

1 Evaluating Metagenome Assembly on a Simple
2 Defined Community with Many Strain Variants

3 Sherine Awad¹, Luiz Irber¹, C. Titus Brown^{1*}
4 ¹**Department of Population Health and Reproduction**
 University of California, Davis
 Davis, CA 95616 USA
 * E-mail: ctbrown@ucdavis.edu

4 June 25, 2017

5 **Abstract**

6 We evaluate the performance of three metagenome assemblers, IDBA,
7 MetaSPAdes, and MEGAHIT, on short-read sequencing of a defined
8 “mock” community containing 64 genomes (Shakya et al. (2013)). We
9 update the reference metagenome for this mock community and detect
10 several additional genomes in the read data set. We show that strain
11 confusion results in significant loss in assembly of reference genomes
12 that are otherwise completely present in the read data set. In agree-
13 ment with previous studies, we find that MEGAHIT performs best
14 computationally; we also show that MEGAHIT tends to recover larger
15 portions of the strain variants than the other assemblers.

16 Introduction

17 Metagenomics refers to sequencing of DNA from a mixture of organisms,
18 often from an environmental or uncultured sample. Unlike whole genome
19 sequencing, metagenomics targets a mixture of genomes, which introduces
20 metagenome-specific challenges in analysis [1]. Most approaches to analyz-
21 ing metagenomic data rely on mapping or comparing sequencing reads to
22 reference sequence collections. However, reference databases contain only a
23 small subset of microbial diversity [2], and much of the remaining diversity
24 is evolutionarily distant and search techniques may not recover it [3].

25 As sequencing capacity increases and sequence data is generated from
26 many more environmental samples, metagenomics is increasingly using *de*
27 *novo* assembly techniques to generate new reference genomes and metagenomes
28 [4]. There are a number of metagenome assemblers that are widely used -
29 see [5] for an overview of the available software, and [1] for a review of the
30 different assembler methodologies. However, evaluating the results of these
31 assemblers is challenging due to the general lack of good quality reference
32 metagenomes.

33 Moya et al. in [6] evaluated metagenome assembly using two simulated
34 454 viral metagenome and six assemblers. The assemblies were evaluated
35 based on several metrics including N50, percentages of reads assembled, ac-
36 curacy when compared to the reference genome. In addition to, chimeras per
37 contigs and the effect of assembly on taxonomic and functional annotations.

38 Mavromatis et al. in [7] provided a benchmark study to evaluate the
39 fidelity of metagenome processing methods. The study used simulated
40 metagenomic data sets constructed at different complexity levels. The datasets
41 were assembled using Phrap v3.57, Arachne v.2 [8] and JAZZ [9]. This study
42 evaluates assembly, gene prediction, and binning methods. However, the
43 study did not evaluate the assembly quality against a reference genome.

44 Rangwala et al. in [10] presented an evaluation study of metagenome
45 assembly. The study used a de Bruijn graph based assembler ABYSS [11] to
46 assemble simulated metagenome reads of 36 bp. The data set is classified at
47 different complexity levels. The study compared the quality of the assembly
48 of the data sets in terms of contig length and assembly accuracy. The
49 study also took into consideration the effect of kmer size and the degree of
50 chimericity. However, the study evaluated the assembly based on only one
51 assembler. Also, both previous studies used simulated data, which may lack
52 confounders of assembly such as sequencing artifacts and GC bias.

53 In a landmark study, Shakya et al. (2013) constructed a synthetic com-

54 munity of organisms by mixing DNA isolated from individual cultures of
55 64 bacteria and archaea, including a variety of strains across a range of
56 nucleotide distances [12]. In addition to performing 16s amplicon analy-
57 sis and doing 454 sequencing, the authors shotgun-sequenced the mixture
58 with Illumina. While the authors concluded that this metagenomic sequenc-
59 ing generally outperformed amplicon sequencing, they did not conduct an
60 assembly based analysis. This data set was also used in several other eval-
61 uation studies, including gbttools for binning [13] and benchmarking of the
62 MEGAHIT assembler [14].

63 More recently, several benchmark studies systematically evaluated metagenome
64 assembly of short reads. The Critical Assessment of Metagenome Interpre-
65 tation (CAMI) collaboration benchmarked a number of metagenome assem-
66 blers on several data sets of varying complexity, evaluating recovery of novel
67 genomes and multiple strain variants [3]. Notably, CAMI concluded that
68 “The resolution of strain-level diversity represents a substantial challenge
69 to all evaluated programs.” Another recent study evaluated eight assem-
70 blers on nine environmental metagenomes and three simulated data sets
71 and provided a workflow for choosing a metagenome assembler based on
72 the biological goal and computational resources available [15]. [5] explored
73 metagenome assembler performance on a pair of real data sets, again con-
74 cluding that the biological goal and computational resources defined the
75 choice of assembler. Also see [16] for an analysis of a previously generated
76 HMP benchmark data set; however, the Illumina reads used for this study
77 are much shorter than current sequencing and are arguably not relevant for
78 future studies.

79 In this study, we extend previous work by delving into questions of
80 chimeric misassembly and strain recovery in the Shakya et al. (2013) data
81 set. First, we update the list of reference genomes for Shakya et al. to in-
82 clude the latest GenBank assemblies along with plasmids. We then compare
83 IDBA [17], MetaSPAdes [18], and MEGAHIT [19] performance on assem-
84 bling this short-read data set, and explore concordance in recovery between
85 the three assemblers. We describe the effects of “strain confusion” between
86 multiple strains. We also detect and analyze several previously unreported
87 strains and genomes in the Shakya et al. data set. We find that in the ab-
88 sence of closely related genomes, all three metagenome assemblers recover
89 95% or more of known reference genomes. However, in the presence of
90 closely related genomes, these three metagenome assemblers vary widely in
91 their performance and, in extreme cases, can fail to recover the majority of
92 some genomes even when they are completely present in the reads. Our re-

93 port provides strong guidance on choice of assemblers and extends previous
94 analyses of this low-complexity metagenome benchmarking data set.

95 **Datasets**

96 We used a diverse mock community data set constructed by pooling DNA
97 from 64 species of bacteria and archaea and sequencing them with Illumina
98 HiSeq. The raw data set consisted of 109,629,496 reads from Illumina HiSeq
99 101 bp paired-end sequencing (2x101) with an untrimmed total length of
100 11.07 Gbp and an estimated fragment size of 380 bp [12].

101 The original reads are available through the NCBI Sequence Read Archive
102 at Accession SRX200676. We updated the 64 reference genomes sets from
103 NCBI GenBank using the latest available assemblies with plasmid content
104 (June 2017); the accession numbers are available as `accession-list-ref.txt`
105 in the Zenodo repository, DOI: 10.5281/zenodo.818050. For convenience, the
106 updated reference genome collection is available for download at the archival
107 URL <https://osf.io/vbhy5/>.

108 **Methods**

109 The analysis code and run scripts for this paper are written in Python and
110 bash, and are available at [https://github.com/dib-lab/2016-metagenome-](https://github.com/dib-lab/2016-metagenome-assembly-eval/)
111 `assembly-eval/` (archived at Zenodo DOI: @DOI: 10.5281/zenodo.818050).
112 The scripts and overall pipeline were examined by the first and senior au-
113 thors for correctness. In addition, the bespoke reference-based analysis
114 scripts were tested by running them on a single-colony *E. coli* MG1655 data
115 set with a high quality reference genome [20].

116 **Quality Filtering**

117 We removed adapters with Trimmomatic v0.30 in paired-end mode with
118 the TruSeq adapters [21], using light quality score trimming (`LEADING:2`
119 `TRAILING:2 SLIDINGWINDOW:4:2 MINLEN:25`) as recommended in MacManes,
120 2014 [22].

121 **Reference Coverage Profile**

122 To evaluate how much of the reference metagenome was contained in the
123 read data, we used `bwa aln` (v0.7.7.r441) to map reads to the reference

124 genome [23]. We then calculated how many reference bases were covered by
125 mapped reads (custom script `coverage-profile.py`).

126 Measuring k-mer inclusion and Jaccard similarity

127 We used MinHashing as implemented in sourmash to estimate k-mer inclu-
128 sion and Jaccard similarity between data sets [24]. MinHash signatures were
129 prepared with `sourmash compute` using `--scaled 10000`. K-mer inclusion
130 was computed by taking the ratio of the number of intersecting hashes with
131 the query over the total number of hashes in the subject MinHash. Jac-
132 card similarity was computed as in [25] by taking the ratio of the number
133 of intersecting hashes between the query and subject over the number of
134 hashes in the union. K-mer sizes for comparison were chosen at 21, 31, or
135 51, depending on the level of taxonomic specificity desired - genus, species,
136 or strain, respectively, as described in [26].

137 Where specified, high-abundance k-mers were selected for counting by
138 using the script `trim-low-abund.py` script with `-C 5` from khmer v2 [27,
139 28].

140 Assemblers

141 We assembled the quality-filtered reads using three different assemblers:
142 IDBA-UD [17], MetaSPAdes [18], and MEGAHIT [19]. For IDBA-UD v1.1.1
143 [17], we used `--pre_correction` to perform pre-correction before assembly
144 and `-r` for the pe files.

145 For MetaSPAdes v3.9.0 [18], we used `--meta --pe1-12 --pe1-s` where
146 `--meta` is used for metagenomic data sets, `--pe1-12` specifies the interlaced
147 reads for the first paired-end library, and `--pe1-s` provides the orphan reads
148 remaining from quality trimming.

149 For MEGAHIT v1.1.1-2-g02102e1 [19], we used `-l 101 -m 3e9 --cpu-only`
150 where `-l` is for maximum read length, `-m` is for max memory in bytes to
151 be used in constructing the graph, and `--cpu-only` to use only the CPU
152 and no GPUs. We also used `--presets meta-large` for large and complex
153 metagenomes, and `--12` and `-r` to specify the interleaved-paired-end and
154 single-end files respectively. MEGAHIT allows the specification of a memory
155 limit and we used `-M 1e+10` for 10 GB.

156 All three assemblies were executed on the same high-memory buy-in
157 node on the Michigan State University High Performance Compute Cluster,
158 and we recorded RAM and CPU time of each assembly job using the `qstat`

159 utility at the end of each run.

160 Unless otherwise mentioned, we eliminated all contigs less than 500 bp
161 from each assembly prior to further analysis.

162 Mapping

163 We aligned all quality-filtered reads to the reference metagenome with `bwa`
164 `aln` (v0.7.7.r441) [23]. We aligned paired-end and orphaned reads separately.
165 We then used `samtools` (v0.1.19) [29] to convert SAM files to BAM files for
166 both paired-end and orphaned reads. To count the unaligned reads, we
167 included only those records with the “4” flag in the SAM files [29].

168 Assembly analysis using NUCmer

169 We used the NUCmer tool from MUMmer3.23 [30] to align assemblies to the
170 reference genome with options `-coords -p`. Then we parsed the generated
171 “.coords” file using a custom script `analyze_assembly.py`, and calculated
172 several analysis metrics across all three assemblies at a 99% alignment iden-
173 tity.

174 Reference-based analysis of the assemblies

175 We conducted reference-based analysis of the assemblies under two condi-
176 tions. “Loose” alignment conditions used all available alignments, including
177 redundant and overlapping alignments. “Strict” alignment conditions took
178 only the longest alignment for any given contig, eliminating all other align-
179 ments.

180 The script `summarize-coords2.py` was used to calculate aligned cov-
181 erage from the loose alignment conditions: each base in the reference was
182 marked as “covered” if it was included in at least one alignment. The script
183 `analyze_ng50.py` was used to calculate NGA 50 for each individual refer-
184 ence genome.

185 Analysis of chimeric misassemblies

186 We analyzed each assembly for chimeric misassemblies by counting the num-
187 ber of contigs that contained matches to two distinct reference genomes. In
188 order to remove secondary alignments from consideration, we included only

189 the longest non-overlapping NUCmer alignments for each contig at a mini-
190 mum alignment identity of 99%. We then used the script `analyze_chimeric2.py`
191 to find individual contigs that matched more than one distinct reference
192 genome. As a negative control on our analysis, we verified that this ap-
193 proach yielded no positive results when applied to the alignments of the
194 reference metagenome against itself.

195 Analysis of unmapped reads

196 We conducted assembly and analysis of unmapped reads with MEGAHIT,
197 NUCmer, and sourmash as above. The new GenBank genomes are listed in
198 the Zenodo archive at the file `accession-list-unmapped.txt` and for con-
199 venience are available for download at the archival URL <https://osf.io/34ef8/>.

200 Results

201 The raw data is high quality.

202 The reads contains 11,072,579,096 bp (11.07 Gbp) in 109,629,496 reads with
203 101.0 average length (2x101bp Illumina HiSeq).

204 Trimming removed 686,735 reads (0.63%). After trimming, we retained
205 108,422,358 paired reads containing 10.94 Gbp with an average length of
206 100.9 bases. A total of 46.56 Mbp remained in 520,403 orphan reads with
207 an average length of 89.5 bases. In total, the quality trimmed data contained
208 10.98 Gbp in 108,942,761 reads. This quality trimmed (“QC”) data set was
209 used as the basis for all further analyses.

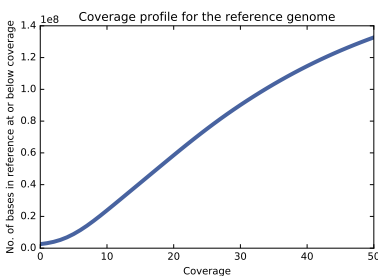


Figure 1: Cumulative coverage profile for the reference metagenome, based on read mapping.

Table 1: Jaccard containment of the reference in the reads

k-mer size	% reference in reads
21	96.8%
31	95.9%
41	94.9%
51	94.1%

210 **The reference metagenome is not completely present in the**
211 **reads.**

212 We next evaluated the fraction of the reference genome covered by at least
213 one read (see Methods for details). Quality filtered reads cover 203,058,414
214 (98.76%) bases of the reference metagenome (205,603,715 bp total size). Fig-
215 ure 1 shows the cumulative coverage profile of the reference metagenome,
216 and the percentage of bases with that coverage. Most of the reference
217 metagenome was covered at least minimally; only 3.33% of the reference
218 metagenome had mapping coverage <5, and 1.24% of the bases in the ref-
219 erence were not covered by any reads in the QC data set.

220 In order to evaluate reconstructability with De Bruijn graph assemblers,
221 we next examined k-mer containment of the reference in the reads for k of
222 21, 31, 41, and 51 (Table 1). The k-mer overlap decreases from 96.8% to
223 94.1% as the k-mer size increases. This could be caused by low coverage of
224 some portions of the reference and/or variation between the reads and the
225 reference.

226 **Some individual reference genomes are poorly represented in**
227 **the reads.**

Table 2: Top uncovered genomes

Genome	Read coverage
<i>Desulfovibrio vulgaris</i> DP4	93.2%
<i>Thermus thermophilus</i> HB27	91.1%
<i>Enterococcus faecalis</i> V583	74.6%
<i>Fusobacterium nucleatum</i>	47.6%

228 To see if specific reference genomes exhibited low coverage, we analyzed
229 read mapping coverage for individual genomes. Of the 64 reference genomes

Table 3: Genomes removed from reference for low 51-mer presence

51-mers in reads	Genome
98.7	<i>Leptothrix cholodnii</i>
98.7	<i>Haloferax volcanii</i> DS2
98.6	<i>Salinispora tropica</i> CNB-440
97.4	<i>Deinococcus radiodurans</i>
97.2	<i>Zymomonas mobilis</i>
97.1	<i>Ruegeria pomeroyi</i>
96.8	<i>Shewanella baltica</i> OS223
95.5	<i>B. bronchiseptica</i> D989
94.5	<i>Burkholderia xenovorans</i>
72.0	<i>Desulfovibrio vulgaris</i> DP4
65.0	<i>Thermus thermophilus</i> HB27
53.4	<i>Enterococcus faecalis</i>
4.7	<i>Fusobacterium nucleatum</i> ATCC 25586

230 used in the metagenome, 60 had a per-base mapping coverage above 95%.
231 The remaining four varied significantly (Table 2), with *F. nucleatum* the
232 lowest – only 47.6% of the bases in the reference genome are covered by one
233 or more mapped reads.

234 We next did a 51-mer containment analysis of each reference genome in
235 the reads; $k=51$ was chosen so as to be specific to strain content [26]. 99%
236 or more of the constituent 51-mers for 51 of the 64 reference genomes were
237 present in the reads, suggesting that each of the 51 genomes was entirely
238 present at some minimal coverage.

239 We excluded the remaining 13 genomes (see Table 3) from any fur-
240 ther reference-based analysis because interpreting recovery and misassembly
241 statistics for these genomes would be confounding; also see the discussion of
242 strain variants, below.

243 **MEGAHIT is the fastest and lowest-memory assembler eval-** 244 **uated**

245 We ran three commonly used metagenome assemblers on the QC data set:
246 IDBA-UD, MetaSPAdes, and MEGAHIT. We recorded the time and mem-
247 ory usage of each (Table 4). In computational requirements, MEGAHIT
248 outperformed both MetaSPAdes and IDBA-UD considerably, producing an
249 assembly in four hours (“wall time”) – approximately 12 times faster than

Table 4: Running Time and Memory Utilization

Assembler	CPU time	Wall time	RAM
MEGAHIT	52hr 25m	4 hr 9m	11.4 GB
IDBA-UD	49h	49h	39.8GB
MetaSPAdes	94hr 43m	94hr 44m	100.7 GB

250 IDBA and 23 times faster than MetaSPAdes. MEGAHIT used only 11.4
251 GB of RAM – 1/3rd to 1/9th the memory used by IDBA and MetaSPAdes,
252 respectively.

253 CPU time measurements (which include processing on multiple CPU
254 cores) show that MEGAHIT and IDBA are competitive in overall process-
255 ing time, but MEGAHIT’s ability to make use of multiple cores results in
256 significantly less overall assembly time; this is particularly relevant given
257 the increasing availability of manycore processors. Despite a variety of con-
258 figuration attempts, we were unable to get MetaSPAdes to use threading
259 effectively; however, we note that even with perfectly parallel processing
260 on 16 cores, MetaSPAdes would take 6 hours and still use approximately 9
261 times as much RAM as MEGAHIT.

262 The assemblies contain most of the raw data

Table 5: Read and high-abundance (> 5) k-mer exclusion from assemblies

Assembly	Unmapped Reads	51-mers omitted
IDBA	3,328,674 (3.05%)	2.4%
MetaSPAdes	3,844,123 (3.52%)	3.2%
MEGAHIT	2,737,640 (2.51%)	2.8%

263 We assessed read inclusion in assemblies by mapping the QC reads to
264 the length-filtered assemblies and counting the remaining unmapped reads.
265 Depending on the assembly, between 2.7 million and 3.9 million reads (2.5-
266 3.5%) did not map to the assemblies (Table 5). All of the assemblies included
267 the large majority of high-abundance 51-mers (more than 96.8% in all cases).

268 Much of the reference is covered by the assemblies.

269 We next evaluated the extent to which the assembled contigs recovered the
270 “known/true” metagenome sequence by aligning each assembly to the ad-

Table 6: Contig coverage of reference with loose alignment conditions.

Assembly	bases aligned	duplication	51-mers
MEGAHIT	94.8%	1.0%	96.7%
MetaSPAdes	93.1%	1.1%	96.2%
IDBA	93.6%	0.98%	97.2%

271 justed reference (Table 6). Each of the three assemblers generates contigs
272 that cover more than 93.1% of the reference metagenome at high identity
273 (99%) with little duplication (approximately 1%). All three assemblies con-
274 tain between 96.2% and 97.2% of the 51-mers in the reference.

275 At 99% identity with the loose mapping approach, approximately 2.5% of
276 the reference is missed by all three assemblers, while 1.7% is uniquely covered
277 by MEGAHIT, 0.74% is uniquely covered by MetaSPAdes, and 0.64% is
278 uniquely covered by IDBA.

279 **The generated contigs are broadly accurate.**

Table 7: Contig accuracy measured by reference coverage with strict alignment.

Assembly	% covered
MEGAHIT	89.3%
IDBA	87.7%
MetaSPAdes	83.4%

280 When counting only the best (longest) alignment per contig at a 99%
281 identity threshold, each of the three assemblies recovers more than 87.3% of
282 the reference, with MEGAHIT recovering the most – 89.3% of the reference
283 (Table 7).

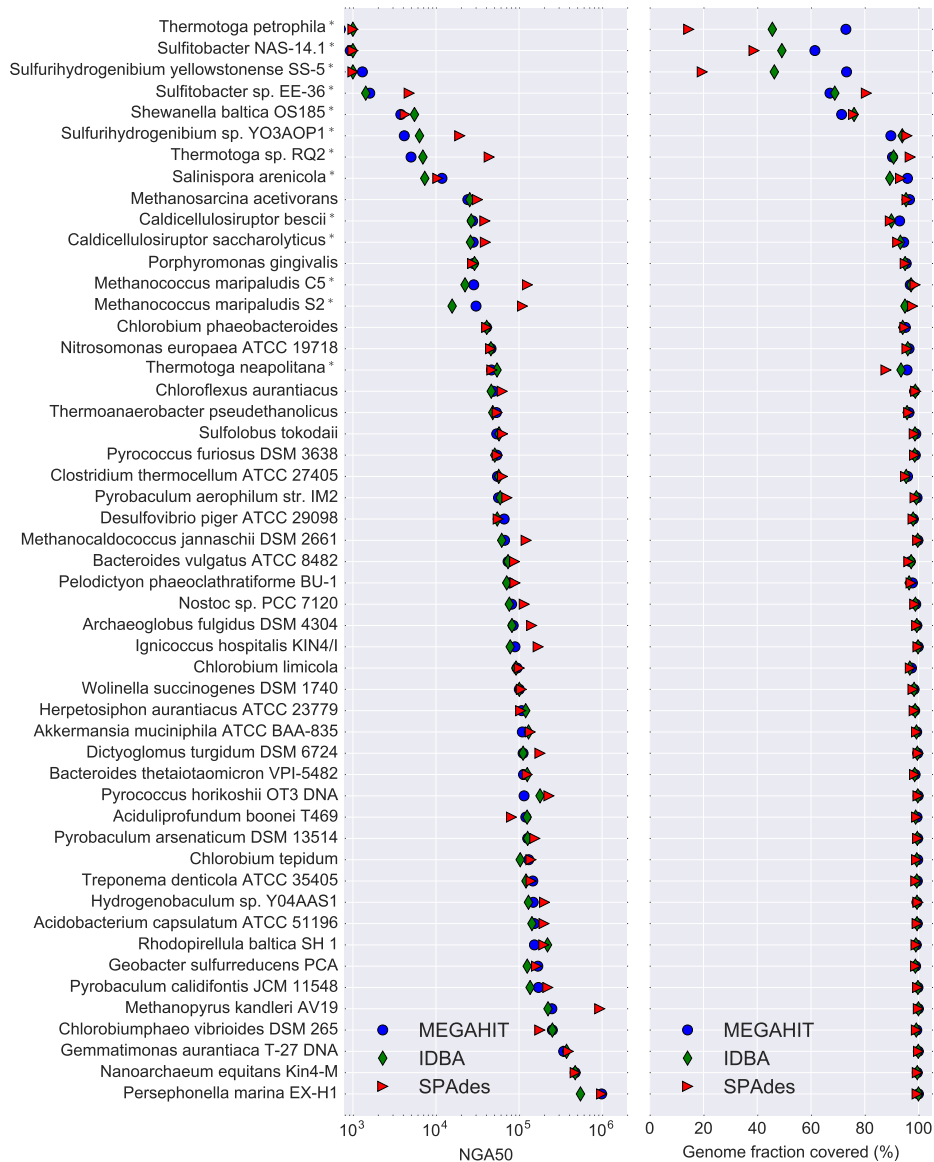


Figure 2: NGA50 and genome fraction covered, by genome and assembler. A '*' after the name indicates the presence of at least one other genome with > 2% Jaccard similarity at k=31 in the community.

284 **Individual genome statistics vary widely in the assemblies.**

285 We computed the NGA50 for each individual genome and assembly in order
286 to compare assembler performance on genome recovery (see left panel of Fig-
287 ure 2). The NGA50 statistics for individual genomes vary widely, but there
288 are consistent assembler-specific trends: IDBA yields the lowest NGA50 for
289 28 of the 51 genomes, while MetaSPAdes yields the highest NGA50 for 32
290 of the 51 genomes.

291 We also evaluated aligned coverage per genome for each of the three
292 assemblies (right panel, Figure 2). We found that 13 of the 51 genomes were
293 missing 5% or more of bases in at least one assembly, despite all 51 genomes
294 having 99% or higher read- and 51-mer coverage.

295 There are 12 genomes with k=31 Jaccard similarity greater than 2%
296 to other genomes in the community, and these (denoted by '*' after the
297 name) typically had lower NGA50 and aligned coverage numbers than other
298 genomes. In particular, these constituted 12 of the 13 genomes missing 5%
299 or more of their content, and the lowest eight NGA50 numbers.

300 **Longer contigs are less likely to be chimeric.**

Table 8: Chimeric contigs by contig length.

Assembly	> 50kb	> 5kb	> 500 bp
IDBA	0	1	7 (0.06%)
MEGAHIT	1	4	14 (0.13%)
MetaSPAdes	0	3	30 (0.48%)

301 Chimerism is the formation of contigs that include sequence from multi-
302 ple genomes. We evaluated the rate of chimerism in contigs at three different
303 contig length cutoffs: 500bp, 5kb, and 50kb (Table 8). We found that the
304 percentage of contigs that match to the genomes of two or more different
305 species drop as the minimum contig size increases, to the point where only
306 the MEGAHIT assembly had a single chimeric contig longer than 50kb.
307 Overall, chimeric misassemblies were rare, with no assembler generating
308 more than 30 chimeric contigs out of thousands of total contigs.

309 **The unmapped reads contain strain variants of reference genomes.**

310 Approximately 4.8 million reads (4.4%) from the QC data set did not map
311 anywhere in the reference provided by the authors of [12]. We extracted

Table 9: GenBank genomes detected in assembly of unmapped reads

match	GenBank genome
44.1%	<i>Fusobacterium sp.</i> OBRC1
23.0%	<i>P. ruminis strain</i> ML2
18.2%	<i>Thermus thermophilus</i> HB8
7.7%	<i>P. ruminis strain</i> CGMCC
8.2%	<i>Enterococcus faecalis</i> M7
7.3%	<i>F. nucleatum</i> 13_3C
3.7%	<i>F. nucleatum subsp. polymorphum</i>
2.9%	<i>Fusobacterium hwasookii</i>
1.0%	<i>E. coli isolate</i> YS
1.7%	<i>F. nucleatum subsp. polymorphum, alt.</i>
1.9%	<i>F. nucleatum subsp. vincentii</i>

312 and assembled these reads in isolation using MEGAHIT, yielding 6.5 Mbp
313 of assembly in 1711 contigs > 500bp in length. We then did a k-mer in-
314 clusion analysis of this assembly against all of the GenBank genomes at
315 k=31, and estimated the fraction of the k-mers that belonged to different
316 species (Table 9). We find that 51.1% of the k-mer content of these contigs
317 positively match to a genome present in GenBank but not in the reference
318 metagenome.

319 To verify these assignments, we aligned the MEGAHIT assembly of un-
320 mapped reads to the GenBank genomes in Table 9 with NUCmer using
321 “loose” alignment criteria. We found that 1.78 Mbp of the contigs aligned
322 at 99% identity or better to these GenBank genomes. We also confirmed
323 that, as expected, there are no matches in this assembly to the full updated
324 reference metagenome.

325 We note that all but the two *P. ruminis* matches and the *E. coli* isolate
326 YS are strain variants of species that are part of the defined community
327 but are not completely present in the reads (see Table 2). For *Proteini-
328 clasticum ruminis*, there is no closely related species in the mock community
329 design, and very little of the MEGAHIT assembly aligns to known *P. ru-
330 minis* genomes at 99%. However, there are many alignments to *P. ruminis*
331 at 94% or higher, for approximately 2.73 Mbp total. This suggests that the
332 unmapped reads contain at least some data from a novel species of *Proteini-
333 clasticum*; this matches the observation in [12] of a contaminating genome
334 from an unknown *Clostridium* spp., as at the time there was no *P. ruminis*
335 genome.

336 Discussion

337 Assembly recovers basic content sensitively and accurately.

338 All three assemblers performed well in assembling contigs from the con-
339 tent that was fully present in reads and k-mers. After length filtering,
340 all three assemblies contained more than 95% of the reference (Table 6);
341 even with removal of secondary alignments, more than 87% was recovered
342 by each assembler (Table 7). About half the constituent genomes had an
343 NGA50 of 50kb or higher (Figure 2), which, while low for current Illumina
344 single-genome sequencing, is sufficient to recover operon-level relationships
345 for many genes.

346 The presence of multiple closely related genomes confounds 347 assembly.

348 In agreement with CAMI, we also find that the presence of closely related
349 genomes in the metagenome causes loss of assembly [3]. This is clearly shown
350 by Figure 2, where 12 of the bottom 14 genomes by NGA50 (left panel)
351 also exhibit poor genome recovery by assembly (right panel). Interestingly,
352 different assemblers handle this quite differently, with e.g. MetaSPAdes
353 failing to recover essentially any of *Thermotoga petrophila*, while MEGAHIT
354 recovers 73%. The presence of nearby genomes is an almost perfect predictor
355 that one or more assembler will fail to recover 5% or more - of the 13/51
356 genomes for which less than 95% is recovered, 12 of them have close genomes
357 in the community. Interestingly, very little similarity is needed - all genomes
358 with Jaccard similarity of 2% or higher at k=31 exhibit these problems.

359 The *Shewanella baltica* OS185 genome is a good example: there are two
360 strain variants, OS185 and OS223, present in the defined community. Both
361 are present at more than 99% in the reads, and more than 98% in 51-mers,
362 but only 75% of *S. baltica* OS185 and 50% of *S. baltica* OS223 are recovered
363 by assemblers. This is a clear case of “strain confusion” where the assemblers
364 simply fail to output contigs for a substantial portion of the two genomes.

365 Another interest of this study was to examine cross-species chimeric as-
366 sembly, in which a single contig is formed from multiple genomes. In Table 8,
367 we show that there is relatively little cross-species chimerism. Surprisingly,
368 what little is present is length-dependent: longer contigs are less likely to
369 be chimeric. This might well be due to the same “strain confusion” effect
370 as above, where contigs that share paths in the assembly graphs are broken
371 in twain.

372 **MEGAHIT performs best by several metrics.**

373 MEGAHIT is clearly the most efficient computationally, outperforming both
374 MetaSPAdes and IDBA by 3-9x in memory and 12-23x in time (Table 4).
375 The MEGAHIT assembly also included more of the reads than either IDBA
376 or MetaSPAdes, and omitted only 0.4% more of the unique 51-mers from
377 the reads than IDBA. MEGAHIT covered more of the reference genome
378 with both loose and strict alignments (Table 6 and Table 7), with little
379 duplication. This is clearly because of MEGAHIT's generally superior per-
380 formance in recovering the genomes of closely related strains (Figure 2, right
381 panel). The sum "fraction of genome recovered" is arguably the most im-
382 portant measure of a metagenome assembler (see [5] in particular) and here
383 MEGAHIT excels for individual genomes even in the presence of strain vari-
384 ation.

385 When comparing details of sequence recovery between the assemblers,
386 the assembly content differs by only a small amount when loose alignments
387 are allowed: all three assemblers miss more content (approximately 2.5% of
388 the reference) than they generate uniquely (1.7% or less). In addition to
389 preferring no one assembler over any other, this suggests that combining as-
390 semblies may have little value in terms of recovering additional metagenome
391 content.

392 **The missing reference may be present in strain variants of the**
393 **intended species.**

394 Several individual genomes are missing in measurable portion from the QC
395 reads (Table 2), and many QC reads (4.4% of 108m) did not map to the
396 full reference metagenome. These appear to be related issues: upon anal-
397 ysis of the unmapped reads against GenBank, we find that many of the
398 contigs assembled from the unmapped reads can be assigned to strain vari-
399 ants of the species in the mock community (Table 9). This suggests that
400 the constructors of the mock community may have unintentionally included
401 strain variants of *Fusobacterium nucleatum*, *Thermus thermophilus* HB27,
402 and *Enterococcus faecalis*; note that the microbes used were sourced from
403 the community rather than the ATCC (M. Podar, pers. communication). In
404 addition, we detect what may be portions of a novel member of the *Proteini-*
405 *clasticum* genus in the assembly of these reads - this is likely the *Clostridium*
406 spp. detected through amplicon sequencing in [12].

407 Without returning to the original DNA samples, it is impossible to con-
408 clusively confirm that unintended strains were used in the construction of the

409 mock community. In particular, our analysis is dependent on the genomes in
410 GenBank: the genomes we detect in the contigs are clearly closely related to
411 GenBank genomes not in the reference metagenome, based on k-mer anal-
412 ysis and contig alignment. However, GenBank is unlikely to contain the
413 exact genomes of the actually included strain variants, rendering conclusive
414 identification impossible.

415 Conclusions

416 Overall, assembly of this mock community works well, with good recovery
417 of known genomic sequence for the majority of genomes. All three assem-
418 blers that we evaluated recover similar amounts of most genomic sequence,
419 but (recapitulating several other studies [3, 5, 15]) MEGAHIT is compu-
420 tationally the most efficient of the three. We note that assembly resolves
421 substantial portions of several previously undetected strain variants, as well
422 as recovering a substantial portion of a novel *Proteiniclasticum* spp. that
423 was detected via amplicon analysis in [12], suggesting that assembly is a
424 useful complement to amplicon or reference-based analyses.

425 The presence of closely related strains is a major confounder of metagenome
426 assembly, and causes assemblers to drop considerable portions of genomes
427 that (based on read mapping and k-mer inclusion) are clearly present. In this
428 relatively simple community, this strain confusion is present but does not
429 dominate the assembly. However, real microbial communities are likely to
430 have many closely related strains and any resulting loss of assembly would
431 be hard to detect in the absence of good reference genomes. While high
432 polymorphism rates in e.g. animal genomes are known to cause duplication
433 or loss of assembly, some solutions have emerged that make use of assump-
434 tions of uniform coverage and diploidy [31]. These solutions cannot however
435 be transferred directly to metagenomes, which have unknown abundance
436 distributions and strain content.

437 An additional concern is that metagenome assemblies are often per-
438 formed after pooling data sets to increase coverage (e.g. [4, 32]); this pooled
439 data is more likely to contain multiple strains, which would then in turn
440 adversely affect assembly of strains. This may not be resolvable within the
441 current paradigm of assembly, which focuses on outputting linear assem-
442 blies that cannot properly represent strain variation. The human genomics
443 community is moving towards using *reference graphs*, which can represent
444 multiple incompatible variants in a single data structure [33]; this approach,
445 however, requires high-quality isolate reference genomes, which are generally

446 unavailable for environmental microbes.

447 Long read sequencing (and related technologies) will undoubtedly help
448 resolve strain variation in the future, but even with highly accurate long-
449 read sequencing, current sequencing depth is still too low to resolve deep
450 environmental metagenomes [34, 35]. It is unclear how well long error-
451 prone reads (such as those output by Pacific Biosciences SMRT [36] and
452 Oxford Nanopore instruments [37]) will perform on complex metagenomes:
453 with high error rates, deep coverage of each individual genome is required
454 to achieve accurate assembly, and this may not be easily obtainable for
455 complex communities. Single-molecule barcoding (e.g. 10X Genomics [38])
456 and HiC approaches [39] show promise but these remain untested on well-
457 defined complex communities and are still challenged by the complexity of
458 complex environmental metagenomes; see [40, 41, 42].

459 Much of our analysis above depends on having a high-quality “mock”
460 metagenome. While computationally constructed synthetic communities
461 and computational “spike-ins” to real data sets can provide valuable controls
462 (e.g. see [15] and [43]) we strongly believe that standardized communities
463 constructed *in vitro* and sequenced with the latest technologies are critical to
464 the evaluation of both canonical and emerging tools, e.g. efforts such as [44].
465 From the perspective of tool evaluation, we must disagree somewhat with
466 Vollmers et al. [5]: good metagenome tool evaluation necessarily depends on
467 mock communities that are as realistic as we can make them. Likewise, from
468 the perspective of bench biologists, actually sequencing real DNA is critical
469 because it can evaluate confounding effects such as kit contamination [45].
470 Large-scale studies of computational approaches systematically applied to
471 mock communities such as CAMI [3] can then provide fair comparisons of
472 entire toolchains (wet + dry) applied to these mock communities.

473 We omitted two important questions in this study: binning and choice
474 of parameters. We chose not to evaluate genome binning because most
475 binning strategies either operate post-assembly (see e.g. [46]), in which
476 case the challenges with assembly discussed above will apply; or require
477 multiple samples (e.g. [47]), which we do not have. We also chose to use
478 only default parameters with all three assemblers, for two reasons. First,
479 we are not aware of any widely used automated approaches for determining
480 the “best” set of parameters or evaluating the output, other than those
481 integrated into the assemblers themselves (e.g. choice of k-mer sizes), and
482 absent such guidance we do not feel comfortable blessing any particular set of
483 parameters; here the choice of default parameters is parsimonious. Second,
484 any parameter exploration pipeline would not only need to be automated

485 but would need to run multiple assemblies, whose time and resource usage
486 should be measured; in this case, any comparison based on runtime of the
487 parameter choice pipeline should naturally favor MEGAHIT because of its
488 substantial advantage in computational efficiency.

489 **Author contributions**

490 SA, LI and CTB developed, tested, and executed the analytical pipeline.
491 SA and CTB created the tables and figures and wrote the paper.

492 **Competing interests**

493 No competing interest to our knowledge.

494 **Grant information**

495 This work is funded by Gordon and Betty Moore Foundation Grant GBMF4551
496 and NIH NHGRI R01 grant HG007513-03, both to CTB.

497 **Acknowledgments**

498 We thank Michael R. Crusoe and Phillip T. Brooks for input on analysis and
499 pipeline development. We thank Migun Shakya, Mircea Podar, Jiarong Guo,
500 Harald R. Gruber-Vodicka, Juliane Wippler, Krista Ternus, and Stephen
501 Turner for valuable comments on drafts of this manuscript.

502 **References**

- 503 [1] Jay Ghurye, Victoria Cepeda-Espinoza, and Mihai Pop. Metagenomic assem-
504 bly: Overview, challenges and applications. *The Yale Journal of Biology and*
505 *Medicine*, 89(3):353–362, 2016.
- 506 [2] Nikos C. Kyrpides, Philip Hugenholtz, Jonathan A. Eisen, Tanja Woyke,
507 Markus Göker, Charles T. Parker, Rudolf Amann, Brian J. Beck, Patrick S. G.
508 Chain, Jongsik Chun, Rita R. Colwell, Antoine Danchin, Peter Dawyndt, Tom
509 Dedeurwaerdere, Edward F. DeLong, John C. Detter, Paul De Vos, Timothy J.
510 Donohue, Xiu-Zhu Dong, Dusko S. Ehrlich, Claire Fraser, Richard Gibbs, Jack
511 Gilbert, Paul Gilna, Frank Oliver Glöckner, Janet K. Jansson, Jay D. Keasling,
512 Rob Knight, David Labeda, Alla Lapidus, Jung-Sook Lee, Wen-Jun Li, Juncai
513 MA, Victor Markowitz, Edward R. B. Moore, Mark Morrison, Folker Meyer,
514 Karen E. Nelson, Moriya Ohkuma, Christos A. Ouzounis, Norman Pace, Julian

- 515 Parkhill, Nan Qin, Ramon Rossello-Mora, Johannes Sikorski, David Smith,
516 Mitch Sogin, Rick Stevens, Uli Stingl, Ken ichiro Suzuki, Dorothea Taylor,
517 Jim M. Tiedje, Brian Tindall, Michael Wagner, George Weinstock, Jean Weis-
518 senbach, Owen White, Jun Wang, Lixin Zhang, Yu-Guang Zhou, Dawn Field,
519 William B. Whitman, George M. Garrity, and Hans-Peter Klenk. Genomic
520 encyclopedia of bacteria and archaea: Sequencing a myriad of type strains.
521 *PLoS Biology*, 12(8):e1001920, aug 2014. doi: 10.1371/journal.pbio.1001920.
522 URL <https://doi.org/10.1371/journal.pbio.1001920>.
- 523 [3] Alexander Sczyrba, Peter Hofmann, Peter Belmann, David Koslicki, Stefan
524 Janssen, Johannes Droege, Ivan Gregor, Stephan Majda, Jessika Fiedler,
525 Eik Dahms, Andreas Bremges, Adrian Fritz, Ruben Garrido-Oter, Tue
526 Sparholt Jorgensen, Nicole Shapiro, Philip D Blood, Alexey Gurevich, Yang
527 Bai, Dmitrij Turaev, Matthew Z DeMaere, Rayan Chikhi, Niranjan Nagara-
528 jan, Christopher Quince, Lars Hestbjerg Hansen, Soren J Sorensen, Burton
529 K H Chia, Bertrand Denis, Jeff L Froula, Zhong Wang, Robert Egan, Dong-
530 wan Don Kang, Jeffrey J Cook, Charles Deltel, Michael Beckstette, Claire
531 Lemaitre, Pierre Peterlongo, Guillaume Rizk, Dominique Lavenier, Yu-Wei
532 Wu, Steven W Singer, Chirag Jain, Marc Strous, Heiner Klingenberg, Peter
533 Meinicke, Michael Barton, Thomas Lingner, Hsin-Hung Lin, Yu-Chieh Liao,
534 Genivaldo Gueiros Z. Silva, Daniel A Cuevas, Robert A Edwards, Surya Saha,
535 Vitor C Piro, Bernhard Y Renard, Mihai Pop, Hans-Peter Klenk, Markus
536 Goeker, Nikos Kyrpides, Tanja Woyke, Julia A Vorholt, Paul Schulze-Lefert,
537 Edward M Rubin, Aaron E Darling, Thomas Rattei, and Alice C McHardy.
538 Critical assessment of metagenome interpretation - a benchmark of compu-
539 tational metagenomics software. *bioRxiv*, 2017. doi: 10.1101/099127. URL
540 <http://biorxiv.org/content/early/2017/01/09/099127>.
- 541 [4] I. Sharon, M. J. Morowitz, B. C. Thomas, E. K. Costello, D. A. Relman,
542 and J. F. Banfield. Time series community genomics analysis reveals rapid
543 shifts in bacterial species, strains, and phage during infant gut colonization.
544 *Genome Research*, 23(1):111–120, aug 2012. doi: 10.1101/gr.142315.112. URL
545 <https://doi.org/10.1101/gr.142315.112>.
- 546 [5] John Vollmers, Sandra Wiegand, and Anne-Kristin Kaster. Compar-
547 ing and evaluating metagenome assembly tools from a microbio-
548 logist’s perspective - not only size matters! *PLOS ONE*, 12
549 (1):e0169662, jan 2017. doi: 10.1371/journal.pone.0169662. URL
550 <https://doi.org/10.1371/journal.pone.0169662>.
- 551 [6] Jorge F Vázquez-Castellanos, Rodrigo García-López, Vicente Pérez-Brocal,
552 Miguel Pignatelli, and Andrés Moya. Comparison of different assembly and
553 annotation tools on analysis of simulated viral metagenomic communities in
554 the gut. *BMC genomics*, 15(1):1, 2014.
- 555 [7] Konstantinos Mavromatis, Natalia Ivanova, Kerrie Barry, Harris Shapiro, Eu-
556 gene Goltsman, Alice C McHardy, Isidore Rigoutsos, Asaf Salamov, Frank

- 557 Korzeniewski, Miriam Land, et al. Use of simulated data sets to evaluate the
558 fidelity of metagenomic processing methods. *Nature methods*, 4(6):495–500,
559 2007.
- 560 [8] David B Jaffe, Jonathan Butler, Sante Gnerre, Evan Mauceli, Kerstin
561 Lindblad-Toh, Jill P Mesirov, Michael C Zody, and Eric S Lander. Whole-
562 genome sequence assembly for mammalian genomes: Arachne 2. *Genome
563 research*, 13(1):91–96, 2003.
- 564 [9] Samuel Aparicio, Jarrod Chapman, Elia Stupka, Nik Putnam, Jer-ming Chia,
565 Paramvir Dehal, Alan Christoffels, Sam Rash, Shawn Hoon, Arian Smit, et al.
566 Whole-genome shotgun assembly and analysis of the genome of *fugu rubripes*.
567 *Science*, 297(5585):1301–1310, 2002.
- 568 [10] Anveshi Charuvaka and Huzefa Rangwala. Evaluation of short read metage-
569 nomic assembly. *BMC genomics*, 12(2):1, 2011.
- 570 [11] Jared T Simpson, Kim Wong, Shaun D Jackman, Jacqueline E Schein,
571 Steven JM Jones, and Inanç Birol. Abyss: a parallel assembler for short read
572 sequence data. *Genome research*, 19(6):1117–1123, 2009.
- 573 [12] Shakya Migun, Christopher Quince, James Campbell, Zamin Yang, Christo-
574 pher Schadt, and Mircea Podar. Comparative metagenomic and rrna microbial
575 diversity characterization using archaeal and bacterial synthetic communities.
576 *Environmental Microbiology*, 15(6):1882–1899, 2013.
- 577 [13] Brandon K. B. Seah and Harald R. Gruber-Vodicka. gbtools: In-
578 teractive visualization of metagenome bins in r. *Frontiers in Mi-
579 crobiology*, 6, dec 2015. doi: 10.3389/fmicb.2015.01451. URL
580 <https://doi.org/10.3389/fmicb.2015.01451>.
- 581 [14] Dinghua Li, Ruibang Luo, Chi-Man Liu, Chi-Ming Leung, Hing-
582 Fung Ting, Kunihiko Sadakane, Hiroshi Yamashita, and Tak-Wah
583 Lam. MEGAHIT v1.0: A fast and scalable metagenome assembler
584 driven by advanced methodologies and community practices. *Meth-
585 ods*, 102:3–11, jun 2016. doi: 10.1016/j.ymeth.2016.02.020. URL
586 <https://doi.org/10.1016/j.ymeth.2016.02.020>.
- 587 [15] Andries Johannes van der Walt, Marc Warwick Van Goethem,
588 Jean-Baptiste Ramond, Thulani Peter Makhwanyane, Oleg Reva,
589 and Don Arthur Cowan. Assembling metagenomes, one com-
590 munity at a time. *bioRxiv*, 2017. doi: 10.1101/120154. URL
591 <http://biorxiv.org/content/early/2017/06/06/120154>.
- 592 [16] William W. Greenwald, Niels Klitgord, Victor Seguritan, Shibu Yooseph,
593 J. Craig Venter, Chad Garner, Karen E. Nelson, and Weizhong Li. Utilization
594 of defined microbial communities enables effective evaluation of meta-genomic

- 595 assemblies. *BMC Genomics*, 18(1), apr 2017. doi: 10.1186/s12864-017-3679-5.
596 URL <https://doi.org/10.1186/s12864-017-3679-5>.
- 597 [17] Yu Peng, Henry C.M. Leung, S.M. Yiu, and Francis Y.L. Chin. Idba-ud: a de
598 novo assembler for single-cell and metagenomic sequencing data with highly
599 uneven depth. *Bioinformatics*, 28:1420–1428, 2012.
- 600 [18] Sergey Nurk, Dmitry Meleshko, Anton Korobeynikov, and Pavel A. Pevzner.
601 metaSPAdes: a new versatile metagenomic assembler. *Genome Re-*
602 *search*, 27(5):824–834, mar 2017. doi: 10.1101/gr.213959.116. URL
603 <https://doi.org/10.1101/gr.213959.116>.
- 604 [19] Dinghua Li, Ruibang Luo, Chi-Man Liu, Chi-Ming Leung, Hing-Fung Ting,
605 Kunihiko Sadakane, Hiroshi Yamashita, and Tak-Wah Lam. Megahit v1. 0:
606 A fast and scalable metagenome assembler driven by advanced methodologies
607 and community practices. *Methods*, 102:3–11, 2016.
- 608 [20] H Chitsaz, JL Yee-Greenbaum, G Tesler, MJ Lombardo, CL Dupont, JH Bad-
609 ger, M Novotny, DB Rusch, LJ Fraser, NA Gormley, O Schulz-Trieglaff,
610 GP Smith, DJ Evers, PA Pevzner, and RS Lasken. Efficient de novo assembly
611 of single-cell bacterial genomes from short-read data sets. *Nat Biotechnol*, 29
612 (10):915–21, 2011.
- 613 [21] Anthony M. Bolger, Marc Lohse, and Bjoern Usadel. Trimmomatic: A flexible
614 trimmer for illumina sequence data. *Bioinformatics*, 30(15):2114–2120, 2014.
- 615 [22] Matthew D MacManes. On the optimal trimming of high-throughput mrna
616 sequence data. *Frontiers in genetics*, 5:13, 2014.
- 617 [23] Heng Li and Richard Durbin. Fast and accurate short read alignment with
618 burrows–wheeler transform. *Bioinformatics*, 25(14):1754–1760, 2009.
- 619 [24] C. Titus Brown and Luiz Irber. sourmash: a library for MinHash sketch-
620 ing of DNA. *The Journal of Open Source Software*, 1(5), sep 2016. doi:
621 10.21105/joss.00027. URL <https://doi.org/10.21105/joss.00027>.
- 622 [25] Brian D. Ondov, Todd J. Treangen, Páll Melsted, Adam B. Mal-
623 lonee, Nicholas H. Bergman, Sergey Koren, and Adam M. Phillippy.
624 Mash: fast genome and metagenome distance estimation using MinHash.
625 *Genome Biology*, 17(1), jun 2016. doi: 10.1186/s13059-016-0997-x. URL
626 <https://doi.org/10.1186/s13059-016-0997-x>.
- 627 [26] David Koslicki and Daniel Falush. Metapalette: a k-mer painting approach
628 for metagenomic taxonomic profiling and quantification of novel strain vari-
629 ation. *mSystems*, 1(3), 2016. doi: 10.1128/mSystems.00020-16. URL
630 <http://msystems.asm.org/content/1/3/e00020-16>.

- 631 [27] Zhang Qingpeng, Awad Sherine, and Brown Titus. Crossing the streams:
632 a framework for streaming analysis of short dna sequencing reads. *PeerJ*
633 *PrePrints* 3:e1100 <https://dx.doi.org/10.7287/peerj.preprints.890v1>, 2015.
- 634 [28] MR Crusoe, HF Alameldin, S Awad, E Boucher, A Caldwell, R Cartwright,
635 A Charbonneau, B Constantinides, G Edverson, S Fay, J Fenton, T Fenzl,
636 J Fish, L Garcia-Gutierrez, P Garland, J Gluck, I Gonzalez, S Guermond,
637 J Guo, A Gupta, JR Herr, A Howe, A Hyer, A Hrpfer, L Irber, R Kidd,
638 D Lin, J Lippi, T Mansour, P McA’Nulty, E McDonald, J Mizzi, KD Mur-
639 ray, JR Nahum, K Nanlohy, AJ Nederbragt, H Ortiz-Zuazaga, J Ory, J Pell,
640 C Pepe-Ranne, ZN Russ, E Schwarz, C Scott, J Seaman, S Sievert, J Simp-
641 son, CT Skennerton, J Spencer, R Srinivasan, D Standage, JA Stapleton,
642 SR Steinman, J Stein, B Taylor, W Trimble, HL Wiencko, M Wright,
643 B Wyss, Q Zhang, e zyme, and CT Brown. The khmer software pack-
644 age: enabling efficient nucleotide sequence analysis [version 1; referees: 2 ap-
645 proved, 1 approved with reservations]. *F1000Research*, 4(900), 2015. doi:
646 10.12688/f1000research.6924.1.
- 647 [29] Heng Li, Bob Handsaker, Alec Wysoker, Tim Fennell, Jue Ruan, Nils Homer,
648 Gabor Marth, Goncalo Abecasis, Richard Durbin, et al. The sequence align-
649 ment/map format and samtools. *Bioinformatics*, 25(16):2078–2079, 2009.
- 650 [30] Stefan Kurtz, Adam Phillippy, Arthur L Delcher, Michael Smoot, Martin
651 Shumway, Corina Antonescu, and Steven L Salzberg. Versatile and open soft-
652 ware for comparing large genomes. *Genome biology*, 5(2):1, 2004.
- 653 [31] J. H. Kim, M. S. Waterman, and L. M. Li. Diploid genome reconstruc-
654 tion of *Ciona intestinalis* and comparative analysis with *Ciona savignyi*.
655 *Genome Research*, 17(7):1101–1110, jun 2007. doi: 10.1101/gr.5894107. URL
656 <https://doi.org/10.1101/gr.5894107>.
- 657 [32] Ping Hu, Lauren Tom, Andrea Singh, Brian C. Thomas, Brett J. Baker,
658 Yvette M. Piceno, Gary L. Andersen, and Jillian F. Banfield. Genome-
659 resolved metagenomic analysis reveals roles for candidate phyla and other
660 microbial community members in biogeochemical transformations in oil reser-
661 voirs. *mBio*, 7(1):e01669–15, jan 2016. doi: 10.1128/mbio.01669-15. URL
662 <https://doi.org/10.1128/mbio.01669-15>.
- 663 [33] Benedict Paten, Adam M Novak, Jordan M Eizenga, and Erik Garrison.
664 Genome graphs and the evolution of genome inference. *Genome research*,
665 27(5):665–676, 2017.
- 666 [34] Itai Sharon, Michael Kertesz, Laura A. Hug, Dmitry Pushkarev, Timothy A.
667 Blauwkamp, Cindy J. Castelle, Mojgan Amirebrahimi, Brian C. Thomas,
668 David Burstein, Susannah G. Tringe, Kenneth H. Williams, and Jillian F.
669 Banfield. Accurate, multi-kb reads resolve complex populations and de-
670 tect rare microorganisms. *Genome Research*, 25(4):534–543, feb 2015. doi:
671 10.1101/gr.183012.114. URL <https://doi.org/10.1101/gr.183012.114>.

- 672 [35] Richard Allen White, Eric M. Bottos, Taniya Roy Chowdhury, Jeremy D.
673 Zucker, Colin J. Brislawn, Carrie D. Nicora, Sarah J. Fansler, Kurt R. Glae-
674 semann, Kevin Glass, and Janet K. Jansson. Moleculo long-read sequenc-
675 ing facilitates assembly and genomic binning from complex soil metagenomes.
676 *mSystems*, 1(3):e00045–16, jun 2016. doi: 10.1128/msystems.00045-16. URL
677 <https://doi.org/10.1128/msystems.00045-16>.
- 678 [36] J. Eid, A. Fehr, J. Gray, K. Luong, J. Lyle, G. Otto, P. Peluso, D. Rank,
679 P. Baybayan, B. Bettman, A. Bibillo, K. Bjornson, B. Chaudhuri, F. Chris-
680 tians, R. Cicero, S. Clark, R. Dalal, A. deWinter, J. Dixon, M. Foquet,
681 A. Gaertner, P. Hardenbol, C. Heiner, K. Hester, D. Holden, G. Kearns,
682 X. Kong, R. Kuse, Y. Lacroix, S. Lin, P. Lundquist, C. Ma, P. Marks,
683 M. Maxham, D. Murphy, I. Park, T. Pham, M. Phillips, J. Roy, R. Se-
684 bra, G. Shen, J. Sorenson, A. Tomaney, K. Travers, M. Trulson, J. Vieceli,
685 J. Wegener, D. Wu, A. Yang, D. Zaccarin, P. Zhao, F. Zhong, J. Korlach,
686 and S. Turner. Real-time DNA sequencing from single polymerase molecules.
687 *Science*, 323(5910):133–138, jan 2009. doi: 10.1126/science.1162986. URL
688 <https://doi.org/10.1126/science.1162986>.
- 689 [37] Gerald M Cherf, Kate R Lieberman, Hytham Rashid, Christopher E Lam,
690 Kevin Karplus, and Mark Akeson. Automated forward and reverse ratcheting
691 of DNA in a nanopore at 5-AA precision. *Nature Biotechnology*, 30(4):344–348,
692 feb 2012. doi: 10.1038/nbt.2147. URL <https://doi.org/10.1038/nbt.2147>.
- 693 [38] Eli Moss, Alex Bishara, Ekaterina Tkachenko, Joyce B Kang,
694 Tessa M Andermann, Christina Wood, Christine Handy, Hanlee
695 Ji, Serafim Batzoglou, and Ami S Bhatt. De novo assembly of
696 microbial genomes from human gut metagenomes using barcoded
697 short read sequences. *bioRxiv*, 2017. doi: 10.1101/125211. URL
698 <http://biorxiv.org/content/early/2017/04/07/125211>.
- 699 [39] Caiti Smukowski Heil, Joshua N. Burton, Ivan Liachko, Anne Friedrich,
700 Noah A. Hanson, Cody L. Morris, Joseph Schacherer, Jay Shendure,
701 James H. Thomas, and Maitreya J. Dunham. Identification of a
702 novel interspecific hybrid yeast from a metagenomic open fermentation
703 sample using hi-c. *bioRxiv*, 2017. doi: 10.1101/150722. URL
704 <http://biorxiv.org/content/early/2017/06/15/150722>.
- 705 [40] Joshua N. Burton, Ivan Liachko, Maitreya J. Dunham, and Jay Shendure.
706 Species-level deconvolution of metagenome assemblies with hi-c–based contact
707 probability maps. *G3*, 4(7):1339–1346, may 2014. doi: 10.1534/g3.114.011825.
708 URL <https://doi.org/10.1534/g3.114.011825>.
- 709 [41] Martial Marbouty, Axel Cournac, Jean-François Flot, Hervé Marie-Nelly,
710 Julien Mozziconacci, and Romain Koszul. Metagenomic chromosome con-
711 formation capture (meta3c) unveils the diversity of chromosome organiza-

- 712 tion in microorganisms. *eLife*, 3, dec 2014. doi: 10.7554/elife.03318. URL
713 <https://doi.org/10.7554/elife.03318>.
- 714 [42] Christopher W. Beitel, Lutz Froenicke, Jenna M. Lang, Ian F. Korf,
715 Richard W. Micheltore, Jonathan A. Eisen, and Aaron E. Darling. Strain- and
716 plasmid-level deconvolution of a synthetic metagenome by sequencing proxim-
717 ity ligation products. *PeerJ*, 2:e415, may 2014. doi: 10.7717/peerj.415. URL
718 <https://doi.org/10.7717/peerj.415>.
- 719 [43] Adina Chuang Howe, Janet K Jansson, Stephanie A Malfatti, Susannah G
720 Tringe, James M Tiedje, and C Titus Brown. Tackling soil diversity with the
721 assembly of large, complex metagenomes. *Proceedings of the National Academy
722 of Sciences*, 111(13):4904–4909, 2014.
- 723 [44] Bonnie L. Brown, Mick Watson, Samuel S. Minot, Maria C.
724 Rivera, and Rima B. Franklin. MinION™ nanopore sequenc-
725 ing of environmental metagenomes: a synthetic approach. *Giga-
726 Science*, 6(3):1–10, feb 2017. doi: 10.1093/gigascience/gix007. URL
727 <https://doi.org/10.1093/gigascience/gix007>.
- 728 [45] Susannah J Salter, Michael J Cox, Elena M Turek, Szymon T Calus,
729 William O Cookson, Miriam F Moffatt, Paul Turner, Julian Parkhill,
730 Nicholas J Loman, and Alan W Walker. Reagent and laboratory
731 contamination can critically impact sequence-based microbiome analyses.
732 *BMC Biology*, 12(1), nov 2014. doi: 10.1186/s12915-014-0087-z. URL
733 <https://doi.org/10.1186/s12915-014-0087-z>.
- 734 [46] Cedric C Laczny, Christina Kiefer, Valentina Galata, Tobias Fehlmann,
735 Christina Backes, and Andreas Keller. Busybee web: metagenomic data anal-
736 ysis by bootstrapped supervised binning and annotation. *Nucleic Acids Re-
737 search*, page gkx348, 2017.
- 738 [47] Brian Cleary, Ilana Lauren Brito, Katherine Huang, Dirk Gevers, Terrance
739 Shea, Sarah Young, and Eric J Alm. Detection of low-abundance bacte-
740 rial strains in metagenomic datasets by eigengenome partitioning. *Nature
741 Biotechnology*, 33(10):1053–1060, sep 2015. doi: 10.1038/nbt.3329. URL
742 <https://doi.org/10.1038/nbt.3329>.