# LSHTM Research Online

https://researchonline.lshtm.ac.uk

1 **Correcting index databases improves metagenomic studies**
2
3 Guillaume Méric[1,2,3]*§, Ryan R. Wick[2]§*, Stephen C. Watts[2], Kathryn E. Holt[2,4], Michael Inouye[1,5,6,7]
4
5
6 [1]Cambridge Baker Systems Genomics Initiative, Baker Heart and Diabetes Institute, 75 Commercial Rd,
7 Melbourne 3004, Victoria, Australia
8 [2]Department of Infectious Diseases, Central Clinical School, Monash University, Melbourne, Victoria 3004,
9 Australia
10 [3]The Milner Centre for Evolution, University of Bath, Claverton Down, Bath, BA2 7AY, UK
11 [4]London School of Hygiene & Tropical Medicine, London WC1E 7HT, UK
12 [5]Cambridge Baker Systems Genomics Initiative, Department of Public Health and Primary Care, University of
13 Cambridge, Cambridge CB1 8RN, UK
14 [6]Cambridge Substantive Site, Health Data Research UK, Wellcome Genome Campus, Hinxton, UK
15 [7]The Alan Turing Institute, London, UK
16
17 * Corresponding authors: guillaume.meric@baker.edu.au; rrwick@gmail.com;
18 § These authors have contributed equally.
19
20
21 Keywords: taxonomic classification, species definitions, metagenomics, GTDB
22 Running title: Larger indices improve metagenomic classification
23
24
25 **Abstract**
26
27 Assessing the taxonomic composition of metagenomic samples is an important first step in
28 understanding the biology and ecology of microbial communities in complex environments. Despite a
29 wealth of algorithms and tools for metagenomic classification, relatively little effort has been put into
30 the critical task of improving the quality of reference indices to which metagenomic reads are
31 assigned. Here, we inferred the taxonomic composition of 404 publicly available metagenomes from
32 human, marine and soil environments, using custom index databases modified according to two
33 factors: the number of reference genomes used to build the databases, and the monophyletic strictness
34 of species definitions. Index databases built following the NCBI taxonomic system were also
35 compared to others using Genome Taxonomy Database (GTDB) taxonomic redefinitions. We
36 observed a considerable increase in the rate of read classification using modified reference index
37 databases as compared to a default NCBI RefSeq database, with up to a 4.4-, 6.4- and 2.2-fold
38 increase in classified reads per sample for human, marine and soil metagenomes, respectively.
39 Importantly, targeted correction for 70 common human pathogens and bacterial genera in the index
40 database increased their specific detection levels in human metagenomes. We also show the choice of
41 index database can influence downstream diversity and distance estimates for microbiome data.
42 Overall, the study shows a large amount of accessible information in metagenomes remains
43 unexploited using current methods, and that the same data analysed using different index databases
44 could potentially lead to different conclusions. These results have implications for the power and
45 design of individual microbiome studies, and for comparison and meta-analysis of microbiome
46 datasets.
47

48 **Introduction**
49
50 For more than 3 billion years, microbes have established complex ecological niches in environments
51 and hosts throughout the planet. This makes them ubiquitous components of biogeochemical cycles
52 on land [1], in the sea [2], the atmosphere [3], and on or inside other living organisms [4, 5] including
53 humans, in which they are important for development and health [6, 7]. However, technical
54 constraints limit our ability to study the ecology of microorganisms, in particular the widespread lack
55 of suitable culturing methods [8]. An important advance in the analysis of microbial communities has
56 been the use of sequence-based, culture-independent methods to study the diversity and composition
57 of clinical and environmental samples and their biological functions. The increasing affordability of
58 high-throughput sequencing has led to an increase in metagenomics studies, in which a sample's total
59 extracted DNA can be sequenced as a whole. Accurately determining and quantifying the taxonomic
60 composition of a metagenome is a critical first step in many analyses, such as the association with
61 host phenotype, host genotype, disease status or environmental properties.
62
63 Metagenomic classification begins with the accurate assignment of sequencing reads to a reference
64 database, or "index", comprising reference genomes and their corresponding taxonomic definitions. A
65 wealth of metagenomic classification algorithms have been developed in the last few years [4, 9-16],
66 mainly focusing on improving classification speed and memory usage, including popular methods
67 such as Kraken [17, 18] or Centrifuge [19]. A given read sequence may be shared among closely-
68 related species, particularly when the read length is short, and so classifiers can assign reads to the last
69 common ancestor (LCA) of all taxa sharing their sequence ("LCA-classification"). Despite the
70 development of ever-more efficient classifying algorithms and tools, comparatively little has been
71 done to improve the quality of the reference indices used to define the taxa to which reads are
72 assigned. Recent efforts showed that the addition of new genomes to NCBI RefSeq could influence
73 metagenomic classification performance, with indices built on most recent releases of NCBI RefSeq
74 able to classify more reads overall, but fewer at the species level [20]. Generally, most methods and
75 studies use a selection of representative, often complete genomes from curated repositories to build
76 indices from all described bacteria and archaea using their reported taxonomic definitions [16, 21],
77 typically NCBI RefSeq [22] for whole representative genomes, and SILVA, Greengenes, or RDP [23]
78 for 16S rRNA-based studies.
79
80 Defining accurate monophyletic bacterial species boundaries has always been a challenge. Bacterial
81 taxonomy has historically been defined using imprecise biochemical or ecological phenotypes, with
82 more recent genotyping studies offering numerous examples of clustered "species" previously thought
83 to be distinct, and *vice versa* [24-27]. As a result, microbial taxonomies in reference repositories are
84 riddled with inconsistencies, with described taxa often forming polyphyletic groupings [28, 29],
85 necessitating reconciliation between microbial systematics and genomics [30]. This has recently been
86 addressed by redefining taxonomic definitions using a phylogenetic depth coefficient inferred from a
87 robust prokaryotic phylogeny [28]. This effort, summarised in the Genome Taxonomy Database
88 (GTDB), aims to define strictly monophyletic species groups of equivalent phylogenetic depth. It
89 produced a wealth of novel definitions at various taxonomic levels of the microbial tree of life,
90 redefining approximately 58% of all previous NCBI-based taxonomic definitions [28].
91
92 Most classification tools will recommend the use of default indices, built using a set of complete
93 genomes from NCBI RefSeq. In this study, we assessed the potential for improvement and addressed
94 the following questions: does the choice of reference index affect the performance of metagenomic
95 classification? Does the addition of draft reference genome sequences improve classification? Should
96 we use default NCBI-based indices, custom human microbiome-enhanced indices; or GTDB-based
97 indices? Is the inclusion of metagenome-assembled genomes (MAGs) beneficial? Is the strict
98 monophyly of taxonomic definitions in indices important for classification performance? To answer
99 these questions, we created seven custom indices (**Table 1**) using NCBI-based and GTDB-based
100 taxonomic definitions, and examined their classification performance on samples from three diverse
101 and representative metagenomic datasets: human body sites, marine and soil. Our work addresses the
102 metagenomic classification bias, whereby sequencing reads for particular taxa are present in

103 metagenomics data but remain unclassified using current methods and recommendations. This has
104 important consequences for the classification of metagenomic datasets and downstream applications
105 such as microbiome-wide association studies.
106
107 **Results**
108
109 **Substantial improvements in classification performance can be achieved using larger indices**
110 To examine the impact of custom indices on metagenomic classification performance, we classified
111 404 metagenomic samples from three different datasets using seven custom indices (**Figure S1, Table**
112 **1**) and quantified the proportion of reads per sample that were classified to any taxon and the
113 proportion that remained unclassified (**Figure 1, Figure S2**). Our custom index databases were
114 corrected for two distinct factors: (a) number of reference genomes for each species used to build the
115 index, and (b) strict monophyletic species definition for these reference genomes (**Table 1**). We
116 observed a drastic improvement in classification performance using custom indices, built with more
117 reference genomes, i.e. the greater the number of reference genomes used to build the index, the
118 greater the proportion of reads classified (**Figure 1A-C**). This effect was not associated with
119 sequencing depth (**Figure S3**). For instance, using the NCBI_r88_Human17k index on human
120 metagenomes, which includes only 1.67-fold more genomes than NCBI_r88 (selectively chosen from
121 70 known human microbiome taxa) and monophyly correction (**Table 1**), the median proportion of
122 classified reads per sample increased from 54.7% to 76.5%. The index built with the largest number
123 of genomes, GTDB_r86_46k, consistently classified the most reads in every sample tested. The
124 increase in the median percentage of classified reads per sample from a default NCBI_r86
125 classification for human metagenomes was from 53.6% to 91.3% (median increase of +69.4%; range
126 of +3.9% to +342.8%) (**Figure S2A, Table S2, Table S3**). Similarly, the increase in classified reads
127 per sample for marine metagenomes was from a median of 14.1% to a median of 55.2% (median
128 increase of +276.2%; range of +94.6% to +536.3%); and in soil metagenomes from 33.2% to 66.3%
129 (median increase of +100.7% reads/sample; range of +85.7% to +120.6%) (**Figure S2B-C, Table S2,**
130 **Table S3**).
131
132 We next show that the number of reference genomes rather than strict monophyly of the index
133 database led to increased classification rate. To do so, we compared two indices built using the same
134 reference genomes with (GTDB_r86_8.6k) and without (NCBI_r86) strict monophyletic definitions.
135 When considering all three datasets together, the proportion of unclassified and classified
136 reads/sample were nearly identical using GTDB_r86_8.6k over NCBI_r86 (median increase of +71
137 classified reads per sample, range of difference from 0 to +2,213 reads/sample, equivalent to a median
138 increase of less than +0.0005% of total reads/sample) (**Figure 1, Table S2, Table S3**), indicating that
139 strict monophyly alone does not substantially affect classification rate. On the other hand, the
140 comparison of classification performance using GTDB_r86 vs GTDB_r86_46k captures the effect of
141 adding reference genomes (28,560 vs 46,006 total genomes, respectively) to two similarly
142 monophyletic indices. When compared to GTDB_r86, GTDB_r86_46k produced more classified
143 reads in almost every (402/404) human and environmental sample tested (**Figure 1**, **Figure S2, Table**
144 **S2, Table S3**), with median percentage change in classified reads/sample of +2.7% in human samples
145 (range of –5.6% to +18.7%), +3.2% reads/sample in marine metagenomes (range of +1.8% to +5.0%)
146 and +2.1% (range of +1.9% to +2.3%) in soil metagenomes.
147
148 In human samples, a median of 8.6% of total reads/sample (range of 0.9% to 31.6%) (**Table S2**)
149 remained unclassified even when using our best corrected index GTDB_r86_46k. To investigate what
150 these remaining unclassified reads are, we reclassified them using a pre-computed index based on the
151 nucleotide (nt) database of NCBI, which excludes any whole genome sequence from the WGS or
152 RefSeq databases but includes sequences from all taxonomic domains of life (results in **Figure S4**).
153 The large majority of these reads remained unclassified (~8.5% of total reads/sample); a substantial
154 proportion were attributed to eukaryotic (~0.86% of total reads/sample) and viral (~0.16% of total
155 reads/sample) taxa, which are not included in the custom indices used in this study (**Figure S4**).
156 Notably, ~1.1% of all reads/sample were still attributed to Bacteria and Archaea (**Figure S4**). As the
157 nt database of NCBI excludes reference genomes from WGS and RefSeq, this classification reflects

158  either the presence of genomic fragments in the nt database that are not in WGS of RefSeq, or that
159  these reads mapped to rarer genomic variants that were not included in the 46,006 representative
160  genomes from the GTDB_r86_46k index. As a very small fraction of these unclassified reads were
161  prokaryotic, this result suggests that the GTDB_r86_46k index is much more likely to capture and
162  classify most accessible prokaryotic reads from human metagenomes than default methods.
163
164  **Classification to lower taxonomic ranks is increased and more accurate using larger indices**
165  The interpretation of metagenomics data often focuses on lower taxonomic levels, typically genus-
166  and species-level. We compared the taxonomic levels of lowest-common-ancestor (LCA) read
167  classification between the different indices (**Figure 1D-F, Table S4-S5**). The observed trend in all
168  three datasets was that as indices included more reference genomes, more reads were classified to
169  genus and species level (**Figure 1D-F, Table S4-S5**). In particular, GTDB_r86_46k index showed a
170  greater proportion of reads from human samples classified to genus (median increase of +387.2% in
171  reads/sample; range of –22.4% to +3914.7%) and species level (median increase of +44.4%
172  reads/sample; range of –31.6% to +371.7%), as compared with the default NCBI_r86 index (**Figure
173  1D, Table S4-S5**). For marine samples, corresponding median increases using GTDB_r86_46k were
174  +503.4% reads classified to genus/sample (range of +68.1% to +1124.9%) and +269.2% reads
175  classified to species/sample (range of +64.3% to +567.2%) (**Figure 1E, Table S4-S5**); and +113.4%
176  reads classified to genus/sample (range of +98.8% to +140.0%) and +90.9% reads classified to
177  species/sample (range of +71.4% to +114.8%), for soil samples (**Figure 1F, Table S4-S5**).
178
179  Interestingly, of the two best performing indices, GTDB_r86 (built with almost 18,000 less reference
180  genomes than GTDB_r86_46k) classified a median of –28.1% less reads/sample (range of –70.3% to
181  +70.7%) at the genus level, but a median of +3.7% more reads/sample (range of –17.1% to +28.7%)
182  more reads at the species level than GTDB_r86_46k in human samples (**Figure 1D-F, Table S4-S5**).
183  A similar trend was observed in marine and soil samples (**Table S4-S5**). This is likely because the
184  larger the index, the greater the likelihood it includes genomes from two different species that share
185  genes via recent horizontal transfer, which renders those gene sequences ambiguous at the species
186  level so that they can be attributed to their LCA only. In this way, the largest index GTDB_r86_46k
187  can be considered to offer a more accurate representation of taxonomic classification, with ambiguous
188  reads being accurately attributed to the LCA rather than erroneously to a single species.
189
190  **The specific composition of corrected indices affects classification performance and detection**
191  **levels of specific taxa**
192  Unsurprisingly, the specific composition of reference genomes in custom indices affected
193  classification performance. To demonstrate this, we expanded the default NCBI_r88 index by
194  increasing the coverage of 70 human-associated bacterial genera (including pathogens) by 6,819
195  reference genomes and also correcting monophyly for these genera (to produce the
196  NCBI_r88_Human17k index; **Table 1**, **File S1**). This expansion of the index from NCBI_r88 had a
197  significantly greater impact on the overall read classification rate for the human metagenomes (mean
198  increase of +44.3% reads/sample; mean range of +0.5% to +249.3%) compared to the environmental
199  metagenomes (mean of +21.2% and +7.0% reads/sample for marine and soil samples, respectively;
200  $p<0.0001$, D=0.3355; Kolmogorov-Smirnov test on human vs. environmental per-sample increase
201  percentage distributions) (**Table S3-S4**). The effect was also clear at both genus and species levels,
202  with mean increases of +181.1% (genus) and +36.4% (species) in human samples compared to mean
203  increases of +28.8% and +10.0% (genus), and +19.2% and +6.5% (species) in marine and soil
204  samples, respectively (**Figure 1D-F, Table S3-S4**).
205
206  Specifically, we also observed that 63/70 (90%) of the genera expanded in the NCBI_r88_Human17k
207  index could be classified and detected in human metagenomes at a higher level using this index
208  compared to the default NCBI_r88 index (**Figure S5**). *Pseudomonas*, *Enterobacter*, *Butyrivibrio*,
209  *Lactobacillus*, *Alistipes*, *Moraxella*, *Parabacteroides* and *Faecalibacterium* were amongst the genera
210  with the most significant improvement in detection levels using the expanded index database in HMP
211  metagenomes (**Figure S5**). The detection levels for common human pathogens, including *Yersinia,*
212  *Clostridium, Helicobacter* or *Acinetobacter*, were also improved when using NCBI_r88_Human17k

213  (**Figure S5B**). In human samples, up to 20% of all reads that remained unclassified using NCBI_r88
214  but that could be classified using NCBI_r88_Human17k belonged to *Prevotella* and *Bacteroides*, the
215  rest being attributed to a variety of other genera (**Figure S5C, S5D**). When examining particular
216  species of interest, the detection of *Lactobacillus crispatus* in vaginal samples, *Haemophilus*
217  *parainfluenzae*, *Campylobacter concisus* and *Campylobacter showae* in buccal and throat samples,
218  were particularly improved by the use of the corrected NCBI_r88_Human17k index, along with
219  numerous distinct species of *Prevotella*, *Bacteroides* and *Alistipes* in samples from various body sites
220  (**Figure S5C, S5D**). Our results demonstrate that increasing the number of reference genomes for
221  specific genera of interest can substantially improve their detection levels.
222
223  **Impact of metagenome-assembled genomes on classification performance**
224  The recently published GTDB taxonomic system (release 86.0) includes 3,087 metagenome-
225  assembled genomes (MAGs) in the taxonomic redefinition of the prokaryotic tree of life [28]. We
226  assessed whether the addition of these potentially new taxa to a reference index improved
227  metagenomic classification on the human, marine and soil test datasets. The addition of 3,087 MAGs
228  to GTDB increased the proportion of reads classified by mean +0.72% (human), +0.63% (marine) and
229  +0.51% (soil) (GTDB_r86_noMAGs vs GTDB_r86 index; **Figure 1**, **Figure S2, Table S2-S3**). These
230  results show that adding MAGs to index databases can in principle increase classification
231  performance. However, this increase was limited in our test, likely because the MAGs included in
232  GTDB release 86.0 do not capture many novel sequences from the microbiomes analysed (human,
233  marine, and soil).
234
235  **GTDB-based species definitions affect taxonomic composition, abundance and diversity metrics,**
236  **downstream analyses and interpretation from metagenomic studies**
237  The use of corrected indices had a substantial effect on downstream metagenomic analyses. We
238  compared the 30 most abundant taxa for HMP samples at the family, genus and species levels, from
239  classifications using NCBI_r88, NCBI_r88_Human17k and GTDB_r86_46k (**Figure 2, Figure S6**).
240  A total of 19 (63%) families, 15 (50%) genera and 7 (23%) species appeared in the top 30 taxa using
241  all three indices. Thus, the higher the taxonomic order examined, the more agreement across index
242  databases (**Figure 2A, Figure S6**). Notably, even for taxa in the top 30 using all three indices, the
243  order of abundance varied substantially (**Figure 2B, Figure S6**). Some of this variation was
244  attributable to many taxa having been reclassified and renamed in the larger, monophyly-corrected
245  databases (particularly GTDB_r86_46k, see yellow bars in **Figure 2A**). The increased taxonomic
246  granularity within the GTDB system sometimes led to previously common taxa being divided and
247  redefined as multiple different sub-lineages, each with a distinct taxon name. However, there were
248  also differences in the relative abundances of top 30 taxa that were not explained by this (**Figure S7**).
249  For example, the relative abundance rank of families *Porphyromonadaceae* and *Corynebacteriaceae*
250  were reversed using NCBI_r88 vs. GTDB_r86_46k, as were the genera *Lactobacillus* and
251  *Bifidobacterium*, and the species *Bacteroides fragilis* and *Bacteroides thetaiotaomicron* (**Figure S7**).
252
253  Alpha diversity (within-sample diversity), which has been associated with various phenotypes in
254  different microbiomes [6, 31-33], is estimated directly from taxonomic composition data and
255  therefore showed significant differences between indices. We compared three alpha-diversity metrics
256  at the genus level (observed genus richness, genus evenness and Shannon index at the genus level),
257  calculated from taxonomic composition tables summarised at the genus level based on classifications
258  of the same test data sets but using seven different index databases (**Figure 3**). As expected, the large
259  GTDB-based indices showed a much higher richness, but also had an effect on the evenness of genus
260  distribution, especially in marine metagenomes, which affected Shannon diversity index distribution
261  (**Figure 3**). Notably the effect of index database on alpha diversity values varied between samples,
262  with some increasing in value and others decreasing. In some cases these differences were substantive
263  enough to alter the results of statistical tests for difference in alpha diversity between samples from
264  different body sites (**Figure 3B, Table S6**). For example, in our subset of the HMP dataset, faecal
265  samples were found to have significantly lower Shannon diversity than buccal samples when using
266  the NCBI_r86 index (median of 1.44 [IQR 1.03-2.25] vs median of 2.41 [IQR 2.02-2.70] respectively,
267  p=0.027) (**Table S6**). A similar result was obtained using NCBI_r88 index (**Table S6**). However no

268 such differences were found between faecal and buccal samples when Shannon diversity was
269 calculated using any of the GTDB-based indices (median of 2.09 [IQR 1.62-2.83] vs median of 2.59
270 [IQR 2.16-2.90] respectively, p=0.999 using GTDB_r86_8.6k) (**Table S6**). The situation was
271 reversed when comparing Shannon diversity for faecal and skin samples, with significant differences
272 obtained using GTDB_r86_8.6k (median of 2.09 [IQR 1.62-2.83] vs median of 0.81 [IQR 0.63-1.12],
273 p=0.001) but not using NCBI_r86 (median of 1.44 [IQR 1.03-2.25] vs median of 0.77 [IQR 0.63-
274 1.12], p=0.965) (**Table S6**).

275
276 We also examined the effect of index database choice on beta-diversity, or between-sample diversity
277 assessed by calculating Bray-Curtis dissimilarity between groups of samples from different sources
278 (**Figure S8, S9, S10, Table S7**). The effect on beta-diversity was more subtle than for alpha diversity,
279 with the large GTDB indices yielding greater distance estimates between groups of samples that were
280 already dissimilar using default methods (dissimilarity above 80%; **Figure S8, S9**), but did not
281 significantly alter the overall clustering patterns (**Figure S10**).

282
283
284 **Discussion**
285
286 Considerable efforts have been made to improve methods for detection of taxonomic and functional
287 markers in complex metagenomic samples, including increasing sequencing depth, optimising
288 classification algorithms and developing more accurate *de novo* metagenome assembly tools. In this
289 study, we showed that the index database is a major source of variation in classification performance
290 and has significant ramifications for downstream analyses, which may be substantive enough to
291 change study conclusions (e.g. alpha diversity). Commonly utilised index databases lead to sub-
292 optimal taxonomic classification, with a minority of some read sets being classified. Increasing the
293 number of phylogenetically consistent reference genomes in an index database in either a broad or
294 targeted manner had consistently positive effects on increasing the proportion of reads classified
295 (sometimes several fold higher) and classification to greater taxonomic resolution. To facilitate
296 metagenomic analyses without the need for deeper sequencing or *de novo* assembly, we make freely
297 available these improved index databases (https://github.com/rrwick/Metagenomics-Index-
298 Correction) for two commonly-used classifiers, Centrifuge and Kraken2 and the tools to construct
299 them as NCBI RefSeq and GTDB expand.

300
301 We found that large indices built using recently developed and largely phylogenetically-coherent
302 taxonomic species definitions, such as GTDB [28], greatly increased the number of classified reads.
303 Our results suggest that more coherent taxonomic definitions and accurate taxonomic boundaries,
304 such as those proposed within GTDB, may improve statistical power and biological interpretation of
305 subsequent results, particularly those for compositional and diversity analyses (summarised in **Figure
306 4**). This results in greater taxon granularity, i.e. smaller, more discrete clades of similar phylogenetic
307 depth than commonly known phylogroups, which increases classification accuracy and may improve
308 downstream applications, such as association analysis for particular traits. For example, in
309 microbiome-wide association studies using large cohorts, a weak association with a poorly-defined
310 lineage may be caused by a strong association with a well-defined subset of the poorly-defined
311 lineage (**Figure 4**). Furthermore, at a fixed confidence level, increasing the classification rate of a
312 metagenomic sample offers a more accurate representation of its microbial diversity and may, as we
313 have shown, affect study conclusions. As such, the approach we propose here facilitates improved
314 metagenomic analysis across the full spectrum of sequencing depths. In particular, our results may
315 facilitate "shallow sequencing" metagenomics [34] by maximising the extraction of taxonomic
316 information from samples sequenced at lower depth, thus enabling more cost-effective comparison of
317 thousands or tens of thousands of samples in large-scale metagenomic and multi-omics studies.
318 Lastly, our study shows the importance of consistency in index database when comparing results
319 across studies. Differences in reference genomes and taxonomic coherence may introduce artefacts
320 when integrating metagenomic data across studies, and therefore care should be taken when
321 performing combined or meta-analyses.
322

323
324 **Material and Methods**
325
326 **Description of corrected index databases**
327
328 Seven different indices, ranging in size from 8674 to 46,006 complete genomes, were built to
329 compare the effect of various factors on metagenomic classification performance (**Table 1**). As the
330 focus of this study was not to compare the performance of specific metagenomic classifier tools
331 themselves, but rather to evaluate the impact of custom indices on metagenomic classification, we
332 picked a recently developed classifier, Centrifuge [19], on the basis of an easy-to-use index
333 customisation pipeline, fast metagenomic classification performance and lower RAM usage than
334 alternative tools. Centrifuge allows for the building of custom indices (via the *centrifuge-build*
335 indexer), taking as input a set of sequences with taxonomic labels and a ranked taxonomic tree
336 describing the relationships between those labels.
337
338 The NCBI_r86 and NCBI_r88 indices were built from the default collections of complete bacterial
339 and archaeal genomes from NCBI RefSeq releases 86 (n=8,674 genomes) and 88 (n=10,089
340 genomes), respectively, using the NCBI taxonomy tree. The NCBI_r88_Human17k index was based
341 on the NCBI_r88 index, with the addition of 6819 further reference genomes from NCBI GenBank
342 plus manual curation of the taxonomy for 70 common human commensal and pathogenic bacterial
343 genera (**File S1**). We used the Bacsort pipeline (https://github.com/rrwick/Bacsort) to manually curate
344 the taxonomy within each of these 70 genera to enforce strict monophyly. We also built four indices
345 using the GTDB taxonomic system (**Table 1**). GTDB is based on curation of >125,000 whole genome
346 sequences sourced from NCBI RefSeq and metagenome-assembled genomes (MAGs); but with
347 taxonomic labels and tree re-defined based on phylogenetic relationships inferred from the
348 concatenation of 120 proteins and enforcing strict monophyly [28]. We built the GTDB_r86 index
349 from the default GTDB release 86 set of 28,941 dereplicated bacterial and archaeal genomes
350 representative of the GTDB taxonomy [28], "dereplication" being defined as the selection of
351 reference genomes representative of phylogenetic similarity clusters [28]. In the original GTDB
352 publication and website, two genomes were found to be "replicates" when a set of conditions were
353 met, typically when their Mash distance was ≤0.05 (~ANI of 95%) [4]. The GTDB_r86_8.6k index
354 was built using the exact same 8,674 complete reference genomes as for NCBI_r86, but using the
355 taxonomic labels and trees assigned by GTDB. The GTDB_r86_noMAGs index was built exactly like
356 GTDB_r86, but excluding all 3,087 metagenome-assembled genomes (MAGs) identified in GTDB
357 release 86. Finally, the GTDB_r86_46k index was built using a lower Mash threshold for
358 dereplication than in the default GTDB_r86 set. Specifically, this index included a total of 46,006
359 reference genomes (18,634 more than GTDB_r86), each representative of similarity clusters defined
360 using a Mash [10] distance threshold of ≤0.005 (~ANI of 99.5%). The tax_from_gtdb.py and
361 dereplicate_assembly.py scripts are available in https://github.com/rrwick/Metagenomics-Index-
362 Correction with instructions.
363
364 **Metagenomic datasets**
365
366 We used a total of 404 publicly available metagenomes representing a variety of commonly-studied
367 environments: human body sites, marine and soil environments (**Table S1, Figure S1**). Human
368 samples were from the WGS-PP1 study of the Human Microbiome Project (HMP) and obtained
369 through the HMPDACC.org website [9]. HMP samples were chosen with the following
370 considerations: we kept a representative proportion of each body site represented in the WGS-PP1
371 study, we did not subset a body site source with less than five samples, and we excluded samples with
372 a high (>90%) proportion of low quality reads and samples with low sequencing depth. A total of 98
373 representative samples were selected, corresponding to ~9.2% (n=98/1067) of the HMP WGS-PP1
374 study. A total of 246 marine metagenomic samples were isolated from a range of locations in
375 epipelagic and mesopelagic waters around the world as part of the TARA Oceans survey [35, 36], and
376 were downloaded from the EBI repository (study MGYS00002008; BioProject PRJEB1787). A total
377 of 60 soil metagenomes were sampled in a recent study from meadows ground at various depths [37],

378 and were obtained from NCBI BioProject PRJNA449266. Accessions for all readsets are listed in
379 **Table S1**.
380
381 **Assessment of metagenomic classification performance**
382
383 For all classifications, we ran Centrifuge version 1.0.4 [19] on a Linux x86 cluster with 16 cores and
384 128 GB of RAM allocated for each sample classification. The run time ranged from 11 to 45 minutes
385 per metagenomic sample, depending on the index used for classification and the sequencing depth of
386 the sample. Classification reports were built from the resulting output files using the *centrifuge-*
387 *kreport* tool, and reports were visualised and exported using Pavian version 0.8.1 [38] and custom-
388 made scripts, available at https://github.com/rrwick/Metagenomics-Index-Correction.
389
390 Classification performance was assessed by first comparing the number of unclassified and classified
391 reads per sample for each index database used. This provides an unambiguous way to measure how
392 much of the total microbial information present in each sample can be classified. We also compared
393 the taxonomic ranks to which reads were assigned using each index. It should be noted that the NCBI
394 prokaryotic taxonomic system includes many additional and ambiguous taxonomic ranks that are not
395 present in GTDB, such as "subphylum", "infraclass", "superclass", "subtribe" or "strain". To make
396 results comparable between taxonomic systems, reads were always attributed and reported to the LCA
397 of the standard ranks: phylum, class, order, family, genus and species.
398
399 A pre-compiled index based on the nucleotide (nt) database is available from the Centrifuge website
400 (http://www.ccb.jhu.edu/software/centrifuge/, compiled on 03/03/2018). This database includes all
401 traditional divisions of GenBank, EMBL and DDBJ, and thus includes eukaryotic and viral sequences
402 in addition to prokaryotes. However, the nt database excludes the WGS section of GenBank, which
403 should have a negative impact on the determination of accurate species-specific microbial markers.
404 Accordingly, we observed that the classification of 10 random HMP metagenomes using nt resulted in
405 more unclassified reads than when using GTDB_r86_46k (data not shown). To investigate the origin
406 of the reads which were unclassified by the GTDB_r86_46k index (the best-performing custom index
407 in this study), we reclassified them using the nt database.
408
409 Finally, we assessed the effect of using different indices on commonly-used ecological diversity
410 metrics. The calculation of alpha and beta diversity estimates (observed genus richness, genus
411 evenness, Shannon diversity and Bray-Curtis dissimilarity at the genus level) was performed using the
412 R package *phyloseq* version 1.24.2 [39].
413
414 **Custom scripts and pre-computed index databases availability**
415
416 A collection of scripts used to prepare, compare and analyse Centrifuge classifications using custom
417 index databases, either based on default NCBI or GTDB taxonomic systems, is available at:
418 https://github.com/rrwick/Metagenomics-Index-Correction with instructions. Pre-computed versions
419 of the NCBI_r88_Human17k and GTDB_r86_46k indices suitable for use with Centrifuge [19],
420 Kraken1 [18], Kraken2 (https://ccb.jhu.edu/software/kraken2/) and their variants (KrakenUniq [17],
421 LiveKraken      [40]),      are      freely      available      to      from:
422 https://monash.figshare.com/projects/Metagenomics_Index_Correction/65534.
423

**Competing interests**

The authors declare that they have no competing interests.

**Funding**

**Authors' contributions**

GM designed the study, generated performed analyses, interpreted results and was the major contributor in writing the manuscript. RRW helped generating scripts and databases, interpreted results and participated in the writing of the manuscript. SCW helped to generate code and databases. KEH and MI were key contributors to the study design, interpretation of results and the writing of the manuscript. All authors read and approved the final manuscript.

443     **Figure and table legends**
444
445     **Table 1**. **Description of the seven classification indices used in this study.** The release numbers
446     correspond to NCBI RefSeq releases of genomes from which the reference genomes used to build
447     indices were obtained.
448
449     **Figure 1. Large index databases substantially improve metagenomic classification performance**
450     **and accuracy, including at lower taxonomic levels.** Sequencing reads from three datasets (HMP
451     samples, n=98; TARA Oceans samples, n=246; meadow soil samples, n=60) were classified using the
452     seven index databases presented in Table 1. Boxplots (Tukey) show the distribution of the proportion
453     of unclassified and classified reads/samples for human samples (A), marine samples (B) and soil
454     samples (C) using seven indices (y-axis) is shown for each index size (x-axis), defined as the number
455     of reference genomes used to build the index. Distributions of the breakdown of read classification to
456     the two lowest taxonomic levels (genus, species) for human samples (D), marine samples (E) and soil
457     samples (F) are shown for the NCBI_r86 default index (light blue), two indices based on NCBI_r88
458     (NCBI_r88 in green and NCBI_r88_Human17k in pink) and two indices based on GTDB_r86
459     (GTDB_r86 in orange and GTDB_r86_46k in purple).
460
461     **Figure 2. Effect of index database correction on metagenomic composition.** (A) Number of shared
462     top 30 most abundant families, genera and species after classification of 98 HMP samples using
463     default NCBI_r88 index and corrected NCBI_r88_Human17k and GTDB_r86_46k indices. (B)
464     Comparisons of relative abundances (-$\log_{10}$ scale) between the default NCBI_r88 classification and
465     NCBI_r88_Human17k classifications (left) and GTDB_r86_46k (right) for taxa in the top 30 most
466     abundant of all three classifications (19 families, 15 genera and 7 species). To assess changes in rank
467     order consistency between the classifications, Spearman's rank correlation coefficient, and the
468     associated p-value are shown for both comparisons of NCBI_r88_Human17k and GTDB_r86_46k
469     classifications with NCBI_r88 for all taxa, and at each taxonomic ranks.
470
471     **Figure 3. Using corrected indices to classify metagenomes affects measures of alpha diversity.**
472     (A) The values of three measures of alpha diversity (observed genus richness, genus evenness and
473     Shannon diversity index at the genus level) for each metagenomic sample from three datasets (HMP
474     subset, TARA Oceans, Meadow soil samples) are shown. Three specific comparisons of values are
475     presented, between NCBI_r86 and GTDB_r86_8.6k, between NCBI_r88 and NCBI_r88_Human17k
476     and between GTDB_r86, GTDB_r86_noMAGs and GTDB_r86_46k. Each sample is represented by a
477     line coloured by isolation phenotype. Statistical comparisons of distributions presented in this panel
478     are shown in Table S6. (B) Effect of classification index on alpha diversity metrics comparisons
479     between groups. The scatter plots compare the significance of ANOVA tests on all alpha-diversity
480     measures for each of three comparisons, *P*-values using Dunn's multiple testing (with Holm's
481     correction). The dotted lines represent proportionality for which p-values are identical for
482     classifications using both indices, the red lines denote the p-value threshold of 0.05 (-$\log_{10}$=1.301) for
483     each index.
484
485     **Figure 4**. **Increased taxonomic granularity in classification indices can improve the**
486     **interpretation of microbiome-wide association studies**. (A) Increased taxonomic granularity is
487     defined here by the accurate redefinition, splitting and merging of phylogenetically-coherent strictly-
488     monophyletic lineages, as performed using GTDB. In this example, taxon A is split into taxa A1, A2,
489     A3 and A4, and taxon B is split into B1 and B2. (B) Example classification using two index
490     databases, a smaller number of reference genomes with polyphyletic definitions (left) and a larger
491     number of reference genomes with monophyletic definitions (right). (C) Example effects of index
492     database correction on downstream analysis involving alpha-diversity metrics (left) or for
493     microbiome-wide association studies (right).
494
495     **Figure S1. Description of 404 metagenomic samples used in this study.** The distribution of the
496     number of reads/sample is shown for 98 human (A), 246 marine (B) and 60 soil (C) samples,

497 according to various sampling information (body site for human samples, and sampling depth for
498 marine and soil samples).
499
500 **Figure S2. Per-sample change (% and fold-change) in unclassified and classified reads/sample**
501 **using seven default and corrected NCBI- and GTDB-based indices.** Per-sample percent and fold-
502 changes are shown for the three metagenomic datasets: (A) human samples (n=98), (B) TARA
503 Oceans samples (n=246) and (C) meadow soil samples (n=60). Values are normalised to the number
504 of reads unclassified and classified using the default NCBI_r86 index.
505
506 **Figure S3. Classification improvement using corrected indices is unaffected by variations in**
507 **sequencing depth.** The total number of reads/sample (a proxy for sequencing depth) was plotted
508 against the number of unclassified reads/sample using 6 default and corrected indices, for human (A),
509 marine (B) and soil (C) metagenomes. The regression line was calculated using a linear model fit
510 ("lm" in ggplot2 *geom_smooth* function)  for each index.
511
512 **Figure S4. Reads from human metagenomes that remained unclassified using the**
513 **GTDB_r86_46k index are mostly unknown and eukaryotic.** (A) Proportion of reads/sample from
514 HMP samples (n=98) that are unclassified using GTDB_r86_46k, according to body site of isolation;
515 (B) and (C) outcome of re-classification of these specific reads using an index based on the NCBI
516 nucleotide database (nt; pre-computed on the 03/03/2018 and available on the Centrifuge website
517 [http://www.ccb.jhu.edu/software/centrifuge/]) in number of reads/sample (B) and in proportion (C);
518 (D) per-sample breakdown of domain re-classification, showing the proportion of reads attributed to
519 Eukaryota, Bacteria, Archaea and Viruses or unclassified.
520
521 **Figure S5. Targeted correction for 70 specific bacterial genera increases their detection levels in**
522 **human metagenomes.** (A) The average number of reads classified to all 1001 different bacterial
523 genera to which at least one read was attributed using the NCBI_r88 and the NCBI_r88_Human17k
524 indices were compared, with the 70 specifically corrected genera were highlighted in blue while the
525 non-corrected genera are shown in red. Any point above the line denotes genera to which more reads
526 were classified using the NCBI_r88 index, while any point below the line denotes genera to which
527 more reads were classified using the NCBI_r88_Human17k index. (B) Number of classified
528 reads/sample using a default index and a targeted correction for 70 specific bacterial genera. The 70
529 corrected genera are shown, along with their corresponding distribution of the number of classified
530 reads/sample using NCBI_r88 (red) or NCBI_r88_Human17k (blue). The column on the right
531 indicates the p-value and significance thresholds after Wilcoxon signed-rank tests comparing the two
532 indices. Genera with the highest significance in difference are shown in red, orange and yellow, and
533 non-significant differences are shown in white. (C and D) From all reads that were unclassified using
534 NCBI_r88 but classified using NCBI_r88_Human17k, the top 30 genera (C) and species (D) to which
535 these reads were attributed, in proportion, are shown. Boxplots of different colours denote different
536 isolation sources, showing how body sites are differently affected by the index correction.
537
538 **Figure S6. Effect of index database correction on metagenomic compositional data and most**
539 **abundant taxa.** The 30 most relative abundant families, genera and species to which reads were
540 attributed using three indices (default NCBI_r88 and corrected indices NCBI_r88_Human17k and
541 GTDB_r86_46k) are shown in boxplots. The colour of the taxon name in the y-axes denotes whether
542 the same taxon label was found to be present in the top 30 most abundant taxa after classification by
543 all three indices (blue), by NCBI_r88 and NCBI_r88_Human17k and not GTDB_r86_46k (orange),
544 by GTDB_r86_46k and either NCBI_r88 or NCBI_r88_Human17k (pink) or only in one index
545 (black). In red are indicated the taxon definitions that are existing only in one index. For the
546 comparison at the genus level, green arrows indicate whether the corresponding genus has been
547 specifically corrected in the NCBI_r88_Human17k index.
548
549 **Figure S7**. **Comparison of metagenomic compositional data and most abundant taxa after**
550 **classification with indices built from the same reference genomes taxonomically defined using**
551 **NCBI- or GTDB-based definitions.** The 30 most relative abundant families (A), genera (B) and

552 species (C) to which reads were attributed using the NCBI_r86 and GTDB_r86_8.6k indices, built
553 with the exact same set of complete reference genomes from NCBI RefSeq release 86, are shown in
554 boxplots. The correspondence between the top 30 most abundant taxa from the two classifications is
555 reflected by lines between the two plots. The colouring of the lines denote taxa with an exact
556 correspondence in both indices (plain blue) or whether the GTDB redefinition of taxa affected the
557 correspondence (dotted red). The taxa written in red were created in GTDB.

558

559 **Figure S8. Effect of using corrected indices to classify metagenomes on Bray-Curtis**
560 **dissimilarity between groups of HMP samples (grouped by body site isolation).** (A) Bray-Curtis
561 dissimilarity distributions are shown for pairwise group comparisons between buccal, throat, skin,
562 faecal and vaginal samples of the HMP dataset subset (n=98) used in this study, using seven different
563 classification indices. Coloured panels denote within-group comparisons, white panels denote
564 between-group comparisons. (B) Visualisation of the same data, but ordered to contrast the effect of
565 index database on pairwise group comparisons of Bray-Curtis dissimilarity. Colours represent body
566 sites, similarly to panel A.

567

568 **Figure S9. Effect of using corrected indices to classify metagenomes on Bray-Curtis**
569 **dissimilarity between groups of TARA Oceans and meadow soil samples (grouped by body site**
570 **isolation).** (A-B) Bray-Curtis dissimilarity distributions are shown for pairwise group comparisons
571 between buccal, throat, skin, faecal and vaginal samples of the TARA Oceans dataset (n=246, panel
572 A) and meadow soil dataset (n=60, panel B) used in this study, using seven different classification
573 indices. Coloured panels denote within-group comparisons, white panels denote between-group
574 comparisons. (C-D) Visualisation of the same data, but ordered to contrast the effect of index
575 database on pairwise group comparisons of Bray-Curtis dissimilarity for TARA Oceans samples
576 (panel C) and meadow soil samples (panel D). Statistical comparisons of distributions presented in
577 this figure are shown in Table S7.

578

579 **Figure S10. Effect of using corrected indices to classify metagenomes on Bray-Curtis**
580 **dissimilarity (ordination plots**). Between-sample diversity was compared by calculating and
581 ordinating Bray-Curtis distance measures for samples classified using NCBI_r88 (C, H, M),
582 NCBI_r88_Human17k (D, I, N) and GTDB_r86_46k (E, J, O). To compare the effect of indices on
583 beta-diversity, we performed permutational multivariate analysis of variance (PERMANOVA) on the
584 Bray-Curtis distances to measure the association between sample information (such as isolation
585 source for the HMP samples, or sampling depth for the marine and soil samples) and variance within
586 the dataset, indicated in bold in panels C-E, H-J, M-O.

587

588 **Table S1. Sample description and accession number for 404 public human, marine and soil**
589 **metagenomes used in this study.**

590

591 **Table S2. Summary of classification outcome for three datasets using seven different index**
592 **databases.** The median, average, minimum and maximum values of the number and proportion of
593 classified and unclassified reads/sample is shown for every condition. Classifications were performed
594 using Centrifuge version 1.0.4. Detailed per-sample values are shown in Table S3.

595

596 **Table S3. Detailed per-sample classification outcome using samples from three datasets, after**
597 **classification using seven different index databases.** The number and proportion of classified and
598 unclassified reads/sample is shown for every condition. Classifications were performed using
599 Centrifuge version 1.0.4. A summary of these results is shown in Table S2.

600

601 **Table S4. Summary of classification outcome at the genus and species level for samples from**
602 **three datasets using five different index databases.** The median, average, minimum and maximum
603 values of the number and proportion of classified reads/sample is shown for every condition.
604 Classifications were performed using Centrifuge version 1.0.4. Detailed per-sample values are shown
605 in Table S5.

606

607 **Table S5. Detailed per-sample classification outcome to genus and species levels using samples**
608 **from three datasets, after classification using five different index databases.** The number and
609 proportion of classified reads/sample is shown for every condition. Classifications were performed
610 using Centrifuge version 1.0.4. A summary of these results is shown in Table S4.
611
612 **Table S6. Influence of index correction on the significance of differences between alpha**
613 **diversity of isolation phenotype groups.** Comparison of significance (p-values) from Dunn's
614 multiple testings with Holm correction after ANOVA on three alpha diversity metrics comparisons
615 (observed richness, evenness and Shannon diversity index) between isolation phenotype groups in
616 three metagenomic datasets (body site for HMP samples, depth of sampling for TARA Oceans and
617 meadow soil samples) after classification with six index databases. Specifically, the alpha diversity
618 metric distribution between isolation group pairs were compared after classifications with NCBI_r86
619 and GTDB_r86_8.6k, NCBI_r88 and NCBI_r88_Human17k and GTDB_r86, GTDB_r86_noMAGs
620 and GTDB_r86_46k.
621
622 **Table S7. Influence of index correction on the significance of differences between beta diversity**
623 **(Bray-Curtis dissimilarity) calculated between isolation phenotype groups**. Comparison of
624 significance (p-values) from Dunn's multiple testings with Holm correction after ANOVA on Bray
625 Curtis dissimilarity comparisons between isolation phenotype groups in three metagenomic datasets
626 (body site for HMP samples, depth of sampling for TARA Oceans and meadow soil samples) after
627 classification with four index databases. Specifically, comparisons were between NCBI_r88 and
628 NCBI_r88_Human17k and GTDB_r86 and GTDB_r86_46k.
629
630 **File S1. Description of the NCBI_r88_Human17k index database creation**. Pre- and post-
631 correction Newick and XML phylogenetic trees built on hybrid FastANI/Mash distances for 9928
632 genomes from 70 genera of interest (*Acinetobacter, Alistipes, Anaerostipes, Atlantibacter,*
633 *Bacteroides, Barnesiella, Bifidobacterium, Blautia, Brenneria, Buttiauxella, Butyrivibrio,*
634 *Campylobacter, Cedecea, Citrobacter, Clostridium, Coprococcus, Dickeya, Dorea, Edwardsiella,*
635 *Enterobacter, Erwinia, Escherichia, Eubacterium, Faecalibacterium, Haemophilus, Hafnia,*
636 *Helicobacter, Intestinimonas, Izhakiella, Klebsiella, Kluyvera, Kosakonia, Lachnoclostridium,*
637 *Lactobacillus, Leclercia, Lelliottia, Mangrovibacter, Moraxella, Morganella, Nissabacter,*
638 *Odoribacter, Oscillibacter, Pantoea, Parabacteroides, Paraprevotella, Pectobacterium,*
639 *Phascolarctobacterium, Phytobacter, Plesiomonas, Prevotella, Proteus, Providencia,*
640 *Pseudescherichia, Pseudocitrobacter, Pseudomonas, Psychrobacter, Rahnella, Raoultella, Roseburia,*
641 *Rosenbergiella, Rouxiella, Ruminiclostridium, Ruminococcus, Salmonella, Serratia, Siccibacter,*
642 *Tatumella, Trabulsiella, Xenorhabdus* and *Yersinia*), suitable for visualisation in Archeopteryx [41].
643 The species definitions for these genomes were corrected and made strictly monophyletic to create the
644 "NCBI_r88_Human17k" index.
645
646

**References**

1.  van der Heijden MG, Bardgett RD, van Straalen NM: **The unseen majority: soil microbes as drivers of plant diversity and productivity in terrestrial ecosystems.** *Ecol Lett* 2008, **11:**296-310.
2.  Stocker R: **Marine microbes see a sea of gradients.** *Science* 2012, **338:**628-633.
3.  Vaitilingom M, Amato P, Sancelme M, Laj P, Leriche M, Delort AM: **Contribution of microbial activity to carbon chemistry in clouds.** *Appl Environ Microbiol* 2010, **76:**23-29.
4.  Hacquard S, Garrido-Oter R, Gonzalez A, Spaepen S, Ackermann G, Lebeis S, McHardy AC, Dangl JL, Knight R, Ley R, Schulze-Lefert P: **Microbiota and Host Nutrition across Plant and Animal Kingdoms.** *Cell Host Microbe* 2015, **17:**603-616.
5.  Zilber-Rosenberg I, Rosenberg E: **Role of microorganisms in the evolution of animals and plants: the hologenome theory of evolution.** *FEMS Microbiol Rev* 2008, **32:**723-735.
6.  Human Microbiome Project C: **Structure, function and diversity of the healthy human microbiome.** *Nature* 2012, **486:**207-214.
7.  Cho I, Blaser MJ: **The human microbiome: at the interface of health and disease.** *Nat Rev Genet* 2012, **13:**260-270.
8.  Epstein SS: **The phenomenon of microbial uncultivability.** *Curr Opin Microbiol* 2013, **16:**636-642.
9.  Group NHW, Peterson J, Garges S, Giovanni M, McInnes P, Wang L, Schloss JA, Bonazzi V, McEwen JE, Wetterstrand KA, et al: **The NIH Human Microbiome Project.** *Genome Res* 2009, **19:**2317-2323.
10. Ondov BD, Treangen TJ, Melsted P, Mallonee AB, Bergman NH, Koren S, Phillippy AM: **Mash: Fast genome and metagenome distance estimation using MinHash.** *Genome Biology* 2016, **17:**1-14.
11. Wainwright M, Wickramasinghe NC, Narlikar JV, Rajaratnam P: **Microorganisms cultured from stratospheric air samples obtained at 41 km.** *FEMS Microbiol Lett* 2003, **218:**161-165.
12. Truong DT, Franzosa EA, Tickle TL, Scholz M, Weingart G, Pasolli E, Tett A, Huttenhower C, Segata N: **MetaPhlAn2 for enhanced metagenomic taxonomic profiling.** *Nat Methods* 2015, **12:**902-903.
13. Sankar A, Malone B, Bayliss SC, Pascoe B, Meric G, Hitchings MD, Sheppard SK, Feil EJ, Corander J, Honkela A: **Bayesian identification of bacterial strains from sequencing data.** *Microb Genom* 2016, **2:**e000075.
14. Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman FD, Costello EK, Fierer N, Pena AG, Goodrich JK, Gordon JI, et al: **QIIME allows analysis of high-throughput community sequencing data.** *Nat Methods* 2010, **7:**335-336.
15. Darling AE, Jospin G, Lowe E, Matsen FAt, Bik HM, Eisen JA: **PhyloSift: phylogenetic analysis of genomes and metagenomes.** *PeerJ* 2014, **2:**e243.
16. McIntyre ABR, Ounit R, Afshinnekoo E, Prill RJ, Henaff E, Alexander N, Minot SS, Danko D, Foox J, Ahsanuddin S, et al: **Comprehensive benchmarking and ensemble approaches for metagenomic classifiers.** *Genome Biol* 2017, **18:**182.
17. Breitwieser FP, Baker DN, Salzberg SL: **KrakenUniq: confident and fast metagenomics classification using unique k-mer counts.** *Genome Biol* 2018, **19:**198.
18. Wood DE, Salzberg SL: **Kraken: ultrafast metagenomic sequence classification using exact alignments.** *Genome Biol* 2014, **15:**R46.
19. Kim D, Song L, Breitwieser FP, Salzberg SL: **Centrifuge: rapid and accurate classificaton of metagenomic sequences.** *Genome Research* 2016, **26:**1-9.
20. Nasko DJ, Koren S, Phillippy AM, Treangen TJ: **RefSeq database growth influences the accuracy of k-mer-based lowest common ancestor species identification.** *Genome Biol* 2018, **19:**165.
21. Martin TC, Visconti A, Spector TD, Falchi M: **Conducting metagenomic studies in microbiology and clinical research.** *Appl Microbiol Biotechnol* 2018, **102:**8629-8646.
22. Tatusova T, Ciufo S, Federhen S, Fedorov B, McVeigh R, O'Neill K, Tolstoy I, Zaslavsky L: **Update on RefSeq microbial genomes resources.** *Nucleic Acids Res* 2015, **43:**D599-605.

702   23.   Balvociute M, Huson DH: **SILVA, RDP, Greengenes, NCBI and OTT - how do these**
703         **taxonomies compare?** *BMC Genomics* 2017, **18:**114.
704   24.   Gomila M, Busquets A, Mulet M, Garcia-Valdes E, Lalucat J: **Clarification of Taxonomic**
705         **Status within the Pseudomonas syringae Species Group Based on a Phylogenomic**
706         **Analysis.** *Front Microbiol* 2017, **8:**2422.
707   25.   Meric G, Mageiros L, Pascoe B, Woodcock DJ, Mourkas E, Lamble S, Bowden R, Jolley
708         KA, Raymond B, Sheppard SK: **Lineage-specific plasmid acquisition and the evolution of**
709         **specialized pathogens in Bacillus thuringiensis and the Bacillus cereus group.** *Mol Ecol*
710         2018, **27:**1524-1540.
711   26.   Pettengill EA, Pettengill JB, Binet R: **Phylogenetic Analyses of Shigella and**
712         **Enteroinvasive Escherichia coli for the Identification of Molecular Epidemiological**
713         **Markers: Whole-Genome Comparative Analysis Does Not Support Distinct Genera**
714         **Designation.** *Front Microbiol* 2015, **6:**1573.
715   27.   Hoffmann M, Monday SR, Fischer M, Brown EW: **Genetic and phylogenetic evidence for**
716         **misidentification of Vibrio species within the Harveyi clade.** *Lett Appl Microbiol* 2012,
717         **54:**160-165.
718   28.   Parks DH, Chuvochina M, Waite DW, Rinke C, Skarshewski A, Chaumeil PA, Hugenholtz P:
719         **A standardized bacterial taxonomy based on genome phylogeny substantially revises the**
720         **tree of life.** *Nat Biotechnol* 2018, **36:**996-1004.
721   29.   Koeppel AF, Wu M: **Surprisingly extensive mixed phylogenetic and ecological signals**
722         **among bacterial Operational Taxonomic Units.** *Nucleic Acids Res* 2013, **41:**5175-5188.
723   30.   Klenk HP, Goker M: **En route to a genome-based classification of Archaea and Bacteria?**
724         *Syst Appl Microbiol* 2010, **33:**175-182.
725   31.   Kostic AD, Gevers D, Siljander H, Vatanen T, Hyotylainen T, Hamalainen AM, Peet A,
726         Tillmann V, Poho P, Mattila I, et al: **The dynamics of the human infant gut microbiome in**
727         **development and in progression toward type 1 diabetes.** *Cell Host Microbe* 2015, **17:**260-
728         273.
729   32.   Hahn A, Fanous H, Jensen C, Chaney H, Sami I, Perez GF, Koumbourlis AC, Louie S, Bost
730         JE, van den Anker JN, et al: **Changes in microbiome diversity following beta-lactam**
731         **antibiotic treatment are associated with therapeutic versus subtherapeutic antibiotic**
732         **exposure in cystic fibrosis.** *Sci Rep* 2019, **9:**2534.
733   33.   Jackson MA, Verdi S, Maxan ME, Shin CM, Zierer J, Bowyer RCE, Martin T, Williams
734         FMK, Menni C, Bell JT, et al: **Gut microbiota associations with common diseases and**
735         **prescription medications in a population-based cohort.** *Nat Commun* 2018, **9:**2655.
736   34.   Hillmann B, Al-Ghalith GA, Shields-Cutler RR, Zhu Q, Gohl DM, Beckman KB, Knight R,
737         Knights D: **Evaluating the Information Content of Shallow Shotgun Metagenomics.**
738         *mSystems* 2018, **3**.
739   35.   Pesant S, Not F, Picheral M, Kandels-Lewis S, Le Bescot N, Gorsky G, Iudicone D, Karsenti
740         E, Speich S, Trouble R, et al: **Open science resources for the discovery and analysis of**
741         **Tara Oceans data.** *Sci Data* 2015, **2:**150023.
742   36.   Sunagawa S, Coelho LP, Chaffron S, Kultima JR, Labadie K, Salazar G, Djahanschiri B,
743         Zeller G, Mende DR, Alberti A, et al: **Ocean plankton. Structure and function of the**
744         **global ocean microbiome.** *Science* 2015, **348:**1261359.
745   37.   Crits-Christoph A, Diamond S, Butterfield CN, Thomas BC, Banfield JF: **Novel soil bacteria**
746         **possess diverse genes for secondary metabolite biosynthesis.** *Nature* 2018, **558:**440-444.
747   38.   Breitwieser FP, Salzberg SL: **Pavian: Interactive analysis of metagenomics data for**
748         **microbiomics and pathogen identification.** *bioRxiv* 2016:**2014-2017.
749   39.   McMurdie PJ, Holmes S: **phyloseq: an R package for reproducible interactive analysis**
750         **and graphics of microbiome census data.** *PLoS One* 2013, **8:**e61217.
751   40.   Tausch SH, Strauch B, Andrusch A, Loka TP, Lindner MS, Nitsche A, Renard BY:
752         **LiveKraken--real-time metagenomic classification of illumina data.** *Bioinformatics* 2018,
753         **34:**3750-3752.
754   41.   Han MV, Zmasek CM: **phyloXML: XML for evolutionary biology and comparative**
755         **genomics.** *BMC Bioinformatics* 2009, **10:**356.
756

**Table 1. Description of 7 classification indices used in this study**

| Index name | Reference database and taxonomic definitions used | Description | Total number of reference genomes included | Strictly monophyletic species definitions |
|---|---|---|---|---|
| NCBI_r86 | NCBI RefSeq release 86 | Complete microbial genomes (r86) | 8674 | N |
| GTDB_r86_8.6k | GTDB release 86 | Same genomes as NCBI_8.6k but with GTDB taxonomic definitions, to compare effect of strict monophyletic definitions | 8674 | Y |
| NCBI_r88 | NCBI RefSeq release 88 | Complete microbial genomes (r88) | 10089 | N |
| NCBI_r88_Human17k | NCBI RefSeq release 88 | Same as NCBI_r88 with the addition of all draft genomes from 70 genera of interest, strictly corrected for monophyly | 16908 | Y (for 70 genera only) |
| GTDB_r86_noMAGs | GTDB release 86 | GTDB r86 without metagenome-assembled genomes (MAGs) | 25660 | Y |
| GTDB_r86 | GTDB release 86 | All dereplicated* bacterial and archaeal genomes used to curate the GTDB taxonomy in the GTDB study | 28560 | Y |
| GTDB_r86_46k | GTDB release 86 | Manual dereplication* of GTDB release 86 to get more bacterial and archaeal reference genomes with GTDB taxonomic definitions | 46006 | Y |

* Dereplication is defined as the threshold-based selection of representative reference genomes for phylogenetically-similar clusters.

# Figure 1.



**Legend:**

| Index name | Taxonomic definitions | # genomes | | Index name | Taxonomic definitions | # genomes |
|---|---|---|---|---|---|---|
| NCBI_r86 | Default NCBI RefSeq r86 | 8,674 | | GTDB_r86_noMAGs | GTDB_r86 without MAGs | 25,660 |
| GTDB_r86_8.6k | NCBI_r86 with GTDB definitions | 8,674 | | GTDB_r86 | Representative genomes from GTDB | 28,560 |
| NCBI_r88 | Default NCBI RefSeq r88 | 10,089 | | GTDB_r86_46k | GTDB_r86 + more genomes / taxon | 46,006 |
| NCBI_r88_Human17k | NCBI_r88 + 70 corrected human-associated genera | 16,908 | | | | |

**Figure 2.**

# Figure 3.

# Figure 4.



**A** Default taxonomic definitions | Increased taxonomic granularity (e.g. GTDB)

**B** Metagenomic classification with index database built using:

Default reference genomes and default taxonomic definitions (e.g. default NCBI)

More reference genomes + increased taxonomic granularity (e.g. GTDB)

→ **Impact on analyses based on:**
- composition data (relative abundance)
- diversity metrics (alpha, beta)

**C** Impact of index database correction on downstream analyses:

*Diversity analyses:*

Richness

Evenness

Shannon diversity

Default    Corrected

*Microbiome-wide association studies (MWAS) Q-Q plot:*

Observed (-log₁₀ *p*-value)

Expected (-log₁₀ *p*-value)

**MWAS using:**
Default classification
Corrected classification

# Figure S1.

**Figure S2.**



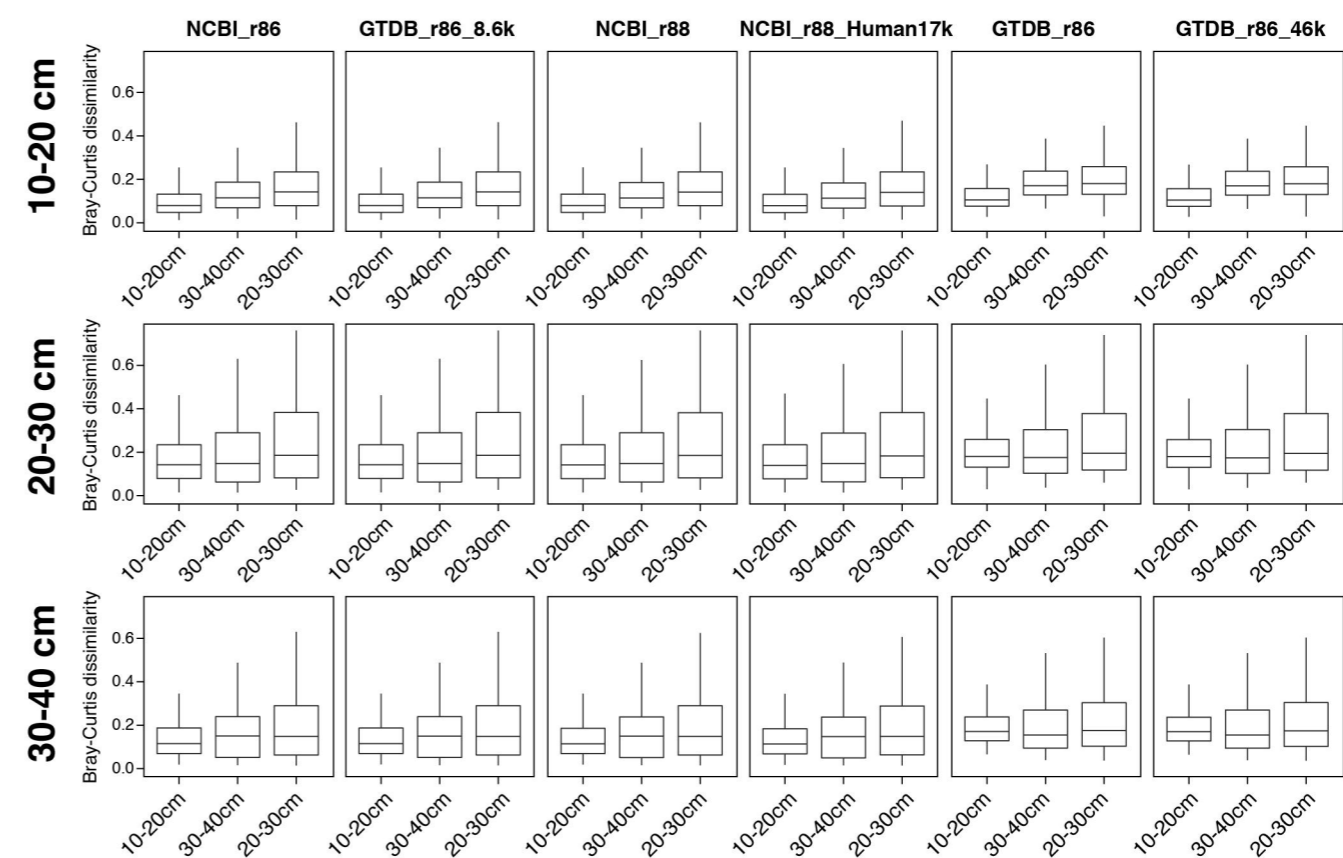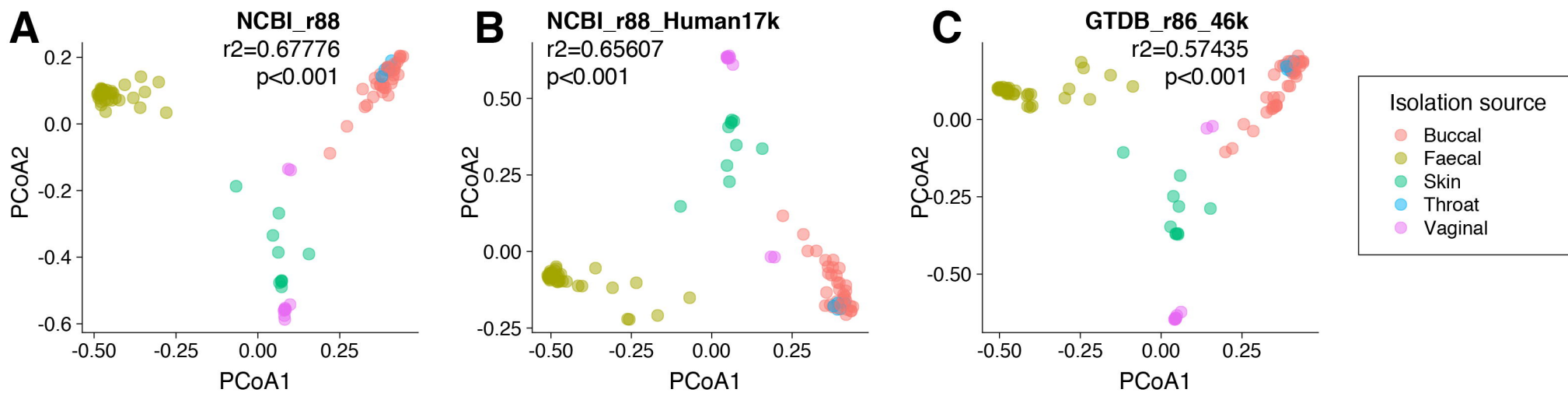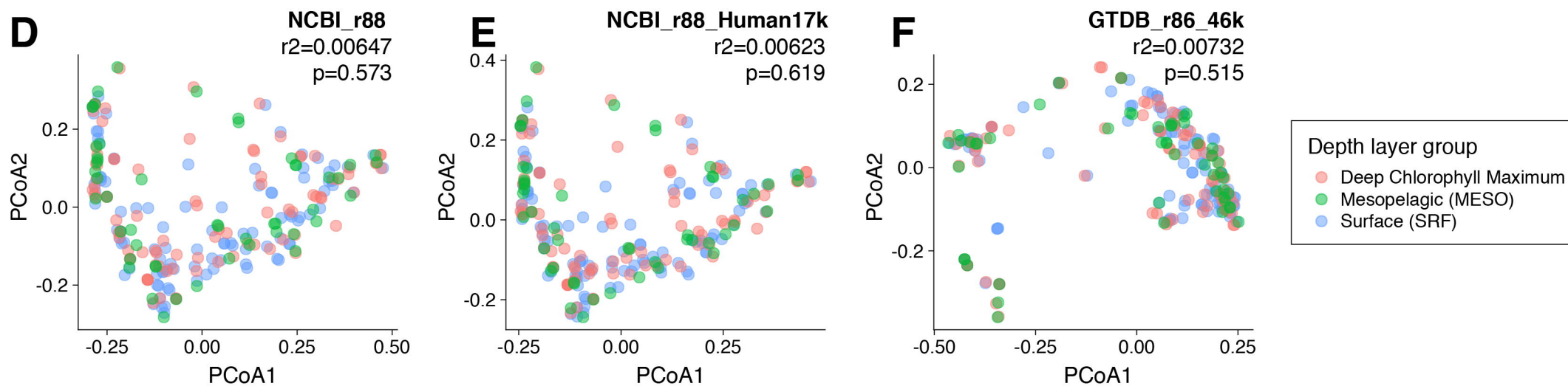**A** HMP samples (n=98)  **B** TARA Oceans samples (n=246)  **C** Meadow soil samples (n=60)

# Figure S3.



**A. HMP samples (n=98)**

**B. TARA Oceans samples (n=246)**

**C. Meadow soil samples (n=60)**

# Figure S4.

# Figure S5.

# Figure S6.

# Figure S7.

# Figure S8.

# Figure S9.

**Figure S10.**