

Evaluating performance of metagenomic characterization algorithms using *in silico* datasets generated with FASTQSim

Anna Shcherbina¹, Darrell O. Ricke*¹, Nelson Chiu¹

*Corresponding author: Darrell Ricke

¹ MIT Lincoln Laboratory
Bioengineering systems and Technologies Group
244 Wood St
Lexington, MA, 02420

Author e-mail addresses:

- Anna Shcherbina - annashch@stanford.edu
- Darrell Ricke – Darrell.Ricke@ll.mit.edu
- Nelson Chiu – Nelson.Chiu@ll.mit.edu

This work is sponsored by the Assistant Secretary of Defense for Research and Engineering under Air Force Contract #FA8721-05-C-0002. Opinions, interpretations, recommendations and conclusions are those of the authors and are not necessarily endorsed by the United States Government.

1 Abstract

2 Background

3 *In silico* bacterial, viral, and human truth datasets were generated to evaluate available metagenomics
4 algorithms. Sequenced datasets include background organisms, creating ambiguity in the true
5 source organism for each read. Bacterial and viral datasets were created with even and staggered
6 coverage to evaluate organism identification, read mapping, and gene identification capabilities of
7 available algorithms. These truth datasets are provided as a resource for the development and
8 refinement of metagenomic algorithms. Algorithm performance on these truth datasets can inform
9 decision makers on strengths and weaknesses of available algorithms and how the results may be best
10 leveraged for bacterial and viral organism identification and characterization.
11

12
13 Source organisms were selected to mirror communities described in the Human Microbiome Project as
14 well as the emerging pathogens listed by the National Institute of Allergy and Infectious Diseases. The
15 six *in silico* datasets were used to evaluate the performance of six leading metagenomics algorithms:
16 MetaScope, Kraken, LMAT, MetaPhlAn, MetaCV, and MetaPhyler.
17

18 Results

19 Algorithms were evaluated on runtime, true positive organisms identified to the genus and species
20 levels, false positive organisms identified to genus and species level, read mapping, relative abundance
21 estimation, and gene calling. No algorithm out performed the others in all categories, and the
22 algorithm or algorithms of choice strongly depends on analysis goals. MetaPhlAn excels for bacteria
23 and LMAT for viruses. The algorithms were ranked by overall performance using a normalized
24 weighted sum of the above metrics, and MetaScope emerged as the overall winner, followed by Kraken
25 and LMAT.
26

27 Conclusions

28 Simulated FASTQ datasets with well-characterized truth data about microbial community composition
29 reveal numerous insights about the relative strengths and weaknesses of the metagenomics algorithms
30 evaluated. The simulated datasets are available to download from the Sequence Read Archive
31 (SRP062063).
32

33 Keywords

34 FASTQsim, Metagenomics, *in silico*, evaluation, Kraken, LMAT, MetaPhlAn, MetaPhyler, MetaScope,
35 MetaCV

36 Background

37 Continuing advances in sequencing technologies are increasing the feasibility of sequencing entire
38 microbial communities rather than individual organisms. This has led to rapid developments in the field
39 of metagenomics aimed at studying genomic material recovered directly from environmental and
40 medical samples. Sequencing the metagenome enables the capture of greater genetic diversity than
41 can be sampled with highly targeted approaches such as microarrays. Metagenomic sequencing has a
42 number of applications for medical diagnostics (i.e. human gut microbiome analysis), environmental
43 profiling (i.e. soil samples), and homeland defense[1-3]. Metagenomic techniques also enable the study
44 of communities of organisms simulated *in vitro*[4].

45 Simultaneously, a number of bioinformatics tools have been developed to analyze metagenomic

46 sample data. They employ a variety of techniques to achieve the opposing goals of high accuracy and
47 low runtime. In this study, the performance of these varied approaches to metagenomic sequence
48 classification was evaluated on a suite of *in silico* datasets with perfectly characterized composition.
49 MetaScope, winner of the Defense Threat Reduction Agency's Grand Challenge[5], relies on sequence
50 analysis using spaced seeds followed by an augmented least common ancestor algorithm to map reads
51 and assign genes for input FASTQ samples[6, 7]. Kraken[8] uses exact alignment of k-mers in
52 combination with an optimized database and another version of the least common ancestors algorithm.
53 MetaPhlAn[9] relies on unique clade-specific marker genes identified from 3000 reference genomes.
54 The Livermore Metagenomic Analysis Toolkit (LMAT) exploits genetic relationships between different
55 organisms by pre-computing the occurrence of each short sequence across the entire reference
56 database and storing the evolutionarily conserved sequence patterns[10-12]. MetaCV translates
57 nucleotide sequences into six frame peptides, which are then decomposed into k-mers. The k-mer
58 frequency is computed in a protein-reference database and used to assign k-mer weights[13]. Finally,
59 MetaPhyler uses a precomputed database of reference phylogenetic marker genes to build a sequence
60 classifier. The classifier, based on BLAST, uses trained thresholds for various combinations of
61 taxonomic ranks, sequence length, and reference genomes[14].

62
63 Simulated *in silico* datasets are a valuable tool for metagenomic research and provide capabilities to
64 evaluate algorithm performance as well as to test hypotheses that cannot be examined through
65 empirical observation. For example, simulated data has revealed biases and heterogeneity in the
66 estimation of diversity metrics from metagenomics samples[15]. Additionally, multiple studies have
67 demonstrated the usefulness of simulated metagenomics datasets for benchmarking sequence
68 assembly and gene prediction pipelines[16-18]. Simulated datasets are also an effective means of
69 parameter optimization for improved algorithm performance and can be used to optimize study design.
70 Sequence simulation can aid with answering questions about coverage requirements, necessary
71 sequence length, and whether paired-end or single-end sequencing should be used. For example, the
72 ART simulator was successfully used by the 1000 Genomes Project Consortium to examine the effects
73 of read length and PE insert size on a read's ability to map to the human genome[19].

74 In this study, six *in silico* datasets were simulated by the FASTQsim tool. **Figure S1** illustrates the
75 composition of each dataset. These datasets contained sequences from reference bacterial and viral
76 genomes, as most human pathogens are members of these taxa. The HMP Even and HMP Staggered
77 datasets were generated to include sequences from the 20 organisms from the Human Microbiome
78 Project[20] (**Supplementary Table 1**). The HMP organisms were selected for inclusion after an attempt
79 to benchmark the performance of MetaScope with the HMP dataset revealed potential contamination
80 in the dataset. As the HMP benchmark dataset was generated by sequencing organisms cultured *in*
81 *vitro*, there was no absolute truth for any background contaminant organisms in the dataset and it was
82 not possible to determine whether the contamination was real or whether MetaScope was calling false
83 positive organisms.

84
85 The bacterial dataset (**Supplementary Table 2**) was designed to test algorithm specificity. Four genera
86 of pathogens were selected from the National Institute of Allergy and Infectious Diseases list of
87 biodefense and emerging infectious disease agents[21] due to their relevance to disease diagnostics
88 from metagenomics samples. These included *Yersinia*, *Coxiella*, *Brucella*, and *Salmonella*. Additionally,
89 the *Escherichia* genus was added to the list due to the high abundance of representative sequences in
90 GenBank[22].

91
92 Two virus datasets were generated with 21 species across 11 representative genera (**Supplementary**

93 **Table 3).** As with the bacterial dataset, candidates were selected due to their inclusion on the NIAID list
94 of emerging pathogens (Marburg virus, Machupo virus, Sudan ebolavirus, Junin virus, Guanarito virus,
95 Chapare virus, Omsk hemorrhagic fever virus) as well as abundance of representative organisms in
96 GenBank (HIV 1, HIV 2, Influenza A virus).

97 Finally, a dataset of human reads from build GRCh38 at 10x (22 million reads) coverage was generated
98 to test host-filtering capabilities of each algorithm. This dataset was generated to measure how well
99 algorithms can overcome the challenges posed by human sequence contamination in public reference
100 databases[23]. For example, endogenous retroviral remnants may be incorrectly classified as belonging
101 to viral genomes in a sample[24-26].

102 **Methods**

103 **Improvements to FASTQsim**

104 The FASTQsim toolkit was augmented to annotate gene information for simulated reads[27]. The
105 “FASTQmapGenes” functionality was added, allowing users to specify NCBI accession ids to use for
106 annotating gene information in simulated reads. The FASTQsim toolkit uses the Entrez and SeqIO
107 libraries from BioPython[28] to download the specified files from GenBank in .gb format. The
108 GenbankParser[29] java application is then used to parse the .gb files in order to extract all information
109 encoded in the CDS and Gene tags. These gene and CDS annotations are appended to the headers
110 within the simulated FASTQ files generated by FASTQsim, such that all reads that fall within a CDS or
111 gene region are annotated with the corresponding CDS and gene information.
112

113 **In silico data generation**

114 The FASTQsim toolkit was used to generate six *in silico* datasets. All were generated with the Illumina
115 error and read length profile included with FASTQsim version 2.0, with no host background added.
116 Specifically, read length of 150 bases was used, with single base mutation, insertion, and deletion rates
117 as specified in the FASTQsim v. 2.0 documentation
118 (<http://sourceforge.net/p/fastqsim/code/ci/master/tree/params/illumina/>). NCBI identifiers for all
119 input data are listed in **Supplementary Tables 1-3**. The Krona toolkit[30] was used to visualize
120 evaluation dataset composition.
121

122 Two *in silico* datasets were generated – “HMP Even” and “HMP Staggered” (**Supplementary Table 1**).
123 For the HMP even dataset, FASTQsim was executed to provide equal number of reads for each species
124 of organism (approximately 60,000 reads per species), with one exception -- 559 reads for *Streptococcus*
125 *agalactiae* were added to simulate a low-level contaminant organism. Version 2.0 of the FASTQsim
126 algorithm probabilistically simulated read counts and error distributions based on a provided model.
127 Due to the probabilistic nature of the algorithm, coverage levels deviated slightly from the specified
128 60,000 reads, with the largest deviation observed for the *E. faecalis* organism (52,290 reads). For the
129 HMP Staggered dataset, coverage levels varied from 11.3x (217,512 reads) for *Actinomyces*
130 *odontolyticus* to 0.001x (2 reads) for *Neisseria meningitidis*. The goal of the staggered dataset was to
131 evaluate the ability of metagenomic algorithms to detect organisms present at very low concentrations,
132 i.e. less than 5 reads.
133

134 The bacterial dataset included reads from the genera *Yersinia*, *Coxiella*, *Brucella*, *Salmonella*, and
135 *Escherichia*. For each of the five genera, several representative species were selected (i.e., *Brucella*
136 *abortus*, *Brucella melitensis*, *Brucella suis*). Next, several representative strains were selected for each
137 species (i.e. *Brucella melitensis* ATCC 23457, *Brucella melitensis* biovar abortus 2308, *Brucella melitensis*

138 biovar 1 strain 16M, and *Brucella melitensis* M28). Organisms were spiked into a FASTQ dataset with
139 coverage levels ranging from 10x to 0.00002x (1 read).

140
141 For the Virus Even dataset, 10x coverage of each organism was simulated. For the Virus Staggered
142 dataset, coverage varied from 100x for Sudan ebolavirus to 0.5x for the Human coronavirus HKU1.

144 Metagenomic algorithm execution

145 Six metagenomic algorithms were selected for execution on the evaluation datasets. These included:

- 146 • MetaScope – winner of the Defense Threat Reduction Agency’s Grand Challenge[7] (version
147 2.0)
- 148 • MetaPhlAn[9] (version 1.7.8, <https://bitbucket.org/nsegata/MetaPhlAn/src/>),
- 149 • MetaCV[13] (version 2.3.0, <http://sourceforge.net/projects/metacv/files/>),
- 150 • MetaPhyler[14] (version 1.13, <http://MetaPhyler.cbcb.umd.edu/#download>),
- 151 • Kraken[8] (v0.10.5, <https://ccb.jhu.edu/software/kraken/>),
- 152 • LMAT[10-12] (v1.2.5, <http://sourceforge.net/projects/lmat/>).

153
154 All algorithms were executed on each of the evaluation datasets using a machine with 512 GB of
155 RAM, 64 cores, 1 TB hard drive, running the Fedora 17 operating system. All algorithms were
156 executed with the default set of databases described in their respective documentation, downloaded
157 on March 1, 2015. Algorithms were evaluated using 60 of the 64 available cores.

158 Attempts were also made to install and run the SURPI (v1.0, <https://github.com/chiulab/surpi>)[31] and
159 compressed BLAST (v0.9, <http://cast.csail.mit.edu/>)[32] algorithms, but these were unsuccessful.

161 Algorithm performance evaluation

162 Runtime in seconds, true positive genus and species calls, false positive genus and species calls, read
163 mapping, and relative abundance results at the species level were computed for all algorithm results.
164 Additionally, correct gene calls were calculated for the set of algorithms that provided gene calling
165 results (MetaScope, MetaCV, LMAT). The Gene ID Conversion function in the DAVID Bioinformatics
166 Database[33] was used to convert across gene representation formats utilized by the three algorithms.
167 Genes were marked as true positives if they matched the gene id, official gene symbol, locus tag,
168 protein id, or specific product name of the truth data.

170 Availability of supporting data

171 The FASTQsim toolkit can be downloaded from SourceForge: <http://sourceforge.net/projects/fastqsim/>
172

173 *In silico* evaluation datasets can be downloaded from the Sequence Read Archive: SRP062063

174 SRR2146185 – Virus Staggered dataset

175 SRR2146184-- Virus Even dataset

176 SRR2146183—Bacterial dataset

177 SRR2146181—HMP Staggered dataset

178 SRR2146182 – HMP Even dataset

179 Results and Discussion

180 Runtime in seconds, true positive genus and species identification, false positive genus and species
181 identification, and false negative species calls were determined for each of the metagenomic algorithms
182 (**Figure 1**). Among the algorithms evaluated, only MetaScope mapped a small number of reads in our
183 datasets to a taxon rank below species. Consequently, although the initial focus of the Bacterial dataset

184 was to assess the ability of the algorithms to distinguish between different strains of the same species,
185 it was decided to evaluate both true and false positives at species and genus level. To determine an
186 overall rank of the algorithms across the datasets, the area occupied by each in the radar plot was
187 computed (**Table 1**). When the polygon area was calculated using the MATLAB polyarea function and
188 summed across all datasets, MetaScope emerges as the winner, with the largest overall area. Kraken
189 and LMAT are the runner-ups, and MetaPhyler performed the worst. In addition to the algorithms'
190 rank overall, several trends can be noted in the individual performance categories.

191 The algorithms diverged in runtime by several orders of magnitude (**Table 2**). Overall, MetaPhlAn had
192 the shortest runtime. The algorithm had the fastest time on the three bacterial datasets – 22.64 s for
193 HMP Even, 53.3 s on HMP staggered, and 220 s. on Bacteria. The second fastest times for these three
194 datasets were 5 to 10 times slower: 233 s (MetaPhlAn), 261 s (MetaScope), and 2,700 s (LMAT),
195 respectively. MetaPhlAn is able to execute quickly partly because it does not perform a host-filtering
196 step. MetaPhlAn came in second for the virus datasets, with a runtime of 11 seconds on both,
197 compared to 9 and 7 seconds for Kraken. MetaPhlAn failed to run on the human dataset. Kraken,
198 MetaScope, and LMAT exhibited similar runtimes on all datasets, averaging 353 s on HMP Even, 354 s
199 on HMP staggered, and 3,595 s on Bacteria. On the other end of the spectrum, MetaPhyler was an
200 outlier for high runtime, requiring 15,480 s on HMP Even, 19,231 s on HMP staggered, and 129,600 s on
201 Bacteria.

202 In addition to its high speed, MetaPhlAn also achieved the highest accuracy, defined as ratio of true
203 positives to false positives, on the bacterial datasets. It identified all 20 species in the HMP even
204 dataset with only a single false positive organism. On HMP staggered, it missed 4 species out of 20 but
205 reported only 2 false positive species. MetaScope, the runner up, reported a single false negative
206 species but 414 false positives. However, the MetaPhlAn reference database is customized for
207 bacteria, and no support exists at the time of this writing for profiling viruses or eukaryotes.
208 MetaScope achieved the second- highest ratio of true positives to false positives, reporting slightly
209 more true positives and approximately half as many false positives as Kraken. LMAT was the least
210 conservative and reported the highest number of false positive organisms. MetaPhyler made highly
211 conservative calls—false positives were low, but so were true positives. Additionally, MetaPhyler, and
212 MetaCV, as well as MetaPhlAn, did not report results for the viral datasets.

213 Algorithm performance on the Human dataset (**Figure 2k**) illustrates the efficacy of the host-filtering
214 step for each algorithm. The human reference genome is incomplete[34, 35] and misses regions
215 specific to individual host subjects. These missed regions show up as false positives on the Human
216 evaluation dataset – algorithms assign them to organisms other than the human host because these
217 reads are not removed during the host filtering step. For example, MetaScope reports 152 organisms,
218 with fewer than 100 reads assigned to each. Kraken has a similar false positive profile; it reports 1,266
219 species that account for <1% of the reads in the dataset. MetaCV reports 2,998 false organisms with
220 low read count, and LMAT reports 1,118 species that account for less than 0.01% of the reads.
221 MetaPhyler does not report results more specific than the Class taxonomy level for the Human
222 dataset, in line with the conservative approach of this algorithm. MetaPhlAn crashes with a
223 segmentation fault on the Human dataset, which most likely is an artifact of the non-host-filtering
224 approach used by this algorithm.

225
226 The algorithms were evaluated based on their ability to correctly map reads and predict relative
227 abundance of the organisms in the data (Figures 2,3). For the bacterial datasets, Kraken and MetaScope

228 classified the highest number of reads correctly for both the genus and species level, and cluster closest
229 to the truth in the dendrogram. However, for the viral datasets, LMAT performed best, classifying the
230 most reads correctly.

231 Although the *Actinomyces odontolyticus* (NZ_DS264586.1) organism had the highest coverage (11.3x,
232 217512 reads) in the HMP staggered dataset, the algorithms on the whole did not perform well on this
233 organism. It was not identified by the Kraken, MetaCV, and MetaPhyler algorithms, and called at a low
234 level by MetaScope (153 reads) (**Figure 2g**) MetaCV mapped the most reads correctly –108,211 (49.7%)
235 and MetaPhlAn was second best, identifying 22,647 (10.4%) of the reads. None of the algorithms
236 identified any of the 2,045 *A. odontolyticus* genes (**Figure 5b**). This poor performance likely results from
237 the fact that *A. odontolyticus* genome annotation in GenBank is incomplete[36]. Conversely, at the
238 species level, five of the six algorithms mapped a high number of reads to *Streptococcus agalactiae* for
239 both the HMP even and HMP staggered datasets (**Figure 2f, 2g**), but only a small number of reads for
240 this organism were present in the truth data. The relative abundance of *Streptococcus mutans* is lower
241 in the algorithm calls as compared to truth, while the relative abundance of *Streptococcus agalactiae*
242 is higher, suggesting that a number of the reads called for *S. agalactiae* are actually from *S. mutans*
243 (Figure 3b, 3d). This implies difficulty distinguishing between closely related species. Similarly, a high
244 number of reads are assigned correctly to the *Yersinia* and *Escherichia* genera by Kraken and
245 MetaScope (**Figure 2c**.) However, the algorithms under-assign reads for *Escherichia albertii* and over-
246 assign reads for *Yersinia pseudotuberculosis*, which indicates difficulty in distinguishing between these
247 species (Figure 2h).

248 Overall, algorithms were equally as able to identify organisms in the staggered datasets as in the even
249 datasets, suggesting that accurate read mapping depends more on the database supplied to the
250 algorithm rather than the abundance of the organism in the dataset. Additionally, for the bacterial
251 datasets, Kraken, MetaScope, LMAT, and MetaPhlAn generally agreed on read mapping assignments.
252 However, for the viral datasets, the algorithms missed different sets of organisms – i.e., in **Figure 3i**,
253 LMAT failed to map reads for HIV1, Influenza A virus, Marburg virus, and Machupo virus, whereas
254 MetaScope and Kraken correctly mapped reads for these organisms. However, MetaScope and Kraken
255 both failed to map reads for Human papillomavirus 5, SARS coronavirus, Human papillomavirus 32, and
256 Canine papillomavirus 3, while LMAT succeeded in mapping reads for these organisms. This suggests
257 that for viral datasets, it might be worthwhile to execute both LMAT and one of Kraken or MetaScope,
258 and calculate the union of the results.

259 The algorithms were also evaluated based on false positive hits (Figure 4). MetaCV and LMAT have
260 diverse error profiles – small numbers of reads are mapped to a high number of false positive
261 organisms. Our past experiences with the MetaScope algorithm suggest that this false positive profile
262 indicates an algorithm has difficulty classifying organisms that are not present in the reference database.
263 Ideally, when an algorithm encounters a novel organism, it should regress up the taxonomic tree until a
264 nearest neighbor for the unknown organism can be established. However, the algorithm may instead
265 report all reference organisms that match the unknown sample to a certain threshold. In contrast,
266 Kraken has a highly concentrated error profiles; fewer than 20 false positive organisms are reported, but
267 several thousand reads are mapped to each of them, suggesting high confidence calls. **Figure 4c** and **4d**
268 summarizes the top 20 organisms in terms number of mapped reads, indicating high agreement
269 between Kraken and MetaScope. On the list of false positive genera are several members of the
270 *Enterobacteriaceae* family, including *Shigella*, *Klebsiella*, and *Enterobacter*. The true positive genera
271 *Salmonella*, *Escherichia*, and *Yersinia* are members of this family as well. More difficult to explain is the

272 presence of the *Methanobolus* genus, which is a member of the kingdom Archaea and is distantly
273 related to the bacteria in the truth data.

274 For the viral datasets, MetaCV returned a high number of false positives and exhibited poor
275 performance. Kyasanur forest disease virus, a close relative to the true positive Omsk hemorrhagic
276 fever virus was the sole false positive for LMAT, and MetaScope did not report any false positive
277 organisms for either viral dataset.

278 Finally, the gene calling capabilities of the algorithms were evaluated (Figure 5). Only MetaScope,
279 LMAT, and MetaCV call genes, so these three were included for analysis. For the HMP
280 Even/Staggered, Bacteria, and Virus Staggered datasets, MetaScope identified the most genes
281 correctly out of the three algorithms. LMAT identified more correct genes on the Virus Even dataset
282 (101, compared to 93 for MetaScope).
283

284 **Conclusions**

285 In summary, *in silico* datasets with known truth data for read and gene distribution across different
286 taxons serve as a valuable tool for evaluating algorithm performance. The HMP Even/Staggered,
287 Bacteria, Virus Even/Staggered, and Human datasets generated with FASTQsim elucidate multiple
288 patterns in performance for leading metagenomics algorithms. No algorithm outperformed the others
289 in all categories, and the algorithm of choice strongly depends on analysis goals. For bacterial datasets,
290 MetaPhlAn is a clear winner, achieving the lowest runtime, highest ratio of true positives to false
291 positives, and the most precise read mapping. However, MetaPhlAn does not assign genes and does not
292 work on taxons other than Bacteria. LMAT is a clear winner for viral datasets in terms of accuracy, and
293 also provides gene calling functionality. The algorithm most closely matched the relative abundance
294 profile of the truth genera and species across all datasets. However, LMAT also reported the highest
295 rate of false positive genera and species calls on the bacterial datasets. Kraken and MetaScope were
296 the runners up in terms of runtime, ratio of true positives to false positives, and read mapping.
297 MetaScope also performed best for gene mapping, which Kraken does not do. These algorithms
298 performed solidly across all categories evaluated and can be applied most universally across versatile
299 metagenomic applications. MetaPhyler and MetaCV came in last for runtime, ratio of true positives to
300 false positives, and read mapping. They also do not provide results out of the box for viral datasets.

301 Although viral, bacterial, and human datasets were simulated for this study, the techniques described
302 here can be extended to evaluate metagenomic algorithm performance for other taxa. For example,
303 fungal contamination incidents at medical facilities such as the 2012 incident at the New England
304 Compounding Center[37] can be contained more quickly and effectively with the aid of metagenomic
305 sequencing. Other potential applications include rapid diagnosis of parasite infections[38].

306 List of abbreviations

307

308 GB - gigabyte

309 RAM – random-access memory

310 s – seconds

311 TB - terabyte

312 x – fold coverage

313

314 Competing interests

315

316 The authors declare that they have no competing interests.

317 Ethics Committee Approval

318 Ethics approval was not required for this study because all data was generated *in silico* using references
319 available in GenBank, as indicated in Supplementary Tables 1-3.

320

321 Authors' contributions

322 AS implemented FASTQSim updates and generated *in silico* datasets. AS and NC benchmarked
323 algorithm performance on evaluation datasets. AS and DR wrote the manuscript. DR conceived of the
324 study. All authors read and approved the final manuscript.

325

326 References

327

- 328 [1] H. G. Martin, N. Ivanova, V. Kunin, F. Warnecke, K. W. Barry, A. C. McHardy, *et al.*, "Metagenomic
329 analysis of two enhanced biological phosphorus removal (EBPR) sludge communities," *Nat*
330 *Biotech*, vol. 24, pp. 1263-1269, 10//print 2006.
- 331 [2] G. W. Tyson, J. Chapman, P. Hugenholtz, E. E. Allen, R. J. Ram, P. M. Richardson, *et al.*,
332 "Community structure and metabolism through reconstruction of microbial genomes from the
333 environment," *Nature*, vol. 428, pp. 37-43, 03/04/print 2004.
- 334 [3] J. C. Venter, K. Remington, J. F. Heidelberg, A. L. Halpern, D. Rusch, J. A. Eisen, *et al.*,
335 "Environmental Genome Shotgun Sequencing of the Sargasso Sea," *Science*, vol. 304, pp. 66-74,
336 2004.
- 337 [4] J. L. Morgan, A. E. Darling, and J. A. Eisen, "Metagenomic Sequencing of an *In Vitro*-Simulated
338 Microbial Community," *PLoS ONE*, vol. 5, p. e10209, 2010.
- 339 [5] InnoCentive. (2013). *Identifying organisms from a stream of DNA sequences*. Available:
340 <https://www.innocentive.com/ar/challenge/9933138>
- 341 [6] L. Ilie and S. Ilie, "Multiple spaced seeds for homology search," *Bioinformatics*, vol. 23, pp. 2969-
342 2977, 2007.
- 343 [7] B. Buchfink, C. Xie, and D. H. Huson, "MetaScope - Fast and accurate identification of microbes in
344 metagenomic sequencing data.," *arXiv.org*, submitted.
- 345 [8] D. E. Wood and S. L. Salzberg, "Kraken: ultrafast metagenomic sequence classification using exact
346 alignments," *Genome Biology*, vol. 15, 2014.
- 347 [9] N. Segata, L. Waldron, A. Ballarini, V. Narasimhan, O. Jousson, and C. Huttenhower,
348 "Metagenomic microbial community profiling using unique clade-specific marker genes," *Nat*
349 *Meth*, vol. 9, pp. 811-814, 08//print 2012.
- 350 [10] S. Ames, J. E. Allen, D. A. Hysom, G. S. Lloyd, and M. B. Gokhale, "Design and Optimization of a
351 Metagenomics Analysis Workflow for NVRAM," in *Parallel & Distributed Processing Symposium*
352 *Workshops (IPDPSW), 2014 IEEE International*, 2014, pp. 556-565.
- 353 [11] S. K. Ames, D. A. Hysom, S. N. Gardner, G. S. Lloyd, M. B. Gokhale, and J. E. Allen, "Scalable
354 metagenomic taxonomy classification using a reference genome database," *Bioinformatics*, vol.
355 29, pp. 2253-2260, 2013.

- 356 [12] B. Van Essen, H. Hsieh, S. Ames, and M. Gokhale, "DI-MMAP: A High Performance Memory-Map
357 Runtime for Data-Intensive Applications," in *High Performance Computing, Networking, Storage
358 and Analysis (SCC), 2012 SC Companion:*, 2012, pp. 731-735.
- 359 [13] J. Liu, H. Wang, H. Yang, Y. Zhang, J. Wang, F. Zhao, *et al.*, "Composition-based classification of
360 short metagenomic sequences elucidates the landscapes of taxonomic and functional enrichment
361 of microorganisms," *Nucleic Acids Research*, vol. 41, 2013.
- 362 [14] B. Liu, T. Gibbons, M. Ghodsi, T. Trengen, and M. Pop, "Accurate and fast estimation of taxonomic
363 profiles from metagenomic shotgun sequences," *BMC Genomics*, vol. 12, 2010.
- 364 [15] G. Bonilla-Rosso, "Lessons learned from simulated metagenomic datasets," *Encyclopedia of
365 Metagenomics*, pp. 1-8, 2014.
- 366 [16] K. Mavromatis, N. Ivanova, K. Barry, H. Shapiro, E. Goltsman, A. McHardy, *et al.*, "Use of simulated
367 data sets to evaluate the fidelity of metagenomic processing methods," *Nature Methods*, vol. 4,
368 pp. 495-500, 2007.
- 369 [17] M. Pignatelli and A. Moya, "Evaluating the fidelity of *de novo* short read metagenomic assembly
370 using simulated data," *PLOS One*, 2011.
- 371 [18] D. Mende, A. Waller, S. Sunagawa, A. Jarvelin, M. Chan, M. Arumugam, *et al.*, "Assessment of
372 metagenomic assembly using simulated next generation sequencing data," *PLOS one*, 2012.
- 373 [19] W. Huang, L. Li, J. Myers, and G. Marth, "ART: a next-generation sequencing read simulator,"
374 *Bioinformatics*, vol. 28, pp. 593-594, 2011.
- 375 [20] The NIH HMP Working Group, "The NIH Human Microbiome Project," *Genome Research*, vol. 19,
376 pp. 2317-2323, 2009.
- 377 [21] N. I. o. A. a. I. Disease. (2015). *NIAID Category A, B, and C Priority Pathogens*. Available:
378 <https://www.niaid.nih.gov/topics/biodefenselated/biodefense/pages/cata.aspx>
- 379 [22] D. Benson, M. Cavanaugh, K. Clark, I. Karsch-Mizrachi, D. Lipman, J. Ostell, *et al.*, "GenBank,"
380 *Nucleic Acids Research*, vol. 41, pp. D36-42, 2013.
- 381 [23] Y. Chen, C. Lin, C. Wang, H. Wu, and P. Hwant, "An optimized procedure greatly improves EST
382 vector contamination removal," *BMC Genomics*, vol. 8, 2007.
- 383 [24] P. Bork and A. Bairoch, "Go hunting in sequence databases but watch out for the traps," *Trends
384 Genet*, vol. 12, pp. 425-427, 1996.
- 385 [25] "Quality control in databanks for molecular biology," *Bioessays*, vol. 22, pp. 1024-1034, 2000.
- 386 [26] G. Seluja, A. Farmer, M. McLeod, C. Harger, and P. Schad, "Establishing a method of vector
387 contamination identification in database sequenes," *Bioinformatics* vol. 15, pp. 106-110, 1999.
- 388 [27] A. Shcherbina, "FASTQSim: platform-independent data characterization and in silico read
389 generation for NGS datasets," *BMC Research Notes*, vol. 7, p. 533, 2014.
- 390 [28] P. Cock, T. Antao, J. Chang, B. Chapman, C. Cox, A. Dalke, *et al.*, "Biopython: freely available
391 Python tools for computational molecular biology and bioinformatics," *Bioinformatics*, vol. 25, pp.
392 1422-1423, 2009.
- 393 [29] D. O. Ricke. (2011). *GenBankParser*. Available:
394 <https://github.com/doricke/BioTools/tree/master/GenBankParser>
- 395 [30] B. Ondov, N. Bergman, and A. Phillippy, "Interactive metagenomic visualization in a Web
396 browser," *BMC Bioinformatics*, vol. 12, 2011.
- 397 [31] S. Naccache, S. Federman, N. Veeraraghavan, M. Zaharia, D. Lee, E. Samayoa, *et al.*, "A cloud-
398 compatible bioinformatics pipeline for ultrarapid pathogen identification from next-generation
399 sequencing of clinical samples," *Genome Research*, vol. 24, pp. 1180-1192, 2014.
- 400 [32] N. Daniels, A. Gallant, J. Peng, L. Cowen, M. Baym, and B. Berger, "Compressive genomics for
401 protein databases," *Bioinformatics*, vol. 29, pp. i283-i290, 2013.
- 402 [33] D. Huang, B. Sherman, and R. Lempicki, "Systematic and integrative analysis of large gene lists
403 using DAVID Bioinformatics Resources," *Nature Protocols*, vol. 4, pp. 44-57, 2009.
- 404 [34] C. Alkan, S. Sajjadian, and E. Eichler, "Limitations of next-generation genome sequence assembly,"
405 *Nature Methods*, vol. 8, pp. 61-65, 2010.
- 406 [35] T. Smith and S. Porter, "Development and role of the human reference sequence in personal
407 genomics," *Wiley Online Library*, 2014.

- 408 [36] N. Sarkonen, "Oral Actinomyces Species in Health and Disease: Identification, Occurrence and
409 Importance of Early Colonization," ed: National Public Health Institute, 2007.
- 410 [37] R. Vijayakumar, T. Sandle, and C. Manoharan, "Review of fungal contamination in pharmaceutical
411 products and phenotypic identification of contaminants by conventional methods," *European*
412 *Journal of Parenteral and Pharmaceutical Sciences*, vol. 17, pp. 4-19, 2012.
- 413 [38] M. Ndao, "Diagnosis of parasitic diseases: old and new approaches," *Interdisciplinary Perspectives*
414 *on Infectious Diseases*, vol. 2009, 2009.
- 415

Figures

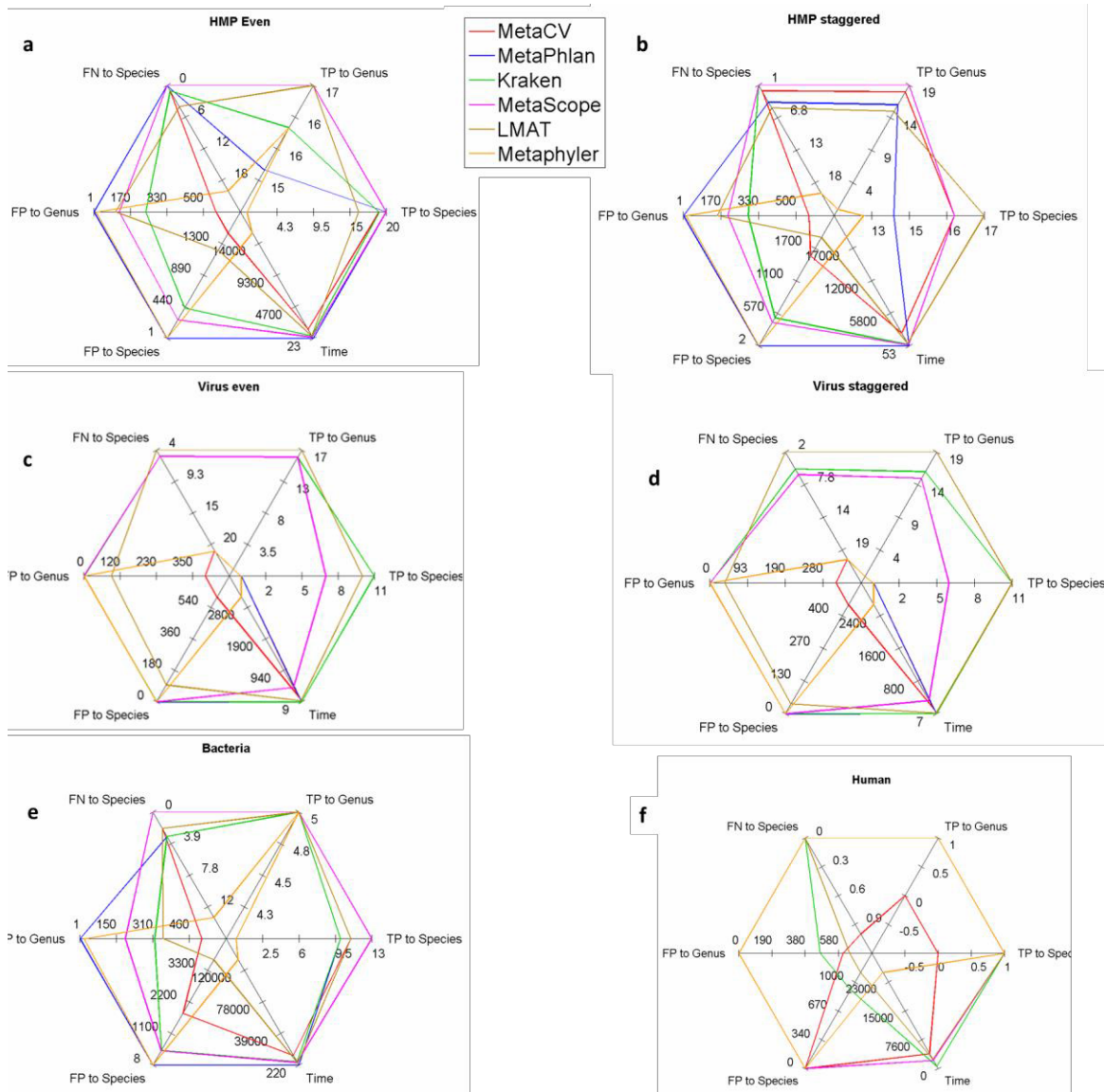
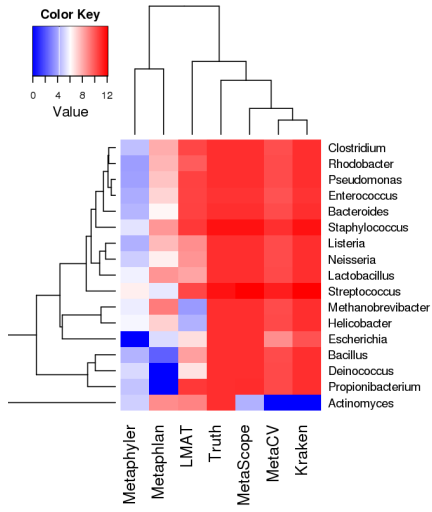
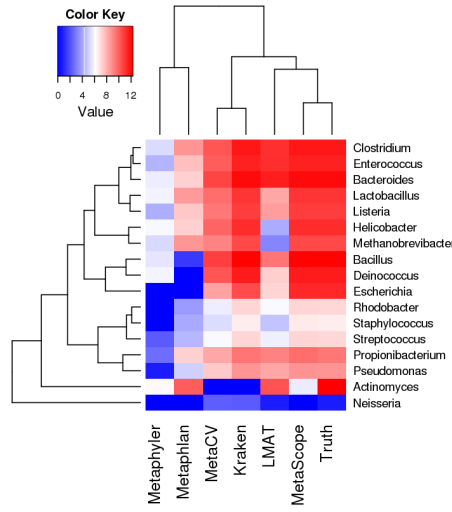


Figure 1. Performance metrics for 6 metagenomic analysis algorithms across the 6 in silico evaluation datasets. Algorithms evaluated include MetaCV (red line), MetaPhlan (blue line), Kraken (green line), MetaScope (pink line), LMAT (brown line), MetaPhyler (orange line). Metrics evaluated include true positives (TP) to genus level, TP to species level, false positives (FP) to genus level, FP to species level, false negatives (FN) to species level, and runtime in seconds. Values indicative of high performance are at the periphery of the radar plot, values indicative of poor performance are at the center of the plot. **a.** HMP dataset with even coverage. **b.** HMP dataset with staggered coverage. **c.** Virus dataset with even coverage. **d.** Virus dataset with staggered coverage. **e.** Bacterial dataset. **f.** Human dataset. The MetaPhlan algorithm failed to run on the human dataset.

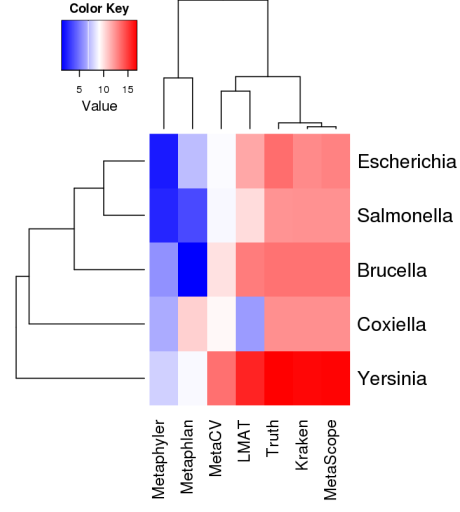
a HMP Even, Genus



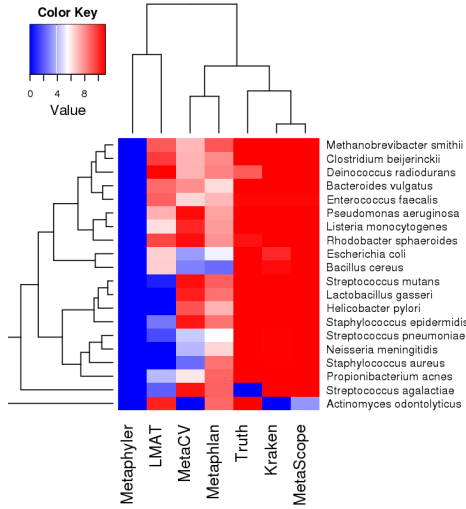
b HMP Staggered, Genus



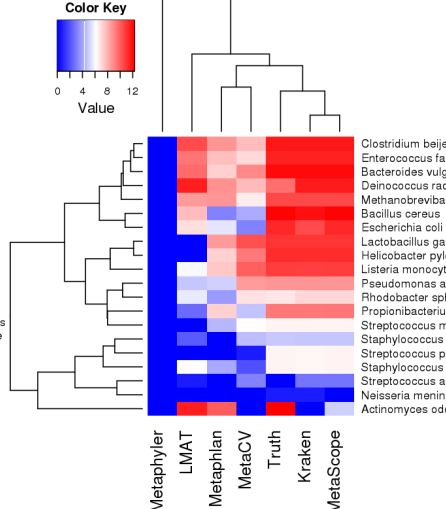
c Bacteria, Genus



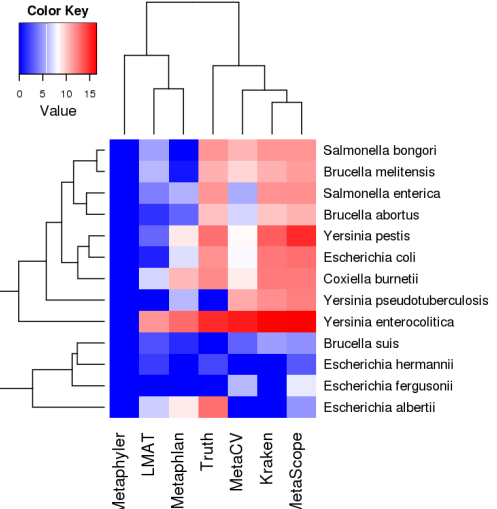
d HMP Even, Species



e HMP Staggered, Species



f Bacteria, Species



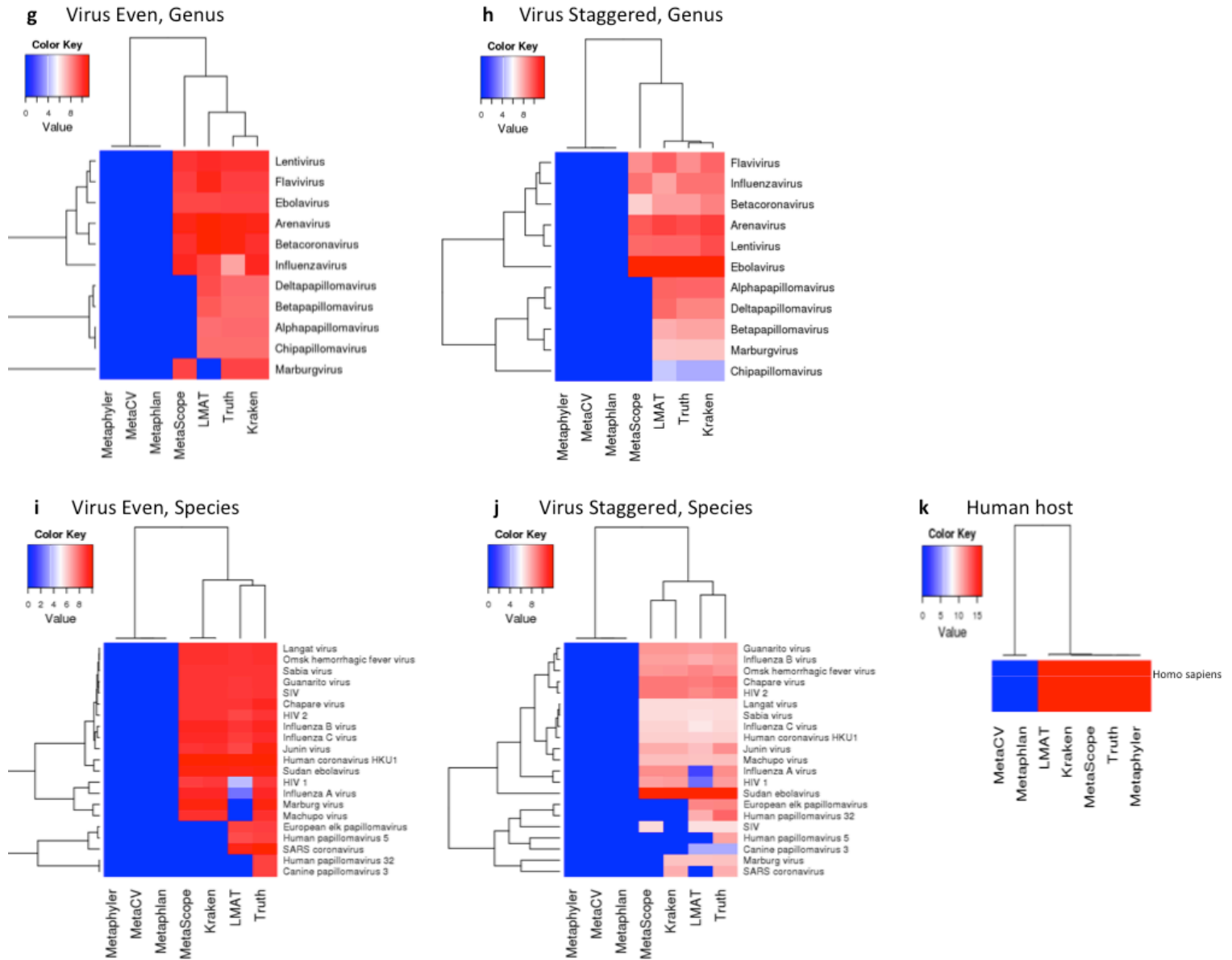


Figure 2. Number of correctly assigned reads to each organism at the genus and species level. Heatmap color scales are log10 (number of correctly assigned reads). The “Truth” column indicates the number of reads spiked into the FASTQ input file for the specified genus or species. **a. – e.** Reads mapped correctly to the genus level for the HMP even, HMP staggered, bacteria, virus even, virus staggered datasets, respectively. **f. – k.** Reads mapped correctly to the species level for HMP even, HMP staggered, bacteria, virus even, virus staggered, and human datasets, respectively.

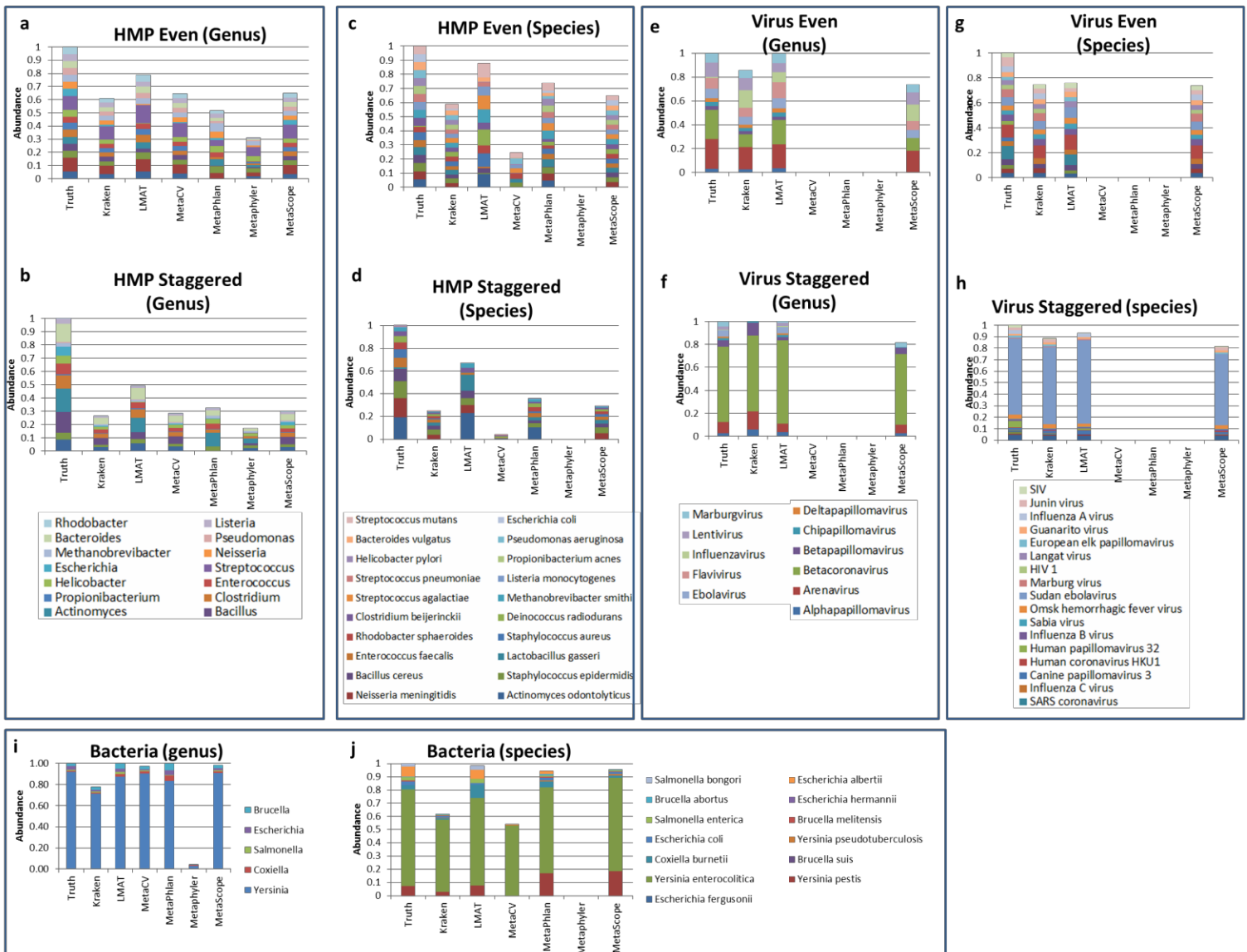
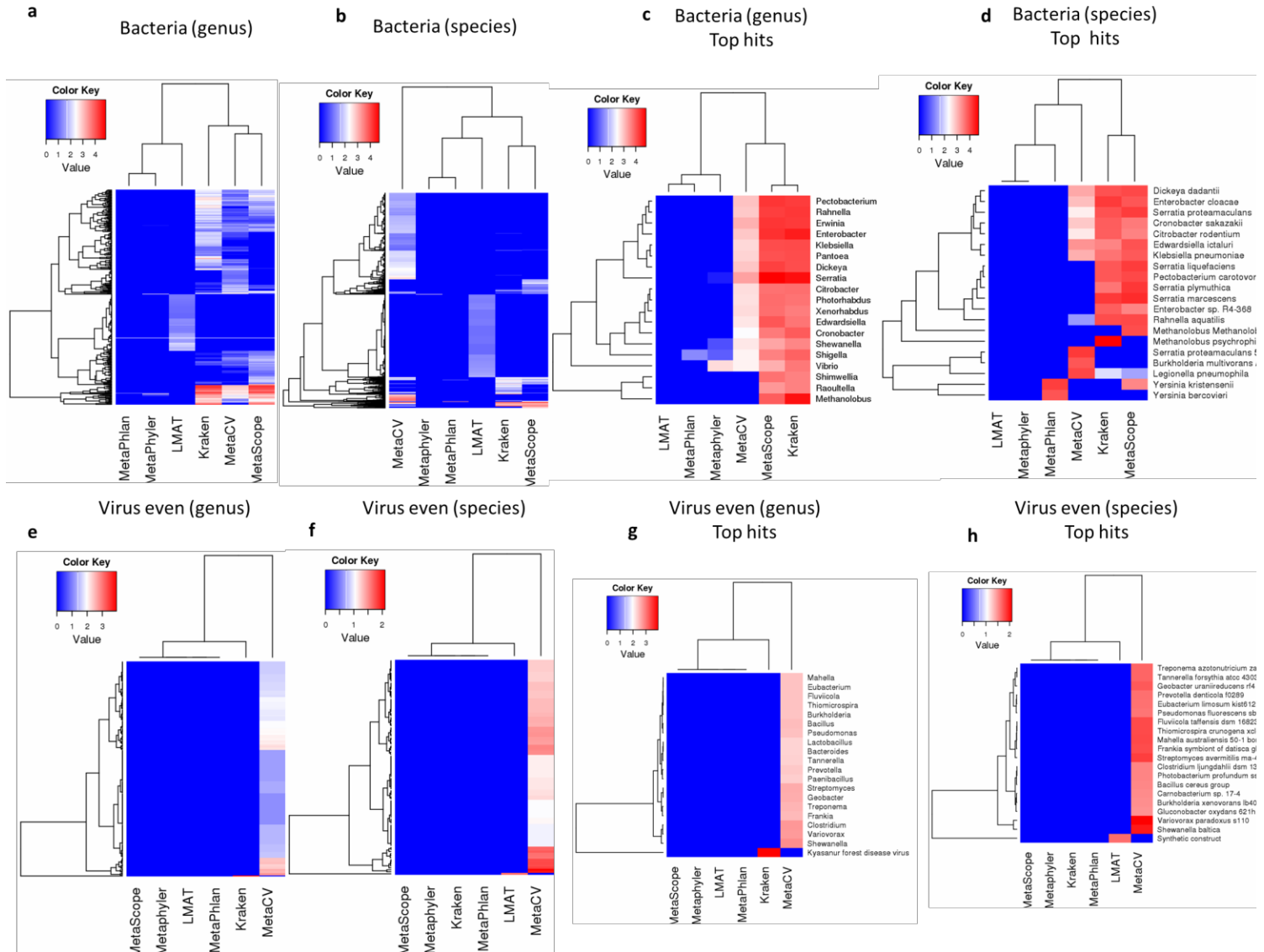


Figure 3. Relative abundance of organisms to the species and genus level. “Truth” column indicates relative abundance of genera and species added to the in silico FASTQ input file.

3



4
5

6

7

8

9

10

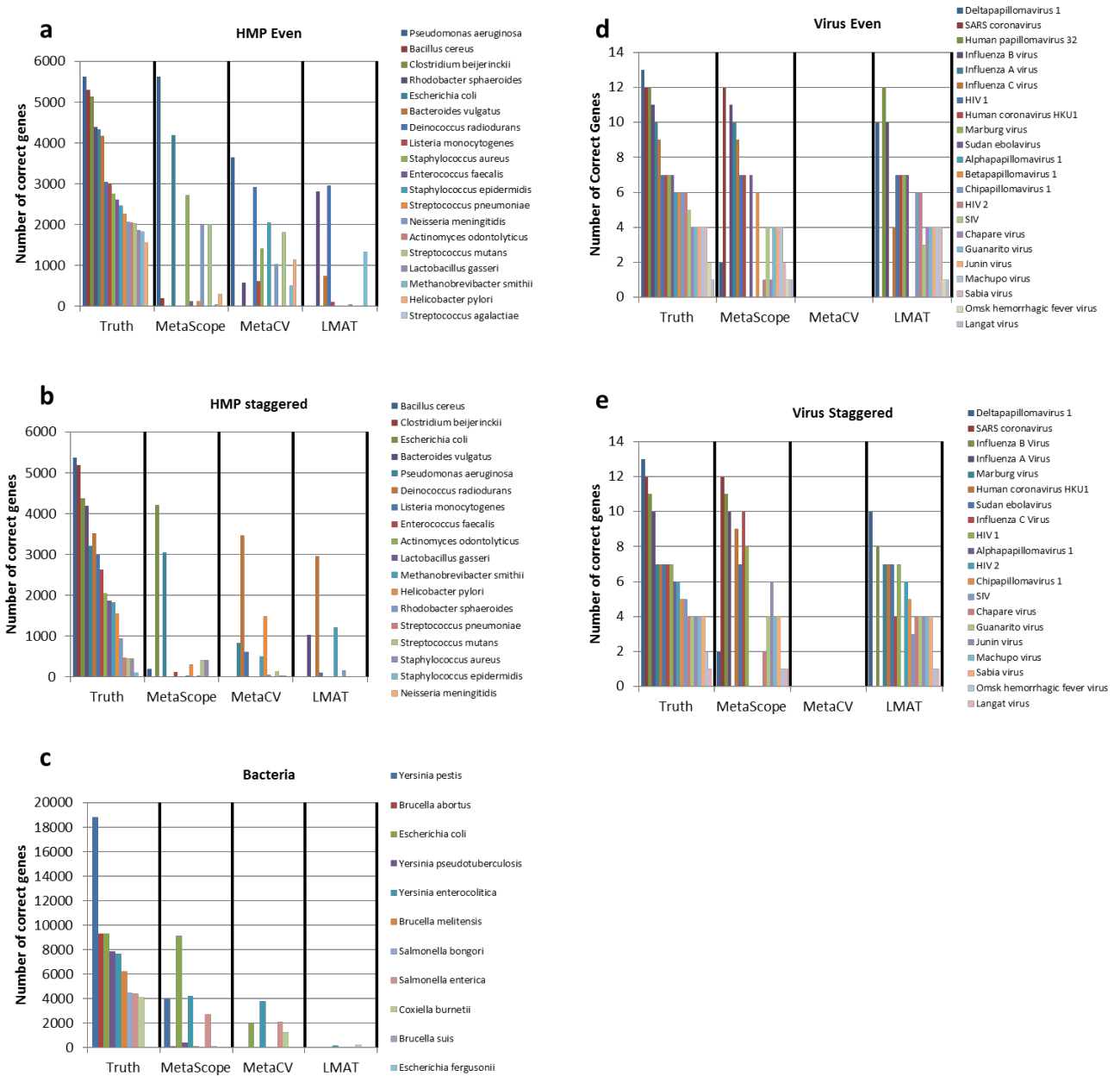
11

12

13

Figure 4. False positive organisms identified to the genus and species level by the 6 metagenomic algorithms. Heatmap color scales are log₁₀ (number of incorrectly assigned false positive reads) for a genus or species. **a.** All false positive genera identified in the bacterial dataset. **b.** All false positive species identified in the bacterial dataset. **c.** 20 false positive genera for the bacterial dataset with the most assigned reads. **d.** 20 false positive species for the bacterial dataset with the most assigned reads. **e.** All false positive genera identified in the virus even dataset. **f.** All false positive species identified in the virus even dataset. **g.** 20 false positive genera for the virus even dataset with the most assigned reads. **h.** 20 false positive species for the virus even dataset with the most assigned reads.

464



465

466

467

468

469

470

Figure 5. Number of genes correctly identified to the species level across the 5 evaluation datasets (the 6th evaluation dataset consisting of human host reads is not shown). “Truth” column indicates the number of genes with non-zero read coverage in the dataset. MetaScope, MetaCV, and LMAT algorithms provide gene assignment capabilities; Kraken, MetaPhyler, and MetaPhlAn do not call genes and were not included in this evaluation.

471 **Tables**

472 **Table 1.** Radar plot area in normalized units across six evaluation datasets. Higher areas, indicative of better
 473 performance, are colored in blue.
 474

Dataset	Human	Virus Staggered	Virus Even	Bacteria	HMP Staggered	HMP Even	Area Sum
MetaScope	2.54	1.90	2.14	2.23	2.03	2.32	13.15
Kraken	1.64	2.36	1.21	1.75	1.90	1.76	10.62
LMAT	1.41	2.45	2.25	1.38	1.46	1.59	10.54
MetaPhlAn		0.48	0.99	2.25	1.88	2.02	7.62
MetaCV	0.82	0.09	0.14	1.43	1.25	1.08	4.80
MetaPhyler	1.88	0.60	0.09	0.67	0.57	0.63	4.44

475 **Table 2.** Algorithm runtime in seconds across six evaluation datasets.
 476

Dataset	Human	Virus Staggered	Virus Even	Bacteria	HMP Staggered	HMP Even
MetaScope	2160	327	427	3686	261	233
Kraken	600	7	9	4400	300	400
LMAT	2428	20	39	2700	502	427
MetaPhlAn	Seg Fault	12	12	220	53	23
MetaCV	3873	120	150	11966	2337	1322
MetaPhyler	25200	2640	3100	129600	19231	15480

477
 478 **Supplementary Materials**
 479

480
 481 **Supplementary Table 1.** Source organisms and coverage levels for HMP Even and HMP Staggered datasets.

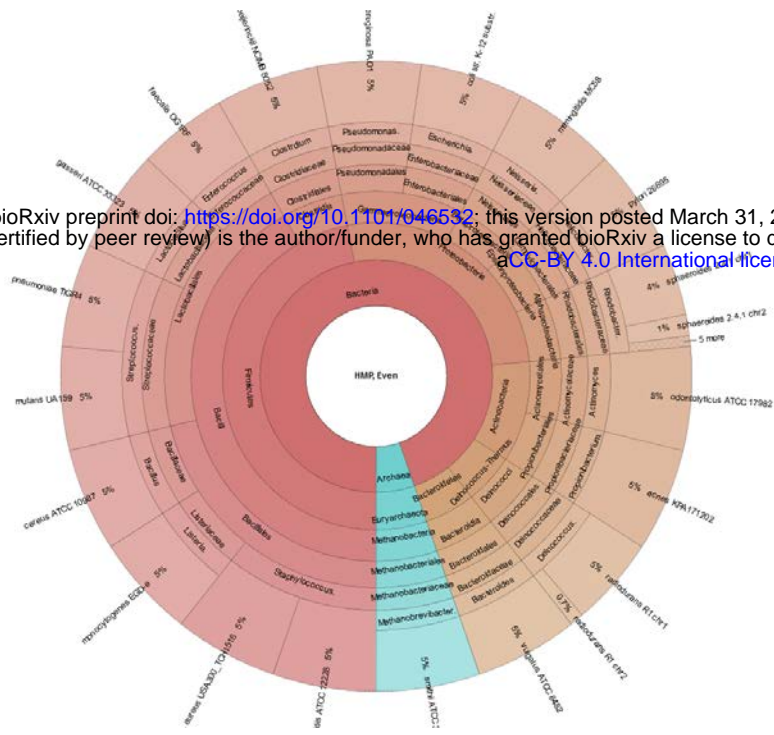
482 **Supplementary Table 2.** Source organisms and coverage levels for Bacterial dataset.

483
 484 **Supplementary Table 3.** Source organisms and coverage levels for Virus Even and Virus Staggered datasets.
 485

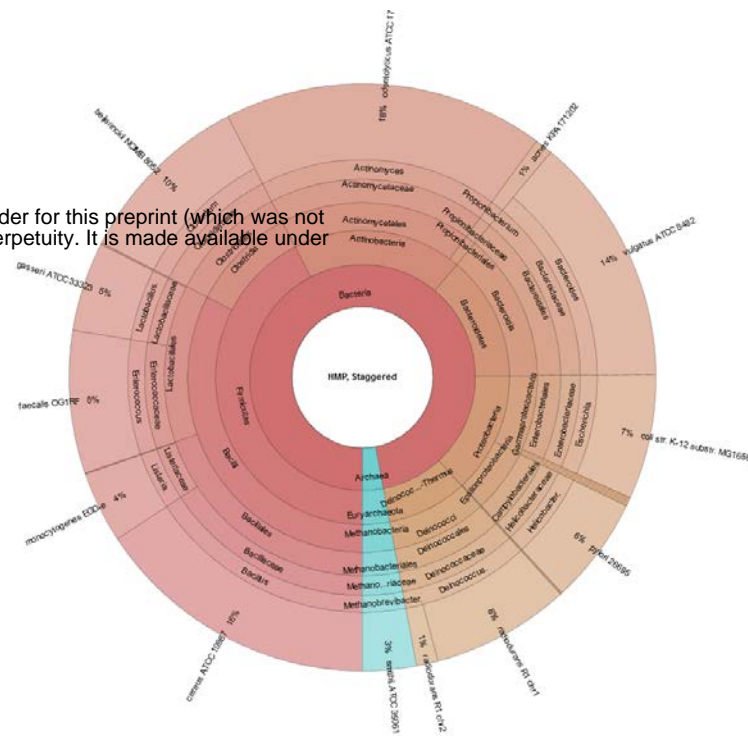
486 **Figure S1.** FASTQSim in silico dataset composition to strain level. **a.** 20 bacteria from the Human Microbiome
 487 Project (HMP), even coverage levels. **b.** Same 20 bacteria from HMP, staggered coverage levels. **c.** 22 species of
 488 viruses across 11 genera, even coverage levels. **d.** Same 22 species of viruses, staggered coverage levels. **e.** 33
 489 strains of bacteria representing 13 species and 5 genera. See Krona HTML files for a-e.
 490
 491

bioRxiv preprint doi: <https://doi.org/10.1101/046532>; this version posted March 31, 2016. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

a



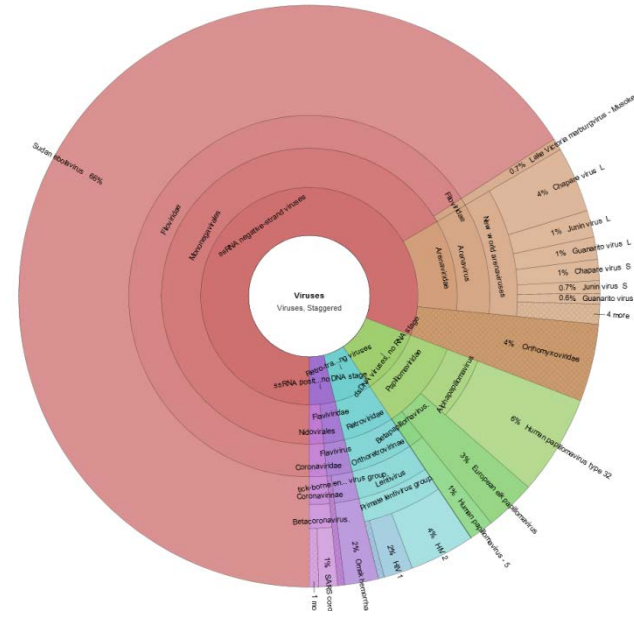
b



c



d



e

