

1 **Assembling metagenomes, one community at a time**

2 Andries J. van der Walt^{1,2†}, Marc W. Van Goethem^{1†}, Jean-Baptiste Ramond¹, Thulani P.
3 Makhalanyane¹, Oleg Reva² & Don A. Cowan¹

4 1 Centre for Microbial Ecology and Genomics (CMEG), Department of Genetics, University of
5 Pretoria, Pretoria, South Africa

6 2 Centre for Bioinformatics and Computational Biology, Department of Biochemistry, University of
7 Pretoria, Pretoria, South Africa

8 † Contributed equally

9 † Corresponding authors: Andries van der Walt or Marc Van Goethem

10 Centre for Microbial Ecology and Genomics (CMEG)

11 Natural Sciences Building 2

12 Lynnwood Road

13 University of Pretoria

14 Pretoria 0028

15 South Africa

16 Email: andriesvanderwalt@gmail.com / mwvangoethem@gmail.com

17 **Abstract**

18 Metagenomics allows unprecedented access to uncultured environmental microorganisms. The
19 assembly of metagenomic sequences facilitates gene prediction and annotation, and enables the
20 assembly of draft genomes of dominant taxa, including uncultured members of a community.
21 However, while numerous platforms have been developed for this critical step, there is currently no
22 strict standard for the assembly of metagenomic sequence data. To assist with selection of an
23 appropriate metagenome assembler we evaluated the capabilities of eight prominent assembly tools
24 on nine publicly-available environmental metagenomes. We found that assembler choice ultimately
25 hinges on the scientific question at hand, the resources available and the bioinformatic competence
26 of the researcher. We provide a concise workflow for the selection of the best assembly tool.

27

28 Introduction

29 The 'science' of metagenomics has greatly accelerated the study of uncultured microorganisms in
30 their natural environments, providing unparalleled insights into microbial community composition and
31 putative functionality¹. Even though shotgun metagenomic sequencing provides comprehensive
32 access to microbial genomic material, many of the encoded functional genes are substantially longer
33 (~1000 bp²) than the length of reads provided by the sequencing platforms³ most commonly used
34 for shotgun metagenomic studies (Illumina HiSeq 3000, 2 x 150 bp; <http://www.illumina.com/>). Thus,
35 raw sequence data alone are typically not sufficient for an in-depth analysis of a communities
36 functional gene repertoire. Moreover, unassembled metagenomic sequence data are fragmented,
37 noisy, error prone and contain uneven sequencing depths⁴.

38 To assist in the accurate and thorough analysis of metagenomes, sequence data can be assembled
39 into larger contiguous segments (contigs)⁵. To this end, numerous metagenome assembly tools
40 (assemblers) have been developed, the vast majority of which assemble sequences in *de novo*
41 fashion. In short, metagenomic sequences are split into predefined segments (*k*-mers), which are
42 overlapped into a network, and paths are traversed iteratively to create longer contigs⁶. *De novo*
43 assembly is advantageous as it allows for more confident gene prediction than is attainable from
44 unassembled data⁷. Furthermore, *de novo* assembled metagenomes facilitate the discovery and
45 reconstruction of novel genomes and/or genomic elements⁸.

46 Improvements to assembly quality have greatly expanded the scope of questions that can be
47 answered using shotgun metagenome sequencing including, for example: determination of microbial
48 community composition and functional capacity⁹, microbial population properties¹⁰, comparisons of
49 microbial communities from various environments¹¹, extraction of full genomes from metagenomes⁵
50 and genomics-informed microorganism isolation¹². Each of these questions require researchers to
51 emphasise specific features of the metagenome. Genome-centric questions^{5,12} require long
52 contigs/scaffolds, while gene-centric questions⁹⁻¹¹ require high confidence contigs and the assembly
53 of a large proportion of the metagenomic dataset.

54 Considering the wealth of available assemblers, it is particularly important that researchers
55 understand assembler performance, especially for investigators who lack appropriate bioinformatic
56 expertise. Firstly, an assembler needs to produce a high proportion of long contigs (>1000 bp). Long
57 contigs allow for more accurate interpretation of full genes within a genomic context and facilitate
58 the reconstruction of single genomes. A good assembler should also utilize most of the raw
59 sequence data to generate the largest assembly span possible. Furthermore, an assembler needs
60 an intuitive and user-friendly interface to enable assembly with minimal effort and rapid processing
61 of the metagenomic data. Finally, tools should be able to assemble metagenomes using the least
62 computational resources possible. Metagenomic assemblers are consistently being developed, this
63 requires regular benchmarking, as with other bioinformatic tools¹³.

64 Here we benchmark seven prominent open-source metagenome assemblers (Velvet v1.2.10¹⁴,
65 MetaVelvet v1.2.02¹⁵, SPAdes v3.9.0¹⁶, metaSPAdes v3.9.0¹⁷, Ray Meta v2.3.1¹⁸, IDBA-UD v1.1.1¹⁹,
66 MEGAHIT v1.0.6²⁰ as well as the commercially-available CLC Genomics Workbench v8.5.1
67 (QIAGEN Bioinformatics; [https://www.qiagenbioinformatics.com/products/clc-genomics-
68 workbench/](https://www.qiagenbioinformatics.com/products/clc-genomics-workbench/); Supplementary Table 1). We compare each assemblers performance on nine complex
69 metagenomes from three distinct environments (i.e., three publicly available metagenomes each
70 from soil, aquatic and human gut niches). While most of the assemblers assessed here have been
71 tested and reviewed extensively²¹⁻²⁵, in this article we provide an elegant reference framework which
72 both experienced and inexperienced researchers can use to determine which assembler is best
73 aligned with their project scope, resources and computational background.

74

75 **Materials and Methods**

76 *Metagenomic datasets*

77 In this study we contrast the assemblies of nine publicly available metagenomic datasets uploaded
78 to the MG-RAST server (<http://metagenomics.anl.gov/>), the SRA (<https://www.ncbi.nlm.nih.gov/sra>)
79 or the JGI IMG/M database (<https://img.jgi.doe.gov/>). The metagenomes are from three distinct
80 environments, namely; soil (Iowa⁸, Oklahoma²⁶, and Permafrost²⁷); aquatic (Kolkata Lake
81 (unpublished data), Arctic Frost Flower²⁸ and Tara Ocean²⁹) and human gut niches (Scandinavian
82 Gut³⁰, European Gut³¹ and Infant Gut³²; Table 1). Each dataset was unique in its complexity and
83 sequencing was performed at different depths. All metagenomes were sequenced using Illumina
84 short read technology producing paired-end reads ranging from 100 to 151 bp in length. Most
85 datasets were sequenced on the Illumina HiSeq 2000 platform, except for the Permafrost
86 metagenome which was sequenced using an Illumina Genome Analyzer II, and the Kolkata Lake
87 metagenome which comprised sequences generated by an Illumina MiSeq. This allowed for
88 comparisons of each assemblers' performance under different coverage and taxonomic diversity.
89 We opted to exclusively evaluate metagenomes sequenced using Illumina platforms due to their
90 popularity and applicability to metagenomic datasets³.

91 Prior to assembling the short read metagenomes, we used Prinseq-lite v0.20.4³³ for read quality
92 control. We removed all reads with mean quality scores of less than 20 [*-min_qual_mean* 20], and
93 removed all sequences contains any ambiguous bases (*N*) [*-ns_max_n* 0].

94 After quality filtering, we assessed the level of coverage of each metagenome using Nonpareil, a
95 statistical program that uses read redundancy to estimate sequence coverage³⁴.

96 *Evaluation of the metagenome assemblers*

97 Most assemblies were performed on a local server (48 Intel® Xeon® CPU E5-2680 v3 @ 2.50 GHz
98 processors, 504 GB physical memory, 15 TB disk space) using 8 threads. However, SPAdes,
99 metaSPAdes and IDBA-UD required more memory, and assembly was performed on the Lengau
100 cluster of the Centre for High Performance Computing (CHPC) for the Iowa and Oklahoma soil
101 datasets. SPAdes, metaSPAdes, IDBA-UD and MEGAHIT iteratively analyse *k*-mer lengths to find
102 the optimal value, and these assemblers were allowed to optimise their own *k*-mer lengths. The other
103 assemblers used *k*-mer values of 55 (Velvet: 51; MetaVelvet: 51; SPAdes: 33, 55, 71; metaSPAdes:
104 33, 55, 71; Ray Meta: 55; IDBA-UD: 20, 30, 40, 50, 60, 70, 71; MEGAHIT: 21, 41, 61, 81, 99; CLC
105 Genomics Workbench: 55). To control for *k*-mer length bias, we compared each assembler's
106 performance at *k*-mer lengths between 50 and 61. Quality of the generated assemblies were
107 assessed using MetaQUAST. This tool calculates basic assembly statistics, including number of
108 contigs above various lengths (500 bp, 1 kbp, 5 kbp and 50 kbp), assembly span above various
109 lengths (500 bp, 1 kbp, 5 kbp and 50 kbp), *N*50 lengths and *L*50 lengths. To assess the amount of
110 data that was used for each assembly, we mapped back the short fragment sequencing reads to the

111 constructed metagenomes. This was performed using Bowtie 2³⁵, using the sensitive setting. Time
112 and memory (RAM) taken to complete assembly was calculated using an in-house bash script.

113 The heatmap was generated using the heatmaply package³⁶ in R. Values were calculated as a mean
114 proportion from the average value obtained for each statistic. This provided ratios of over- or under-
115 performance relative to the average assembly statistic (-2 to +2). The radial plots were generated
116 using log-transformed data for each assembly statistic of relevance to ensure concise representation
117 of the data.

118 Data availability is provided in Supplementary Tables 1 and 2. A link to each to assembler
119 benchmarked is provided, as are the accession numbers for all nine metagenomes assessed.

120

121 **Results**

122 *Metagenome data and dataset complexity*

123 Using Nonpareil, we confirmed that the soil metagenomes were more complex (less redundant) than
124 the aquatic and human guts metagenomes, which were the least complex (highly redundant; Figure
125 1)³⁷⁻³⁹. All the human gut metagenomes came close to sequencing saturation (with at least 75% of
126 the diversity sequenced; Figure 1). The infant gut metagenome was sequenced to above 90%
127 estimated average coverage (~94%). However, all the sequencing depths reached were insufficient
128 to describe the complete spectrum of microbial members in the samples assessed. For example,
129 the largest metagenome assessed here, the Iowa soil metagenome, only described 48.8% of the
130 total microbial diversity despite the utilization of 47 Gbp of sequence data. We note the published
131 predictions that 2-5 Gbp of sequence data would fully capture an entire natural microbial
132 community⁴⁰.

133 Estimates of the number of microbial species per gram of soil still vary substantially, with values
134 ranging from 2000⁴¹ to more than 830000³⁷. These estimates do not include eukaryotic microbes,
135 which generally possess much larger genomes and are much more difficult to fully sequence⁴². We
136 therefore propose that the sequencing depth required to provide comprehensive coverage of soil
137 metagenomes should be increased by an order of magnitude, to ~100 Gbp. This is a function of the
138 extreme taxonomic heterogeneity of soil microbial communities, and highlights the challenges
139 associated with assembling low coverage metagenomes.

140 *Strategy and approaches of the current research*

141 We defined four criteria to assess the performance of each metagenomic assembler tested; (1) ease
142 of use and assembler attributes, (2) quality of assemblies generated and computational
143 requirements, (3) influence of sequencing depth and coverage and (4) suitability to different
144 environments.

145 *1. Ease of use and assembler attributes*

146 Many researchers entering the field of metagenomics are inexperienced in the use of intricate
147 bioinformatic tools, and may lack extensive computational resources. To assess the ease of use for
148 inexperienced computational biologists we evaluated the availability of a web application or graphical
149 user interface (GUI), ease of installation, availability and completeness of manuals, MPI compatibility
150 and programming language.

151 Seven of the assemblers tested here use command-line interfaces (CLI) and are open-source
152 freeware (Velvet, MetaVelvet, SPAdes, metaSPAdes, Ray Meta, IDBA-UD and MEGAHIT). Only the
153 commercial software CLC Genomics Workbench (Qiagen) implements a GUI (Supplementary Table
154 1). CLC is easily installed on most Linux, Windows or MacOS computers, whereas all other
155 assemblers are limited to Unix-based operating systems. The GUI is intuitive, and users can

156 assemble simply by using a point-and-click interface. CLC provides substantial support (via manuals
157 and web based tutorials) and was the most user-friendly assembler tested here.

158 Unix-based assemblers are inherently more difficult to use and must be installed or compiled from
159 source code using the CLI. All assemblers that are CLI-based can be downloaded from GitHub, while
160 some tools (SPAdes, metaSPAdes, Ray Meta, Velvet and MetaVelvet) provide download links from
161 their respective parent websites. All tools, barring SPAdes and metaSPAdes, provide Message
162 Passing Interface (MPI) compatibility, allowing parallelization which reduces computational time. All
163 tools assessed here provide manuals or 'readme' files either on their websites or GitHub repositories,
164 although others, such as IDBA-UD, MetaVelvet, are not comprehensive and lack information on
165 installation or implementation. Tools with more complete manuals (MEGAHIT and Ray Meta) feature
166 extensive wiki pages and frequently asked questions. The number of citations, websites,
167 programming languages and MPI compatibility of all the tools assessed are provided in
168 Supplementary Table 1.

169 2. *Benchmarking quality of assemblies generated and computational requirements*

170 Evaluating metagenome assembly quality is challenging without the use of known reference
171 genomes for diverse microbial communities. We compared assembly quality using many standard
172 metrics, including the total number of contigs longer than 500 bp, 1 kbp (referred to as long contigs
173 throughout) and 50 kbp (referred to as ultra-long contigs throughout), maximum contig length, *N50*
174 length of the contigs (length of the median contig, representing the length of the smallest contig at
175 which half of the assembly is represented), mapping rate and assembly span (total length assembled
176 using contigs > 500 bp). We used MetaQUAST to evaluate these assembly quality statistics²².

177 We selected the Tara Ocean metagenome²⁹ for a comparison of each assembler at *k*-mer lengths
178 between 50 and 61. We selected this range as the assemblers which automatically optimize *k*-mer
179 values generally set sizes within this range. We set the other non-optimizing assemblers to 55.
180 Compared to the other metagenomes, the Tara Ocean metagenomic dataset is of intermediate
181 complexity and sequencing depth (Figure 1, Supplementary Table 2). This metagenome was
182 sequenced on an Illumina HiSeq instrument, which is currently the most widely used shotgun
183 metagenome sequencing technology³. This 5.4 Gbp metagenome comprised more than 27 million
184 high-quality read pairs with a mean read pair length of 200.3 bp (Supplementary Table 2).

185 SPAdes (1415), Ray Meta (1329), IDBA-UD (1166) and metaSPAdes (1124) provided assemblies
186 with high *N50* values (> 1000 bp), while the assemblies generated using CLC, MEGAHIT, Velvet
187 and MetaVelvet produced *N50* statistics below 1000 bp (Figure 2; Table 1; Supplementary Figure
188 1). Overall, the assembly spans varied considerably with SPAdes (275.9 Mbp), MEGAHIT (210.6
189 Mbp), metaSPAdes (202.8 Mbp) and IDBA-UD (179.7 Mbp) assembling the largest metagenomes.
190 Assembly span was correlated with the number of reads mapping back to the assemblies ($R^2 = 0.92$;
191 Supplementary Figure 2), with SPAdes and metaSPAdes having the highest values (Table 1). Both

192 IDBA-UD and MEGAHIT mapped back more than 50% of the sequence reads to the assemblies.
193 SPAdes also produced the most contigs over 1 kbp (70711), while MEGAHIT, IDBA-UD and
194 metaSPAdes created fewer contigs in that size range, but all were comparable to each other
195 (between 48640 and 56243 contigs). The largest contig was assembled by SPAdes (197 kbp),
196 followed by metaSPAdes (142 kbp) and IDBA-UD (101 kbp). These three assemblers also produced
197 the most 'ultra-long' contigs (> 50 kbp); with 54, 37 and 2 contigs, respectively.

198 The computational requirements of an assembly tool should be a major consideration when selecting
199 an assembler. We evaluated all assemblers in relation to the time taken to assemble the Tara Ocean
200 metagenome (Supplementary Figure 3; Table 1) using the same number of threads ($n=8$;
201 Supplementary Figure 3A). Velvet, MetaVelvet and CLC assembled the metagenome in less than
202 an hour, while MEGAHIT and Ray Meta were substantially slower, assembling over multiple hours.
203 IDBA-UD, SPAdes and metaSPAdes required considerably more time to complete assembly, taking
204 approximately 24 hours, or more. In terms of memory requirements, SPAdes was the most 'memory
205 expensive' (157 GB of RAM), followed by Velvet and MetaVelvet (both 109 GB), which is
206 substantially more RAM than is available on an average desktop computer (16 GB). By contrast,
207 MEGAHIT (11 GB) and CLC (16 GB) were the most memory efficient assemblers (Supplementary
208 Figure 1B and 3; Table 1).

209 Overall, SPAdes, metaSPAdes, IDBA-UD and MEGAHIT displayed the best performances in
210 assembling this metagenome of intermediate size and complexity, as they produced very high $N50$
211 values, a high proportion of long contigs and the widest assembly spans. While SPAdes was the
212 best assembler overall, MEGAHIT was the most memory efficient, as it produced an assembly
213 comparable to the best performing assemblers while using only a fraction of computational
214 resources.

215 3. *Benchmarking influence of sequencing depth and coverage*

216 Temperate soil communities are generally more diverse than extreme counterparts (e.g., permafrost;
217 Supplementary Table 3, Figure 1)¹¹. Subsequently, high levels of diversity within these biomes
218 require much deeper sequencing effort. Differences in microorganism abundances and strain level
219 heterogeneity introduce complications during metagenome assembly, resulting in increased memory
220 requirements and longer computational run-times, which may challenge assemblers. The two
221 temperate soil metagenomes assessed here have vastly different sequencing depths, thus providing
222 us with the scope to assess the influence of sequencing depth on the performance of each
223 assembler. The Oklahoma soil metagenome²⁶ had a low sequencing depth (9 Gbp) and estimated
224 coverage (11%), c.f. the Iowa soil metagenome⁸, which had a very high sequencing depth (47 Gbp)
225 and 49% estimated coverage (Figure 1, Supplementary Table 2). We predicted that deeper
226 sequencing effort would be correlated with an increase in metagenome coverage³⁴.

227 All assemblers successfully assembled the Oklahoma metagenome, although SPAdes required
228 considerably more memory (up to 1 TB RAM, Supplementary Table 3). Nevertheless, SPAdes
229 produced the best assembly statistics for most categories (9548 long contigs and an assembly span
230 of 54.3 Mbp; Supplementary Table 3; Figure 2). IDBA-UD and MEGAHIT used less than 500 GB of
231 RAM and were comparable in performance (3828 and 3416 long contigs, and assembly spans of
232 17.2 Mbp and 20.2 Mbp, respectively; Supplementary Table 3; Figure 2). It is noteworthy that while
233 metaSPAdes was one of the best performing assemblers for the Tara Ocean metagenome (Figure
234 2), it performed poorly here (Supplementary Table 3; Figure 2), suggesting that metaSPAdes is ill-
235 suited to assembling low coverage metagenomes.

236 The massive Iowa soil metagenome could not be assembled by either SPAdes or IDBA-UD using
237 our available computing resources (1 TB of RAM). This is in agreement with the methodology
238 described by the authors who generated this dataset, who digitally normalized and partitioned the
239 Iowa metagenome to allow for assembly using Velvet⁸. Remarkably, MEGAHIT and CLC assembled
240 the Iowa metagenome using less than 500 GB of RAM. MEGAHIT performed best across most
241 categories tested (assembly span of 1036.5 Mbp, largest contig of 104841 bp, and 277623 long
242 contigs; Figure 2), while CLC produced the third-best assembly (assembly span of 432.7 Mbp,
243 largest contig of 70207 and 114196 long contigs), using less than 64GB of memory. MetaSPAdes
244 performed comparably to MEGAHIT but had much higher computational resource requirements to
245 assemble the Iowa soil metagenome, using up to 1TB of RAM (assembly span of 873.8 Mbp, largest
246 contig of 188499 bp, and 225046 long contigs).

247 Overall, we found that sequencing depth greatly influenced the performance of the assemblers,
248 although the most memory-efficient tools, MEGAHIT and CLC, performed well irrespective of
249 sequencing coverage. SPAdes and IDBA-UD produced good assemblies for the Oklahoma soil
250 metagenome, but were extremely expensive in terms of memory and failed to assemble the Iowa
251 soil metagenome. We found that metaSPAdes produced a better assembly for the Iowa soil
252 metagenome than the Oklahoma dataset. MetaSPAdes performed optimally for the assembly of the
253 high-coverage metagenome, but was less efficient in the assembly of the low-coverage
254 metagenome.

255 *4. Benchmarking suitability to various environments*

256 Environmental samples are widely dissimilar in microbial community complexity and have distinct
257 taxonomic compositions. In this study, we assembled metagenomes from three environmental
258 biomes of different phylotypic complexities. Overall, SPAdes, MEGAHIT, IDBA-UD and metaSPAdes
259 assembled most of the metagenomes well, according to the parameters we evaluated
260 (Supplementary Tables 3-5). SPAdes consistently provided the largest contigs and the widest
261 assembly spans. MEGAHIT demanded far fewer computational resources, and yet produced similar
262 assemblies to metaSPAdes and IDBA-UD. CLC provided assemblies of moderate to high quality,
263 was the easiest to use and performed particularly well on large metagenomes. Together, these

264 results indicate that no single assembler performs best across all sequencing platforms and
265 datasets.

266 *How to select a metagenome assembler*

267 Bioinformatics projects can be limited by memory (RAM) requirements. SPAdes, metaSPAdes,
268 IDBA-UD, Velvet and MetaVelvet all have large memory requirements during the assembly of
269 massive datasets. MEGAHIT and CLC are extremely memory efficient, as they required less than
270 500 GB and 64 GB of RAM, respectively, to assemble the massive Iowa soil metagenome.
271 MEGAHIT, for example, generates succinct *de Bruijn* graphs to achieve efficient memory usage²⁰.

272 Our results indicate that although many assemblers perform comparably, their applicability is defined
273 by the research question at hand. SPAdes, for example, generated good assemblies with multiple
274 long and ultra-long contigs for most datasets. These are ideal characteristics for genome-centric
275 studies, which require the binning of draft genomes from community sequence data⁴³. By contrast,
276 metaSPAdes considers read coverage during assembly, making it more applicable for microbial
277 community profiling¹⁷. While SPAdes and metaSPAdes produced the best assemblies in general,
278 MEGAHIT performed comparably and emerged as a rapid and memory efficient alternative
279 assembler.

280 In conclusion, we argue that when selecting an assembler, the primary consideration should be the
281 research question. Selecting an appropriate assembler is essential to make full use of metagenomic
282 sequence dataset. The primary objectives of the project, whether gene- or genome-centric, for
283 example, should dictate the choice of assembler. We suggest that a secondary consideration should
284 be the computational resources available to the researcher. Some assemblers are very memory
285 efficient, while others sacrifice computational resources for improved assembly quality. Finally, as
286 most assemblers use a CLI (and are more flexible than those constrained by a GUI), the GUI-based
287 CLC platform is an excellent alternative if bioinformatic skill level is a consideration.

288 **Other analyses**

289 In additional analyses (Figure 2), we compared the performance of each assembler on a low diversity
290 soil metagenome (Supplementary Table 3), other aquatic metagenomes (Supplementary Table 4)
291 and human gut microbiomes (Supplementary Table 5).

292 Discussion

293 Over the last decade, high throughput sequencing has revolutionised the field of microbial ecology⁴⁴.
294 Amplicon-based technologies have allowed for near-complete classification of whole microbial
295 communities, including populations of bacteria, archaea and fungi⁴⁵. The emergence of two key
296 platforms for analysing amplicon sequencing data, mothur⁴⁶ and QIIME⁴⁷, has allowed for
297 methodological standards to be set, which enables robust comparisons between studies⁴⁸.

298 While whole community shotgun metagenome sequencing has facilitated the in-depth description of
299 microbial communities from diverse environments, such as the human gut⁴⁹ and acid mine drainage
300 systems⁵⁰, no standards exist with regard to assembly platforms or their use. While numerous
301 reviews on strategies to analyse metagenomic data have been published⁵¹, there are currently no
302 standard assembly procedures implemented to enable thorough comparative analyses between
303 projects. Numerous pipelines for processing metagenomic sequence data are available. These
304 typically integrate existing tools into a single workflow for rapid, standardized analysis (e.g., MG-
305 RAST, MetAMOS, and IMG/M)⁵²⁻⁵⁴. However, few of these pipelines are as widely used as mothur
306 or QIIME in barcoding studies. This is partly because integrated metagenome analysis tools, such
307 as MetAMOS, do not achieve the flexibility afforded by using each tool individually (e.g., using
308 separate tools for assembly, binning and taxonomic assignment).

309 Consequently, investigators can analyse unassembled reads¹¹, optimize their assembly parameters
310 or even develop their own tools to assemble their data prior to further analysis⁵⁵. However, within
311 the scope of metagenome assembly, essential details are often omitted when describing methods⁵⁶.
312 This leads to methodological discrepancies, and severely limits the possibility of making routine,
313 robust comparisons between studies. This issue was recently highlighted by Vollmers, et al. ²¹, who
314 reported that the taxonomic diversity patterns of microbial communities differed substantially,
315 depending on the assembler used.

316 In our comparative analyses of the most popular assembly platforms, SPAdes produced the most
317 long contigs, independent of the metagenome origin. SPAdes is ideal for genome-centric research
318 questions that require long and ultra-long contigs, such as those that aim to bin and reconstruct
319 single genomes from shotgun metagenomes¹⁶. By contrast, MEGAHIT, IDBA-UD and metaSPAdes
320 provided very large assembly spans and consider sequence coverage during assembly, reducing
321 the number of misassemblies generated. These tools are thus more appropriate for research
322 questions related to taxonomic profiling of natural microbial communities, for functionally annotating
323 microbial communities, for the analysis of population scale dynamics or for comparison of microbial
324 communities across biomes^{17,20}. Overall, MEGAHIT produced some of the best assemblies
325 throughout this study, while only using a fraction of the computational resources required by other
326 assemblers. We strongly recommend MEGAHIT for researchers who do not have access to large
327 computational resources. Finally, the CLC assembler is ideal for researchers who lack a depth of
328 bioinformatic knowledge, or who prefer to use a GUI and are willing to invest in software which is

329 easier to use. CLC is easy to install, has an intuitive interface and provides a compromise in which
330 assembly quality may be sacrificed for ease of use. Strikingly, the most widely cited assembler
331 assessed here (Velvet cited 5621 times; Supplementary Table 1) did not perform well across most
332 metagenomes, while scarcely cited platforms (MEGAHIT, metaSPAdes cited 82 and 8 times,
333 respectively; Supplementary Table 1) performed well across most statistics assessed here.

334 No assembler tested here consistently provided superior assemblies across the different
335 metagenomes. Consequently, we propose a viable methodology for the selection of an appropriate
336 assembler, dictated by (1) the scientific research question posed, then by (2) the computational
337 resources available, and (3) the bioinformatics skill level of the researcher (Figure 4). In light of the
338 above proposed framework, we urge researchers to carefully consider the assembler used (as well
339 as the entire bioinformatics pipeline followed) while specifically bearing in mind their research
340 question and what feature of the dataset they want accentuated.

341

342 **Acknowledgements**

343 We thank the National Research Foundation under the grant numbers 102910 (AJvdW) and 97891
344 (MWVG) and the University of Pretoria Genomics Research Institute for financial assistance. We
345 thank Dr Surendra Vikram for helping with the installation of assemblers, as well as Dane Kennedy
346 from the CHPC, South Africa with his assistance in accessing the large memory nodes of the CHPC.
347 We thank S. de Scally for support and helpful discussions.

348 **Author contributions**

349 A.J.v.d.W. contributed research, analysis, formulation of ideas and writing. M.W.V.G. contributed
350 research, analysis, formulation of ideas and writing. J.B.R. contributed formulation of ideas, direction
351 and writing. T.P.M. contributed formulation of ideas and writing. O.R. contributed direction and
352 writing. D.A.C. contributed formulation of ideas, direction and writing.

353 **Competing financial interests**

354 The authors declare no competing financial interests.

355

356 **References**

- 357 1 Riesenfeld, C. S., Schloss, P. D. & Handelsman, J. Metagenomics: genomic analysis of microbial
358 communities. *Annu. Rev. Genet.* **38**, 525-552 (2004).
- 359 2 Xu, L. *et al.* Average gene length is highly conserved in prokaryotes and eukaryotes and diverges
360 only between the two kingdoms. *Molecular biology and evolution* **23**, 1107-1108 (2006).
- 361 3 Clooney, A. G. *et al.* Comparing Apples and Oranges?: Next Generation Sequencing and Its Impact
362 on Microbiome Analysis. *PLoS one* **11**, e0148028 (2016).
- 363 4 Nagarajan, N. & Pop, M. Sequence assembly demystified. *Nature Reviews Genetics* **14**, 157-167
364 (2013).
- 365 5 Anantharaman, K. *et al.* Thousands of microbial genomes shed light on interconnected
366 biogeochemical processes in an aquifer system. *Nature communications* **7**, 13219 (2016).
- 367 6 Compeau, P. E., Pevzner, P. A. & Tesler, G. How to apply de Bruijn graphs to genome assembly.
368 *Nature biotechnology* **29**, 987-991 (2011).
- 369 7 Loman, N. J. *et al.* A culture-independent sequence-based metagenomics approach to the
370 investigation of an outbreak of Shiga-toxigenic *Escherichia coli* O104: H4. *Jama* **309**, 1502-1510
371 (2013).
- 372 8 Howe, A. C. *et al.* Tackling soil diversity with the assembly of large, complex metagenomes.
373 *Proceedings of the National Academy of Sciences* **111**, 4904-4909 (2014).
- 374 9 Handelsman, J. Metagenomics: application of genomics to uncultured microorganisms.
375 *Microbiology and molecular biology reviews* **68**, 669-685 (2004).
- 376 10 Tyson, G. W. *et al.* Community structure and metabolism through reconstruction of microbial
377 genomes from the environment. *Nature* **428**, 37-43 (2004).
- 378 11 Fierer, N. *et al.* Cross-biome metagenomic analyses of soil microbial communities and their
379 functional attributes. *Proceedings of the National Academy of Sciences* **109**, 21390-21395 (2012).
- 380 12 Wurch, L. *et al.* Genomics-informed isolation and characterization of a symbiotic Nanoarchaeota
381 system from a terrestrial geothermal environment. *Nature communications* **7** (2016).
- 382 13 Baruzzo, G. *et al.* Simulation-based comprehensive benchmarking of RNA-seq aligners. *Nature*
383 *Methods* (2016).
- 384 14 Zerbino, D. R. & Birney, E. Velvet: algorithms for de novo short read assembly using de Bruijn
385 graphs. *Genome research* **18**, 821-829 (2008).
- 386 15 Namiki, T., Hachiya, T., Tanaka, H. & Sakakibara, Y. MetaVelvet: an extension of Velvet assembler to
387 de novo metagenome assembly from short sequence reads. *Nucleic acids research* **40**, e155-e155
388 (2012).
- 389 16 Bankevich, A. *et al.* SPAdes: a new genome assembly algorithm and its applications to single-cell
390 sequencing. *Journal of Computational Biology* **19**, 455-477 (2012).
- 391 17 Nurk, S., Meleshko, D., Korobeynikov, A. & Pevzner, P. metaSPAdes: a new versatile de novo
392 metagenomics assembler. *arXiv preprint arXiv:1604.03071* (2016).
- 393 18 Boisvert, S., Raymond, F., Godzaridis, É., Laviolette, F. & Corbeil, J. Ray Meta: scalable de novo
394 metagenome assembly and profiling. *Genome biology* **13**, 1 (2012).
- 395 19 Peng, Y., Leung, H. C., Yiu, S.-M. & Chin, F. Y. IDBA-UD: a de novo assembler for single-cell and
396 metagenomic sequencing data with highly uneven depth. *Bioinformatics* **28**, 1420-1428 (2012).
- 397 20 Li, D., Liu, C.-M., Luo, R., Sadakane, K. & Lam, T.-W. MEGAHIT: an ultra-fast single-node solution for
398 large and complex metagenomics assembly via succinct de Bruijn graph. *Bioinformatics*, btv033
399 (2015).
- 400 21 Vollmers, J., Wiegand, S. & Kaster, A.-K. Comparing and Evaluating Metagenome Assembly Tools
401 from a Microbiologist's Perspective-Not Only Size Matters! *PLoS one* **12**, e0169662 (2017).
- 402 22 Mikheenko, A., Saveliev, V. & Gurevich, A. MetaQUAST: evaluation of metagenome assemblies.
403 *Bioinformatics*, btv697 (2015).
- 404 23 Vázquez-Castellanos, J. F., García-López, R., Pérez-Brocal, V., Pignatelli, M. & Moya, A. Comparison
405 of different assembly and annotation tools on analysis of simulated viral metagenomic
406 communities in the gut. *BMC genomics* **15**, 37 (2014).
- 407 24 Pignatelli, M. & Moya, A. Evaluating the fidelity of de novo short read metagenomic assembly using
408 simulated data. *PLoS one* **6**, e19984 (2011).

- 409 25 Charuvaka, A. & Rangwala, H. Evaluation of short read metagenomic assembly. *BMC genomics* **12**,
410 S8 (2011).
- 411 26 Luo, C. *et al.* Soil microbial community responses to a decade of warming as revealed by
412 comparative metagenomics. *Applied and environmental microbiology* **80**, 1777-1786 (2014).
- 413 27 Hultman, J. *et al.* Multi-omics of permafrost, active layer and thermokarst bog soil microbiomes.
414 *Nature* (2015).
- 415 28 Bowman, J. S., Berthiaume, C. T., Armbrust, E. V. & Deming, J. W. The genetic potential for key
416 biogeochemical processes in Arctic frost flowers and young sea ice revealed by metagenomic
417 analysis. *FEMS microbiology ecology* **89**, 376-387 (2014).
- 418 29 Sunagawa, S. *et al.* Structure and function of the global ocean microbiome. *Science* **348**, 1261359
419 (2015).
- 420 30 Tremaroli, V. *et al.* Roux-en-Y gastric bypass and vertical banded gastroplasty induce long-term
421 changes on the human gut microbiome contributing to fat mass regulation. *Cell metabolism* **22**,
422 228-238 (2015).
- 423 31 Karlsson, F. H. *et al.* Gut metagenome in European women with normal, impaired and diabetic
424 glucose control. *Nature* **498**, 99-103 (2013).
- 425 32 Bäckhed, F. *et al.* Dynamics and stabilization of the human gut microbiome during the first year of
426 life. *Cell host & microbe* **17**, 690-703 (2015).
- 427 33 Schmieder, R. & Edwards, R. Quality control and preprocessing of metagenomic datasets.
428 *Bioinformatics* **27**, 863-864 (2011).
- 429 34 Rodriguez-R, L. M. & Konstantinidis, K. T. Nonpareil: a redundancy-based approach to assess the
430 level of coverage in metagenomic datasets. *Bioinformatics* **30**, 629-635 (2014).
- 431 35 Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nature methods* **9**, 357-
432 359 (2012).
- 433 36 Galili, T. heatmaply: interactive heat maps (with R). *Month* (2016).
- 434 37 Gans, J., Wolinsky, M. & Dunbar, J. Computational improvements reveal great bacterial diversity
435 and high metal toxicity in soil. *Science* **309**, 1387-1390 (2005).
- 436 38 Torsvik, V., Øvreås, L. & Thingstad, T. F. Prokaryotic diversity--magnitude, dynamics, and controlling
437 factors. *Science* **296**, 1064-1066 (2002).
- 438 39 Tringe, S. G. *et al.* Comparative metagenomics of microbial communities. *Science* **308**, 554-557
439 (2005).
- 440 40 Venter, J. C. *et al.* Environmental genome shotgun sequencing of the Sargasso Sea. *science* **304**, 66-
441 74 (2004).
- 442 41 Schloss, P. D. & Handelsman, J. Toward a census of bacteria in soil. *PLoS computational biology* **2**,
443 e92 (2006).
- 444 42 Lynch, M. & Conery, J. S. The origins of genome complexity. *science* **302**, 1401-1404 (2003).
- 445 43 Becraft, E. D. *et al.* Single-Cell-Genomics-Facilitated Read Binning of Candidate Phylum EM19
446 Genomes from Geothermal Spring Metagenomes. *Applied and environmental microbiology* **82**, 992-
447 1003 (2016).
- 448 44 Shokralla, S., Spall, J. L., Gibson, J. F. & Hajibabaei, M. Next-generation sequencing technologies for
449 environmental DNA research. *Molecular ecology* **21**, 1794-1805 (2012).
- 450 45 Taberlet, P., Coissac, E., Pompanon, F., Brochmann, C. & Willerslev, E. Towards next-generation
451 biodiversity assessment using DNA metabarcoding. *Molecular ecology* **21**, 2045-2050 (2012).
- 452 46 Schloss, P. D. *et al.* Introducing mothur: open-source, platform-independent, community-supported
453 software for describing and comparing microbial communities. *Applied and environmental*
454 *microbiology* **75**, 7537-7541, doi:10.1128/AEM.01541-09 (2009).
- 455 47 Caporaso, J. G. *et al.* QIIME allows analysis of high-throughput community sequencing data. *Nature*
456 *methods* **7**, 335-336 (2010).
- 457 48 Plummer, E., Twin, J., Bulach, D. M., Garland, S. M. & Tabrizi, S. N. A Comparison of Three
458 Bioinformatics Pipelines for the Analysis of Preterm Gut Microbiota using 16S rRNA Gene
459 Sequencing Data. *Journal of Proteomics & Bioinformatics* **8**, 283 (2015).
- 460 49 Gill, S. R. *et al.* Metagenomic analysis of the human distal gut microbiome. *science* **312**, 1355-1359
461 (2006).

- 462 50 Kantor, R. S. *et al.* Genome-resolved meta-omics ties microbial dynamics to process performance in
463 biotechnology for thiocyanate degradation. *Environmental Science & Technology* (2017).
- 464 51 Scholz, M. B., Lo, C.-C. & Chain, P. S. Next generation sequencing and bioinformatic bottlenecks: the
465 current state of metagenomic data analysis. *Current opinion in biotechnology* **23**, 9-15 (2012).
- 466 52 Markowitz, V. M. *et al.* IMG/M 4 version of the integrated metagenome comparative analysis
467 system. *Nucleic Acids Research* **42**, D568-D573 (2014).
- 468 53 Glass, E. M., Wilkening, J., Wilke, A., Antonopoulos, D. & Meyer, F. Using the metagenomics RAST
469 server (MG-RAST) for analyzing shotgun metagenomes. *Cold Spring Harbor Protocols* **2010**, pdb.
470 prot5368 (2010).
- 471 54 Treangen, T. J. *et al.* MetAMOS: a modular and open source metagenomic assembly and analysis
472 pipeline. *Genome Biol* **14**, R2 (2013).
- 473 55 Oh, J. *et al.* Biogeography and individuality shape function in the human skin metagenome. *Nature*
474 **514**, 59-64 (2014).
- 475 56 Brown, C. T. *et al.* Unusual biology across a group comprising more than 15% of domain Bacteria.
476 *Nature* **523**, 208-211 (2015).
- 477

478 **Figure and Table legends.**

479 **Table 1. Computational requirements for assembly of the Tara Oceans metagenome.** Time
480 required is given in seconds, minutes and hours for illustrative purposes and memory in GB of RAM
481 required.

482 **Figure 1. Nonpareil estimates of sequence coverage (redundancy) for the 9 metagenomes**
483 **studied.** Metagenomes are grouped according to their environmental niche.

484 **Figure 2. Radial plots showing the assembly statistics for two soil metagenomes.** A) Oklahoma
485 soil metagenome. B) Iowa soil metagenome

486 **Figure 3. Proposed workflow to select a metagenome assembler based on the research**
487 **question, the computational resources available and the bioinformatic expertise of the**
488 **researcher**

489 **Table 1. Assembly statistics and computational requirements of the Tara Oceans metagenome.** Time required is given in seconds, minutes and
 490 hours for illustrative purposes and memory is provided in GB of RAM required.

	Tara Ocean							
	CLC	IDBA-UD	MEGAHIT	metaSPAdes	MetaVelvet	Ray Meta	SPAdes	Velvet
Number of contigs (≥ 500 bp)	50,716	163,815	216,938	185,419	67,161	6,128	220,178	57,816
Total length	46,069,409	179,686,756	210,621,485	202,770,058	55,972,515	7,277,214	275,920,632	45,425,460
No. of long contigs (≥ 1 kbp)	10,720	50,498	56,243	48,640	12,590	2,179	70,711	8,802
No. of ultra-long contigs (≥ 50 kbp)	0	2	1	37	0	0	54	0
Largest contig	39,748	101,400	62,649	141,519	30,177	41,443	197,381	21,980
<i>N50</i>	880	1,166	982	1124	805	1,329	1415	749
<i>L50</i>	14,113	38,236	58,246	39,033	21,544	1,345	39,617	19,631
Mapping rate (%)	38.98	52.24	55.92	64.03	41.17	8.25	64.46	48.19
Time (seconds)	3,527	69,782	10,455	125,862	2,527	16,419	80,039	2342
Time (minutes)	58.78	1,163.03	174.25	2,097.70	42.12	273.65	1,333.98	39.03
Time (hours)	0.98	19.38	2.90	34.96	0.70	4.56	22.23	0.65
Memory required (GB)	16.23	42.84	10.58	66.53	109.37	42	157.75	109.37

Figure 1. Nonpareil estimates of sequence coverage (redundancy) for the 9 metagenomes studied. Metagenomes are grouped according to their environmental niche.

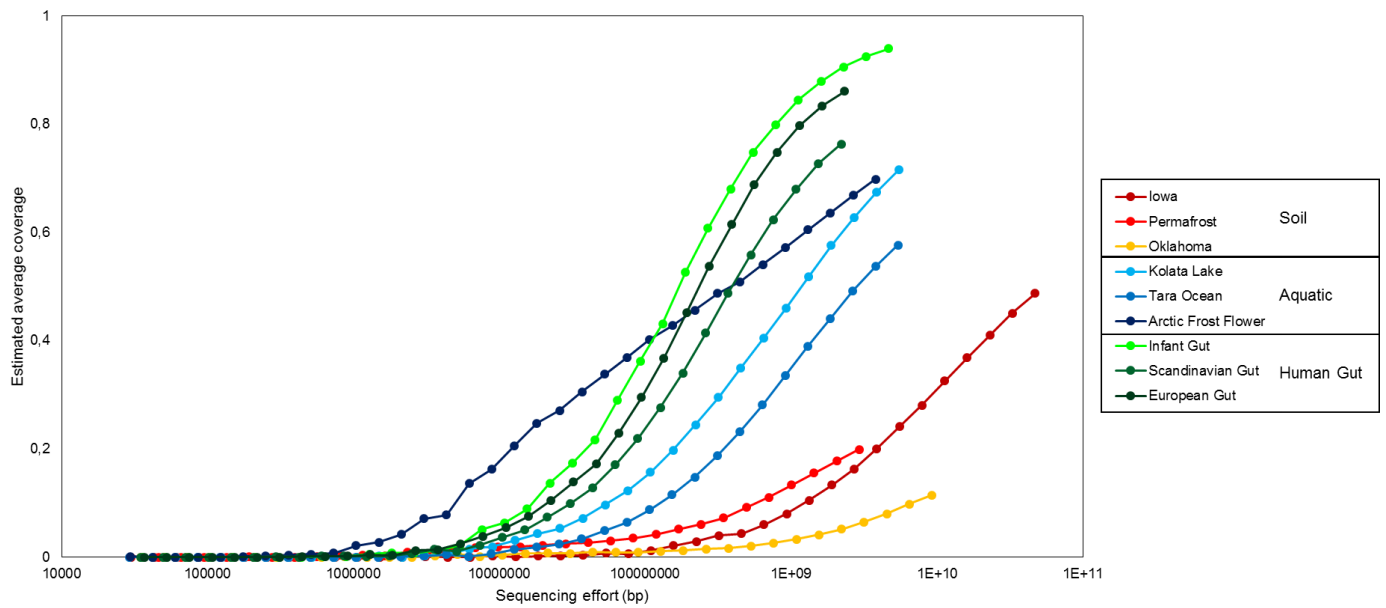


Figure 2. Radial plots showing scaled assembly statistics for the all the metagenomes assessed. Metagenomes are labelled above the respective radial plots.

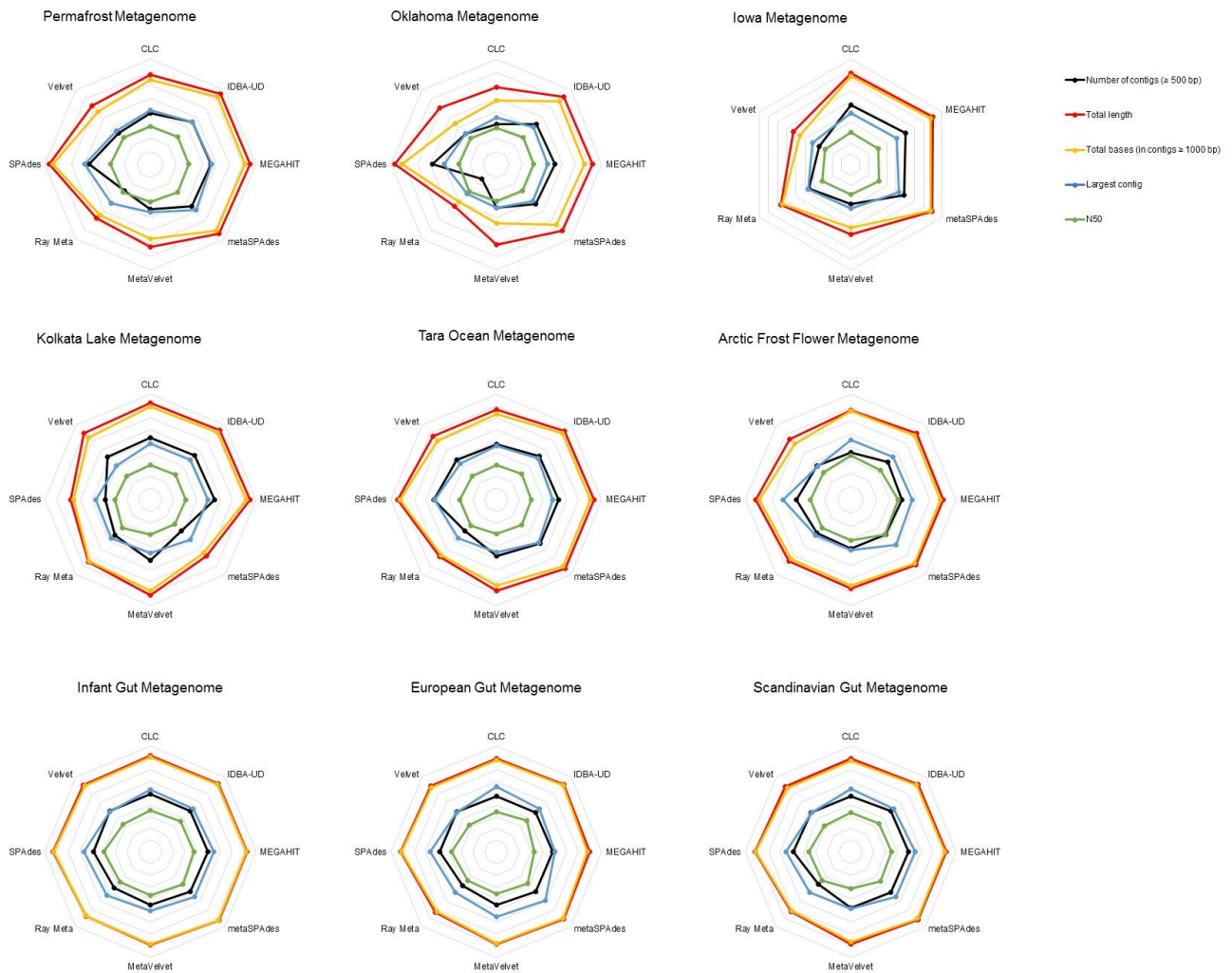


Figure 3. Proposed workflow to select a metagenome assembler based on the research question, the computational resources available and the bioinformatic expertise of the researcher.

