

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

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

4 June 17, 2017

5 **Abstract**

6 We evaluate the performance of three metagenome assemblers, IDBA,
7 SPAdes, and MEGAHIT, on short-read sequencing of a defined “mock”
8 community from Shakya et al. (2013) containing 64 genomes. We up-
9 date 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 of reference genomes that are oth-
12 erwise completely present in the read data set. In agreement with
13 previous studies, we find that MEGAHIT performs best computation-
14 ally; we also show that MEGAHIT tends to recover larger portions of
15 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. Most approaches to analyzing
21 metagenomic data rely on mapping or comparing sequencing reads to refer-
22 ence sequence collections. However, reference databases contain only a small
23 subset of microbial diversity [1], and the much of the remaining diversity is
24 evolutionarily distant and search techniques may not access it.

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 There are a number of metagenome assemblers that are widely used. How-
29 ever, evaluating the results of these assemblers is challenging due to the
30 general lack of good quality reference metagenomes.

31 Moya et al. in [2] evaluated metagenome assembly using two simulated
32 454 viral metagenome and six assemblers. The assemblies were evaluated
33 based on several metrics including N50, percentages of reads assembled, ac-
34 curacy when compared to the reference genome. In addition to, chimeras per
35 contigs and the effect of assembly on taxonomic and functional annotations.

36 Mavromatis et al. in [3] provided a benchmark study to evaluate the
37 fidelity of metagenome process methods. The study used simulated metage-
38 nomic data sets constructed at different complexity levels. The datasets were
39 assembled using Phrap v3.57, Arachne v.2 [4] and JAZZ. [5] This study eval-
40 uates assembly, gene prediction, and binning methods. However, the study
41 did not evaluate the assembly quality against a reference genome.

42 Rangwala et al. in [6] presented an evaluation study of metagenome
43 assembly. The study used a de Bruijn graph based assembler ABYSS [7] to
44 assemble simulated metagenome reads of 36 bp. The data set is classified at
45 different complexity levels. The study compares the quality of the assembly
46 of the data sets in terms of quality measures of contigs length, assembly
47 accuracy. The study also took into consideration the effect of kmer size and
48 the degree of chimericity. However, the study evaluated the assembly based
49 on one assembler, and did not evaluate assembly against several assemblers.
50 Also, both previous studies used simulated data, which may lack confounders
51 of assembly such as sequencing artifacts and GC bias.

52 Shakya et al. (2013) constructed a synthetic community of organisms by
53 mixing DNA isolated from individual cultures of 64 bacteria and archaea,

54 including a variety of strains across a range of nucleotide distances [8]. In
55 addition to performing 16s amplicon analysis and doing 454 sequencing, the
56 authors shotgun-sequenced the mixture with Illumina. While the authors
57 concluded that this metagenomic sequencing generally outperformed ampli-
58 con sequencing, they did not conduct an assembly based analysis.

59 More recently, several benchmark studies systematically evaluated metagenome
60 assembly of short reads. The Critical Assessment of Metagenome Interpre-
61 tation (CAMI) collaboration benchmarked a number of metagenome assem-
62 blers on several data sets of varying complexity, evaluating recovery of novel
63 genomes and multiple strain variants [9]. Notably, CAMI concluded that
64 “The resolution of strain-level diversity represents a substantial challenge to
65 all evaluated programs.” Another recent study evaluated eight assemblers
66 on nine environmental metagenomes and three simulated data sets [10].

67 In this study, we extend previous work by delving into questions of
68 chimeric misassembly and strain recovery in the Shakya et al. (2013) data
69 set. First, we update the list of reference genomes for Shakya et al. to in-
70 clude the latest Genbank assemblies. We then compare IDBA [11], SPAdes
71 [12], and MEGAHIT [13] performance on assembling this short-read data
72 set, and explore concordance in recovery between the three assemblers. We
73 describe the effects of “strain confusion” between multiple strains. We also
74 detect and analyze several previously unreported strains and genomes in
75 the Shakya et al. data set. We find that in the absence of closely related
76 genomes, all three metagenome assemblers recover 95% or more of known
77 reference genomes. However, in the presence of closely related genomes,
78 metagenome assemblers vary widely in their performance and can fail to re-
79 cover the majority of some genomes even when they are completely present
80 in the reads. Our report provides strong guidance on choice of assemblers
81 and extends previous analyses of this low-complexity metagenome bench-
82 marking data set.

83 Datasets

84 We used a diverse mock community data set constructed by pooling DNA
85 from 64 species of bacteria and archaea and sequencing them with Illumina
86 HiSeq. The raw data set consisted of 109,629,496 reads from Illumina HiSeq
87 101 bp paired-end sequencing (2x101) with an untrimmed total length of
88 11.07 Gbp and an estimated fragment size of 380 bp [8].

89 The original reads are available through the NCBI Sequence Read Archive
90 at Accession SRX200676. We updated the 64 reference genomes sets from

91 NCBI Genbank using the latest available assemblies (June 2017); updated
92 data is available for download at <https://osf.io/8uxj9/>.

93 **Methods**

94 The analysis code and run scripts for this paper are available at: [https://github.com/dib-](https://github.com/dib-lab/2015-metagenome-assembly/)
95 [lab/2015-metagenome-assembly/](https://github.com/dib-lab/2015-metagenome-assembly/). The scripts and overall pipeline were ex-
96 amined by the first and senior authors for correctness. In addition, the
97 bespoke reference-based analysis scripts were tested by running them on a
98 single-colony *E. coli* MG1655 data set with a high quality reference genome
99 [14].

100 **Quality Filtering**

101 We removed adapters with Trimmomatic v0.30 in paired-end mode with the
102 Truseq adapters [15], using light quality score trimming as recommended in
103 MacManes, 2014 [16].

104 **Reference Coverage Profile**

105 To evaluate how much of the reference metagenome was contained in the
106 read data, we used `bwa aln` (v0.7.7.r441) to map reads to the reference
107 genome [17]. We then calculated how many reference bases were covered by
108 mapped reads (custom script `coverage-profile.py`).

109 **Measuring k-mer inclusion and Jaccard similarity**

110 We used MinHashing as implemented in sourmash to estimate k-mer inclu-
111 sion and Jaccard similarity between data sets [18]. MinHash signatures were
112 prepared with ‘sourmash compute’ using ‘-scaled 10000’. K-mer inclusion
113 was computed by taking the ratio of the number of intersecting hashes with
114 the query over the total number of hashes in the subject MinHash. Jac-
115 card similarity was computed as in [19] by taking the ratio of the number
116 of intersecting hashes between the query and subject over the number of
117 hashes in the union. K-mer sizes for comparison were chosen at 21, 31, or
118 51, depending on the level of taxonomic specificity desired - genus, species,
119 or strain, as described in [20].

120 When specified, high-abundance k-mers were selected for counting by
121 using the script `trim-low-abund.py` script with `-C 5` from khmer 2.x [21,

122 22].

123 Assemblers

124 We assembled the quality-filtered reads using three different assemblers:
125 IDBA-UD [11], MetaSPAdes [12], and MEGAHIT [13]. For IDBA-UD v1.1.1
126 [11], we used `--pre_correction` to perform pre-correction before assembly
127 and `-r` for the pe files.

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

132 For MEGAHIT v1.1.1-2-g02102e1 [13], we used `-l 101 -m 3e9 --cpu-only`
133 where `-l` is for maximum read length, `-m` is for max memory in bytes to
134 be used in constructing the graph, and `--cpu-only` to use only the CPU
135 and no GPUs. We also used `--presets meta-large` for large and complex
136 metagenomes, and `--12` and `-r` to specify the interleaved-paired-end and
137 single-end files respectively. MEGAHIT allows the specification of a memory
138 limit and we used `-M 1e+10` for 10 GB.

139 All three assemblies were executed on the same high-memory buy-in
140 node on the Michigan State University High Performance Compute Cluster,
141 and we recorded RAM and CPU time of each assembly job using the `qstat`
142 utility at the end of each run.

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

145 Mapping

146 We aligned all quality-filtered reads to the reference metagenome with `bwa`
147 `aln` (v0.7.7.r441) [17]. We aligned paired-end and orphaned reads separately.
148 We then used `samtools` (v0.1.19) [23] to convert SAM files to BAM files for
149 both paired-end and orphaned reads. To count the unaligned reads, we
150 included only those records with the “4” flag in the SAM files [23].

151 Assembly analysis using Nucmer

152 We used the `NUCmer` tool from `MUMmer3.23` [24] to align assemblies to the
153 reference genome with options `-coords -p`. Then we parsed the generated
154 “.coords” file using a custom script `analyze_assembly.py`, and calculated

several analysis metrics across all three assemblies at a 99% alignment identity.

Reference-based analysis of the assemblies

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

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

Analysis of chimeric misassemblies

We analyzed each assembly for chimeric misassemblies by counting the number of contigs that contained matches to two distinct reference genomes. In order to remove secondary alignments from consideration, we included only the longest non-overlapping NUCmer alignments for each contig at a minimum alignment identity of 99%. We then used the script `analyze_chimeric2.py` to find individual contigs that matched more than one distinct reference genome. As a negative control on our analysis, we verified that this approach yielded no positive results when applied to the alignments of the reference metagenome against itself.

Results

The raw data is high quality.

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

Trimming removed 686,735 reads (0.63%). After trimming, we retained 108,422,358 paired reads containing 10.94 Gbp with an average length of 100.9 bases. A total of 46.56 Mbp remained in 520,403 orphan reads with an average length of 89.5 bases. In total, the quality trimmed data contained

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%

186 10.98 Gbp in 108,942,761 reads. This quality trimmed (“QC”) data set was
 187 used as the basis for all further analyses.

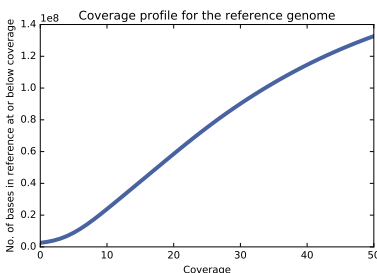


Figure 1: Cumulative coverage profile for the reference metagenome, based on read mapping.

188 **The reference metagenome is not completely present in the**
 189 **reads.**

190 We next evaluated the fraction of the reference genome covered by at least
 191 one read (see Methods for details). Quality filtered reads cover 203,058,414
 192 (98.76%) bases of the reference metagenome (205,603,715 bp total size). Fig-
 193 ure 1 shows the cumulative coverage profile of the reference metagenome,
 194 and the percentage of bases with that coverage. Most of the reference
 195 metagenome was covered at least minimally; only 3.33% of the reference
 196 metagenome had mapping coverage <5 , and 1.24% of the bases in the ref-
 197 erence were not covered by any reads in the QC data set.

198 In order to evaluate reconstructability with De Bruijn graph assemblers,
 199 we next examined k-mer containment of the reference in the reads for k of
 200 21, 31, 41, and 51 (Table 1). The k-mer overlap decreases from 96.8% to
 201 94.1% as the k-mer size increases. This could be caused by low coverage of
 202 some portions of the reference and/or variation between the reads and the

203 reference.

204 **Some individual reference genomes are poorly represented in**
205 **the reads.**

Table 2: Top uncovered genomes

Genome	Read coverage	21-mer presence
<i>B. bronchiseptica</i>	98.2%	97.3%
<i>D. vulgaris</i> DP4	93.2%	82.5%
<i>T. thermophilus</i> HB27	91.1%	79.7%
<i>E. faecalis</i> V583	74.6%	65.6%
<i>F. nucleatum</i>	47.6%	18.2%

206 To see if specific reference genomes exhibited low coverage, we analyzed
207 read mapping coverage and 21-mer containment for individual genomes.
208 Of the 64 reference genomes used in the metagenome, 59 had a per-base
209 mapping coverage above 95% and a 21-mer containment in the QC reads
210 above 95%. The remaining five varied significantly in both metrics (Table 4),
211 with *F. nucleatum* the lowest – only 47.6% of the bases in the reference
212 genome are covered by one or more mapped reads, and only 18.2% of the
213 21-mers in the *F. nucleatum* reference genome are present in the reads at
214 any abundance.

215 We next did a 51-mer containment analysis of each reference genome
216 in the reads, mimicking the analysis done in [20]. 99% or more of the
217 constituent 51-mers for 51 of the 64 reference genomes were present in the
218 reads, suggesting that each of the 51 genomes was entirely present at some
219 minimal coverage.

220 We excluded the remaining 13 genomes (see Table 3) from any compar-
221 ative analysis of assembly quality, because interpreting coverage and misas-
222 sembly analysis for these genomes would be challenging.

223 **MEGAHIT is the fastest and lowest-memory assembler eval-**
224 **uated**

225 We ran three commonly used metagenome assemblers on the QC data set:
226 IDBA-UD, SPAdes, and MEGAHIT. We recorded the time and memory
227 usage of each (Table 4). In computational requirements, MEGAHIT out-
228 performed both SPAdes and IDBA-UD considerably, producing an assembly

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

Table 4: Running Time and Memory Utilization

Assembler	CPU time	Wall time	RAM
MEGAHIT	52hr 25m	4 hr 9m	11.4 GB
IDBA-UD	17h		149.1 GB
SPAdes	94hr 43m	94hr 44m	100.7 GB

in four hours – approximately 4 times faster than IDBA and 8 times faster than SPAdes. MEGAHIT used only 11.4 GB of RAM – 1/13th to 1/9th the memory used by IDBA and SPAdes, respectively.

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%
SPAdes	3,844,123 (3.52%)	3.2%
MEGAHIT	2,737,640 (2.51%)	2.8%

We assessed read inclusion in assemblies by mapping the QC reads to the length-filtered assemblies and counting the remaining unmapped reads.

235 Depending on the assembly, between 2.7 million and 3.9 million reads (2.5-
 236 3.5%) did not map to the assemblies (Table 5). All of the assemblies included
 237 the large majority of high-abundance 51-mers (more than 96.8% in all cases).

238 **Much of the reference is covered by the assemblies.**

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

Assembly	bases aligned	duplication	51-mers
MEGAHIT	96.2%	0.72%	96.7%
SPAdes	95.8%	0.99%	96.2%
IDBA	95.6%	0.88%	97.2%

239 We next evaluated the extent to which the assembled contigs recovered
 240 the “known/true” metagenome sequence by aligning each assembly to the
 241 adjusted reference (Table 6). Each of the three assemblers generates contigs
 242 that cover more than 95.6% of the reference metagenome at high identity
 243 (99%) with little duplication (0.72-0.99%). All three assemblies contain
 244 between 96.2% and 97.2% of the 51-mers in the reference.

245 At 99% identity with the loose mapping approach, approximately 1.8%
 246 of the reference is missed by all three assemblers, while 0.9% is uniquely
 247 covered by MEGAHIT, 0.6% is uniquely covered by SPAdes, and 0.4% is
 248 uniquely covered by IDBA.

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

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

Assembly	% covered
MEGAHIT	93.8%
IDBA	89.5%
SPAdes	87.3%

250 When counting only the best (longest) alignment per contig at a 99%
 251 identity threshold, each of the three assemblies recovers more than 87.3% of
 252 the reference, with MEGAHIT recovering the most – 93.8% of the reference
 253 (Table 7).

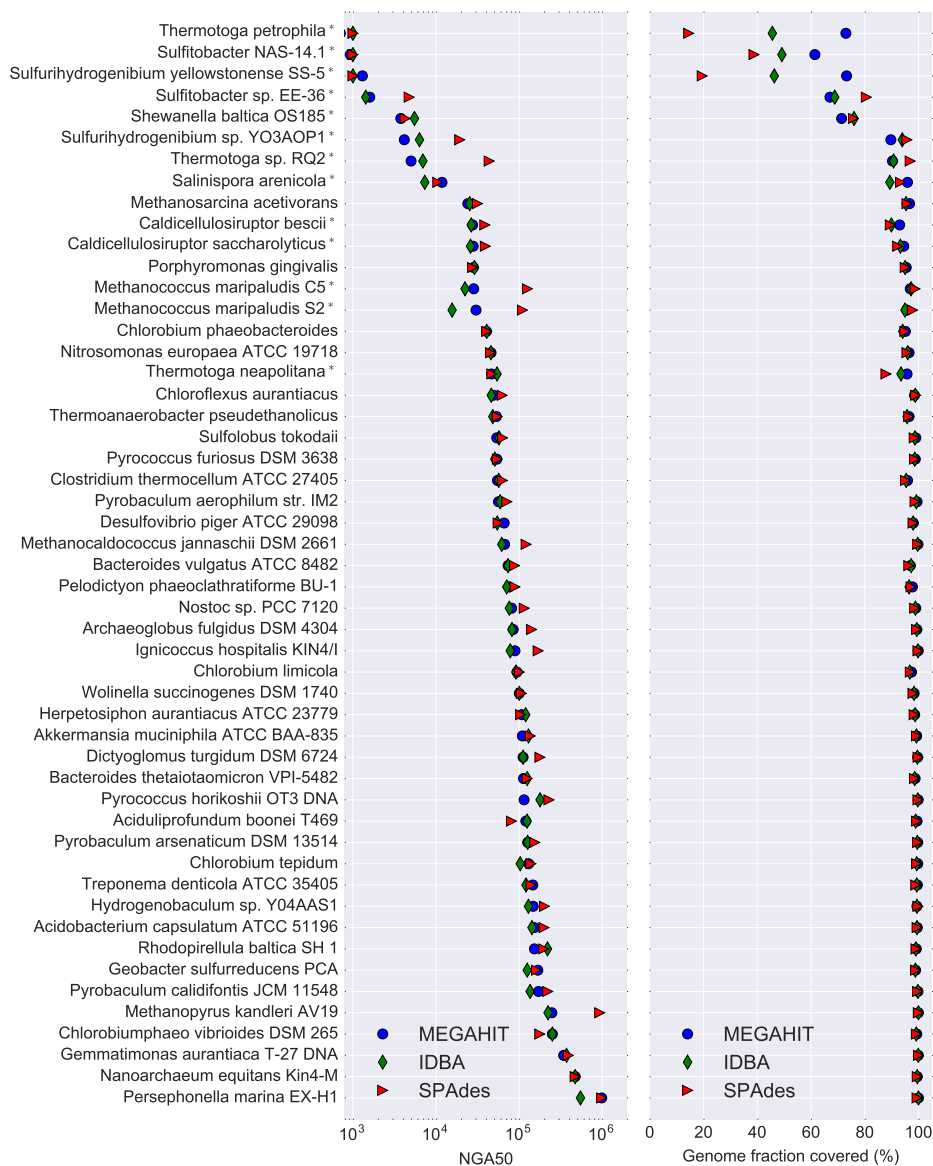


Figure 2: NGA50 by genome and assembler. A '*' after the name indicates the presence of at least one other genome with > 2% Jaccard similarity at k=31 in the community.

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

255 We computed the NGA50 for each individual genome and assembly in order
 256 to compare assembler performance on genome recovery (see left panel of Fig-
 257 ure 2). The NGA50 statistics for individual genomes vary widely, but there
 258 are consistent assembler-specific trends: IDBA yields the lowest NGA50 for
 259 28 of the 51 genomes, while SPAdes yields the highest NGA50 for 32 of the
 260 51 genomes.

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

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

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

Table 8: Chimeric contigs by contig length.

Assembly	> 50kb	> 5kb	> 500 bp
IDBA	0	1	7
MEGAHIT	1	4	14
SPAdes	0	3	30

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

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

280 Approximately 4.8 million reads (4.4%) from the QC data set did not map
 281 anywhere in the reference provided by the authors of [8]. We extracted and

Table 9: Genbank genomes detected in assembly of unmapped reads

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

282 assembled these reads in isolation using MEGAHIT, yielding 6.5 Mbp of
 283 assembly in 1711 contigs > 500bp in length. We then did a k-mer inclusion
 284 analysis of this assembly against all of the Genbank genomes at k=31, and
 285 estimated the fraction of the k-mers that belonged to different species (Ta-
 286 ble 9). We find that 51.1% of the k-mer content of these contigs positively
 287 match to a genome present in Genbank but not in the reference metagenome.

288 To verify these assignments, we aligned the MEGAHIT assembly of un-
 289 mapped reads to the Genbank genomes in Table 9 with nucmer using “loose”
 290 alignment criteria. We found that 1.78 Mbp of the contigs aligned at 99%
 291 identity or better to these Genbank genomes. We also confirmed that, as
 292 expected, there are no matches in this assembly to the full updated reference
 293 metagenome.

294 We note that all but the two *P. ruminis* matches and the *E. coli* isolate
 295 YS are strain variants of species that are part of the defined community
 296 but are not completely present in the reads (see Table 2). For *Proteiniclas-*
 297 *ticum ruminis*, there is no closely related species in the mock community
 298 design, and very little of the MEGAHIT assembly aligns to known *P. ru-*
 299 *minis* genomes at 99%. However, there are many alignments to *P. ruminis*
 300 at 94% or higher, for approximately 2.73 Mbp total. This suggests that the
 301 unmapped reads contain at least some data from a novel species of *Proteini-*
 302 *clasticum*.

303 Discussion

304 Assembly recovers basic content sensitively and accurately.

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

313 The presence of multiple closely related genomes confounds 314 assembly.

315 As reported by CAMI, we also find that the presence of closely related
316 genomes in the metagenome causes many assembly problems. This is clearly
317 shown by Figure 2, where 12 of the bottom 14 genomes by NGA50 (left
318 panel) also exhibit poor genome recovery by assembly (right panel). Inter-
319 estingly, different assemblers handle this quite differently, with e.g. SPAdes
320 failing to recover essentially any of *Thermotoga petrophila*, while MEGAHIT
321 recovers 73%. The presence of nearby genomes is an almost perfect predic-
322 tor that one or more assembler will fail to recover 5% or more - of the 13/51
323 genomes for which less than 95% is recovered, 12 of them have close genomes
324 in the community. Interestingly, very little similarity is needed - all genomes
325 with Jaccard similarity of 2% or higher at k=31 exhibited these problems.

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

332 Another interest of this study was to examine cross-species chimeric
333 assembly, in which a single contig is formed from multiple genomes. In
334 Table 8, we show that there is relatively little cross-species chimerism.

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

336 MEGAHIT is clearly the most efficient computationally, outperforming both
337 SPAdes and IDBA by 5-10x in memory and 17-42x in time (Table 4). The
338 MEGAHIT assembly also included more of the reads than either IDBA or
339 SPAdes, and omitted only 0.4% more of the unique 51-mers from the reads
340 than IDBA. MEGAHIT covered more of the reference genome with both
341 loose and strict alignments (Table 6 and Table 7), with little duplication.
342 This is clearly because of MEGAHIT's superior performance in recovering
343 the genomes of closely related strains (Figure 2, right panel).

344 Between the assemblers, the assembly content differs by only a small
345 amount when loose alignments are allowed: all three assemblers miss more
346 content (approximately 1.8% of the reference) than they generate uniquely
347 (0.9% or less). In addition to preferring no one assembler over any other,
348 this suggests that combining assemblies may have little value in terms of
349 recovering additional metagenome content.

350 **The missing reference may be present in strain variants of the**
351 **intended species.**

352 Several individual genomes are missing in measurable portion from the QC
353 reads (Table 2), and many QC reads (4.4% of 108m) did not map to the
354 full reference metagenome. These appear to be related issues: upon analysis
355 of the unmapped reads against Genbank, we find that many of the contigs
356 assembled from the unmapped reads can be assigned to strain variants of
357 the species in the mock community (Table 9). This suggests that the con-
358 structors of the mock community may have unintentionally included strain
359 variants of *Fusobacterium nucleatum*, *Thermus thermophilus* HB27, and *En-*
360 *terococcus faecalis*. In addition, we detect what may be portions of a novel
361 member of the *Proteiniclasticum* genus in the assembly of these reads.

362 Without returning to the original DNA samples, it is impossible to con-
363 clusively confirm that unintended strains were used in the construction
364 of the mock community. In particular, our analysis is dependent on the
365 genomes in Genbank: the genomes we detect in the contigs are clearly more
366 closely related to Genbank genomes other than the species in the reference
367 metagenome, based on k-mer analysis and contig alignment. However, Gen-
368 bank is unlikely to contain the exact genomes of the included strain variants,
369 rendering conclusive identification impossible.

370 Conclusions

371 Overall, assembly of this mock community works well, with good recovery of
372 known genomic sequence for the majority of genomes. All three assemblers
373 that we evaluated recover similar amounts of most genomic sequence, but
374 (recapitulating several other studies) MEGAHIT is computationally most
375 efficient.

376 The presence of closely related strains is a major confounder of metagenome
377 assembly, and causes assemblers to drop considerable portions of genomes
378 that (based on read mapping and k-mer inclusion) are clearly present. In
379 this relatively simple community, this strain confusion is present but does
380 not dominate the assembly. However, real microbial communities are likely
381 to have many closely related strains and any resulting loss of assembly will
382 be hard to detect in the absence of good reference genomes. While high
383 polymorphism rates in e.g. animal genomes is known to cause duplication
384 or loss of assembly, some solutions have emerged that make use of assump-
385 tions of uniform coverage and diploidy. These solutions cannot however
386 be transferred directly to metagenomes, which have unknown abundance
387 distributions and strain content.

388 An additional concern is that metagenome assemblies are often per-
389 formed after pooling data sets to increase coverage; this pooled data is more
390 likely to contain multiple strains, which would then in turn adversely af-
391 fect assembly of strains. This may not be resolvable within the current
392 paradigm of assembly, which focuses on outputting linear assemblies that
393 cannot properly represent strain variation.

394 Long read sequencing (and related technologies) may help resolve strain
395 variants in the future, but even with highly accurate long-read sequencing,
396 sequencing depth is still too low to resolve deep metagenomes [25]. It is
397 unclear how well long error-prone reads (such as those output by Pacific
398 Biosciences SMRT and Oxford Nanopore instruments) will perform on com-
399 plex metagenomes; with high error rates, deep coverage of each individual
400 genome is required to achieve accurate assembly, and this may not be easily
401 obtainable for complex communities. Single-molecule barcoding (e.g. 10X
402 Genomics) and HiC approaches may work better but these remain untested
403 on well-defined communities.

404 Author contributions

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

407 Competing interests

408 No competing interest to our knowledge.

409 Grant information

410 This work is funded by Moore and NIH.

411 Acknowledgments

412 We thank Michael R. Crusoe and Phillip T. Brooks for input on analysis
413 and pipeline development.

414 References

- 415 [1] Nikos C. Kyrpides, Philip Hugenholtz, Jonathan A. Eisen, Tanja Woyke,
416 Markus Göker, Charles T. Parker, Rudolf Amann, Brian J. Beck, Patrick S. G.
417 Chain, Jongsik Chun, Rita R. Colwell, Antoine Danchin, Peter Dawyndt, Tom
418 Dedeurwaerdere, Edward F. DeLong, John C. Detter, Paul De Vos, Timothy J.
419 Donohue, Xiu-Zhu Dong, Dusko S. Ehrlich, Claire Fraser, Richard Gibbs, Jack
420 Gilbert, Paul Gilna, Frank Oliver Glöckner, Janet K. Jansson, Jay D. Keasling,
421 Rob Knight, David Labeda, Alla Lapidus, Jung-Sook Lee, Wen-Jun Li, Juncai
422 MA, Victor Markowitz, Edward R. B. Moore, Mark Morrison, Folker Meyer,
423 Karen E. Nelson, Moriya Ohkuma, Christos A. Ouzounis, Norman Pace, Julian
424 Parkhill, Nan Qin, Ramon Rossello-Mora, Johannes Sikorski, David Smith,
425 Mitch Sogin, Rick Stevens, Uli Stingl, Ken ichiro Suzuki, Dorothea Taylor,
426 Jim M. Tiedje, Brian Tindall, Michael Wagner, George Weinstock, Jean Weis-
427 senbach, Owen White, Jun Wang, Lixin Zhang, Yu-Guang Zhou, Dawn Field,
428 William B. Whitman, George M. Garrity, and Hans-Peter Klenk. Genomic
429 encyclopedia of bacteria and archaea: Sequencing a myriad of type strains.
430 *PLoS Biology*, 12(8):e1001920, aug 2014. doi: 10.1371/journal.pbio.1001920.
431 URL <https://doi.org/10.1371/journal.pbio.1001920>.
- 432 [2] Jorge F Vázquez-Castellanos, Rodrigo García-López, Vicente Pérez-Brocal,
433 Miguel Pignatelli, and Andrés Moya. Comparison of different assembly and
434 annotation tools on analysis of simulated viral metagenomic communities in
435 the gut. *BMC genomics*, 15(1):1, 2014.

- 436 [3] Konstantinos Mavromatis, Natalia Ivanova, Kerrie Barry, Harris Shapiro, Eu-
437 gene Goltsman, Alice C McHardy, Isidore Rigoutsos, Asaf Salamov, Frank
438 Korzeniewski, Miriam Land, et al. Use of simulated data sets to evaluate the
439 fidelity of metagenomic processing methods. *Nature methods*, 4(6):495–500,
440 2007.
- 441 [4] David B Jaffe, Jonathan Butler, Sante Gnerre, Evan Mauceli, Kerstin
442 Lindblad-Toh, Jill P Mesirov, Michael C Zody, and Eric S Lander. Whole-
443 genome sequence assembly for mammalian genomes: Arachne 2. *Genome*
444 *research*, 13(1):91–96, 2003.
- 445 [5] Samuel Aparicio, Jarrod Chapman, Elia Stupka, Nik Putnam, Jer-ming Chia,
446 Paramvir Dehal, Alan Christoffels, Sam Rash, Shawn Hoon, Arian Smit, et al.
447 Whole-genome shotgun assembly and analysis of the genome of *fugu rubripes*.
448 *Science*, 297(5585):1301–1310, 2002.
- 449 [6] Anveshi Charuvaka and Huzefa Rangwala. Evaluation of short read metage-
450 nomic assembly. *BMC genomics*, 12(2):1, 2011.
- 451 [7] Jared T Simpson, Kim Wong, Shaun D Jackman, Jacqueline E Schein,
452 Steven JM Jones, and Inanç Birol. Abyss: a parallel assembler for short read
453 sequence data. *Genome research*, 19(6):1117–1123, 2009.
- 454 [8] Shakya Migun, Christopher Quince, James Campbell, Zamin Yang, Christo-
455 pher Schadt, and Mircea Podar. Comparative metagenomic and rrna microbial
456 diversity characterization using archaeal and bacterial synthetic communities.
457 *Enivromental Microbiology*, 15(6):1882–1899, 2013.
- 458 [9] Alexander Sczyrba, Peter Hofmann, Peter Belmann, David Koslicki, Stefan
459 Janssen, Johannes Droege, Ivan Gregor, Stephan Majda, Jessika Fiedler,
460 Eik Dahms, Andreas Bremges, Adrian Fritz, Ruben Garrido-Oter, Tue
461 Sparholt Jorgensen, Nicole Shapiro, Philip D Blood, Alexey Gurevich, Yang
462 Bai, Dmitrij Turaev, Matthew Z DeMaere, Rayan Chikhi, Niranjan Nagara-
463 jan, Christopher Quince, Lars Hestbjerg Hansen, Soren J Sorensen, Burton
464 K H Chia, Bertrand Denis, Jeff L Froula, Zhong Wang, Robert Egan, Dong-
465 wan Don Kang, Jeffrey J Cook, Charles Deltel, Michael Beckstette, Claire
466 Lemaitre, Pierre Peterlongo, Guillaume Rizk, Dominique Lavenier, Yu-Wei
467 Wu, Steven W Singer, Chirag Jain, Marc Strous, Heiner Klingenberg, Peter
468 Meinicke, Michael Barton, Thomas Lingner, Hsin-Hung Lin, Yu-Chieh Liao,
469 Genivaldo Gueiros Z. Silva, Daniel A Cuevas, Robert A Edwards, Surya Saha,
470 Vitor C Piro, Bernhard Y Renard, Mihai Pop, Hans-Peter Klenk, Markus
471 Goeker, Nikos Kyrpides, Tanja Woyke, Julia A Vorholt, Paul Schulze-Lefert,
472 Edward M Rubin, Aaron E Darling, Thomas Rattei, and Alice C McHardy.
473 Critical assessment of metagenome interpretation - a benchmark of compu-
474 tational metagenomics software. *bioRxiv*, 2017. doi: 10.1101/099127. URL
475 <http://biorxiv.org/content/early/2017/01/09/099127>.

- [10] Andries Johannes van der Walt, Marc Warwick Van Goethem, Jean-Baptiste Ramond, Thulani Peter Makhalanyane, Oleg Reva, and Don Arthur Cowan. Assembling metagenomes, one community at a time. *bioRxiv*, 2017. doi: 10.1101/120154. URL <http://biorxiv.org/content/early/2017/06/06/120154>.
- [11] Yu Peng, Henry C.M. Leung, S.M. Yiu, and Francis Y.L. Chin. Idba-ud: a de novo assembler for single-cell and metagenomic sequencing data with highly uneven depth. *Bioinformatics*, 28:1420–1428, 2012.
- [12] Sergey Nurk, Dmitry Meleshko, Anton Korobeynikov, and Pavel A. Pevzner. metaSPAdes: a new versatile metagenomic assembler. *Genome Research*, 27(5):824–834, mar 2017. doi: 10.1101/gr.213959.116. URL <https://doi.org/10.1101/gr.213959.116>.
- [13] Dinghua Li, Ruibang Luo, Chi-Man Liu, Chi-Ming Leung, Hing-Fung Ting, Kunihiro Sadakane, Hiroshi Yamashita, and Tak-Wah Lam. Megahit v1. 0: A fast and scalable metagenome assembler driven by advanced methodologies and community practices. *Methods*, 102:3–11, 2016.
- [14] H Chitsaz, JL Yee-Greenbaum, G Tesler, MJ Lombardo, CL Dupont, JH Badger, M Novotny, DB Rusch, LJ Fraser, NA Gormley, O Schulz-Trieglaff, GP Smith, DJ Evers, PA Pevzner, and RS Lasken. Efficient de novo assembly of single-cell bacterial genomes from short-read data sets. *Nat Biotechnol*, 29(10):915–21, 2011.
- [15] Anthony M. Bolger, Marc Lohse, and Bjoern Usadel. Trimmomatic: A flexible trimmer for illumina sequence data. *Bioinformatics*, 30(15):2114–2120, 2014.
- [16] Matthew D MacManes. On the optimal trimming of high-throughput mrna sequence data. *Frontiers in genetics*, 5:13, 2014.
- [17] Heng Li and Richard Durbin. Fast and accurate short read alignment with burrows–wheeler transform. *Bioinformatics*, 25(14):1754–1760, 2009.
- [18] C. Titus Brown and Luiz Irber. sourmash: a library for MinHash sketching of DNA. *The Journal of Open Source Software*, 1(5), sep 2016. doi: 10.21105/joss.00027. URL <https://doi.org/10.21105/joss.00027>.
- [19] Brian D. Ondov, Todd J. Treangen, Páll Melsted, Adam B. Mallonee, Nicholas H. Bergman, Sergey Koren, and Adam M. Phillippy. Mash: fast genome and metagenome distance estimation using MinHash. *Genome Biology*, 17(1), jun 2016. doi: 10.1186/s13059-016-0997-x. URL <https://doi.org/10.1186/s13059-016-0997-x>.
- [20] David Koslicki and Daniel Falush. Metapalette: a k-mer painting approach for metagenomic taxonomic profiling and quantification of novel strain variation. *mSystems*, 1(3), 2016. doi: 10.1128/mSystems.00020-16. URL <http://msystems.asm.org/content/1/3/e00020-16>.

- 515 [21] Zhang Qingpeng, Awad Sherine, and Brown Titus. Crossing the streams:
516 a framework for streaming analysis of short dna sequencing reads. *PeerJ*
517 *PrePrints* 3:e1100 <https://dx.doi.org/10.7287/peerj.preprints.890v1>, 2015.
- 518 [22] MR Crusoe, HF Alameldin, S Awad, E Boucher, A Caldwell, R Cartwright,
519 A Charbonneau, B Constantinides, G Edverson, S Fay, J Fenton, T Fenzl,
520 J Fish, L Garcia-Gutierrez, P Garland, J Gluck, I Gonzlez, S Guermond,
521 J Guo, A Gupta, JR Herr, A Howe, A Hyer, A Hrpfer, L Irber, R Kidd,
522 D Lin, J Lippi, T Mansour, P McA’Nulty, E McDonald, J Mizzi, KD Mur-
523 ray, JR Nahum, K Nanlohy, AJ Nederbragt, H Ortiz-Zuazaga, J Ory, J Pell,
524 C Pepe-Ramney, ZN Russ, E Schwarz, C Scott, J Seaman, S Sievert, J Simp-
525 son, CT Skennerton, J Spencer, R Srinivasan, D Standage, JA Stapleton,
526 SR Steinman, J Stein, B Taylor, W Trimble, HL Wiencko, M Wright,
527 B Wyss, Q Zhang, e zyme, and CT Brown. The khmer software pack-
528 age: enabling efficient nucleotide sequence analysis [version 1; referees: 2 ap-
529 proved, 1 approved with reservations]. *F1000Research*, 4(900), 2015. doi:
530 10.12688/f1000research.6924.1.
- 531 [23] Heng Li, Bob Handsaker, Alec Wysoker, Tim Fennell, Jue Ruan, Nils Homer,
532 Gabor Marth, Goncalo Abecasis, Richard Durbin, et al. The sequence align-
533 ment/map format and samtools. *Bioinformatics*, 25(16):2078–2079, 2009.
- 534 [24] Stefan Kurtz, Adam Phillippy, Arthur L Delcher, Michael Smoot, Martin
535 Shumway, Corina Antonescu, and Steven L Salzberg. Versatile and open soft-
536 ware for comparing large genomes. *Genome biology*, 5(2):1, 2004.
- 537 [25] Itai Sharon, Michael Kertesz, Laura A. Hug, Dmitry Pushkarev, Timothy A.
538 Blauwkamp, Cindy J. Castelle, Mojgan Amirebrahimi, Brian C. Thomas,
539 David Burstein, Susannah G. Tringe, Kenneth H. Williams, and Jillian F.
540 Banfield. Accurate, multi-kb reads resolve complex populations and de-
541 tect rare microorganisms. *Genome Research*, 25(4):534–543, feb 2015. doi:
542 10.1101/gr.183012.114. URL <https://doi.org/10.1101/gr.183012.114>.