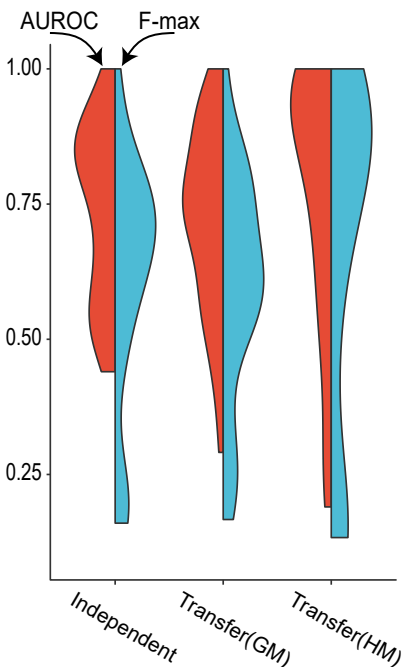
b. Estimated source proportions by different models (queries: samples of 12 months)



c. Distribution of infant gut microbiome during first year of age using PCoA

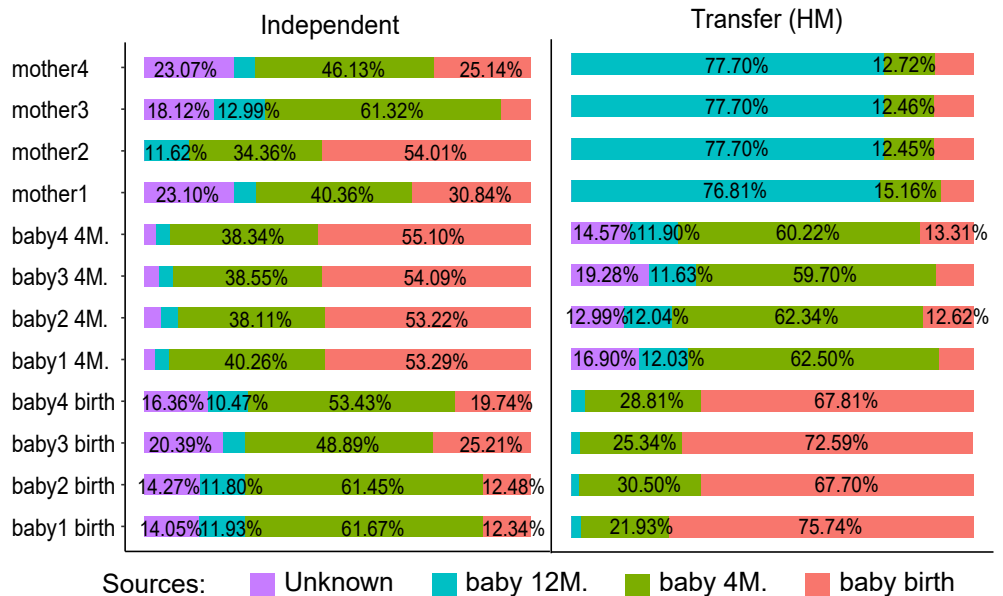


d. Overall performances

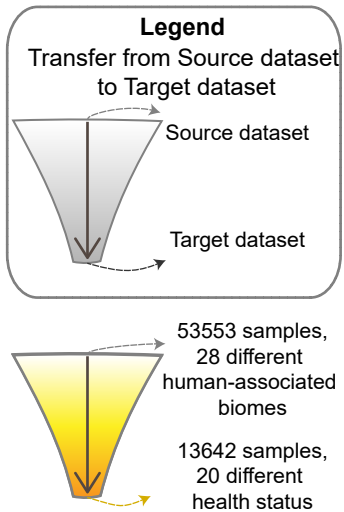


e.

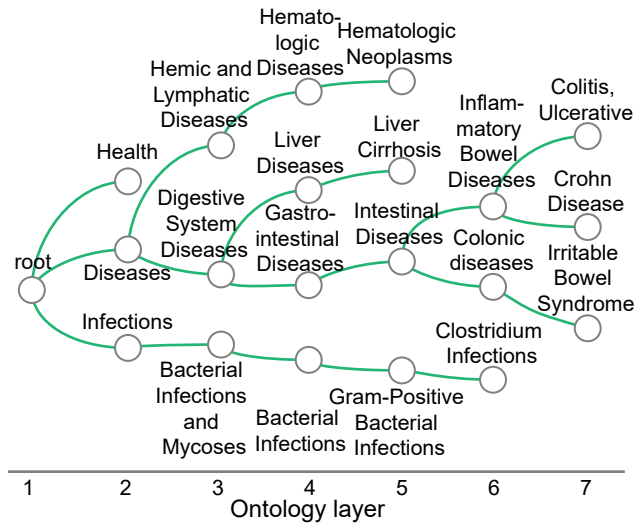
Comparison of estimated contributions of source biomes for different samples



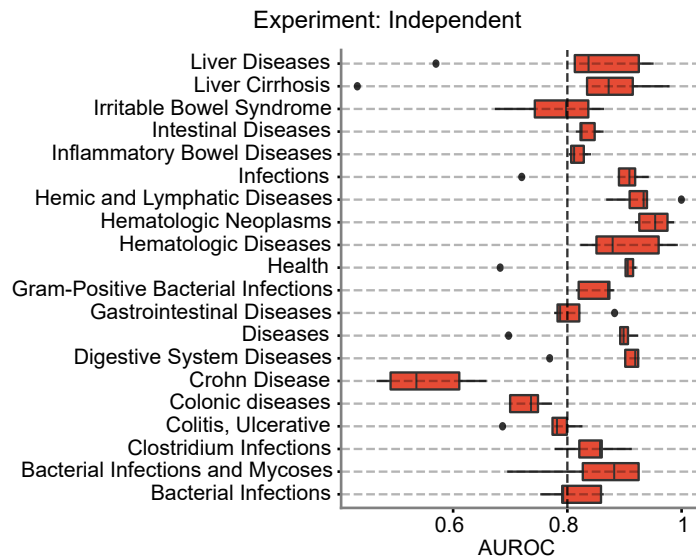
a. Transfer process



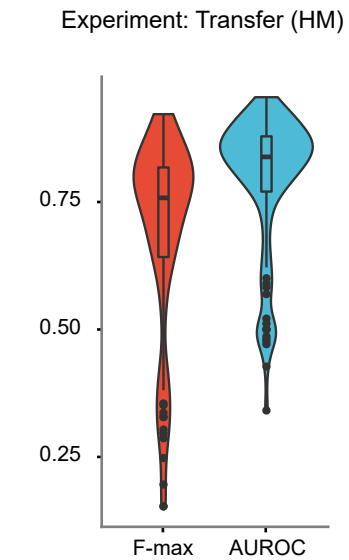
b. Disease Ontology



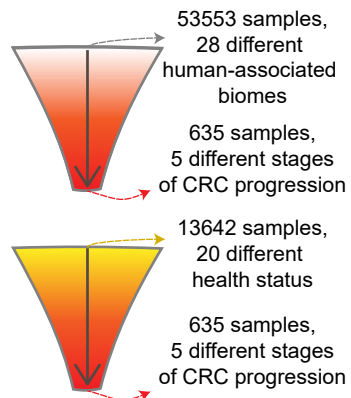
c. Distinguishing disease phenotypes



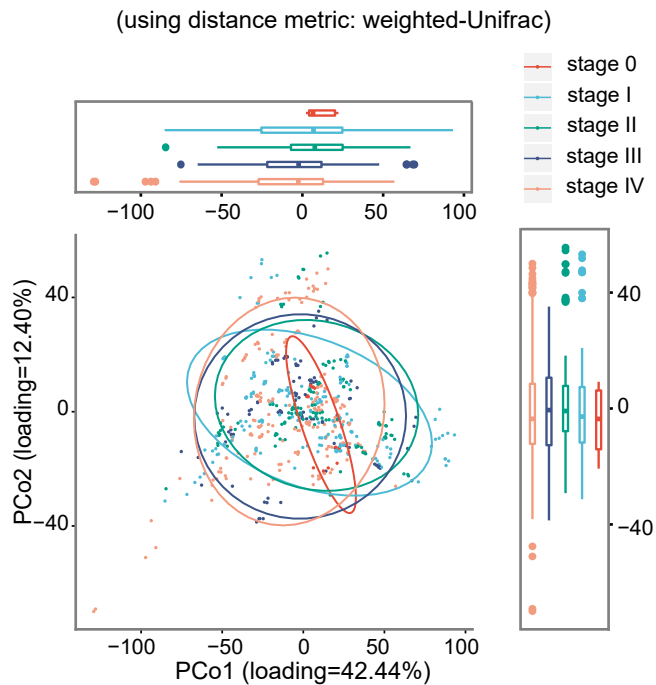
d. Distinguishing diseases



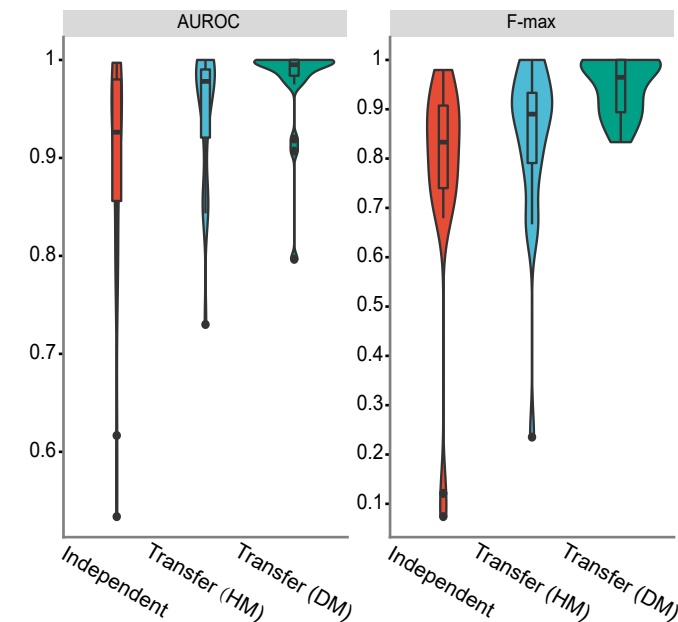
e. Transfer process



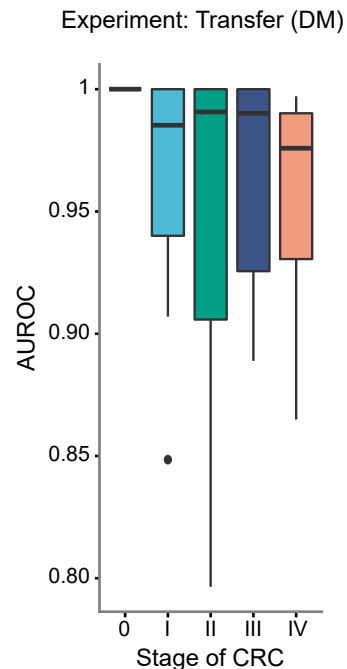
g. Distribution of gut microbiome across CRC stages



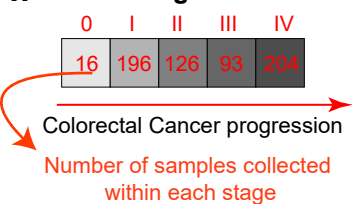
h. Differentiating stages for CRC



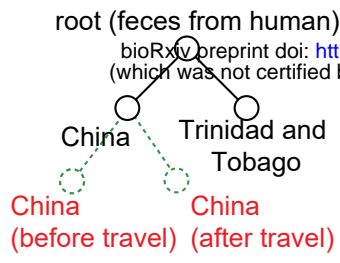
i. Transfer (DM)



f. The stage of CRC

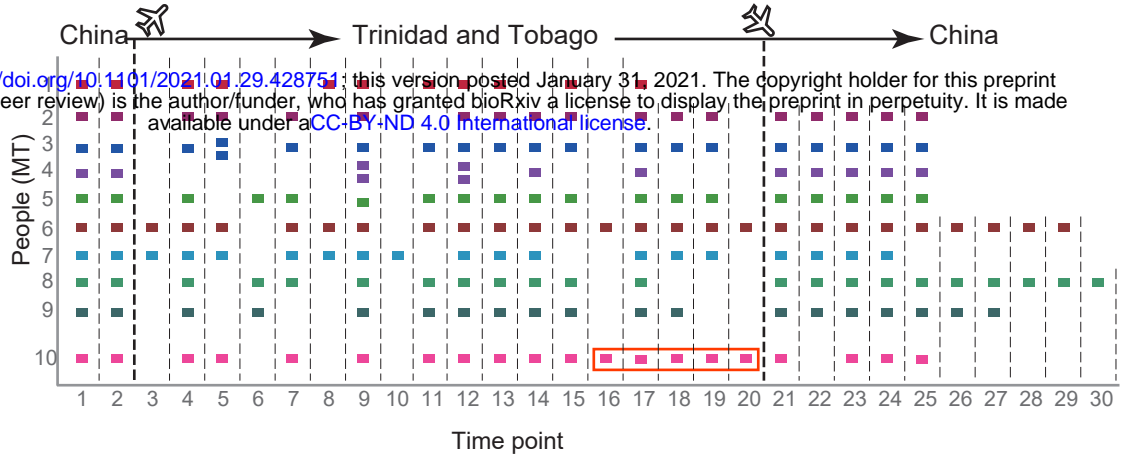


a. Biome ontology
for travel longitudinal study

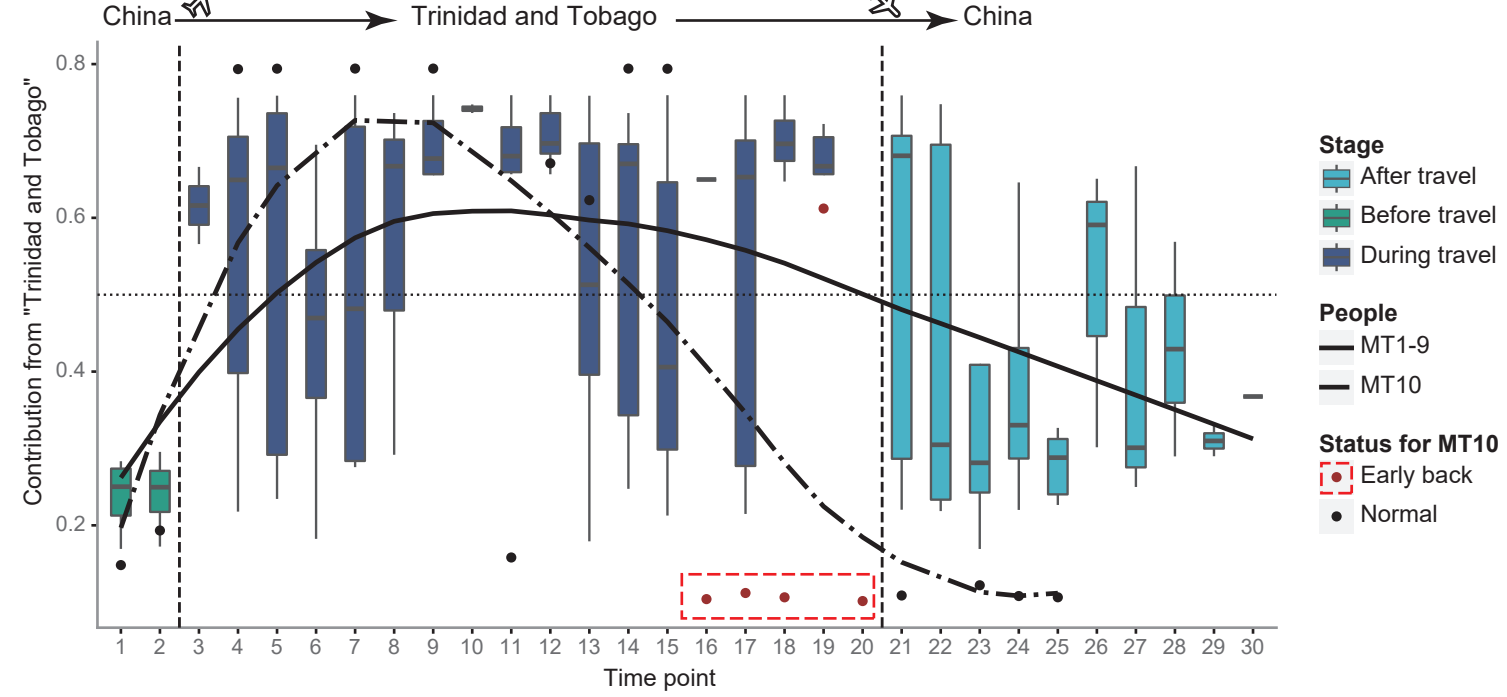


Could be inferred based on the prediction and additional metadata

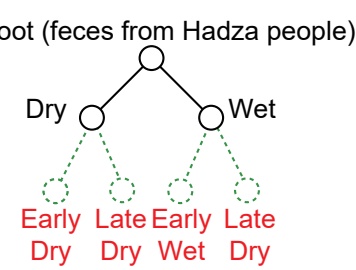
b. Sampling along the timeline



c. Estimated contributions over the timeline



d. Biome ontology
for seasonal pattern analysis



Could be inferred based on the prediction and additional metadata

e. Estimated contributions over the timeline

