1 **Ontology-Aware Deep Learning Enables Ultrafast, Accurate and Interpretable**

2 **Source Tracking among Sub-Million Microbial Community Samples from**

3 **Hundreds of Niches**

4 Yuguo Zha[1,$], Hui Chong[1,$], Hao Qiu[1], Kai Kang[1], Yuzheng Dun[2], Zhixue Chen[3,4], Xuefeng

5 Cui[3,4,*], Kang Ning[1,*]

6 [1]Key Laboratory of Molecular Biophysics of the Ministry of Education, Hubei Key Laboratory of

7 Bioinformatics and Molecular-imaging, Center of AI Biology, Department of Bioinformatics and

8 Systems Biology, College of Life Science and Technology, Huazhong University of Science and

9 Technology, Wuhan 430074, Hubei, China

10 [2]School of Mathematics and Statistics, Huazhong University of Science and Technology,   Wuhan

11 430074, Hubei, China

12 [3]School of Computer Science and Technology, Shandong University, Qingdao 250101, Shandong,

13 China

14 [4]Institute for Interdisciplinary Information Sciences, Tsinghua University, Beijing 100084, China

15 [$]These authors contribute equally to this work

16 [*]Correspondence should be addressed to K.N (Email: ningkang@hust.edu.cn) and X.C (Email:

17 xfcui@email.sdu.edu.cn)

18

## Abstract

20 The taxonomical structure of microbial community sample is highly habitat-specific,

21 making it possible for source tracking niches where samples are originated. Current

22 methods face challenges when the number of samples and niches are magnitudes

23 more than current in use, under which circumstances they are unable to accurately

24 source track samples in a timely manner, rendering them difficult in knowledge

25 discovery from sub-million heterogeneous samples. Here, we introduce a deep

26  learning method based on Ontology-aware Neural Network approach, ONN4MST

27  (https://github.com/HUST-NingKang-Lab/ONN4MST),  which  takes  into

28  consideration the ontology structure of niches and the relationship of samples from

29  these ontologically-organized niches. ONN4MST's superiority in accuracy, speed and

30  robustness have been proven, for example with an accuracy of 0.99 and AUC of 0.97

31  in a microbial source tracking experiment that 125,823 samples and 114 niches were

32  involved. Moreover, ONN4MST has been utilized on several source tracking

33  applications, showing that it could provide highly-interpretable results from samples

34  with previously less-studied niches, detect microbial contaminants, and identify

35  similar samples from ontologically-remote niches, with high fidelity.

36  **Keywords**: Ontology-aware Neural Network (ONN), microbial source tracking

37  (MST), deep learning, ultrafast, niches

38

39

40

## 41  Introduction

42  With the rapid accumulation of microbial community samples from various niches

43  (biomes) around the world, as well as the huge volume of sequencing data deposited

44  into public databases, such as those from the "Human Microbiome Project"[1,2] and the

45  "Earth Microbiome Project"[3,4], knowledge about microbial communities and their

46  influence on environment and human health has grown rapidly[5,6]. Such massive

47  amount of microbial community samples provide the opportunity to study the

48  inconspicuous evolution and ecological patterns among microbial communities,

49  especially habitat-specific patterns.

50

51  One key challenge faced with such a paramount number of heterogeneous samples is

52  to track potential origin of a microbial community sample, as well as distinguishing

53  samples from different health conditions or diverse environments, calling for fast and

54      accurate source tracking[7-9]. Taxonomical composition of a microbial community

55      sample is usually represented by hierarchically-structured taxa and their relative

56      abundances (also referred to as the community structure), and these taxa are

57      functioning in concert to maintain the stability of the microbial community and its

58      adaptation to the specific environment (also referred to as the niche or biome)[10,11].

59      Biomes are organized in an ontology structure with six different layers (simply

60      referred to as the biome ontology). Layer one is the highest layer containing only one

61      biome "Root" and layer six is the lowest (bottom) layer containing biomes such as

62      "Fecal". Biomes of lower layers such as "Human gut" belong to those of higher layers

63      such as "Human digestive system", whereby EBI MGnify currently contain the most

64      up-to-date biome structure[11] with more than one hundred biomes as of January 2020.

65      In general, microbial community samples from the same biome tend to have similar

66      community structures, while such similarities are highly dependent on the biome

67      layers. Source tracking the microbial community samples, especially among a

68      massive amount of samples, remains a challenging problem today.

69

70      Several methods for microbial community source tracking have already been

71      proposed[9,12-15]. They can generally be divided into two categories, namely

72      distance-based methods such as Jensen-Shannon Divergence (JSD)[16], Striped

73      UniFrac[13] and Meta-Prism[17], and unsupervised machine learning methods such as

74      SourceTracker based on Bayesian algorithm[15] and FEAST based on

75      Expected-Maximization algorithm[9]. However, the limitations of these methods are

76      apparent: Firstly, due to the nature of distance-based method and unsupervised

77      method, they are relatively slow, especially when the number of samples exceed tens

78      of thousands[9], hindering them from identifying potential source environments in a

79      timely manner. Secondly, there is still a lack of method for accurate source tracking

80      from more than a hundred biomes, largely due to the resolution limitation of both

81      methods[9,15]. Thirdly, current methods are not suitable for knowledge discovery of

82      samples from previously less studied or unknown biomes.

83

84 To address these limitations, we developed ONN4MST, an Ontology-aware Neural

85 Network (ONN) computational model for microbial source tracking. It is a supervised

86 learning method, and has utilized the biome ontology information. It has provided an

87 ultrafast (less than 0.1 seconds) and accurate (AUC higher than 0.97 in most cases)

88 solution for searching a sample against dataset containing more than a hundred

89 potential biomes and sub-million samples, and also out-performed state-of-the-art

90 methods in scalability and stability. The ability of ONN4MST on knowledge

91 discovery is also demonstrated by utilization in various source tracking applications:

92 it enables source tracking of samples whose niches are previously less studied or

93 unknown, detection of microbial contaminants, as well as identification of similar

94 samples from ontologically-remote biomes, showing the unique importance of

95 ONN4MST in knowledge discovery from huge amount of microbial community

96 samples of heterogeneous biomes.

97

## Results

99 **Ontology-aware Neural Network**

100 ONN4MST uses an Ontology-aware Neural Network (ONN) model for source

101 tracking. When training the model, all training samples' community structures are

102 decoded, each converted to a matrix containing the taxa at different taxonomical

103 levels and their relative abundances (simply referred to as the Matrix). The ONN

104 model uses the Matrix as input and reshapes it into tensors which point to biomes at

105 every different layer of the biome ontology. To fit the structure of biome ontology, the

106 ONN model uses multiple ontology units, each belonging to one of the six specific

107 layers of biome ontology (**Fig. 1a**). The conceptual modules, the training procedure

108 and the evaluation procedure of the ONN model are illustrated in **Supplementary Fig.**

109 **1** and described in **Methods**.

110

111 The source tracking procedure of ONN4MST is illustrated in **Fig. 1b**. Since

112 ONN4MST is the first method available that could source track the samples at

113    different layers of biome ontology, the search scheme of ONN4MST is completely

114    different from other methods (**Fig. 1b**). While ONN4MST goes through the biome

115    ontology to find the best possible source along different layers, other methods such as

116    FEAST and SourceTracker treat all biomes as anarchically equal. The overall scheme

117    of building the ONN model and using ONN4MST for source tracking is illustrated in

118    **Supplementary Fig. 2**. Note that the contributions of every known biome would be

119    estimated by the ONN model respectively.

120

121    **General model enables accurate source tracking with high scalability and**

122    **stability**

123    We constructed five datasets, representing sample collections with different numbers

124    of biomes and samples, covering more than 100,000 real microbial community

125    samples (**Supplementary Tables 1 and 2**). These five datasets contain samples from

126    different niches including "Host_associated", "Environmental" and "Engineered" as

127    top biomes, which are representative of high-quality microbial community samples in

128    public resources (**Supplementary Table 2**, **Methods**). Since these five datasets were

129    designed to have varied complexities, each including different number of samples

130    from different number of biomes, they could serve well for the evaluation of

131    ONN4MST and other methods (**Fig. 2a**): The Combined dataset contains 125,823

132    samples and 114 biomes, which represents the largest datasets, as well as the largest

133    model (the general model), used in this study. The FEAST dataset contains only

134    10,270 samples and 3 biomes. While the Human dataset, Water dataset, Soil datasets

135    are respectively with moderate sample sizes (**Supplementary Tables 1**).

136

137    First and foremost, the performances of ONN4MST on all five datasets were

138    evaluated. Results showed that the predicted biomes by ONN4MST were very close

139    to the actual biomes, regardless of the datasets used for evaluation. For example,

140    ONN4MST could achieve an accuracy of 0.99 and AUC of 0.97 on searching the

141    Combined dataset with 125,823 samples from 114 biomes. When we applied

142    ONN4MST on Human, Soil, Water and FEAST datasets, the accuracy and AUC of

143    ONN4MST were also higher than 0.98 and 0.96 for these datasets (**Table 1**,

144    **Supplementary Fig. 3**).

145

146    ONN4MST based on selected features performed equally well or better than that

147    based on all features. There are 44,668 taxa (or features) in total used in ONN4MST,

148    while ONN4MST_FS (ONN4MST based on selected features) has utilized only 1,462

149    selected features (see **Methods** and **Supplementary Table 3**). Results showed that

150    based on 1,462 selected features, ONN4MST_FS could attain slightly higher accuracy

151    (0.997 vs. 0.995, on Combined dataset), AUC and $F_{max}$ compared to ONN4MST

152    using all features (**Table 1**, **Supplementary Fig. 3**), which means that there is a

153    certain degree of redundancy among all 44,668 features, and we can achieve the same

154    accuracy with just 1,462 features compared with that using all 44,668 features. These

155    results have emphasized the scalability and stability of the general model built based

156    on the Combined dataset, either based on using all features, or using selected features.

157

158    Furthermore, we evaluated the universality of the general model built based on the

159    Combined dataset, by applying it directly on the Human, Water, Soil, and FEAST

160    datasets. It was found that the source tracking by using the general model was

161    successful on those datasets which are composed of samples mostly from the

162    Combined dataset's samples (**Supplementary Table 4**, results on Human, Water, Soil

163    datasets). However, when we applied the general model on datasets in which most of

164    the samples were not previously observed in the general model or have more detailed

165    biome ontology compared to the biome ontology used in general model, the general

166    model would not perform well (**Supplementary Table 4**, results on FEAST dataset).

167    Besides, results showed that it was unsuccessful when we applied the human model

168    (the model built based on Human dataset) for source tracking on Soil and Water

169    datasets (**Supplementary Table 5**).

170

171    **Comparison of ONN4MST and other source tracking methods**

172    We then compared all six source tracking methods on all five datasets with different

173    complexities (**Fig. 2a**). Results on all five datasets were evaluated seperately (**Fig.**

174    **2b,d**). Among the four datasets excluding FEAST dataset, ONN4MST was superior to

175    other methods: ONN4MST reached an AUC of 0.97, while other methods only

176    reached a maximum of 0.89 (**Fig. 2d**). As for the FEAST dataset, ONN4MST reached

177    an AUC of 0.99, while other methods only reached a maximum of 0.96.

178

179    The performances of these methods on five datasets depend on the datasets'

180    complexities (**Fig. 2c**). While Soil dataset and Water dataset are among those with the

181    highest Shannon diversity, the AUCs on these two datasets are also lower than those

182    on Human dataset and Combined dataset. The high AUC on FEAST dataset is largely

183    due to the small number of biomes used in FEAST dataset (**Supplementary Table 1**).

184    On the other hand, the performance of ONN4MST on each dataset did not depend

185    heavily on the number of samples in that dataset (provided that there are at least

186    10,000 samples in the dataset) (**Fig. 2c**, **Supplementary Table 1**). Furthermore, the

187    prediction accuracies were not biased for certain biomes (provided that there are at

188    least 100 samples in each biome) (**Supplementary Table 6**).

189

190    We further analyzed ONN4MST's performances at different biome layers (**Fig. 2e,f**).

191    Since it is the only method available that could source track samples at different

192    layers of biome ontology, we have remolded other methods' search scheme into a

193    hierarchical prediction scheme (see **Methods**), so that their results are comparable to

194    ONN4MST's. Results have clearly shown that ONN4MST and ONN4MST_FS

195    reached an AUC of 0.97 in minimum at all layers for the Combined dataset and these

196    were noticeably superior to other methods (**Fig. 2e,f**). Thus, ONN4MST is not just the

197    only method available that could source track the samples at different layers, but also

198    the best method even when other methods were remoulded for such purpose.

199

200    **Running time and memory utilization benchmark**

201    We evaluated the time and memory cost of all methods using a computational

202    platform comprising Quadruplex E7-4809 v3 CPU with 315 GB RAM, Nvidia Tesla

203   K80 GPU with 12 GB RAM. For time cost comparison, all actual times (search time,

204   excluding I/O time) were converted to the equivalent time on a single core.

205

206   ONN4MST is superior to other methods in search time and memory utilization where

207   the superiority expands as the number of source samples increases (**Fig. 3**). First of all,

208   we tested the time cost by searching a single query against the five datasets

209   respectively. For the Combined dataset including 125,823 source samples,

210   ONN4MST and ONN4MST_FS took 0.18 seconds and 0.04 seconds, respectively,

211   while distance-based methods took at least 1 second for a query. And FEAST took

212   more than 100,000 seconds, and SourceTracker took even more time (**Fig. 3a**, on the

213   Combined dataset, as also verified in Shenhav *et al.*[9]). Interestingly, though the time

214   spent by FEAST and Source Tracker per thousand of source samples were both less

215   than those reported in Shenhav *et*

216   *al.*[9], these two methods costed magnitudes more time than ONN4MST     (**Fig.    3a**).

217   When we linearly extrapolated the number of source samples to one million in the

218   dataset to be searched, the advantage of ONN4MST over other methods still held (**Fig.**

219   **3a**, hollow bars). When searching different number of queries against the Combined

220   dataset, we observed the time cost follows this trend: supervised methods

221   (ONN4MST and ONN4MST_FS) ≤ distance-based methods (JSD, Meta-Prism and

222   Striped UniFrac) < unsupervised methods (FEAST and SourceTracker) (**Fig. 3b**).

223   Again, when we linearly extrapolated the number of queries to one million in a batch,

224   the advantage of ONN4MST over other methods still held (**Fig. 3b**, hollow bars).

225

226   When memory utilization was evaluated, we have also observed the superiority of

227   ONN4MST over most of the other methods. Specifically, when searching a single

228   query against the Combined dataset, ONN4MST and ONN4MST_FS needed 22 GB

229   and 2 GB of memory, respectively; while FEAST and SourceTracker needed 84 GB

230   and 18 GB of memory, respectively; and JSD needed 47 GB of memory. Striped

231   UniFrac and Meta-Prism (https://github.com/HUST-NingKang-Lab/Meta-Prism-2.0)

232   were comparable with ONN4MST_FS in memory utilization, since they have

233    optimized the data structure for sample comparison. When the number of queries in a

234    batch exceeded 10,000, or the size of dataset to be searched varies, ONN4MST and

235    ONN4MST_FS remain the ones that needed the least memory (**Fig. 3c,d**). Details

236    about running time and memory utilization are presented in **Supplementary Tables**

237    **7-10**.

238

239    **Utility of ONN4MST in various source tracking applications**

240    The objective of microbial community sample source tracking is knowledge discovery

241    from the huge amount of microbial community samples of heterogeneous sources.

242    Thus, we showcased the ability of ONN4MST in knowledge discovery from several

243    perspectives: firstly, it can ensure accurate and interpretable source tracking, even on

244    distinguishing samples from ontologically-close biomes; secondly, when samples'

245    biomes are previously less studied or unknown, ONN4MST could provide accurate

246    clues for possible biome at higher layers, supplementing the information about such

247    less-studies biome; thirdly, ONN4MST could help for accurate microbial contaminant

248    detection; finally, "open search" of sample among the source samples with almost all

249    possible biomes could identify similar samples from ontologically-remote biomes,

250    leading to novel knowledge discovery.

251

252    **Centenarians share similar gut microbiota with young individuals**

253    ONN4MST can distinguish samples from ontologically-close biomes, thus offers a

254    quantitative way to characterize the development of human gut microbial community.

255    In this context, we leveraged external sources of young individuals (30 years old on

256    average) to understand the unique properties of gut microbiota in centenarians

257    (persons over 100 years old). To demonstrate this capability, we first built a

258    self-defined ONN model with two layers of biome ontology: "human gut" as first

259    layer, while "Young human gut" and "Others or unknown" at second layer, through

260    using a training set which contains 5,000 randomly selected human gut samples from

261    the Combined dataset (**Supplementary Table 1**), together with 800 randomly selected

262    human gut samples from young individuals in published studies[18,19]. Then, samples

263    from centenarians (30 from Italy, and 51 from China)[18,19] were used as queries for

264    performing source tracking with the self-defined ONN model. Results revealed a

265    significantly larger "Young human gut" contribution (Wilcoxon-test, $p < 1e-3$) in

266    centenarians (**Supplementary Fig. 4**), regardless of the locations where these samples

267    were collected, which were consistent with the results of published studies[18,19]. To

268    prove that these gut microbiota properties were unique in centenarians, we have

269    further collected 770 samples of normal seniors from another published study[20] as

270    queries for comparison. However, we could not observe the same phenomenon in

271    these normal seniors (**Supplementary Fig. 4**).

272

273    Several other case studies that distinguish samples from ontologically-close biomes

274    have also been conducted, with details in **Supplementary Note**, **Supplementary**

275    **Figures 5 and 6**.

276

277    **Detecting microbial contamination in built environment**

278    To validate ONN4MST's ability on microbial contamination detection, we analyzed

279    microbial community data collected by Lax *et al.*[21] In this analysis, we investigated

280    microbial contamination at several indoor house surfaces. We used skin samples from

281    several body parts (skin, foot, hand and nose) and additional environmental, plants

282    and mammal samples from the Combined dataset (**Supplementary Table 1**) as source

283    samples, and samples from indoor house surfaces ("Bathroom Door Knob", "Front

284    Door Knob", "Kitchen Counter", "Kitchen Floor" and "Kitchen Light Switch") as

285    queries. Our analysis results by using ONN4MST have shown that microbial

286    communities on these surfaces mostly originated from humans (**Fig. 4a**), largely in

287    agreement with the original analyses of Lax *et al.*[21] using SourceTracker, and differs

288    slightly from the results of Shenhav *et al.*[9] These results were reasonable considering

289    the strong influence of skin microbial communities on indoor house surfaces[22], while

290    they have again emphasized the challenge of source disambiguation for methods that

291    do not consider ontology structure of the biomes. That is, treating each individual

292    sample as an independent potential source would make differentiation of tiny sample

293 differences among ontologically-close biomes impossible, thus underestimating the

294 contributions of known sources at higher layers. We further investigated the

295 composition of the unknown sources existed in **Fig. 4a**. In addition to the contribution

296 of human, we found evidence for contributions from barley and bean product

297 (0.6-1.1%) and marine product (0.2-0.4%) for kitchen environments, and potential

298 evidence for contributions from agricultural (0.7-1.1%) and coastal (0.2-0.6%) for

299 door knobs (**Fig. 4b,c**).

300

**Source tracking of environmental samples from less studied biomes**

302 This investigation was based on searching 11 groundwater samples from another

303 published study[23] (the biome "Groundwater" is less studied, with a handful of samples

304 in the MGnify database, **Supplementary Table 2**) against the Combined dataset.

305 ONN4MST could successfully identify the actual biome for the majority of these

306 samples at different biome layers, such as "Aquatic" at the third layer and

307 "Freshwater" at the fourth layer (**Fig. 5a-c**) (results at the fifth and sixth layers were

308 shown in **Supplementary Fig. 7**). In contrast, FEAST and SourceTracker could not

309 identify any source near "Groundwater", while they only identified "Nutrient

310 (Wastewater)" with the meaning marginally related with groundwater (**Fig. 5d,e**).

311 Such differences in identification of actual biome are largely due to the fact that

312 ONN4MST could screen the whole biome ontology, and identify possible sources at

313 different layers, enabling it to at least identify the higher biome under which the

314 actual biome belongs to, with high fidelity. Whereas FEAST and SourceTracker were

315 designed without considering the biome ontology, they would assign "Unknown" for

316 many of these samples. These results indicated that when the actual biome of sample

317 was previously less studied, ONN4MST could provide accurate clues for possible

318 biome at higher layers in the biome ontology, and such clues would become valuable

319 assets in guiding the manual curation of these samples.

320

**Discovery of similar samples from ontologically-remote biomes**

322 Another advantage of ONN4MST in source tracking is its ability for "open search"

323  without any *a priori* knowledge about possible biomes where the query might be from,

324  enabling it for novel knowledge discovery. We tested ONN4MST's "open search"

325  results, and found that it could discover similar samples among ontologically-remote

326  biomes "Engineered", "Host_associated" and "Environmental" (**Supplementary**

327  **Table 11**). While some of the samples from the biome

328  "Root-Environmental-Aquatic-Marine-Intertidal_zone" share similar environments

329  (Baltic Sea) with the query sample from the biome

330  "Root-Engineered-Wastewater-Industrial_wastewater-Petrochemical", the literature

331  has also verified that this query sample was marine-sourced "MGYS00005175" (from

332  MGnify database). Such examples were plentiful (**Supplementary Fig. 8**), and many

333  had very high contributions ($> 0.8$). However, there were also examples which might

334  indicate possible mis-annotation or possible contaminations of samples in the MGnify

335  database. For instance, more than 10 samples from the study "MGYS00001610"

336  (from MGnify database) with annotated biome

337  "Root-Engineered-Wastewater-Water_and_sludge" have been identified by

338  ONN4MST as from biome

339  "Root-Host_associated-Mammals-Digestive_system-Large_intestine-Fecal"

340  (**Supplementary Fig. 8**). These results have verified our hypothesis that open search

341  of sample among the source samples with almost all possible biomes could reveal

342  remotely-similar samples, leading to novel knowledge that is never identified or

343  interpreted before.

344

## Discussion

346  ONN4MST was designed to address the urgent need for fast, accurate and

347  interpretable microbial community source tracking. It has been built based on an

348  Ontology-aware Neural Network model, which has provided a solution for source

349  tracking among sub-million samples and hundreds of biomes, outperforming

350  state-of-the-art methods, thus enabling knowledge discovery from these

351  heterogeneous samples. Microbial community sample source tracking has become

352   increasingly important, largely due to the needs of source tracking in multiple areas.

353   The requirements for high accuracy, high speed and high interpretability have thus

354   become critical considerations for a successful source tracking method, especially

355   when faced with the ever more complex situation where sub-million microbial

356   community samples from hundreds of biomes are provided as possible sources for

357   search.

358

359   The superiority of ONN4MST is established in several contexts. Firstly, ONN4MST

360   is very robust against dataset heterogeneity: from a dataset with the number of biomes

361   ranging from a handful to more than a hundred, as well as with the number of samples

362   ranging from a few thousand to sub-million, it always provides the highest accuracies

363   (AUC > 0.97) among state-of-the-art methods compared, making it the most scalable

364   source tracking method. Secondly, based on the Human, Water and Soil datasets, the

365   source tracking accuracies are all near-perfect (AUC > 0.97), indicating that

366   ONN4MST could provide reliable insights for downstream analysis on implicating

367   taxonomical or functional differences between healthy and diseased phenotypes, or on

368   illuminating tiny differences among environmental samples from even slightly

369   different niches. Furthermore, even when source tracking a sample against a database

370   of sub-million samples, only less than 0.1 seconds is needed when we conduct

371   ONN4MST search based on selected features, which is several orders of magnitude

372   faster than other contemporary methods. Finally, the ability of ONN4MST for 'open

373   search', without any *a priori* knowledge about possible biomes where the query might

374   be from, enables it for interpretable knowledge discovery.

375

376   The advantage of ONN4MST over other state-of-the-art source tracking methods is

377   essentially dependent on two technical advancements: the deep learning model, and

378   the ontology structure. Though the currently ongoing shift towards supervised

379   learning methods is not surprising for the source tracking research, the superior

380   performance of ONN4MST over existing methods is still quite pronounced.

381   ONN4MST's advantage also stems from its consideration of the ontology structure of

382　the biomes: by embedding the ontology considerations into the ONN learning model,

383　ONN4MST naturally becomes suitable for solving the ontology relationships among

384　biomes.

385

386　ONN4MST is not without limitations. Most importantly, the accuracy of ONN4MST

387　is heavily dependent on the ONN model built based on existing biome ontology

388　information. If there comes a new biome ontology with more detailed biomes

389　involved (for example, if we need to refine the source tracking results to human gut

390　down, to differentiate niches such as adult's gut from infant's gut), or simply with

391　more biome relationships involved, then the ONN model should be re-trained for

392　accurate source tracking. Such biome ontology-wide scalability problem could

393　potentially be solved by Transfer Learning approaches.

394

395　In summary, ONN4MST is an ontology-aware deep learning method that has pushed

396　the envelope of microbial source tracking, enabling near-optimal accurate, ultrafast

397　and interpretable source tracking. ONN4MST has enabled in-depth pattern and

398　function discoveries among sub-million microbial community samples, allowing for

399　tracking the potential origin of microbial community with diverse niche background,

400　as well as distinguishing samples from different health conditions or diverse

401　environments. Thus, it could have a broader area of application, such as

402　contamination screening, novel or refined biome discovery, new functional

403　microbiome discovery, and even source tracking of biomes from which protein

404　sequences could be supplemented for computational protein 3D structure

405　prediction[24,25].

406

407　**References**

408　1　Turnbaugh, P. J. *et al.* The human microbiome project. *Nature* **449**, 804-810
409　　　(2007).
410　2　Proctor, L. M. *et al.* The Integrative Human Microbiome Project. *Nature* **569**,
411　　　641-648 (2019).
412　3　Gilbert, J. A., Jansson, J. K. & Knight, R. The Earth Microbiome project:

413    successes and aspirations. *BMC Biol* **12**, 69-69 (2014).

414  4   Thompson, L. R. *et al.* A communal catalogue reveals Earth's multiscale
415    microbial diversity. *Nature* **551**, 457-463 (2017).

416  5   Dominguez-Bello, M. G. *et al.* Partial restoration of the microbiota of
417    cesarean-born infants via vaginal microbial transfer. *Nat Med* **22**, 250-253
418    (2016).

419  6   Thomas, S. *et al.* The Host Microbiome Regulates and Maintains Human
420    Health: A Primer and Perspective for Non-Microbiologists. *Cancer Res* **77**,
421    1783-1812 (2017).

422  7   Lladó, S., López-Mondéjar, R. & Baldrian, P. Drivers of microbial community
423    structure in forest soils. *Applied Microbiology and Biotechnology* **102**,
424    4331-4338 (2018).

425  8   Grond, K., Guilani, H. & Hird, S. M. Spatial heterogeneity of the shorebird
426    gastrointestinal microbiome. *R Soc Open Sci* **7**, 191609-191609 (2020).

427  9   Shenhav, L. *et al.* FEAST: fast expectation-maximization for microbial source
428    tracking. *Nature Methods* **16**, 627-632 (2019).

429  10  Tokeshi, M. Species Abundance Patterns and Community Structure. *advances
430    in ecological research* **24**, 111-186 (1993).

431  11  Mitchell, A. L. *et al.* MGnify: the microbiome analysis resource in 2020.
432    *Nucleic Acids Research* **48**, D570-D578 (2019).

433  12  Simpson, J. M., Santo Domingo, J. W. & Reasoner, D. J. Microbial Source
434    Tracking: State of the Science. *Environmental Science & Technology* **36**,
435    5279-5288 (2002).

436  13  Lozupone, C. & Knight, R. UniFrac: a new phylogenetic method for
437    comparing microbial communities. *Appl Environ Microbiol* **71**, 8228-8235
438    (2005).

439  14  Smith, A., Sterba-Boatwright, B. & Mott, J. Novel application of a statistical
440    technique, Random Forests, in a bacterial source tracking study. *Water
441    research* **44**, 4067-4076 (2010).

442  15  Knights, D. *et al.* Bayesian community-wide culture-independent microbial
443    source tracking. *Nature methods* **8**, 761-763 (2011).

444  16  Lin, J. Divergence measures based on the Shannon entropy. *IEEE
445    Transactions on Information Theory* **37**, 145-151 (1991).

446  17  Zhu, M., Kang, K. & Ning, K. Meta-Prism: Ultra-fast and highly accurate
447    microbial community structure search utilizing dual indexing and parallel
448    computation. *Briefings in bioinformatics* (2020).

449  18  Bian, G. *et al.* The Gut Microbiota of Healthy Aged Chinese Is Similar to That
450    of the Healthy Young. *mSphere* **2**, e00327-00317 (2017).

451  19  Biagi, E. *et al.* Through ageing, and beyond: gut microbiota and inflammatory
452    status in seniors and centenarians. *PLoS One* **5**, e10667-e10667 (2010).

453  20  Jeffery, I. B., Lynch, D. B. & O'Toole, P. W. Composition and temporal
454    stability of the gut microbiota in older persons. *The ISME Journal* **10**, 170-182
455    (2016).

456  21  Lax, S. *et al.* Longitudinal analysis of microbial interaction between humans

457       and the indoor environment. *Science* **345**, 1048-1052 (2014).

458   22   Timmis, K., Jebok, F., Rohde, M. & Molinari, G. Microbiome Yarns:
459       microbiome of the built environment, paranormal microbiology, and the power
460       of single cell genomics1,2,3,4. *Microb Biotechnol* **11**, 575-587 (2018).

461   23   Alsalah, D., Al-Jassim, N., Timraz, K. & Hong, P.-Y. Assessing the
462       Groundwater Quality at a Saudi Arabian Agricultural Site and the Occurrence
463       of Opportunistic Pathogens on Irrigated Food Produce. *Int J Environ Res*
464       *Public Health* **12**, 12391-12411 (2015).

465   24   Ovchinnikov, S. *et al.* Protein structure determination using metagenome
466       sequence data. *Science* **355**, 294-298 (2017).

467   25   Wang, Y. *et al.* Fueling ab initio folding with marine metagenomics enables
468       structure and function predictions of new protein families. *Genome Biol* **20**,
469       229-229 (2019).

470

# Methods

**Datasets**

We evaluated the performances of ONN4MST and other source tracking methods based on five different datasets (**Supplementary Table 1**). These five datasets comprise samples from different niches, which are representative of high-quality samples in public resources.

The "Combined dataset" consists of 125,823 microbial community samples collected from EBI MGnify database (https://www.ebi.ac.uk/metagenomics/), accessed as of January 2020 (**Supplementary Table 1**). This is a comprehensive dataset containing samples from 114 biomes (**Supplementary Table 2**), and the 125,823 microbial community samples represent more than half of the samples in EBI MGnify (as of January 1st, 2020). These samples contain taxonomical information for 225 phyla, 6,232 families, 16,081 genera and 45,477 species.

The "Human dataset" consists of 53,553 microbial community samples selected from the Combined dataset, representing a subset of samples from the human niches (**Supplementary Table 1**). Specifically, these samples are collected under these biomes: "Root-Host_associated-Human-Skin", "Root-Host_associated-Human-Circulatory_system", "Root-Host_associated-Human-Digestive_system" and "Root-Host_associated-Human-Reproductive_system" (biomes at higher layer). This dataset contains 53,553 samples from a total of 25 biomes. These samples contain taxonomical information for 204 phyla, 2,801 families, 6,523 genera and 16,135 species.

The "Water dataset" consists of 27,667 microbial community samples selected from the Combined dataset, representing a subset of samples from the water niches (**Supplementary Table 1**). Specifically, these samples are collected under these

500     biomes:                                "Root-Environmental-Aquatic-Freshwater",

501     "Root-Environmental-Aquatic-Marine"                        and

502     "Root-Environmental-Aquatic-Non-marine_Saline_and_Alkaline" (biomes at higher

503     layer). This dataset contains 27,667 samples from a total of 44 biomes. These samples

504     contain taxonomical information for 222 phyla, 6,040 families, 15,261 genera and

505     36,406 species.

506

507     The "Soil dataset" consists of 11,528 microbial community samples selected from the

508     Combined dataset, representing a subset of samples from the soil niches

509     (**Supplementary Table 1**). Specifically, these samples are collected under these

510     biomes:                    "Root-Environmental-Terrestrial-Soil",           and

511     "Root-Host_associated-Plants-Rhizosphere" (biomes at higher layer). This dataset

512     contains 11,528 samples from a total of 16 biomes. These samples contain

513     taxonomical information for 201 phyla, 2,962 families, 6,753 genera and 12,769

514     species.

515

516     These three datasets (Human, Water and Soil datasets) were designed with several

517     reasons in consideration. Firstly, these three datasets are representative enough and

518     frequently-used subsets[11] from the Combined dataset. Secondly, these three datasets

519     are also distinct, since the Alpha diversity of samples from each of these datasets is

520     significantly different from the other two: while samples from soil niches are

521     considered more complicated, those from human and water niches are considered less

522     so. Finally, samples from these niches are more comprehensively explored than other

523     less studied niches, and they are of relatively higher quality of samples from these

524     three niches.

525

526     The "FEAST dataset" consists of 10,270 microbial community samples selected from

527     the datasets used in the Lax *et al.*[9] (**Supplementary Table 1**). Specifically, these

528     samples are all collected from three biomes ("Root-Host_associated-Human",

529     "Root-Host_associated-Human-Digestive_system-Large_intestine-Fecal"        and

530   "Root-Mixed"). These samples contain taxonomical information for 133 phyla, 1,118

531   families, 3,389 genera and 5,762 species. The "FEAST dataset" is the smallest dataset

532   used in this study, and it is the simplest dataset with regard to the number of biomes

533   involved. Yet it is a dataset of unique importance, as the source tracking methods

534   evaluated in this study could be benchmarked on this medium-sized and credible

535   human gut dataset[9,15] for fair assessment of accuracy and efficiency.

536

537   **Data representation**

538   we generated the Matrix for each microbial community sample, so that the

539   abundances for all taxa at seven taxonomical levels including super-kingdom,

540   kingdom, phylum, class, order, family, and genus (simply referred to as "sk", "k", "p",

541   "c", "o", "f", and "g") can be retained. The abundance of taxa at different levels were

542   filled in the Matrix (**Figure 1**). Within the Matrix, seven columns respectively

543   represent seven taxonomical levels. And 44,668 rows respectively represent relative

544   abundance for 44,668 taxa (also referred to as features). For a detailed description and

545   an example of the data representation, see **Supplementary Note** and **Supplementary**

546   **Table 3**.

547

548   **Feature selection**

549   To improve the efficiency and accuracy of ONN4MST, we conducted feature

550   selection by using a random forest regression model (Python-3.7.4 and

551   Scikit-learn-0.22.1). An abundance-based pre-filtering and an importance-based

552   selection were performed in sequential order. In doing so, we treated each row

553   (representing the abundances of a taxon, see **Supplementary Table 3**) of the Matrix

554   as a feature. Then, a series of adaptive thresholds ($C\overline{R}_l$ and $C\overline{I}_l$) were applied to

555   different taxon levels, in which $\overline{R}_l$ and $\overline{I}_l$ stand respectively for the relative

556   abundance and the feature importance. $level \in \{sk, k, p, c, o, f, g\}$ and the

557   coefficient $C$ was set to $0.001$. As a result, we have selected 1,462 features with

558   relative abundance and feature importance above the thresholds from all 44,668

559   features involved in this study.

560

**Biome ontology**

We constructed a comprehensive biome ontology using 114 biomes (**Supplementary Table 2**) collected from EBI MGnify database (https://www.ebi.ac.uk/metagenomics/biomes). In this process, we organized the biome ontology as a tree, by treating a biome with multiple parent biomes in the higher layer (e.g. "Human-Digestive_system" and "Mammal-Digestive_system") as seperate biomes. Next, the ontology tree containing 6 layers and 133 nodes (representing 114 biomes) was constructed, by using Python-3.7.4 and Treelib-1.5.5. As a result, each biome was represented by at least one node in the ontology tree. The ontology tree has "Root" at the first layer, biomes (nodes) including "Environmental", "Host_associated", and "Engineered" at the second layer, and 7, 22, and 56 biomes (nodes) at the third to fifth layers respectively, with 43 biomes (nodes) including "Coral reef", "Fecal" and "Saliva" at the bottom (sixth) layer (**Supplementary Table 2).**

**Sample Labeling**

In all experiments, we used microbial samples each with a label annotated by using 6-layers biome ontology to validate our model. For example, there are 22 samples labeled as "Root-Host_associated-Human-Digestive_system-Oral-Throat" in the Combined dataset (by separating different layers with the "-" symbol).

**Building ONN model**

We used Tensorflow-1.14[26] to build and train our Ontology-aware Neural Network model. Our model was trained on a computational platform comprising Quadruplex E7-4809 v3 CPU with 315 GB RAM and Nvidia Tesla K80 GPU with 12 GB RAM.

Ontology-aware Neural Network has four conceptual modules in total: a feature extraction module for basic feature extraction, a feature encoding module for layer-specific feature encoding, a feature integration module for inter-layer

590    information integration, and an ontology prediction module for ontology walk through

591    and source contribution calculation (**Supplementary Fig. 1a**). The feature extraction

592    module accepts a sample represented by the Matrix, extracts the feature information

593    from the Matrix and deliver them to the feature encoding module. The feature

594    encoding module consists of a series of fully-connected layers. It accepts the output of

595    feature extraction module, and encodes layer-specific feature information for each of

596    the six biome ontology layers. The feature integration module consists of several

597    fully-connected layers, which serves for inter-layer information integration. The

598    ontology prediction module consists of five sigmoid layers (corresponding to the $2^{nd}$,

599    $3^{rd}$, $4^{th}$, $5^{th}$ and $6^{th}$ biome ontology layers), each sigmoid layer accepts the output of

600    feature encoding module and computes the contribution of all biome sources on its

601    corresponding biome ontology layer.

602

603    We chose 8-fold cross validation for model training and testing (**Supplementary Fig.**

604    **1c**). For each dataset, we randomly split it into 8 folds, each fold including a training

605    set (87.5%) and a testing set (12.5%). For each fold, the model was trained (in batches

606    of 512 samples) for 30,000 iterations or until training accuracy converged, and the

607    model with the highest accuracy on the training set was selected for testing. The

608    results on the testing set are organized in the form of a hierarchical prediction (with

609    prediction results from $2^{nd}$ to $6^{th}$ layers), which would then be evaluated.

610

611    **Other methods used in this study**

612    Three distance-based methods: JSD, Striped UniFrac and Meta-Prism, two

613    unsupervised machine learning methods: Expected-Maximization based method

614    FEAST and Bayesian based method SourceTracker; as well as our supervised deep

615    learning method (ONN4MST), were applied for microbial source tracking. In this

616    study, the source tracking results (predicted biomes) of multiple methods were

617    compared against the microbial community samples' actual source (actual biomes).

618

619     The distance-based methods are based on pair-wise calculation of sample distances,

620     and such methods depend heavily on the presence of species and their relative

621     abundance for individual samples, regardless of weighted or unweighted scoring

622     functions used. Among distance-based methods, JSD does not consider the

623     phylogenetic relationships among species, while methods such as Striped UniFrac and

624     Meta-Prism do (we have used Meta-Prism 2.0 for comparison in this study). However,

625     distance-based methods have a binomial increase in time cost with the increase of the

626     number of samples.

627

628     Unsupervised methods for microbial community sample comparison are based on

629     profile-based statistical models, either the Bayesian model used in the SourceTracker

630     method, or the Expected-Maximization (EM) model used in the FEAST method.

631     Unsupervised methods are typically more accurate than distance-based methods.

632     However, since unsupervised methods still do not consider the intricate but important

633     patterns of a set of samples from similar niches, their tolerance to noisy signals in

634     samples is not high, hence potentially would lead to biased mismatches. Details about

635     the source tracking methods other than ONN4MST used in this study are provided in

636     **Supplementary Note**.

637

638     **Hierarchical prediction**

639     In order to carry out comparison of ONN4MST against other methods at different

640     layers of biome ontology, all other methods were remolded, so that the prediction

641     results of these methods (excluding ONN4MST) at different layers could be produced.

642     Based on the source contributions of biomes at the sixth (bottom) layer, the source

643     contributions of biomes for other layers were computed using $P_f = \sum_{f_c \in C_f} P_{f_c}$ . Where

644     $P_f$ is a source contribution for $f$, $C_f$ is a set of children biomes for biome source $f$

645     in the biome ontology. $f_c$ is a child biome of $f$. We used NumPy-1.18.1 and

646     Treelib-1.5.5 in the process.

647

648 **Benchmarking measures**

649 To benchmark and compare the results based on ONN4MST and the other five

650 methods, we used these measures:

651

$$TP_f(t) = \sum_i I(f \in P_i(t) \wedge f \in T_i) \tag{1}$$

$$TN_f(t) = \sum_i I(f \notin P_i(t) \wedge f \notin T_i) \tag{2}$$

$$FP_f(t) = \sum_i I(f \in P_i(t) \wedge f \notin T_i) \tag{3}$$

$$FN_f(t) = \sum_i I(f \notin P_i(t) \wedge f \in T_i) \tag{4}$$

$$TPR_f(t) = \frac{TP_f(t)}{TP_f(t) + FN_f(t)} \tag{5}$$

$$FPR_f(t) = \frac{FP_f(t)}{FP_f(t) + TN_f(t)} \tag{6}$$

$$TPR(t) = \frac{1}{F}\sum_{f=1}^{F} TPR_f(t) \tag{7}$$

$$FPR(t) = \frac{1}{F}\sum_{f=1}^{F} FPR_f(t) \tag{8}$$

660 where $f$ is a biome source, $P_i(t)$ is a set of predicted biomes for a microbial

661 community sample $i$ and threshold $t \in [0,1]$ with a step size of 0.01, $T_i$ is a set of

662 actual biomes for a sample $i$, $F$ is the total number of biomes, and $I$ is a logical

663 operation function, the value of $I$ is 1 when the result of logical operation is TRUE,

664 else 0.

665

666 Four evaluation metrics ($Accuracy$, $Precision$, $Recall$ and $F_{max}$) were introduced.

667 These evaluation metrics are computed with the following formulas:

$$Accuracy(t) = \frac{TP_f(t) + TN_f(t)}{TP_f(t) + FP_f(t) + TN_f(t) + FN_f(t)} \tag{9}$$

$$Precision_f(t) = \frac{TP_f(t)}{TP_f(t) + FP_f(t)} \tag{10}$$

$$Recall_f(t) = \frac{TP_f(t)}{TP_f(t) + FN_f(t)} \tag{11}$$

671 where $TP$ is true positive, $TN$ is true negative, $FP$ is false positive, $FN$ is false

672 negative. Subsequently, we compute $F1$ for threshold $t \in [0,1]$ with a step size of

673 0.01 by using the average precision and average recall for all actual biomes that we

674 predicted at least one time. Then, we select the maximum $F1$ as $F_{max}$. These

675 evaluation metrics are computed with the following formulas:

676
$$AvgPrecision(t) = \frac{1}{F}\sum_{f=1}^{F} Precision_f(t) \tag{12}$$

677
$$AvgRecall(t) = \frac{1}{F}\sum_{f=1}^{F} Recall_f(t) \tag{13}$$

678
$$F_{max} = max_t \left\{ \frac{2 \cdot AvgPrecision(t) \cdot AvgRecall(t)}{AvgPrecision(t) + AvgRecall(t)} \right\} \tag{14}$$

679

680 Then, ROC (Receiver Operating Characteristic) curves, which are based on

681 contrasting the true positive rate (TPR) against the false positive rate (FPR), were

682 plotted. AUC (Area Under the Curve) reflects the ability of model to correctly predict

683 the biomes (sources) of microbial community samples. AUC is calculated with the

684 following formula:

685
$$AUC = \int_0^1 TPR(t)\big(-FPR'(t)\big)dt \tag{15}$$

686

## Data availability

688 The selected samples from Combined dataset, which were assigned to Human dataset,

689 Water dataset, Soil dataset respectively, were annotated with their respective

690 assignments in **Supplementary Table 2**. Data download links are provided in

691 **Supplementary Table 12**.

692

## Code availability

694 All source codes have been uploaded to the website at:

695 https://github.com/HUST-NingKang-Lab/ONN4MST. Detailed parameters of

696 software and package we used in this study are provided in **Supplementary Table 13**.

697

## References

699 26    Abadi, M. *et al.* Tensorflow: a system for large-scale machine learning.
700       *Operating Systems Design and Implementation*, 265-283 (2016).

701

## Acknowledgments

## Author contributions

KN conceived of and proposed the idea, and designed the study. YGZ, HC, HQ, KK, YZD, ZXC performed the experiments and analyzed the data. YGZ, HC, KN and XC contributed to editing and proof-reading the manuscript. All authors read and approved the final manuscript.

## Competing interests

The authors declare that they have no competing interests.

## Ethics approval and consent to participate

Not applicable

## Figures

### Figure 1

**Fig. 1: Building and using the Ontology-aware Neural Network model for microbial source tracking. a.** The sample data representation and training process of ONN model. **i.** Sample data are transformed into the Matrix. With the Matrix, each column represents a taxonomical level and each row represents a feature; **ii.** In parallel, samples are mapped to biome ontology according to their niches; **iii.** The model is built and updated according to both samples' abundance matrices and biome ontology information. More details about building, testing and using the ONN model for source tracking are illustrated in **Supplementary Fig. 1** and **Supplementary Fig. 2**. **b.** An illustrated example of microbial source tracking procedure using ONN4MST. **i.** The input is the community structure of a real microbial community sample (this sample is from the biome "Root-Host_associated-Human-Digestive_system-Oral-Saliva") that has been preprocessed and the Matrix has been provided into the model; **ii.** Source tracking process at different layers. The red arrows indicate the search process from layer 1 to layer 6, accompanied with source contribution annotated in red. To compare with the procedure of ONN4MST, the yellow and blue arrows indicated the source tracking results (among the overall top 5 sources) of FEAST and Source Tracker, together with their source contributions, respectively. The actual biome is annotated by a red check mark; **iii.** The predicted biomes (with source contributions) by ONN4MST, FEAST and SourceTracker.

739  **Figure 2**

740  **Fig. 2: ONN4MST's prediction accuracies are among the best on different datasets and**

741  **different biome layers, while the performance of ONN4MST does not depend heavily on the**

742  **number of biomes or number of samples in the dataset. a.** The five datasets with varied

743  complexities have provided source tracking tasks with different difficulties. **b.** The ROC curve of

744  ONN4MST and other methods on all five datasets. **c.** The number of samples, the Shannon

745  diversity and the source tracking results by different methods for the five datasets. The samples

746  involved in each dataset are shown with blue bars, the Shannon diversity of each dataset is shown

747  with red boxes, the AUC of several methods on each dataset is shown with dash lines. **d.** The AUC

748  of all methods on all five datasets. **e.** The number of biomes and the source tracking results by

749  different methods at different layers for the Combined dataset. The samples involved in each

750  biome ontology layer are shown with blue bars, the AUC of different methods on each layer is

751  shown with dash lines. **f.** The AUC of all methods at different layers. (**Abbreviations**.

752  ONN4MST_FS: ONN4MST using selected features).

753 **Figure 3**

754 **Fig. 3: ONN4MST is superior to other methods in search time and memory utilization. a.**

755 Running time of different methods when search one query against different datasets. **b.** Running

756 time of different methods when search queries of different sizes against Combined dataset. **c.**

757 Memory utilization of all methods when search one query against different datasets. **d.** Memory

758 utilization of all methods when search queries of different sizes against Combined dataset. **Note**: a

759 hollow bar means that the value represent by this bar is the result of linearly extrapolation, both

760 for running time and for memory utilization. (**Abbreviations**. ONN4MST_FS: ONN4MST using

761 selected features, 1M: Results of linearly extrapolation with one million samples in use).

762 **Figure 4**

763 **Fig. 4: The contribution of the unknown sources in indoor house surface samples using**

764 **ONN4MST. a.** Mean source contributions considering 4 human skin sources (hand, foot, nose and

765 skin-other across all inhabitants) using data from Lax *et al.*[21] **b,c.**

766 Further decomposition of the unknown sources existed in **Fig. 4a** has revealed other microbial con

767 taminates in built environment.

768 **Figure 5**

769 **Fig. 5: Successful source tracking of environmental samples from a less studied biome by**

770 **using ONN4MST.** Results were based on using 11 samples from groundwater environment,

771 which represented a biome previously less studied. **a-c.** Source tracking results by using

772 ONN4MST at the second, third and fourth layers; **d.** Source tracking results by using FEAST; **e.**

773 Source tracking results by using SourceTracker. Actual biome of query sample:

774 "Root-Environmental-Aquatic-Freshwater-Groundwater". A_1, A_2: two samples collected from a

775 single well; B_1, B_2: two samples collected from another single well; C_1, C_2: two samples

776 collected from the third single well; D-H: samples collected from other five wells, respectively.

777 **Tables**

778 **Table 1**

779 **Table 1. Evaluation of ONN4MST on all five datasets.** ONN4MST achieved the accuracy

780 higher than 0.98 for all five datasets, and the AUC higher than 0.97 for all five datasets. **Note:** For

781 each dataset, we used the model trained on that dataset for evaluation. The evaluation procedure of

782 the ONN model is illustrated in **Supplementary Fig. 1c** and described in **Methods**. ONN4MST

783 based on all features and selected features were both evaluated at the bottom (sixth) layer with a

784 threshold of 0.5. (**Abbreviations**. Pr: Precision, Rc: Recall, Acc: Accuracy).

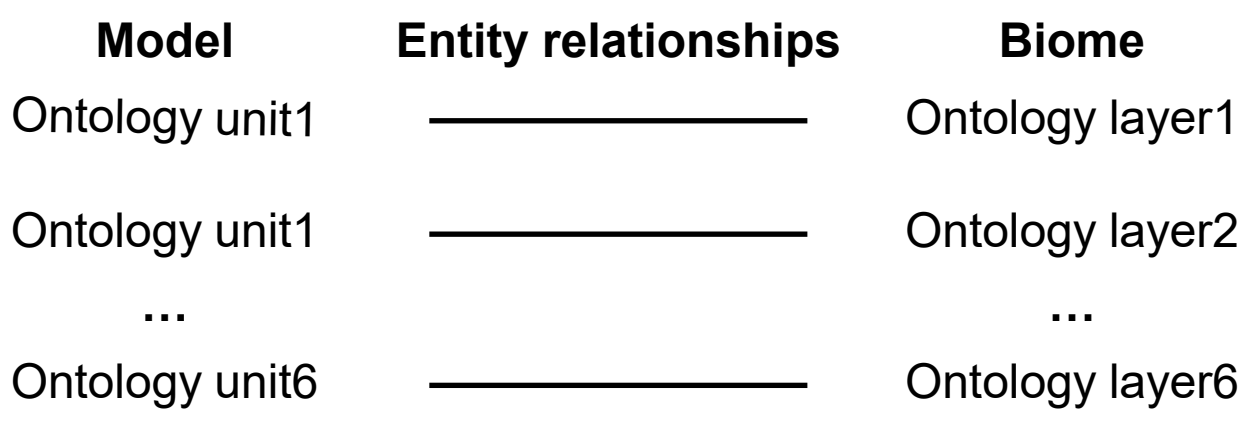| Dataset | #Biomes | #Samples | All features | | | | | Selected features | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | Pr | Rc | Acc | $F_{max}$ | AUC | Pr | Rc | Acc | $F_{max}$ | AUC |
| Combined | 114 | 125,823 | 0.826 | 0.662 | 0.995 | 0.740 | 0.971 | 0.868 | 0.774 | 0.997 | 0.820 | 0.977 |
| Human | 25 | 53,553 | 0.822 | 0.521 | 0.984 | 0.695 | 0.972 | 0.894 | 0.826 | 0.991 | 0.863 | 0.984 |
| Water | 44 | 27,667 | 0.842 | 0.766 | 0.992 | 0.803 | 0.966 | 0.854 | 0.764 | 0.992 | 0.813 | 0.971 |
| Soil | 16 | 11,528 | 0.915 | 0.778 | 0.986 | 0.850 | 0.974 | 0.892 | 0.881 | 0.989 | 0.890 | 0.982 |
| FEAST | 3 | 10,270 | 0.793 | 0.795 | 0.984 | 0.803 | 0.980 | 0.895 | 0.812 | 0.989 | 0.862 | 0.991 |

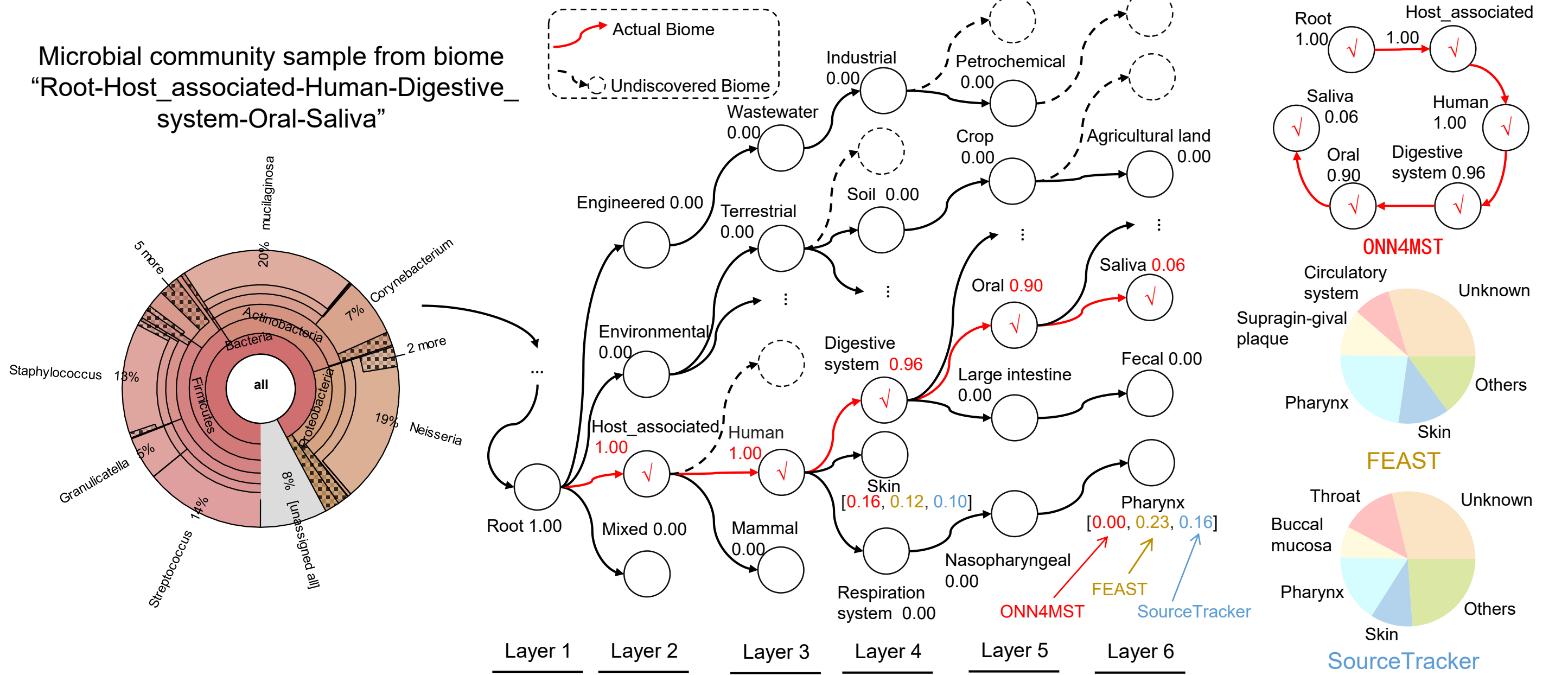785

**a**

**i. Sample data representation**
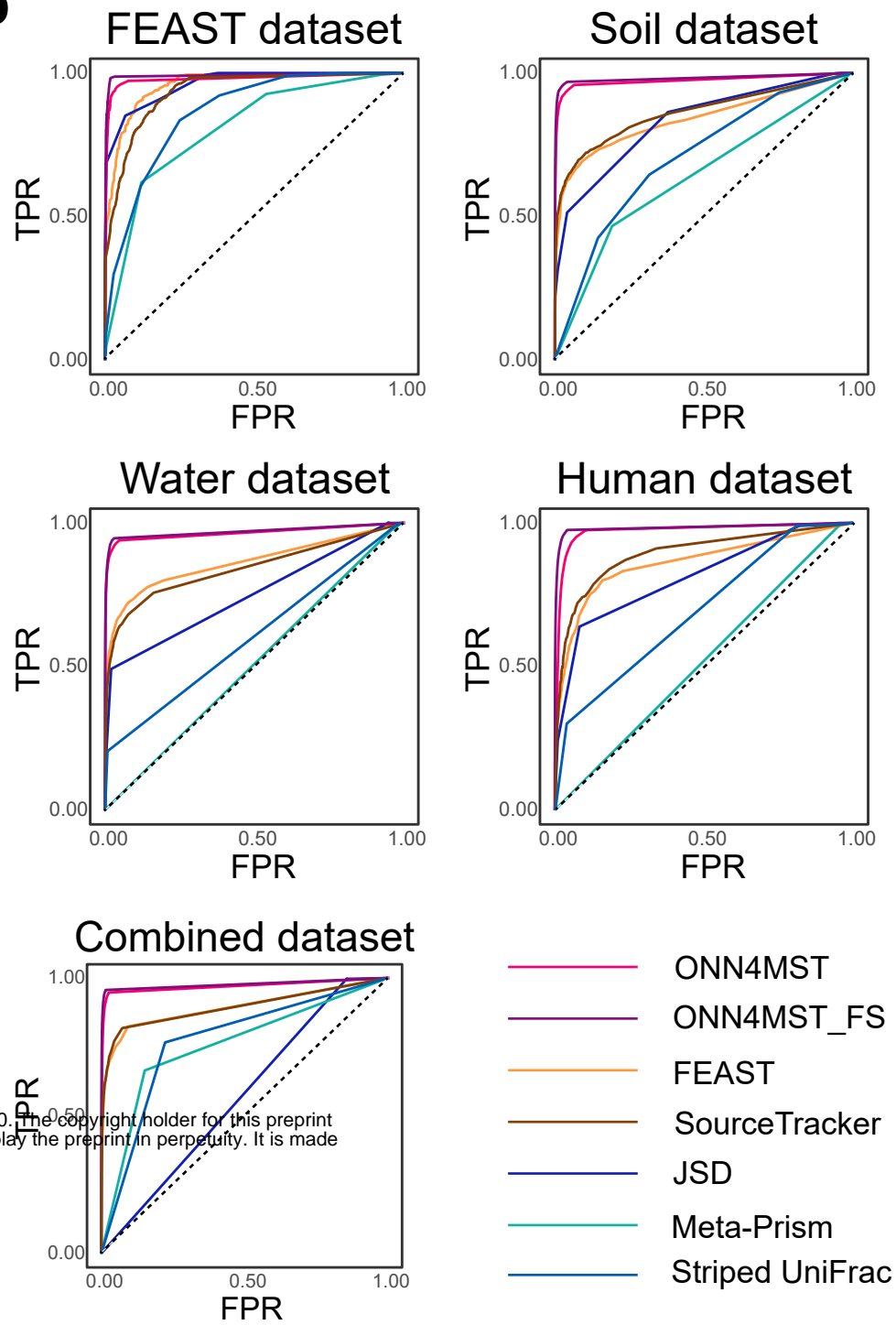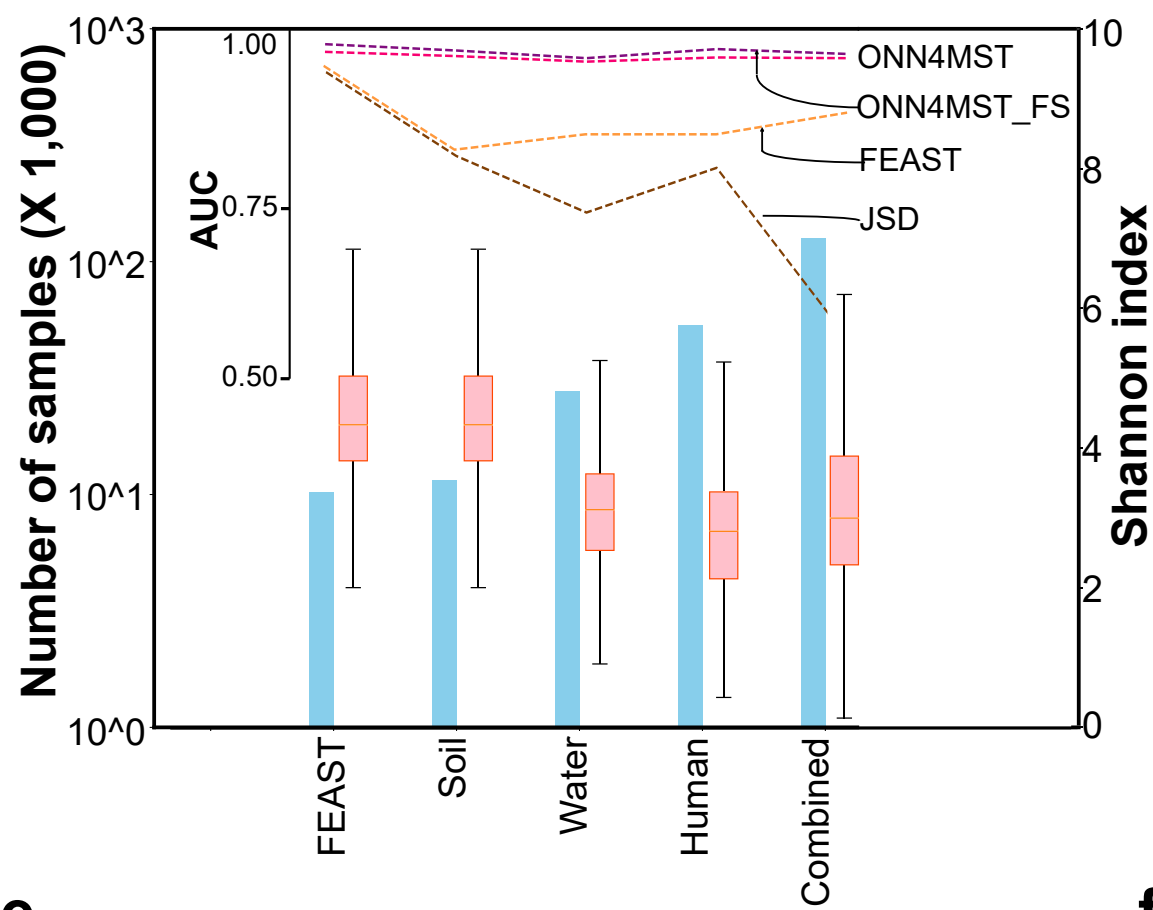
**ii. Biome ontology (with samples mapped)**
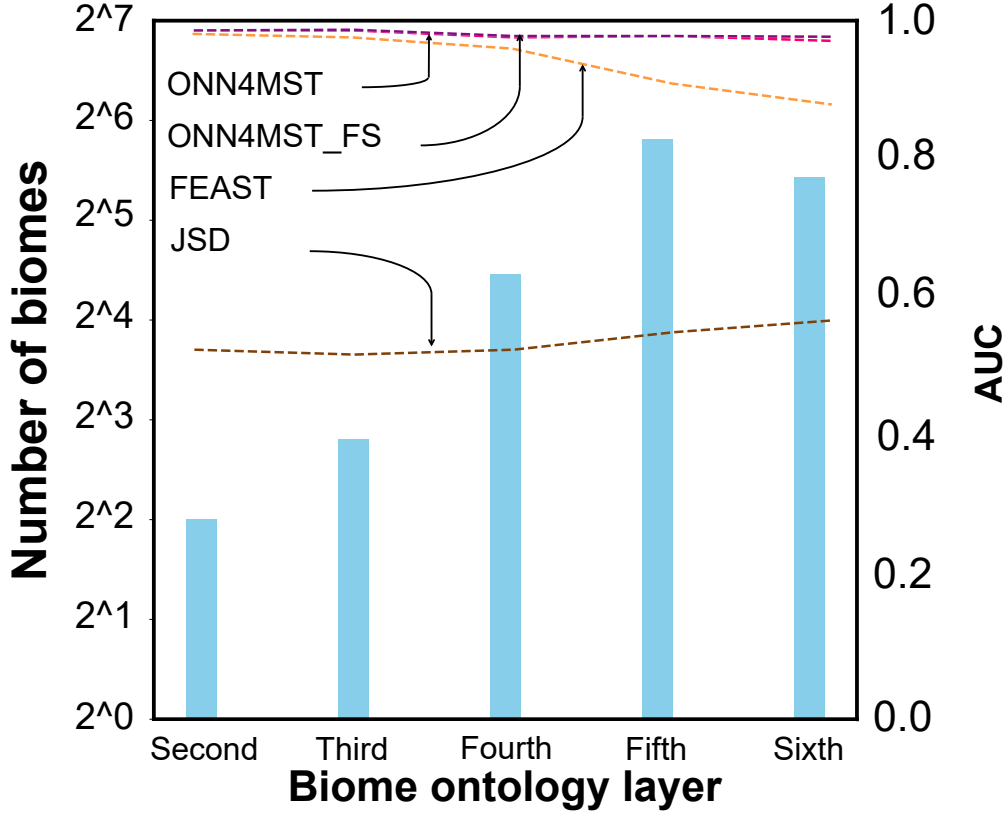
**iii. Ontology-aware Neural Network**

**b**

**i. Input**

**ii. Source tracking process**

**iii. Output**

**a**

Running time (s) vs Datasets to be searched (FEAST, Soil, Water, Human, Combined, 1M)

**b**

Running time (s) vs Number of queries (1, 100, 10,000, 1,000,000)

**c**

Memory utilization (GB) vs Datasets to be searched (FEAST, Soil, Water, Human, Combined, 1M)

**d**

Memory utilization (GB) vs Number of queries (1, 100, 10,000, 1,000,000)

ONN4MST   FEAST   JSD   Striped UniFrac
ONN4MST_FS   SourceTracker   Meta-Prism

**a** ONN4MST Layer-2  **b** ONN4MST Layer-3  **c** ONN4MST Layer-4

**d** FEAST  **e** SourceTracker

Legend a:
Unknown / Engineered / Environmental / Host_associated / Mixed / Others

Legend b:
Unknown / Aquatic / Terrestrial / Human / Plants / Others

Legend c:
Unknown / Freshwater / Marine / Soil / Rhizosphere / Others

Legend d:
Unknown / Industrial / Nutrient(Wastewater) / Contaminated(Soil) / Mixed / Others

Legend e:
Unknown / Petrochemical / Dissolved_organics_(aerobic) / Contaminated(Soil) / Mixed / Others

Actual biome: "Root-Environmental-Aquatic-Freshwater-Groundwater"

◆ Source: "Environmental"     ▼ Source: "Aquatic"

★ Source: "Freshwater"     ⬟ Source: "Nutrient (Wastewater)"

**a**

Feature Integration Module

I2 64 I3 64 I4 64 I5 64

Input 44668 x 7

B1 512

E2 128 E3 128 E4 128 E5 128 E6 128

Feature Encoding Module

Feature Extraction Module

P2 4 P3 7 P4 22 P5 56 P6 43

Ontology Prediction Module

**b**

Model training

Initial state   Intermediate states   End state

Layers
L1 L2 ⋯ L6

Biomes
B 1
B 2
⋮
B n

Low
Enrichment
High

Final model

**c**

Train   Test
Train   Test   Train
⋯
Test   Train

8-fold cross validation

"Root–Environmental"
0.14          0.05
"Root"          1.00   0.86   0.82
"Root-Host_associated"   0.13

Hierarchical predictions

TPR
FPR

Evaluation

**a** Microbial community structures

Samples from biome "Environmental"
Samples from biome "Host-associated"
Samples from biome "Mammals"
Samples from biome "Human-gut"
Training samples
Sample(s) from unknown source
Query sample(s)

**b** Ontology-aware profiles

kingdom
phylum
...
species
1 2 3 4 ...... n

kingdom
phylum
...
species
1 2 3 4 ...... n
Abundance matrix

**c** Building the ONN prediction model

"Human-gut"
"Host-associated"
"Environmental"
"Engineered"
"Mixed"

Biome ontology

Samples
Parameter updates

Biome ontology (with samples mapped)

**d** Source tracking model

L1 L2 ... L6
B 1
B 2
⋮
B n

Enrichment
Low
High

**e** Hierarchical predictions

| Host_associated | Engineered | Environmental | Mixed | Layer 2 |
|---|---|---|---|---|
| 0.95 | 0.02 | 0.01 | 0.01 | |

| Mammals | Plants | Waste-water | Terrestrial | Layer 3 |
|---|---|---|---|---|
| 0.92 | 0.03 | 0.03 | 0.02 | |

| Fecal | Cecum | Rumen | Throat | Layer 6 |
|---|---|---|---|---|
| 0.91 | 0.03 | 0.03 | 0.02 | |

Query sample(s) from "Root-Host_associated-Mammals- …… -Fecal"

**a** **ONN4MST**

**b** **ONN4MST_FS**

Dataset

— Combined
— Human
— Water
— Soil
— FEAST

Panel a (ONN4MST):
- AUC=0.971
- AUC=0.972
- AUC=0.966
- AUC=0.974
- AUC=0.980

Panel b (ONN4MST_FS):
- AUC=0.977
- AUC=0.984
- AUC=0.971
- AUC=0.982
- AUC=0.991

**a** Centenarians in Italy

p = 1.6e−15

#Samples = 30
(100 ≤ age ≤ 120)

**b** Centenarians in China

p < 2.22e−16

#Samples = 51
(100 ≤ age ≤ 120)

**c** Seniors

p < 2.22e−16

#Samples = 770
(age ≥ 64, age = 78 ± 8)

Source contribution

Young          Others or Unknown

**a** JSD distribution

**b**

**Possible remote similarities**

| Sink | & | Source |
|------|---|--------|
| Intestine | X | Skin |
| Intestine | X | Respiratory_system |
| Intestine | X | Reproductive_system |
| ...... | X | ...... |
| Oral | X | Circulatory_system |
| ...... | X | ...... |

**c**

| Method | ONN4MST | | | | FEAST |
|--------|---------|---|---|---|-------|
| Cutoff | Layer2 | Layer3 | Layer4 | Layer5 | |
| 40 | **0.977** | 0.963 | 0.963 | 0.716 | 0.597 |
| 70 | **0.957** | 0.913 | 0.923 | 0.583 | 0.350 |
| 90 | **0.933** | 0.867 | 0.830 | 0.403 | 0.150 |

**d**

Sample ID: MGYS00001248-SRR2761086
Actual biome:"Root-Host-associated-Human-Digestive_system-Large_intestine"

ONN4MST : {Layer2 | "Root-Host_associated": 0.999,    √
            Layer3 | "Root-Host_associated-Human": 0.999,    √
            Layer4 | "Root-Host_associated-Human-Digestive_system": 0.999,    √
            Layer5 | "Root-Host_associated-Human-Digestive_system-Large_intestine": 0.968    √

FEAST :    {"Root-Host-associated-Human-Skin": 0.163,
            "Unknown": 0.837}

ONN4MST Layer-5 | ONN4MST Layer-6

Legend (Layer-5): Unknown, Lake, Oceanic, Alkaline, Permafrost(Soil), Others

Legend (Layer-6): Unknown, Photic_zone(Oceanic), Sediment(Oceanic), Sediment(Alkaline), Agricultural_land, Others

**a**

Sink ID: MGYS00005175-SRR6319590     Represent Source ID: MGYS00002650-SRR3589592

Actual biome: Root-Engineered-Wastewater-Industrial_wastewater-Petrochemical     Predicted biome: Root-Environmental-Aquatic-Marine-Intertidal_zone

Planctomycetes: 2.36 %
Verrucomicrobia: 4.13 %
Unassigned Bacteria: 4.94 %
Bacteroidetes: 9.81 %
Actinobacteria: 10.80 %
Proteobacteria: 64.46 %

Other: 1.26 %
Verrucomicrobia: 1.45 %
Unassigned Bacteria: 2.63 %
Bacteroidetes: 17.26 %
Proteobacteria: 50.69 %
Actinobacteria: 23.92 %

**b**

Sink ID: MGYS00004521-SRR6901946     Represent Source ID: MGYS00002650-SRR3589534

Actual biome: Root-Engineered-Wastewater-Industrial_wastewater-Agricultural_wastewater     Predicted biome: Root-Environmental-Aquatic-Marine-Intertidal_zone

Verrucomicrobia: 1.48 %
Planctomycetes: 4.62 %
Actinobacteria: 5.05 %
Unassigned Bacteria: 6.69 %
Bacteroidetes: 7.54 %
Proteobacteria: 69.95%

Other: 1.12 %
Unassigned Bacteria: 2.82 %
Actinobacteria: 6.78 %
Bacteroidetes: 13.07 %
Proteobacteria: 74.42 %

**c**

Sink ID: MGYS00001610-ERR982889     Represent Source ID: MGYS00004714-ERR3258060

Actual biome: Root-Engineered-Wastewater-Water_and_sludge     Predicted biome: Root-Host_associated-Mammals-Digestive_system-Large_intestine-Fecal

Other: 3.99 %
Deltaproteobacteria: 1.72 %
Bacilli: 5.75 %
Bacteroidia: 11.06 %
Unassigned: 54.21 %
Clostridia: 18.85 %

Other: 2.47 %
Kiritimatiellaeota: 2.13 %
Proteobacteria: 2.52 %
Verrucomicrobia: 5.41 %
Firmicutes: 54.47 %
Bacteroidetes: 27.11 %