

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

2 Andries J. van der Walt<sup>1,2</sup>, Marc W. Van Goethem<sup>1</sup>, Jean-Baptiste Ramond<sup>1</sup>, Thulani P.  
3 Makhalanyane<sup>1</sup>, Oleg Reva<sup>2</sup> & Don A. Cowan<sup>1\*</sup>

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 \* Corresponding author: Prof. Don A. Cowan

9 Centre for Microbial Ecology and Genomics (CMEG)

10 Natural Sciences Building 2

11 Lynnwood Road

12 University of Pretoria

13 Pretoria 0028

14 South Africa

15 Tel: +27 12 420 5873

16 Email: [don.cowan@up.ac.za](mailto:don.cowan@up.ac.za)

17

18 Andries Johannes van der Walt Email: [andriesvanderwalt@gmail.com](mailto:andriesvanderwalt@gmail.com)

19 Marc Warwick Van Goethem Email: [mwvangoethem@gmail.com](mailto:mwvangoethem@gmail.com)

20 Jean-Baptiste Ramond Email: [jbaptiste.ramond@gmail.com](mailto:jbaptiste.ramond@gmail.com)

21 Thulani Peter Makhalanyane Email: [Thulani.Makhalanyane@up.ac.za](mailto:Thulani.Makhalanyane@up.ac.za)

22 Oleg Reva Email: [oleg.reva@up.ac.za](mailto:oleg.reva@up.ac.za)

23 Don Arthur Cowan Email: [Don.Cowan@up.ac.za](mailto:Don.Cowan@up.ac.za)

24

## 25 **Abstract**

26 **Background:** Metagenomics allows unprecedented access to uncultured environmental  
27 microorganisms. The analysis of metagenomic sequences facilitates gene prediction and annotation,  
28 and enables the assembly of draft genomes, including uncultured members of a community.  
29 However, while several platforms have been developed for this critical step, there is currently no  
30 clear framework for the assembly of metagenomic sequence data.

31 **Results:** To assist with selection of an appropriate metagenome assembler we evaluated the  
32 capabilities of nine prominent assembly tools on nine publicly-available environmental  
33 metagenomes, as well as three simulated datasets. Overall, we found that SPAdes provided the  
34 largest contigs and highest *N50* values across 6 of the 9 environmental datasets, followed by  
35 MEGAHIT and metaSPAdes. MEGAHIT emerged as a computationally inexpensive alternative to  
36 SPAdes, assembling the most complex dataset using less than 500 GB of RAM and within 10 hours.

37 **Conclusions:** We found that assembler choice ultimately depends on the scientific question, the  
38 available resources and the bioinformatic competence of the researcher. We provide a concise  
39 workflow for the selection of the best assembly tool.

40 **Keywords:** metagenome assembly; microbial ecology; Illumina HiSeq; assembler; bioinformatics

41

## 42 **Background**

43 The 'science' of metagenomics has greatly accelerated the study of uncultured microorganisms in  
44 their natural environments, providing unparalleled insights into microbial community composition and  
45 putative functionality [1]. Even though shotgun metagenomic sequencing provides comprehensive  
46 access to microbial genomic material, many of the encoded functional genes are substantially longer  
47 (~1000 bp [2]) than the length of reads provided by the sequencing platforms [3] most commonly  
48 used for shotgun metagenomic studies (Illumina HiSeq 3000, 2 x 150 bp; <http://www.illumina.com/>).  
49 Thus, raw sequence data alone are typically not sufficient for an in-depth analysis of a communities  
50 functional gene repertoire. Moreover, unassembled metagenomic sequence data are fragmented,  
51 noisy, error prone and contain uneven sequencing depths [4].

52 To assist in the accurate and thorough analysis of metagenomes, sequence data can be assembled  
53 into larger contiguous segments (contigs) [5]. To this end, numerous metagenome assembly tools  
54 (assemblers) have been developed, the vast majority of which assemble sequences in *de novo*  
55 fashion. In short, metagenomic sequences are split into predefined segments (*k*-mers), which are  
56 overlapped into a network, and paths are traversed iteratively to create longer contigs [6]. *De novo*  
57 assembly is advantageous as it allows for more confident gene prediction than is attainable from  
58 unassembled data [7]. Furthermore, *de novo* assembled metagenomes facilitate the discovery and  
59 reconstruction of novel genomes and/or genomic elements [8].

60 Improvements to assembly quality have greatly expanded the scope of questions that can be  
61 answered using shotgun metagenome sequencing including, for example: determination of microbial  
62 community composition and functional capacity [9], microbial population properties [10],  
63 comparisons of microbial communities from various environments [11], extraction of full genomes  
64 from metagenomes [5] and genomics-informed microorganism isolation [12]. Each of these  
65 questions require researchers to emphasise specific features of the metagenome. Genome-centric  
66 questions [5, 12] require long contigs/scaffolds, while gene-centric questions [9-11] require high  
67 confidence contigs and the assembly of a large proportion of the metagenomic dataset.

68 Considering the wealth of available assemblers, it is particularly important that researchers  
69 understand assembler performance, especially for investigators who lack appropriate bioinformatic  
70 expertise. Firstly, an assembler needs to produce a high proportion of long contigs (>1000 bp). Long  
71 contigs allow for more accurate interpretation of full genes within a genomic context and facilitate  
72 the reconstruction of single genomes. A good assembler should also utilize most of the raw  
73 sequence data to generate the largest assembly span possible. Furthermore, an assembler needs  
74 an intuitive and user-friendly interface to enable assembly with minimal effort and rapid processing  
75 of the metagenomic data. Finally, tools should be able to assemble metagenomes using the least  
76 computational resources possible. Metagenomic assemblers are consistently being developed, this  
77 requires regular benchmarking, as with other bioinformatic tools [13].

78 Here we benchmark eight prominent open-source metagenome assemblers (Velvet v1.2.10 [14],  
79 MetaVelvet v1.2.02 [15], SPAdes v3.9.0 [16], metaSPAdes v3.9.0 [17], Ray Meta v2.3.1 [18], IDBA-  
80 UD v1.1.1 [19], MEGAHIT v1.0.6 [20] and Omega v1.4 as well as the commercially-available CLC  
81 Genomics Workbench v8.5.1 (QIAGEN Bioinformatics;  
82 <https://www.qiagenbioinformatics.com/products/clc-genomics-workbench/>; Supplementary Table 1).  
83 We compare each assemblers performance on nine complex metagenomes from three distinct  
84 environments (i.e., three publicly available metagenomes each from soil, aquatic and human gut  
85 niches) as well as three simulated datasets. While most of the assemblers assessed here have been  
86 tested and reviewed extensively [21-25], in this article we provide an elegant reference framework  
87 which both experienced and inexperienced researchers can use to determine which assembler is  
88 best aligned with their project scope, resources and computational background.

89

## 90 **Methods**

### 91 *Metagenomic datasets*

92 In this study we contrast the assemblies of nine publicly available metagenomic datasets uploaded  
93 to the MG-RAST server (<http://metagenomics.anl.gov/>), or the sequence read archive (SRA)  
94 (<https://www.ncbi.nlm.nih.gov/sra>). The metagenomes are from three distinct environments, namely;  
95 soil (Iowa [8], Oklahoma [26], and Permafrost [27]); aquatic (Kolkata Lake (unpublished data), Arctic

96 Frost Flower [28] and Tara Ocean [29]) and human gut niches (Scandinavian Gut [30], European  
97 Gut [31] and Infant Gut [32]; Table 1). Each dataset was unique in its complexity and sequencing  
98 was performed at different depths. All metagenomes were sequenced using Illumina short read  
99 technology producing paired-end reads ranging from 100 to 151 bp in length. Most datasets were  
100 sequenced on the Illumina HiSeq 2000 platform, except for the Permafrost metagenome which was  
101 sequenced using an Illumina Genome Analyzer II, and the Kolkata Lake metagenome which  
102 comprised sequences generated by an Illumina MiSeq. This allowed for comparisons of each  
103 assemblers' performance under different coverage and taxonomic diversity. We opted to exclusively  
104 evaluate metagenomes sequenced using Illumina platforms due to their popularity and applicability  
105 to metagenomic datasets [3].

106 Prior to assembling the short read metagenomes, we used Prinseq-lite v0.20.4 [33] for read quality  
107 control. We removed all reads with mean quality scores of less than 20 [-min\_qual\_mean 20], and  
108 removed all sequences contains any ambiguous bases (N) [-ns\_max\_n 0].

109 After quality filtering, we assessed the level of coverage of each metagenome using Nonpareil, a  
110 statistical program that uses read redundancy to estimate sequence coverage [34].

#### 111 *Evaluation of the metagenome assemblers*

112 Most assemblies were performed on a local server (48 Intel® Xeon® CPU E5-2680 v3 @ 2.50 GHz  
113 processors, 504 GB physical memory, 15 TB disk space) using 8 threads. However, SPAdes,  
114 metaSPAdes and IDBA-UD required more memory, and assembly was performed on the Lengau  
115 cluster of the Centre for High Performance Computing (CHPC) for the Iowa and Oklahoma soil  
116 datasets. SPAdes, metaSPAdes, IDBA-UD and MEGAHIT iteratively analyse *k*-mer lengths to find  
117 the optimal value, and these assemblers were allowed to optimise their own *k*-mer lengths. The other  
118 assemblers used *k*-mer values of 55 (Velvet: 51; MetaVelvet: 51; SPAdes: 33, 55, 71; metaSPAdes:  
119 33, 55, 71; Ray Meta: 55; IDBA-UD: 20, 30, 40, 50, 60, 70, 71; MEGAHIT: 21, 41, 61, 81, 99; CLC  
120 Genomics Workbench: 55). In contrast to the above *de Bruijn* graph assemblers, Omega uses  
121 overlap-layout-consensus graphs to generate assemblies. Read pairs are first aligned, followed by  
122 read error correction, hash-table construction, overlap graph construction before generating contigs.  
123 We selected a minimum overlap length of 60. To control for *k*-mer length bias, we compared each  
124 assembler's performance at *k*-mer lengths between 50 and 61. Quality of the generated assemblies  
125 were assessed using MetaQUAST. This tool calculates basic assembly statistics, including number  
126 of contigs above various lengths (500 bp, 1 kbp, 5 kbp and 50 kbp), assembly span above various  
127 lengths (500 bp, 1 kbp, 5 kbp and 50 kbp), *N*50 lengths and *L*50 lengths. To assess the accuracy  
128 and specificity of each assembler, the included synthetic metagenomes were assessed against their  
129 respective constitutive reference genomes in MetaQUAST.

130 To assess the volume of sequencing data that was used for each assembly, we mapped back the  
131 short fragment sequencing reads to the constructed metagenomes. This was performed using

132 Bowtie 2 [35], using the sensitive setting. Time and memory (RAM) taken to complete assembly  
133 were calculated using an in-house bash script.

134 All tables and figures were drawn in R v3.4.0 or Microsoft Excel. Figure 1 was generated using the  
135 freely-available tool Nonpareil. Nonpareil estimates the percentage sequence coverage of  
136 metagenomes (as a fraction of 1) using either the forward or reverse sequence reads. These values  
137 are then plotted using a scatter plot function. Figure 2 was generated using the *heatmaply* package  
138 [36], and clustered using the *hclust* hierarchical clustering package in R. Values were calculated as  
139 a mean over- or under-representation relative to the average value obtained for all the assemblers  
140 assessed here. This provided ratios of over- or under-performance relative to the average assembly  
141 statistic (-1 to +4). Figure 3 was generated using log-transformed data for each assembly statistic of  
142 relevance to ensure concise representation of the data.

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

145

## 146 **Results**

### 147 *Metagenome data and dataset complexity*

148 Using Nonpareil, we confirmed that the soil metagenomes were more complex (less redundant) than  
149 the aquatic and human guts metagenomes, which were the least complex (highly redundant; Figure  
150 1) [37-39]. All the human gut metagenomes came close to sequencing saturation (with at least 75%  
151 of the diversity sequenced; Figure 1). The infant gut metagenome was sequenced to above 90%  
152 estimated average coverage (~94%). However, all the sequencing depths reached were insufficient  
153 to describe the complete spectrum of microbial members in the samples assessed. For example,  
154 the largest metagenome assessed here, the Iowa soil metagenome, only described 48.8% of the  
155 total microbial diversity despite the utilization of 47 Gbp of sequence data.

156 Estimates of the number of microbial species per gram of soil still vary substantially, with values  
157 ranging from 2000 [41] to more than 830000 [37]. These estimates do not include eukaryotic  
158 microbes, which generally possess much larger genomes and are much more difficult to fully  
159 sequence [42]. We note the published predictions that 2-5 Gbp of sequence data would fully capture  
160 an entire natural microbial community [40]. Based on our analysis, we propose that the sequencing  
161 depth required to provide comprehensive coverage of soil metagenomes should be increased by an  
162 order of magnitude, to ~100 Gbp. This is a function of the extreme taxonomic heterogeneity of soil  
163 microbial communities, and highlights the challenge of assembling low coverage metagenomes.

### 164 *Strategy and approaches of the current research*

165 We defined five measures to assess the performance of each metagenomic assembler tested; (1)  
166 ease of use and assembler attributes, (2) quality of assemblies generated and computational  
167 requirements, (3) influence of sequencing depth and coverage, (4) suitability to different  
168 environments and (5) their performance on metagenomes of known composition.

169 *1. Ease of use and assembler attributes*

170 Many researchers entering the field of metagenomics are inexperienced in the use of intricate  
171 bioinformatic tools, and may lack extensive computational resources. To assess the ease of use for  
172 inexperienced computational biologists we evaluated the availability of a web application or graphical  
173 user interface (GUI), ease of installation, availability and completeness of manuals, Message  
174 Passing Interface (MPI) compatibility and programming language.

175 Eight of the assemblers tested here use command-line interfaces (CLI) and are open-source  
176 freeware (Velvet, MetaVelvet, SPAdes, metaSPAdes, Ray Meta, IDBA-UD, MEGAHIT and Omega).  
177 Only the commercial software CLC Genomics Workbench (Qiagen) implements a GUI  
178 (Supplementary Table 1). CLC is easily installed on most Linux, Windows or MacOS computers,  
179 whereas all other assemblers are limited to Unix-based operating systems. The GUI is intuitive, and  
180 users can assemble simply by using a point-and-click interface. CLC provides substantial support  
181 (via manuals and web based tutorials) and was the most user-friendly assembler tested here.

182 Unix-based assemblers are inherently more difficult to use and must be installed or compiled from  
183 source code using the CLI. All assemblers that are CLI-based can be downloaded from GitHub, while  
184 some tools (SPAdes, metaSPAdes, Ray Meta, Velvet, MetaVelvet and Omega) provide download  
185 links from their respective parent websites. All tools, barring SPAdes, metaSPAdes and IDBA-UD,  
186 provide MPI compatibility, allowing parallelization which reduces computational time. All tools  
187 assessed here provide manuals or 'readme' files either on their websites or GitHub repositories,  
188 although others, such as IDBA-UD, MetaVelvet and Omega, are not comprehensive and lack  
189 information on installation or implementation. Tools with more complete manuals (MEGAHIT and  
190 Ray Meta) feature extensive wiki pages and frequently asked questions. The number of citations,  
191 websites, programming languages and MPI compatibility of all the tools assessed are provided in  
192 Supplementary Table 1.

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

194 Evaluating metagenome assembly quality is challenging without the use of known reference  
195 genomes for diverse microbial communities. We compared assembly quality using many standard  
196 metrics, including the total number of contigs longer than 500 bp, 1 kbp (referred to as long contigs  
197 throughout) and 50 kbp (referred to as ultra-long contigs throughout), maximum contig length, *N50*  
198 length of the contigs (length of the median contig, representing the length of the smallest contig at

199 which half of the assembly is represented), mapping rate and assembly span (total length assembled  
200 using contigs > 500 bp). We used MetaQUAST to evaluate these assembly quality statistics [22].

201 We selected the Tara Ocean metagenome [29] for a comparison of each assembler at *k*-mer lengths  
202 between 50 and 61. We selected this range as the assemblers which automatically optimize *k*-mer  
203 values generally set sizes within this range. We set the other non-optimizing assemblers to 55.  
204 Compared to the other natural metagenomes, the Tara Ocean metagenomic dataset is of  
205 intermediate complexity and sequencing depth (Figure 1, Supplementary Table 2). This  
206 metagenome was sequenced on an Illumina HiSeq instrument, which is currently the most widely  
207 used shotgun metagenome sequencing technology [3]. This 5.4 Gbp metagenome comprised more  
208 than 27 million high-quality read pairs with a mean read pair length of 200.3 bp (Supplementary  
209 Table 2).

210 Omega (2691), SPAdes (1415), Ray Meta (1329), IDBA-UD (1166) and metaSPAdes (1124)  
211 provided assemblies with high *N50* values (> 1000 bp), while the assemblies generated using CLC,  
212 MEGAHIT, Velvet and MetaVelvet produced *N50* statistics below 1000 bp (Figure 2; Table 1).  
213 Overall, the assembly spans varied considerably with SPAdes (275.9 Mbp), MEGAHIT (210.6 Mbp),  
214 metaSPAdes (202.8 Mbp) and IDBA-UD (179.7 Mbp) assembling the largest metagenomes.  
215 Assembly span was correlated with the number of reads mapping back to the assemblies ( $R^2 = 0.83$ ;  
216 Supplementary Figure 3), with SPAdes and metaSPAdes having the highest values (Table 1). Both  
217 IDBA-UD and MEGAHIT mapped back more than 50% of the sequence reads to the assemblies.  
218 SPAdes also produced the most contigs over 1 kbp (70711), while MEGAHIT, IDBA-UD and  
219 metaSPAdes created fewer contigs in that size range, but all were comparable to each other  
220 (between 48640 and 56243 contigs). The largest contig was assembled by SPAdes (197 kbp),  
221 followed by metaSPAdes (142 kbp), Omega (102 kbp) and IDBA-UD (101 kbp). These three  
222 assemblers also produced the most 'ultra-long' contigs (> 50 kbp); with 54, 37 and 2 contigs,  
223 respectively.

224 The computational requirements of an assembly tool should be a major consideration when selecting  
225 an assembler. We evaluated all assemblers in relation to the time taken to assemble the Tara Ocean  
226 metagenome (Supplementary Figure 2; Table 1) using the same number of threads ( $n=8$ ;  
227 Supplementary Figure 2A). Velvet, MetaVelvet and CLC assembled the metagenome in less than  
228 an hour, while MEGAHIT and Ray Meta were substantially slower, assembling over multiple hours.  
229 IDBA-UD, SPAdes and metaSPAdes required considerably more time to complete assembly, taking  
230 approximately 24 hours, or more. Omega required the most time to assemble the metagenome,  
231 taking approximately 48 hours. In terms of memory requirements, SPAdes was the most 'memory  
232 expensive' (157 GB of RAM), followed by Velvet and MetaVelvet (both 109 GB), which is  
233 substantially more RAM than is available on an average desktop computer (16 GB). By contrast,  
234 MEGAHIT (11 GB) and CLC (16 GB) were the most memory efficient assemblers (Figure 2 and  
235 Supplementary Figure 3; Table 1).

236 Overall, SPAdes, metaSPAdes, IDBA-UD and MEGAHIT displayed the best performances in  
237 assembling this metagenome of intermediate size and complexity, as they produced very high *N50*  
238 values, a high proportion of long contigs and the widest assembly spans. While SPAdes was the  
239 best assembler overall, MEGAHIT was the most memory efficient, as it produced an assembly  
240 comparable to the best performing assemblers while using only a fraction of computational  
241 resources.

### 242 3. *Benchmarking influence of sequencing depth and coverage*

243 Temperate soil communities are generally more diverse than extreme counterparts (e.g., permafrost;  
244 Supplementary Table 3, Figure 1) [11]. Subsequently, high levels of diversity within these biomes  
245 require much deeper sequencing effort. Differences in microorganism abundances and strain level  
246 heterogeneity introduce complications during metagenome assembly, resulting in increased memory  
247 requirements and longer computational run-times, which may challenge assemblers. The two  
248 temperate soil metagenomes assessed here have vastly different sequencing depths, thus providing  
249 us with the scope to assess the influence of sequencing depth on the performance of each  
250 assembler. The Oklahoma soil metagenome [26] had a low sequencing depth (9 Gbp) and estimated  
251 coverage (11%), c.f. the Iowa soil metagenome [8], which had a very high sequencing depth (47  
252 Gbp) and 49% estimated coverage (Figure 1, Supplementary Table 2). We predicted that deeper  
253 sequencing effort would be correlated with an increase in metagenome coverage [34].

254 All assemblers successfully assembled the Oklahoma metagenome, although SPAdes required  
255 considerably more memory (up to 1 TB RAM, Supplementary Table 3). Nevertheless, SPAdes  
256 produced the best assembly statistics for most categories (9548 long contigs and an assembly span  
257 of 54.3 Mbp; Supplementary Table 3; Figure 3). IDBA-UD and MEGAHIT used less than 500 GB of  
258 RAM and were comparable in performance (3828 and 3416 long contigs, and assembly spans of  
259 17.2 Mbp and 20.2 Mbp, respectively; Supplementary Table 3; Figure 3). It is noteworthy that while  
260 metaSPAdes was one of the best performing assemblers for the Tara Ocean metagenome (Figure  
261 2), it performed poorly here (Supplementary Table 3; Figure 3), suggesting that metaSPAdes is ill-  
262 suited to assembling low coverage metagenomes.

263 The massive Iowa soil metagenome could not be assembled by either SPAdes or IDBA-UD using  
264 our available computing resources (1 TB of RAM). This is in agreement with the methodology  
265 described by the authors who generated this dataset, who digitally normalized and partitioned the  
266 Iowa metagenome to allow for assembly using Velvet [8]. Remarkably, MEGAHIT and CLC  
267 assembled the Iowa metagenome using less than 500 GB of RAM. MEGAHIT performed best across  
268 most categories tested (assembly span of 1036.5 Mbp, largest contig of 104841 bp, and 277623  
269 long contigs; Figure 3), while CLC produced the third-best assembly (assembly span of 432.7 Mbp,  
270 largest contig of 70207 and 114196 long contigs), using less than 64GB of memory. MetaSPAdes  
271 performed comparably to MEGAHIT but had much higher computational resource requirements to



272 assembly the Iowa soil metagenome, using up to 1TB of RAM (assembly span of 873.8 Mbp, largest  
273 contig of 188499 bp, and 225046 long contigs).

274 Overall, we found that sequencing depth greatly influenced the performance of the assemblers,  
275 although the most memory-efficient tools, MEGAHIT and CLC, performed well irrespective of  
276 sequencing coverage. SPAdes and IDBA-UD produced good assemblies for the Oklahoma soil  
277 metagenome, but were extremely expensive in terms of memory and failed to assemble the Iowa  
278 soil metagenome. We found that metaSPAdes produced a better assembly for the Iowa soil  
279 metagenome than the Oklahoma dataset. MetaSPAdes performed optimally for the assembly of the  
280 high-coverage metagenome, but was less efficient in the assembly of the low-coverage  
281 metagenome.

#### 282 *4. Benchmarking suitability to various environments*

283 Environmental samples are widely dissimilar in microbial community complexity and have distinct  
284 taxonomic compositions. In this study, we assembled metagenomes from three environmental  
285 biomes of different phylotypic complexities. Overall, SPAdes, MEGAHIT, IDBA-UD and metaSPAdes  
286 assembled most of the metagenomes well, according to the parameters we evaluated  
287 (Supplementary Tables 3-5). SPAdes consistently provided the largest contigs and the widest  
288 assembly spans. MEGAHIT demanded far fewer computational resources, and yet produced similar  
289 assemblies to metaSPAdes and IDBA-UD. CLC provided assemblies of moderate to high quality,  
290 was the easiest to use and performed particularly well on large metagenomes. Together, these  
291 results indicate that no single assembler performs best across all sequencing platforms and  
292 datasets.

#### 293 *5. Benchmarking on synthetic metagenomes*

294 As previously indicated, assessing metagenome assembler performance is complicated due to the  
295 unknown composition of environmental microbial communities. To overcome this challenge, we  
296 included three synthetic metagenomes of known composition to assess the error rates (such as  
297 number of indels, misassemblies, and ambiguous bases) generated by each assembler. These three  
298 metagenomes represented three discreet complexities (low, medium and high; Supplementary  
299 Figure 1), in order to challenge the assemblers with the unique properties of each dataset.

300 Our analysis show that more complex metagenomes led to higher error rates in the resultant  
301 assemblies (Figure 4). Notably, SPAdes produced the most misassemblies (643, 4928 and 77264  
302 for the assemblies of low, medium and high complexity synthetic metagenomes, respectively) and  
303 the highest unaligned lengths (46 kbp, 891 kbp and 19 Mbp, respectively). IDBA-UD produced a  
304 high number of misassemblies while Omega consistently produced the most mismatches in all  
305 synthetic datasets (more than 1500 mismatches per 100 kbp for all synthetic metagenomes). CLC  
306 and Ray Meta consistently produced more than 100 ambiguous bases (N's) per 100 kbp in each of

307 the generated synthetic assemblies. Finally, CLC also incorporated the most indels per 100 kbp in  
308 all complexity classes (more than double the number of indels produced by any other assembler).

309

### 310 *How to select a metagenome assembler*

311 Bioinformatics projects can be limited by memory (RAM) requirements. SPAdes, metaSPAdes,  
312 IDBA-UD, Velvet and MetaVelvet all have large memory requirements during the assembly of  
313 massive datasets. MEGAHIT, Omega and CLC are extremely memory efficient, as they required  
314 less than 500 GB of RAM to assemble the massive Iowa soil metagenome. MEGAHIT, for example,  
315 generates succinct *de Bruijn* graphs to achieve efficient memory usage [20].

316 Our results indicate that although many assemblers perform comparably, their applicability is defined  
317 by the research question at hand. SPAdes, for example, generated good assemblies with the most  
318 long and ultra-long contigs for most datasets. These are ideal characteristics for genome-centric  
319 studies, which require the binning of draft genomes from community sequence data [43]. By contrast,  
320 metaSPAdes considers read coverage during assembly, making it more applicable for microbial  
321 community profiling [17]. While SPAdes and metaSPAdes produced the best assemblies in general,  
322 MEGAHIT performed comparably and emerged as a rapid and memory efficient alternative  
323 assembler.

324 However, it should be noted that SPAdes and IDBA-UD generate high numbers of misassemblies  
325 and contigs that do not align to the reference genomes. Other assemblers such as Omega, CLC and  
326 Ray Meta each have unique error profiles, which should be considered in light of the research  
327 questions asked. For example, when assessing strain level genomic variations (SNP's), assemblers  
328 that generate high numbers of indels and mismatches should be avoided. In addition, while SPAdes  
329 generates many mismatches, if the aim is to extract single genomes from a metagenome, manual  
330 curation of the newly re-constructed draft genomes will identify and correct such misassemblies.

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

### 339 **Other analyses**

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

343

## 344 **Discussion**

345 Over the last decade, high throughput sequencing has revolutionised the field of microbial ecology  
346 [44]. Amplicon-based technologies have allowed for near-complete classification of whole microbial  
347 communities, including populations of bacteria, archaea and fungi [45]. The emergence of two key  
348 platforms for analysing amplicon sequencing data, mothur [46] and QIIME [47], has allowed for  
349 methodological standards to be set, which enables robust comparisons between studies [48].

350 While whole community shotgun metagenome sequencing has facilitated the in-depth description of  
351 microbial communities from diverse environments, such as the human gut [49] and acid mine  
352 drainage systems [50], no standards exist with regard to assembly platforms or their use. While  
353 numerous reviews on strategies to analyse metagenomic data have been published [51], there are  
354 currently no standard assembly procedures implemented to enable thorough comparative analyses  
355 between projects. Numerous pipelines for processing metagenomic sequence data are available.  
356 These typically integrate existing tools into a single workflow for rapid, standardized analysis (e.g.,  
357 MG-RAST, MetAMOS, and IMG/M) [52-54]. However, few of these pipelines are as widely used as  
358 mothur or QIIME in barcoding studies. This is partly because integrated metagenome analysis tools,  
359 such as MetAMOS, do not achieve the flexibility afforded by using each tool individually (e.g., using  
360 separate tools for assembly, binning and taxonomic assignment).

361 Consequently, investigators can analyse unassembled reads [11], optimize their assembly  
362 parameters or even develop their own tools to assemble their data prior to further analysis [55].  
363 However, within the scope of metagenome assembly, essential details are often omitted when  
364 describing methods [56]. This leads to methodological discrepancies, and severely limits the  
365 possibility of making routine, robust comparisons between studies. This issue was recently  
366 highlighted by J Vollmers, S Wiegand and A-K Kaster [21] and WW Greenwald, N Klitgord, V  
367 Seguritan, S Yooseph, JC Venter, C Garner, KE Nelson and W Li [57] who reported that the  
368 taxonomic diversity patterns of microbial communities differed substantially, depending on the  
369 assembler used. While some recent studies have applied single cell sequencing [58] or chromosome  
370 capture [59] approaches to enhance metagenome assembly, these techniques remain inaccessible  
371 to most researchers. We provide an evaluation of commonly-used assemblers on standard shotgun  
372 sequenced metagenomes.

373 In our comparative analyses of the most popular assembly platforms, SPAdes produced the most  
374 long contigs, independent of the metagenome origin. However, this assembler introduced a large  
375 number of misassemblies in high complexity datasets. SPAdes is ideal for genome-centric research

376 questions that require long and ultra-long contigs, such as those that aim to bin and reconstruct  
377 single genomes from shotgun metagenomes [16]. By contrast, MEGAHIT and metaSPAdes provided  
378 very large assembly spans and consider sequence coverage during assembly, reducing the number  
379 of misassemblies generated. IDBA-UD also produced large assembly spans and a high number of  
380 contigs, but at the cost of generating misassemblies for complex datasets. These tools are thus more  
381 appropriate for research questions related to taxonomic profiling of natural microbial communities,  
382 for functionally annotating microbial communities, for the analysis of population scale dynamics or  
383 for comparison of microbial communities across biomes [17, 20]. By analysing metagenomes of  
384 known composition and complexity, we found that each assembler tested here generated a unique  
385 error profile (e.g., IDBA-UD produces many misassemblies, CLC produces many indels and Omega  
386 produces many mismatches). As mentioned above, this excluded some assemblers from specific  
387 research objectives (e.g., using CLC for variant calling). This reiterates the fact that the research  
388 question should be the primary consideration when selecting the appropriate assembler, and that  
389 these assembler-specific drawbacks should also be considered.

390 Overall, MEGAHIT produced some of the best assemblies throughout this study, while only using a  
391 fraction of the computational resources required by other assemblers. We strongly recommend  
392 MEGAHIT for researchers who do not have access to large computational resources. Finally, the  
393 CLC assembler is ideal for researchers who lack a depth of bioinformatic knowledge, or who prefer  
394 to use a GUI and are willing to invest in software which is easier to use. CLC is easy to install, has  
395 an intuitive interface and provides a compromise in which assembly quality may be sacrificed for  
396 ease of use. Strikingly, the most widely cited assembler assessed here (Velvet cited 5974 times;  
397 Supplementary Table 1) did not perform well across most metagenomes, while scarcely cited  
398 platforms (MEGAHIT, metaSPAdes cited 114 and 18 times, respectively; Supplementary Table 1)  
399 performed well across most statistics assessed here.

400

## 401 **Conclusions**

402 No assembler tested here consistently provided superior assemblies across the different  
403 metagenomes. Consequently, we propose a viable methodology for the selection of an appropriate  
404 assembler, dictated by (1) the scientific research question posed, then by (2) the computational  
405 resources available, and (3) the bioinformatics skill level of the researcher (Figure 5). In light of the  
406 above proposed framework, we urge researchers to carefully consider the assembler used (as well  
407 as the entire bioinformatics pipeline followed) while specifically bearing in mind their research  
408 question and what feature of the dataset they want accentuated.

409

## 410 **List of abbreviations**

411 contigs - contiguous segments  
412 SRA - sequence read archive  
413 MG-RAST - Metagenomics Rapid Annotation Server  
414 CHPC - Centre for High Performance Computing  
415 RAM random access memory  
416 GUI - graphical user interface  
417 MPI - Message Passing Interface  
418 CLI - command-line interface  
419 SNP – single nucleotide polymorphism

420

## 421 **Declarations**

### 422 **Ethics approval and consent to participate**

423 Not applicable

### 424 **Consent for publication**

425 Not applicable

### 426 **Availability of data and material**

427 The datasets analysed during the current study are available in the Sequence Read Archive  
428 (SRA) repository, under accession numbers SRR351474, SRR958082, ERR598950,  
429 ERR526087, SRR341725 and ERR732914. As well as the MG-RAST repository, under  
430 accession numbers 4470009.3 - 4470010.3, 4673644.3 - 4673645.3 and 4537104.3 -  
431 4537105.3. Synthetic simulated metagenomes are available from [https://data.cami-  
432 challenge.org](https://data.cami-challenge.org).

### 433 **Competing interests**

434 The authors declare that they have no competing interests.

### 435 **Funding**

436 The research was funded by the National Research Foundation under the grant numbers  
437 102910 (AJvdW) and 97891 (MWVG) and the University of Pretoria (Genomics Research  
438 Institute). Funding bodies had no influence on the study design, data collection, analysis,  
439 interpretation of data and writing of the manuscript.

440 **Authors' contributions**

441 AJvdW contributed most to research, analysis, formulation of ideas and writing. MWVG  
442 contributed research, analysis, formulation of ideas and writing. JBR contributed formulation  
443 of ideas, direction and writing. TPM contributed formulation of ideas and writing. OR  
444 contributed direction and writing. DAC contributed formulation of ideas, direction and writing.

445 **Acknowledgements**

446 We thank Dr Surendra Vikram for helping with the installation of assemblers, as well as Dane  
447 Kennedy from the CHPC, South Africa with his assistance in accessing the large memory  
448 nodes of the CHPC. We gratefully acknowledge support from Ms Amanda van der Walt help  
449 with improving the quality of our figures. Finally, we thank S. de Scally and A. E. Visser for  
450 their support and helpful discussions.

451

452

453 **References**

- 454 1. Riesenfeld CS, Schloss PD, Handelsman J: **Metagenomics: genomic analysis of microbial**  
455 **communities.** *Annu Rev Genet* 2004, **38**:525-552.
- 456 2. Xu L, Chen H, Hu X, Zhang R, Zhang Z, Luo Z: **Average gene length is highly conserved in**  
457 **prokaryotes and eukaryotes and diverges only between the two kingdoms.** *Molecular biology and*  
458 *evolution* 2006, **23**(6):1107-1108.
- 459 3. Clooney AG, Fouhy F, Sleator RD, O'Driscoll A, Stanton C, Cotter PD, Claesson MJ: **Comparing**  
460 **Apples and Oranges?: Next Generation Sequencing and Its Impact on Microbiome Analysis.** *PLoS*  
461 *one* 2016, **11**(2):e0148028.
- 462 4. Nagarajan N, Pop M: **Sequence assembly demystified.** *Nature Reviews Genetics* 2013, **14**(3):157-  
463 167.
- 464 5. Anantharaman K, Brown CT, Hug LA, Sharon I, Castelle CJ, Probst AJ, Thomas BC, Singh A, Wilkins  
465 MJ, Karaoz U: **Thousands of microbial genomes shed light on interconnected biogeochemical**  
466 **processes in an aquifer system.** *Nature Communications* 2016, **7**:13219.
- 467 6. Compeau PE, Pevzner PA, Tesler G: **How to apply de Bruijn graphs to genome assembly.** *Nature*  
468 *biotechnology* 2011, **29**(11):987-991.
- 469 7. Loman NJ, Constantinidou C, Christner M, Rohde H, Chan JZ-M, Quick J, Weir JC, Quince C, Smith  
470 GP, Betley JR: **A culture-independent sequence-based metagenomics approach to the**  
471 **investigation of an outbreak of Shiga-toxigenic Escherichia coli O104: H4.** *Jama* 2013,  
472 **309**(14):1502-1510.
- 473 8. Howe AC, Jansson JK, Malfatti SA, Tringe SG, Tiedje JM, Brown CT: **Tackling soil diversity with the**  
474 **assembly of large, complex metagenomes.** *Proceedings of the National Academy of Sciences* 2014,  
475 **111**(13):4904-4909.
- 476 9. Handelsman J: **Metagenomics: application of genomics to uncultured microorganisms.**  
477 *Microbiology and molecular biology reviews* 2004, **68**(4):669-685.
- 478 10. Tyson GW, Chapman J, Hugenholtz P, Allen EE, Ram RJ, Richardson PM, Solovyev VV, Rubin EM,  
479 Rokhsar DS, Banfield JF: **Community structure and metabolism through reconstruction of**  
480 **microbial genomes from the environment.** *Nature* 2004, **428**(6978):37-43.
- 481 11. Fierer N, Leff JW, Adams BJ, Nielsen UN, Bates ST, Lauber CL, Owens S, Gilbert JA, Wall DH,  
482 Caporaso JG: **Cross-biome metagenomic analyses of soil microbial communities and their**

- 483 **functional attributes. *Proceedings of the National Academy of Sciences* 2012, **109**(52):21390-**  
484 **21395.**
- 485 12. Wurch L, Giannone RJ, Belisle BS, Swift C, Utturkar S, Hettich RL, Reysenbach A-L, Podar M:  
486 **Genomics-informed isolation and characterization of a symbiotic Nanoarchaeota system from a**  
487 **terrestrial geothermal environment. *Nature Communications* 2016, **7**.**
- 488 13. Baruzzo G, Hayer KE, Kim EJ, Di Camillo B, FitzGerald GA, Grant GR: **Simulation-based**  
489 **comprehensive benchmarking of RNA-seq aligners. *Nature Methods* 2016.**
- 490 14. Zerbino DR, Birney E: **Velvet: algorithms for de novo short read assembly using de Bruijn graphs.**  
491 ***Genome research* 2008, **18**(5):821-829.**
- 492 15. Namiki T, Hachiya T, Tanaka H, Sakakibara Y: **MetaVelvet: an extension of Velvet assembler to de**  
493 **novu metagenome assembly from short sequence reads. *Nucleic acids research* 2012,**  
494 ****40**(20):e155-e155.**
- 495 16. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin VM, Nikolenko SI, Pham  
496 S, Prjibelski AD: **SPAdes: a new genome assembly algorithm and its applications to single-cell**  
497 **sequencing. *Journal of Computational Biology* 2012, **19**(5):455-477.**
- 498 17. Nurk S, Meleshko D, Korobeynikov A, Pevzner P: **metaSPAdes: a new versatile de novo**  
499 **metagenomics assembler. *arXiv preprint arXiv:160403071* 2016.**
- 500 18. Boisvert S, Raymond F, Godzaridis É, Laviolette F, Corbeil J: **Ray Meta: scalable de novo**  
501 **metagenome assembly and profiling. *Genome biology* 2012, **13**(12):1.**
- 502 19. Peng Y, Leung HC, Yiu S-M, Chin FY: **IDBA-UD: a de novo assembler for single-cell and**  
503 **metagenomic sequencing data with highly uneven depth. *Bioinformatics* 2012, **28**(11):1420-1428.**
- 504 20. Li D, Liu C-M, Luo R, Sadakane K, Lam T-W: **MEGAHIT: an ultra-fast single-node solution for large**  
505 **and complex metagenomics assembly via succinct de Bruijn graph. *Bioinformatics* 2015:btv033.**
- 506 21. Vollmers J, Wiegand S, Kaster A-K: **Comparing and Evaluating Metagenome Assembly Tools from a**  
507 **Microbiologist's Perspective-Not Only Size Matters! *PloS one* 2017, **12**(1):e0169662.**
- 508 22. Mikheenko A, Saveliev V, Gurevich A: **MetaQUAST: evaluation of metagenome assemblies.**  
509 ***Bioinformatics* 2015:btv697.**
- 510 23. Vázquez-Castellanos JF, García-López R, Pérez-Brocal V, Pignatelli M, Moya A: **Comparison of**  
511 **different assembly and annotation tools on analysis of simulated viral metagenomic communities**  
512 **in the gut. *BMC genomics* 2014, **15**(1):37.**
- 513 24. Pignatelli M, Moya A: **Evaluating the fidelity of de novo short read metagenomic assembly using**  
514 **simulated data. *PloS one* 2011, **6**(5):e19984.**
- 515 25. Charuvaka A, Rangwala H: **Evaluation of short read metagenomic assembly. *BMC genomics* 2011,**  
516 ****12**(2):S8.**
- 517 26. Luo C, Rodriguez-R LM, Johnston ER, Wu L, Cheng L, Xue K, Tu Q, Deng Y, He Z, Shi JZ: **Soil microbial**  
518 **community responses to a decade of warming as revealed by comparative metagenomics.**  
519 ***Applied and environmental microbiology* 2014, **80**(5):1777-1786.**
- 520 27. Hultman J, Waldrop MP, Mackelprang R, David MM, McFarland J, Blazewicz SJ, Harden J, Turetsky  
521 MR, McGuire AD, Shah MB: **Multi-omics of permafrost, active layer and thermokarst bog soil**  
522 **microbiomes. *Nature* 2015.**
- 523 28. Bowman JS, Berthiaume CT, Armbrust EV, Deming JW: **The genetic potential for key**  
524 **biogeochemical processes in Arctic frost flowers and young sea ice revealed by metagenomic**  
525 **analysis. *FEMS microbiology ecology* 2014, **89**(2):376-387.**
- 526 29. Sunagawa S, Coelho LP, Chaffron S, Kultima JR, Labadie K, Salazar G, Djahanschiri B, Zeller G, Mende  
527 DR, Alberti A: **Structure and function of the global ocean microbiome. *Science* 2015,**  
528 ****348**(6237):1261359.**
- 529 30. Tremaroli V, Karlsson F, Werling M, Ståhlman M, Kovatcheva-Datchary P, Olbers T, Fändriks L, le  
530 Roux CW, Nielsen J, Bäckhed F: **Roux-en-Y gastric bypass and vertical banded gastroplasty induce**  
531 **long-term changes on the human gut microbiome contributing to fat mass regulation. *Cell***  
532 ***metabolism* 2015, **22**(2):228-238.**
- 533 31. Karlsson FH, Tremaroli V, Nookaew I, Bergström G, Behre CJ, Fagerberg B, Nielsen J, Bäckhed F: **Gut**  
534 **metagenome in European women with normal, impaired and diabetic glucose control. *Nature***  
535 **2013, **498**(7452):99-103.**

- 536 32. Bäckhed F, Roswall J, Peng Y, Feng Q, Jia H, Kovatcheva-Datchary P, Li Y, Xia Y, Xie H, Zhong H:  
537 **Dynamics and stabilization of the human gut microbiome during the first year of life.** *Cell host &*  
538 *microbe* 2015, **17**(5):690-703.
- 539 33. Schmieder R, Edwards R: **Quality control and preprocessing of metagenomic datasets.**  
540 *Bioinformatics* 2011, **27**(6):863-864.
- 541 34. Rodriguez-R LM, Konstantinidis KT: **Nonpareil: a redundancy-based approach to assess the level of**  
542 **coverage in metagenomic datasets.** *Bioinformatics* 2014, **30**(5):629-635.
- 543 35. Langmead B, Salzberg SL: **Fast gapped-read alignment with Bowtie 2.** *Nature methods* 2012,  
544 **9**(4):357-359.
- 545 36. Galili T: **heatmaply: interactive heat maps (with R).** *Month* 2016.
- 546 37. Gans J, Wolinsky M, Dunbar J: **Computational improvements reveal great bacterial diversity and**  
547 **high metal toxicity in soil.** *Science* 2005, **309**(5739):1387-1390.
- 548 38. Torsvik V, Øvreås L, Thingstad TF: **Prokaryotic diversity--magnitude, dynamics, and controlling**  
549 **factors.** *Science* 2002, **296**(5570):1064-1066.
- 550 39. Tringe SG, Von Mering C, Kobayashi A, Salamov AA, Chen K, Chang HW, Podar M, Short JM, Mathur  
551 EJ, Detter JC: **Comparative metagenomics of microbial communities.** *Science* 2005, **308**(5721):554-  
552 557.
- 553 40. Venter JC, Remington K, Heidelberg JF, Halpern AL, Rusch D, Eisen JA, Wu D, Paulsen I, Nelson KE,  
554 Nelson W: **Environmental genome shotgun sequencing of the Sargasso Sea.** *science* 2004,  
555 **304**(5667):66-74.
- 556 41. Schloss PD, Handelsman J: **Toward a census of bacteria in soil.** *PLoS computational biology* 2006,  
557 **2**(7):e92.
- 558 42. Lynch M, Conery JS: **The origins of genome complexity.** *science* 2003, **302**(5649):1401-1404.
- 559 43. Becraft ED, Dodsworth JA, Murugapiran SK, Ohlsson JI, Briggs BR, Kanbar J, De Vlaminck I, Quake  
560 SR, Dong H, Hedlund BP: **Single-Cell-Genomics-Facilitated Read Binning of Candidate Phylum**  
561 **EM19 Genomes from Geothermal Spring Metagenomes.** *Applied and environmental microbiology*  
562 2016, **82**(4):992-1003.
- 563 44. Shokralla S, Spall JL, Gibson JF, Hajibabaei M: **Next-generation sequencing technologies for**  
564 **environmental DNA research.** *Molecular ecology* 2012, **21**(8):1794-1805.
- 565 45. Taberlet P, Coissac E, Pompanon F, Brochmann C, Willerslev E: **Towards next-generation**  
566 **biodiversity assessment using DNA metabarcoding.** *Molecular ecology* 2012, **21**(8):2045-2050.
- 567 46. Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, Lesniewski RA, Oakley BB,  
568 Parks DH, Robinson CJ *et al*: **Introducing mothur: open-source, platform-independent,**  
569 **community-supported software for describing and comparing microbial communities.** *Appl*  
570 *Environ Microbiol* 2009, **75**(23):7537-7541.
- 571 47. Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman FD, Costello EK, Fierer N, Pena AG,  
572 Goodrich JK, Gordon JI: **QIIME allows analysis of high-throughput community sequencing data.**  
573 *Nature methods* 2010, **7**(5):335-336.
- 574 48. Plummer E, Twin J, Bulach DM, Garland SM, Tabrizi SN: **A Comparison of Three Bioinformatics**  
575 **Pipelines for the Analysis of Preterm Gut Microbiota using 16S rRNA Gene Sequencing Data.**  
576 *Journal of Proteomics & Bioinformatics* 2015, **8**(12):283.
- 577 49. Gill SR, Pop M, DeBoy RT, Eckburg PB, Turnbaugh PJ, Samuel BS, Gordon JI, Relman DA, Fraser-  
578 Liggett CM, Nelson KE: **Metagenomic analysis of the human distal gut microbiome.** *science* 2006,  
579 **312**(5778):1355-1359.
- 580 50. Kantor RS, Huddy RJ, Iyer RM, Thomas BC, Brown CT, Anantharaman K, Tringe SG, Hettich RL,  
581 Harrison ST, Banfield JF: **Genome-resolved meta-omics ties microbial dynamics to process**  
582 **performance in biotechnology for thiocyanate degradation.** *Environmental Science & Technology*  
583 2017.
- 584 51. Scholz MB, Lo C-C, Chain PS: **Next generation sequencing and bioinformatic bottlenecks: the**  
585 **current state of metagenomic data analysis.** *Current opinion in biotechnology* 2012, **23**(1):9-15.
- 586 52. Markowitz VM, Chen I-MA, Chu K, Szeto E, Palaniappan K, Pillay M, Ratner A, Huang J, Pagani I,  
587 Tringe S: **IMG/M 4 version of the integrated metagenome comparative analysis system.** *Nucleic*  
588 *Acids Research* 2014, **42**(D1):D568-D573.



- 589 53. Glass EM, Wilkening J, Wilke A, Antonopoulos D, Meyer F: **Using the metagenomics RAST server**  
590 **(MG-RAST) for analyzing shotgun metagenomes.** *Cold Spring Harbor Protocols* 2010, **2010**(1):pdb.  
591 prot5368.
- 592 54. Treangen TJ, Koren S, Sommer DD, Liu B, Astrovskaya I, Ondov B, Darling AE, Phillippy AM, Pop M:  
593 **MetAMOS: a modular and open source metagenomic assembly and analysis pipeline.** *Genome*  
594 *Biol* 2013, **14**(1):R2.
- 595 55. Oh J, Byrd AL, Deming C, Conlan S, Kong HH, Segre JA, Program NCS: **Biogeography and**  
596 **individuality shape function in the human skin metagenome.** *Nature* 2014, **514**(7520):59-64.
- 597 56. Brown CT, Hug LA, Thomas BC, Sharon I, Castelle CJ, Singh A, Wilkins MJ, Wrighton KC, Williams KH,  
598 Banfield JF: **Unusual biology across a group comprising more than 15% of domain Bacteria.** *Nature*  
599 2015, **523**(7559):208-211.
- 600 57. Greenwald WW, Klitgord N, Seguritan V, Yooseph S, Venter JC, Garner C, Nelson KE, Li W:  
601 **Utilization of defined microbial communities enables effective evaluation of meta-genomic**  
602 **assemblies.** *BMC genomics* 2017, **18**(1):296.
- 603 58. Ji P, Zhang Y, Wang J, Zhao F: **MetaSort untangles metagenome assembly by reducing microbial**  
604 **community complexity.** *Nature communications* 2017, **8**.
- 605 59. Marbouty M, Baudry L, Cournac A, Koszul R: **Scaffolding bacterial genomes and probing host-virus**  
606 **interactions in gut microbiome by proximity ligation (chromosome capture) assay.** *Science*  
607 *Advances* 2017, **3**(2):e1602105.
- 608

609 **Figure and Table legends.**

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

613 **Figure 1.** Nonpareil estimates of sequence coverage (redundancy) for the 9 metagenomes studied.  
614 Metagenomes are grouped according to their environmental niche, red colours indicate soil  
615 metagenomes, blue colours indicate aquatic metagenomes and green colours are used for human  
616 gut metagenomes. Sequencing effort is indicated in base pairs on a log scale and the estimated  
617 coverage achieved is shown as a fraction of 1.

618 **Figure 2.** Heatmap displaying the assembly statistics measured and computational resources used  
619 by the nine tested assemblers on the Tara Ocean metagenome. Well performing statistics are shown  
620 in yellow, while dark blue regions indicate poor performance. Clustering of assemblers and assembly  
621 statistics was done using an hierarchical clustering method in R (hclust).

622 **Figure 3.** Radial plots showing assembly statistics for all metagenomes assessed as measured by  
623 the number of contigs larger than 500 bp, the total length of the assembly, the number of contigs  
624 larger than 1 kbp, the total bases calculated using only contigs larger than 1 kbp, the largest contigs,  
625 the *N50* value and for the synthetic datasets the fraction of contigs which aligned to the reference  
626 genomes provided. Metagenomes are labelled above the respective radial plots, where the first row  
627 represents the soils metagenomes, followed by aquatic, human gut and synthetic metagenomes.

628 **Figure 4.** Assembler performance on synthetics simulated datasets, measured by (a) number of  
629 misassemblies, (b) unaligned length, (c) number of unassigned bases (N's) per 100 kbp, (d) number  
630 of mismatches per 100 kbp and the number of indels per 100 kbp. These statistics represent negative  
631 assembly statistics and are a reflection of poor performance. Each assembler is indicated by different  
632 colors and the complexity of the synthetic dataset is indicated on the x-axis.

633 **Figure 5.** Proposed workflow to select a metagenome assembler based on the research question,  
634 the computational resources available and the bioinformatic expertise of the researcher.

635

636 **Additional Files**

637 **Supplementary Tables and Figures**

638 **.pdf**

639 **Title of data:**

640 **Supplementary Table 1. Attributes of *de novo* assemblers used in this study.** Included in this  
641 table are the versions of each assembler used in this study, along with the release date of each

642 version. We provide a link to each assemblers' website accompanied by its reference and number  
643 of citations. We gauge ease of use by providing the programming language and MPI compatibility of  
644 each tool as well as assessing the completeness of each tools' available documentation.

645 **Supplementary Table 2. Characteristics of the metagenomic datasets used in this study.**

646 Three metagenomes from three distinct environments (Soil, Aquatic and Human gut) were selected,  
647 and we provide accession numbers, sequencing platforms used and basic sequence characteristics  
648 (pre- and post-filtering) of each metagenome.

649 **Supplementary Table 3. Assembly statistics for the assembled aquatic metagenomes.**

650 **Supplementary Table 4. Assembly statistics for the assembled soil metagenomes.**

651 **Supplementary Table 5. Assembly statistics for the assembled human gut metagenomes.**

652 **Supplementary Table 6. Assembly statistics for the synthetic metagenomes.**

653 **Supplementary Figure 1. Nonpareil estimates of sequence coverage (redundancy) for the 3**  
654 **synthetic metagenomes studied.**

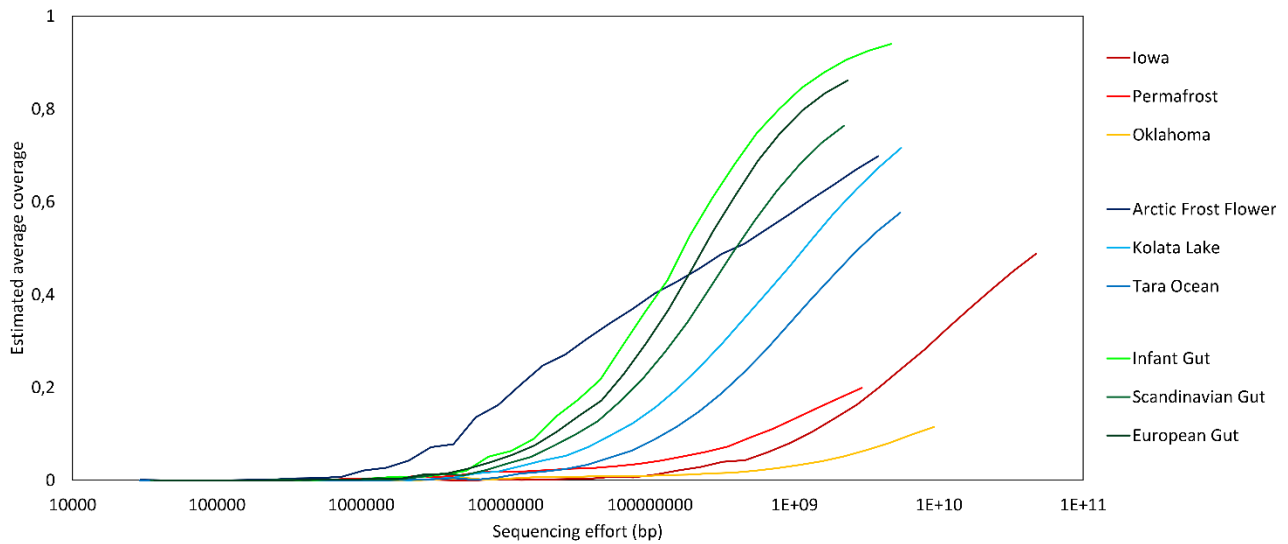
655 **Supplementary Figure 2.** Computational requirements for the Tara Ocean metagenome. A) Total  
656 assembly span proportional to wall time required. B) Total assembly span in relation to peak memory  
657 usage.

658 **Supplementary Figure 3.** Correlation between assembly span and mapping rate. The exponential  
659 trendline indicates a very strong positive correlation between the amount of data utilized and the size  
660 of the generated assembly ( $R^2 = 0.83$ ).

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

	Tara Ocean								
	CLC	IDBA-UD	MEGAHIT	metaSPAdes	MetaVelvet	Omega	Ray Meta	SPAdes	Velvet
Number of contigs ( $\geq 500$ bp)	50,716	163,815	216,938	185,419	67,161	15,982	6,128	220,178	57,816
Total length	46,069,409	179,686,756	210,621,485	202,770,058	55,972,515	34,861,819	7,277,214	275,920,632	45,425,460
No. of long contigs ( $\geq 1$ kbp)	10,720	50,498	56,243	48,640	12,590	13,305	2,179	70,711	8,802
No. of ultra-long contigs ( $\geq 50$ kbp)	0	2	1	37	0	9	0	54	0
Largest contig	39,748	101,400	62,649	141,519	30,177	102,255	41,443	197,381	21,980
N50	880	1,166	982	1,124	805	2,691	1,329	1,415	749
L50	14,113	38,236	58,246	39,033	21,544	2,737	1,345	39,617	19,631
Mapping rate (%)	38.98	52.24	55.92	64.03	4,117	13.64	8.25	64.46	48.19
Time (seconds)	3,527	69,782	10,455	125,862	2,527	168,213	16,419	80,039	2342
Time (minutes)	58.78	1,163.03	174.25	2,097.70	42.12	2803.55	273.65	1,333.98	39.03
Time (hours)	0.98	19.38	2.90	34.96	0.70	46.73	4.56	22.23	0.65
Memory required (GB)	16.23	42.84	10.58	66.53	109.37	30.7	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, red colours indicate soil metagenomes, blue colours indicate aquatic metagenomes and green colours are used for human gut metagenomes. Sequencing effort is indicated in base pairs on a log scale and the estimated coverage achieved is shown as a fraction of 1.





**Figure 3.** Radial plots showing assembly statistics for all metagenomes assessed as measured by the number of contigs larger than 500 bp, the total length of the assembly, the number of contigs larger than 1 kbp, the total bases calculated using only contigs larger than 1 kbp, the largest contigs, the *N50* value and for the synthetic datasets the fraction of contigs which aligned to the reference genomes provided. Metagenomes are labelled above the respective radial plots, where the first row represents the soils metagenomes, followed by aquatic, human gut and synthetic metagenomes.







**Figure 5.** Proposed workflow to select a metagenome assembler based on the research question, the computational resources available and the bioinformatic expertise of the researcher.

