

1 **Title**

2 **Metatranscriptome of human fecal microbial communities in a cohort**
3 **of adult men**

4 **Supplementary information – SI**

5
6 Galeb S. Abu-Alt^{a,c,1}, Raaj S. Mehta^{d,1}, Jason Lloyd-Price^{a,c}, Himel Mallick^{a,c}, Toby Branck^{a,g},
7 Kerry L. Ivey^{b,h}, David A. Drew^{d,e}, Casey DuLong^a, Eric Rimm^{b,i}, Jacques Izard^f, Andrew T.
8 Chan^{c,d,e,i,1,2}, Curtis Huttenhower^{a,c,1,2}

9 ^aBiostatistics Department and ^bDepartment of Nutrition, Harvard T. H. Chan School of Public
10 Health, Boston, MA 02115; ^cThe Broad Institute, Cambridge, MA 02142; ^dClinical and
11 Translational Epidemiology Unit, ^eDivision of Gastroenterology, Massachusetts General Hospital
12 and Harvard Medical School, Boston, MA 02114; ^fUniversity of Nebraska, Lincoln, 1901 North
13 21 Street, Lincoln, NE 68588; ^gU.S. Army Natick Soldier Systems Center in Natick, MA 01760,
14 ^hSouth Australian Health and Medical Research Institute, Infection and Immunity Theme, School
15 of Medicine, Flinders University, Adelaide, Australia, 5000; ⁱChanning Division of Network
16 Medicine, Department of Medicine, Brigham and Women's Hospital
17

18 ¹Contributed equally to this work

19 ²To whom correspondence should be addressed:

20 Curtis Huttenhower
21 Harvard T. H. Chan School of Public Health
22 Department of Biostatistics
23 655 Huntington Avenue
24 Building 1, Room 413
25 Boston, Massachusetts 02115
26 Phone: 617.432.4912
27 chuttenh@hsph.harvard.edu

28 Andrew T. Chan
29 Massachusetts General Hospital
30 Clinical and Translational Epidemiology Unit
31 55 Fruit St, GRJ-825C
32 Boston MA 02114
33 achan@mgh.harvard.edu

34

35 **Supplementary Notes**

36 **Viruses detected from non-enriched metagenomes and metatranscriptomes do not reflect** 37 **variation of bacterial taxa**

38 A small number of both DNA and RNA viruses were quantified confidently by MetaPhlan2, which
39 is likely an underestimate of the gut virome diversity since our extraction protocol did not enrich
40 for virus-like particles. Of the 30 DNA viral species detected in the cohort, 29 were bacteriophage
41 belonging primarily to the Siphoviridae (19) family, with few Myoviridae (6), Podoviridae (2), and
42 Inoviridae (1) members, and one *Escherichia* phage of unknown taxonomy. The most common
43 putative phage identifications were *C2likevirus* (in 116 participants), *Epsilon15likevirus* (in 16
44 participants), and *Lactobacillus Lc Nu* (in 16 participants), however, no correlation was found
45 between the abundance distribution of these Caudovirales phages and their natural host genera
46 *Lactococcus*, *Escherichia*, and *Lactobacillus*; although power for detecting this is low given non-
47 virally-enriched detection rates. In metatranscriptomes, we identified 88 RNA viral species
48 belonging to 19 families. Apart from rare Leviviridae and Iflaviridae members, these were all plant
49 viruses, in agreement with previous studies¹ and possibly due to dietary ingestion. Most RNA
50 viruses were detected in <30 (10%) metatranscriptomes with the exception of *Pepper mild mottle*
51 *virus* (48% prevalence) and *Tomato mosaic virus* (39% prevalence), which together accounted
52 for 31% of viral RNA on average (when present). Although gut viral ecology is more difficult to
53 analyze than that of the bacteriome due to inadequate viral reference sequences², these methods
54 allow for some incidental analysis of DNA phage and RNA plant viruses in human gut
55 metagenomes and metatranscriptomes.

56 **Effect of GC content and ORF length on transcription ratios.**

57 We analytically evaluated the effect of GC content and ORF length on transcription ratios, finding
58 no interaction (**Supplementary Fig. 5**). The 430 structured MetaCyc pathways analyzed here
59 were quantified from 808,694 UniRef90 gene families (3.4% of the total UniRef90 database) that
60 had detectable DNA in at least one sample in our study. Of those, 89,991 and 44,792 UniRef90s
61 had non-zero DNA and RNA abundance in at least one sample, respectively. This resulted in
62 37,085 pathway-associated UniRef90 gene families for which RNA/DNA ratios were calculated.
63 Among these genes, there were no significant differences in ratio when stratified by either %GC
64 or length of nucleotide sequences; when tested continuously, there was no significant correlation
65 of either length or %GC with RNA abundances. Intriguingly, a very low effect size (Pearson's $r <$
66 0.04), but significant, correlation was observed between sequence length and DNA abundance.
67 Given the comparability of DNA and RNA protocols in this study, it is not clear why this might
68 have arisen, but at such a minimal effect size it does not affect the study's conclusions.

69 **Core and variable fecal metatranscriptomes differ from the metagenome.**

70 The distribution of transcript abundance ranged over four orders of magnitude among 210
71 pathways that were transcribed in >10% of samples (**Supplementary Fig. 6F**). The highest
72 transcription ratios consistently arose from pathways that were both low prevalence and
73 taxonomically restricted, e.g. archaeal methanogenesis and coenzyme 420 biosynthesis, as
74 previously suggested by our pilot study³. Following energy metabolism and fermentation, which
75 tended to dominate in both prevalence and expression levels, the highest metatranscription was
76 observed for saturated and unsaturated fatty acid elongation pathways, albeit in less than one-
77 third of samples. Fatty acids are generated from acetyl-CoA, which in turn is produced mainly
78 during glycolytic energy release, and together this may explain the concerted metatranscription
79 of glycolysis and energy-expensive fatty acid elongation. As the primary role for bacterial fatty
80 acids is to serve as precursors for cell membrane building blocks (e.g. phospholipids), this likely

81 signals widespread cell growth in the typical fecal microbiota⁴. On the other end of the spectrum,
82 pathways with the lowest metatranscription had mean RNA abundances below their
83 corresponding DNA relative abundances, with prevalence of metatranscription ranging between
84 15% (sulfate assimilation/cysteine synthesis) and 95% (peptidoglycan synthesis) of samples. This
85 low tail of metatranscription included several amino acid synthesis pathways, including
86 methionine, homoserine, aromatic and seleno-amino acids, followed by cofactor biosynthesis,
87 including thiamine (and variants), tetrapyrrole, etc. Prevalent metatranscription of degradation of
88 stachyose (PWY-6527), a legume-derived non-digestible tetrasaccharide that promotes SCFA
89 producers, may reflect diet preferences. Together, these findings would underline that the fecal
90 microbiome does not prioritize *de novo* synthesis of amino acids or widespread activation of
91 specialized functions, yet displays high dynamic range and milieu activities such as transformation
92 of phenolics, stress adaptation, and secondary metabolism.

93 **Genetic divergence patterns of stool-associated bacterial strains is species-specific and** 94 **preserved among host populations**

95 Nucleotide substitution rates within and between cohorts were strikingly similar for the compared
96 species, indicating that species' evolutionary strategies within the stool niche were comparable
97 between these host populations (**Fig. 6C**). The amount of genetic change was higher for
98 Firmicutes than Bacteroidetes and did not appear to be simply a function of species prevalence
99 in the two cohorts. For example, *Bacteroides dorei* and *uniformis*, and *Alistipes putredinis* had
100 comparable prevalence with *Ruminococcus bromii*, *Dialister invisus*, and *Eubacterium rectale*, yet
101 appreciably fewer nucleotide substitutions between strains. This may be due to *Bacteroides*
102 species' more restrictive definitions by systematics⁵, serving as a reminder that culture-based
103 isolate information and culture-independent microbial profiling may need further resolution as
104 strain and transcriptional meta'omics are explored.

105 **Species-function relationships in fecal meta'omes.**

106 We quantified how tightly was each pathway coexpressed - that is, the extent to which the multiple
107 enzymes making up each pathway were expressed at similar abundances within each organism
108 and meta'ome (**Supplementary Dataset 4**). This was assessed using the average variance of
109 gene families' transcription log ratios across samples (see **Methods**), here termed the EC
110 dispersion. The distribution of dispersions from all pathways' ECs was significantly below 1 (one-
111 sided t-test $P=1.1 \times 10^{-16}$), with a mean of 0.89, indicating that functionally-related genes are co-
112 expressed on average. Tightly coexpressed pathways (low dispersion) included methanogenesis
113 (dispersion 0.26), two pathways for L-histidine degradation (0.38, 0.39), and degradation of the
114 glutaryl-CoA (0.49) intermediate of tryptophan and lysine metabolism. Tryptophan and histidine
115 are among the energetically most expensive amino acids to synthesize^{6,7}, for which tight co-
116 expression of degradation pathway is not surprising. No evidence was found for a relationship
117 between EC dispersion and the number of species that transcribed the pathway (Spearman rho -
118 0.01). Differences between pathways that were considered a part of the core or variably
119 expressed metatranscriptome were also not detected (Wilcoxon rank-sum test p-value 0.10).

120

121 **Supplementary Discussion**

122 We briefly review here the current literature on the topic of microbiome sample stabilization with
123 RNAlater. The reported minor effect of choice of sample handling method on microbiome
124 composition⁸ lacks testing for statistical significance of any variance, suggesting that between-
125 condition variation in that study was comparable to replicate variation (and much smaller than
126 population variability). In addition to our own validation work³, which indicates a negligible effect
127 of RNAlater on microbial community composition, there are numerous reports on the evaluation
128 of methods for storage and handling of microbiome samples in a cohort setting. These studies
129 reveal a lack of significant alteration in community structure between samples preserved with
130 RNAlater, ethanol, lyophilization, fecal occult blood test cards, and freezing at -80°C⁹⁻¹⁴.

131 The recent Choo et al. study¹⁵ reports a statistically significant effect of storage method on
132 microbial community composition, based on assessing the variation among differently stored
133 samples collected from one individual sampled three times over 30 days. However, Choo et al.
134 report that the variation attributable to storage method is markedly smaller than the variation
135 explained by different sample time points (i.e. smaller than intra-individual difference, which in
136 turn is far smaller than inter-individual difference). Notably, the variation introduced by RNAlater,
137 relative to freezing at -80°C, was comparable to that introduced by sample storage in
138 OMNIGene.Gut, another popular sample stabilization kit. Furthermore, significant but localized
139 differences in taxon abundance relative to freezing were comparable among RNAlater and
140 OmniGene.Gut. Finally, Choo et al evaluated stool samples from a single subject, which are not
141 likely to be representative at the population level.

142

143

144

145 **Supplementary Figures**

146 Please note that **Supplementary Figures 6-10** are multi-panel, multi-page figures that were
147 submitted as separate files due to their length and size.

148 **Supplementary Figure 1. Taxonomic profiles of gut metagenome ecology and stability.**
149 Summaries of taxonomic membership and population diversity in the MLVS, which broadly agree
150 with previous comparable gut metagenome profiles¹⁶⁻¹⁸. **A)** Inter-individual variation of major
151 phyla. Seven out of fourteen phyla were found present at a relative abundance >0.1% with a
152 prevalence >10%. Panels indicate collection time points, with the number of participants in
153 parenthesis, and samples ordered by decreasing mean abundance of the most abundant phylum
154 (Firmicutes). **B)** Relative abundances of most abundant (when present) species (rows) across
155 913 samples. Columns in the heatmap were ordered based on average linkage clustering on a
156 Euclidean distance matrix of log₁₀ relative abundances. The grey color indicates that the species
157 was not detected. **C)** Principal coordinates ordination of 307 subjects on Bray-Curtis dissimilarity
158 between abundance profiles for 139 species (detected in ≥10% samples at ≥0.01% abundance)
159 averaged over time points and colored by sequencing depth, and each point in the ordination is
160 one participant. The ordination of taxonomic profiles averaged for individual time points shows no
161 bias from variability in input sequencing depth of samples. Labels t1-t4 represent sampling time
162 points; samples were self-collected in two pairs (t1-t2 and t3-t4), six months apart, with each pair
163 spanning 2-3 days.

164 **Supplementary Figure 2. Variation in genus composition among MLVS, HMP1-II, and**
165 **ELDERMET cohorts.** PCoA of Bray-Curtis dissimilarities in genus composition of samples from
166 **(A)** the ELDERMET cohort (participant age >65 yrs) reported by Claesson *et al*¹⁹, **(B)** the
167 ELDERMET, MLVS (age range 65-81 yrs) and HMP1-II (age range 20-40 yrs) cohorts based on
168 the intersection of genera detected in all three cohorts, and **(C)** the ELDERMET, MLVS and
169 HMP1-II cohorts based on the union of *all* genera detected in the three cohorts. Cohort sample
170 numbers are in parenthesis. Only genera with abundances ≥0.01% in ≥10% of samples in
171 respective cohorts were used for ordination analysis; i.e. 42, 57, and 52 genera in ELDERMET,
172 MLVS, and HMP1-II cohorts. **A)** Taxonomic composition of the ELDERMET cohort reported in¹⁹
173 was recapitulated in our analysis indicating that 72% samples were dominated by Bacteroidetes
174 with 56% average abundance across all 192 samples; Firmicutes averaged 39% abundance per
175 sample. In contrast to Claesson *et al*¹⁹, though, we found that the control samples (9 young
176 adults) in the ELDERMET cohort also contained Bacteroidetes at a slightly higher average
177 abundance (45%) than Firmicutes (43%). Such differences in resulting taxonomic profiles are
178 likely a consequence of different OTU calling pipelines. Claesson *et al* used the RDP tool suite
179 version 10.16 whereas we used UPARSE version 9.0.2132 for *de novo* OTU clustering and the
180 RDP classifier version 2.2 for taxonomic classification of OTU centroid sequences against the
181 Greengenes 13_8 database. **B)** Twenty-seven genera represented the taxonomic intersection of
182 all three cohorts. In our MLVS data, taxonomic profiling reinforced Firmicutes and Bacteroidetes
183 as the dominant provenance of bacterial clades, and, unlike in young adults of the HMP1-II
184 cohort^{20,17}, the proportion was tilted in favor of Firmicutes (50.6% ± 14.3%; mean ± s.d.) over
185 Bacteroidetes (40.4% ± 14.3%) in the MLVS. One potential confounder in this comparison is the
186 different DNA extraction protocol between MLVS and HMP, which was more efficient in extracting
187 Bacteroidetes DNA from HMP samples^{21,22}. **C)** Seventy genera represented the taxonomic union
188 of all three datasets and, in addition to sample processing, the comparison is biased by
189 differences in taxonomic assignment strategies for metagenomic and amplicon sequencing reads.
190 Explicit comparison with data from an ELDERMET publication including whole metagenome
191 shotgun (WMS) sequence data¹⁸ was not possible as neither WMS nor 16S sequencing reads
192 from that study are available from the MG-RAST server where the data were deposited

193 (<http://metagenomics.anl.gov/mgmain.html?mgpage=project&project=mgp154>). Taken together,
194 these larger metagenomic and new metatranscriptomic data showed a greater enrichment for
195 Firmicutes, clearer resolution at the species level, and fewer signs of instability or directly age-
196 linked configurations, possibly due to the comparably high level of population health despite
197 participant age. All numbers in parenthesis are stool metagenome sample counts from a total of
198 307 MLVS, 253 HMP (male and female), and 170 ELDERMET (male and female) participants.

199 **Supplementary Figure 3. Feature detection as a function of sequencing depth.** Effect of
200 sample sequencing depth on the ability to detect microbiome functional features in metagenomic
201 and metatranscriptomic sequence data. HUMAnN2 functional profiling of DNA and RNA quality
202 filtered reads was performed on individual samples in species-specific mode, i.e. nucleotide
203 alignment against pangenomes of species identified in the sample with MetaPhlAn2, and in
204 combined species-specific and -agnostic mode, in which reads not matching any pangenome
205 reference sequences were subjected to DIAMOND²³ translated searching against the UniRef90
206 database. Each sample is represented by a green and blue point in each plot. Linear regression
207 trends with 95% confidence intervals are represented by straight lines and grey shading in each
208 plot. Four plots per row from left to right show read alignment rates, and counts of detected
209 UniRef90 gene families, enzymes, and pathways as a function aligned read counts, for 913 DNA
210 samples (A) and 347 RNA samples (B). The number of gene families detected in metagenomic
211 samples increased by less than half a log over a log difference in sequencing depth (A), but well
212 over one log for metatranscriptomes (B) indicating great transcriptional capacity of the gut
213 metagenome. Detection of UniRef90 transcripts Increasing sequencing depth would improve
214 feature detection from RNA samples whereas feature detection was saturated with these input
215 DNA read counts. C) Species rarefaction curve for samples with total counts above the 1st decile
216 (836 samples). The vertical reference line is set at 50,000; at rarified count of 50,000 the median
217 ratio of rarified to observed number of species of samples nears one (boxplot inset). On average,
218 3.5% of input reads per sample were considered by MetaPhlAn2, implying per sample saturation
219 at <1.5M input reads which is roughly seven times less than the average per sample sequencing
220 depth after quality control (9.3M paired-end reads). The ratio of the Chao extrapolated richness
221 from all samples ('specpool' function in the R/vegan package) to the observed number of species
222 in all samples, indicated 89.9% species saturation in the MLVS cohort. D) Rarefaction curves of
223 UniRef90 gene family abundances, using data from samples with total counts above the 9th decile
224 (93 samples), plateaus at a count of 5M. This was also indicated by boxplot summaries of rarefied
225 to observed UniRef90 ratios (inset). The per-sample average read usage rate by HUMAnN2 was
226 60%, implying per sample saturation at 8.3M quality filtered reads. E) For UniRef90 transcript
227 abundances, curves plateau at similar rarefaction levels based on analysis of 184 samples with
228 highest total counts. Average sequencing depth for RNA samples was 6.7M paired-end reads,
229 after quality control. Boxplot whiskers represent 1.5 times the inter-quartile range from the first
230 and third quartiles. RPKs – reads per kilobase.

231 **Supplementary Figure 4. Definition of core metatranscriptome that is robust to sequencing**
232 **depth.** Number of pathways (from a total of 340) with prevalence exceeding the given threshold,
233 calculated from 341 samples with RNA sequencing depth greater than 1M, 2M, 4M, and 8M reads.
234 A change in slope is observed at 81 pathways, which is robust to changes in sequencing depth.
235 These pathways were thus defined as “core”.

236 **Supplementary Figure 5. UniRef90 gene and DNA-normalized transcript abundance is not**
237 **biased by GC content and ORF length.** A) RNA/DNA ratios for gene families (UniRef90s, total
238 n=37,085) do not vary significantly by GC content, plotted as deciles from the lowest to highest
239 %GC in gene families analyzed in 341 metagenome-metatranscriptome paired samples from 96
240 MLVS participants. GC content was calculated as an average across a single representative

241 nucleotide sequence per UniRef90 family. Boxplot whiskers represent 1.5 times the inter-quartile
242 range from the first and third quartiles. Half-open interval labels for x-axis ticks include only the
243 second endpoint, e.g. (0,1] includes values greater than zero and less than or equal to 1. **B)** When
244 analyzed continuously (rather than quantized), neither DNA abundance nor RNA abundance of
245 each UniRef90 gene family (one per point) is strongly correlated with %GC. Additionally, this small
246 degree of %GC bias does not differ between DNA vs. RNA abundances. **C)** As above, RNA/DNA
247 ratios for gene families do not vary significantly by gene length, as deciles from highest to lowest
248 lengths across gene families. Length was again calculated using a single representative
249 sequence per UniRef90. Boxplot whiskers represent 1.5 times the inter-quartile range from the
250 first and third quartiles. **D)** As above, neither DNA nor RNA abundance of UniRef90 gene families
251 were strongly correlated with length. A slightly greater association was detected between greater
252 metagenomic (but not metatranscriptomic) abundance of shorter genes, but even this weak
253 association was of extremely low effect size (maximum absolute Pearson correlation <0.04).

254 **Supplementary Figure 6. Core and variable metatranscriptomes of the stool microbiome,**
255 **with pathway definitions and distribution range of pathway transcript abundances.** DNA-
256 normalized transcript abundances for 239 gut microbiome pathways with detectable RNA in >10
257 of the 341 metatranscriptomes, collected from 96 MLVS participants. **A)** Core metatranscriptome
258 pathways (transcribed in >80% of samples) with RNA:DNA transcription ratio >1. **B)** Low-
259 expression core metatranscriptome pathways with transcript abundance detectable in >80% of
260 samples but an RNA:DNA ratio <1. **C)** Variably metatranscribed pathways detected in DNA but
261 below detection in at least half of RNA samples, and **D)** variably metatranscribed pathways below
262 detection in DNA (and matching RNA) in 30%-80% of the 341 samples. **E)** Thirty-eight pathways
263 that do not fall into any of the categories depicted in **A-D**. **F)** Pathways with the 30 highest and 30
264 lowest mean DNA-normalized transcript abundances among the 341 metatranscriptome samples.
265 Points indicate individual samples with medians overlaid per pathway, with prevalence in
266 parenthesis; see Supplementary Notes for supplementary results text.

267 **Supplementary Figure 7. Per pathway species contributions to metagenomes and**
268 **metatranscriptomes.** Each point in a given pathway plot is a contributing species, and species
269 contributions to DNA and RNA are expressed as relative abundances; i.e. the average
270 abundances from 341 metagenome-metatranscriptome sample pairs from 96 participants. For
271 example, if DNA or RNA for a given pathway is contributed by a single species, based on species-
272 specific HUMAnN2 functional profiling, then the corresponding log₁₀ value along the x or y axis,
273 respectively, is 1. Color scheme: red - species (points) that contributed more RNA than DNA for
274 a given pathway; blue - species (points) that contributed more DNA than RNA for a given pathway;
275 grey - species (points) that contributed equal levels of RNA and DNA for a given pathway. Number
276 of species within each plot (n) and Spearman correlation coefficients (Rho) between species'
277 contributions to DNA and RNA abundances of a pathway are provided in plot titles.

278 **Supplementary Figure 8. Species-stratified distributions of metagenomic potential (DNA)**
279 **and metatranscriptomic activity (RNA) for all pathways with non-zero abundance in at least**
280 **10% of samples.** The 40 most transcriptionally-active species are shown (additional species are
281 grouped as "other"). Abundances were normalized within each pathway for 189 subject-week
282 pairs, from 96 participants. For each pathway, the number of samples with non-zero RNA and
283 DNA is given in the x-axis label. Subjects were ordered to emphasize blocks of subjects with
284 similar metatranscriptomic profiles (see **Methods**). Pathways are sorted in decreasing order of
285 their Weighted Spearman coefficients (see **Fig. 4B**).

286 **Supplementary Figure 9. Ecological interactions in the gut microbiome for individual time**
287 **points.** Significant co-variation and co-exclusion relationships among 104 species in stool
288 metagenomes of MLVS participants. Each node represents a species and edges correspond to

289 significant interactions inferred by BAnOCC (see **Methods**). Stool microbiome taxonomic profiles
 290 were averaged within each participant for the first (215 participants) and second (258 participants)
 291 collection pairs (separated by 6 months). 95% credible interval criteria was used to assess
 292 significance, and only estimated absolute correlations with effect sizes ≥ 0.15 are reported.

293 **Supplementary Figure 10: Strain-level diversity is robust across cohorts.** Principal
 294 coordinate analysis of pairwise nucleotide substitution rates among strains of 21 species identified
 295 in stool metagenomes from MLVS and HMP1-II cohorts. Nucleotide substitution rates were
 296 calculated from multiple sequence alignments using the Kimura Two-Parameter distance²⁴. All
 297 numbers in plot titles are sample counts in which indicated strains were above limit of detection;
 298 from a total of 913 MLVS stool metagenomes and 564 HMP stool metagenomes (from 253 male
 299 and female HMP participants) that were analyzed with StrainPhlAn.

300

301

302

303

304 **Supplementary Tables**

305 **Supplementary Table 1. Functional profiling of MLVS metagenomes and**
 306 **metatranscriptomes.** UniRef90 gene families identified from DNA and RNA, plus those
 307 characterizable to enzymes and pathways per sample and in the entire cohort.

	Metagenome (<i>n</i> = 913)		Metatranscriptome (<i>n</i> = 347)	
	Features per sample (mean \pm s.d.)	Unique in cohort*	Features per sample (mean \pm s.d.)	Unique in cohort
UniRef90	173,609 \pm 36,157	1,569,171	32,279 \pm 21,537	602,896
ECs	1045 \pm 128	1,909	623 \pm 149	1,570
Pathways	253 \pm 40	429	129 \pm 48	340

308 * - Number of unique non-redundant features in entire cohort.

309 **SI References**

310
311 1 Zhang, T. *et al.* RNA viral community in human feces: prevalence of plant pathogenic
312 viruses. *PLoS Biol* **4**, e3, doi:10.1371/journal.pbio.0040003 (2006).
313 2 Virgin, H. W. The virome in mammalian physiology and disease. *Cell* **157**, 142-150,
314 doi:10.1016/j.cell.2014.02.032 (2014).
315 3 Franzosa, E. A. *et al.* Relating the metatranscriptome and metagenome of the human
316 gut. *Proceedings of the National Academy of Sciences of the United States of America*
317 **111**, E2329-2338, doi:10.1073/pnas.1319284111 (2014).
318 4 Zhang, Y. M. & Rock, C. O. Membrane lipid homeostasis in bacteria. *Nature reviews.*
319 *Microbiology* **6**, 222-233, doi:10.1038/nrmicro1839 (2008).
320 5 SHAH, H. N. & COLLINS, M. D. Proposal To Restrict the Genus *Bacteroides* (Castellani
321 and Chalmers) to *Bacteroides fragilis* and Closely Related Species. *International Journal*
322 *of Systematic and Evolutionary Microbiology* **39**, 85-87, doi:doi:10.1099/00207713-39-1-
323 85 (1989).
324 6 Akashi, H. & Gojobori, T. Metabolic efficiency and amino acid composition in the
325 proteomes of *Escherichia coli* and *Bacillus subtilis*. *Proceedings of the National*
326 *Academy of Sciences of the United States of America* **99**, 3695-3700,
327 doi:10.1073/pnas.062526999 (2002).
328 7 Bender, R. A. Regulation of the histidine utilization (hut) system in bacteria. *Microbiology*
329 *and molecular biology reviews : MMBR* **76**, 565-584, doi:10.1128/MMBR.00014-12
330 (2012).
331 8 Reck, M. *et al.* Stool metatranscriptomics: A technical guideline for mRNA stabilisation
332 and isolation. *BMC Genomics* **16**, 494, doi:10.1186/s12864-015-1694-y (2015).
333 9 Song, S. J. *et al.* Preservation Methods Differ in Fecal Microbiome Stability, Affecting
334 Suitability for Field Studies. *mSystems* **1**, doi:10.1128/mSystems.00021-16 (2016).
335 10 Voigt, A. Y. *et al.* Temporal and technical variability of human gut metagenomes.
336 *Genome biology* **16**, 73, doi:10.1186/s13059-015-0639-8 (2015).
337 11 Vlckova, K., Mrazek, J., Kopečný, J. & Petrzalkova, K. J. Evaluation of different storage
338 methods to characterize the fecal bacterial communities of captive western lowland
339 gorillas (*Gorilla gorilla gorilla*). *J Microbiol Methods* **91**, 45-51,
340 doi:10.1016/j.mimet.2012.07.015 (2012).
341 12 Blekhman, R. *et al.* Common methods for fecal sample storage in field studies yield
342 consistent signatures of individual identity in microbiome sequencing data. *Sci Rep* **6**,
343 31519, doi:10.1038/srep31519 (2016).
344 13 Flores, R. *et al.* Collection media and delayed freezing effects on microbial composition
345 of human stool. *Microbiome* **3**, 33, doi:10.1186/s40168-015-0092-7 (2015).
346 14 Vogtmann, E. *et al.* Comparison of Fecal Collection Methods for Microbiota Studies in
347 Bangladesh. *Applied and environmental microbiology* **83**, doi:10.1128/AEM.00361-17
348 (2017).
349 15 Choo, J. M., Leong, L. E. & Rogers, G. B. Sample storage conditions significantly
350 influence faecal microbiome profiles. *Sci Rep* **5**, 16350, doi:10.1038/srep16350 (2015).
351 16 Qin, J. *et al.* A human gut microbial gene catalogue established by metagenomic
352 sequencing. *Nature* **464**, 59-65, doi:10.1038/nature08821 (2010).
353 17 Consortium, H. M. P. Structure, function and diversity of the healthy human microbiome.
354 *Nature* **486**, 207-214, doi:10.1038/nature11234 (2012).
355 18 Claesson, M. J. *et al.* Gut microbiota composition correlates with diet and health in the
356 elderly. *Nature* **488**, 178-184, doi:10.1038/nature11319 (2012).
357 19 Claesson, M. J. *et al.* Composition, variability, and temporal stability of the intestinal
358 microbiota of the elderly. *Proceedings of the National Academy of Sciences of the*

359 *United States of America* **108 Suppl 1**, 4586-4591, doi:10.1073/pnas.1000097107
360 (2011).

361 20 Lloyd-Price, J. *et al.* Strains, functions and dynamics in the expanded Human
362 Microbiome Project. *Nature* **550**, 61-66, doi:10.1038/nature23889 (2017).

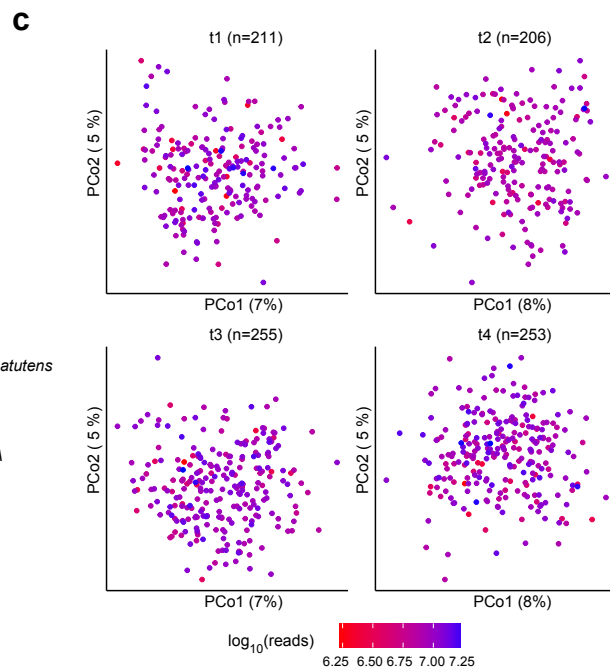
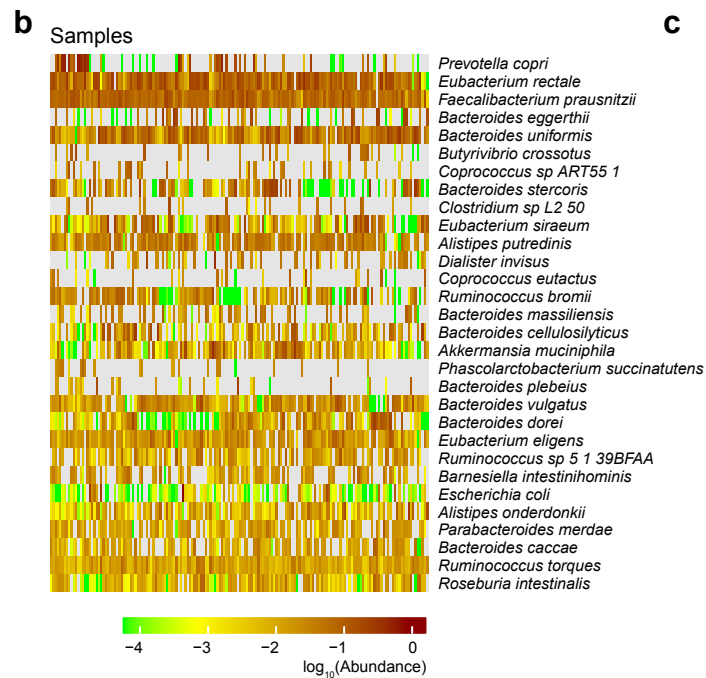
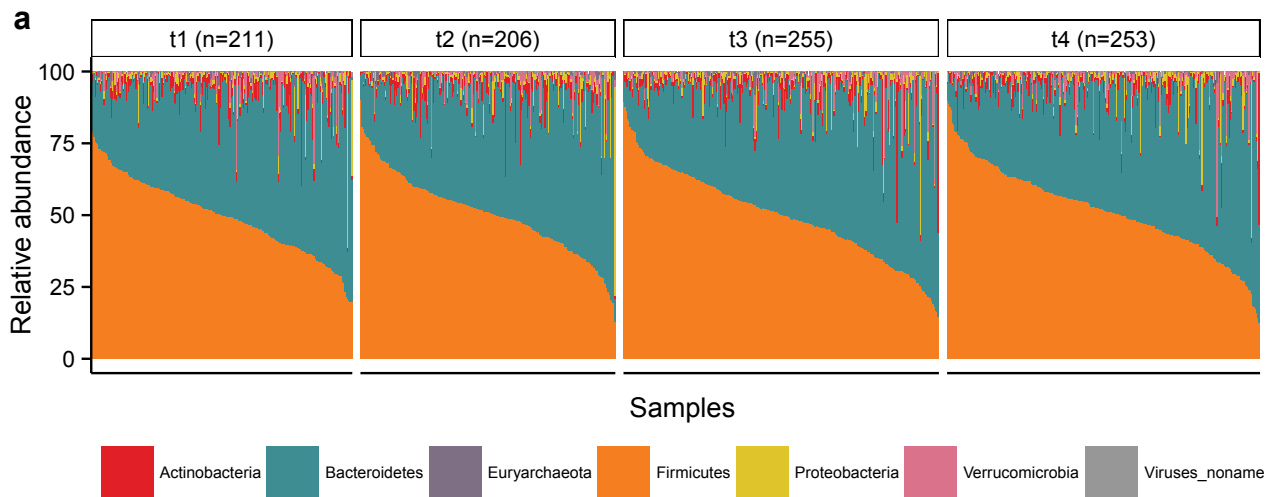
363 21 Wesolowska-Andersen, A. *et al.* Choice of bacterial DNA extraction method from fecal
364 material influences community structure as evaluated by metagenomic analysis.
365 *Microbiome* **2**, 19, doi:10.1186/2049-2618-2-19 (2014).

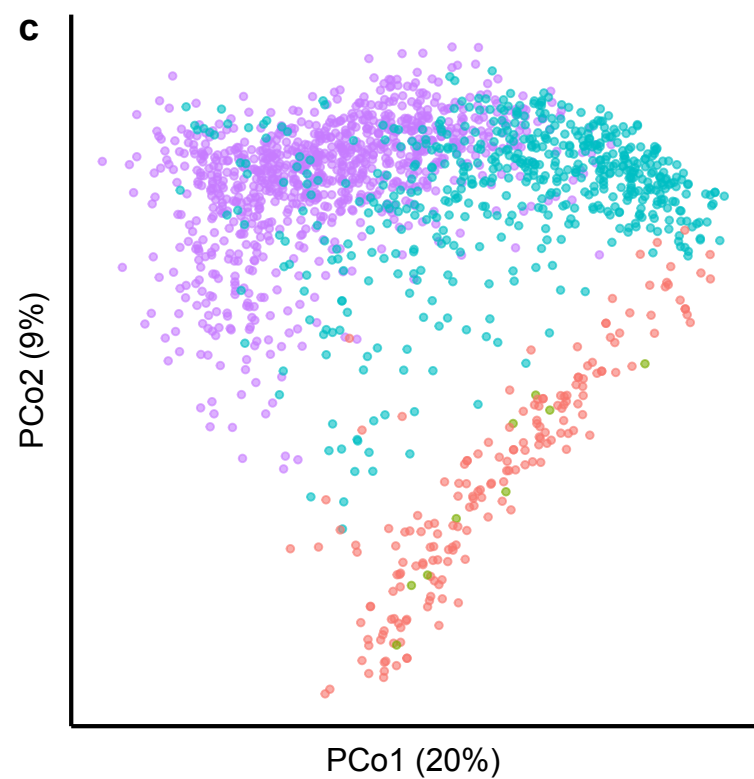
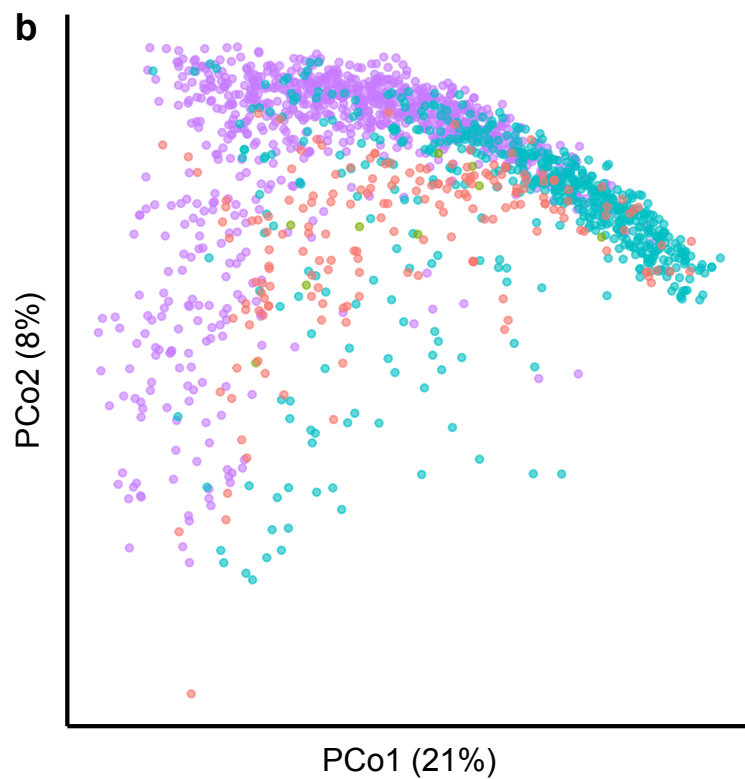
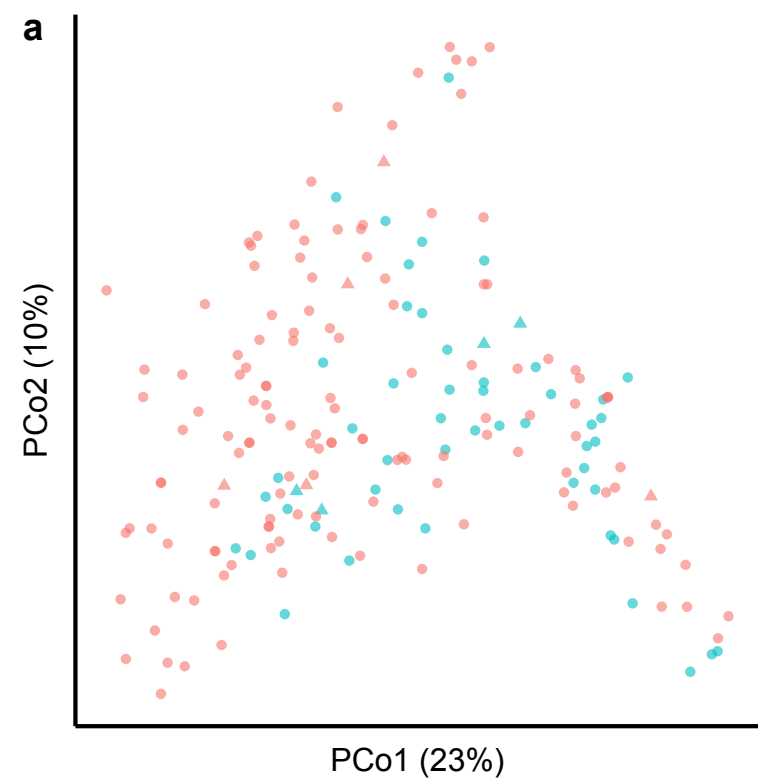
366 22 Shafquat, A., Joice, R., Simmons, S. L. & Huttenhower, C. Functional and phylogenetic
367 assembly of microbial communities in the human microbiome. *Trends in microbiology*
368 **22**, 261-266, doi:10.1016/j.tim.2014.01.011 (2014).

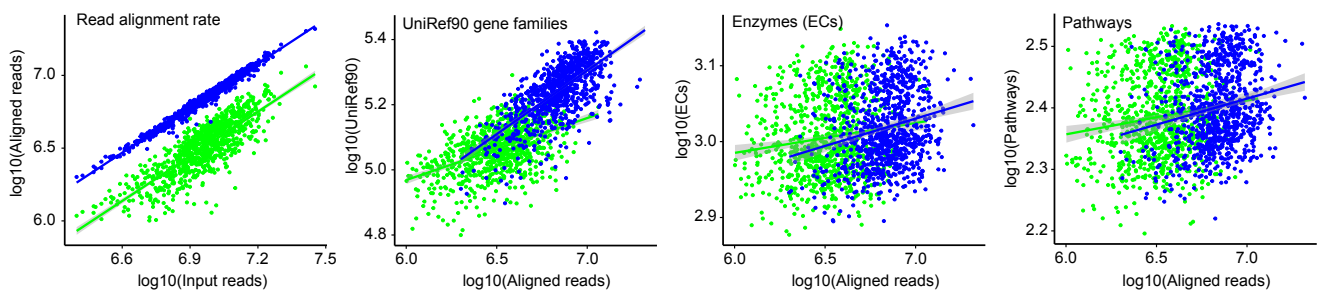
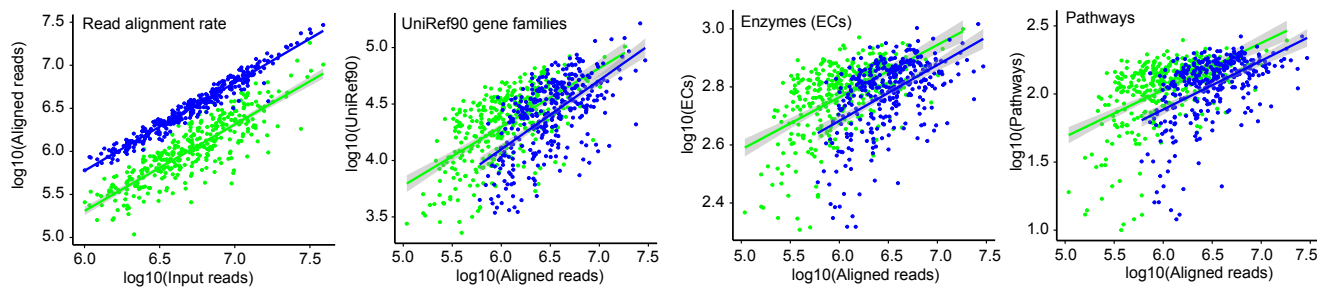
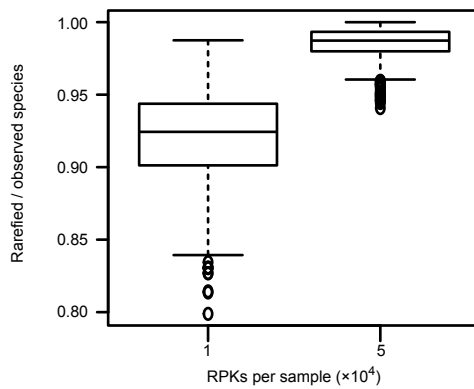
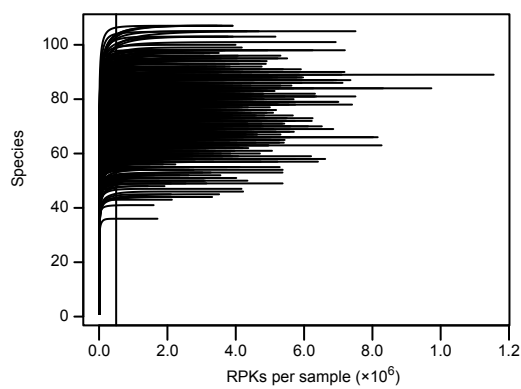
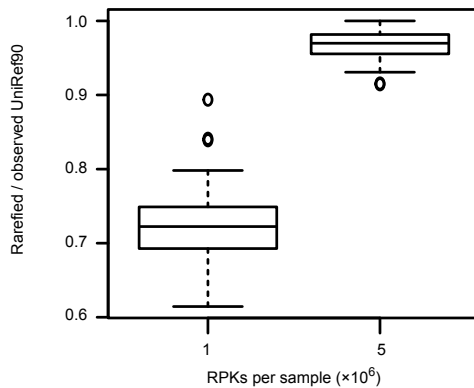
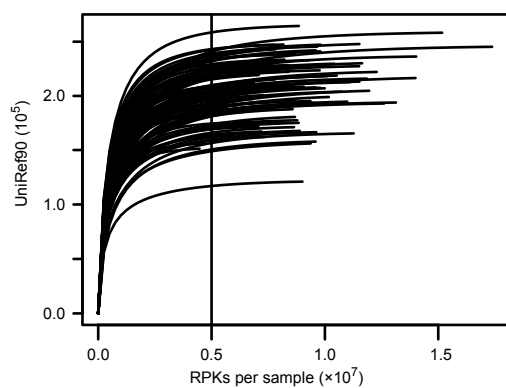
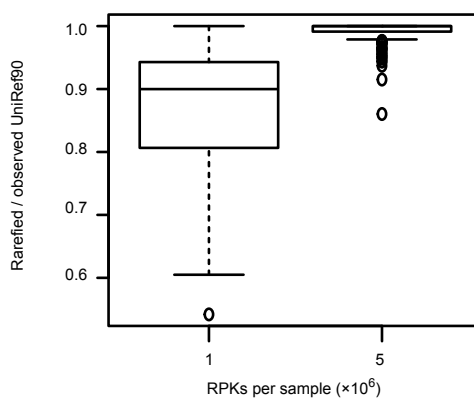
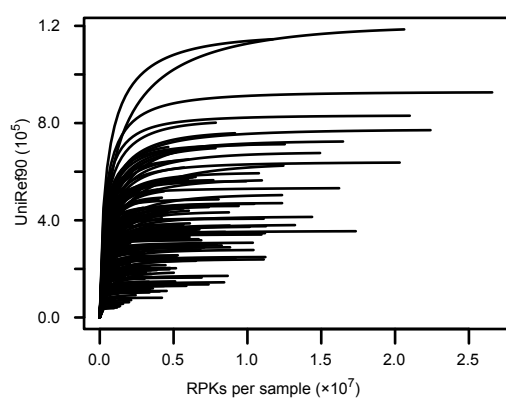
369 23 Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using
370 DIAMOND. *Nature methods* **12**, 59-60, doi:10.1038/nmeth.3176 (2015).

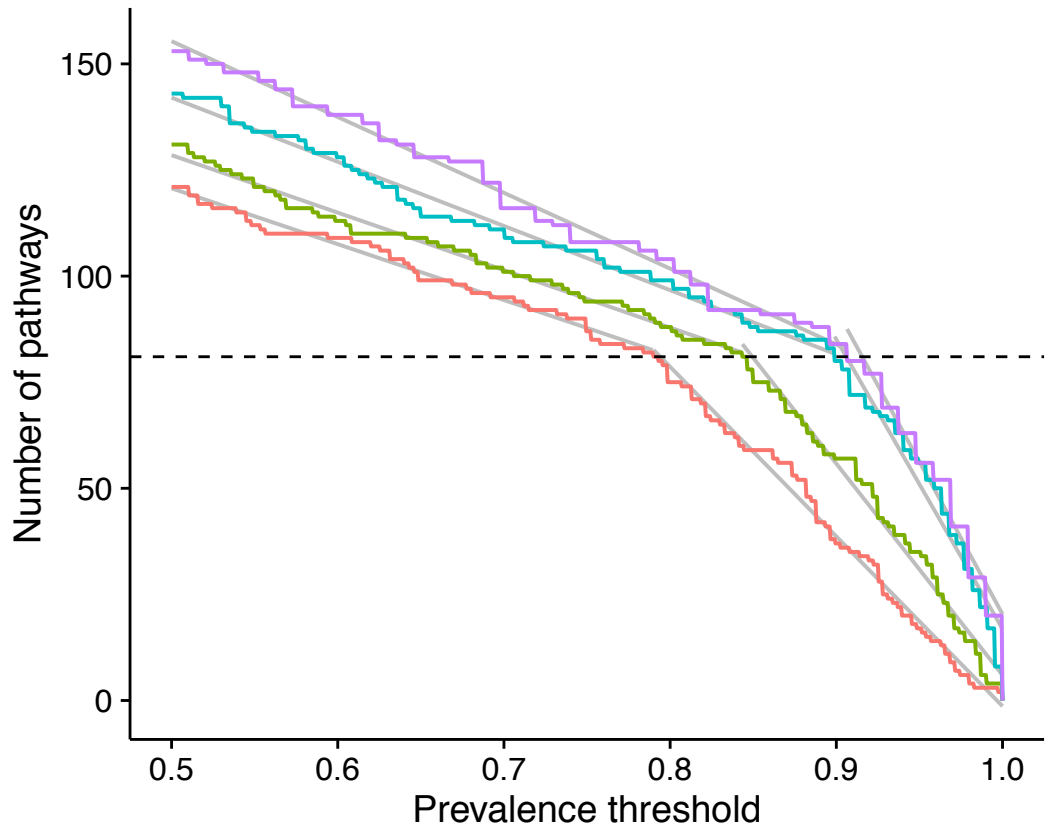
371 24 Kimura, M. A simple method for estimating evolutionary rates of base substitutions
372 through comparative studies of nucleotide sequences. *J Mol Evol* **16**, 111-120 (1980).

373





a Metagenome features**b** Metatranscriptome features**c** Species rarefaction**d** UniRef90 DNA rarefaction**e** UniRef90 RNA rarefaction



Minimum RNA reads — 1M — 2M — 4M — 8M

