

1 ***In Silico* Benchmarking of Metagenomic Tools for Coding Sequence Detection Reveals**
2 **the Limits of Sensitivity and Precision**

3

4 Jonathan Louis Golob¹, Samuel Schwartz Minot²

5

6 ¹ Infectious Diseases, Internal Medicine, Michigan Medicine, University of Michigan, Ann Arbor,
7 Michigan, USA

8 ² Microbiome Research Initiative, Fred Hutchinson Cancer Research Center, Seattle,
9 Washington, USA

10

11 Corresponding Author:

12 Samuel Schwartz Minot

13 1100 Fairview Ave N, E4-100

14 Seattle, WA. 98109-1024

15 206-667-2884

16 sminot@fredhutch.org

17

18 **Abstract**

19 High-throughput sequencing can establish the functional capacity of a microbial community by
20 cataloging the protein-coding sequences (CDS) present in the metagenome of the community.
21 The relative performance of different computational methods for identifying CDS from whole-
22 genome shotgun sequencing (WGS) is not fully established.

23
24 Here we present an automated benchmarking workflow, using synthetic shotgun sequencing
25 reads for which we know the true CDS content of the underlying communities, to determine the
26 relative performance (sensitivity, positive predictive value or PPV, and computational efficiency)
27 of different metagenome analysis tools for extracting the CDS content of a microbial community.

28
29 Assembly-based methods are limited by coverage depth, with poor sensitivity for CDS at < 5X
30 depth of sequencing, but have excellent PPV. Mapping-based techniques are more sensitive at
31 low coverage depths, but can struggle with PPV. We additionally describe an expectation
32 maximization based iterative algorithmic approach which we show to successfully improve the
33 PPV of a mapping based technique while retaining improved sensitivity and computational
34 efficiency.

35

36

37

38

39 **Introduction**

40 High throughput (or “next-generation”) sequencing has uncovered the hidden complexity of
41 microbial communities living within and upon the human body, as well as the link between the
42 human microbiome and health [1–4]. The taxonomic composition of a microbial community can
43 be inferred by sequencing PCR amplicons spanning variable regions of a taxonomically
44 informative gene (i.e. the 16S rRNA gene or the CPN60 gene)[5–8]. Alternatively, DNA
45 recovered from a sample can be put through Whole-Genome Sequencing (WGS), which
46 samples the complete genomic content of a sample via random fragmentation and sequencing
47 [9]. WGS differs from amplicon sequencing by (1) providing genomic data from all organisms in
48 a sample—not limited to any single domain of life; (2) enabling a high degree of taxonomic
49 resolution which identifies the subspecies and strains present in a sample; and (3) generating a
50 “functional” metagenomic profile of the protein coding sequences (CDS) that are present in a
51 sample in addition to the organisms which contain those genes [10]. While the term “functional”
52 can often be used to describe predicted metabolic pathways, in this case are limiting our scope
53 to the identification of CDS without presupposing knowledge of any annotations.

54

55 There are three broad computational approaches used to generate an estimate of the functional
56 metagenome (CDS content) of a microbial community from WGS reads: (1) The inferred
57 taxonomic composition can be used to construct a custom database of protein-coding genes
58 from the set of reference organisms detected in the sample (e.g. HUMAnN2, MIDAS) [11,12].
59 (2) *De novo* assembly, in which the WGS reads are combined into contigs, which can be further
60 used to identify open reading frames (e.g. metaSPAdes, IDBA-UD) [13,14]. (3) The WGS reads
61 can be directly mapped (aligned) to a closed reference of protein coding sequences (which is
62 also a downstream component of HUMAnN2 and MIDAS).

63

64 Proteins can evolve by duplication events, truncation, homologous recombination, and other
65 means that result in the sharing of highly conserved domains between otherwise distinct CDS
66 [15]. As a result, mapping of reads to a closed reference of CDS is challenged by the fact that
67 some reads may align equally well to multiple references: “multi-mapping” reads.

68

69 Metagenomic tools have been benchmarked extensively for their ability to determine the
70 taxonomic composition of a microbial community [16–19]. The relative ability of metagenomic
71 analysis approaches and tools to accurately infer the CDS catalog of a microbial community has
72 yet to be established. Additionally, benchmarking efforts are often limited in their long-term utility
73 by the practical challenges of repeating the computational analysis with the addition of newly
74 available tools. We address this core challenge of benchmarking by implementing our analysis
75 within a workflow management tool, Nextflow [20], which achieves a high degree of
76 reproducibility by executing each component task within Docker containers, a portable and fixed
77 computational environment.

78

79 Here we establish sensitivity and positive predictive value (PPV) of computational tools for
80 determining the CDS content of a microbial community metagenome, using synthetic
81 communities and reads generated *in silico* for which we know the true CDS content of the
82 community. We establish that assembly-based approaches achieve a near-perfect PPV, but
83 struggle with sensitivity for CDS at a low sequencing coverage depth. Mapping-based
84 approaches are more sensitive, particularly at low coverage depths, but struggle with PPV. We
85 introduce an expectation-maximization based approach for mapping based metagenomics that
86 retains the sensitivity and improves the PPV of CDS calls close to that of assembly-based
87 approaches.

88

89 **Materials & Methods**

90 **Evaluating Computational Tools**

91 All of the analytical steps for analyzing computational tools for CDS detection from
92 metagenomes were executed within a single analytical workflow ('evaluate-gene-detection.nf')
93 which can be downloaded from [95](https://github.com/FredHutch/evaluate-gene-level-
94 <u>metagenomics-tools</u> and executed via Nextflow. That analytical workflow follows this approach:</p></div><div data-bbox=)

96 1. Simulate metagenomes (n=100)

97 a. Randomly select host-associated genomes from NCBI/RefSeq (n=20). (A list of
98 genomes from host-associated organisms is available in the supplemental
99 materials.)

100 b. Make a file with all of the CDS records from those genomes

101 c. Assign sequencing depths for each genome from a log-normal distribution
102 (mean=5x, std=1 log), with a maximum possible depth of 100x

103 d. Make a file with the depth of sequencing for each CDS from step (1b) above

104 e. Simulate reads from whole genome sequences via ART (paired-end read length
105 250bp, mean fragment length 1kb +/- 300bp)

106 f. Interleave paired end FASTQ data

107 2. Run tools

108 a. For assembly-based tools, perform assembly from paired end FASTQ data and
109 predict CDS records from the resulting contigs

110 b. For mapping-based tools, run the tool and then extract the FASTA for all
111 detected CDS records

112 3. Perform evaluation

113 a. For each tool, align the FASTA with all detected CDS records against the set of
114 truly present CDS records (from step 1b)

- 115 i. Prior to alignment, both sets of FASTAs are clustered at 90% amino acid
116 identity to account for sets of homologous genes in the simulated
117 metagenome
- 118 b. Filter to the top hit for each detected CDS
- 119 c. Assign each detected CDS as:
- 120 i. True positive: The detected CDS is the mutual best hit for a truly present
121 CDS
- 122 ii. False positive: The detected CDS does not align against any truly present
123 CDS
- 124 iii. Duplicate: The detected CDS aligns against a truly present CDS, but is
125 not the best hit (i.e. there are multiple non-overlapping detected CDS
126 records that each align against a single truly present CDS).
- 127 d. Calculate accuracy metrics:
- 128 i. Sensitivity is calculated as the number of true positives (TP) divided by
129 the number of true positives and false negatives (FN): $TP / (TP + FN)$
- 130 ii. Positive Predictive Value is calculated as the number of true positives
131 (TP) divided by the number of true positives and false positives (FP): $TP /$
132 $(TP + FP)$
- 133 iii. Uniqueness is calculated as the number of true positives
134 (TP) divided by the number of true positives and duplicates (DUP): $TP /$
135 $(TP + DUP)$

136 **FAMLI Implementation**

137 FAMLI is available as an open source software package on GitHub at
138 <https://github.com/FredHutch/FAMLI>. In addition, Docker images are provided at
139 <https://quay.io/repository/fhcrc-microbiome/famli> to facilitate easy usage by the research

140 community with a high degree of computational reproducibility. FAMLI can be run with the single
141 executable "famli", which encompasses:

- 142 1. Downloading reference data and query FASTQ files (supporting local paths, FTP, and
143 Amazon Web Service (AWS) object storage)
- 144 2. Aligning query FASTQ files in amino acid space with DIAMOND
- 145 3. Parsing the translated alignments
- 146 4. Running the FAMLI algorithm to filter unlikely reference peptides and assign multi-
147 mapping query reads to a single unique reference.
- 148 5. Summarizing the results in a single output file
- 149 6. Copying the output file to a remote directory (supporting local paths, FTP, and AWS
150 object storage)

151 The help flag ("-h" or "--help") can be used to print a complete list of options, including the flags
152 used to run the filtering process starting from step 4 above..

153 **FAMLI Overall Approach**

- 154 1. Align all input nucleotide reads in against a reference database of peptides; UniRef 90
155 was used for this study [21].
- 156 2. Calculate the coverage depth (CD) across the length of each reference, representing the
157 number of reads aligning to each amino acid position of the reference.
- 158 3. Filter out any reference sequences with highly uneven coverage:

$$159 \frac{CD_{STD}}{CD_{Mean}} < 1.0 \quad (1)$$

160 Where STD is standard deviation of per-base coverage values.

- 161 4. Calculate initial score for a given query coming from a subject using the alignment
162 bitscores to weight the relative possibilities for a given query, normalizing the scores to
163 total to 1 for a given query.

- 164 5. Iteratively, until no further references are pruned or a maximum number of iterations is
165 reached:
- 166 i. WEIGHTING and RENORMALIZING: The score of queries being from a subject
167 from the prior iteration are weighted by the sum of scores for a given subject, and
168 then renormalized to sum to 1 for each query.
 - 169 ii. PRUNING. Determine the maximum likelihood for each query. Prune away all
170 other likelihoods for the query below a threshold.
- 171 6. Repeat filtering steps 2-3 using the set of deduplicated alignments resulting from step 4.
172

173 Here are some examples:

- 174 ☐ For reference A and reference B that both have some aligning query reads, if **there is**
175 **uneven depth for reference A** but relatively even depth across reference B, then
176 **reference A is removed from the candidate list** while reference B is kept as a
177 candidate.
- 178 ☐ If query **read #1 aligns equally-well to reference A and reference C**, but **there is 2x**
179 **more query read depth for reference A as compared to reference C** across the
180 entire sample, then **reference C's alignment is removed from the list of candidates**
181 **for query read #1.**

182 A more detailed description of the method is available in the supplemental materials. An
183 interactive demonstration of our algorithm is available as a Jupyter notebook is available at
184 <https://github.com/FredHutch/FAMLI/blob/master/schematic/FAMLI-schematic-figure-GB.ipynb>
185

186 **Comparison of FAMLI to HUMAnN2, SPAdes, Top Hit, and Top 20**

187 The version of FAMLI presented in this paper is v1.3, which can be found at
188 <https://github.com/FredHutch/FAMLI/releases/tag/v1.3>. FAMLI was executed in this analysis

189 using a Docker image hosted at <https://quay.io/repository/fhcr-microbiome/famli> with the tag
190 v1.3 (sha256:25c34c73964f).

191 The version of DIAMOND used for translated nucleotide alignments in this analysis is
192 DIAMOND v0.9.10 using a Docker image compiled from [https://github.com/FredHutch/docker-](https://github.com/FredHutch/docker-diamond)
193 [diamond](https://github.com/FredHutch/docker-diamond) and available at <https://quay.io/repository/fhcr-microbiome/docker-diamond> as
194 v0.9.23--0 (sha256: 0f06003c4190).

195 Comparative analysis of the simulated communities used HUMAnN2 v0.11.1--py27_1, and all
196 code used to run HUMAnN2 can be found in the GitHub repository
197 <https://github.com/FredHutch/docker-humann2> (v0.11.1--6), which is based on the BioBakery
198 Docker image quay.io/biocontainers/humann2:0.11.1--py27_1. The Docker image used to run
199 HUMAnN2 is available at <https://quay.io/repository/fhcr-microbiome/humann2> as v0.11.2--1
200 (sha256:d6426bda36ca).

201 The code used to run SPAdes is maintained by BioContainers and is available at
202 <https://quay.io/repository/biocontainers/spades> as 3.13.0--0 (sha256:9f097c5d6d79).

203 The code used to run megahit is maintained by BioContainers and is available at
204 <https://quay.io/repository/biocontainers/megahit> as 1.1.3--py36_0 (sha256:8c9f17dd0fb1).

205 The code used to run IDBA is maintained by BioContainers and is available at
206 <https://quay.io/repository/biocontainers/idba> as 1.1.3--1 (sha256:51291ffeecc).

207 CDS were predicted from assembled contigs using Prokka as maintained by BioContainers
208 (<https://quay.io/repository/biocontainers/prokka>) 1.12--pl526_0 (sha256:600512072486).

209 The reference database used for the alignment-based analysis was UniRef90
210 (www.uniprot.org/uniref/) [16], downloaded on January 30th, 2018.

211

212 **Simulation of microbial communities**

213 Synthetic microbial communities were simulated using ART

214 (<https://quay.io/repository/biocontainers/art>) 2016.06.05--h869255c_2 (sha256:1cd93ed9f680)

215 with paired-end reads, a read length of 250, mean fragment length of 1000, and fragment size
216 standard deviation of 300. The abundance of each member of a given community was
217 simulated from a log-normal distribution with a mean of 5x, standard deviation of 1-log, and
218 maximum of 100x. Each community contains 20 distinct genomes.

219

220 **Results**

221 **Sensitivity and specificity of metagenomics approaches**

222 For each synthetic community, we cataloged the CDS present and compared these true
223 positives to the reported CDS by each analytic method. For mapping-based methods, we
224 allowed for duplicate calls (i.e. similar but distinct CDS sequences determined by the method to
225 be roughly equally likely to be present). Comparing these CDS catalogs (true and inferred) we
226 were able to calculate a positive predictive value (PPV; true positive / true positive + false
227 positive), sensitivity (true positive / true positive + false negative), and uniqueness (true positive
228 / true positive + duplicates). As shown in Figure 1, mapping-based approaches were more
229 sensitive, particularly when the CDS has low coverage depth, at a cost of PPV and uniqueness.

230

231 The mapping all-hits approach is the simplest approach, accepting as present any CDS that had
232 at least one aligning short-read sequence. While very sensitive, this approach had dismal PPV
233 and uniqueness. A related mapping method is to restrict to CDS with at least one short read that
234 maps uniquely to that CDS: Mapping - unique hits; this approach yielded balanced sensitivity
235 and PPV. FAMLI uses an expectation maximization-based iterative approach (considering
236 evenness of coverage and total coverage depth) and achieves somewhat superior sensitivity
237 and PPV as compared to the Mapping - unique hit approach.

238

239 HUMAAaN2 uses a hybrid approach, combining taxonomic identification, mapping of reads to
240 reference genomes, and then using a mapping - all-hits like approach for the remainder of short

241 reads that do not map to a genome. Our experimental set-up biases in favor of organisms with
242 reference genomes. In this favorable set of circumstances, HUMAaN2 performs well with
243 regards to PPV (superior to any of the tested mapping based approaches), sensitivity (similar at
244 all depths and low-coverage depths, slightly inferior to mapping approaches) and with
245 uniqueness.

246

247 Assembly based approaches have the advantage of near perfect uniqueness (with the
248 assembly process itself resulting in convergence on a single CDS), and the best PPV.

249 Sensitivity was inferior to mapping-based approaches, and varied by the coverage depth for a
250 given CDS (Figure 2).

251

252 **Short reads align equally well to multiple CDS**

253 To better understand why mapping approaches, particularly mapping with acceptance of all hits,
254 has poor sensitivity, we explored the role multi-mapping reads may be playing. To do so, three
255 random unique CDS were selected and 120 simulated reads were generated for each CDS,
256 resulting in a total of 360 simulated reads. These simulated reads were aligned against the
257 UniRef100 database. Each read has only one true origin CDS.

258

259 To account for sequencing errors and poor representation in the reference database, we
260 accepted alignments within a certain percentage of the best alignment for a given read. When
261 we accept all CDS with an alignment within 10% identity ('top-10') of the best alignment for a
262 read, 100,468 CDS are recruited for the 360 reads, an average of 279 (median of 268, minimum
263 of 77 and maximum of 537) CDS recruited from UniRef100 per read (figure XXX D, Start).

264

265 When taking a more restricted approach, only recruiting CDS with an alignment to a read
266 equivalent to the best hit, a total of 57,983 CDS are recruited, an average of 161 (median of

267 165, minimum of 1 and maximum of 384) equally well aligning reference CDS for each
268 simulated read.

269

270 **The FAMLI approach can successfully cull multi-mapped reads**

271 To establish the extent of the multimapping read problem, three random CDS were selected
272 from UniRef100. One hundred and twenty simulated reads were generated from each CDS, and
273 combined into one set of 360 paired reads; each of these reads had one true origin coding
274 sequence.

275

276 We then used Diamond to align these 360 reads against UniRef100. Even after limiting to only
277 alignments within equal in quality to the best hit, there were an average of 161 (median 165, min
278 1, and max 384) reference sequences tied with the best hit per read pair; when limited to
279 alignments within 10% of the best identity, there was a mean of 279, median 268, minimum 77,
280 and max 537 aligning subjects (references) per read pair.

281

282 To filter these alignments, we developed an iterative expectation maximization-based approach
283 that considered both the evenness of coverage and total depth of coverage (weighted by
284 alignment quality) of a candidate CDS in order to cull the vast excess of recruited CDS by the
285 mapping approach, the FAMLI algorithm. Figure 3 shows the FAMLI algorithm applied to the
286 top-10 alignments. Figure 3A shows the coverage (or read depth by base pair) for the three true
287 positive CDS. After filtering for coverage evenness, Fig 3B shows the read-depth of some
288 successfully filtered away references, as well as some references not present in the simulated
289 sample that pass this evenness test. Figure 3C depicts the iterative pruning of alignments by
290 likelihood, showing the candidate references for one query being successfully filtered down over
291 ten iterations to a single reference CDS for the read (the true origin reference for this read).

292

293 By the conclusion of the first evenness filtering, 908 references remain (for the true three); the
294 360 reads remain with an average of 271 (median of 267, minimum of 77, and maximum of 398)
295 equally-well aligning reference CDS. By the conclusion of ten iterations, all reads are
296 successfully assigned now to their true origin CDS (one reference CDS per read) (Figure 3D).

297

298 **Discussion**

299 Randomly fragmented shotgun sequencing of the metagenome of a microbial community offers
300 the promise of inferring the functional capacity of the community by establishing the protein
301 coding sequencing (CDS) present. CDS or gene-level metagenomics offers a more reproducible
302 and mechanistic means of associating the state of the microbiome with functional outcomes in a
303 host or environment [22]. Realizing this promise is predicated on having a reliable set of analytic
304 tools for determining the CDS catalog of a microbial community.

305

306 Here we introduce and employ an approach for benchmarking the performance of different
307 metagenome analysis tools for determining the CDS content of the metagenome. This
308 benchmarking approach is implemented within a reproducible Nextflow workflow, and therefore
309 should be relatively straightforward for other researchers to reproduce and augment as
310 additional tools for CDS detection become available.

311

312 We found that assembly-based tools are limited by sensitivity, particularly at low read coverage.
313 The association between the sensitivity to detect a CDS and the read coverage depth of the
314 CDS is worrisome; the ability of these tools to detect a protein coding sequence is dependent
315 upon community factors, including the relative abundance of the hosting organism, more so
316 than other approaches.

317

318 Mapping-based approaches must address the problem of short-reads from metagenomes
319 aligning equally well to large numbers of distinct CDS sequences. As evident in our simulated
320 communities, the ratio of true to false positive alignments can be in the hundreds to one,
321 resulting in dismal precision unless the alignments are culled or filtered. We suspect some of
322 the limitations experienced by software attempting to use short reads to identify the functional
323 genes encoded by microbial communities, described by [23], may be due to this multi-mapping
324 read problem.

325

326 Here we demonstrate the magnitude of the problem of multiple-mapping of short reads to
327 peptides, revealing a large number of equally-scored alignments; if one simply includes all
328 peptides for which there is at least one short read that aligns equally as well as to any other
329 peptide, the false positives outnumber true positives by an average of about 160:1.

330

331 We describe an algorithmic approach to correctly assign these multiply aligned WGS reads to
332 the proper reference sequence, implemented as the open source software package FAML I
333 (Functional Analysis of Metagenomes by Likelihood Inference). With FAML I, we are able to
334 improve our precision (number of true positives divided by the sum of false and true positives) to
335 about 80%; this performance is consistent over a range of community types. FAML I is more
336 efficient than *de novo* assembly at identifying protein-coding sequences present in a community
337 with regards to both read depth and computational time. While FAML I can be used as a
338 standalone tool to identify protein-coding genes, it could also easily be used to enhance the
339 precision of existing bioinformatics tools (e.g. HUMAnN2).

340

341 The hybrid approach of establishing which taxa are present and first mapping to reference
342 genomes (e.g. HUMAnN2, MIDAS) has merit, and performed well from a sensitivity and positive
343 predictive value perspective in our benchmarking approach. We note that our approach limits

344 our synthetic communities to being those with reference genomes. This biases in favor of this
345 hybrid approach. In the context of microbial communities with a high degree of novelty, we
346 suspect performance would be poorer.

347

348 Thinking about the relative merits of reference-based (e.g. UniRef90) or reference-free (e.g. *de*
349 *novo* assembly) analysis methods, one of our primary considerations was the efficiency of
350 comparing results across large numbers of samples. While reference-free approaches are free
351 by definition from the bias inherent in reference databases, that lack of common reference
352 makes it extremely challenging to compare results between samples. For example, comparing a
353 set of genes between N samples is an $O(N^2)$ problem that scales exponentially with the
354 number of samples. In contrast, by identifying proteins from a reference database (UniRef90),
355 all results are inherently comparable without any additional computation (e.g. sequence
356 alignment), in other words the complexity is $O(1)$. Put simply, with *de novo* assembly (SPAdes)
357 it is *much* more difficult to compare the results for 1,000 samples in contrast to just 10 samples,
358 while for FAMLI or HUMAnN2 it is about the same.

359

360 **Acknowledgements**

361 We would like to acknowledge funding support the Microbiome Research Initiative at the Fred
362 Hutch, lead by David N Fredricks, for supporting this work and Dan Tenenbaum from the
363 scientific computing group at the Fred Hutch for help with establishing computational resources
364 for this manuscript.

365

366 **Conflicts of Interest**

367 The authors have no conflicts of interest to disclose.

368

369 **References**

- 370 1. NIH HMP Working Group, Peterson J, Garges S, Giovanni M, McInnes P, Wang L, et al. The
371 NIH Human Microbiome Project. *Genome Res.* 2009;19:2317–23.
- 372 2. Qin J, Li R, Raes J, Arumugam M, Burgdorf KS, Manichanh C, et al. A human gut microbial
373 gene catalogue established by metagenomic sequencing. *Nature.* 2010;464:59–65.
- 374 3. Cho I, Blaser MJ. The human microbiome: at the interface of health and disease. *Nat Rev*
375 *Genet.* 2012;13:260–70.
- 376 4. Human Microbiome Project Consortium. Structure, function and diversity of the healthy
377 human microbiome. *Nature.* 2012;486:207–14.
- 378 5. Human Microbiome Project Consortium. A framework for human microbiome research.
379 *Nature.* 2012;486:215–21.
- 380 6. Golob JL, Margolis E, Hoffman NG, Fredricks DN. Evaluating the accuracy of amplicon-based
381 microbiome computational pipelines on simulated human gut microbial communities. *BMC*
382 *Bioinformatics.* 2017;18:283.
- 383 7. Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, et al. Introducing
384 mothur: open-source, platform-independent, community-supported software for describing and
385 comparing microbial communities. *Appl Environ Microbiol.* 2009;75:7537–41.
- 386 8. Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman FD, Costello EK, et al. QIIME
387 allows analysis of high-throughput community sequencing data. *Nat Methods.* 2010;7:335–6.
- 388 9. Quince C, Walker AW, Simpson JT, Loman NJ, Segata N. Shotgun metagenomics, from
389 sampling to analysis. *Nat Biotechnol.* 2017;35:833–44.
- 390 10. Scholz MB, Lo C-C, Chain PSG. Next generation sequencing and bioinformatic bottlenecks:
391 the current state of metagenomic data analysis. *Curr Opin Biotechnol.* 2012;23:9–15.
- 392 11. Abubucker S, Segata N, Goll J, Schubert AM, Izard J, Cantarel BL, et al. Metabolic
393 reconstruction for metagenomic data and its application to the human microbiome. *PLoS*
394 *Comput Biol.* 2012;8:e1002358.
- 395 12. Nayfach S, Rodriguez-Mueller B, Garud N, Pollard KS. An integrated metagenomics

- 396 pipeline for strain profiling reveals novel patterns of bacterial transmission and biogeography.
397 *Genome Res.* 2016;26:1612–25.
- 398 13. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, et al. SPAdes: A New
399 Genome Assembly Algorithm and Its Applications to Single-Cell Sequencing. *Journal of*
400 *Computational Biology.* 2012;19:455–77.
- 401 14. Peng Y, Leung HCM, Yiu SM, Chin FYL. IDBA-UD: a de novo assembler for single-cell and
402 metagenomic sequencing data with highly uneven depth. *Bioinformatics.* 2012;28:1420–8.
- 403 15. Fitch WM. Homology a personal view on some of the problems. *Trends Genet.*
404 2000;16:227–31.
- 405 16. McIntyre ABR, Ounit R, Afshinnekoo E, Prill RJ, Hénaff E, Alexander N, et al.
406 Comprehensive benchmarking and ensemble approaches for metagenomic classifiers. *Genome*
407 *Biol.* 2017;18:182.
- 408 17. Sczyrba A, Hofmann P, Belmann P, Koslicki D, Janssen S, Dröge J, et al. Critical
409 Assessment of Metagenome Interpretation—a benchmark of metagenomics software. *Nat*
410 *Methods.* 2017;14:1063–71.
- 411 18. Lindgreen S, Adair KL, Gardner PP. An evaluation of the accuracy and speed of
412 metagenome analysis tools. *Sci Rep.* 2016;6:19233.
- 413 19. Petersen TN, Lukjancenko O, Thomsen MCF, Maddalena Sperotto M, Lund O, Møller
414 Aarestrup F, et al. MGmapper: Reference based mapping and taxonomy annotation of
415 metagenomics sequence reads. *PLoS ONE.* 2017;12:e0176469.
- 416 20. Di Tommaso P, Chatzou M, Floden EW, Barja PP, Palumbo E, Notredame C. Nextflow
417 enables reproducible computational workflows. *Nat Biotechnol.* 2017;35:316–9.
- 418 21. Suzek BE, Wang Y, Huang H, McGarvey PB, Wu CH, the UniProt Consortium. UniRef
419 clusters: a comprehensive and scalable alternative for improving sequence similarity searches.
420 *Bioinformatics.* 2015;31:926–32.
- 421 22. Minot SS, Willis AD. Clustering co-abundant genes identifies components of the gut

422 microbiome that are reproducibly associated with colorectal cancer and inflammatory bowel
423 disease. bioRxiv [Internet]. 2019 [cited 2019 Mar 21]; Available from:
424 <http://biorxiv.org/lookup/doi/10.1101/567818>

425 23. Carr R, Borenstein E. Comparative analysis of functional metagenomic annotation and the
426 mappability of short reads. PLoS ONE. 2014;9:e105776.

427

428 **Figure Legends**

429

430 **Figure 1: Positive predictive value (PPV), sensitivity, and uniqueness of CDS calls by**
431 **metagenomic analysis approaches.** The positive predictive value (true positive over true
432 positive plus false positive), sensitivity (true positive over true positive plus false negative) both
433 overall and subsetted to CDS with 0-5x coverage, and uniqueness (true positive over true
434 positive plus duplicates) on a per-CDS basis with different analysis approaches.

435

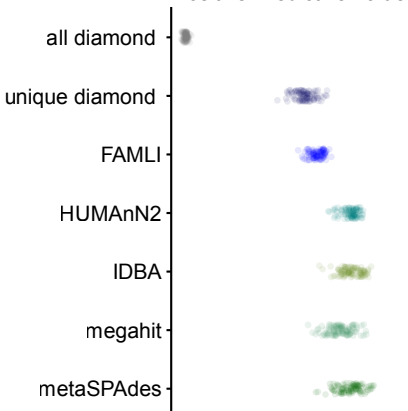
436 **Figure 2: Sensitivity and uniqueness of CDS calls with respect to CDS coverage depth.**
437 Mapping based approaches are both more sensitive, and achieve a plateau of sensitivity at a
438 lower coverage depth as compared to assembly-based methods.

439

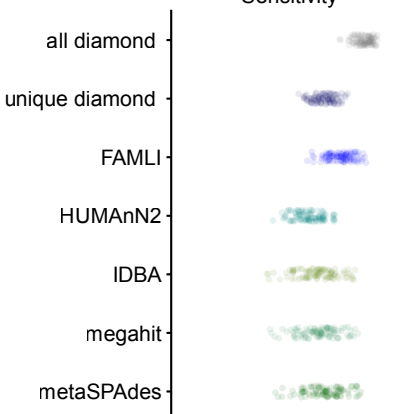
440 **Figure 3: The problem of multiply-mapping short-reads, and the FAML algorithm**
441 **schematized.** Three hundred and sixty simulated reads were generated from three CDS. These
442 simulated read was aligned against the UniRef100 database, and all CDS with an alignment
443 within 10% identity of the best match were retained. A) The read-depth coverage of the three
444 true peptides (top) B) Evenness filtering is used to remove the least likely to be present
445 references from being considered. The left column is three randomly selected references that
446 are successfully filtered at this step, the right three false references that are not filtered. C) The
447 iterative likelihood-based filtering of one randomly selected synthetic read. Each circle

448 represents one remaining aligned reference CDS for this read; the true positive origin reference
449 is in dark green. The length of each line is proportional to the calculated score at this iteration.
450 D)The number of CDS per read as a violin plot. After the tenth iteration, only one reference CDS
451 (the correct) remains for this read.

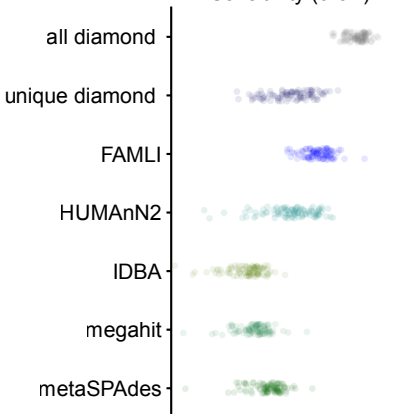
Positive Predictive Value



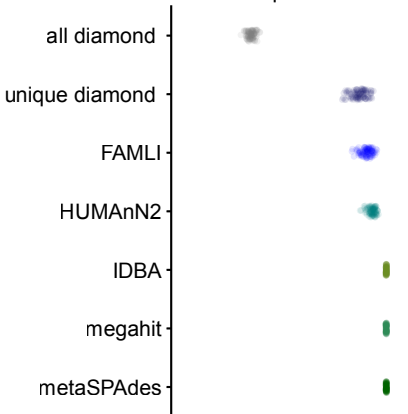
Sensitivity



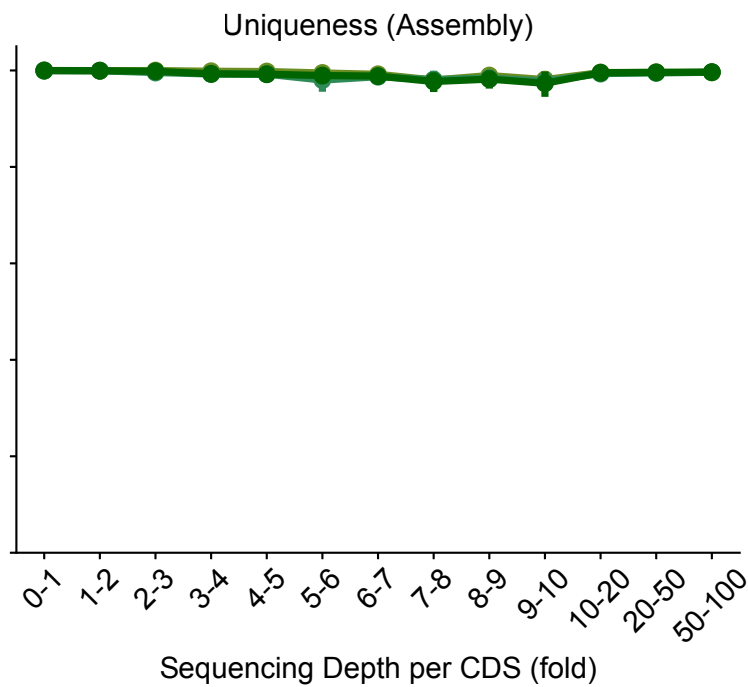
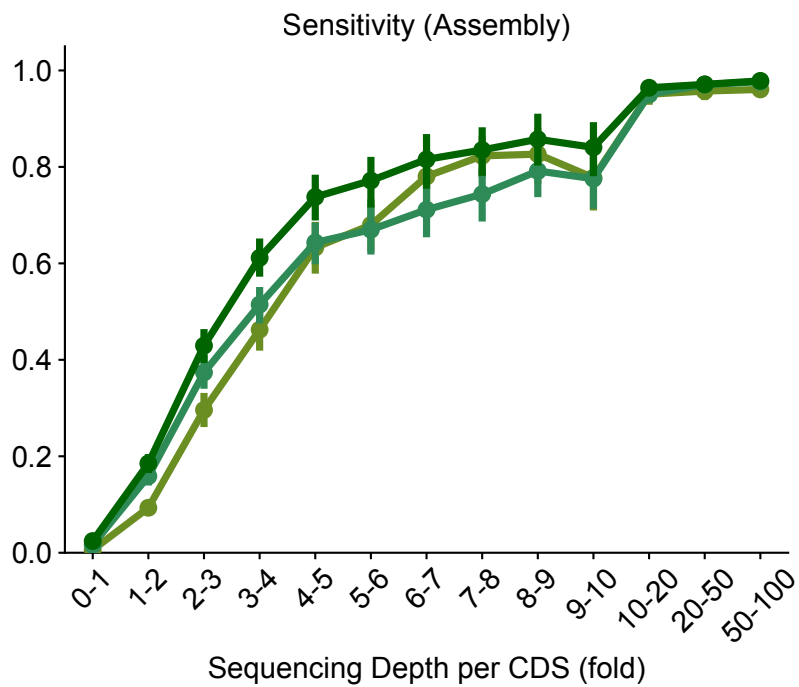
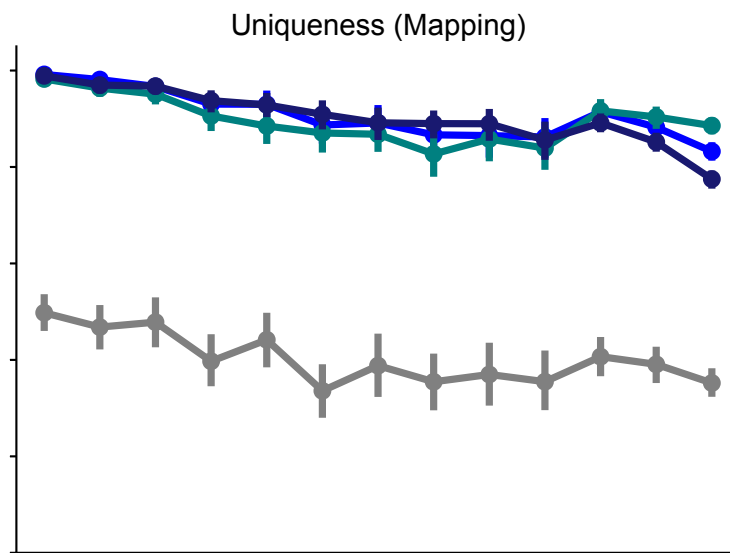
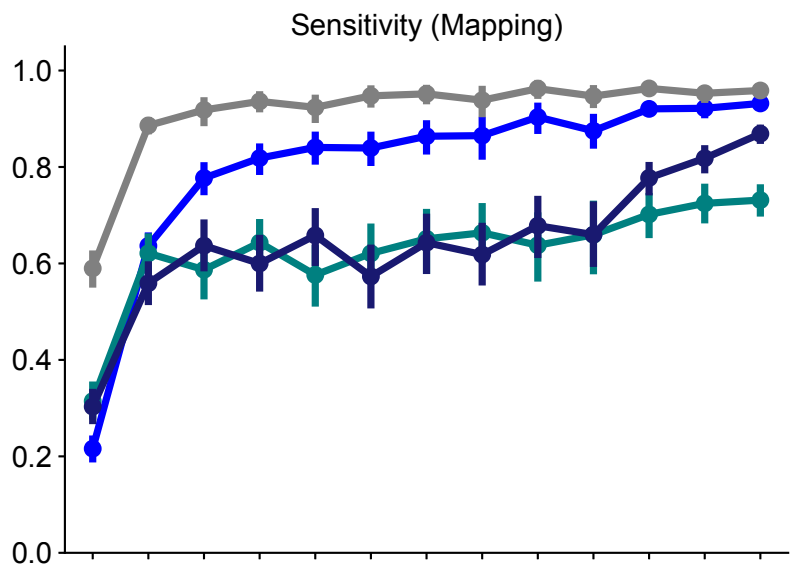
Sensitivity (0-5X)



Uniqueness

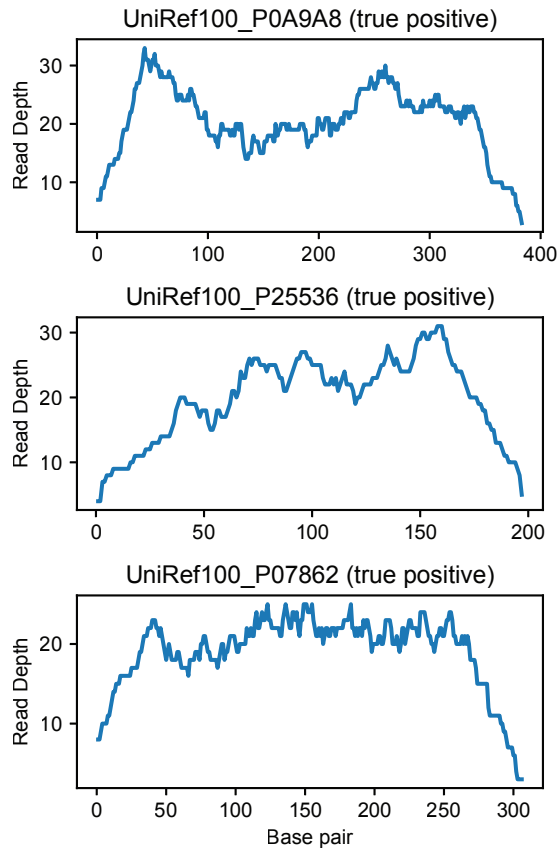


0.0 0.5 1.0

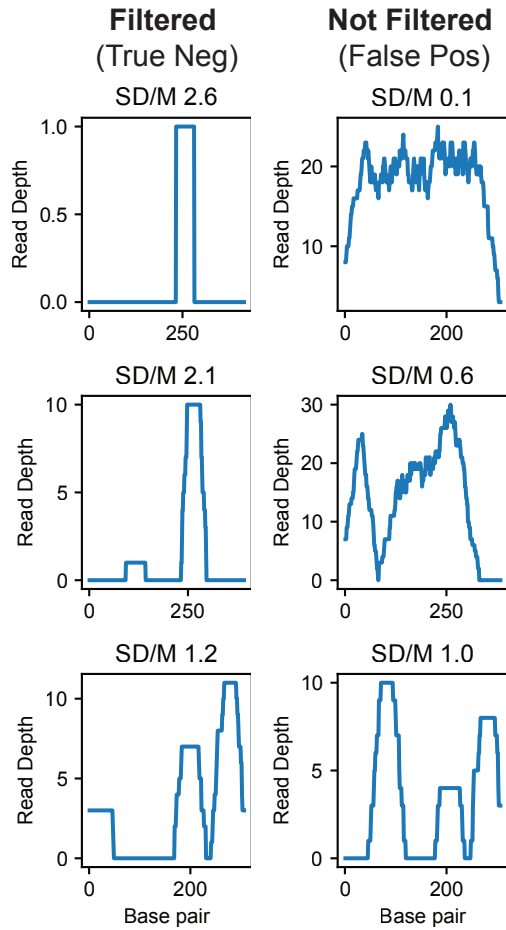


- method
- FAMLI
 - HUMAnN2
 - IDBA
 - all diamond
 - megahit
 - metaSPAdes
 - unique diamond

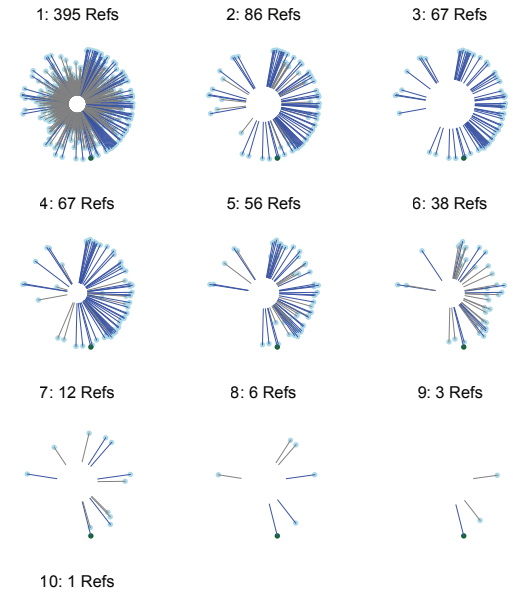
A. Start



B. Evenness Filtering



C. Iterative Alignment Filtering



D. CDS per read by step

