

1 The advantages and disadvantages of short- and long-  
2 read metagenomics to infer bacterial and eukaryotic  
3 community composition

4 Keywords: metagenomics, Nanopore, Illumina, long read, community composition

5 William S. Pearman<sup>1</sup>, Nikki E. Freed<sup>1</sup>, Olin K. Silander<sup>1</sup>

6 1 School of Natural and Computational Sciences, Massey University, Auckland, New  
7 Zealand

8

9 Corresponding Authors:

10 William S Pearman

11 Olin K Silander

12 Private Bag 102904, North Shore, 0745 Auckland, New Zealand

13 Email addresses: [wpearman1996@gmail.com](mailto:wpearman1996@gmail.com); [olinsilander@gmail.com](mailto:olinsilander@gmail.com)

## 14 Abstract

### 15 Background

16 The first step in understanding ecological community diversity and dynamics is quantifying  
17 community membership. An increasingly common method for doing so is through  
18 metagenomics. Because of the rapidly increasing popularity of this approach, a large  
19 number of computational tools and pipelines are available for analysing metagenomic data.  
20 However, the majority of these tools have been designed and benchmarked using highly  
21 accurate short read data (i.e. illumina), with few studies benchmarking classification  
22 accuracy for long error-prone reads (PacBio or Oxford Nanopore). In addition, few tools have  
23 been benchmarked for non-microbial communities.

### 24 Results

25 Here we use simulated error prone Oxford Nanopore and high accuracy Illumina read sets to  
26 systematically investigate the effects of sequence length and taxon type on classification  
27 accuracy for metagenomic data from both microbial and non-microbial communities. We  
28 show that very generally, classification accuracy is far lower for non-microbial communities,  
29 even at low taxonomic resolution (e.g. family rather than genus).

### 30 Conclusions

31 We then show that for two popular taxonomic classifiers, long error-prone reads can  
32 significantly increase classification accuracy, and this is most pronounced for non-microbial  
33 communities. This work provides insight on the expected accuracy for metagenomic  
34 analyses for different taxonomic groups, and establishes the point at which read length  
35 becomes more important than error rate for assigning the correct taxon.

## 36 Introduction

## 37 Applying Metagenomic Methods to Quantify Community Composition

38 To understand ecological community diversity, it is essential to quantify taxon frequency.

39 The most common method of quantifying taxa frequencies is through metabarcoding (Ji et

40 al. 2013). In this method, conserved genomic regions (often 16S rRNA in the case of

41 bacterial and archaeal species; 18S rRNA or Cytochrome c oxidase I for eukaryotic species)

42 are amplified from the sample of interest, sequenced (most often using high-throughput

43 methods such as Illumina), and then classified using one of several available pipelines (e.g.

44 QIIME, MEGAN, Mothur) (Caporaso et al. 2010; Huson et al. 2016; Schloss et al. 2009).

45 Many of these pipelines have been designed around the analysis of bacterial datasets.

46 In contrast to metabarcoding, metagenomic approaches do not rely on the amplification of

47 specific genomic sequences, which can introduce bias. Instead, they aim to quantify

48 community composition based on the recovery and sequencing of all DNA from community

49 samples. Not only do metagenomic methods profile taxon composition in a less biased way

50 than metabarcoding, but they can also yield insight into the functional diversity present in

51 ecosystems (Schloss and Handelsman 2005; Keeling et al. 2014).

52 While metabarcoding approaches have been widely applied to both microbial and eukaryotic

53 taxa, the vast majority of metagenomic studies have focused only on microbial communities.

54 Unsurprisingly, the various advantages and disadvantages of using metagenomic analyses

55 for microbial communities are well-documented (Roumpeka et al. 2017; Thomas, Gilbert,

56 and Meyer 2012; Temperton and Giovannoni 2012). There are likely several factors driving

57 this microbe-centric application of metagenomics, including (1) the greater level of diversity

58 of microbial taxa; (2) the considerable number of microbial taxa that are “unculturable,”

59 making it difficult to collect the requisite amount of DNA for genomic sequencing; (3) the

60 availability of a multitude of non-molecular methods for quantifying multicellular taxa; and (4)

61 the relative paucity of genomic sequence for multicellular organisms in databases (Escobar-

62 Zepeda, Vera-Ponce de León, and Sanchez-Flores 2015) (Supp. Fig.1). This latter factor is

63 perhaps the single largest factor in driving the bias toward microbial metagenomics.

64 However, the amount and diversity of eukaryotic genomic sequence data is rapidly  
65 increasing. Although multicellular metabarcoding databases are currently far more complete  
66 relative to genomic databases, this gap is closing quickly. For example, the Earth  
67 BioGenome project aims to sequence the genomes of upwards of one million eukaryotic  
68 species within the next decade (Lewin et al. 2018). Regardless of the success of this effort,  
69 there are a host of ongoing eukaryotic sequencing projects, including Bat 1K (Teeling et al.  
70 2018), Bird 10K (10,000 bird genomes (OBrien, Haussler, and Ryder 2014)), G10K (10,000  
71 vertebrate genomes (10K Community of Scientists 2009)), and i5K (5000 arthropod  
72 genomes (Robinson et al. 2011)), among others. This suggests that within the next five  
73 years, most multicellular organisms will have at least one member of their family present in  
74 genomic databases, with some groups of multicellular organisms being completely  
75 represented at the genus level.

76 This would increase the utility of metagenomics for assessing membership in plant and  
77 animal communities, especially for cases in which organisms are difficult to observe or  
78 degraded. This is frequently the case for diet studies (Pearman et al. 2018), many  
79 invertebrate communities such as in treeholes (Gossner et al. 2016) or algal holdfasts  
80 (Ojeda and Santelices 1984).

## 81 Analysis of Short-read Metagenomic Data

82 Many metagenomic classification analyses rely on first pass classifiers to assign reads to  
83 one or more taxa, followed by second pass classifiers that can improve on the initial  
84 classification by taking into account the number and relationship of taxa identified in the first  
85 pass. This second step often relies on a lowest common ancestor algorithm (Wood and  
86 Salzberg 2014; Kim et al. 2016; Huson et al. 2016), or by refining taxonomic representation  
87 by examining the results from the first pass classifier (Lu et al. 2016).

88 The most widely used first pass classifier is BLAST, and it is considered gold standard  
89 (McIntyre et al. 2017). However, BLAST is not computationally efficient enough to deal with

90 tens or hundreds of millions of reads. Thus, algorithms for fast metagenomic classification  
91 have been the subject of intense research over the last few years, and include k-mer based  
92 approaches such as CLARK (Ounit et al. 2015), Kraken and related tools (Kraken, Kraken2,  
93 and KrakenUniq) (Wood and Salzberg 2014), Centrifuge (Kim et al. 2016), EnSVMB (Jiang  
94 et al. 2017), and Kaiju (Menzel, Ng, and Krogh 2016), as well as reduced alphabet amino  
95 acid based approaches such as DIAMOND (Buchfink, Xie, and Huson 2015). In almost all  
96 cases these have been designed and benchmarked using short read data (McIntyre et al.  
97 2017).

## 98 Analysis of Long-Read Data Metagenomic data

99 The advent of “third generation” single molecule long read technologies (PacBio and Oxford  
100 Nanopore) has significant implications for metagenomic analyses, most notably for genome  
101 assembly (Frank et al. 2016; Nicholls et al. 2019). These technologies allow read lengths of  
102 10 kilobase pairs (Kbp) and beyond, in strong contrast with the approximately 300 base pairs  
103 (bp) limit of Illumina. However, both PacBio and Nanopore technologies have far higher error  
104 rates (88-94% accuracy for Nanopore (Wick, Judd, and Holt 2018) and 85-87% for PacBio  
105 (Ardui et al. 2018)). The lower accuracy of Nanopore and PacBio (non-circular consensus)  
106 sequence reads may affect the success of current classification methods, and there are few  
107 algorithms designed to exploit long-read data.

108 As a first approach toward determining the use of long-read technologies for metagenomic  
109 applications, we would like to understand the relative advantages and disadvantages of  
110 using short accurate reads versus long error-prone reads. Recent work has shown that  
111 relatively high genus level classifications of approximately 93% have been achieved using  
112 Nanopore-based metagenomic analyses of a mock bacterial community (Brown et al. 2017).  
113 Here we expand this analysis to allow direct comparison between short and long read  
114 approaches. In addition, we compare metagenomic classification success in microbial  
115 communities as compared to communities of multicellular organisms. We find that longer

116 reads, despite their higher error rate, can considerably improve classification accuracy  
117 compared to shorter reads, and that this is especially true for specific taxa.

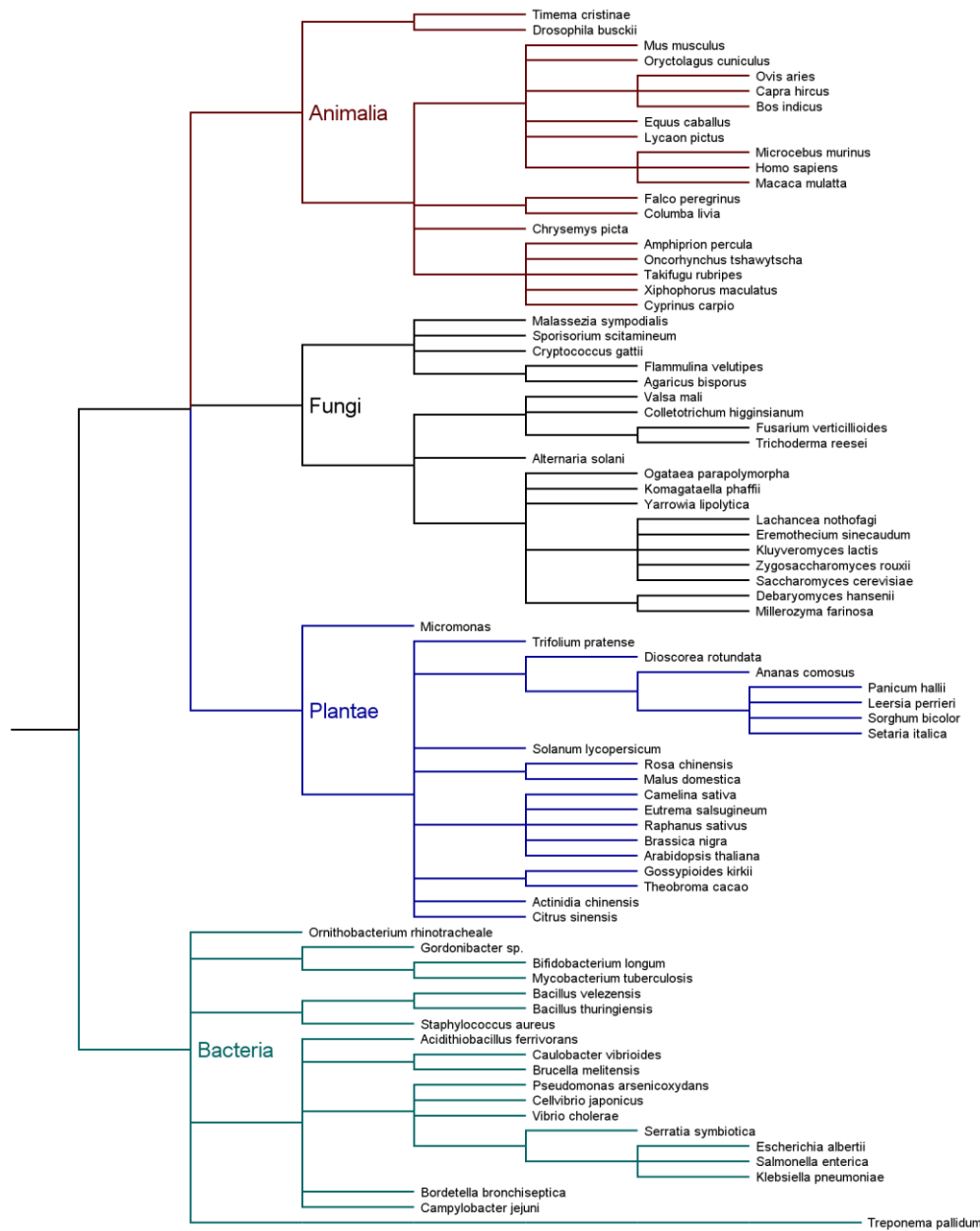
## 118 Methods

### 119 Genomic data

120 For each of four major taxonomic divisions (bacteria, fungi, animals, and plants), we  
121 downloaded 20 genomes from GenBank (Benson et al. 2013). Within each of these  
122 divisions, we included genomes from a total of 22 classes, 46 orders, and 58 families (Figure  
123 1).

### 124 Read simulation

125 We simulated Nanopore reads using NanoSim 2.0.0 (Yang et al. 2017) with the default error  
126 parameters for *E. coli* R9 1D data. This method uses a mixture model to produce simulated  
127 reads with indel and error rates similar to real datasets. The error model is applied equally to  
128 all parts of a read, and the read lengths are drawn from a distribution approximating real  
129 data. To create simulated read data of specific lengths, we truncated the simulated reads  
130 after the relevant number of basepairs using a custom perl script (i.e. to simulate 100bp  
131 Nanopore reads, we truncated all reads in a simulated dataset to 100bp). We did this for  
132 read lengths varying from 100 bp to 4,000 bp at 100 bp intervals, simulating 1,000 reads per  
133 interval for all taxa (a total of 40,000 reads for each taxon, and 3.2 million reads for all taxa  
134 and read lengths).



142 genome, at three read lengths: 100 bp, 150 bp, and 300 bp (a total of 240,000 reads across  
143 all taxa and lengths), and used only single end reads for all analyses.

## 144 Sequence Classification

145 We used BLAST 2.7.1 (Madden 2013) and Kraken2 (Wood and Salzberg 2014) for  
146 sequence classification. We created a local custom database consisting of the NCBI nt  
147 database (downloaded on Feb 8 2019) and the genomes of the 80 taxa that we used to test  
148 classification success. We used the default alignment parameters for BLAST, except for  
149 implementing a maximum e-value of 0.1. We used the match with the highest bit score for all  
150 downstream analyses. For Kraken2 analyses we used the default parameters (in which the  
151 k-mer length is 35 bp and default minimiser length is 31 bp). For Kraken2 we used the taxon  
152 assigned by the lowest common ancestor (LCA) algorithm employed in Kraken2.

## 153 Accuracy metrics

154 To assess the effects of read length on classification accuracy we focus our analysis only on  
155 how often a read is assigned to the correct taxon. For our simulated reads there are three  
156 possible outcomes when querying a database (**Table 1**).

157 We expect that taxa that are well represented in the database, and which have few closely  
158 related taxa, will have high rates of true matches. Taxa with many close relatives in the  
159 database will have many false matches. Taxa that are poorly represented in the database  
160 will have high rates of failed queries. Both of these latter results are in a class usually  
161 referred to as false negatives: we falsely infer taxon A is absent. However, they largely arise  
162 from different mechanisms. Importantly, as genomic databases become more complete, we  
163 expect the fraction of failed queries will decrease. At the same time we expect that the  
164 fraction of false matches may increase, as more and more closely related taxa become  
165 present in the database. The exact nature of this tradeoff is not well explored. Novel  
166 statistical approaches, such as Bayesian re-estimation of species frequencies, may mitigate



167 the problem (Lu et al. 2016); however, improved methods are required address this problem  
168 (Nasko et al. 2018).

169

170 **Table 1. Description of outcomes for database queries.**

Description of outcome	Metric	Notation
A read query from taxon A returns a match from taxon A	True match (we correctly infer taxon A is present)	$M_{\text{true}}$
A read query from taxon A returns a match from a taxon that is not A	False match (we infer taxon A is absent due to a secondary match)	$M_{\text{false}}$
A read query from taxon A returns no hit at all	Failed query (we infer taxon A is absent due to database paucity).	$M_{\text{fail}}$

171

172 There are other aspects of classification success that we do not focus on here. The first of  
173 these is the notion of a true negative: a sequence that is known to *not* arise from any taxa,  
174 should not return a match to any taxa. This is not a biologically realistic situation (all  
175 sequences arise from a taxon), although this aspect is useful when trying to assess the  
176 performance of different classifiers [ref Gardner] and presenting the full truth table. The  
177 second aspect we do not consider here are false positives: if a read query matches taxon A,  
178 but does not arise from taxon A. We would thus falsely interpret taxon A as being present in  
179 a community. This metric is intrinsic to the composition of the community rather than just  
180 each taxon and the database. For example, if taxon A dominates the community, then it  
181 cannot have high rates of false positives relative to true positives simply because the vast  
182 majority of read queries from the community will be from taxon A and thus true positives.  
183 Conversely if taxon B is extremely rare, there will be a large number of false positives  
184 relative to true positives, as very few read queries will be from taxon B, resulting in a very  
185 small number of true positives.

186 Thus, we use a simplified set of metrics (see **Table 1**) that are not intrinsically related to  
187 community composition: true matches, false matches, and failed queries. We used our

188 simulated genomic sequence reads from 80 taxa to quantify these three outcomes at both  
189 the genus and family level. To assign genus and family from species, we used the NCBI  
190 taxonomy database (Federhen 2012) (which is used by BLAST as the default taxon  
191 classifier).

192 We calculate two ratios from the three metrics in **Table 1**. The first is the fraction of true  
193 positives classified correctly (i.e. recall):

$$194 \text{ Recall} = M_{\text{true}} / (M_{\text{true}} + M_{\text{false}} + M_{\text{fail}})$$

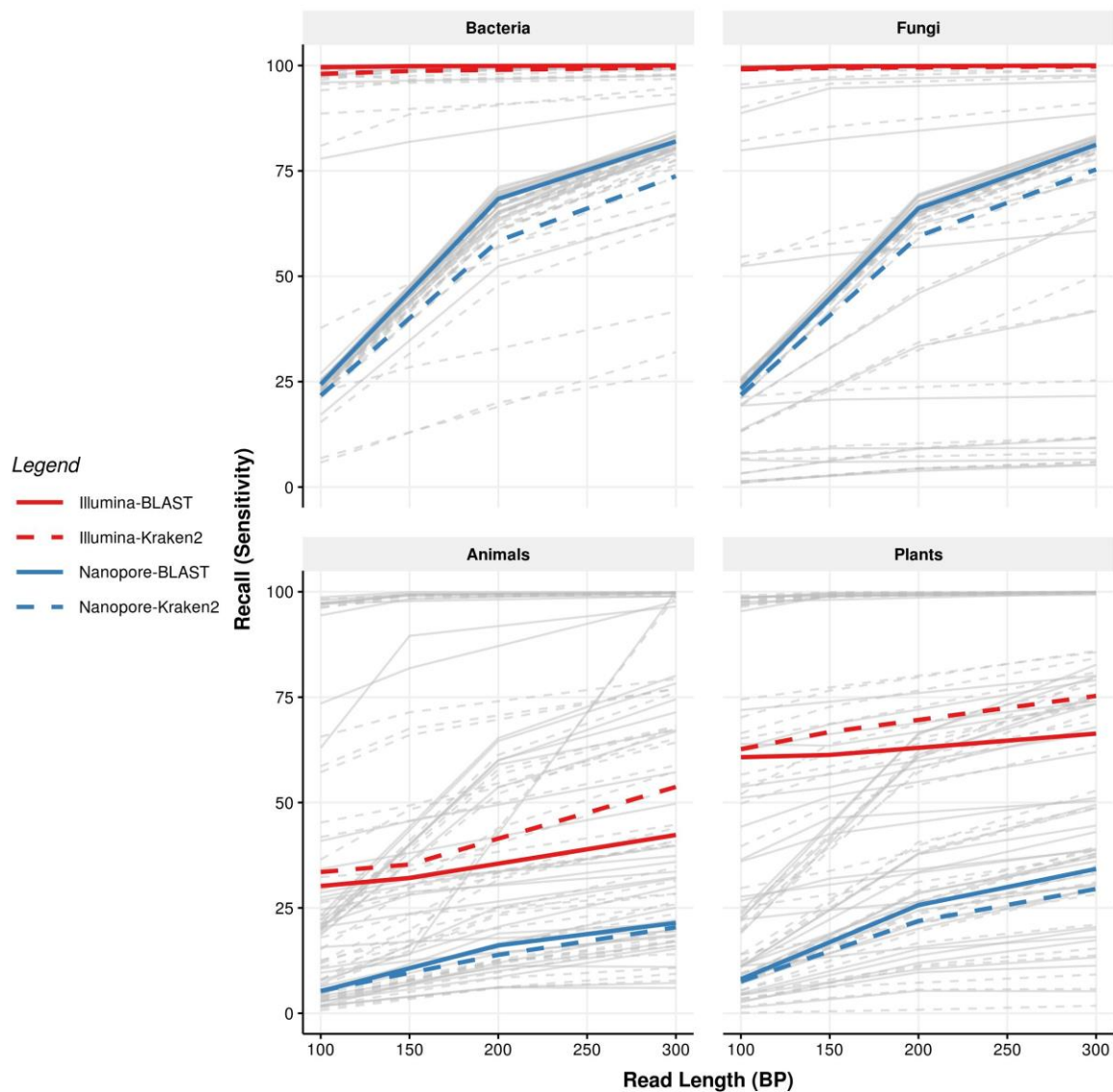
195 The second is the ratio of true matches to false matches. This simply excludes failed queries  
196 from the equation. We term this second metric classification success.

$$197 \text{ Classification Success} = M_{\text{true}} / (M_{\text{true}} + M_{\text{false}})$$

198 The critical difference between these metrics is that taxa which are poorly represented in the  
199 database may nevertheless have high rates of classification success, although recall will  
200 necessarily be low. However, as the fraction of failed queries approaches zero (which we  
201 expect as genomic databases grow), these two metrics become equivalent.

## 202 Results

203 We first looked only at short read lengths to quantify the effects of sequencing technology  
204 and classifier (BLAST or Kraken2) on recall at the level of genus. For both bacteria and  
205 fungi, we found that recall was at or above 99.9% for Illumina reads of any length (100bp,  
206 150bp, or 300bp), for both BLAST and Kraken2 (**Fig. 2**). In strong contrast, for Nanopore  
207 data, recall was far lower; approximately 25% for 100bp reads and increasing to 75% at  
208 300bp. In general, Kraken2 had slightly lower recall than BLAST.



209

210 **Figure 2. Recall is consistently higher in bacteria and fungi than plants or animals for**  
211 **both short Illumina and Nanopore reads.** Each panel shows recall for the different  
212 kingdoms. Recall for individual taxa is indicated in grey, with median recall shown by dashed  
213 (Nanopore) or solid (Illumina). Blue lines indicate recall rates for reads classified using  
214 Kraken2; red for reads classified using BLAST. Illumina reads exhibit consistently higher  
215 recall; bacteria and fungi exhibit higher recall than plants or animals.

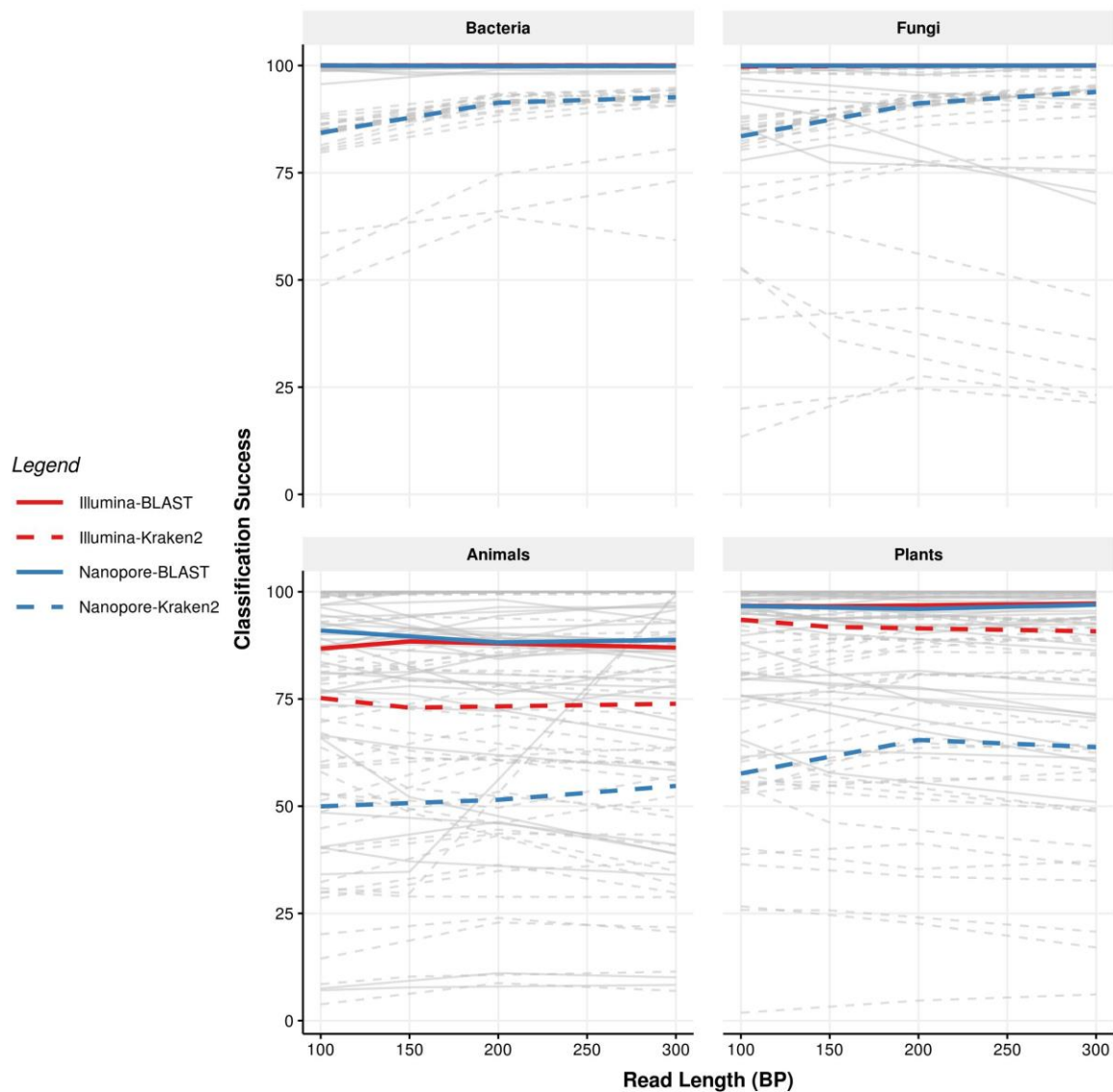
216

217 However, for plants and animals, average recall was low regardless of sequencing  
218 technology. Average recall for Illumina reads peaked at approximately 55% and 75% for  
219 animals and plants, respectively (**Fig 2**, light blue lines). Nanopore recall rates peaked at just  
220 over 20% and 35% for animals and plants, respectively. However, this was highly taxon-  
221 dependent, with some taxa consistently having recall near 100%, while others remained

222 close to 0% regardless of sequencing technology or read length (**Fig. 2**, grey lines). Perhaps  
223 surprisingly, on average Kraken2 outperformed BLAST for Illumina reads for both plant and  
224 animal taxa.

225 We next quantified differences in classification success (the proportion of all classified reads  
226 that were correctly classified), again considering only short read lengths. For bacteria and  
227 fungi, both Illumina and Nanopore reads exhibited high classification success, with the  
228 exception of Kraken2 classification of Nanopore reads (**Fig. 3**). For each sequencing  
229 method and classifier, classification success for plants and animals was low relative to  
230 bacteria and fungi. For both Illumina and Nanopore, BLAST resulted in approximately 87%  
231 and 97% of reads being correctly classified, for animals and plants respectively. However,  
232 Kraken2 success was far lower, especially for Nanopore reads, peaking at 54% in animals  
233 (**Fig. 3**). Over this range of read lengths, we found only a weak relationship between read  
234 length and classification success, in contrast to the results for recall.

235 It is perhaps expected that highly accurate Illumina reads would result in more accurate  
236 taxonomic classification than long error-prone Nanopore reads. However, it is possible to  
237 obtain Nanopore reads far in excess of 300bp (reads up to 2 megabase pairs have been  
238 sequenced), so we next quantified recall and classification success for reads with lengths up  
239 to 4,000 bp. Because such read lengths are not currently possible to obtain using Illumina  
240 technology, we did not measure recall and classification success for Illumina reads of similar  
241 lengths.



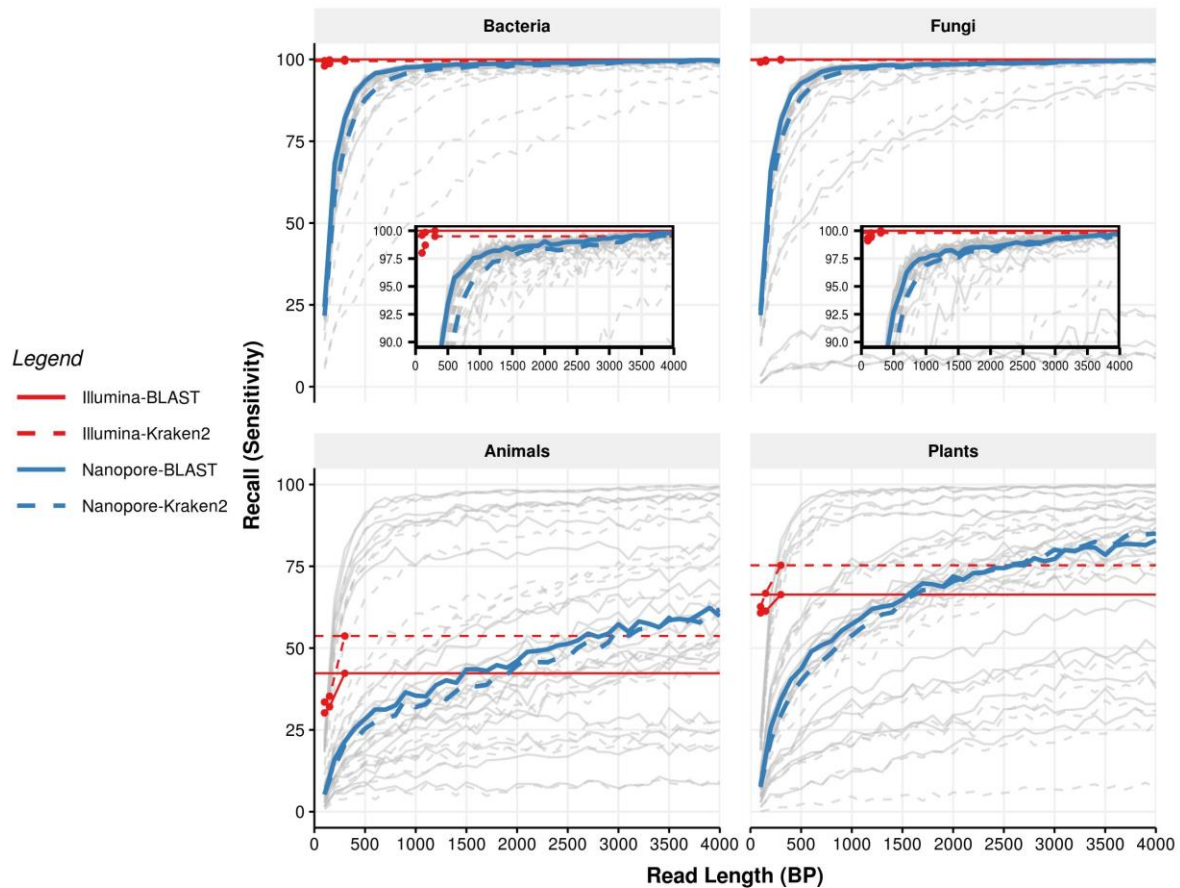
242

243 **Figure 3. Classification success for short reads is weakly related to read length and**  
244 **strongly dependent on classification method.** Each panel shows recall for the different  
245 kingdoms. Classification success for individual taxa is indicated in grey, with median  
246 classification success shown by solid lines (Illumina) or dashed lines (Nanopore). Blue lines  
247 indicate recall rates for reads classified using Kraken2; red for reads classified using BLAST.  
248 For bacteria and fungi, median classification rates of Illumina-BLAST, Illumina-Kraken2, and  
249 Nanopore-BLAST are almost exactly 100% for all read lengths.

250

251 We observed similar relationship between read length and recall for both BLAST and  
252 Kraken2. For bacteria and fungi, recall increased from ~20% using 100 bp reads to almost  
253 100% when using 1500 bp reads. For animals and plants we observed similar trends,  
254 although at no point did recall approach 100%. However, long Nanopore reads surpassed

255 the recall of even the longest Illumina reads (300 bp) classified with Kraken, with crossover  
256 points at approximately 3000 bp for animals and 2500 bp for plants (**Fig. 4**, red and blue  
257 solid lines).



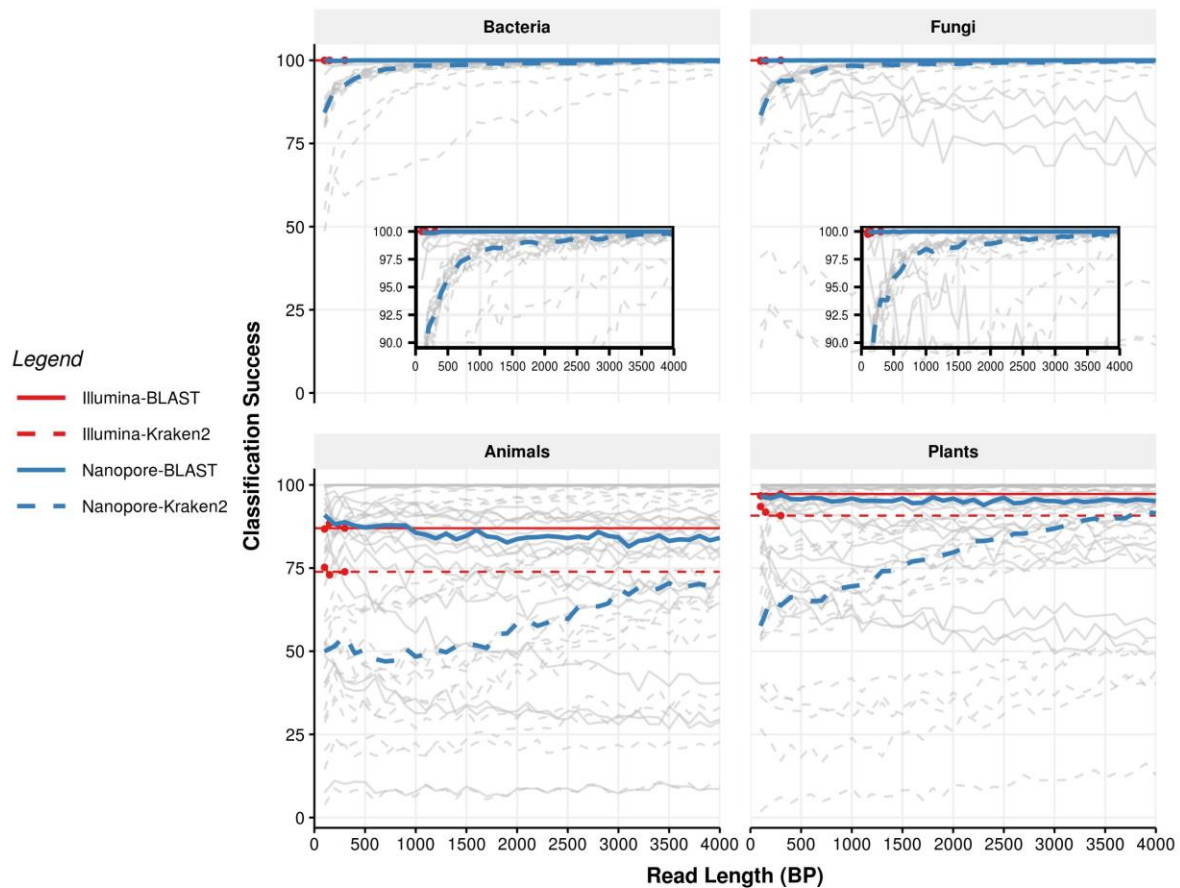
258

259 **Figure 4. Long Nanopore reads equal or surpass the recall of the longest Illumina**  
260 **reads for both BLAST and Kraken2.** Each panel shows recall for the different kingdoms.  
261 Recall for Nanopore reads for individual taxa is indicated in grey, with median recall  
262 indicated by dashed lines, either blue (Kraken2) or red (BLAST). The recall rates for 300 bp  
263 Illumina reads are shown as thin solid lines, again either blue (Kraken2) or red (BLAST).  
264 Coloured points show the recall for all Illumina reads of all lengths (100 bp, 150 bp, and 300  
265 bp).

266

267 We also considered this metric at the level of family. In this case found that for animals,  
268 Nanopore reads surpassed Illumina reads only at lengths close to 4000 bp, reaching  
269 approximately 70% recall at this point (**Supp Fig. 2**). However, for plants Nanopore recall  
270 surpassed Illumina recall at 2500 bp, with 4000 bp reads yielding a recall of approximately

271 90%. We found again that for both animals and plants, Kraken2 recall surpassed BLAST  
272 when relying on Illumina reads.



273

274 **Figure 5. Classification success for long reads is dependent on read length only for**  
275 **Kraken2 classification.** Each panel shows classification success for the different kingdoms.  
276 Classification success for individual taxa is indicated in grey, with median classification  
277 success shown by solid lines (Illumina) or dashed lines (Nanopore). Blue indicates  
278 classification success rates for reads classified using Kraken2, while red indicates those  
279 classified using BLAST. For animal and plants, the classification success of Kraken2  
280 depends strongly on read length, and never surpasses BLAST or Illumina at any length.

281

282 We next examined classification success at longer read lengths. For BLAST we observed no  
283 relationship between classification success and read length for any taxon (**Fig 5**). Bacteria  
284 and fungi both had consistently high classification success (median 100%), while animals  
285 and plants had lower classification success (median 82% and 96%, respectively). However,  
286 for Kraken2 we observed a consistent increase in classification success as read length

287 increased. However, this never exceeded the classification success we observed for BLAST,  
288 nor did it succeed the classification success we observed for short accurate Illumina reads.  
289 Finally, we tested classification success at the level of Family. In this case, we observed that  
290 for BLAST, the classification success for plants was approximately 99% overall read lengths,  
291 while for Kraken2 only 4000 bp reads reached this level. For animals, BLAST classification  
292 success was approximately 95% over all read lengths, but for Kraken2 reached a maximum  
293 of 85% at the longest read lengths (**Supp. Fig. 3**).

## 294 Discussion

295 Here we have compared the relative accuracy of taxon classification using simulated short  
296 accurate reads (Illumina) and long, error-prone reads (Nanopore) with known ground truth.  
297 We have used two simple metrics of success: recall (the ratio of correctly classified reads to  
298 all reads) and classification success (the ratio of correctly classified reads to all classified  
299 reads). We have tested taxon classification using a broad range of taxa, including bacteria,  
300 fungi, animals, and plants.

301 Recall for both BLAST and Kraken2 was improved by the use of long reads, especially in the  
302 case of animals and plants, for which recall improved almost three-fold as read length  
303 increased from 300 bp to 4,000 bp. Generally both Kraken2 and BLAST achieved similar  
304 levels of recall. The exception was for short reads for animals and plants, for which Kraken2  
305 was more accurate than BLAST.

306 We found no relationship between classification success and read length for BLAST. This  
307 implies that the ratio of correctly classified reads to all classified reads remains relatively  
308 constant over different read lengths. However, the number of reads that are classified *at all*  
309 increases with read length (causing an increase in recall). These observations are in line  
310 with what has been observed by others (McHardy et al. 2007). The exception to this lack of  
311 relationship between classification success and read length was for Kraken2, for which the



312 proportion of correctly classified reads increases with read length by more than 50% for both  
313 plants and animals.

314 Our results also indicate that recall for long Nanopore reads was equal to or higher than  
315 short Illumina reads. This was true regardless of kingdom, or classification method, with  
316 Nanopore surpassing 300 bp Illumina reads at approximately 1500 bp for plants and  
317 animals, and surpassing 150 bp Illumina reads at between 1500 bp and 3000 bp for bacteria  
318 and fungi, depending on the methodology (**Fig. 4**). Even the longest Illumina reads, at 300  
319 bp, were outclassed by Nanopore at between 3500 and 4000 bp, depending on  
320 methodology. These results do suggest that one approach to improve Nanopore  
321 classification accuracy is to impose minimum read lengths. This can be achieved by  
322 performing size selection during library preparation or during computational analyses.

323 At first glance, then, there appears to be a clear trade-off between short read Illumina and  
324 long read Nanopore sequencing for metagenomic analyses. While Nanopore allows higher  
325 recall at long read lengths, this advantage is offset by the fact that Illumina generally  
326 provides more reads per run. At most, recall for Nanopore improves 50% beyond 300 bp  
327 Illumina reads, while classification success is similar (using BLAST). Thus, if the read  
328 capacity of Illumina runs is 50% or more than Nanopore, the number of classified reads will  
329 be maximised using Illumina technology - on a per sequencing run basis. However, for many  
330 researchers the more relevant metric is cost per read. In this case, MinION read yields are  
331 approximately equal to MiSeq, and only HiSeq or NovaSeq provides a clear cost advantage  
332 over Nanopore MinION. On the other hand, cost per read for PromethION are not far from  
333 NovaSeq. Thus, we find no clear advantage in using Illumina over Nanopore given the  
334 observed classification accuracy for long inaccurate Nanopore reads.

### 335 Differences in accuracy between bacteria, fungi, animals, and plants

336 We find very large differences in classification accuracy (mostly in terms of recall) for  
337 bacteria and fungi versus plants and animals. The discrepancy between taxonomic groups

338 likely arises from a variety of factors. Among these are the higher degree of divergence  
339 between bacterial species relative to animal and plant species, and the complexity of  
340 bacterial genomes compared to eukaryotic genomes. We discuss these factors below.

341 Bacterial taxa are often considered separate species once they have diverged by 6% ANI  
342 (Average Nucleotide Identity) on a genomic level (Stackebrandt and Goebel 1994;  
343 Konstantinidis and Tiedje 2005). The degree of nucleotide divergence between eukaryotic  
344 species is not standardised (Cognato 2006), and species are generally designated as such  
345 based on the biological species concept put forward by Mayr (Mayr 1999). Although  
346 divergence levels differ substantially between loci (as for bacteria), for some loci general  
347 ranges for eukaryotic species have emerged. For example, for mitochondrial COI, between-  
348 species divergence is usually greater than 3% (Song et al. 2008; Lefébure et al. 2006).

349 These loci are among the fastest diverging loci in plant and animal genomes, and many  
350 other loci may differ by far less than 1% between species. Due to this low level of  
351 divergence, metagenomic classifiers may frequently classify animal and plant genera with  
352 lower accuracy than bacterial genera.

353 A second explanation for the increased classification success in bacteria and fungi is that  
354 these genomes contain fewer repetitive elements than animals or plants (Treangen et al.  
355 2009). Although such repetitive regions are usually masked from classifiers (including  
356 BLAST and Kraken2), this masking may not be complete.

357 A third reason is that the genomic databases for plants and animals are far less complete  
358 than for bacteria and fungi. There is a large difference in the number of genomes and  
359 sequences available for different Kingdoms, with bacteria having significantly more species  
360 present than the next closest kingdom (See Supp. Fig.1). However, we expect this factor will  
361 be mitigated in the future as genomic databases continue to expand and computational  
362 search methods continue to improve.

## 363 Differences in accuracy between Kraken2 and BLAST

364 We observed similar levels of recall for BLAST and Kraken2 over most reads lengths.  
365 However, there were strong differences in classification success. For short reads, Kraken2  
366 classification success was far lower than BLAST. As read lengths increased, Kraken2  
367 classification success approached BLAST. Part of this is likely due to longer reads allowing  
368 multiple k-mer matches, decreasing the probability of a false positive classification. One  
369 perhaps underappreciated advantage of Kraken2 over BLAST is that Kraken2 has reduced  
370 sensitivity to structural variation within reads. As Kraken2 allows multiple k-mers to match  
371 within a read, structural changes (e.g. inversions) are less likely to influence the outcome of  
372 Kraken2 matching. Such structural changes may influence BLAST due to the matching and  
373 extend algorithm. Thus for long reads, classifiers that are insensitive to synteny may be  
374 more successful, especially for taxa in which structural rearrangements are common.

## 375 Conclusions

376 Here we have shown despite being error-prone, Nanopore reads are useful for metagenomic  
377 classification due to their increased length, and that for plant and animal communities, the  
378 classification accuracy of long Nanopore reads exceeds that of Illumina. We found that  
379 classification accuracy is more dependent on the set of taxa being considered than on the  
380 metagenomic classifier being used (Kraken2 or BLAST), and that this was true for both short  
381 accurate (Illumina) and long error-prone (Nanopore) sequence data. Together these data  
382 suggest that one consideration in selecting a metagenomic sequencing method (i.e. long or  
383 short read) is the taxonomic group of interest.

## 384 Declarations

385 **Ethics approval and consent to participate** – not applicable

386 **Consent for publication** – approved by all authors

387 **Availability of data and material** – data for this manuscript will be uploaded to DataDryad if  
388 manuscript is accepted

389 **Competing interests** – the authors have no competing interests to declare

390 **Funding** – Some of this work was supported by a Massey University Strategic Research  
391 Excellence Fund awarded to NF.

392 **Authors' contributions** - WP, NF, and OS conceived the project, WP and OS simulated  
393 and generated the data. WP and OS analysed the data. WP, NF, and OS wrote the paper.

394 **Acknowledgements** - Thanks to Paul Gardner for his helpful and insightful comments on  
395 the manuscript.

## 396 References

397 10K Community of Scientists, Genome. 2009. "Genome 10K: A Proposal to Obtain Whole-  
398 Genome Sequence for 10 000 Vertebrate Species." *The Journal of Heredity*.

399 <https://academic.oup.com/jhered/article-abstract/100/6/659/839176>.

400 Ardui, Simon, Adam Ameer, Joris R. Vermeesch, and Matthew S. Hestand. 2018. "Single  
401 Molecule Real-Time (SMRT) Sequencing Comes of Age: Applications and Utilities for  
402 Medical Diagnostics." *Nucleic Acids Research* 46 (5): 2159–68.

403 Benson, Dennis A., Mark Cavanaugh, Karen Clark, Ilene Karsch-Mizrachi, David J. Lipman,  
404 James Ostell, and Eric W. Sayers. 2013. "GenBank." *Nucleic Acids Research* 41  
405 (Database issue): D36–42.

406 Brown, Bonnie L., Mick Watson, Samuel S. Minot, Maria C. Rivera, and Rima B. Franklin.  
407 2017. "MinION nanopore Sequencing of Environmental Metagenomes: A Synthetic  
408 Approach." *GigaScience* 6 (3): 1–10.

409 Buchfink, Benjamin, Chao Xie, and Daniel H. Huson. 2015. "Fast and Sensitive Protein  
410 Alignment Using DIAMOND." *Nature Methods* 12 (1): 59–60.

411 Caporaso, J. Gregory, Justin Kuczynski, Jesse Stombaugh, Kyle Bittinger, Frederic D.

412 Bushman, Elizabeth K. Costello, Noah Fierer, et al. 2010. "QIIME Allows Analysis of

- 413 High-Throughput Community Sequencing Data.” *Nature Methods* 7 (5): 335–36.
- 414 Cognato, Anthony I. 2006. “Standard Percent DNA Sequence Difference for Insects Does  
415 Not Predict Species Boundaries.” *Journal of Economic Entomology* 99 (4): 1037–45.
- 416 Escobar-Zepeda, Alejandra, Arturo Vera-Ponce de León, and Alejandro Sanchez-Flores.  
417 2015. “The Road to Metagenomics: From Microbiology to DNA Sequencing  
418 Technologies and Bioinformatics.” *Frontiers in Genetics* 6 (December): 348.
- 419 Federhen, Scott. 2012. “The NCBI Taxonomy Database.” *Nucleic Acids Research* 40  
420 (Database issue): D136–43.
- 421 Frank, J. A., Y. Pan, A. Tooming-Klunderud, V. G. H. Eijsink, A. C. McHardy, A. J.  
422 Nederbragt, and P. B. Pope. 2016. “Improved Metagenome Assemblies and Taxonomic  
423 Binning Using Long-Read Circular Consensus Sequence Data.” *Scientific Reports* 6  
424 (May): 25373.
- 425 Homer, Nils. 2017. *DWGSIM* (version 0.1.12). Github. <https://github.com/nh13/DWGSIM>.
- 426 Huson, Daniel H., Sina Beier, Isabell Flade, Anna Górska, Mohamed El-Hadidi, Suparna  
427 Mitra, Hans-Joachim Ruscheweyh, and Rewati Tappu. 2016. “MEGAN Community  
428 Edition - Interactive Exploration and Analysis of Large-Scale Microbiome Sequencing  
429 Data.” *PLoS Computational Biology* 12 (6): e1004957.
- 430 Jiang, Yuan, Jun Wang, Dawen Xia, and Guoxian Yu. 2017. “EnSVMB: Metagenomics  
431 Fragments Classification Using Ensemble SVM and BLAST.” *Scientific Reports* 7 (1):  
432 9440.
- 433 Ji, Yinqiu, Louise Ashton, Scott M. Pedley, David P. Edwards, Yong Tang, Akihiro  
434 Nakamura, Roger Kitching, et al. 2013. “Reliable, Verifiable and Efficient Monitoring of  
435 Biodiversity via Metabarcoding.” *Ecology Letters* 16 (10): 1245–57.
- 436 Keeling, Patrick J., Fabien Burki, Heather M. Wilcox, Bassem Allam, Eric E. Allen, Linda A.  
437 Amaral-Zettler, E. Virginia Armbrust, et al. 2014. “The Marine Microbial Eukaryote  
438 Transcriptome Sequencing Project (MMETSP): Illuminating the Functional Diversity of  
439 Eukaryotic Life in the Oceans through Transcriptome Sequencing.” *PLoS Biology* 12 (6):  
440 e1001889.

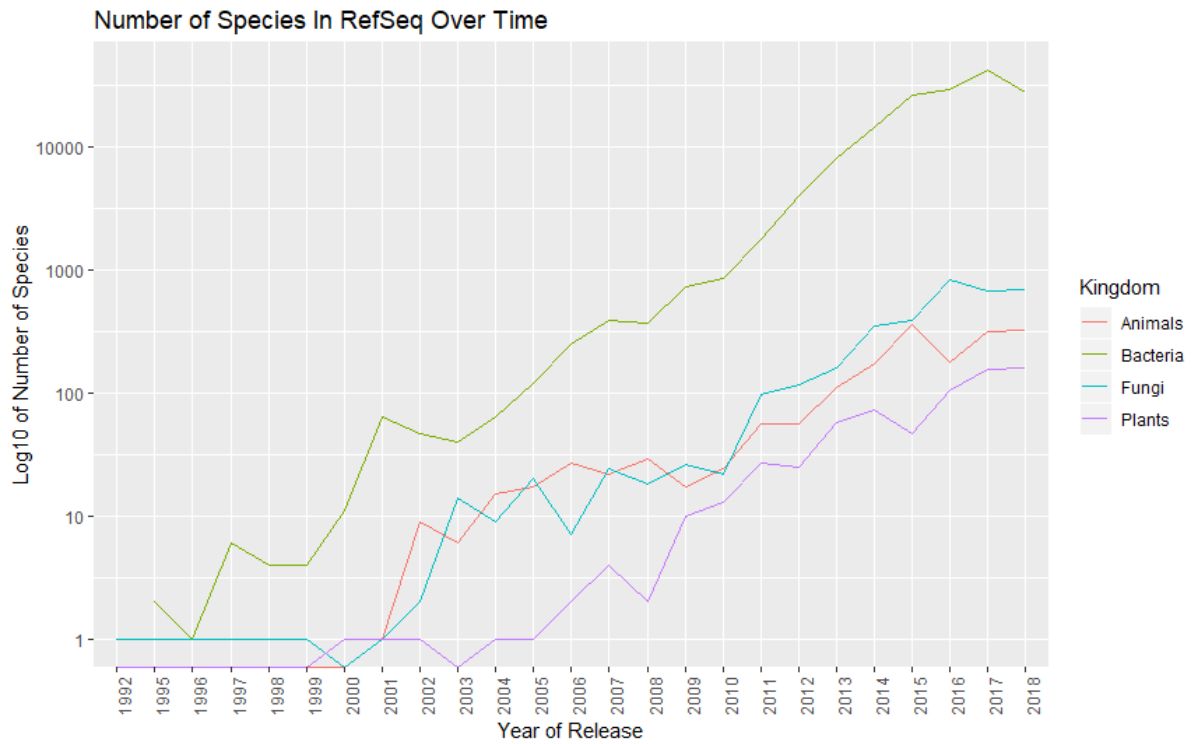
- 441 Kim, Daehwan, Li Song, Florian P. Breitwieser, and Steven L. Salzberg. 2016. "Centrifuge:  
442 Rapid and Sensitive Classification of Metagenomic Sequences." *Genome Research* 26  
443 (12): 1721–29.
- 444 Konstantinidis, Konstantinos T., and James M. Tiedje. 2005. "Genomic Insights That  
445 Advance the Species Definition for Prokaryotes." *Proceedings of the National Academy  
446 of Sciences of the United States of America* 102 (7): 2567–72.
- 447 Lefébure, T., C. J. Douady, M. Gouy, and J. Gibert. 2006. "Relationship between  
448 Morphological Taxonomy and Molecular Divergence within Crustacea: Proposal of a  
449 Molecular Threshold to Help Species Delimitation." *Molecular Phylogenetics and  
450 Evolution* 40 (2): 435–47.
- 451 Lewin, Harris A., Gene E. Robinson, W. John Kress, William J. Baker, Jonathan Coddington,  
452 Keith A. Crandall, Richard Durbin, et al. 2018. "Earth BioGenome Project: Sequencing  
453 Life for the Future of Life." *Proceedings of the National Academy of Sciences of the  
454 United States of America* 115 (17): 4325–33.
- 455 Lu, Jennifer, Florian P. Breitwieser, Peter Thielen, and Steven L. Salzberg. 2016. "Bracken:  
456 Estimating Species Abundance in Metagenomics Data." <https://doi.org/10.1101/051813>.
- 457 Madden, Thomas. 2013. *The BLAST Sequence Analysis Tool*. National Center for  
458 Biotechnology Information (US).
- 459 Mayr, Ernst. 1999. *Systematics and the Origin of Species, from the Viewpoint of a Zoologist*.  
460 Harvard University Press.
- 461 McHardy, Alice Carolyn, Héctor García Martín, Aristotelis Tsirigos, Philip Hugenholtz, and  
462 Isidore Rigoutsos. 2007. "Accurate Phylogenetic Classification of Variable-Length DNA  
463 Fragments." *Nature Methods* 4 (1): 63–72.
- 464 McIntyre, Alexa B. R., Rachid Ounit, Ebrahim Afshinnikoo, Robert J. Prill, Elizabeth Hénaff,  
465 Noah Alexander, Samuel S. Minot, et al. 2017. "Comprehensive Benchmarking and  
466 Ensemble Approaches for Metagenomic Classifiers." *Genome Biology* 18 (1): 182.
- 467 Menzel, Peter, Kim Lee Ng, and Anders Krogh. 2016. "Fast and Sensitive Taxonomic  
468 Classification for Metagenomics with Kaiju." *Nature Communications* 7 (April): 11257.

- 469 Nasko, Daniel J., Sergey Koren, Adam M. Phillippy, and Todd J. Treangen. 2018. "RefSeq  
470 Database Growth Influences the Accuracy of K-Mer-Based Lowest Common Ancestor  
471 Species Identification." *Genome Biology* 19 (1): 165.
- 472 Nicholls, Samuel M., Joshua C. Quick, Shuiquan Tang, and Nicholas J. Loman. 2019. "Ultra-  
473 Deep, Long-Read Nanopore Sequencing of Mock Microbial Community Standards."  
474 *GigaScience* 8 (5). <https://doi.org/10.1093/gigascience/giz043>.
- 475 OBrien, Stephen J., David Haussler, and Oliver Ryder. 2014. "The Birds of Genome10K."  
476 *GigaScience* 3 (1): 32.
- 477 O'Leary, Nuala A., Mathew W. Wright, J. Rodney Brister, Stacy Ciuffo, Diana Haddad, Rich  
478 McVeigh, Bhanu Rajput, et al. 2016. "Reference Sequence (RefSeq) Database at NCBI:  
479 Current Status, Taxonomic Expansion, and Functional Annotation." *Nucleic Acids*  
480 *Research* 44 (D1): D733–45.
- 481 Ounit, Rachid, Steve Wanamaker, Timothy J. Close, and Stefano Lonardi. 2015. "CLARK:  
482 Fast and Accurate Classification of Metagenomic and Genomic Sequences Using  
483 Discriminative K-Mers." *BMC Genomics* 16 (March): 236.
- 484 Robinson, Gene E., Kevin J. Hackett, Mary Purcell-Miramontes, Susan J. Brown, Jay D.  
485 Evans, Marian R. Goldsmith, Daniel Lawson, Jack Okamuro, Hugh M. Robertson, and  
486 David J. Schneider. 2011. "Creating a Buzz about Insect Genomes." *Science* 331  
487 (6023): 1386.
- 488 Roumpeka, Despoina D., R. John Wallace, Frank Escalettes, Ian Fotheringham, and Mick  
489 Watson. 2017. "A Review of Bioinformatics Tools for Bio-Prospecting from Metagenomic  
490 Sequence Data." *Frontiers in Genetics* 8 (March): 23.
- 491 Schloss, Patrick D., and Jo Handelsman. 2005. "Metagenomics for Studying Unculturable  
492 Microorganisms: Cutting the Gordian Knot." *Genome Biology* 6 (8): 229.
- 493 Schloss, Patrick D., Sarah L. Westcott, Thomas Ryabin, Justine R. Hall, Martin Hartmann,  
494 Emily B. Hollister, Ryan A. Lesniewski, et al. 2009. "Introducing Mothur: Open-Source,  
495 Platform-Independent, Community-Supported Software for Describing and Comparing  
496 Microbial Communities." *Applied and Environmental Microbiology* 75 (23): 7537–41.

- 497 Song, Hojun, Jennifer E. Buhay, Michael F. Whiting, and Keith A. Crandall. 2008. "Many  
498 Species in One: DNA Barcoding Overestimates the Number of Species When Nuclear  
499 Mitochondrial Pseudogenes Are Coamplified." *Proceedings of the National Academy of  
500 Sciences of the United States of America* 105 (36): 13486–91.
- 501 Stackebrandt, E., and B. M. Goebel. 1994. "Taxonomic Note: A Place for DNA-DNA  
502 Reassociation and 16S rRNA Sequence Analysis in the Present Species Definition in  
503 Bacteriology." *International Journal of Systematic and Evolutionary Microbiology* 44 (4):  
504 846–49.
- 505 Teeling, Emma C., Sonja C. Vernes, Liliana M. Dávalos, David A. Ray, M. Thomas P.  
506 Gilbert, Eugene Myers, and Bat1K Consortium. 2018. "Bat Biology, Genomes, and the  
507 Bat1K Project: To Generate Chromosome-Level Genomes for All Living Bat Species."  
508 *Annual Review of Animal Biosciences* 6 (February): 23–46.
- 509 Temperton, Ben, and Stephen J. Giovannoni. 2012. "Metagenomics: Microbial Diversity  
510 through a Scratched Lens." *Current Opinion in Microbiology* 15 (5): 605–12.
- 511 Thomas, Torsten, Jack Gilbert, and Folker Meyer. 2012. "Metagenomics - a Guide from  
512 Sampling to Data Analysis." *Microbial Informatics and Experimentation* 2 (1): 3.
- 513 Treangen, Todd J., Anne-Laure Abraham, Marie Touchon, and Eduardo P. C. Rocha. 2009.  
514 "Genesis, Effects and Fates of Repeats in Prokaryotic Genomes." *FEMS Microbiology  
515 Reviews* 33 (3): 539–71.
- 516 Wick, Ryan, Louise M. Judd, and Kathryn E. Holt. 2018. *Comparison of Oxford Nanopore  
517 Basecalling Tools*. <https://doi.org/10.5281/zenodo.1188469>.
- 518 Wood, Derrick E., and Steven L. Salzberg. 2014. "Kraken: Ultrafast Metagenomic Sequence  
519 Classification Using Exact Alignments." *Genome Biology* 15 (3): R46.
- 520 Yang, Chen, Justin Chu, René L. Warren, and Inanç Birol. 2017. "NanoSim: Nanopore  
521 Sequence Read Simulator Based on Statistical Characterization." *GigaScience* 6 (4): 1–  
522 6.



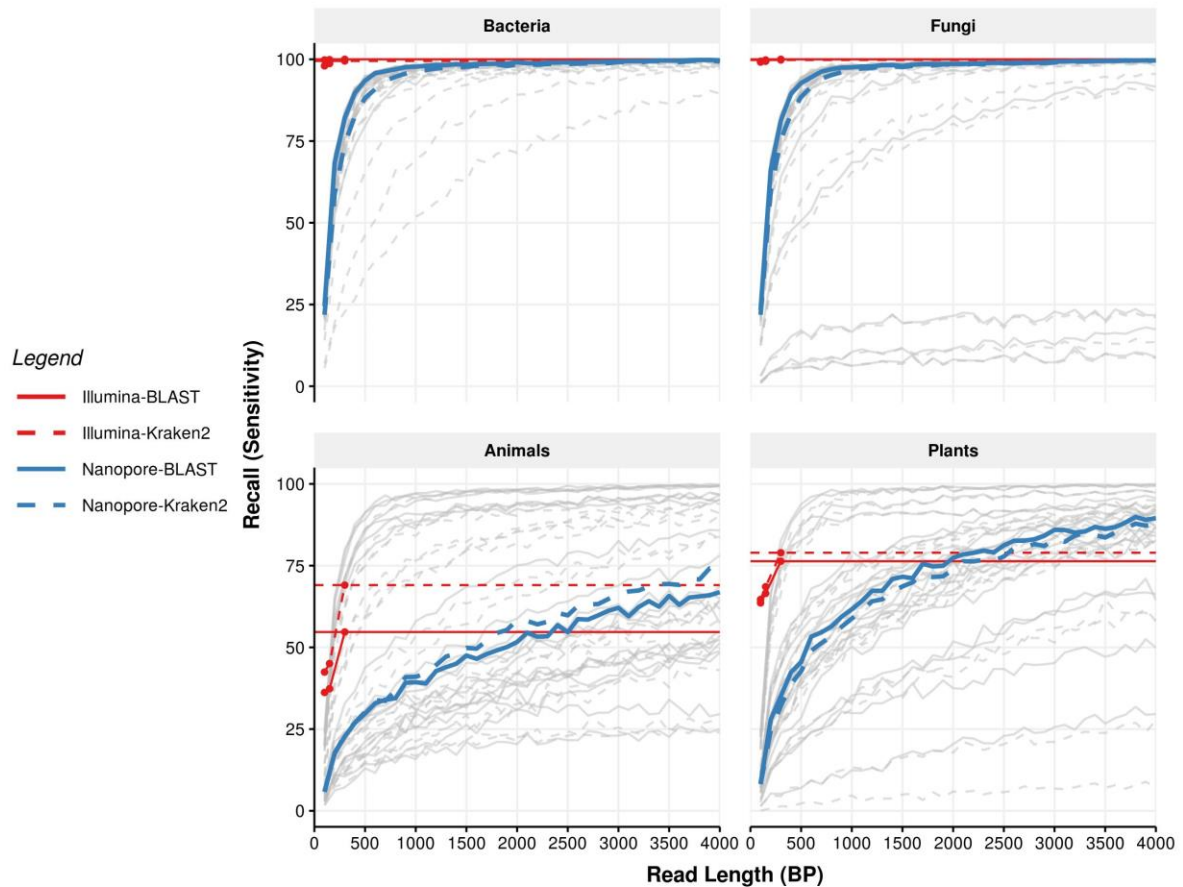
## 523 Supplementary Materials



524

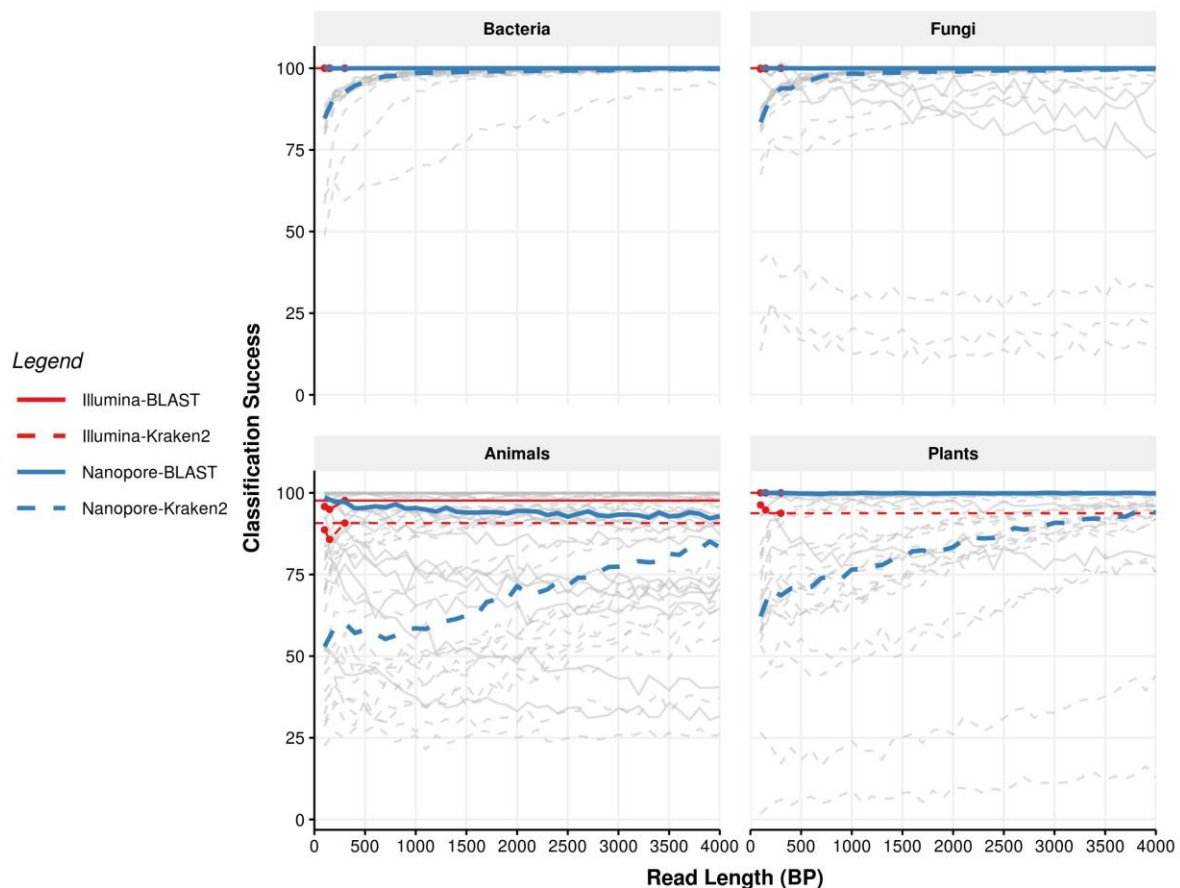
525 **Supplementary Figure 1. The number of species present in the NCBI RefSeq database has**  
526 **grown roughly exponentially over time.** Note that the y-axis is plotted on a log scale. Data were  
527 retrieved from the RefSeq database (O'Leary et al. 2016):  
528 [https://ftp.ncbi.nlm.nih.gov/genomes/GENOME\\_REPORTS/eukaryotes.txt](https://ftp.ncbi.nlm.nih.gov/genomes/GENOME_REPORTS/eukaryotes.txt) and  
529 [https://ftp.ncbi.nlm.nih.gov/genomes/GENOME\\_REPORTS/prokaryotes.txt](https://ftp.ncbi.nlm.nih.gov/genomes/GENOME_REPORTS/prokaryotes.txt)

530



531

532 **Supplementary Figure 2 Recall at the family level.** Each panel shows recall for the  
533 different kingdoms. Recall for Nanopore reads for individual taxa is indicated in grey,  
534 with median recall indicated by dashed lines, either blue (Kraken2) or red (BLAST). The recall  
535 rates for 300 bp Illumina reads are shown as thin solid lines, again either blue (Kraken2) or  
536 red (BLAST). Coloured points show the recall for all Illumina reads of all lengths (100 bp,  
537 150 bp, and 300 bp).  
538



539

540 **Supplementary Figure 3 Classification success at the family level.** Each panel shows  
 541 classification success for the different kingdoms. Classification success for individual taxa is  
 542 indicated in grey, with median classification success shown by solid lines (Illumina) or  
 543 dashed lines (Nanopore). Blue indicates classification success rates for reads classified  
 544 using Kraken2, while red indicates those classified using BLAST. For animal and plants, the  
 545 classification success of Kraken2 depends strongly on read length, and never surpasses  
 546 BLAST or Illumina at any length.  
 547

548 **Supplementary Table 1. List of species include in the *in silico* mock community, with  
 549 associated Kingdom and NCBI**

Species	Kingdom	NCBI Accession
<i>Actinidia chinensis</i>	Plantae	CM009654.1
<i>Ananas comosus</i>	Plantae	CM003813.1
<i>Arabidopsis thaliana</i>	Plantae	CP002684.1
<i>Brassica nigra</i>	Plantae	CM004491.1
<i>Camelina sativa</i>	Plantae	CM002729.1

<i>Citrus sinensis</i>	Plantae	CM001701.1
<i>Dioscorea rotundata</i>	Plantae	BDMI01000001.1
<i>Eutrema salsugineum</i>	Plantae	CM001778.1
<i>Gossypioides kirkii</i>	Plantae	CM008980.1
<i>Leersia perrieri</i>	Plantae	CM002476.1
<i>Malus domestica</i>	Plantae	CM007867.1
<i>Micromonas sp.</i>	Plantae	CP001574.1
<i>Panicum hallii</i>	Plantae	CM008046.2
<i>Raphanus sativus</i>	Plantae	CM007999.1
<i>Rosa chinensis</i>	Plantae	CM009582.1
<i>Setaria italica</i>	Plantae	CM004364.1
<i>Solanum lycopersicum</i>	Plantae	CM001064.3
<i>Sorghum bicolor</i>	Plantae	CM000760.3
<i>Theobroma cacao</i>	Plantae	LT594788.1
<i>Trifolium pratense</i>	Plantae	LT555306.1
<i>Amphiprion percula</i>	Animalia	CM009708.1
<i>Bos indicus</i>	Animalia	CM003021.1
<i>Capra hircus</i>	Animalia	CM001710.2
<i>Chrysemys picta</i>	Animalia	CM002655.1
<i>Columba livia</i>	Animalia	CM007525.1
<i>Cyprinus carpio</i>	Animalia	LN590701.1
<i>Drosophila busckii</i>	Animalia	CP012523.1
<i>Equus caballus</i>	Animalia	CM000377.2
<i>Falco peregrinus</i>	Animalia	CM007505.1
<i>Homo sapiens</i>	Animalia	CM004593.1
<i>Lycaon pictus</i>	Animalia	CM007565.1
<i>Macaca mulatta</i>	Animalia	CM000308.1
<i>Microcebus murinus</i>	Animalia	CM007661.1

<i>Mus musculus</i>	Animalia	CM004154.1
<i>Oncorhynchus tshawytscha</i>	Animalia	CM009202.1
<i>Oryctolagus cuniculus</i>	Animalia	CM000790.1
<i>Ovis aries</i>	Animalia	CM008472.1
<i>Takifugu rubripes</i>	Animalia	HE602535.1
<i>Timema cristinae</i>	Animalia	CM007794.2
<i>Xiphophorus maculatus</i>	Animalia	CM008938.1
<i>Agaricus bisporus</i>	Fungi	CP015470.1
<i>Alternaria solani</i>	Fungi	CP022024.1
<i>Colletotrichum higginsianum</i>	Fungi	CM004455.1
<i>Cryptococcus gattii</i>	Fungi	CP025759.1
<i>Debaryomyces hansenii</i>	Fungi	CR382133.2
<i>Eremothecium sinicaudum</i>	Fungi	CP014242.1
<i>Flammulina velutipes</i>	Fungi	CM002695.1
<i>Fusarium verticillioides</i>	Fungi	CM000578.1
<i>Kluyveromyces lactis</i>	Fungi	CR382121.1
<i>Komagataella phaffii</i>	Fungi	LT962476.1
<i>Lachancea nothofagi</i>	Fungi	LT598449.1
<i>Malassezia sympodialis</i>	Fungi	LT671813.1
<i>Millerozyma farinosa</i>	Fungi	FO082059.1
<i>Ogataea parapolyomorpha</i>	Fungi	CM002300.1
<i>Saccharomyces cerevisiae</i>	Fungi	BK006935.2
<i>Sporisorium scitamineum</i>	Fungi	CP010913.1
<i>Trichoderma reesei</i>	Fungi	CP016232.1
<i>Valsa mali</i>	Fungi	CM003098.1
<i>Yarrowia lipolytica</i>	Fungi	HG934059.1
<i>Zygosaccharomyces rouxii</i>	Fungi	CU928173.1
<i>Acidithiobacillus ferrivorans</i>	Bacteria	LT841305.1

<i>Bacillus thuringiensis</i>	Bacteria	CP015250.1
<i>Bacillus velezensis</i>	Bacteria	CP025939.1
<i>Bifidobacterium longum</i>	Bacteria	CP013673.1
<i>Bordetella bronchiseptica</i>	Bacteria	CM002881.1
<i>Brucella melitensis</i>	Bacteria	CP018494.1
<i>Campylobacter jejuni</i>	Bacteria	CP012689.1
<i>Caulobacter crescentus</i>	Bacteria	AE005673.1
<i>Cellvibrio japonicus</i>	Bacteria	CP000934.1
<i>Escherichia albertii</i>	Bacteria	AP014855.1
<i>Gordonibacter sp.</i>	Bacteria	LT827128.1
<i>Klebsiella pneumoniae</i>	Bacteria	CP025088.1
<i>Mycobacterium tuberculosis</i>	Bacteria	CP023640.1
<i>Ornithobacterium rhinotracheale</i>	Bacteria	CP006828.1
<i>Pseudomonas arsenicoxydans</i>	Bacteria	LT629705.1
<i>Salmonella enterica</i>	Bacteria	CP007400.2
<i>Serratia symbiotica</i>	Bacteria	LN890288.1
<i>Staphylococcus aureus</i>	Bacteria	CP012974.1
<i>Treponema pallidum</i>	Bacteria	CP020366.1
<i>Vibrio cholerae</i>	Bacteria	LT907989.1