

1 **Ecological Dynamics Imposes Fundamental** 2 **Challenges in Microbial Source Tracking**

3 Xu-Wen Wang¹, Lu Wu², Lei Dai^{2,3}, Xiaole Yin⁴, Tong Zhang⁴, Scott T. Weiss¹ & Yang-Yu
4 Liu¹

5 *¹Channing Division of Network Medicine, Department of Medicine, Brigham and Women's*
6 *Hospital and Harvard Medical School, Boston, MA, 02115, USA.*

7 *²CAS Key Laboratory of Quantitative Engineering Biology, Shenzhen Institute of Synthetic*
8 *Biology, Shenzhen Institute of Advanced Technology, Shenzhen 518055, China.*

9 *³University of Chinese Academy of Sciences, Beijing 100049, China.*

10 *⁴Environmental Microbiome Engineering and Biotechnology Laboratory, Department of Civil*
11 *Engineering, The University of Hong Kong, Hong Kong, China.*

12 **ABSTRACT**

13 Quantifying the contributions of possible environmental sources (“sources”) to a specific microbial
14 community (“sink”) is a classical problem in microbiology known as microbial source tracking
15 (MST). Solving the MST problem will not only help us understand how microbial communities
16 were formed, but also have far-reaching applications in pollution control, public health, and
17 forensics. Numerous computational methods, referred to as MST solvers hereafter, have been
18 developed in the past and applied to various real datasets to demonstrate their utility across
19 different contexts. Yet, those MST solvers do not consider microbial interactions and priority
20 effects in microbial communities. Here, we revisit the performance of several representative MST
21 solvers. We show compelling evidence that solving the MST problem using existing MST solvers
22 is impractical when ecological dynamics plays a role in community assembly. In particular, we
23 clearly demonstrate that the presence of either microbial interactions or priority effects will render
24 the MST problem mathematically unsolvable for any MST solver. We further analyze data from
25 fecal microbiota transplantation studies, finding that the state-of-the-art MST solvers fail to
26 identify donors for most of the recipients. Finally, we perform community coalescence
27 experiments to demonstrate that the state-of-the-art MST solvers fail to identify the sources for
28 most of the sinks. Our findings suggest that ecological dynamics imposes fundamental challenges
29 in solving the MST problem using computational approaches.

30 INTRODUCTION

31 Estimating the contributions or mixing proportions of different source microbial communities
32 (“sources”) to a specific microbial community (“sink”) is known as the microbial source tracking
33 (MST) problem¹⁻³. Historically, MST was framed in the context of quantifying the input of various
34 sources of fecal contamination to manage and remediate water pollution⁴. Recently, MST has been
35 used in many other contexts such as healthcare^{5,6} and forensics⁷. This is largely due to the advances
36 in metagenomics and next-generation sequencing technologies, which have enabled us to collect
37 microbiome data at an unprecedented speed⁸⁻¹¹ and provide deep insights into the roles of microbes
38 in the integrity of their environments or the well-being of their hosts¹²⁻¹⁴. Despite these advances,
39 much remains unclear regarding how the microbial communities were formed in the first place
40 and how microbes migrate across different habitats. Understanding the origins of microbial
41 communities by solving the MST problem is crucial for us to reveal their assembly rules, prevent
42 future instances of contamination, and inform disease prevention.

43 Mathematically, the MST problem can be formalized as follows. Consider a sink
44 community represented by a composition vector \mathbf{x} , where x_j corresponds to the relative abundance
45 of species- j , $1 \leq j \leq N$. Let K be the number of known sources to this sink community. Each
46 known source is represented by a composition vector $\mathbf{y}^{(a)}$, where $y_j^{(a)}$ is the relative abundance of
47 species- j in source- a ($1 \leq a \leq K$). In addition to the K known sources, we assume there is an
48 unobserved source labeled as $(K + 1)$. Our goal is to estimate the contributions or mixing
49 proportions of the $(K + 1)$ source communities to form the sink community, i.e., inferring m_a
50 ($a = 1, \dots, K + 1$) that satisfy $\sum_{a=1}^{K+1} m_a \mathbf{y}^{(a)} = \mathbf{x}$ and $\sum_{a=1}^{K+1} m_a = 1$.

51 Previous MST studies typically aimed at defining source-specific indicators (microbial or
52 chemical) with appropriate detection techniques^{1,3}. Recently, numerous computational methods
53 based on machine learning or Bayesian modeling, referred to hereafter as MST solvers, have been
54 developed to infer the contributions of different sources to a sink community^{2,4}. Here we introduce
55 three representative MST solvers. The first solver is based on the classification analysis in machine
56 learning, e.g., using the Random Forest (RF) classifier¹⁵. In this case, each source represents a
57 distinct class and RF will classify the sink into different classes with different probabilities. The
58 probabilities of the sink belonging to the different classes can be naturally interpreted as the mixing
59 proportions or contributions of those sources to the sink. Beyond the simple classification analysis,
60 more advanced statistical methods based on Bayesian modeling have been developed. For example,
61 SourceTracker is a Bayesian MST solver that explicitly models the sink as a convex mixture of

62 sources and infers the mixing proportions via Gibbs sampling¹⁶. Due to its computational
63 complexity, SourceTracker is only applicable to small- or medium-size datasets with a small
64 number of sources. FEAST (fast expectation-maximization for microbial source tracking¹⁷) is a
65 more recent statistical method. FEAST also assumes each sink is a convex combination of sources.
66 But it infers the model parameters via fast expectation-maximization, which is much more scalable
67 than Markov Chain Monte Carlo used by SourceTracker.

68 Both SourceTracker and FEAST have shown promising performance in synthetic datasets
69 and offered biologically meaningful interpretations when applied to real datasets under certain
70 contexts. Yet, the synthetic datasets used to validate these MST solvers were all generated from
71 statistical distributions, rather than dynamics models in community ecology. Hence, the ecological
72 dynamics driving the community assembly is completely ignored. We hypothesize that, after
73 considering the ecological dynamics, the power of those MST solvers might be significantly
74 restricted.

75 Here we consider two factors that heavily affect the ecological dynamics and community
76 assembly: (1) microbial interactions; (2) priority effects. Microbial interactions are ubiquitous.
77 They can be mediated by direct secretion of substances such as bacteriocins^{18,19}, ecological
78 competition between the microbes²⁰, metabolite exchange²¹, or the host's immune system
79 modulation²²⁻²⁴. In the presence of microbial interactions, the final composition of the sink
80 community will in general be fundamentally different from its initial one, i.e., the one right after
81 the source mixing, which is typically not available to us (see Fig.1). Consequently, the source
82 contributions (or mixing proportions) estimated by applying MST solvers to the final sink
83 community will be significantly different from the source contributions estimated by applying
84 MST solvers to the initial sink community.

85 Ecological theory suggests that the establishment of new species in a community can
86 depend on the order and/or timing of their arrival, a phenomenon known as *priority effects*²⁵⁻²⁸.
87 This phenomenon is actually ubiquitous in animal^{29,30}, plant³¹, and microbial communities^{28,32,33}.
88 Mechanisms of priority effects and evidence for their importance have been heavily studied for
89 microbial communities inhabiting a range of environments, including the mammalian gut³⁴⁻³⁷, the
90 plant phyllosphere³⁸⁻⁴⁰ and rhizosphere^{41,42}, soil⁴³, freshwaters⁴⁴ and oceans^{45,46}. For example, it
91 has been pointed out that priority effects probably shape the human gut microbiome during early
92 childhood⁴⁷. In particular, the infant's exposure history and the patterns of dispersal from various
93 sites in or on their mother could mediate the observed mutual exclusion between *Bacteroides spp.*,
94 *Escherichia spp.* and lactic acid producers such as *Bifidobacterium spp.* and *Lactobacillus spp.*⁴⁷.

95 In the presence of priority effects, even if the mixing proportions (source contributions) are exactly
96 the same, sink communities resulting from mixing the same set of sources but with different mixing
97 orders could be drastically different (see Fig.1). Thus, for the different sink communities, the
98 source contributions estimated by MST solvers will also be quite different, contradicting the truth.

99 To test our hypothesis, in this work we first systematically examined the impact of
100 microbial interactions and priority effects on the performance of existing MST solvers using
101 synthetic data generated by a classical population dynamics model in community ecology. We
102 found that those solvers fail in the presence of microbial interactions or priority effects. We offered
103 mathematical explanations for the failures. We then applied FEAST and SourceTracker, the two
104 state-of-the-art MST solvers, to analyze data from two fecal microbiota transplantation (FMT)
105 studies, finding that it fails to identify donors for most of the recipients. To experimentally validate
106 our hypothesis, we performed community coalescence experiments, where fecal samples from 24
107 healthy individuals (i.e., sources) were mixed and cultured *ex vivo* to form 481 sink communities.
108 We found that FEAST and SourceTracker fail to identify sources for most of the sinks. These
109 results underscore the fundamental challenges imposed by ecological dynamics in solving the
110 MST problem using computational approaches.

111

112 RESULTS

113 Impact of microbial interactions on MST.

114 To illustrate the impact of microbial interactions on MST, we simulated source and sink
115 communities as the steady states of a classical population dynamics model in community ecology
116 --- the Generalized Lotka-Volterra (GLV) model: $dX_i/dt = X_i(r_i + \sum_{j=1}^N a_{ij} X_j)$, $i = 1, \dots, N$. Here
117 X_i is the abundance (or biomass) of species- i and r_i is its intrinsic growth rate. The microbial
118 interaction matrix $\mathbf{A} = (a_{ij}) \in \mathbb{R}^{N \times N}$ can be represented by an ecological network $\mathcal{G}(\mathbf{A})$: there is
119 a directed edge ($j \rightarrow i$) in the network if and only if $a_{ij} \neq 0$. And $a_{ij} > 0$ (< 0 , or $= 0$) means
120 that species- j promotes (inhibits or does not affect) the growth of species- i , respectively. To
121 generate the matrix \mathbf{A} , we first generate the underlying network $\mathcal{G}(\mathbf{A})$ using a random graph
122 model⁴⁸ with N nodes (species) and connectivity C (representing the probability of randomly
123 connecting two nodes). Then for each link ($j \rightarrow i$) $\in \mathcal{G}(\mathbf{A})$ with $j \neq i$, we draw a_{ij} from a normal
124 distribution $\mathbb{N}(0, \sigma^2)$. Here, the standard deviation σ of this normal distribution can be considered
125 as the characteristic inter-species interaction strength. Despite its simplicity, the GLV model has

126 been successfully applied to describe the population dynamics of various microbial communities,
127 from the soil⁴⁹ and lakes⁵⁰ to the human gut^{51,52}.

128 We generated three source communities, S_1 , S_2 and S_3 , each with 30 species drawn from a
129 pool of $N = 90$ species. To simplify the MST problem, we ensured the three sources do not share
130 any common species, and the intrinsic growth rates of all species were set to be identical ($r_i = 0.5$
131 for $i = 1, \dots, N$). The composition vectors of S_1 , S_2 and S_3 (denoted as $\mathbf{y}^{(1)}$, $\mathbf{y}^{(2)}$, $\mathbf{y}^{(3)}$, respectively)
132 were obtained by running the GLV model until a steady state was reached and then normalizing
133 the steady-state abundance of each species by the total biomass of the community (see SI Sec.1
134 for details).

135 To systematically examine the impact of microbial interactions on MST, we tuned the
136 connectivity C of the ecological network $\mathcal{G}(\mathbf{A})$ and the characteristic inter-species interaction
137 strength σ in the GLV model. For a given pair of (C, σ) , we simulated 100 sink communities with
138 the initial composition vector $\mathbf{x}(0)$ given by a random mixture of the three source communities,
139 i.e., $\mathbf{x}(0) = m_1\mathbf{y}^{(1)} + m_2\mathbf{y}^{(2)} + m_3\mathbf{y}^{(3)}$, where m_a 's were drawn from uniform distribution
140 $\mathcal{U}(0,1)$ with the constraint that $\sum_a m_a = 1$. The final composition of each sink was obtained by
141 running the GLV model until a steady state. Note that to disentangle the impacts of microbial
142 interactions and priority effects on MST, here we assume a simultaneous mixing, i.e., all the
143 sources (and their species) are available at the same time to avoid priority effects.

144 We found that, with identical intrinsic species growth rates, both FEAST and
145 SourceTracker can achieve very high accuracy (with the coefficients of determination of the
146 estimated proportions $R^2 = 1$) in the absence of microbial interactions: $C = 0$ (Fig.2a) or $\sigma = 0$
147 (Fig.2b). This can be explained as follows. First, in the absence of microbial interactions and with
148 identical intrinsic species growth rates, the final composition of each sink will be identical to its
149 initial composition (right after the mixture of the three sources). Second, the three sources do not
150 share any common species, hence the MST problem becomes trivial for those solvers that assume
151 each sink is a convex combination of sources. Note that even in this ideal case, the classification-
152 based MST solver (i.e., RF) does not perform very well. This is because, as the combination of
153 different sources, the sink community's composition does not necessarily need to be similar to the
154 composition of any source.

155 Interestingly, with a nonzero C or σ , none of the three MST solvers can successfully
156 estimate the source contributions (indicated by $R^2 \approx 0$). This implies that the existing MST solvers
157 will completely fail as long as microbial interactions are present, and even in the absence of priority
158 effects (see Fig.2a,b).

159 The unsolvability of the MST problem in the presence of microbial interactions can be
160 conceptually explained as follows. Any microbial interactions will drive the sink community to
161 evolve from its initial state to its final state (Fig.2c,d). The final state will be generally different
162 from the initial one. There are two exceptions. First, the initial sink community is already at its
163 steady state and hence will not change over time. This case almost never happens, because the
164 initial sink is obtained by mixing multiple sources. Even though the sources are at their respective
165 steady states, simply mixing them will not lead to another steady state. The interactions among the
166 species across different sources will affect the assembly of the sink community. Some source-
167 specific species might even die out due to competition. Second, the system has a periodic
168 trajectory in the state space, and the initial and final states happen to be identical. This coincidence
169 generally will not happen for an unspecified time interval between the initial and final states. (See
170 SI Sec.2 for a more mathematical explanation on the difference between the initial and final states
171 of the sink community, using generic population dynamics models.) Since the initial and final
172 states of the sink community are different, the source contributions estimated by applying any
173 MST solver to the final sink community will also be different from that estimated by applying the
174 MST solver to the initial sink community. We can avoid this issue by inferring the initial state
175 from the final state. But this is impossible if the system is globally stable, i.e., any feasible initial
176 state will result in the same final state. Even if such global stability does not exist, inferring the
177 initial state from the final one would typically require detailed knowledge of the ecological
178 dynamics, which is not known *a priori*. All these factors suggest that without *a priori* knowledge
179 on the ecological dynamics, the MST problem is mathematically unsolvable in the presence of
180 microbial interactions.

181

182 [Impact of priority effects on MST.](#)

183 To examine the impact of priority effects on MST, we again simulated three source communities
184 S_1 , S_2 and S_3 whose species collections do not have any overlap (30 species for each source). The
185 final compositions of sources were obtained by running the GLV model until reaching a steady
186 state and then normalizing the steady-state abundance of each species by the total biomass of the
187 community (see SI Sec. 1 for details). For each of the $3! = 6$ mixing orders, we generated a sink
188 by mixing three sources with equal proportion $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$, then ran the GLV model to obtain its final
189 composition. For comparison purposes, we also generated a sink through simultaneous mixing of
190 the three sources with equal proportion $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$. We visualized the compositions of the three

191 sources and the seven sinks using the t-distributed stochastic neighbor embedding (t-SNE) method,
192 finding that the compositions of the seven sinks are clearly different (see Fig.3a). We then ran
193 FEAST, the fastest MST solver, to estimate the contributions of the three sources to each sink,
194 finding that the contributions are different for different sinks, despite the true mixing proportions
195 being exactly the same (Fig.3b). In the above simulations we set the network connectivity $C = 0.5$
196 and the characteristic interaction strength $\sigma = 1$.

197 The above results make us wonder the solvability of the MST problem in the presence of
198 priority effects. Here we offer an outline of proof that the MST problem is mathematically
199 unsolvable in the presence of priority effects. Consider a set of source communities. If we mix
200 them in different orders (but using the same set of mixing proportions), this will generally lead to
201 different sink communities due to priority effects. The between-sink dissimilarity can be as large
202 as the between-source dissimilarity (see Fig.3c). We emphasize that different mixing orders
203 generally result in different sink communities even in the absence of any microbial interactions
204 (see SI Sec.3 for a mathematical explanation). For different sink communities, the source
205 contributions estimated by any computational method (i.e., MST solver) will also be different,
206 contradicting the fact that the source contributions (i.e., mixing proportions) are exactly the
207 same. This proof by contradiction clearly illustrates that the MST problem is mathematically
208 unsolvable in the presence of priority effects.

209

210 [Evaluation of MST solvers using data from FMT studies.](#)

211 During FMT, fecal microbiota from a carefully screened, healthy donor is introduced to a recipient
212 through either the lower or upper gastrointestinal tract. It is a “natural” mixing experiment that can
213 be used to evaluate the performance of MST solvers. To achieve that, we applied FEAST and
214 SourceTracker to analyze data from two FMT studies^{53,54}.

215 In the first study, recurrent *Clostridioides difficile* infection (rCDI) patients were treated
216 with encapsulated donor material for FMT (cap-FMT)⁵³. Fig.4a shows the donor-recipient
217 relationship between 7 healthy donors and 88 rCDI patients (i.e., recipients). Each trajectory
218 represents a donor and one of its recipients with fecal samples collected at (up to) five different
219 time points: pre-FMT, 2–6 days post FMT, weeks (7–20 days) post FMT, months (21–60 days)
220 post FMT, and long term (>60 days). The Principal Coordinate Analysis (PCoA) plot of all the
221 microbiome samples is shown in Fig.4b. We tested if FEAST can correctly identify the donor of a
222 recipient. To achieve that, we considered each post-FMT sample of each recipient as a sink

223 community and considered the fecal samples of all the 7 donors, as well as the recipient's pre-
224 FMT sample as potential source communities. Then we applied FEAST to solve the MST problem.
225 For each sink community, among all the 7 donors, we referred to the one whose fecal sample has
226 the highest contribution estimated by FEAST as the "predicted donor" (green squares, Fig.4c,
227 Fig.S1). Interestingly, we found that for a large portion (61%) of the sink communities, FEAST
228 failed to identify the true donor (red circles, Fig.4c, Fig.S1), though the average Jensen-Shannon
229 divergence among those donors is higher enough (0.63). Similar results were found for
230 SourceTracker (see Fig.S2). These results clearly demonstrate the limitation of existing MST
231 solvers.

232 In the second FMT study, the gut microbiota of human donors with autism spectrum
233 disorder (ASD) or typically-developing (TD) controls were transplanted into germ-free mice⁵⁴.
234 The dataset includes 8 donors, 13 recipients, and in total 106 post-FMT sink communities. We
235 again examined whether FEAST can correctly identify the true donor of each sink community. For
236 each sink community, among the 8 donors, we refer to the one whose fecal sample has the highest
237 contribution predicted by FEAST as the "predicted donor" (green squares, Fig.S3). We found that
238 for 40% of the sink communities, FEAST failed to identify the true donor (red circles, Fig.S3).
239 Similar results were observed for SourceTracker (see Fig.S4).

240

241 [Evaluation of MST solvers using data from community coalescence experiments.](#)

242 To further evaluate MST solvers using real data, we performed community coalescence
243 experiments, where fecal microbiota from 24 healthy individuals (i.e., sources) were mixed and
244 cultured *ex vivo* to form 481 sink communities (see SI Sec.4 for details). Among the 481 sinks,
245 256 sinks were obtained by mixing two different sources (pair-wise mixing), and the remaining
246 225 sinks were obtained by mixing four different sources (quadruple-wise mixing). After
247 inoculation, the sink communities were transferred into fresh medium every 24 hours (1:200
248 dilution) for 10 transfers⁵⁵ (see Fig.5a). Samples collected at the final time point were sequenced
249 and the resulting taxonomic profiles were considered as the steady-state composition of sinks (see
250 Methods). As expected, we found that the source and sink communities had distinct taxonomic
251 profiles (Fig.S5-S6).

252 To examine the performance of FEAST in community coalescence experiments, we first
253 applied FEAST to analyze the compositions of the 256 sinks obtained in the pair-wise mixing
254 experiments. We ranked the estimated contributions of 24 potential sources to each sink and
255 selected the top-two as the predicted sources. We found that the predicted sources (green squares)

256 are different from the true sources (red circles) for most of the 256 sinks (Fig.5b and Fig.S7). This
257 is also true for the cases of quadruple-wise mixing (Fig.S9). Similar results were observed for
258 SourceTracker (see Fig.S8, S10).

259 Note that some donor samples (e.g., S0820B, S0814D) were predicted as sources for many
260 sinks. We found this is due to the high abundance of common ASVs shared by sinks and those
261 particular sources (Fig.S11).

262

263 DISCUSSION

264 Many computational methods have been developed to solve the MST problem. Yet, those methods
265 ignored the underlying ecological dynamics that drive the assembly of microbial communities. For
266 example, as a Bayesian MST solver, SourceTracker explicitly models the sink as a convex mixture
267 of sources and infers the mixing proportions via Gibbs sampling¹⁶. This approach was inspired by
268 the “analogy” between quantifying the proportion of different source environments to a sink
269 microbial community and inferring the mixing proportions of conversation topics in a test
270 document^{56,57}. Here we point out that this analogy is inappropriate. In topic modeling, which is a
271 specific research area in natural language processing, the goal is to discover the abstract “topics”
272 that occur in a collection of documents. In a sense, those documents are static or “dead”. By
273 contrast, in MST we are typically dealing with alive (or even flourishing) microbial communities,
274 where ecological dynamics plays an important role in community assembly and determining their
275 state, i.e., the microbial composition. In the presence of ecological dynamics, a sink community
276 cannot be simply considered as a convex mixture of known and unknown sources. In this work,
277 through numerical simulations, analytical calculations, and real data analysis, we presented
278 compelling evidence that ecological dynamics impose fundamental challenges in MST. In
279 particular, we clearly demonstrate that the presence of either microbial interactions or priority
280 effects will render the MST problem mathematically unsolvable for any MST solver.

281 MST solvers have been applied to various real datasets and demonstrated their utility across
282 two fundamentally different contexts. First, as originally intended, they were used to quantify the
283 contribution of different source environments to a sink microbial community. For example,
284 SourceTracker was used to estimate the contributions of bacteria from ‘gut’, ‘oral’, ‘skin’, ‘soil’
285 and ‘unknown’ sources to several indoor sink environments (e.g., office buildings, hospitals, and
286 research laboratories)¹⁶. It was found that wet-lab surface communities tended to be composed
287 mainly of bacteria from ‘skin’ and ‘unknown’, while neonatal intensive care units and office

288 communities were typically dominated by skin bacteria. FEAST was used to estimate if taxa in the
289 infant gut originate from the birth canal, or if they are derived from some other external source at
290 a later time point¹⁷. By treating samples taken from the infants at age 12 months as sinks,
291 considering respective earlier time points and maternal samples as sources, a significantly larger
292 maternal contribution in vaginally delivered infants over cesarean-delivered infants was found.
293 Moreover, biological mothers were more likely to be identified as sources of their infant's
294 microbiome than other potential source communities. Although these results seem reasonable and
295 agree well with our intuition, we suggest that the whole community of microbiome research should
296 be very cautious when interpreting the results of existing MST solvers in this context. The source
297 contributions estimated by MST solvers might be quite different from the true contributions due
298 to complex ecological dynamics. This is particularly important for microbial communities living
299 in nutrient-rich environments such as the human gut. For microbial communities living in
300 oligotrophic environments (e.g., soil, ocean, etc.), the growth rates of bacteria and assembly
301 process of communities are relatively slow^{58,59,60} and the impact of ecological dynamics on MST
302 might be relatively low⁶¹. But even in this case, interpreting the results of existing MST solvers
303 should be done with great caution.

304 Second, MST solvers have been used as a metric of similarity¹⁷. In this context, instead of
305 quantifying the contribution of different sources to a sink, we aim for capturing the similarities
306 between the sink and its characteristic environments using mixing proportions estimated by MST
307 solvers. Each sink can be represented by a similarity feature vector, characterizing its similarity to
308 each of its characteristic environments. For example, FEAST has been used in this context to
309 distinguish patients in ICU from healthy adults, and capture shifts in microbial community
310 composition that may underlie differences between pathogenic and neutral phenotypes¹⁷. We think
311 this is a much more meaningful and practical way of using MST solvers to analyze real data.

312 A recent study has shown that the strain tracking approach⁶² can predict whether two
313 metagenomics samples originate from the same donor via counting the number of species that
314 share closely-related strains. Yet, the contribution of different sources to a given sink remains
315 unknown. More importantly, challenges imposed by ecological dynamics are still there, which do
316 not rely on a particular sequencing method. For example, in the presence of microbial interactions
317 and priority effects, those source-specific microbial strains may not be able to survive in the sink
318 community at all. This actually raises a serious concern on any approaches based on indicator
319 species in solving the MST problem.

320

321 References

- 322 1. Simpson, J. M., Santo Domingo, J. W. & Reasoner, D. J. Microbial source tracking: state
323 of the science. *Environ. Sci. Technol.* **36**, 5279–5288 (2002).
- 324 2. Hagedorn, C., Blanch, A. R. & Harwood, V. J. *Microbial source tracking: methods,*
325 *applications, and case studies.* (Springer Science & Business Media, 2011).
- 326 3. Rock, C., Rivera, B. & Gerba, C. P. Microbial source tracking. in *Environmental*
327 *Microbiology* 309–317 (Elsevier, 2015).
- 328 4. Devane, M. L., Weaver, L., Singh, S. K. & Gilpin, B. J. Fecal source tracking methods to
329 elucidate critical sources of pathogens and contaminant microbial transport through New
330 Zealand agricultural watersheds—a review. *J. Environ. Manage.* **222**, 293–303 (2018).
- 331 5. McDonald, D. *et al.* Extreme dysbiosis of the microbiome in critical illness. *Mosphere* **1**,
332 e00199-16 (2016).
- 333 6. Dominguez-Bello, M. G. *et al.* Partial restoration of the microbiota of cesarean-born
334 infants via vaginal microbial transfer. *Nat. Med.* **22**, 250–253 (2016).
- 335 7. Teaf, C. M., Flores, D., Garber, M. & Harwood, V. J. Toward forensic uses of microbial
336 source tracking. *Microbiol. Spectr.* **6**, 6.1. 05 (2018).
- 337 8. Qin, J. *et al.* A human gut microbial gene catalogue established by metagenomic
338 sequencing. *nature* **464**, 59–65 (2010).
- 339 9. Franzosa, E. A. *et al.* Species-level functional profiling of metagenomes and
340 metatranscriptomes. *Nat. Methods* **15**, 962–968 (2018).
- 341 10. Gill, S. R. *et al.* Metagenomic analysis of the human distal gut microbiome. *science* **312**,
342 1355–1359 (2006).
- 343 11. Goodrich, J. K. *et al.* Conducting a microbiome study. *Cell* **158**, 250–262 (2014).
- 344 12. Louca, S. *et al.* Function and functional redundancy in microbial systems. *Nat. Ecol.*
345 *Evol.* **2**, 936–943 (2018).
- 346 13. Sunagawa, S. *et al.* Structure and function of the global ocean microbiome. *Science* **348**,
347 (2015).
- 348 14. The Human Microbiome Project Consortium. Structure, function and diversity of the
349 healthy human microbiome. *Nature* **486**, 207–214 (2012).
- 350 15. Smith, A., Sterba-Boatwright, B. & Mott, J. Novel application of a statistical technique,
351 Random Forests, in a bacterial source tracking study. *Water Res.* **44**, 4067–4076 (2010).
- 352 16. Knights, D. *et al.* Bayesian community-wide culture-independent microbial source
353 tracking. *Nat. Methods* **8**, 761–763 (2011).
- 354 17. Shenhav, L. *et al.* FEAST: fast expectation-maximization for microbial source tracking.
355 *Nat. Methods* **16**, 627–632 (2019).
- 356 18. Bucci, V., Nadell, C. D. & Xavier, J. B. The evolution of bacteriocin production in
357 bacterial biofilms. *Am. Nat.* **178**, E162–E173 (2011).
- 358 19. Levy, R. & Borenstein, E. Metabolic modeling of species interaction in the human
359 microbiome elucidates community-level assembly rules. *Proc. Natl. Acad. Sci.* **110**, 12804–
360 12809 (2013).
- 361 20. Bucci, V., Bradde, S., Biroli, G. & Xavier, J. B. Social interaction, noise and antibiotic-
362 mediated switches in the intestinal microbiota. *PLoS Comput. Biol.* **8**, e1002497 (2012).
- 363 21. Uehling, J. K. *et al.* Microfluidics and metabolomics reveal symbiotic bacterial–fungal
364 interactions between *Mortierella elongata* and *Burkholderia* include metabolite exchange. *Front.*
365 *Microbiol.* **10**, 2163 (2019).
- 366 22. Dorrestein, P. C., Mazmanian, S. K. & Knight, R. Finding the missing links among
367 metabolites, microbes, and the host. *Immunity* **40**, 824–832 (2014).
- 368 23. Round, J. L. & Mazmanian, S. K. The gut microbiota shapes intestinal immune responses
369 during health and disease. *Nat. Rev. Immunol.* **9**, 313–323 (2009).

- 370 24. Buffie, C. G. & Pamer, E. G. Microbiota-mediated colonization resistance against
371 intestinal pathogens. *Nat. Rev. Immunol.* **13**, 790–801 (2013).
- 372 25. Grainger, T. N., Letten, A. D., Gilbert, B. & Fukami, T. Applying modern coexistence
373 theory to priority effects. *Proc. Natl. Acad. Sci.* **116**, 6205–6210 (2019).
- 374 26. Ke, P.-J. & Letten, A. D. Coexistence theory and the frequency-dependence of priority
375 effects. *Nat. Ecol. Evol.* **2**, 1691–1695 (2018).
- 376 27. Zhao, N., Saavedra, S. & Liu, Y.-Y. Impact of colonization history on the composition of
377 ecological systems. *Phys. Rev. E* **103**, 052403 (2021).
- 378 28. Debray, R. *et al.* Priority effects in microbiome assembly. *Nat. Rev. Microbiol.* 1–13
379 (2021).
- 380 29. Alford, R. A. & Wilbur, H. M. Priority effects in experimental pond communities:
381 competition between *Bufo* and *Rana*. *Ecology* **66**, 1097–1105 (1985).
- 382 30. Rasmussen, N. L., Van Allen, B. G. & Rudolf, V. H. Linking phenological shifts to
383 species interactions through size-mediated priority effects. *J. Anim. Ecol.* **83**, 1206–1215 (2014).
- 384 31. Kardol, P., Souza, L. & Classen, A. T. Resource availability mediates the importance of
385 priority effects in plant community assembly and ecosystem function. *Oikos* **122**, 84–94 (2013).
- 386 32. Clay, P. A., Dhir, K., Rudolf, V. H. & Duffy, M. A. Within-host priority effects
387 systematically alter pathogen coexistence. *Am. Nat.* **193**, 187–199 (2019).
- 388 33. Louette, G. & De Meester, L. Predation and priority effects in experimental zooplankton
389 communities. *Oikos* **116**, 419–426 (2007).
- 390 34. Chng, K. R. *et al.* Metagenome-wide association analysis identifies microbial
391 determinants of post-antibiotic ecological recovery in the gut. *Nat. Ecol. Evol.* **4**, 1256–1267
392 (2020).
- 393 35. Lee, S. M. *et al.* Bacterial colonization factors control specificity and stability of the gut
394 microbiota. *Nature* **501**, 426–429 (2013).
- 395 36. Martínez, I. *et al.* Experimental evaluation of the importance of colonization history in
396 early-life gut microbiota assembly. *Elife* **7**, e36521 (2018).
- 397 37. Furman, O. *et al.* Stochasticity constrained by deterministic effects of diet and age drive
398 rumen microbiome assembly dynamics. *Nat. Commun.* **11**, 1–13 (2020).
- 399 38. Seybold, H. *et al.* A fungal pathogen induces systemic susceptibility and systemic shifts
400 in wheat metabolome and microbiome composition. *Nat. Commun.* **11**, 1–12 (2020).
- 401 39. Carlström, C. I. *et al.* Synthetic microbiota reveal priority effects and keystone strains in
402 the *Arabidopsis* phyllosphere. *Nat. Ecol. Evol.* **3**, 1445–1454 (2019).
- 403 40. Halliday, F. W. *et al.* Facilitative priority effects drive parasite assembly under
404 coinfection. *Nat. Ecol. Evol.* **4**, 1510–1521 (2020).
- 405 41. Wei, Z. *et al.* Trophic network architecture of root-associated bacterial communities
406 determines pathogen invasion and plant health. *Nat. Commun.* **6**, 1–9 (2015).
- 407 42. Kennedy, P. G., Peay, K. G. & Bruns, T. D. Root tip competition among ectomycorrhizal
408 fungi: are priority effects a rule or an exception? *Ecology* **90**, 2098–2107 (2009).
- 409 43. Ferrero, A. F. Effect of compaction simulating cattle trampling on soil physical
410 characteristics in woodland. *Soil Tillage Res.* **19**, 319–329 (1991).
- 411 44. Fukami, T. *et al.* Assembly history dictates ecosystem functioning: evidence from wood
412 decomposer communities. *Ecol. Lett.* **13**, 675–684 (2010).
- 413 45. Enke, T. N. *et al.* Modular assembly of polysaccharide-degrading marine microbial
414 communities. *Curr. Biol.* **29**, 1528–1535. e6 (2019).
- 415 46. Svoboda, P., Lindström, E. S., Ahmed Osman, O. & Langenheder, S. Dispersal timing
416 determines the importance of priority effects in bacterial communities. *ISME J.* **12**, 644–646
417 (2018).

- 418 47. Sprockett, D., Fukami, T. & Relman, D. A. Role of priority effects in the early-life
419 assembly of the gut microbiota. *Nat. Rev. Gastroenterol. Hepatol.* **15**, 197–205 (2018).
- 420 48. Erdos, P. & Rényi, A. On the evolution of random graphs. *Publ Math Inst Hung Acad Sci*
421 **5**, 17–60 (1960).
- 422 49. Moore, J. C., de Ruiter, P. C., Hunt, H. W., Coleman, D. C. & Freckman, D. W.
423 Microcosms and soil ecology: critical linkages between fields studies and modelling food webs.
424 *Ecology* **77**, 694–705 (1996).
- 425 50. Dam, P., Fonseca, L. L., Konstantinidis, K. T. & Voit, E. O. Dynamic models of the
426 complex microbial metapopulation of lake mendota. *NPJ Syst. Biol. Appl.* **2**, 1–7 (2016).
- 427 51. Stein, R. R. *et al.* Ecological modeling from time-series inference: insight into dynamics
428 and stability of intestinal microbiota. *PLoS Comput. Biol.* **9**, e1003388 (2013).
- 429 52. Buffie, C. G. *et al.* Precision microbiome reconstitution restores bile acid mediated
430 resistance to *Clostridium difficile*. *Nature* **517**, 205–208 (2015).
- 431 53. Staley, C. *et al.* Predicting recurrence of *Clostridium difficile* infection following
432 encapsulated fecal microbiota transplantation. *Microbiome* **6**, 166 (2018).
- 433 54. Sharon, G. *et al.* Human Gut Microbiota from Autism Spectrum Disorder Promote
434 Behavioral Symptoms in Mice. *Cell* **177**, 1600-1618.e17 (2019).
- 435 55. Aranda-Díaz, A. *et al.* Establishment and characterization of stable, diverse, fecal-derived
436 in vitro microbial communities that model the intestinal microbiota. *Cell Host Microbe* (2022).
- 437 56. Blei, D. M., Ng, A. Y. & Jordan, M. I. Latent Dirichlet Allocation. *J. Mach. Learn. Res.*
438 **3**, 993–1022 (2003).
- 439 57. Finding scientific topics. <https://www.pnas.org/doi/10.1073/pnas.0307752101>
440 doi:10.1073/pnas.0307752101.
- 441 58. Vadia, S. & Levin, P. A. Growth rate and cell size: a re-examination of the growth law.
442 *Curr. Opin. Microbiol.* **24**, 96–103 (2015).
- 443 59. Bastviken, D. *et al.* Chemosynthesis☆. *Ref. Module Earth Syst. Environ. Sci.* (2014).
- 444 60. Harris, D. & Paul, E. A. Measurement of bacterial growth rates in soil. *Appl. Soil Ecol.* **1**,
445 277–290 (1994).
- 446 61. Li, L.-G., Huang, Q., Yin, X. & Zhang, T. Source tracking of antibiotic resistance genes
447 in the environment — Challenges, progress, and prospects. *Water Res.* **185**, 116127 (2020).
- 448 62. Zhao, S., Dai, C. L., Evans, E. D., Lu, Z. & Alm, E. J. *Tracking strains predicts personal*
449 *microbiomes and reveals recent adaptive evolution.*
450 <http://biorxiv.org/lookup/doi/10.1101/2020.09.14.296970> (2020)
451 doi:10.1101/2020.09.14.296970.
- 452 63. Gonzalez, A. *et al.* Qiita: rapid, web-enabled microbiome meta-analysis. *Nat. Methods*
453 **15**, 796–798 (2018).

454
455
456 **Data and code availability.** The sequencing data from the first FMT study is available at Sequence
457 Read Archive at the National Center for Biotechnology Information under BioProject accession
458 number SRP070464. The sequencing data from the second FMT study is available at Qiita⁶³ with
459 ID: 11809. The raw sequencing data from the community coalescence experiments is available at
460 European Nucleotide Archive (ENA) under study accession number PRJEB51290. The code used
461 to generate the simulated data is available at: <https://github.com/spxuw/MST>.

462

463 **Author Contributions.** Y.-Y.L conceived and designed the project. X.-W.W and Y.-Y.L did the
464 analytical calculations. X.-W.W did all the numerical calculations and analyzed all the simulated
465 and real datasets. L.W. and L.D. designed and performed the community coalescence experiments.
466 X.-W.W. and Y.-Y.L wrote the manuscript. L.W., L.D., X.Y., T.Z., and S.T.W interpreted the
467 results, reviewed and edited the manuscript. All authors approved the manuscript.

468
469 **Acknowledgement.** Y.-Y.L. acknowledges grants from National Institutes of Health
470 (R01AI141529, R01HD093761, RF1AG067744, UH3OD023268, U19AI095219 and
471 U01HL089856).

472

473

474

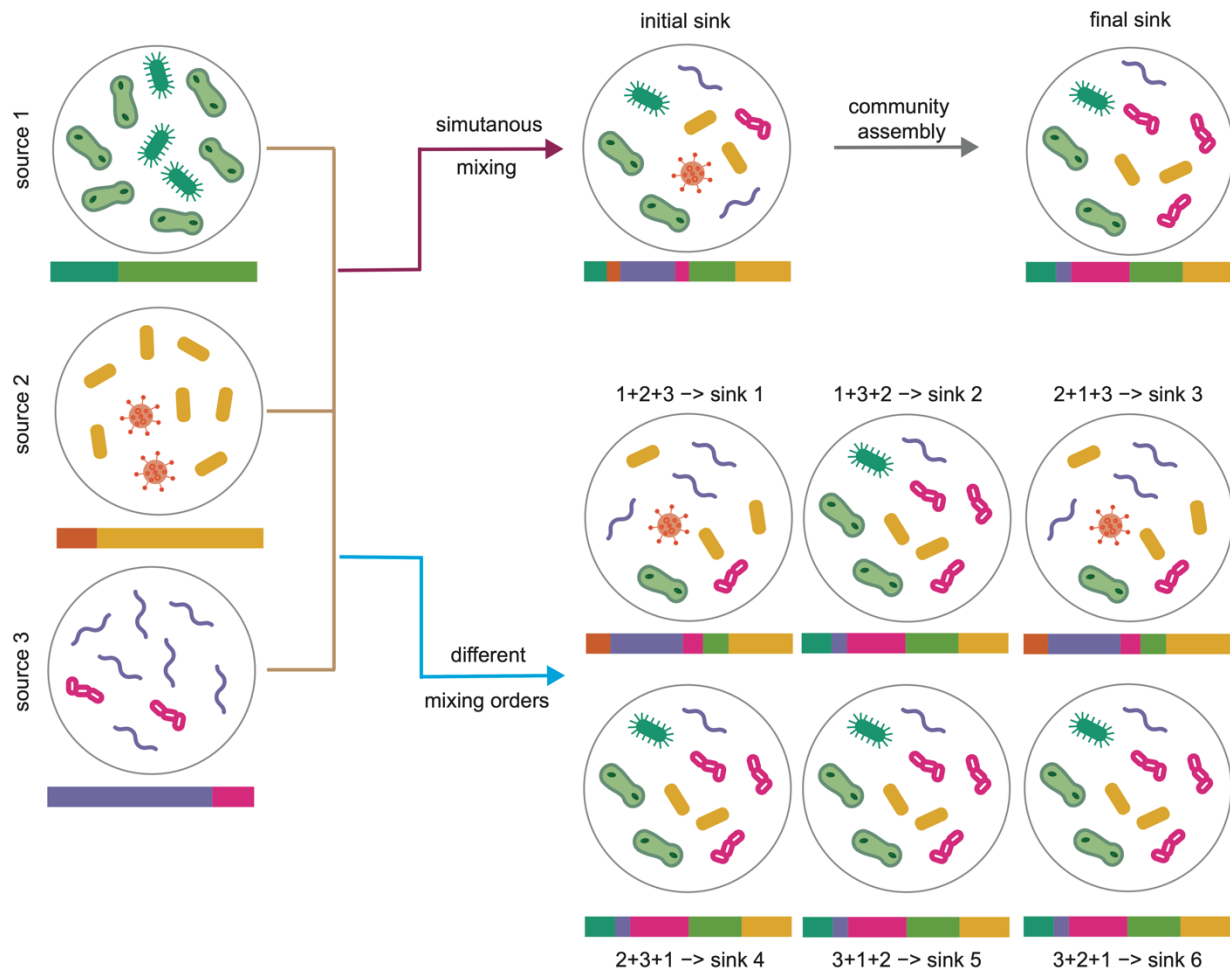
475

476

477

478

479

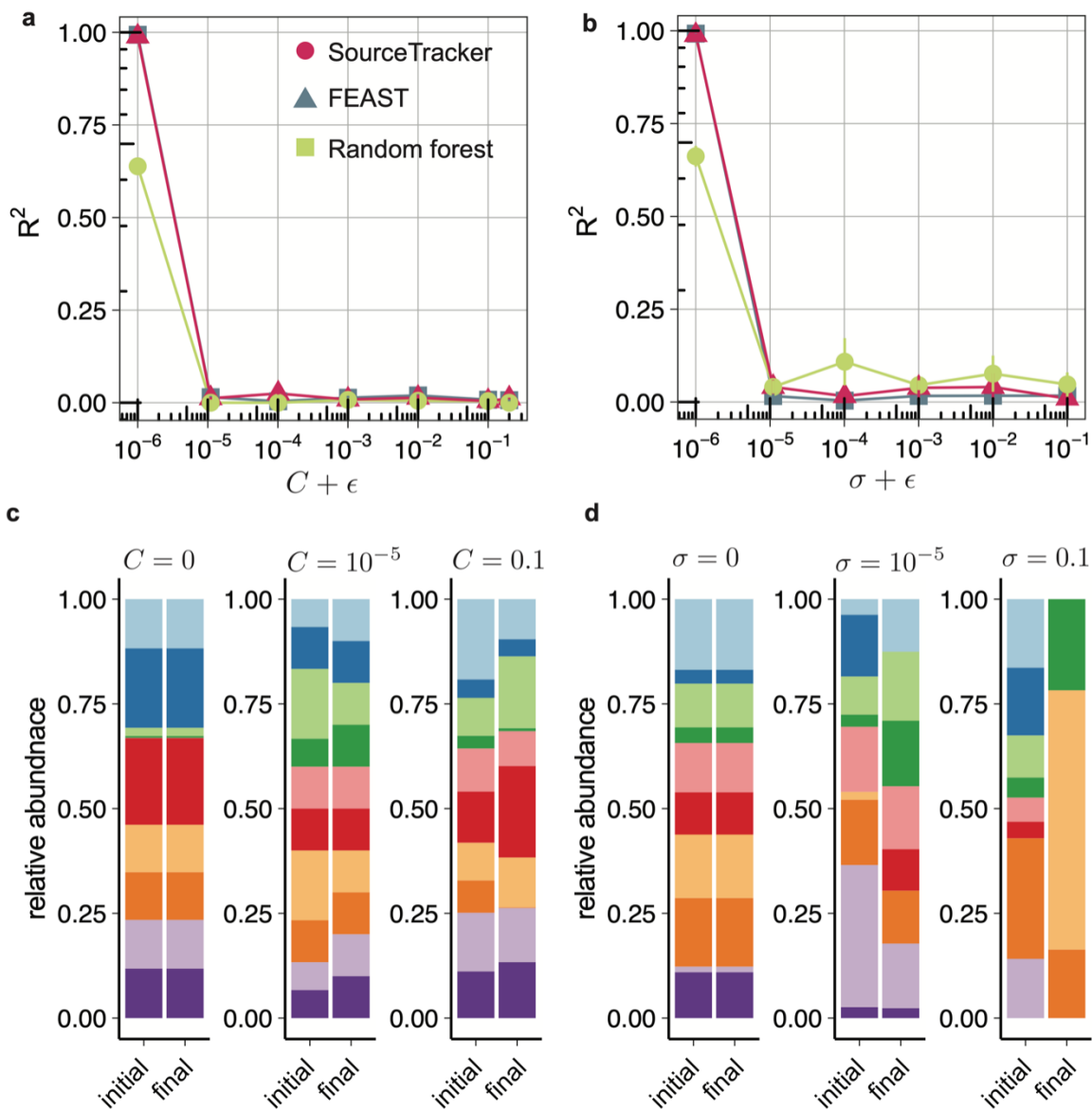


480

481 **Figure 1: Ecological dynamics imposes fundamental challenges in microbial source tracking.**

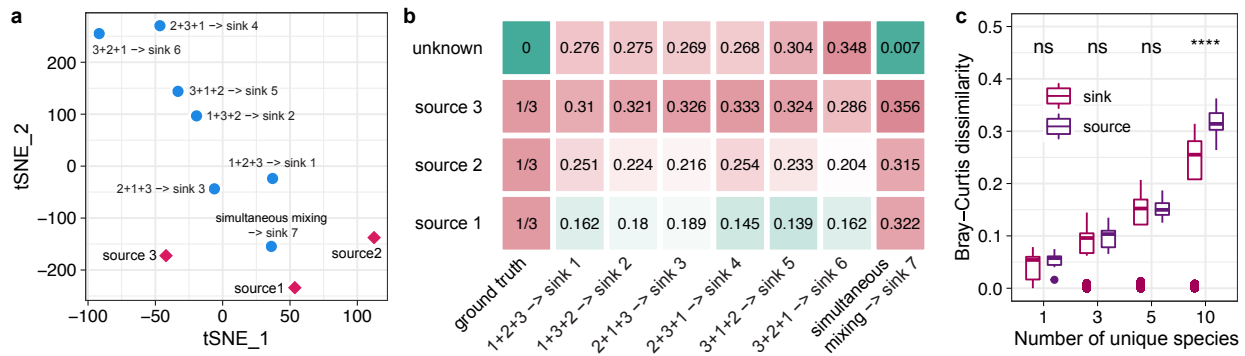
482 **(Top)** A sink is obtained by simultaneously mixing three sources (without any species overlap)
 483 with mixing proportions (1/3,1/3,1/3). Due to the presence of microbial interactions, the initial
 484 composition of the sink community (right after the mixing, which is typically not available for
 485 MST) can be significantly different from the final composition (which is the input of MST solvers).
 486 Applying any MST solver to the final sink composition will yield different results from applying
 487 the MST solver to the initial sink composition. **(Bottom)** Due to the priority effects, three sources
 488 mixed with different orders can result in total $3! = 6$ different sinks with different compositions,
 489 even if the mixing proportions of the sources are exactly the same for the different mixing orders.

490



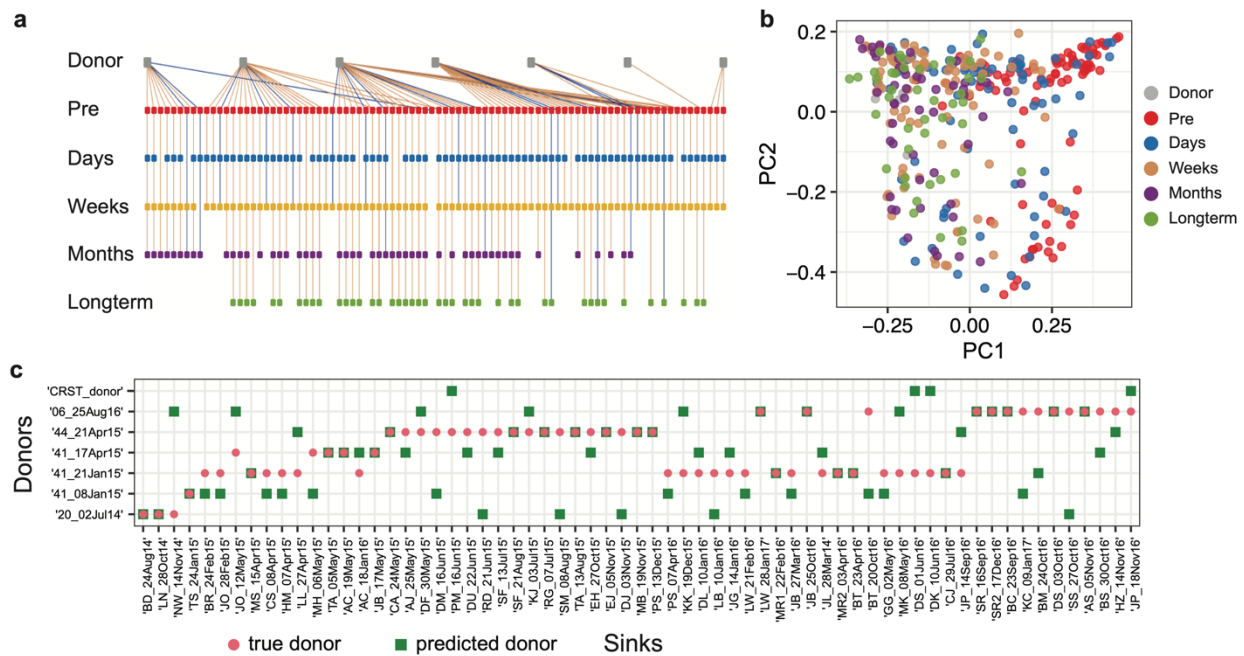
491

492 **Figure 2: Impact of microbial interactions on MST.** a-b, Performance of SourceTracker (red),
 493 FEAST (blue) and Random Forest (green) in simulated sinks with different network connectivity
 494 C (a) and characteristic interaction strengths σ (b). Each simulation was performed using 3
 495 synthetic sources and 100 synthetic sinks. Accuracy of each method is measured as the coefficients
 496 of determination (R^2) of the estimated proportions. Each point represents the mean R^2 for three
 497 independent source sets; error bars show s.e.m ($n = 3$) of the mean of R^2 over three sources. c-d,
 498 Initial and final steady compositions (we only show the relative abundance of the first 10 species
 499 for visualization purpose) of a sink with different network connectivity (c) and characteristic
 500 interaction strengths (d). In (a,c), the diagonal elements of the interaction matrix \mathbf{A} are set to be
 501 $a_{ii} = -5C$ to ensure the stability of the community, and the characteristic interaction strength $\sigma =$
 502 0.1 . In (b,d), we set $a_{ii} = -5\sigma$ to ensure the stability, and the network connectivity $C = 0.5$. In
 503 all the simulations, we set the intrinsic growth rate $r = 0.5$ for all the species. We added a pseudo
 504 number $\epsilon = 10^{-6}$ to the x-axis for visualization purpose.



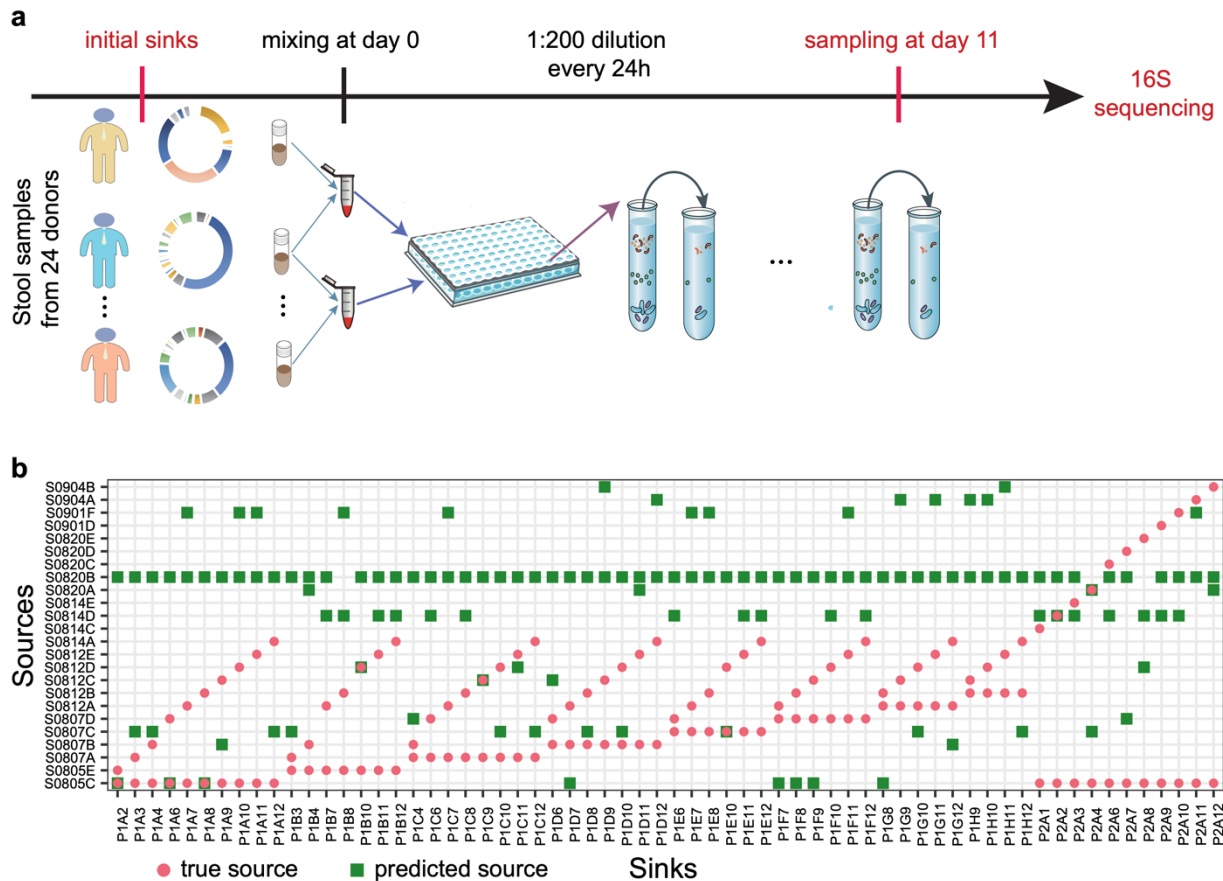
505
 506 **Figure 3: Impact of priority effects on MST. a-b,** We synthesized three sources S_1 , S_2 and S_3
 507 whose species collections do not have any overlap (30 species for each source). We mixed these
 508 three sources using six different mixing orders but with the same mixing proportions $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$,
 509 rendering six sinks. We set the network connectivity $C = 0.5$, the characteristic interaction
 510 strength $\sigma = 1$, the intrinsic growth rate $r = 0.5$ for each species. We set the diagonal elements of
 511 interaction matrix \mathbf{A} to be $a_{ii} = -5$ to ensure the stability. **a,** Dimensionality reduction using t-
 512 SNE shows the variations among the six sinks generated from the six different mixing orders. **b,**
 513 Contribution of each source to the six simulated sinks estimated by FEAST. **c,** Between-sink and
 514 between-source Bray-Curtis dissimilarity. We synthesized five sources. The species collection of
 515 each source includes N_u unique species and the remaining $(90 - 5N_u)$ species are shared by all
 516 the sources. We mixed these five sources with the same mixing proportions $(\frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5})$ in 100
 517 different mixing orders randomly chosen from the total $5! = 120$ mixing orders. We set the
 518 network connectivity $C = 0.5$, the characteristic interaction strength $\sigma = 1$, the intrinsic growth
 519 rate $r = 0.5$ for each species. We set the diagonal elements of interaction matrix \mathbf{A} to be $a_{ii} =$
 520 -10 to ensure the stability. P-values were calculated using one-sided Wilcoxon test.

521
 522
 523
 524
 525
 526



527

528 **Figure 4: Evaluation of FEAST using FMT data from Staley et al.**⁵³ **a**, Donor-recipient
 529 relationship. Each trajectory represents a donor and its corresponding recipients at up to 5 time
 530 points. Trajectories of recipients who responded to FMT (i.e., responders) are colored in yellow.
 531 Trajectories of non-responders are colored in blue. **b**, Principal Coordinates Analysis (PCoA) plot
 532 based on the Bray-Curtis dissimilarity. **c**, True donor (red cycle) vs. predicted donor (green square)
 533 of each recipient. For each post-FMT community (sink), among all the 7 donors, we referred to
 534 the one whose fecal sample has the highest contribution estimated by FEAST as the “predicted
 535 donor”. Here, we only showed the results for the first 65 sinks for the visualization purpose (see
 536 Fig.S1 for results of the remaining 194 sinks). Sources: microbiome samples of donors and the
 537 pre-FMT samples of recipients; Sinks: post-FMT samples of recipients.



1 **Ecological Dynamics Imposes Fundamental Challenges in** 2 **Microbial Source Tracking**

3 *Supplementary Information*

4 Xu-Wen Wang¹, Lu Wu², Lei Dai^{2,3}, Xiaole Yin⁴, Tong Zhang⁴, Scott T. Weiss¹ & Yang-Yu
5 Liu¹

6 ¹*Channing Division of Network Medicine, Department of Medicine, Brigham and Women's*
7 *Hospital and Harvard Medical School, Boston, MA, 02115, USA.*

8 ²*CAS Key Laboratory of Quantitative Engineering Biology, Shenzhen Institute of Synthetic*
9 *Biology, Shenzhen Institute of Advanced Technology, Shenzhen 518055, China.*

10 ³*University of Chinese Academy of Sciences, Beijing 100049, China.*

11 ⁴*Environmental Microbiome Engineering and Biotechnology Laboratory, Department of Civil*
12 *Engineering, The University of Hong Kong, Hong Kong, China.*

13 **Table of Contents**

| | | |
|----|--|----------|
| 14 | 1. Using an ecological model to generate synthetic microbiome data. | 2 |
| 15 | 2. Microbial interactions affect the assembly of the sink community. | 4 |
| 16 | 3. Priority effects affect the assembly of the sink community. | 5 |
| 17 | 4. Community coalescence experiments. | 6 |
| 18 | References. | 8 |

34 1. Using an ecological model to generate synthetic microbiome data.

35 To systematically reveal the impacts of the microbial interactions and priority effects on MST, we
36 generated synthetic data using the classical Generalized Lotka-Volterra (GLV) model¹:

$$37 \quad \frac{dX_i(t)}{dt} = X_i(t) \left[r_i + \sum_{j=1}^N a_{ij} X_j(t) \right], i = 1, \dots, N.$$

38 Here $X_i(t)$ represents the absolute abundance of species- i at time $t \geq 0$, r_i is its intrinsic growth
39 rate, which is randomly drawn from a uniform distribution $\mathcal{U}(0,1)$, if not specified otherwise. The
40 inter-species interactions are encoded in the interaction matrix $\mathbf{A} = (a_{ij}) \in \mathbb{R}^{N \times N}$, with $a_{ij} > 0$
41 (< 0 , or $= 0$) means that species- j promotes (inhibits or does not affect) the growth of species- i ,
42 respectively. To generate the matrix \mathbf{A} , we first generate the underlying ecological network $\mathcal{G}(\mathbf{A})$
43 using an Erdős-Rényi random graph model² with N nodes (species) and connectivity C (the
44 probability of randomly connecting two nodes). Then for each link $(j \rightarrow i) \in \mathcal{G}(\mathbf{A})$ with $j \neq i$, we
45 draw a_{ij} from a normal distribution $\mathbb{N}(0, \sigma^2)$. The standard deviation σ of this normal distribution
46 represents the characteristic inter-species interaction strength. To ensure the stability of the system,
47 the diagonal elements of \mathbf{A} are set to be $a_{ii} = -dC$ in tuning C or $a_{ii} = -d\sigma$ in tuning σ . Here d
48 is a positive constant. All other entries of \mathbf{A} are set to be zero.

49
50 We generated k source communities, S_1, S_2, \dots, S_k , each with N_s species drawn from a pool of
51 $N = 90$ species. To simplify the MST problem, the intrinsic growth rates of all species were set
52 to be identical ($r = 0.5$). The composition vectors of S_1, S_2, \dots, S_k (denoted as $\mathbf{y}^{(1)}, \mathbf{y}^{(2)}, \dots, \mathbf{y}^{(k)}$,
53 respectively) were obtained by running the GLV model (i.e., numerically solving the ordinary
54 differential equations (ODEs) in the GLV model) with initial species abundances randomly chosen
55 from a uniform distribution $\mathcal{U}(0,1)$, until a steady state was reached and then normalizing the
56 steady-state abundance of each species by the total biomass of the community.

57
58 The sink obtained by simultaneously mixing the k sources was simulated as follows:

- 59 1) The mixing proportions of k sources were randomly drawn from a uniform distribution
60 with constraint $\sum_{a=1}^k m_a = 1$.
- 61 2) The initial composition of the sink community is calculated as: $\mathbf{x}(0) = m_1 \mathbf{y}^{(1)} +$
62 $m_2 \mathbf{y}^{(2)} + \dots + m_k \mathbf{y}^{(k)}$. And the initial (absolute) abundance vector is chosen to be $\mathbf{X}(0) =$
63 $\mathbf{x}(0)$.

64 3) Run the GLV model until it reaches a steady state and normalize the steady-state abundance
65 vector by the total biomass of the sink community to get the final composition of the sink
66 community.

67
68 Consider a particular mixing order π among the total $k!$ mixing orders. Let $\pi(a)$ denote the label
69 of the a -th source in the mixing order. $a, \pi(a) \in \{1, \dots, k\}$. The sink obtained by mixing the k
70 sources in the order π was simulated as follows:

- 71 1) The mixing proportions of the k sources were set to be equal: $m = \frac{1}{k}, a = 1, \dots, k$.
- 72 2) The initial abundance vector of the sink community is determined by the composition of
73 the first source in the order π , i.e., $\pi(1)$, as $\mathbf{X}_0^{(1)} = m \mathbf{y}^{(\pi(1))}$. Then we run the GLV model
74 until it reaches a steady state. Denote the steady-state abundance vector as $\mathbf{X}_{ss}^{(1)}$.
- 75 3) Then the second source $\pi(2)$ arrives. Right after the mixing, the abundance vector of the
76 sink community becomes $\mathbf{X}_0^{(2)} = \mathbf{X}_{ss}^{(1)} + m \mathbf{y}^{(\pi(2))}$. Then we run the GLV model until it
77 reaches a steady state. Denote the steady-state abundance vector as $\mathbf{X}_{ss}^{(2)}$.
- 78 4) Repeat step-3 until all the k sources have been added to the sink. Note that right after the
79 arrival of the k -th source, the abundance vector of the sink community becomes $\mathbf{X}_0^{(k)} =$
80 $\mathbf{X}_{ss}^{(k-1)} + m \mathbf{y}^{(\pi(k))}$. Then we run the GLV model until it reaches a steady state. Denote the
81 steady-state abundance vector as $\mathbf{X}_{ss}^{(k)}$.
- 82 5) Normalize the final steady-state abundance vector $\mathbf{X}_{ss}^{(k)}$ by the total biomass of the sink
83 community to get the final composition of the sink community.

84
85 Since the input data of MST solvers is the OTU count table, for both sink and source communities,
86 we converted the species relative abundances into counts by multiplying the absolute abundances
87 and a fix number (1,000 in all the simulations) and rounding to the nearest integers as the synthetic
88 count data generated by the GLV model.

89
90
91
92
93

94 2. Microbial interactions affect the assembly of the sink community.

95 The deep reason why the existing MST solvers are almost doomed to fail in the presence of
96 microbial interactions is that the true contributions of different sources are only reflected in the
97 sink's initial composition, which will evolve to a final composition following complex ecological
98 dynamics. In general, the final composition will be quite different from the initial one. Here we
99 sketch a proof.

100

101 Let us consider a sink generated by mixing K non-overlapping sources with compositions given
102 by $\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(K)}$, respectively. The initial abundance vector of the sink is denoted as $\mathbf{X}(0) =$
103 $(X_1(0), \dots, X_N(0))$, and its initial composition is given by $\mathbf{x}(0) = (x_1(0), \dots, x_N(0))$ with $x_i(0) =$
104 $X_i(0) / \sum_{i=1}^N X_i(0)$ representing the relative abundance of species i . Note that $\mathbf{x}(0) = \sum_a m_a \mathbf{y}^{(a)}$.
105 Let's assume the population dynamics of the sink community can be represented by a set of
106 ordinary differential equations:

$$107 \quad \dot{\mathbf{X}} = \mathbf{f}(\mathbf{X}; \boldsymbol{\theta}), \quad (1)$$

108 where $\mathbf{X}(t) = (X_1(t), \dots, X_N(t))$ represents the abundance vector at time t , \mathbf{f} is an unspecified
109 nonlinear function with $\boldsymbol{\theta}$ encoding all the ecological parameters, i.e., intrinsic growth rates, and
110 intra- and inter-species interaction strengths. After a small time-step δt , the abundance vector of
111 the sink can be approximated as $\mathbf{X}(\delta t) = \mathbf{X}(0) + \delta t \mathbf{f}(\mathbf{X}(0); \boldsymbol{\theta})$. The ratio of relative abundance
112 for any species pair (i, j) in the initial community is $\alpha(0) = \frac{x_i(0)}{x_j(0)} = \frac{X_i(0)}{X_j(0)}$, while after δt the ratio
113 becomes:

$$114 \quad \alpha(\delta t) = \frac{x_i(\delta t)}{x_j(\delta t)} = \frac{X_i(\delta t)}{X_j(\delta t)} = \frac{X_i(0) + \delta t f_i(\mathbf{X}(0); \boldsymbol{\theta})}{X_j(0) + \delta t f_j(\mathbf{X}(0); \boldsymbol{\theta})}. \quad (2)$$

115 If $X_i(0) = X_j(0)$ and $f_i(\mathbf{X}(0); \boldsymbol{\theta}) = f_j(\mathbf{X}(0); \boldsymbol{\theta})$, then we have $\alpha(\delta t) = \alpha(0)$. But the condition
116 $X_i(0) = X_j(0)$ is too strong to be true. If $X_i(0) \neq X_j(0)$, but $f_i(\mathbf{X}(0); \boldsymbol{\theta}) = X_i(0)g_i(\boldsymbol{\theta})$ and
117 $g_i(\boldsymbol{\theta}) = g_j(\boldsymbol{\theta})$, then we have $\alpha(\delta t) = \alpha(0)$. For a general population dynamics model, this
118 requirement means that there are no inter-species interactions and the intrinsic growth rates of
119 different species are identical, which is also too strong to be true. Hence, in general $\alpha(\delta t) \neq \alpha(0)$,
120 and the final composition of the sink community will be quite different from its initial composition.

121

122

123 3. Priority effects affect the assembly of the sink community.

124 Consider three source communities S_1 , S_2 and S_3 . Let's assume species- i is only present in the
 125 source S_1 and species- j is only present in the source S_2 . We mix the source communities in 6
 126 different orders but with identical mixing proportions $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$. For each mixing order, we assume
 127 the arrival time of the three sources are 0, τ , 2τ , respectively, where τ is a constant (and is large
 128 enough for the resulting sink community to reach a steady state). Suppose we use the composition
 129 of the final sink community taken at 3τ to estimate the contribution of each source. We want to
 130 prove that at time $t = 3\tau$, the sink communities resulting from different mixing orders will have
 131 different compositions, even in the absence of any microbial interactions. To achieve that, let's
 132 compute the ratio between the relative abundance of species- i and that of species- j in the final sink
 133 community at time $t = 3\tau$, i.e., $\alpha_{ij}(3\tau) = \frac{x_i(3\tau)}{x_j(3\tau)} = \frac{X_i(3\tau)}{X_j(3\tau)}$.

134
 135 Consider a particular mixing order $S_1 \rightarrow S_2 \rightarrow S_3$. In the absence of any inter- or intra-species
 136 interactions, species will grow exponentially. Hence, at time $t = 3\tau$, the abundance of a species- i
 137 (which is only present in the source S_1) is given by: $X_i(3\tau) = m_1 X_i(0) \exp(3\tau r_i)$, where $m_1 = \frac{1}{3}$
 138 is the mixing proportion (contribution) of the source S_1 , $X_i(0)$ is the initial abundance of species-
 139 i in the source S_1 , r_i is the intrinsic growth rate of species- i . Similarly, at time $t = 3\tau$, the
 140 abundance of species- j (which is assumed to be only present in the source S_2) is given by:
 141 $X_j(3\tau) = m_2 X_j(0) \exp(2\tau r_j)$. So, we have

$$142 \quad \alpha_{ij}^{123}(3\tau) = \frac{m_1 X_i(0) e^{3\tau r_i}}{m_2 X_j(0) e^{2\tau r_j}} = \alpha_{ij}(0) e^{\tau(3r_i - 2r_j)},$$

143 where the superscript '123' indicates the mixing order $S_1 \rightarrow S_2 \rightarrow S_3$. We can repeat the above
 144 calculation for different mixing orders. The results are summarized here:

$$145 \quad \alpha_{ij}^{132}(3\tau) = \frac{m_1 X_i(0) e^{3\tau r_i}}{m_2 X_j(0) e^{\tau r_j}} = \alpha_{ij}(0) e^{\tau(3r_i - r_j)},$$

$$146 \quad \alpha_{ij}^{213}(3\tau) = \frac{m_1 X_i(0) e^{2\tau r_i}}{m_2 X_j(0) e^{3\tau r_j}} = \alpha_{ij}(0) e^{\tau(2r_i - 3r_j)},$$

$$147 \quad \alpha_{ij}^{231}(3\tau) = \frac{m_1 X_i(0) e^{\tau r_i}}{m_2 X_j(0) e^{3\tau r_j}} = \alpha_{ij}(0) e^{\tau(r_i - 3r_j)},$$

$$148 \quad \alpha_{ij}^{312}(3\tau) = \frac{m_1 X_i(0) e^{2\tau r_i}}{m_2 X_j(0) e^{\tau r_j}} = \alpha_{ij}(0) e^{\tau(2r_i - r_j)},$$

149
$$\alpha_{ij}^{321}(3\tau) = \frac{m_1 X_i(0) e^{\tau r_i}}{m_2 X_j(0) e^{2\tau r_j}} = \alpha_{ij}(0) e^{\tau(r_i - 2r_j)}.$$

150 Note that if the three sources were mixed simultaneously, then we have

151
$$\alpha_{ij}^{\text{simultaneous}}(3\tau) = \frac{m_1 X_i(0) e^{3\tau r_i}}{m_2 X_j(0) e^{3\tau r_j}} = \alpha_{ij}(0) e^{3\tau(r_i - r_j)}.$$

152 Therefore, even in the absence of any microbial interactions, different mixing patterns will result
153 in different final compositions of the sink community, which are also different from that obtained
154 by simultaneous mixing.

155

156 4. Community coalescence experiments.

157 Stool samples from healthy human donors were collected and immediately transferred into the
158 anaerobic workstation (85% N₂, 10% H₂ and 5% CO₂, COY). 10g stool samples were suspended
159 into 50 mL 20% glycerol (in sterile phosphate-buffered saline, with 0.1% L-cysteine
160 hydrochloride). The samples were homogenized by vortexing and then filtered with sterile nylon
161 mesh to remove large particles in fecal matter. Aliquots of the suspension were placed in sterile
162 cryogenic vials and frozen at -80 °C for long-term storage until use.

163

164 Stool samples of 24 individuals were used for the community coalescence experiments. To
165 generate 481 sink communities, samples from two, three or four different individuals were mixed
166 with equal volume. 20 uL stool mixture was inoculated into 980 uL medium in 96-well plates
167 (PCR-96-SG-C, Axygen) for static culturing at 37 °C in the anaerobic workstation. The medium
168 used for *ex vivo* culture was modified from previous studies, which comprises: peptone water
169 (2.0 g/L, CM0009, Thermo Fisher), yeast extract (2.0 g/L, LP0021B, Thermo Fisher), L-cysteine
170 hydrochloride (1 g/L), Tween 80 (2 mL/L), hemin (5 mg/L), vitamin K1 (10 µL/L), NaCl (1.0 g/L),
171 K₂HPO₄ (0.4 g/L), KH₂PO₄ (0.4 g/L), MgSO₄·7H₂O (0.1 g/L), CaCl₂·2H₂O (0.1 g/L), NaHCO₃
172 (4 g/L), porcine gastric mucin (4 g/L, M2378, Sigma-Aldrich), sodium cholate (0.25 g/L) and
173 sodium chenodeoxycholate (0.25 g/L)³. *Ex vivo* culture of gut microbial communities was
174 transferred into fresh medium every 24h (1:200 dilution), for a total of 10 transfers. After each
175 transfer, samples were centrifuged to remove the supernatant and the pellets were stored at -80°C
176 with a plastic seal until DNA extraction.

177

178 The initial stool and *ex vivo*-cultured samples after 10 passages were sequenced. For stool samples,
179 DNA was extracted using the QIAamp Power Fecal Pro DNA Kit (Qiagen) according to the
180 manufacturer's instructions. For cultured samples, DNA extraction (DNeasy UltraClean 96
181 Microbial Kit, Qiagen) and 16S amplicon library preparation were performed by an automated
182 protocol at Tecan Freedom EVO 200. V3-V4 region of 16S rRNA gene was amplified using
183 primers 341F 5'-CCTACGGGNGGCWGCAG -3' and 805R 5'-
184 GACTACHVGGGTATCTAATCC-3' with custom barcodes⁴. Libraries were further pooled
185 together at equal molar ratios and sequenced by Illumina NovaSeq (250 bp paired-end reads) at
186 Novogene Technology (Tianjin, China).

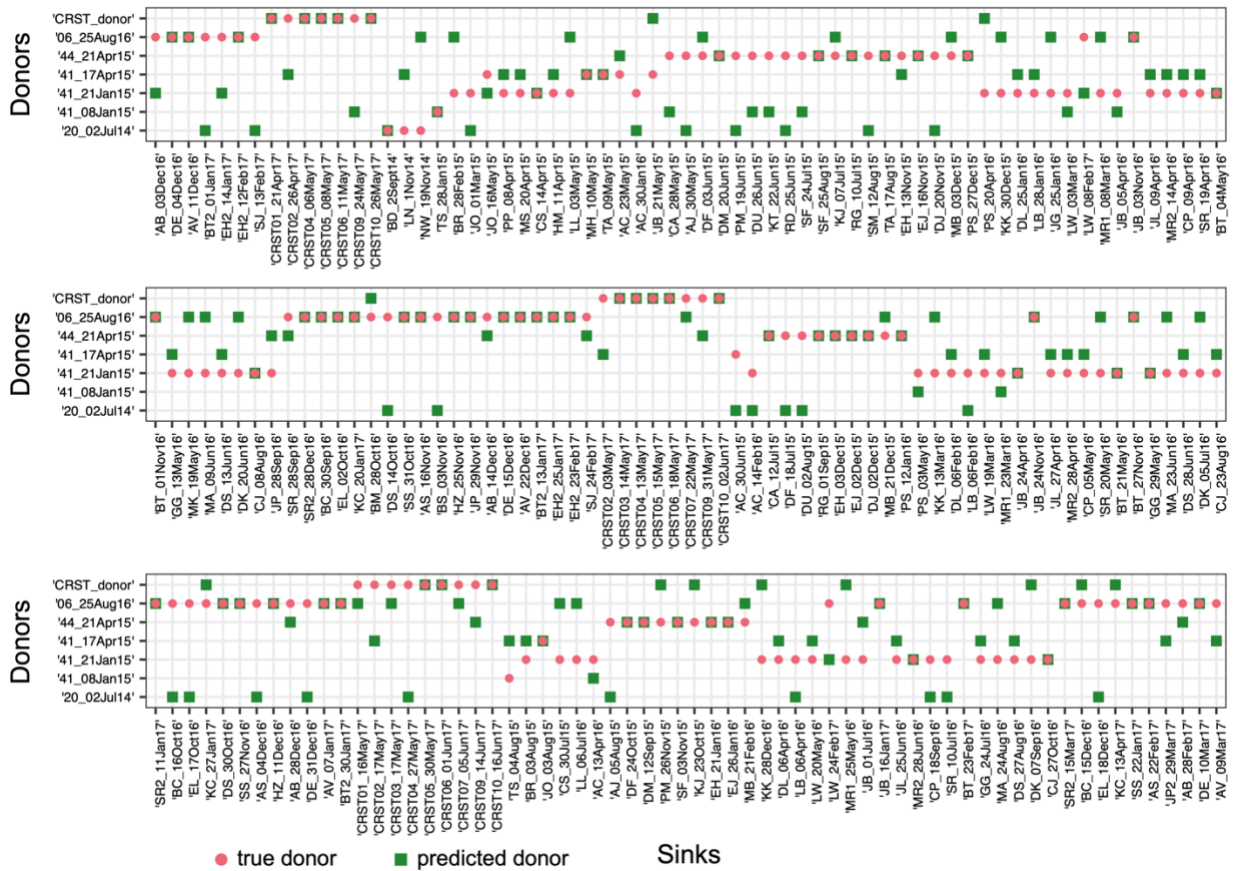
187
188 16S amplicon sequencing data were analyzed by QIIME2 (version 2020.2)⁵. Primers of the raw
189 sequence data were cut with Cutadapt (via q2-cutadapt)⁶. Quality control was performed by
190 DADA2 (via q2-dada2)⁷. All amplicon sequence variants (ASVs) from DADA2 were used to
191 construct a phylogenetic tree with fasttree2 (via q2-phylogeny)⁸. The ASVs were assigned to
192 taxonomy with naïve Bayes classifier (via q2-feature-classifier)⁹ against the SILVA database
193 (SILVA_132_SSURef_Nr99). The ASV table was normalized, and rare ASVs (all features with a
194 total abundance of less than 10 and present in only a single sample) were filtered out.

195
196
197
198
199
200
201
202
203
204
205
206
207
208
209

210 References

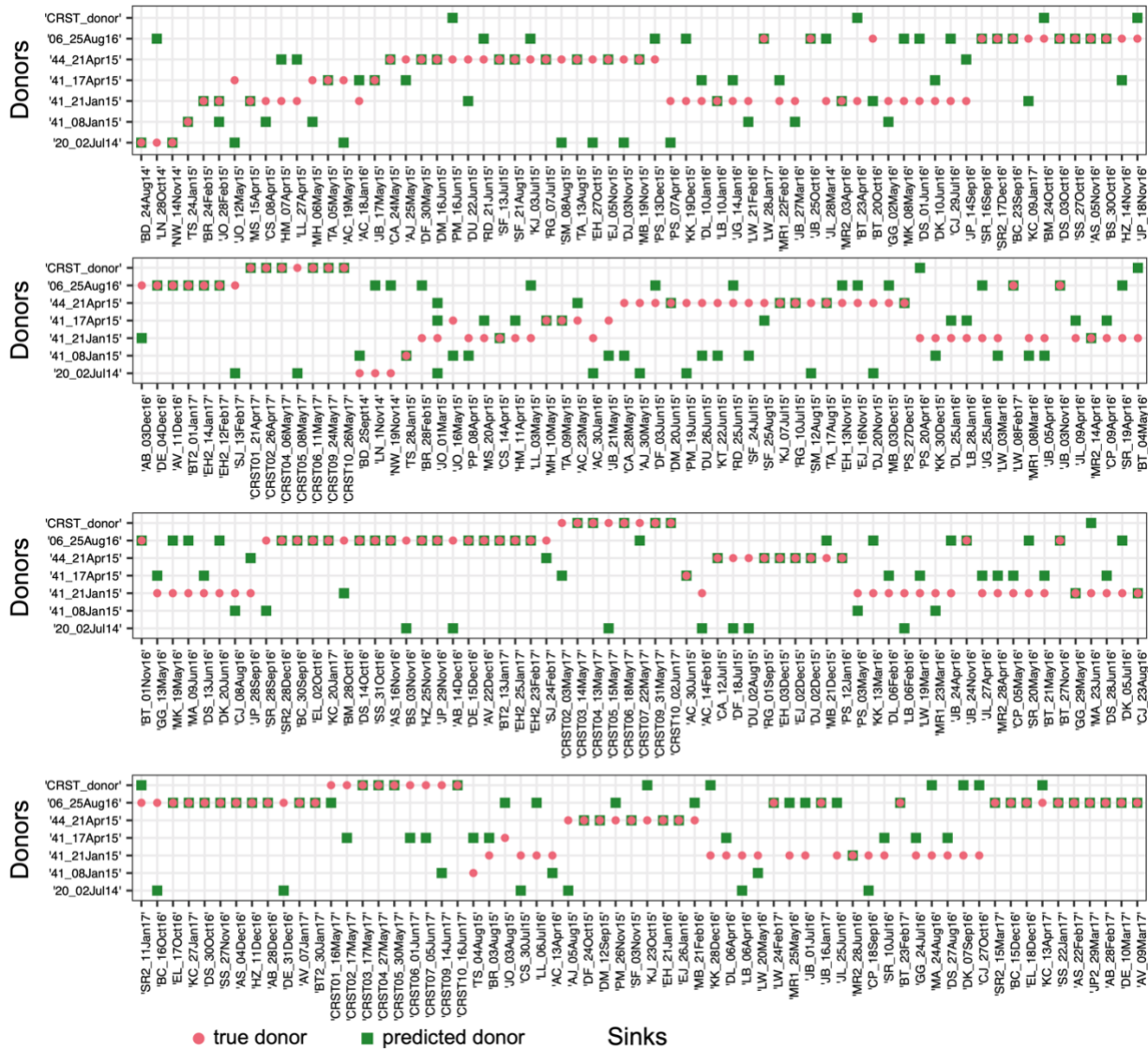
- 211 1. Case, T. J. Illustrated guide to theoretical ecology. *Ecology* **80**, 2848–2848 (1999).
- 212 2. Erdos, P. & Rényi, A. On the evolution of random graphs. *Publ Math Inst Hung Acad Sci* **5**, 17–
- 213 60 (1960).
- 214 3. Li, L. *et al.* An in vitro model maintaining taxon-specific functional activities of the gut
- 215 microbiome. *Nat. Commun.* **10**, 1–11 (2019).
- 216 4. Mizrahi-Man, O., Davenport, E. R. & Gilad, Y. Taxonomic classification of bacterial 16S rRNA
- 217 genes using short sequencing reads: evaluation of effective study designs. *PloS One* **8**,
- 218 e53608 (2013).
- 219 5. Bolyen, E. *et al.* Reproducible, interactive, scalable and extensible microbiome data science
- 220 using QIIME 2. *Nat. Biotechnol.* **37**, 852–857 (2019).
- 221 6. Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing reads.
- 222 *EMBnet J.* **17**, 10–12 (2011).
- 223 7. Callahan, B. J. *et al.* DADA2: high-resolution sample inference from Illumina amplicon data.
- 224 *Nat. Methods* **13**, 581–583 (2016).
- 225 8. Price, M. N., Dehal, P. S. & Arkin, A. P. FastTree 2—approximately maximum-likelihood trees
- 226 for large alignments. *PloS One* **5**, e9490 (2010).
- 227 9. Bokulich, N. A. *et al.* Optimizing taxonomic classification of marker-gene amplicon sequences
- 228 with QIIME 2's q2-feature-classifier plugin. *Microbiome* **6**, 1–17 (2018).
- 229 10. Staley, C. *et al.* Predicting recurrence of *Clostridium difficile* infection following
- 230 encapsulated fecal microbiota transplantation. *Microbiome* **6**, 166 (2018).
- 231 11. Sharon, G. *et al.* Human Gut Microbiota from Autism Spectrum Disorder Promote
- 232 Behavioral Symptoms in Mice. *Cell* **177**, 1600–1618.e17 (2019).

233
234



235
 236 **Figure S1: Evaluation of FEAST using FMT data from Staley et al.¹⁰** True donor (red cycle)
 237 vs. predicted donor (green square) of each recipient. For each post-FMT community
 238 all the 7 donors, we referred to the one whose fecal sample has the highest contribution estimated
 239 by FEAST as the “predicted donor”. In Fig.4c, we presented results of the first 65 sinks. Here, we
 240 showed the results of the remaining 194 sinks. Sources: microbiome samples of donors and the
 241 pre-FMT samples of recipients; Sinks: post-FMT samples of recipients.

242
 243
 244
 245
 246
 247
 248
 249
 250
 251
 252



253

254 **Figure S2: Evaluation of SourceTracker using FMT data from Staley et al.**¹⁰ True donor (red

255 cycle) vs. predicted donor (green square) of each recipient. For each post-FMT community (sink),

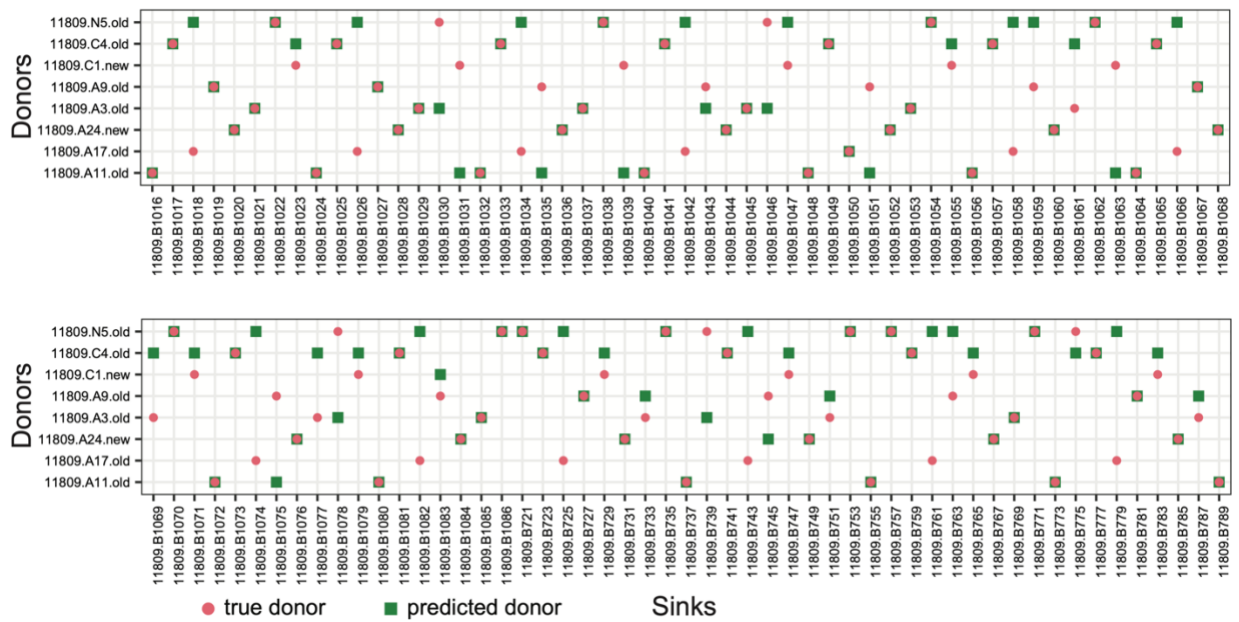
256 among all the 7 donors, we referred to the one whose fecal sample has the highest contribution

257 estimated by SourceTracker as the “predicted donor”. Sources: microbiome samples of donors and

258 the pre-FMT samples of recipients; Sinks: post-FMT samples of recipients.

259

260



261

262 **Figure S3: Evaluation of FEAST using FMT data from Sharon et al.**¹¹ True donors (red cycle)

263 vs. the predicted donor (green square) of each recipient sink given by FEAST using the source and

264 sink compositions as the input. For each post-FMT community (sink), among all the 8 donors, we

265 referred to the one whose fecal sample has the highest contribution estimated by FEAST as the

266 “predicted donor”. Sources: microbiome samples of donors and the pre-FMT samples of recipients;

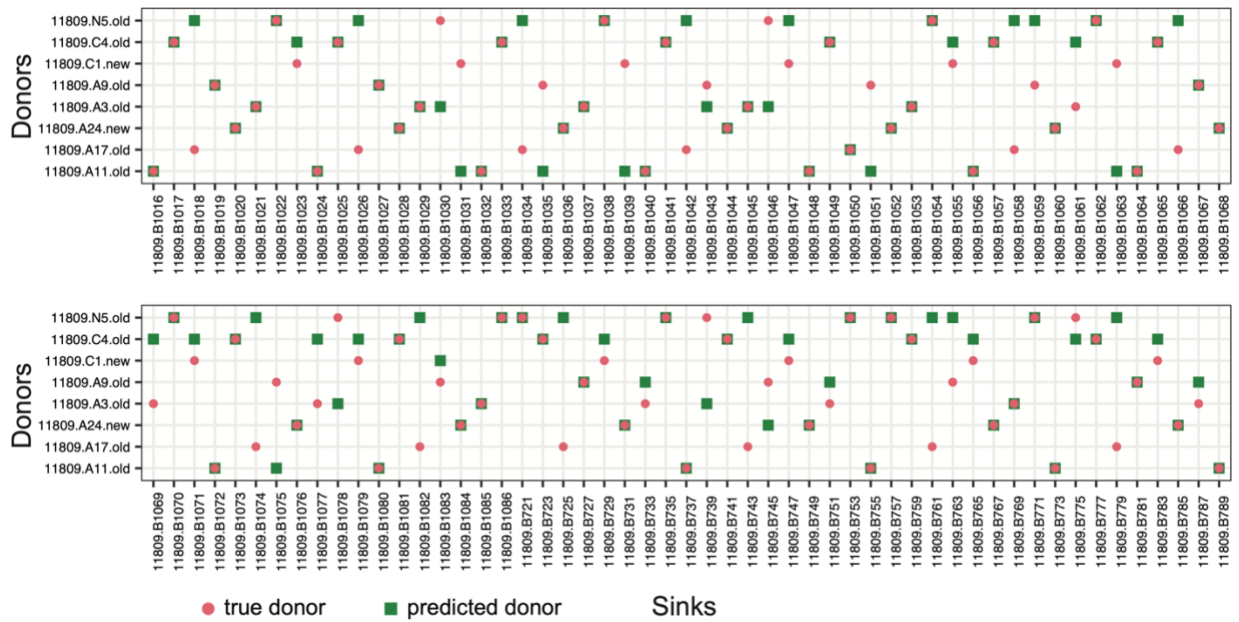
267 Sinks: post-FMT samples of recipients. In total, there are 106 sinks.

268

269

270

271



272

273 **Figure S4: Evaluation of SourceTracker using FMT data from Sharon et al.¹¹ True donors**
274 (red cycle) vs. the predicted donor (green square) of each recipient sink given by SourceTracker
275 using the source and sink compositions as the input. For each post-FMT community (sink), among
276 all the 8 donors, we referred to the one whose fecal sample has the highest contribution estimated
277 by SourceTracker as the “predicted donor”. Sources: microbiome samples of donors and the pre-
278 FMT samples of recipients; Sinks: post-FMT samples of recipients. In total, there are 106 sinks.

279

280

281

282

283

284

285

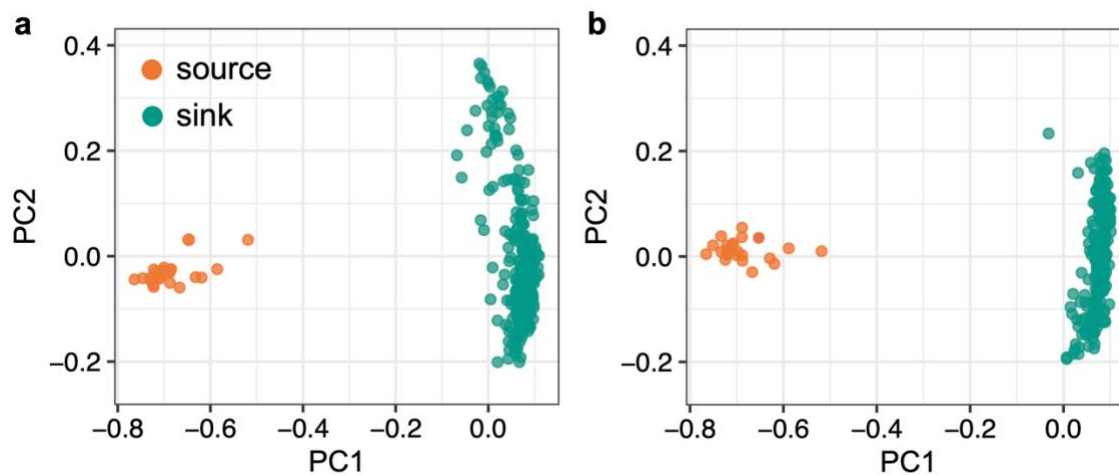
286

287

288

289

290



291

292 **Figure S5:** Principal Coordinate Analysis (PCoA) plot of the sinks and sources in the community
293 coalescence experiments. **a**, Pairwise mixing. **b**, Quadruple-wise mixing.

294

295

296

297

298

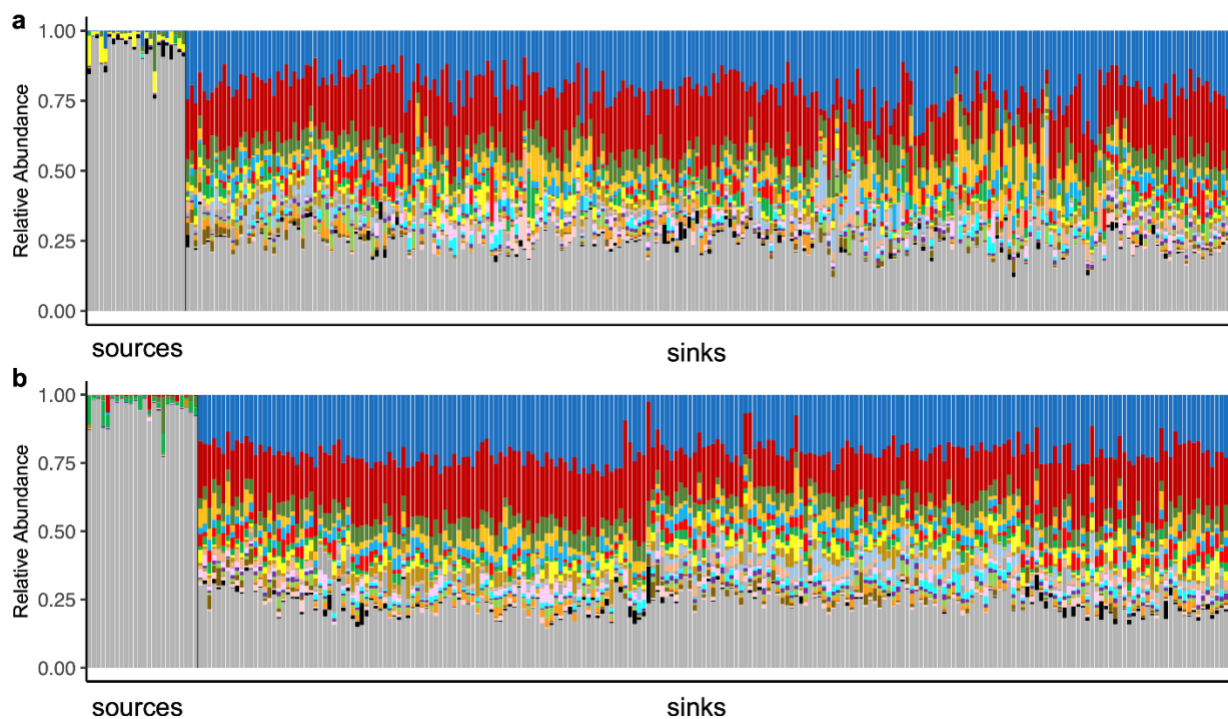
299

300

301

302

303

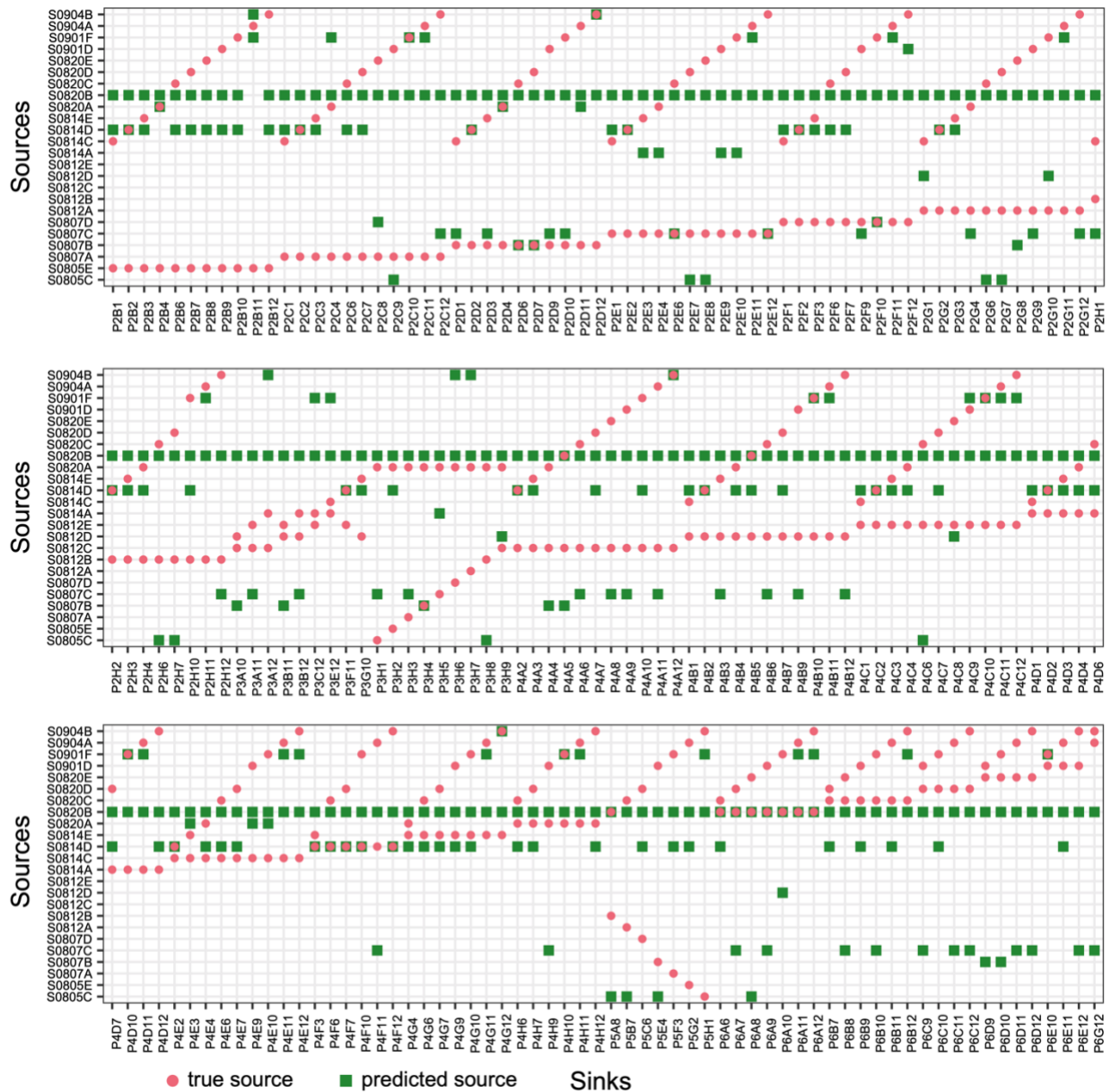


304

305 **Figure S6: Taxonomic profiles of sources and sinks in community coalescence experiments.**

306 **a**, Pairwise mixing. **b**, Quadruple-wise mixing. For visualization purposes, we only showed the
307 top-20 abundant ASVs. All other ASVs were grouped together and shown in gray. We found that
308 some highly abundant ASVs in the source communities have very low abundances in the sink
309 communities, whereas some low-abundance ASVs in source communities flourish in the sink
310 communities. Also, a few ASVs in the sinks were not detected in the sources, indicating that either
311 their relative abundances were below the detection limit or there was potential contamination.

312



313

314 **Figure S7: Performance of FEAST in identifying sources in pairwise community coalescence**

315 **experiments.** True sources (red cycles) vs. predicted sources (green squares) of each sink. For

316 each sink, among the 24 known sources, the two sources with the top-two largest contributions

317 predicted by FEAST were referred to as the predicted sources. In Fig.5, we showed the results of

318 the first 64 sinks. Here we showed the results of the remaining 192 sinks.

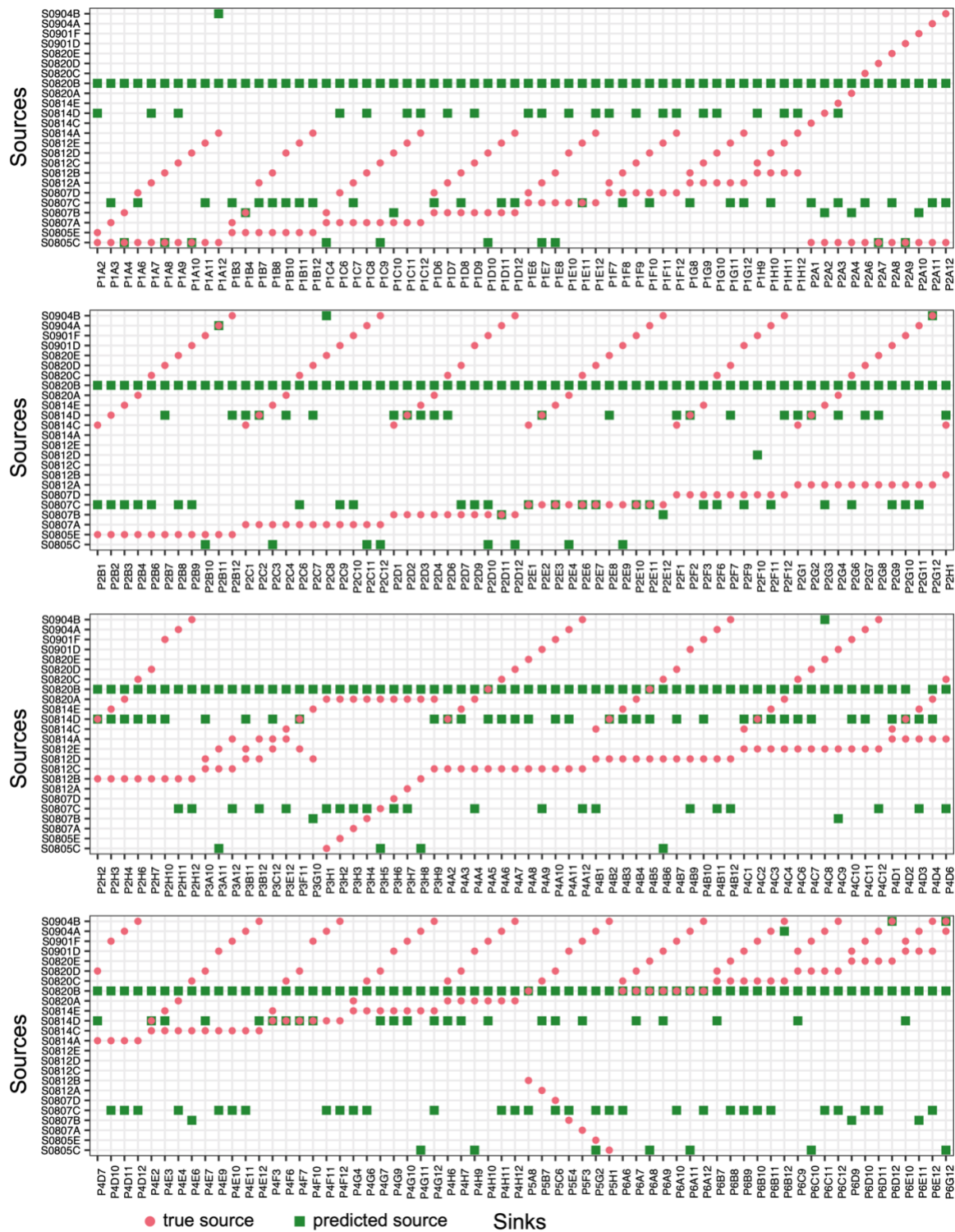
319

320

321

322

323



324

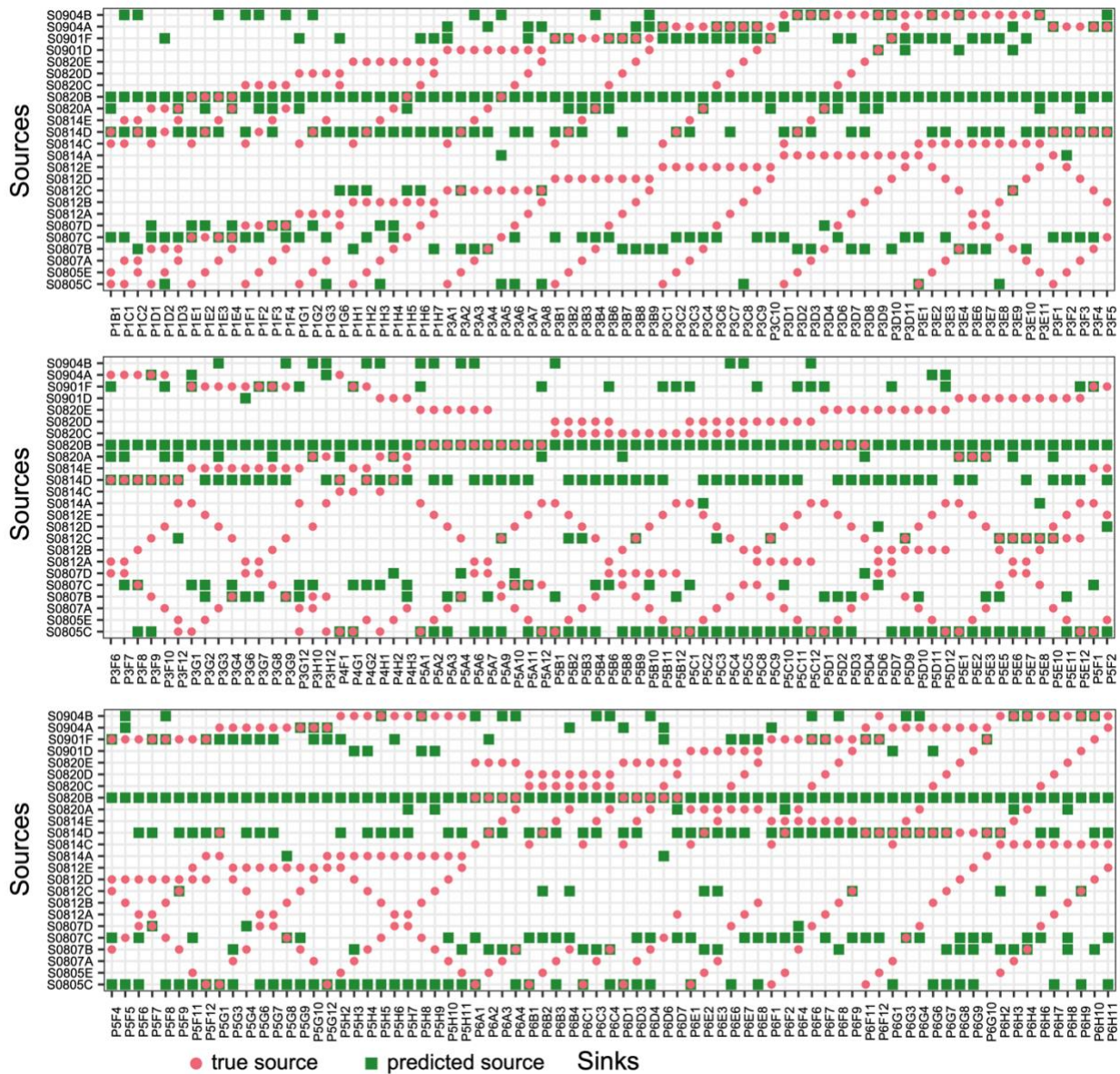
325 **Figure S8: Performance of SourceTracker in identifying sources in pairwise community**

326 **coalescence experiments.** True sources (red cycles) vs. predicted sources (green squares) of each

327 sink. For each sink, among the 24 known sources, the sources with the top-two largest

328 contributions predicted by SourceTracker were referred to as the predicted sources.

329



330

331 **Figure S9: Performance of FEAST in identifying sources in quadruple-wise community**

332 **coalescence experiments.** There are 24 source communities (stool samples from 24 healthy

333 individuals). Each sink community is obtained by mixing four different source communities *ex*

334 *vivo*. The final composition of each sink was obtained from metagenomic sequencing of samples

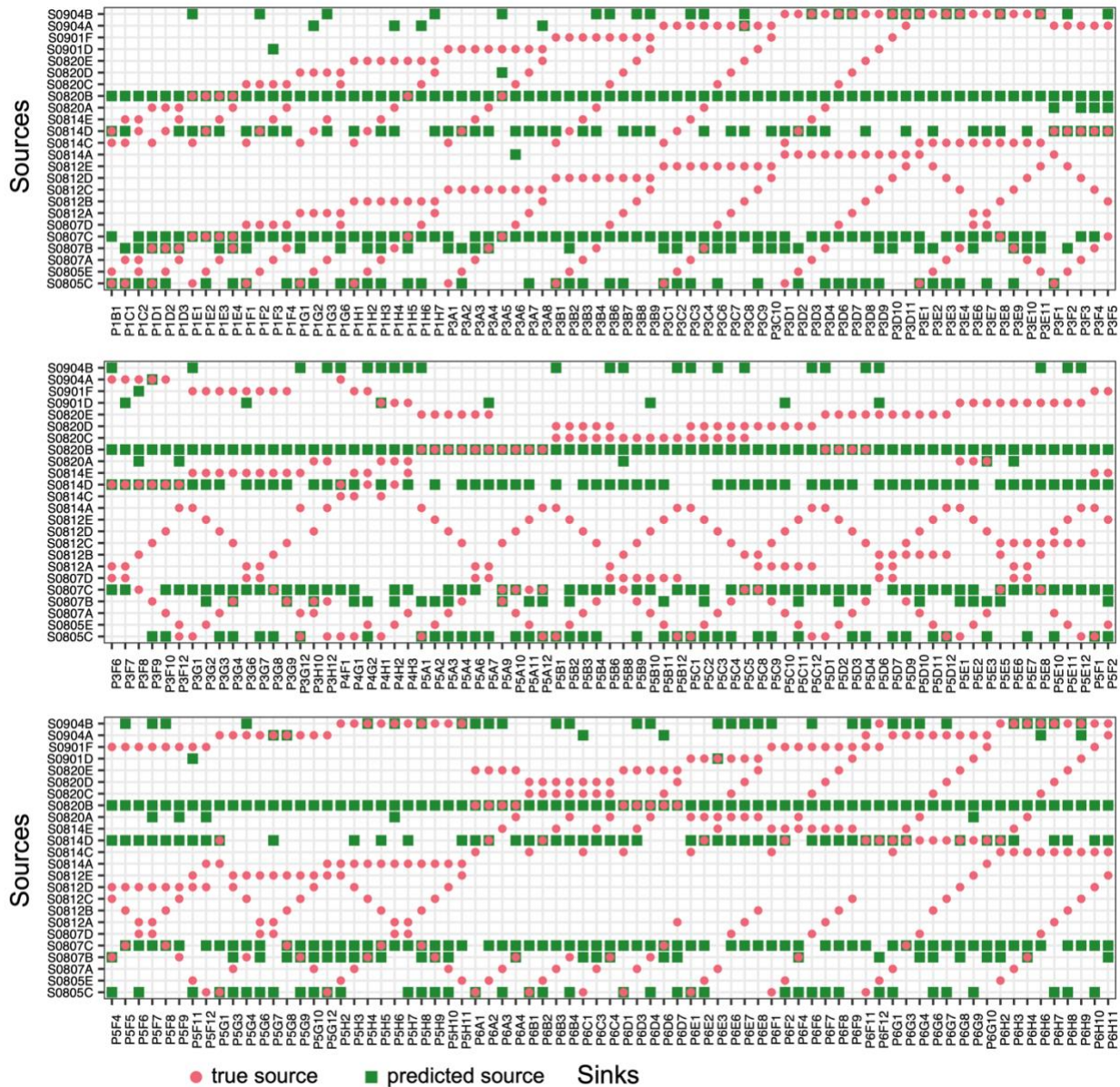
335 collected after 11 days of the *ex vivo* mixing. True sources (red circles) vs. predicted sources (green

336 squares) of each sink. (Each row includes 75 sinks). For each sink, among the 24 known sources,

337 the four sources with the top-four largest contributions predicted by FEAST were referred to as

338 the predicted sources.

339



340

341 **Figure S10: Performance of SourceTracker in identifying sources in quadruple-wise**

342 **community coalescence experiments.** There are 24 sources communities (stool samples from 24

343 healthy individuals). Each sink community is obtained by mixing four different source

344 communities *ex vivo*. The final composition of each sink was obtained from metagenomics

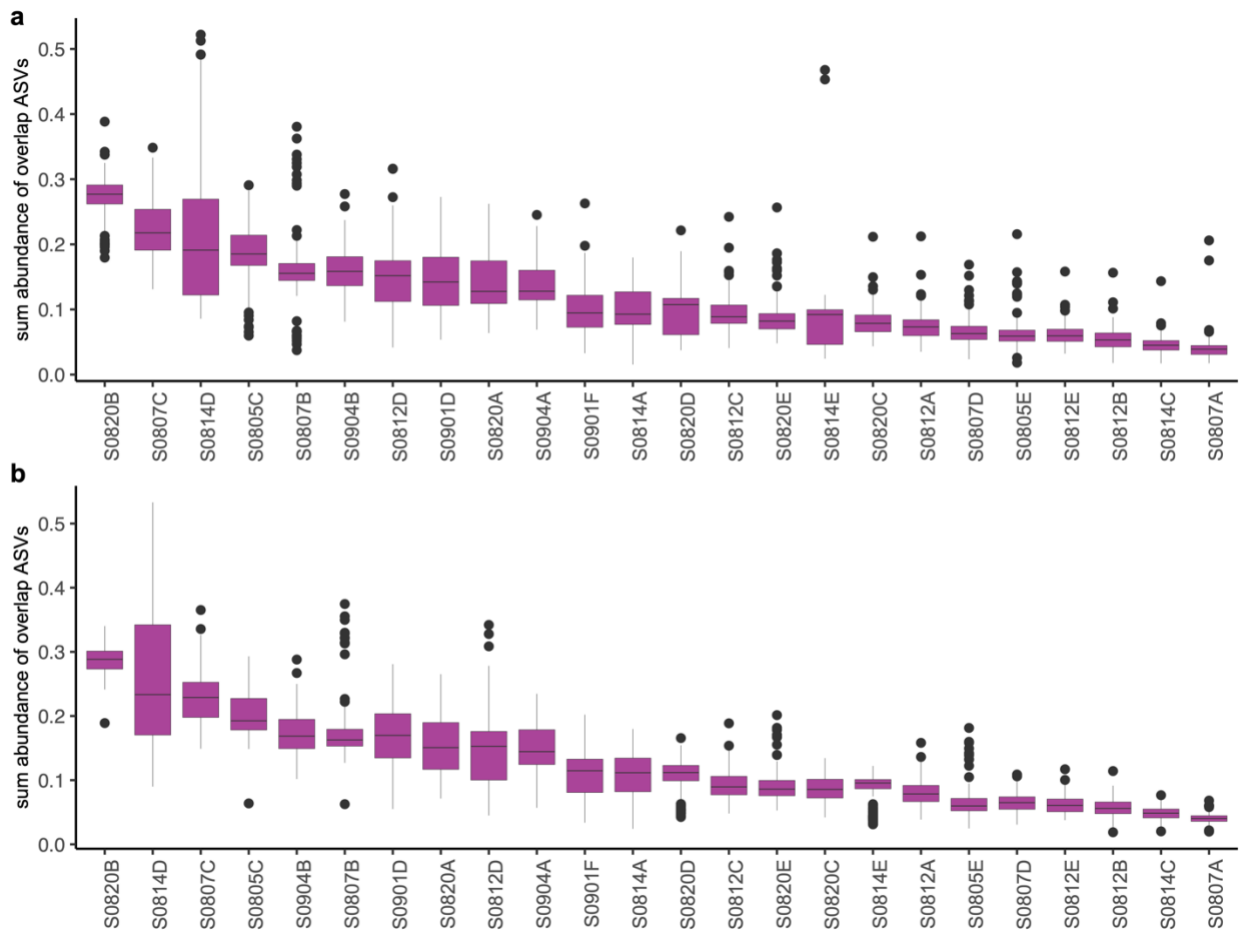
345 sequencing of samples collected after 11 days of the *ex vivo* mixing. True sources (red cycles) vs.

346 predicted sources (green squares) of each sink. (Each row includes 75 sinks). For each sink, among

347 the 24 known sources, the four sources with the top-four largest contributions predicted by

348 SourceTracker were referred to as the predicted sources.

349



350

351 **Figure S11: Relative abundances of common ASVs shared by sources and sinks.** For each
352 sink and source pair, we identified their common ASVs and calculated the total relative abundance
353 of those common ASVs. Each boxplot represents the total relative abundance of common ASVs
354 shared by this source and each of the 256 sinks in the pairwise community coalescence experiments
355 (a); and 225 sinks in quadruple-wise community coalescence experiments (b).

356

357

358

359