

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**
5 University of California, Davis
6 Davis, CA 95616 USA
7 * E-mail: ctbrown@ucdavis.edu

8 June 15, 2017

9 **Abstract**

10 Foo!

7 Introduction

8 Metagenomics refers to sequencing of DNA from a mixture of organisms,
9 often from an environmental or uncultured sample. Unlike whole genome
10 sequencing, metagenomics targets a mixture of genomes, which introduces
11 metagenome-specific challenges in analysis. Most approaches to analyzing
12 metagenomic data rely on mapping or comparing sequencing reads to refer-
13 ence sequence collections. However, reference databases contain only a small
14 subset of microbial diversity (cite: geba), and the much of the remaining
15 diversity is evolutionarily distant and search techniques may not access it.

16 As sequencing capacity increases and sequence data is generated from
17 many more environmental samples, metagenomics is increasingly using de
18 novo assembly techniques to generate new reference genomes and metagenomes.
19 There are a number of metagenome assemblers that are widely used. How-
20 ever, evaluating the results of these assemblers is challenging due to the
21 general lack of good quality reference metagenomes. Below, we evaluate
22 three commonly assemblers - SPAdes, IDBA, and MEGAHIT - on a mock
23 community containing 64 species of microbes with known genomes.

24 Moya et al. in [1] evaluated metagenome assembly using two simulated
25 454 viral metagenome and six assemblers. The assemblies were evaluated
26 based on several metrics including N50, percentages of reads assembled, ac-
27 curacy when compared to the reference genome. In addition to, chimeras per
28 contigs and the effect of assembly on taxonomic and functional annotations.

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

35 Rangwala et al. in [5] presented an evaluation study of metagenome
36 assembly. The study used a de Bruijn graph based assembler ABYSS [6] to
37 assemble simulated metagnome reads of 36 bp. The data set is classified at
38 different complexity levels. The study compares the quality of the assembly
39 of the data sets in terms of quality measures of contigs length, assembly
40 accuracy. The study also took into consideration the effect of kmer size and
41 the degree of chimericity. However, the study evaluated the assembly based
42 on one assembler, and did not evaluate assembly against several assemblers.
43 Also, both previous studies used simulated data, which may lack confounders
44 of assembly such as sequencing artifacts and GC bias.

45 Shakya et al. (2013) constructed a synthetic community of organisms by
46 mixing DNA isolated from individual cultures of 64 bacteria and archaea,
47 including a variety of strains across a range of nucleotide distances. In
48 addition to performing 16s amplicon analysis and doing 454 sequencing, the
49 authors shotgun sequenced the mixture with Illumina (@cite). While the
50 authors concluded that this metagenomic sequencing generally outperformed
51 amplicon sequencing, they did not conduct an assembly based analysis. a
52 mapping based analysis rather than an assembly based analysis.

53 More recently, several benchmark studies systematically evaluated metagenome
54 assembly of short reads. The Critical Assessment of Metagenome Interpretation (CAMI) collaboration benchmarked a number of metagenome assemblers on several data sets of varying complexity, evaluating recovery of novel genomes and multiple strain variants (@cite). Notably, CAMI concluded
58 that “The resolution of strain-level diversity represents a substantial challenge to all evaluated programs.” Another recent study evaluated eight
59 assemblers on nine environmental metagenomes and three simulated data
60 sets (@cite).
61

62 In this paper, we evaluate metagenome assembly on the Illumina data
63 set from Shakya et al. (2013) using three assemblers; IDBA-UD [7], SPAdes
64 [8], and MEGAHIT [9]. These three assemblers were chosen because they
65 are actively used and highly cited, and typically perform well.

66 Below, we evaluate the performance of these three assemblers using the
67 mock community data from the Shakya et al. study. The performance of
68 each assembler is compared in terms of resource utilization, covered genome
69 fraction, duplication ratio, gene recovery, contig misassembly, and contig
70 length.

71 In this report, we extend the CAMI study by delving into questions of
72 chimeric misassembly and strain recovery. First, we update the list of reference genomes for Shakya et al. to include updated assemblies as well as plasmids. We then compare IDBA, SPAdes, and MEGAHIT performance
75 on assembling this short-read data set, and explore concordance in recovery
76 between the three assemblers. We also evaluate inter-strain chimerism in
77 the assemblies and explore the poor assemblies caused by “strain confusion”
78 between two *Shewanella baltica* strain. We detect and analyze several previously unreported strains and genomes in the Shakya et al. data set. Our
80 report provides strong guidance on choice of assemblers and significantly extends previous analyses of this low-complexity metagenome benchmarking
81 data set.
82

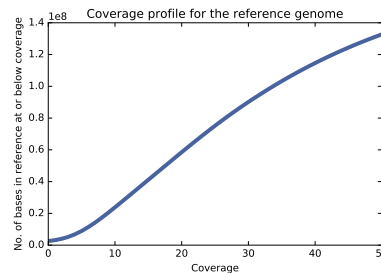


Figure 1: Cumulative coverage profile for the reference metagenome, based on read mapping.

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 [10].

89 The original reads are available through the NCBI Sequence Read Archive
 90 at Accession SRX200676. We received the 64 reference genomes from the
 91 original authors. They consist of 205.6 Mbp of assembled genomes in 64 con-
 92 tigs, and are available for download at <https://dx.doi.org/10.6084/m9.figshare.1506873.v2>.

93 We updated the data sets from NCBI etc. etc. The following genomes
 94 were updated. Updated data is available for download here (OSF).

95 Methods

96 The analysis code and run scripts for this paper are available at: [https://github.com/dib-](https://github.com/dib-lab/2015-metagenome-assembly/)
 97 [lab/2015-metagenome-assembly/](https://github.com/dib-lab/2015-metagenome-assembly/). The scripts and overall pipeline were ex-
 98 amined by the first and senior authors for correctness. In addition, the
 99 bespoke reference-based analysis scripts were tested by running them on a
 100 single-colony *E. coli* MG1655 data set with a high quality reference genome
 101 [11].

102 Quality Filtering

103 We removed adapters with Trimmomatic v0.30 in paired-end mode with the
 104 Truseq adapters [12], using light quality score trimming as recommended in

105 MacManes, 2014 [13].

106 **Reference Coverage Profile**

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

111 **Assemblers**

112 We assembled the quality-filtered reads using three different assemblers:
113 IDBA-UD [7], MetaSPAdes [8], and MEGAHIT [9]. For IDBA-UD v1.1.1
114 [7], we used `--pre_correction` to perform pre-correction before assembly
115 and `-r` for the pe files.

116 For MetaSPAdes v3.9.0 [8], we used `--meta --pe1-12 --pe1-s` where
117 `--meta` is used for metagenomic data sets, `--pe1-12` specifies the interlaced
118 reads for the first paired-end library, and `--pe1-s` provides the orphan reads
119 remaining from quality trimming.

120 For MEGAHIT v1.1.1-2-g02102e1 [9], we used `-l 101 -m 3e9 --cpu-only`
121 where `-l` is for maximum read length, `-m` is for max memory in bytes to
122 be used in constructing the graph, and `--cpu-only` to use only the CPU
123 and no GPUs. We also used `--presets meta-large` for large and complex
124 metagenomes, and `--12` and `-r` to specify the interleaved-paired-end and
125 single-end files respectively. MEGAHIT allows the specification of a memory
126 limit and we used `-M 1e+10` for 10 GB.

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

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

133 **Mapping**

134 We aligned all quality-filtered reads to the reference metagenome with `bwa`
135 `aln` (v0.7.7.r441) [14]. We aligned paired-end and orphaned reads separately.
136 We then used `samtools` (v0.1.19) [15] to convert SAM files to BAM files for

137 both paired-end and orphaned reads. To count the unaligned reads, we
138 included only those records with the “4” flag in the SAM files [15].

139 To extract the reads that contribute to unaligned contigs, we mapped
140 the quality filtered reads to the unaligned contigs using `bwa aln` (v0.7.7.r441)
141 [14]. Then we used `samtools` to retrieve the reads that mapped to the un-
142 aligned contigs.

143 **k-mer Presence**

144 In order to examine k-mer presence for a k-mer size of 20, we built a k-mer
145 counting table from the given quality filtered reads using `load-into-counting.py`
146 from `khmer` [?]. Then we calculate abundance distribution of the k-mers
147 in the quality filtered reads using the pre-made k-mer counting table us-
148 ing `abundance-dist.py`. We followed the same approach to examine k-mer
149 presence in assemblies.

150 **Assembly analysis using Nucmer**

151 We used the NUCmer tool from MUMmer3.23 [16] to align assemblies to the
152 reference genome with options `-coords -p`. Then we parsed the generated
153 “.coords” file using a custom script `analyze_assembly.py`, and calculated
154 several analysis metrics across all three assemblies at two alignment identi-
155 ties, 95% and 99%.

156 **Reference-based analysis of the assemblies**

157 We analyzed metrics for three different sets of contigs, based on the NUCmer
158 alignments. We used the unfiltered NUCmer alignments for the analyses
159 termed “ambiguous.” We also subjected the alignments to two different fil-
160 tering criteria, “best-hit” and “no-misassemblies.” In the best-hit approach,
161 among all alignments of a contig, we took into consideration the longest
162 alignment with an identity above a specified identity threshold (either 95%
163 or 99%). In the no-misassemblies approach, we only counted contigs that
164 have precisely one alignment within the reference.

165 In all approaches, we flag a base in the reference genome as “covered” if
166 it is contained in a kept alignment. We define the duplication ratio as the
167 percentages of bases in the reference covered by two or more kept alignments.
168 We define misassemblies as those contigs that are divided into different parts
169 when mapped to the reference. The number of misassembled contigs is equal

170 to the number of aligned contigs (both totally and partially) in the ambigu-
171 ous approach, minus the number of aligned contigs in the no-misassemblies
172 approach.

173 All approaches have a non-zero duplication ratio within the reference
174 because we do not explicitly discard contigs that map to the same location
175 in the reference.

176 Analysis of chimeric misassemblies

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

186 Results

187 The raw data is high quality.

188 The reads contains 11,072,579,096 bp (11.07 Gbp) in 109,629,496 reads with
189 101.0 average length (2x101bp Illumina HiSeq).

190 Trimming removed 686,735 reads (0.63%). After trimming, we retained
191 108,422,358 paired reads containing 10.94 Gbp with an average length of
192 100.9 bases. A total of 46.56 Mbp remained in 520,403 orphan reads with
193 an average length of 89.5 bases. In total, the quality trimmed data contained
194 10.98 Gbp in 108,942,761 reads. This quality trimmed ("QC") data set was
195 used as the basis for all further analyses.

196 The reference metagenome is not completely present in the 197 reads.

198 We next evaluated the fraction of the reference genome covered by at least
199 one read (see Methods for details). Quality filtered reads cover 203,058,414
200 (98.76%) bases of the reference metagenome (205,603,715 bp total size). Fig-
201 ure 1 shows the cumulative coverage profile of the reference metagenome,

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%

and the percentage of bases with that coverage. Most of the reference metagenome was covered at least minimally; only 3.33% of the reference metagenome had mapping coverage <5 , and 1.24% of the bases in the reference were not covered by any reads in the QC data set.

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

Some individual reference genomes are poorly represented in 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%

To see if specific reference genomes exhibited low coverage, we analyzed read mapping coverage and 21-mer containment for individual genomes. Of the 64 reference genomes used in the metagenome, 59 had a per-base mapping coverage above 95% and a 21-mer containment in the QC reads above 95%. The remaining five varied significantly in both metrics (Table 3), with *F. nucleatum* the lowest – only 47.6% of the bases in the reference genome are covered by one or more mapped reads, and only 18.2% of the 21-mers in the *F. nucleatum* reference genome are present in the reads at

any abundance.

We next did a 51-mer containment analysis of each reference genome in the reads. 99% or more of the constituent 51-mers for 51 of the 64 reference genomes were present in the reads, suggesting that each of the 51 genomes was entirely present at some minimal coverage.

We excluded the remaining 13 genomes from any comparative analysis of assembly quality, because interpreting coverage and misassembly analysis for these genomes would be impossible. (@CTB list or table?)

MEGAHIT is the fastest and lowest-memory assembler evaluated

Table 3: 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

We ran three commonly used metagenome assemblers on the QC data set: IDBA-UD, SPAdes, and MEGAHIT. We recorded the time and memory usage of each (Table 3). In computational requirements, MEGAHIT outperformed both SPAdes and IDBA-UD considerably, producing an assembly 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 4: 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	3,879,573 (3.56%)	3.2%
MEGAHIT	5,848,494 (5.37%)	2.8%

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

242 Depending on the assembly, between 3.3 million and 5.9 million reads (3.0-
 243 5.4%) did not map to the assemblies (Table 4). Here, the MEGAHIT as-
 244 sembly was distinguished by representing 2 million fewer reads than the
 245 IDBA and SPAdes assemblies. K-mer inclusion, however, was more closely
 246 matched across the assemblies, with all three assemblies containing the large
 247 majority of high-abundance 51-mers.

248 Much of the reference is covered by the assemblies.

Table 5: 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%

249 We next evaluated the extent to which the assembled contigs recovered
 250 the “known/true” metagenome sequence by aligning each assembly to the
 251 adjusted reference (Table 5). Each of the three assemblers generates contigs
 252 that cover more than 95.6% of the reference metagenome at high identity
 253 (99%) with little duplication (0.72-0.99%). All three assemblies contain
 254 between 96.2% and 97.2% of the 51-mers in the reference.

255 At 99% identity with the loose mapping approach, approximately 1.8%
 256 of the reference is missed by all three assemblers, while 0.9% is uniquely
 257 covered by MEGAHIT, 0.6% is uniquely covered by SPAdes, and 0.4% is
 258 uniquely covered by IDBA.

259 The generated contigs are broadly accurate.

Table 6: Contig accuracy measured by reference coverage with strict alignment.

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

260 When counting only the best (longest) alignment per contig at a 99%
 261 identity threshold, each of the three assemblies recovers more than 87.3% of

262 the reference, with MEGAHIT recovering the most – 93.8% of the reference
263 (Table 6).

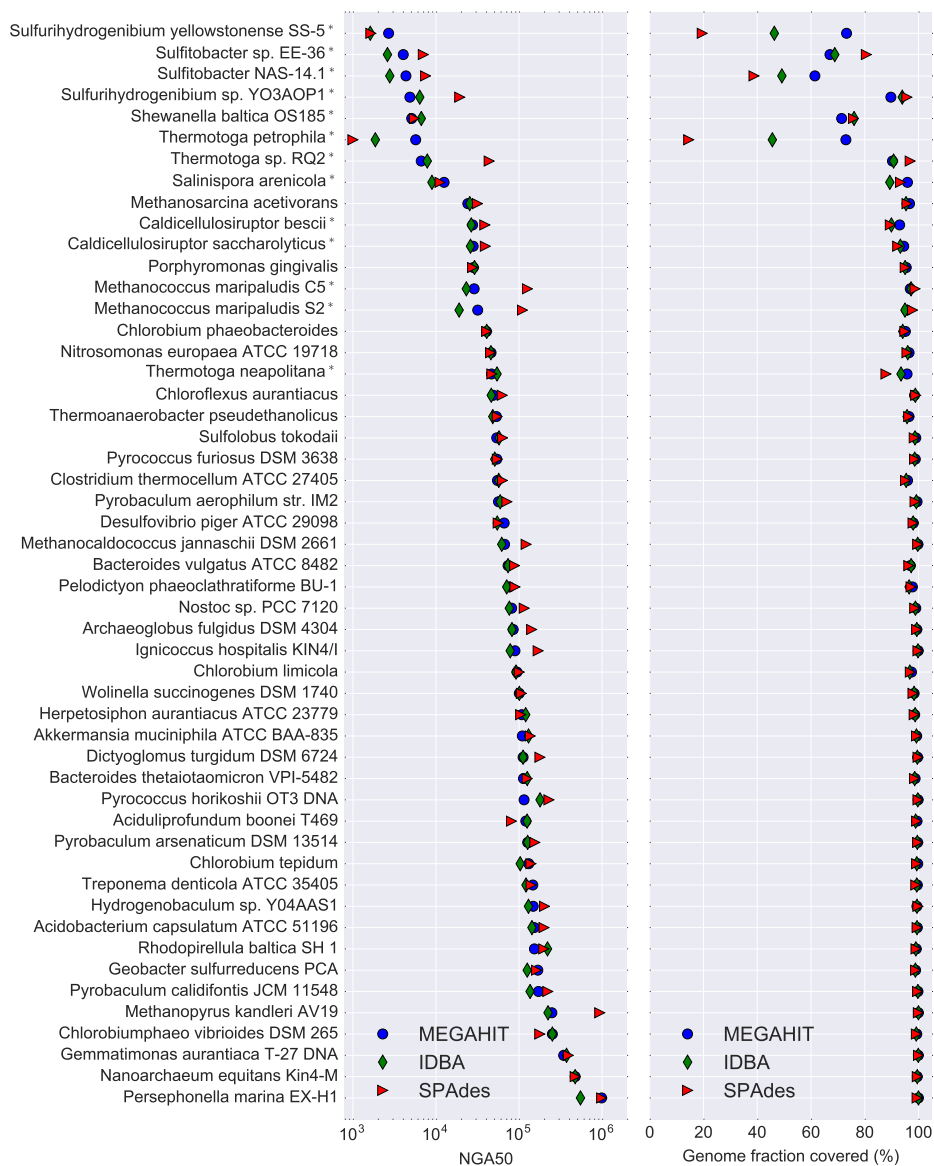


Figure 2: NGA50 by genome and assembler. A '*' after the name indicates the presence of at least one other strain variant in the community.

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

265 We computed the NGA50 for each individual genome and assembly in order
 266 to compare assembler performance on each genome (see left panel of Fig-
 267 ure 2). The NGA50 statistics for individual genomes vary widely, but there
 268 are consistent assembler-specific trends: IDBA yields the lowest NGA50 for
 269 30 of the 51 genomes, while SPAdes yields the highest NGA50 for 33 of the
 270 51 genomes. Genomes with close strain variants in the defined community
 271 (indicated by a '*' suffix in the genome name) tended strongly towards lower
 272 NGA50s.

273 We also evaluated aligned coverage per genome for each of the three
 274 assemblies (right panel, Figure 2. We found that a 13 of the 51 genomes
 275 were missing 5% or more of bases in at least one assembly, despite all 51
 276 genomes having 99% or higher read- and 51-mer coverage. 12 of the 13
 277 genomes missing 5% or more of their content also had at least one close
 278 strain variant in the defined community.

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

Table 7: Chimeric contigs by contig length.

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

280 Chimerism is the formation of contigs that include sequence from multi-
 281 ple genomes. We evaluated the rate of chimerism in contigs at three different
 282 contig length cutoffs: 500bp, 5kb, and 50kb (Table 7). We found that the
 283 percentage of contigs that match to the genomes of two or more different
 284 species drop as the minimum contig size increases, to the point where only
 285 the MEGAHIT assembly had a single chimeric contig longer than 50kb.
 286 Overall, chimeric misassemblies were rare, with no assembler generating
 287 more than 30 chimeric contigs out of thousands of total contigs.

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

289 Approximately 4.8 million reads (4.4%) from the QC data set did not map
 290 anywhere in the reference provided by the authors of @cite. We extracted
 291 and assembled these reads in isolation using MEGAHIT, yielding 6.5 Mbp of

Table 8: 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>

assembly in 1711 contigs > 500bp in length. We then did a k-mer analysis of this assembly against all of the Genbank genomes at k=31 using sourmash @cite, and estimated the fraction of the k-mers that belonged to different species (Table 8). We find that 51.1% of the k-mer content of these contigs positively match to a genome present in Genbank but not in the reference metagenome.

To verify these assignments, we aligned the MEGAHIT assembly of unmapped reads to the Genbank genomes in Table 8 with nucmer using “loose” alignment criteria. We found that 1.78 Mbp of the contigs aligned at 99% identity or better to the Genbank genomes. We also confirmed that, as expected, there are no matches to the reference metagenome.

Interestingly, all but the two *P. ruminis* matches and the *E. coli* isolate YS are strain variants of species that are part of the defined community but are not completely present in the reads (see Table 2). The presence of so many 31-mers from the *Proteiniclasticum ruminis* species is intriguing, since there is no closely related species in the mock community design. Very little of the MEGAHIT assembly aligns to known *P. ruminis* genomes via nucmer at 99%, suggesting that this is an unknown species – there are many alignments to *P. ruminis* at 94% or higher, for approximately 2.73 Mbp total.

312 Discussion

313 Assembly recovers basic content sensitively and accurately.

314 All three assemblers performed well in assembling contigs from the content
315 that was fully present in reads and k-mers. After length filtering, all three
316 assemblies contained more than 95% of the reference (Table 5); even with
317 removal of secondary alignments, more than 87% was recovered by each
318 assembler (Table 6). About half the constituent genomes had an NGA50
319 of 50kb or higher (Figure 2), which, while low for current Illumina single-
320 genome sequencing @cite, is sufficient to recover operon-level relationships
321 for many genes.

322 The presence of multiple strain variants confounds assembly.

323 As reported elsewhere, we also find that strain variation causes many as-
324 sembly problems. This is clearly shown by Figure 2, where 12 of the bottom
325 14 genomes by NGA50 (left panel) also exhibit poor genome recovery by
326 assembly (right panel). Interestingly, different assemblers handle this quite
327 differently, with e.g. SPAdes failing to recover essentially any of *Thermotoga*
328 *petrophila*, while MEGAHIT recovers 73%. The presence of strain variation
329 is an almost perfect predictor that one or more assembler will fail to recover
330 5% or more - of the 13/51 genomes for which less than 95% is recovered, 12
331 of them have close strain variants in the community.

332 The *Shewanella baltica* OS185 genome is a case in point: there are two
333 strain variants, OS185 and OS223, present in the defined community. Both
334 are present at more than 99% in the reads, and more than 98% in 51-mers,
335 but only 75% of *S. baltica* OS185 and 50% of *S. baltica* OS223 are recovered
336 by assemblers. This is a clear case of “strain confusion” where the assemblers
337 simply fail to output contigs for a substantial portion of the two genomes.

338 Another interest of this study was to examine cross-species chimeric
339 assembly, in which a single contig is formed from multiple genomes. In
340 Table 7, we show that there is relatively little cross-species chimerism and
341 that it always occurs between members of the same genus.

342 MEGAHIT performs best by several metrics.

343 MEGAHIT is clearly the most efficient computationally, outperforming both
344 SPAdes and IDBA by 5-10x in memory and 17-42x in time (Table 3). While
345 the MEGAHIT assembly had 2m fewer reads mapping to it than the other

assemblies (Table 1), MEGAHIT covered more of the reference genome with both loose and strict alignments (Table 5 and Table 6), with little duplication. This is clearly because of MEGAHIT’s superior performance in recovering the genomes of closely related strains (Figure 2, right panel).

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

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

Several individual genomes are missing in measurable portion from the QC reads (Table 2), and many QC reads (4.4% of 108m) did not map to the full reference metagenome. These appear to be related issues: upon analysis of the unmapped reads against Genbank, we find that many of the contigs assembled from the unmapped reads can be assigned to strain variants of the species in the mock community (Table 8). This suggests that the constructors of the mock community may have unintentionally included strain variants of *Fusobacterium nucleatum*, *Thermus thermophilus* HB27, and *Enterococcus faecalis*. In addition, we detect what may be portions of a novel member of the *Proteiniclasticum* genus in the assembly of these reads.

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

Conclusions

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

381 efficient.

382 The presence of closely related strains is a major confounder of metagenome
383 assembly, and causes assemblers to drop considerable portions of genomes
384 that (from a read and k-mer perspective) are present. In this relatively
385 simple community, this strain confusion is clearly present but not domi-
386 nant. However, real microbial communities are likely to have many closely
387 related strains and any resulting loss of assembly will be hard to detect in
388 the absence of good reference genomes. While high polymorphism rates in
389 e.g. animal genomes is known to cause duplication or loss of assembly, some
390 solutions have emerged that make use of assumptions of uniform coverage
391 and diploidy. These solutions cannot however be transferred directly to
392 metagenomes, which have unknown abundance distributions and numbers
393 of strains.

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

400 It is unclear how well long error-prone reads (such as those output by
401 Pacific Biosciences SMRT and Oxford Nanopore instruments) will be able
402 to resolve strain variants in metagenomes. High coverage of each individual
403 genome is required to achieve accurate assembly, which may not be easily ob-
404 tainable for complex communities. Fosmid, single-molecule barcoding, and
405 HiC approaches may work better but these remain untested on well defined
406 communities. (note phase genomics, 10x metagenomes papers.) (Mention
407 Sharon@ paper and moleculo).

408 **Author contributions**

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

411 **Competing interests**

412 No competing interest to our knowledge.

413 Grant information

414 This work is funded by Moore and NIH.

415 Acknowledgments

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

418 References

- 419 [1] Jorge F Vázquez-Castellanos, Rodrigo García-López, Vicente Pérez-Brocal,
420 Miguel Pignatelli, and Andrés Moya. Comparison of different assembly and
421 annotation tools on analysis of simulated viral metagenomic communities in
422 the gut. *BMC genomics*, 15(1):1, 2014.
- 423 [2] Konstantinos Mavromatis, Natalia Ivanova, Kerrie Barry, Harris Shapiro, Eu-
424 gene Goltsman, Alice C McHardy, Isidore Rigoutsos, Asaf Salamov, Frank
425 Korzeniewski, Miriam Land, et al. Use of simulated data sets to evaluate the
426 fidelity of metagenomic processing methods. *Nature methods*, 4(6):495–500,
427 2007.
- 428 [3] David B Jaffe, Jonathan Butler, Sante Gnerre, Evan Mauceli, Kerstin
429 Lindblad-Toh, Jill P Mesirov, Michael C Zody, and Eric S Lander. Whole-
430 genome sequence assembly for mammalian genomes: Arachne 2. *Genome*
431 *research*, 13(1):91–96, 2003.
- 432 [4] Samuel Aparicio, Jarrod Chapman, Elia Stupka, Nik Putnam, Jer-ming Chia,
433 Paramvir Dehal, Alan Christoffels, Sam Rash, Shawn Hoon, Arian Smit, et al.
434 Whole-genome shotgun assembly and analysis of the genome of *fugu rubripes*.
435 *Science*, 297(5585):1301–1310, 2002.
- 436 [5] Anveshi Charuvaka and Huzefa Rangwala. Evaluation of short read metage-
437 nomic assembly. *BMC genomics*, 12(2):1, 2011.
- 438 [6] Jared T Simpson, Kim Wong, Shaun D Jackman, Jacqueline E Schein,
439 Steven JM Jones, and Inanç Birol. Abyss: a parallel assembler for short read
440 sequence data. *Genome research*, 19(6):1117–1123, 2009.
- 441 [7] Yu Peng, Henry C.M. Leung, S.M. Yiu, and Francis Y.L. Chin. Idba-ud: a de
442 novo assembler for single-cell and metagenomic sequencing data with highly
443 uneven depth. *Bioinformatics*, 28:1420–1428, 2012.
- 444 [8] Anton Bankevich, Sergey Nurk, Dmitry Antipov, Alexey A. Gurevich, Mikhail
445 Dvorkin, Alexander S. Kulikov, Valery M. Lesin, Sergey I. Nikolenko, Son

- 446 Pham, Andrey D. Prjibelski, Alexey V. Pyshkin, Alexander V. Sirotkin, Niko-
447 lay Vyahhi, Glenn Tesler, Max A. Alekseyev, and Pavel A. Pevzner. Spades: A
448 new genome assembly algorithm and its applications to single-cell sequencing.
449 *Journal of Computational Biology*, 19(5):455–477, 2012.
- 450 [9] Dinghua Li, Ruibang Luo, Chi-Man Liu, Chi-Ming Leung, Hing-Fung Ting,
451 Kunihiko Sadakane, Hiroshi Yamashita, and Tak-Wah Lam. Megahit v1. 0:
452 A fast and scalable metagenome assembler driven by advanced methodologies
453 and community practices. *Methods*, 102:3–11, 2016.
- 454 [10] Shakya Migun, Christopher Quince, James Campbell, Zamin Yang, Christo-
455 pher Schadt, and Mircea Podar. Comparative metagenomic and rna microbial
456 diversity characterization using archaeal and bacterial synthetic communities.
457 *Enivromental Microbiology*, 15(6):1882–1899, 2013.
- 458 [11] H Chitsaz, JL Yee-Greenbaum, G Tesler, MJ Lombardo, CL Dupont, JH Bad-
459 ger, M Novotny, DB Rusch, LJ Fraser, NA Gormley, O Schulz-Trieglaff,
460 GP Smith, DJ Evers, PA Pevzner, and RS Lasken. Efficient de novo assembly
461 of single-cell bacterial genomes from short-read data sets. *Nat Biotechnol*, 29
462 (10):915–21, 2011.
- 463 [12] Anthony M. Bolger, Marc Lohse, and Bjoern Usadel. Trimmomatic: A flexible
464 trimmer for illumina sequence data. *Bioinformatics*, 30(15):2114–2120, 2014.
- 465 [13] Matthew D MacManes. On the optimal trimming of high-throughput mrna
466 sequence data. *Frontiers in genetics*, 5:13, 2014.
- 467 [14] Heng Li and Richard Durbin. Fast and accurate short read alignment with
468 burrows–wheeler transform. *Bioinformatics*, 25(14):1754–1760, 2009.
- 469 [15] Heng Li, Bob Handsaker, Alec Wysoker, Tim Fennell, Jue Ruan, Nils Homer,
470 Gabor Marth, Goncalo Abecasis, Richard Durbin, et al. The sequence align-
471 ment/map format and samtools. *Bioinformatics*, 25(16):2078–2079, 2009.
- 472 [16] Stefan Kurtz, Adam Phillippy, Arthur L Delcher, Michael Smoot, Martin
473 Shumway, Corina Antonescu, and Steven L Salzberg. Versatile and open soft-
474 ware for comparing large genomes. *Genome biology*, 5(2):1, 2004.