

Characterization and Simulation of Metagenomic Nanopore Sequencing Data with Meta-NanoSim

Chen Yang (✉ cheny@bcgsc.ca)

Canada's Michael Smith Genome Sciences Centre <https://orcid.org/0000-0002-5144-7748>

Theodora Lo

Canada's Michael Smith Genome Sciences Centre

Ka Ming Nip

Canada's Michael Smith Genome Sciences Centre

Saber Hafezqorani

Canada's Michael Smith Genome Sciences Centre

René L Warren

Canada's Michael Smith Genome Sciences Centre

Inanc Birol

Canada's Michael Smith Genome Sciences Centre

Software article

Keywords: Metagenomics, Oxford nanopore sequencing, microbial abundance estimation, sequence simulation, chimeric reads

Posted Date: December 28th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-1125389/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **Characterization and simulation of metagenomic nanopore sequencing data with**

2 **Meta-NanoSim**

3 Chen Yang¹, Theodora Lo^{1,2}, Ka Ming Nip^{1,2}, Saber Hafezqorani^{1,2}, René L Warren¹, Inanc Birol^{1,3*}

4

5 1. 570 W 7th Ave, Canada's Michael Smith Genome Sciences Centre, BC Cancer, V5Z 4S6,

6 Vancouver, BC, Canada

7 2. Bioinformatics Graduate Program, University of British Columbia, Vancouver, BC, Canada

8 3. Department of Medical Genetics, University of British Columbia, Vancouver, BC, Canada

9 * Corresponding author

10

11 Chen Yang: cheny@bcgsc.ca

12 Theodora Lo: tlo@bcgsc.ca

13 Ka Ming Nip: kmnip@bcgsc.ca

14 Saber Hafezqorani: shafezqorani@bcgsc.ca

15 René L Warren: rwarren@bcgsc.ca

16 Inanc Birol: ibirol@bcgsc.ca

17

18

19 **ABSTRACT**

20 **Background:** Nanopore sequencing is crucial to metagenomic studies as its kilobase-long reads
21 can contribute to resolving genomic structural differences among microbes. However,
22 sequencing platform-specific challenges, including high base-call error rate, non-uniform read
23 lengths, and the presence of chimeric artifacts, necessitate specifically designed analytical tools,
24 such as microbial abundance estimation and metagenome assembly algorithms. When
25 developing and testing bioinformatics tools and pipelines, the use of simulated datasets with
26 characteristics that are true to the sequencing platform under evaluation is a cost-effective way
27 to provide a ground truth and assess the performance in a controlled environment.

28 **Results:** Here, we present Meta-NanoSim, a fast and versatile utility that characterizes and
29 simulates the unique properties of nanopore metagenomic reads. It improves upon state-of-the-
30 art methods on microbial abundance estimation through a base-level quantification algorithm.
31 Meta-NanoSim can simulate complex microbial communities composed of both linear and
32 circular genomes, and can stream reference genomes from online servers directly. Simulated
33 datasets showed high congruence with experimental data in terms of read length, error profiles,
34 and abundance levels. We demonstrate that Meta-NanoSim simulated data can facilitate the
35 development of metagenomic algorithms and guide experimental design through a metagenome
36 assembly benchmarking task.

37 **Conclusions:** The Meta-NanoSim characterization module investigates read features including
38 chimeric information and abundance levels, while the simulation module simulates large and
39 complex multi-sample microbial communities with different abundance profiles. All trained
40 models and the software are freely accessible at Github: <https://github.com/bcgsc/NanoSim>.

41 **KEYWORDS**

42 Metagenomics, Oxford nanopore sequencing, microbial abundance estimation, sequence
43 simulation, chimeric reads

44

45 **BACKGROUND**

46 Empowered by the rapid development of next-generation sequencing technologies,
47 metagenomic analysis has enabled comprehensive investigation of the genetic composition and
48 abundance of microbial communities. In investigating microbiomes, metagenomic sequencing
49 bypasses the need to culture each individual species by extracting DNA directly from their natural
50 habitat, making it feasible to study microbes that cannot be isolated or cultured in the laboratory
51 (1,2). Within the past few decades, the improved throughput and reduced cost of next-
52 generation DNA/RNA sequencing platforms has enabled a wide range of metagenomic studies of
53 environmental, pharmaceutical, and medical relevance (3–5).

54

55 Until recently, Illumina short-read sequencing (Illumina Inc., San Diego, CA) has been the
56 technology of choice for metagenomic sequencing projects due to its high throughput, low cost,
57 and low error rate. However, the reads generated by Illumina instruments are often too short
58 (<250 bp) to span inter- and intra-chromosomal homologous regions and suffer from intrinsic
59 biases, thus complicating downstream assembly and taxonomic analysis (6). As a third-
60 generation long-read sequencing technology, nanopore sequencing from Oxford Nanopore
61 Technologies Ltd. (ONT, Oxford, UK), is gaining traction in metagenomic research efforts, due
62 largely to the long read lengths it generates, as well as the portability of some of their sequencing

63 platforms (7). The N50 length metric (the shortest read length to be included for covering 50% of
64 the total number of bases sequenced) in a typical run is over 5 kbp (8) and the reported maximum
65 read length exceeds 2 Mbp. At the high end, whole bacterial or viral genomes may be captured
66 by a few sequencing reads (9,10), making it possible to disambiguate between even closely-
67 related strains. Since it was first announced, metagenomic sequencing by ONT has been playing
68 an essential role in real-time pathogen identification and clinical diagnosis, including during the
69 recent coronavirus disease 2019 pandemic (11–14).

70

71 Although a plethora of metagenomic analysis tools have been developed for short-read
72 sequencing data, the challenges associated with ONT reads, such as high error rate, non-uniform
73 error distributions, and chimeric read artifacts (8,15–17), call for analytical tools designed
74 specifically for long reads. For example, estimation of microbial abundance levels is traditionally
75 computed by counting the number of mapped reads followed by fine-tuning of ambiguous
76 mappings (18,19). This approach has been proven to be cost-effective for Illumina short reads
77 because of their uniform lengths. However, the accuracy of these tools would be understandably
78 impacted when applied on ONT reads, especially for lowly represented genomes, because of the
79 variable lengths and relatively high error rates (5 - 15% depending on the flowcell chemistry and
80 basecalling algorithm) compared to that of Illumina reads (typically less than 1%). In addition,
81 ONT sequencing projects on genomes, transcriptomes, and metagenomes, from prokaryotes to
82 eukaryotes, were all reported to contain certain problematic reads with gapped or chimeric
83 alignments, likely generated due to library preparation or sequencing artifacts (17,20–25). The
84 accuracy of reference-based abundance estimation using merely primary alignments may further

85 be affected by the presence of these chimeric reads, as well as reads that span the start position
86 of a circular genome. To the best of our knowledge, even the state-of-the-art program,
87 MetaMaps, does not account for chimeric reads, but simply uses an Expectation-Maximization
88 (EM) algorithm to disambiguate multi-mapped reads (26). In this work, we show that there is still
89 room for improving metagenomic abundance estimation, a proposition attainable by quantifying
90 aligned bases instead of reads, while leveraging chimeric read information.

91

92 In the process of tool development and benchmarking, a metagenomic ONT read simulator and
93 associated simulated datasets with known ground truth can save time and money. Ideally, such
94 a read simulator should reflect the true characteristics of the ONT platform and allow effective
95 evaluation of bioinformatics tools. In return, the evaluation results can guide the experimental
96 design of metagenomics projects, to help determine the desired sequencing depths and number
97 of replicates (27).

98

99 Currently, the only simulator that specifically simulates ONT metagenomic datasets is CAMISIM
100 (28). The workflow of CAMISIM is focused on the composition design of microbial community
101 given a taxonomy profile, while the abundance levels are drawn from a lognormal distribution.
102 The obvious drawback of this approach is that users cannot request the abundance levels as they
103 need. CAMISIM uses NanoSim (version 1) (15) as its engine to simulate ONT reads for each
104 genome separately once the composition of the community is determined. Following the same
105 idea, one can also use other existing ONT genomic simulators naively to simulate each composite
106 genome separately and then aggregate the reads according to the desired abundance. However,

107 it is impractical to simulate a large microbial community with hundreds or more genomes with
108 this approach, not to mention that the existing simulators for ONT reads are not designed to
109 model metagenomic specific features, such as chimeric reads and deviations in abundance levels.
110 More importantly, the simulation of abundance levels should be consistent with the
111 quantification method, thus merely mixing the reads from different genomes will yield a
112 compromised abundance profile. Taken together, we note that the previous version of NanoSim
113 can be upgraded to capture and simulate read properties specific to metagenomics, especially
114 the microbial abundance levels and chimeric reads – two key factors that may influence
115 metagenome assembly, taxonomy binning, and abundance estimation. Further, in real world
116 scenarios, viruses, bacteria, and fungi co-exist in complex microbial communities, hence it is
117 desirable to have the ability to simulate complex metagenomes comprising both circular and
118 linear genomes.

119

120 Here, we introduce Meta-NanoSim (released within NanoSim version 3), an ONT metagenome
121 simulator for complex microbial communities. Given a training dataset, Meta-NanoSim
122 characterizes read length distributions, error profiles, and alignment ratio models. Optionally, it
123 also detects chimeric reads and estimates microbial abundance levels. In our benchmarks, the
124 performance of the metagenomic abundance estimation feature of Meta-NanoSim surpasses the
125 current state-of-the-art methods. The chimeric read detection feature also improves the read
126 length modelling, and thus simulating this artifact of the technology may challenge metagenomic
127 analytical tools with more realistic erroneous reads. Through benchmarking experiments
128 comparing simulated reads with empirical datasets, we show that Meta-NanoSim preserves the

129 key characteristics of ONT metagenomic reads. Further, we showcase the usability and utility of
130 Meta-NanoSim with a metagenomic assembly task.

131

132 **IMPLEMENTATION**

133 **Meta-NanoSim general design**

134 Meta-NanoSim is implemented in Python as two sub-modules in the NanoSim suite: `meta` in
135 characterization and simulation stages respectively. It learns the technical and metagenomic-
136 specific features of ONT reads in the characterization stage, builds statistical models, and applies
137 them in the simulation stage (Fig. 1). In the characterization stage, it takes ONT metagenomic
138 reads and a reference metagenome as input to infer the ground truth through sequence
139 alignments. Based on those alignments, it models the read length distributions (aligned and
140 unaligned part) via kernel density estimation and basecall events via mixture statistical models.
141 In addition to existing NanoSim features, we introduce two new analyses in its characterization
142 pipeline; chimeric read analysis for genome/metagenomes and abundance estimation for
143 metagenomic datasets.

144

145 To simulate a metagenomic dataset, besides the pre-trained model from the characterization
146 stage, a list of reference genomes of users' choice is required as input to use for simulation,
147 together with their abundance levels and DNA topology (i.e., linear or circular). The tool can
148 optionally stream reference genome sequences from either RefSeq (29) or Ensembl (30)
149 automatically without requiring extra disk storage, which facilitates large microbial community
150 simulations. Since microbial sequencing projects are often carried out in a multi-sample or multi-

151 replicate fashion, Meta-Nanosim is designed to simulate multiple samples in one batch with user-
152 defined abundance level profiles as input.

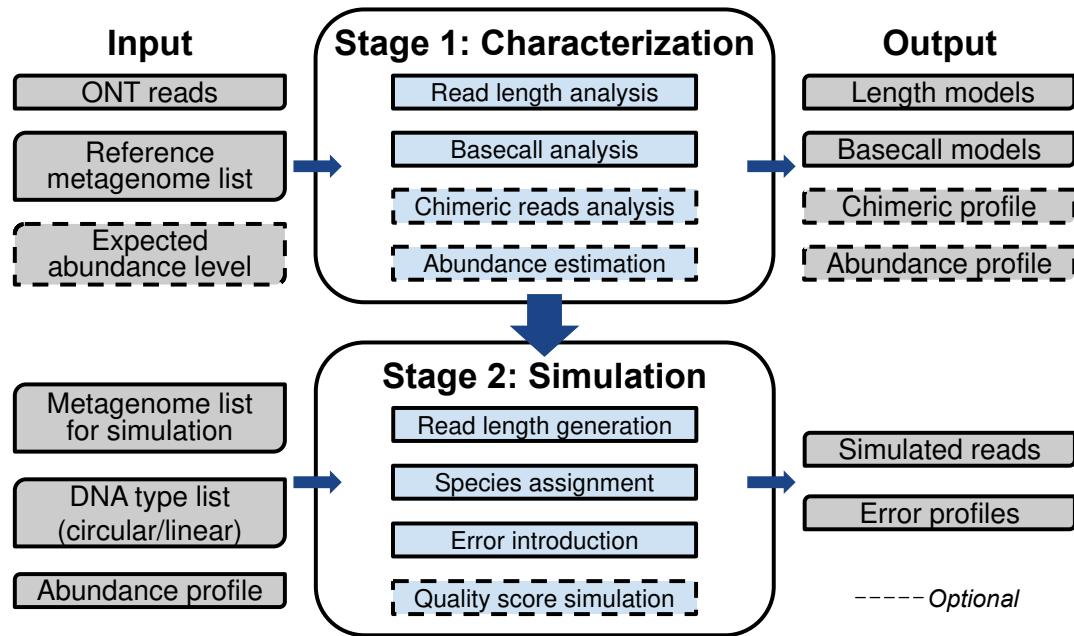


Fig. 1 Meta-NanoSim workflow. Meta-NanoSim consists of two stages: the characterization stage and the simulation stages. In the characterization stage, given a training dataset and reference metagenome, Meta-NanoSim builds models for the read length distributions and basecall events. It optionally profiles chimeric read artifacts and quantifies an abundance profile. It can also calculate the deviation between expected and estimated abundance levels. In the simulation stage, Meta-NanoSim takes four inputs: 1) The genome list for simulation (local or ftp path), 2) DNA type list, 3) abundance profiles to be simulated, and 4) the models generated from the characterization stage. Meta-NanoSim outputs simulated reads and error profiles for ground truth.

153

154 **Chimeric read detection and simulation**

155 During alignment, one ONT read may have multiple sub-alignments to different *loci* on the
156 reference genome/metagenome. When the query and reference coordinates of two or more sub-
157 alignments do not overlap, we define them as a set of compatible alignments. Finding the best
158 compatible alignment set problem is akin to finding the largest compatible interval set in
159 computer science. For each read, we exhaustively search for all compatible alignments for each
160 sub-alignment to generate a list of compatible alignment sets (Fig. S1 in Additional File 1). We

161 then select the best element from the list for downstream analysis, based on alignment quality
162 and total aligned length. If, for a given read, the best element contains two or more compatible
163 alignments, the read is considered as chimeric and its aligned length, gap length, and source
164 species (specific to metagenome mode) are modelled for simulation. One exception here is reads
165 that span the start position of a circular genome; these reads are detected but not considered
166 chimeric. The two sub-alignments in such cases are concatenated as one long alignment to count
167 towards the aligned length and source species.

168

169 As for the source species of each segment in chimeric reads, we observed that the source species
170 of the subsequent segment has a higher probability to be the same as the previous species, and
171 consequently the probabilities of it being other species are decreased, while their relative
172 abundance levels to each other remain the same. We build a hidden Markov model to simulate
173 this observation, where start probability is the input abundance, the emission probability
174 represents which species the next segment is coming from given the previous one, and the
175 transitional probability of species is the change of abundance levels in the underlying Markov
176 chain. Since species to be simulated, namely the states in a Markov model, may be different
177 between the training and simulation metagenome, we generalize the emission probability as a
178 single value called shrinkage rate s ($0 < s \leq 1$). This parameter describes the reduction of
179 abundances (probabilities) of other species, while maintaining the relative abundances among
180 them. Assuming the input abundance is $\{p_A, p_B, p_C, \dots, p_N\}$ for n species, when the first segment
181 comes from species A, the transitional probabilities for the other species would become $\{s \times p_B,$
182 $s \times p_C \dots s \times p_N\}$ and the transitional probability for A would be inflated as $1 - s \times \sum_{i=B}^N p_i$. To learn

183 s , all segments in chimeric reads are divided into overlapping pairs, and the probability for the
184 source species of the second segment being different from the first one is recorded. In this way,
185 we can calculate the reduction of abundance for every species. The average reduction is the
186 shrinkage rate, and the inflated abundance for being from the same species can be inferred as
187 well. The shrinkage rate can also be adjusted by the user, to 1 for example, if one assumes the
188 source species for each segment depends solely on the input abundance levels, and is
189 independent from the previous segment.

190

191 To summarize the simulation of chimeric reads, Meta-NanoSim first determines the number of
192 segments to be simulated based on a Geometric distribution. If the number of segments is
193 greater than or equal to two, it is a chimeric read. Then, Meta-NanoSim generates the lengths of
194 each segment and gap between them using kernel density estimation learnt from empirical reads.
195 The source species of the first segment is randomly picked based on the input abundance level.
196 Starting from the second segment, the abundance levels are re-computed based on the previous
197 species and s . The source species is determined one after another, and then sequences are
198 extracted, mutated with purposely introduced errors, concatenated to each other, and then
199 unaligned regions are added to the reads in the same process as non-chimeric reads.

200

201 **Abundance estimation**

202 Existing abundance estimation methods generally quantify the number of mapped reads or k -
203 mers, under the presumption that all reads have equal lengths. However, since the ONT read
204 length varies several orders of magnitude, it is likely that the mean read length for each species

205 would be different. When all species are equally and deeply sequenced, according to the central
206 limit theorem, the standard deviation of mean lengths would scale with $1/\sqrt{n}$, where n is the
207 number of species. In reality, low-abundant species may have a higher standard deviation
208 because there are fewer sequences representing it. We observe the mean read lengths of
209 uniquely aligned reads to vary substantially in datasets where species abundance levels are
210 logarithmically distributed, necessitating base-level instead of read-level quantification
211 algorithms (Fig. S2 in Additional File 1).

212

213 Another key challenge that confounds short-read metagenomic analysis is ambiguously aligned
214 reads. In ONT datasets, however, most reads are long enough to span inter- and intra-species
215 homologous regions, and the chimeric read detection feature can resolve the estimation for
216 reads that have multiple sub-alignments and for reads that span across the start site of a circular
217 genome. For the remaining small fraction of multi-aligned reads between closely-related species,
218 the estimation for them can be optimized using the EM algorithm.

219

220 The EM algorithm (Algorithm 1) first processes uniquely aligned reads to calculate a baseline
221 abundance profile. Then it starts the expectation step, which is to assign multi-aligned bases
222 proportionally to their respective species based on their relative abundances. Next, in the
223 maximization step, these assigned multi-aligned bases are used to update the abundance profile.
224 The algorithm then goes back to the expectation step to update the fractions of multi-aligned
225 bases based on the new abundance profile. The E and M steps alternate until the difference in
226 abundances between two rounds is lower than a threshold (default: 1%). Note that the

227 abundance levels are in the units of relative genomic DNA weight, and they can be used to
228 calculate genome copy numbers when divided by the respective genome sizes.

229 **Algorithm 1** EM for metagenome abundance estimation

230 abundance_list = {species1: abundance1; species2: abundance2; ...}
231 base_count = {species1: count1; species2: count2; ...}
232
233 Start processing uniquely aligned reads:
234 for each uniquely aligned read and its source species:
235 base_count[species] += aligned bases
236 abundance_list = {species: base_count[species]/ sum(base_count[species])}
237
238 Start processing multi-aligned reads
239 while diff >= min(abundance_list.values()) * 0.01:
240 **E-Step:**
241 for each multi-aligned read:
242 read_abun = the sum of abundances for all possible species for that read
243 for each possible species:
244 fraction = aligned bases * abundance_list[species] / read_abun
245 base_count[species] += fraction
246 **M-Step:**
247 abundance_list = {species: base_count[species]/ sum(base_count[species])}
248 diff = |abundance_list - prev_abundance_list|
249

250 Meta-NanoSim offers two estimation methods, one with the chimeric read detection and one
251 without. When chimeric read detection is enabled, all subalignments are used for computing;
252 when it is disabled, only the primary alignments are used. Meta-NanoSim records the aligned
253 bases for each sub-alignment towards their source genome, and then uses EM algorithm to assign
254 multi-aligned segments proportionally to their putative source genomes iteratively.

255

256 **Abundance deviation simulation**

257 Due to library preparation and sequencing depth, there is a deviation between the sequenced
258 abundance levels and expected ones. Meta-NanoSim offers a means to simulate the abundance

259 deviation with user-defined lower and upper deviation boundaries. We noticed a weak positive
260 correlation between genome size and abundance deviation in our analysis. During simulation we
261 first randomly draw a list of relative error between the deviation boundaries. Next we assign
262 these errors to each genome based on their sizes, namely larger deviations are assigned to larger
263 genomes and smaller ones are assigned to smaller genomes. After applying errors, the
264 genome/species abundances are renormalized to have a total abundance of 100%.

265

266 **RESULTS**

267 Here, we first present the design logic and the performance of two key features, chimeric read
268 detection and abundance estimation. To show the similarity between simulated reads and
269 experimental reads, we generated two simulated datasets using models learned from
270 experimental data, and we compared the performance of Meta-Nanosim with that of CAMISIM.
271 We illustrate how Meta-NanoSim is capable of simulating large complex microbial community by
272 simulating a dataset with 125 species with abundance levels estimated from a human saliva
273 sample. Finally, we showcase a potential use case of Meta-NanoSim with a set of simulated data,
274 benchmarking the metagenomic assembler MetaFlye (31) and demonstrating its scalability with
275 increasing sequencing depth. The input experimental data used for the assembly benchmarks are
276 two sets of publicly available mock community ONT sequencing reads, each dataset containing
277 the same 10 microbial species (eight bacteria and two fungi) but with different abundance
278 distributions. We denote the dataset with evenly distributed abundance levels as the *Even*
279 dataset and the one with logarithmically distributed abundance levels as the *Log* dataset
280 (detailed in Supplementary Methods in Additional File 1).

281 **Chimeric read characterization and simulation**

282 Chimeric reads are also called split reads because one read is split into two or more sub-
283 alignments that are aligned to distinct regions of the genome/metagenome. They may be created
284 due to sequencing artifacts or structural variants when the reference metagenome is not
285 comprehensive. Previous studies reported that chimeric reads represent a non-negligible fraction
286 of ONT sequencing datasets ranging from 1.7% to 8.17% depending on the sequencing kits and
287 identification thresholds (17,20,22,23). In the metagenome datasets used in our study, after
288 ruling out structural variants, we have identified a similar fraction of reads in this category: 2.17%
289 (75,628 reads) in the *Even* dataset and 1.67% (68,444 reads) in the *Log* datasets. These reads are
290 free of known adapters, so their presence may impact downstream analyses, such as assembly,
291 taxonomy binning, and quantification, even after adapter trimming. When aligned to their
292 respective reference genome sequence(s), ONT reads may contain unaligned or soft-clipped
293 regions. In our tests, the chimeric read detection feature of Meta-NanoSim significantly reduced
294 the length of these unaligned regions, which explained why some of the reads have over 1 kbp
295 long unaligned portions (Fig. 2A). As seen in Fig. 2B, the length distributions of the gaps between
296 split alignments follow multi-modal distributions. Meta-NanoSim uses kernel density estimation
297 to model them, with results exhibiting strong similarity between the length distributions of
298 simulated and experimental sequences. We also noticed that the number of segments each read
299 contains can be described as a geometric distribution and the mean probability can be learnt
300 from experimental data (Fig. 2C). On average, each read contains 1.03 segments for both data
301 sets under study.

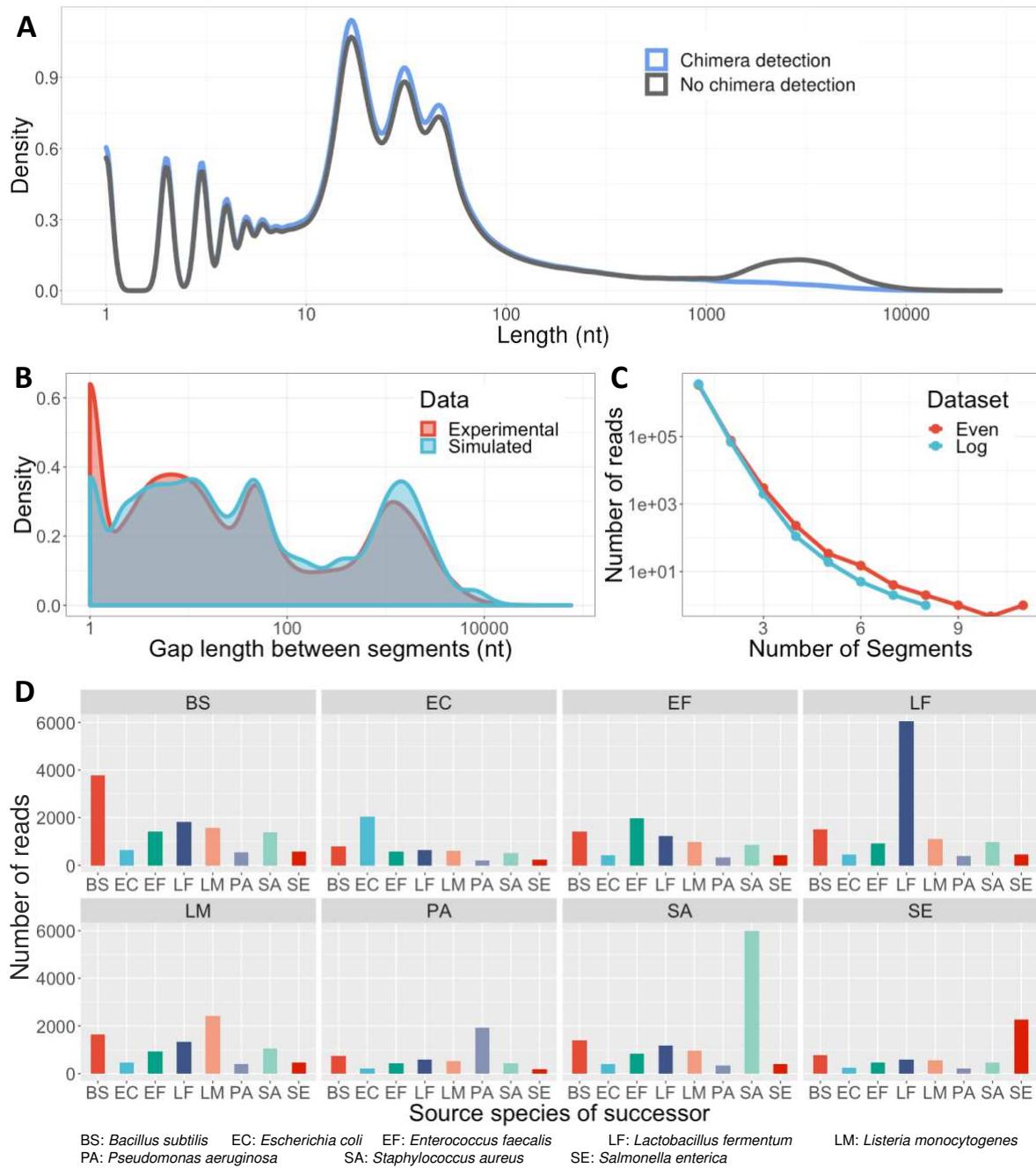


Fig. 2 Chimeric read detection and simulation. **A.** The length distribution of the unaligned regions of reads with or without chimeric read detection for the *Log* dataset (x-axis in logarithmic scale). **B.** The performance of gap length simulation for the *Log* dataset (x-axis in logarithmic scale). **C.** The number of segments each read contains for the *Even* and *Log* datasets. **D.** All segments in chimeric reads in the *Even* dataset are converted into overlapping pairs. Each facet represents one source species of the first segment and the x-axis represents the source species of the second segment. Each facet shows the probability of the second segment given the source species of the first one. *Cryptococcus neoformans* and *Saccharomyces cerevisiae* are excluded here due to their low abundances.

303 Based on the source species to which each split alignment belongs, chimeric reads can be
304 classified as “intra-species-chimeric” or “inter-species-chimeric”. It is observed that the source
305 species of the first segment is controlled by the abundance level, and the one of the subsequent
306 segment is influenced by the identity of the previous species (Fig. 2D). For all eight species, the
307 source species of the subsequent segment more likely to be the same as the previous one, while
308 the probabilities of being other species are proportional to the abundance level. We postulate
309 that this is because DNA molecules of the same species are more likely to gather near the
310 nanopore than being homogeneously dispersed in the buffer. Regardless of the actual cause, this
311 phenomenon can be approximated as a simplified hidden Markov model with a generalized
312 emission probability, which is defined as shrinkage rate s here. To our calculation, s is equal to
313 0.77 for the *Even* dataset and 0.73 for the *Log* dataset, suggesting that its value may be stable
314 across datasets.

315

316 **Abundance estimation**

317 To benchmark, we compared the performance of four abundance estimation methods: Meta-
318 NanoSim estimation with chimeric read detection, Salmon quantification with --meta option
319 (Salmon) (32), the base-level estimation reported in the paper that released the dataset (denoted
320 as “Data Note” from hereon) (8), and MetaMaps. For Meta-NanoSim estimation, we performed
321 an ablation study that removes key components of the algorithm step by step, including
322 estimation on read-level with chimeric read detection (Meta-NanoSim CR) or with EM algorithm
323 (Meta-NanoSim ER), estimation on base-level (Meta-NanoSim B), base-level with chimeric read
324 detection (Meta-NanoSim CB), base-level with EM algorithm (Meta-NanoSim EB), and base-level

325 with chimeric read detection fine-tuned by EM algorithm (Meta-NanoSim ECB). All compared
326 methods, except for MetaMaps, are computed based on Minimap2 alignments. We compared
327 the estimated abundances to the expected ones provided by the manufacturer and reported R-
328 squared, standard deviation, and percent error to assess the performance.

329

330 In general, all base-level quantification methods performed better than read-level quantification
331 methods (Salmon, MetaMaps, Meta-NanoSim CR and ER), and Meta-NanoSim base-level
332 estimations have the highest correlation with the expected abundances (Table 1, Fig. S3 in
333 Additional File 1). For the *Even* dataset, all four Meta-NanoSim base-level methods performed
334 similarly; the stand-alone base-level quantification has the highest R-squared value for the *Even*
335 dataset, while the chimeric read detection helped reduce the percent error, mainly for low-
336 abundance species *Cryptococcus neoformans* (Fig. S3 in Additional File 1). MetaMaps, as a read-
337 level quantification method designed specifically for ONT metagenomic data, although ranked
338 highest among this category, showed a big discrepancy compared to base-level methods with
339 almost twice the percent error. For the *Log* dataset, Meta-NanoSim base-level estimations also
340 had similar R-squared values, higher than other methods. Since the metrics are close to each
341 other for the *Log* dataset and the performance on low-abundance species may be overshadowed
342 by high-abundance species, we also computed the coefficient of correlation and error between
343 log-transformed estimated and expected abundances. After log-transformation, Salmon, Meta-
344 NanoSim ECB, and Meta-NanoSim EB showed similar performances, with Salmon having the
345 highest correlation and Meta-NanoSim ECB having the lowest percent error. The metrics for
346 Meta-NanoSim estimation without EM, on the other hand, dropped significantly due to difficulties

347 in differentiating multi-mapped reads for low-abundance species. When estimating the
 348 abundance levels for the *Log* dataset, Minimap2 incorrectly assigned 18,212 reads to the *E.*
 349 *faecalis* genome as primary alignments, when they can also be aligned to an inter-species
 350 homologous region in the *L. monocytogenes* genome. In fact, *E. faecalis* is a low-abundance
 351 species in the *Log* dataset with only 33 unique alignments. Therefore, the methods with EM
 352 algorithm resolved the multi-aligned reads problem, indicating that EM can be advantageous for
 353 datasets consisting of similar genomes but with large variances in abundance levels.

Table 1 Statistical analysis of the abundance estimation results compared to expected abundances.

Tool	Algorithm				Even dataset			Log dataset			
	E	C	B	R	R ²	Std	PE	R ²	Log R ^{2*}	Std*	PE
Meta-NanoSim	✓	✓	✓		0.7463	0.0225	144.5317	1.0000	0.9920	0.1818	256.7856
		✓	✓		0.7465	0.0225	144.0877	0.9999	0.7899	0.9295	53349.95
	✓		✓		0.7498	0.0224	145.6229	1.0000	0.9917	0.1843	260.6973
			✓		0.7499	0.0224	145.4656	1.0000	0.7895	0.9304	53443.94
		✓		✓	0.4305	0.0337	326.3432	0.9980	0.7778	0.9560	57527.07
	✓			✓	0.4396	0.0335	313.1466	0.9980	0.7776	0.9565	57651.72
Salmon	✓			✓	0.4269	0.0339	318.1502	0.9978	0.9955	0.1366	261.1368
Data Note			✓		0.6702	0.0257	181.1667	0.9998	0.9863	0.2374	359.9314
MetaMaps	✓			✓	0.4420	0.0334	258.4563	0.9979	0.9652	0.3781	1366.705

R²: R-squared, Std: standard deviation, PE: summation of percent error

E: EM algorithm, B: base-level quantification, R: read-level quantification, C: chimeric read detection

* The expected and estimated abundances are log-transformed before calculating R-squared value and standard deviation.

354
 355 To recapitulate our findings, we chose a second logarithmically distributed mock microbial
 356 community, denoted as the *Adp* dataset (33) (Detailed in Supplementary Method in Additional
 357 File 1), and repeated the quantifications with Meta-NanoSim base-level methods, Salmon, and
 358 MetaMaps (Table S1 in Additional File 1). Again, all four Meta-NanoSim methods performed
 359 similar to each other with the highest correlation to the expected values. MetaMaps
 360 quantification, although with the lowest percent error among all compared methods, showed a

361 much lower correlation in terms of R-squared and standard deviation. Taken all together, Meta-
362 NanoSim base-level quantification after chimeric read detection and fine-tuned by EM algorithm
363 balanced correlation and percent error performed robustly well, making it preferable for
364 naturally occurring microbial communities with various abundances.

365
366 Because of the deviation between expected and estimated abundance levels, we introduce a
367 feature that can simulate this observation. We compared the deviation between expected
368 abundances and experimental data, and the simulation results of NanoSim and CAMISIM (Fig. 3).
369 The distribution of abundance deviations of experimental data and Meta-NanoSim simulated
370 reads are statistically the same (Kolmogorov-Smirnov test p -value = 0.787), while for CAMISIM
371 simulated reads, the distribution is pronounced as it does not provide this feature.

372

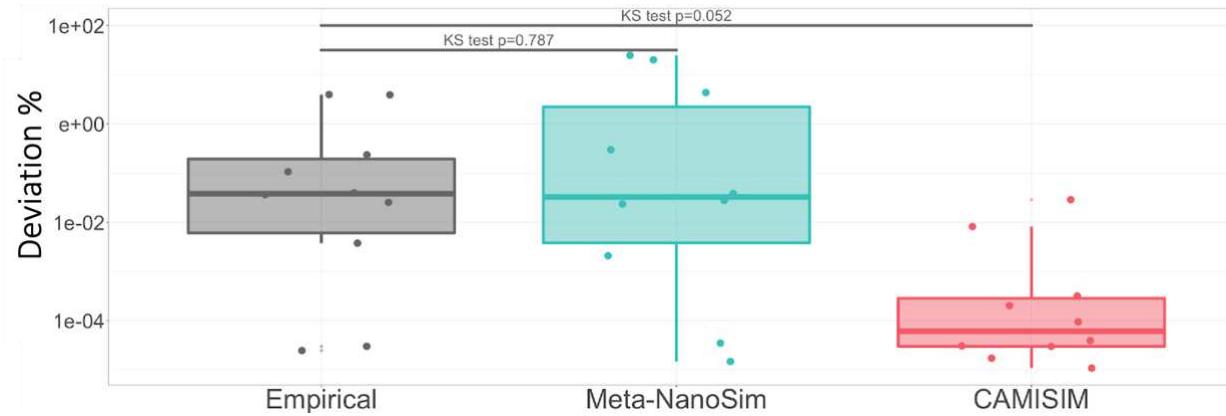


Fig. 3 Abundance level deviations between experimental and simulated metagenomic reads. In this plot, each dot represents a microbial genome, and the y-axis represents the deviation in percentage between the expected values and experimental/simulated values.

373
374 **Comparison between simulated and experimental datasets**
375 To demonstrate the performance of Meta-NanoSim, we trained it with the *Log* dataset and

376 compared the simulated datasets against the result of CAMISIM. With eight processors,
377 simulation of one million reads took under 20 minutes (or under 160 CPU-minutes) for Meta-
378 NanoSim, while CAMISIM required more than six hours to complete.

379

380 The read lengths of simulated datasets from Meta-NanoSim follow the empirical length
381 distribution closely, with a median read length peak at 4,040 nt (3,994 nt for experimental reads)
382 (Fig. 4A). In contrast, the lengths of CAMISIM-simulated reads deviate far from those of the
383 experimental data. Because CAMISIM is only compatible with an old version of NanoSim and uses
384 hard coded pre-trained models learnt from a genomic sequencing data, it was not possible to
385 test it coupled with Meta-NanoSim. Moreover, the length distribution of unaligned regions on
386 Meta-NanoSim simulated reads captures the patterns in experimental reads well, with multiple
387 peaks below 100 nt. In contrast, the lengths of unaligned part in CAMISIM reads are inflated as it
388 does not detect nor simulate chimeric reads. Both Meta-NanoSim and CAMISIM mimic the
389 mismatch and deletion events well when compared to the experimental dataset (Fig. 4B), which
390 demonstrates the robustness of NanoSim mixture statistical models. However, Meta-NanoSim
391 simulates insertion and match events better than CAMISIM, also due to the change of model.

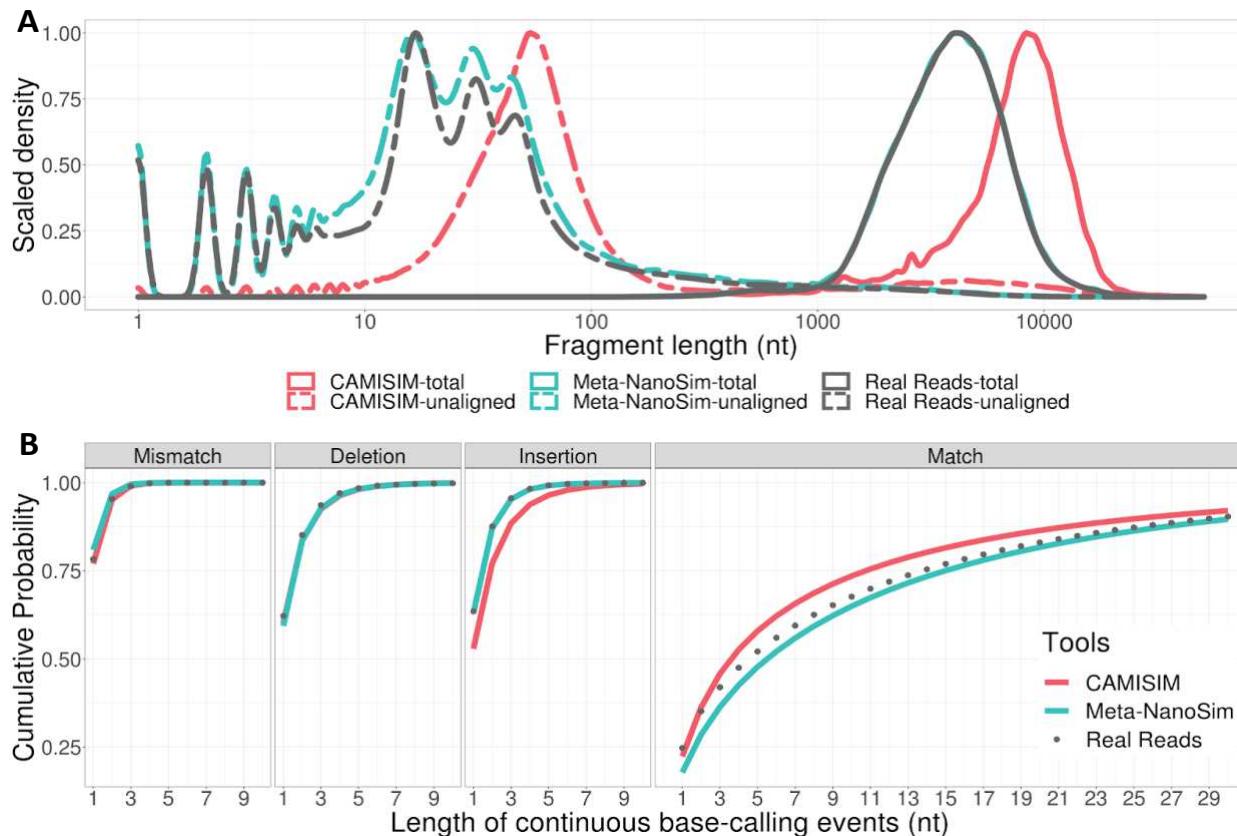


Fig. 4 Performance of Meta-NanoSim and CAMISIM in simulating one million reads from the *Log* dataset. A. Comparison of read length distributions in the experimental vs. simulated reads (x-axis in logarithmic scale). Unaligned length represents the length of unaligned part of each aligned read. B. Cumulative probability distributions of the lengths of matches/errors in experimental and simulated reads.

392

393 Additionally, we challenged Meta-NanoSim with two simulation tasks to mimic real-world use
 394 cases. First, we simulated two samples at the same time with the pre-trained model from the *Log*
 395 dataset. Each sample contained one million reads from the seven species from the *Adp* dataset
 396 with different abundance levels. Meta-NanoSim simulation finished within 51 min with eight
 397 processors. Although the metagenome to be simulated is completely different from the training
 398 one, simulated reads exhibited similar read features as the experimental data (Fig. S4 in
 399 Additional File 1), and the abundance levels are on average 99.93% in accordance with the
 400 expected values. Next, we randomly picked a saliva sample from the Human Microbiome Project

401 (HMP) and simulated ONT reads using the same microbial composition (34). The abundance
402 levels of the 125 different bacteria strains range between 4.28% to 12.49%. By streaming
403 reference genomes from RefSeq directly, it took Meta-NanoSim less than three hours to simulate
404 10 million reads in total.

405

406 **Showcase application in assembly benchmarking**

407 Meta-NanoSim simulated reads can be used to benchmark and assess the performance of
408 bioinformatics algorithms. To demonstrate an application of Meta-NanoSim, we simulated four
409 sets of data with 1, 2, 4, 10 million reads to assess the correctness, scalability and robustness of
410 metaFlye, a long-read metagenomic assembler. We used the models learnt from the *Log* dataset
411 and the abundance levels of the 10 species are shown in Fig. 5. With 128 threads, runtimes
412 ranged from one hour to seven hours (Table S2 in Additional File 1). The maximum resident set
413 size for 10 million reads dataset was 212 GB, but intermediate files occupied over 10 TB of disk
414 space during consensus-building stage. We also tried to assemble a larger dataset of 20 million
415 reads, however the assembly failed after 30 days with an out-of-memory error on a 1 TB RAM
416 server. Before completing, the maximum resident set size was 1.75 TB, and the most time-
417 consuming stage was the graph simplification stage (two weeks).

418

419 In total, MetaFlye reconstructed 27.81 Mbp sequences, adding up to 43.72% of the total
420 reference genome for the simulated dataset with 4M reads. These metrics are similar to the
421 reported assemblies using the original training dataset with 3.48 M reads (28.20 Mbp assembled
422 length that covered 46.00% of the reference metagenome) (31). Generally speaking, the average

423 genome coverage is positively correlated with the number of sequencing reads and abundance
424 levels, and accordingly, the genome reconstruction fraction and NGA50 length are positively
425 correlated with the coverage (Fig. 5). As expected, sub-1x coverage genomes have very poor
426 reconstructions. Between 1x and 10x, the positive correlation is mirrored in multiple species,
427 including *B. subtilis*, *S. cerevisiae*, *E. coli* and *S. enterica*. When the coverage reaches 10x,
428 metaFlye is able to reconstruct the genome to nearly 100% (*S. cerevisiae* in the 4M dataset, Fig.
429 5). When the coverage reaches 30x, the NGA50 length can cover the whole genome size (*B.*
430 *subtilis* in the 2M dataset). Similarly, the contigs reconstructed for these two species decrease as
431 the number of reads increase, showing how increasing sequencing depth can help assembling
432 genomes into one contig for *B. subtilis* and nearly one contig per chromosome for *S. cerevisiae*.
433 In contrast, when the coverage is not high enough, the results may be fragmented or even mis-
434 assembled in the case of *E. coli* and *S. enterica*. However, a higher coverage does not necessarily
435 lead to a better assembly quality. The reconstruction of *L. monocytogenes* deteriorates with more
436 reads when the fold coverage exceeds 1000x. Although the genome fraction remains 100%, the
437 NGA50 length is only half or less of the genome size for the 2M, 4M, and 10M datasets. The drop
438 in NGA50 length can be explained by the increasing number of reconstructed contigs and mis-
439 assemblies (Fig. 5). With one million reads, there are only four contigs that can be mapped to the
440 *L. monocytogenes* genome, and there are no mis-assemblies detected. However, we think that
441 the higher sequencing depth above 1000X has led to many more misassemblies, adversely
442 effecting the assembly contiguity as measured by the NGA50 length metric.

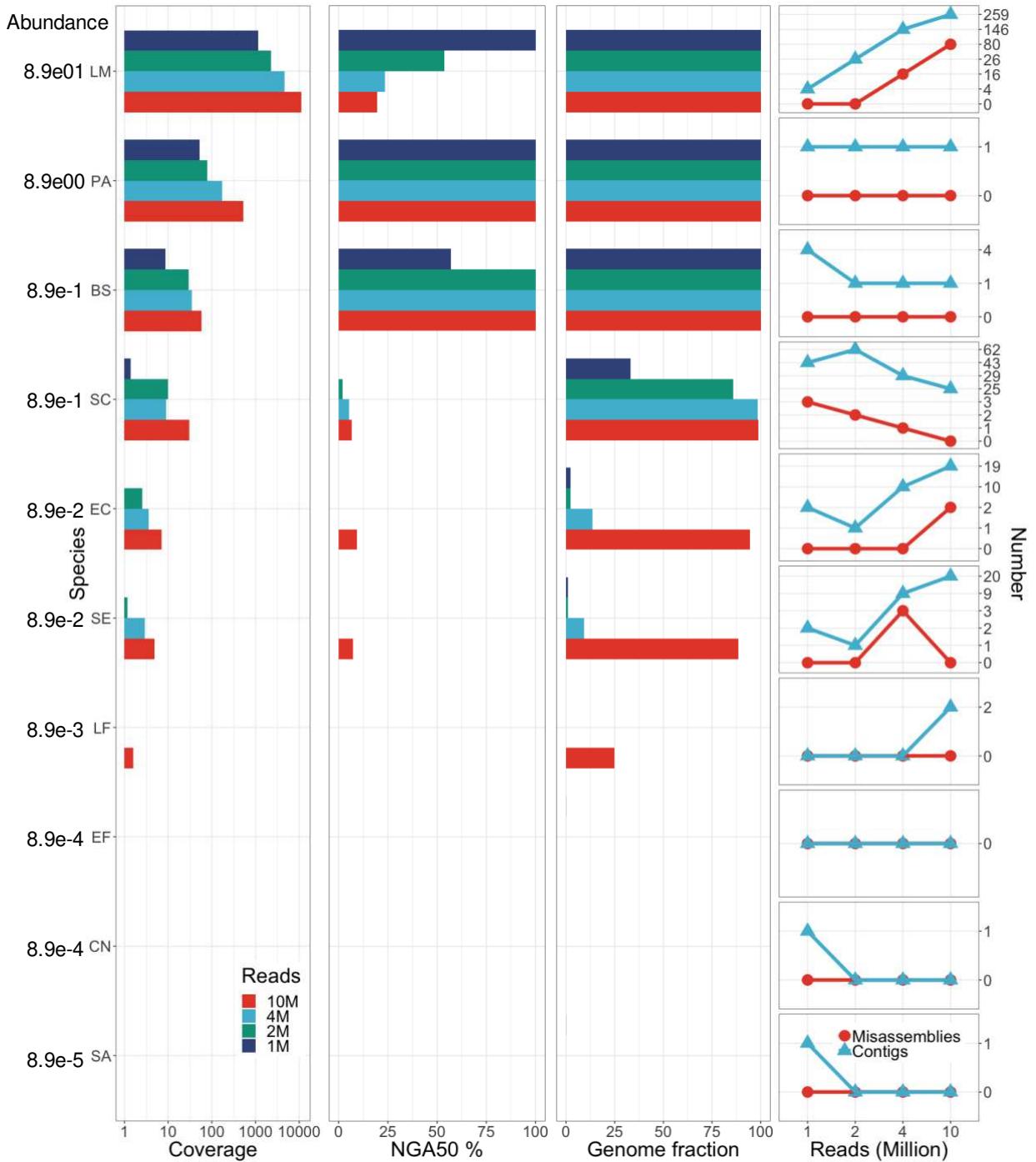


Fig. 5 metaFlye assemblies with four sets of simulated metagenome sequencing data. The four sets of simulated datasets include 1, 2, 4 and 10 million reads, respectively. The abundance is the expected abundance level during simulation. The coverage panel shows the average read depth including plasmids (x-axis in logarithmic scale). NGA50 % represents the NGA50 length divided by the reference genome size. Genome fraction is a proportion between the assembled sequences and each corresponding genome. The right-most panel shows the number of misassemblies and assembled contigs as the number of simulated reads increases. BS: *Bacillus subtilis*, CN: *Cryptococcus neoformans*, EC: *Escherichia coli*, EF: *Enterococcus faecalis*, LF: *Lactobacillus fermentum*, LM: *Listeria monocytogenes*, PA: *Pseudomonas aeruginosa*, SA: *Staphylococcus aureus*, SC: *Saccharomyces cerevisiae*, SE: *Salmonella enterica*.

443 **DISCUSSION**

444 The applications of nanopore sequencing on metagenomic projects are rapidly expanding,
445 motivating the development of metagenomic analysis tools tailored for this specific data type. In
446 this work, we bring two main contributions to ONT metagenomic analysis tasks: (1) introduction
447 of a new base-level quantification method for metagenomic abundance estimation; and (2) an
448 upgrade of NanoSim for metagenomic characterization and simulation.

449

450 Reference-based metagenomic abundance estimation is key to investigating the microbial
451 composition of an environment, enabled by emerging sequencing technologies. The long read
452 length of ONT reads provides an opportunity to resolve homologous regions between species or
453 strains, but complications arise due to their high error rates and non-uniform read lengths.
454 Existing methods, primarily developed for short read technologies, generally assume uniform
455 read lengths and therefore only need to count the number of mapped reads or k -mers. However,
456 it is unreasonable to treat, say, a 100 bp long read and a 1000 bp long read the same when
457 calculating their contributions to the genome abundance. Our results echoed that it is necessary
458 to quantify microbial abundances on a base-level rather than read-level to better leverage this
459 data type. In addition to their higher error rates, a small yet substantial fraction of ONT reads are
460 chimeras, which may obscure the accuracy of estimates. The chimeric read detection feature in
461 Meta-NanoSim works by searching for best compatible alignments and helps reduce the percent
462 error in microbial abundance estimation. For multi-aligned segments where source species
463 cannot be differentiated, we adopted an EM algorithm to optimize the proportional
464 contributions of these segments to each potential source species. Our benchmarking results

465 demonstrate that the combination of these three components can improve correlations with
466 expected abundances. We note that Meta-NanoSim quantification performs better when the
467 abundance levels are more uniform or when low-abundance microbes do not share large
468 homologous regions with high-abundance microbes. Depending on whether the user wishes to
469 achieve a higher correlation or lower percent error, they can choose to disable or enable chimeric
470 read detection, respectively. Although our work is limited to reference-based quantification, we
471 expect it to inspire the design of reference-free methods and eventually have a broader
472 application.

473

474 Built on top of abundance estimation, Meta-NanoSim is able to simulate datasets with desired
475 abundance profiles. It can also recapitulate the abundance level deviation from expected values,
476 an especially useful feature for designing sequencing projects. When the rough abundance of a
477 microbial community is known, it is essential to know how deep the sequencing needs to be, so
478 that each species can be covered by a sufficient number of sequencing reads. However, when
479 the sequenced abundance differs from the expected value, simulated data with abundance
480 variations true to the platform can inform the relationship between sequencing depth and
481 abundance levels.

482

483 The general workflow of characterization and simulation of Meta-NanoSim follows the same
484 paradigm of the previous versions of NanoSim. The chimeric read detection in characterization
485 stage provides a means to profile all chimeric reads in a library regardless of its root cause. When
486 the reference metagenome is inclusive, the chimeric reads are likely introduced by library

487 preparation and sequencing artifacts; while in reality, since the detection relies on alignment,
488 some chimeric reads may also be attributed to structural variants when the source genome is not
489 present in the reference. In this case, the output of the characterization stage can be used to
490 further investigate such events with specifically designed statistically models and algorithms.

491

492 The three new main features added to Meta-NanoSim are 1) chimeric read simulation, 2) the
493 ability to stream reference genomes from online servers, and 3) the simulation of a metagenome
494 composed of a mixture of both linear and circular genomes. As chimeric reads may interfere with
495 downstream analyses, simulated datasets with these artifacts are needed for more accurate
496 performance assessment. Characterizing this feature and introducing it to simulated reads will
497 also diversify error types in the reads, helping to improve the robustness of related algorithms.

498 Reference genome streaming is uniquely advantageous when simulating a large metagenome
499 with hundreds of species. It eliminates the reference genome downloading step for users and
500 saves disk space while keeping the runtime reasonable. Similarly, since metagenomes are
501 naturally composed of both linear and circular genome topologies, having a simulated dataset
502 supporting this important characteristic will add credibility to benchmarking results and better
503 forecast performance with experimental data.

504

505 The benchmarking on a metagenome assembly task showcased that Meta-NanoSim can facilitate
506 relevant tool development as well as guide sequencing projects. The resulting assembly quality
507 of Meta-NanoSim simulated reads is comparable to that of the experimental data with similar
508 coverage. Although publicly available mock community sequencing data provide a more realistic

509 training and test set, simulated data provide a ground truth and has no limit in size, making them
510 perfect for testing the accuracy and scalability of algorithms. Through the use of simulated
511 datasets, we demonstrated that metaFlye assembler performs best when the species coverage
512 is between 10 and 1000-fold. To ensure a successful assembly of low abundance species, it is
513 suggested to calculate the number of reads needed given an estimated abundance first to ensure
514 just enough coverage without wasting resources. For example, it takes 10 million reads to achieve
515 10-fold coverage for a 0.1% abundant species with a genome size of 5Mb. When assembling real
516 microbial communities with highly variable abundance levels, we recommend multiple rounds of
517 assembly with different sample sizes to achieve the best performance for both high- and low-
518 abundance microbes. In addition, developers may analyze in depth the mis-assemblies and errors
519 in assembled contigs with the ground truth provided by Meta-NanoSim to improve their
520 algorithms. The effect of chimeric reads, as a common source of mis-assemblies, can be easily
521 evaluated with simulated reads. Meta-NanoSim also has a perfect read simulation feature,
522 which would allow users to learn the relationship between assembly and coverage without
523 interference from read base errors.

524

525 CONCLUSIONS

526 Meta-NanoSim is an ONT metagenomic simulator that simulates complex microbial communities
527 with read features true to the platform. Given a training dataset, Meta-NanoSim generates read
528 length distributions, error profiles, and alignment ratio models by default. Optionally, it also
529 detects chimeric reads and quantifies species abundance levels. Meta-NanoSim is aimed to
530 capture the platform-specific features and can be adopted to profile datasets from any ONT

531 sequencing chemistry and basecallers tested to date. The performance of metagenomic
532 quantification of Meta-NanoSim surpasses the performance of the current state-of-the-art.
533 Meta-NanoSim is the first ONT metagenomic read simulator that can simulate chimeric reads and
534 abundance levels at base-level. Chimeric read detection improves the read length modelling and
535 helps reproduce such feature in simulated reads to challenge metagenomic assemblers,
536 taxonomy binners, and abundance quantification tools. The tool also supports multiprocessing
537 to speed up simulations and streams reference genomes from online servers to save disk space
538 when hundreds or thousands of genomes are to be simulated in a microbial community. By
539 comparing simulated reads with empirical datasets, we show that Meta-NanoSim preserves
540 some key characteristics of ONT metagenomic reads well. Further, our metagenomic assembly
541 benchmarks demonstrate a use case and utility of Meta-NanoSim. We expect Meta-NanoSim to
542 have broad utility in helping develop, test and improve upon such applications.

543

544 **AVAILABILITY AND REQUIREMENTS**

545 **Project name:** Meta-NanoSim

546 **Project home page:** <https://github.com/bcgsc/NanoSim>

547 **Operating systems:** Platform independent

548 **Programming language:** Python

549 **Other requirements:** <https://github.com/bcgsc/NanoSim/blob/master/README.md>

550 **License:** GNU General Public License

551 **Any restrictions to use by non-academics:** Please contact the authors.

552

553 **LIST OF ABBREVIATIONS**

- 554 bp : basepairs
- 555 EM : Expectation-Maximization
- 556 GB : gigabytes
- 557 GPU : graphics processing unit
- 558 M : million
- 559 NGA50 : length of the shortest alignment block for which longer or equal length alignment blocks
- 560 cover 50% of the reference genome size
- 561 nt : nucleotides
- 562 ONT : Oxford Nanopore Technologies
- 563 TB : terabytes
- 564

565 **DECLARATIONS**

566 **Ethics approval and consent to participate**

567 Not applicable

568

569 **Consent for publication**

570 Not applicable

571

572 **Availability of data and material**

573 Meta-NanoSim is implemented in Python within the NanoSim suite. The source code and pre-

574 trained models used in this study are available on Github: <https://github.com/bcgsc/NanoSim>.

575 NanoSim version 3.0.2 is used for this work. Meta-NanoSim is platform independent, and is
576 released under the GNU GPL license. The data analysed during this study is described in the
577 manuscript and supplementary methods in Additional File 1.

578

579 **Competing interests**

580 The authors declare that they have no competing interests.

581

582 **Funding**

583 This work was supported by Genome Canada and Genome BC [281ANV]; and by the National
584 Human Genome Research Institute of the National Institutes of Health [R01HG007182].
585 Scholarship funding was provided by the University of British Columbia, and the Natural Sciences
586 and Engineering Research Council of Canada. The content is solely the responsibility of the
587 authors and does not necessarily represent the official views of the funding organizations.

588

589 **Authors' contributions**

590 IB and CY conceived and designed the study. CY designed and implemented the software with
591 the help of TL, SH, and KMN. KMN and SH provided additional help with the software
592 maintainence. CY drafted the manuscript, and all authors were involved in its revision. All authors
593 read and approved the final manuscript.

594

595 **Acknowledgements**

596 Not Applicable

597 **REFERENCE**

- 598 1. Handelsman J. Metagenomics: Application of Genomics to Uncultured Microorganisms. *Microbiol*
599 *Mol Biol Rev.* 2004;
- 600 2. Chen K, Pachter L. Bioinformatics for whole-genome shotgun sequencing of microbial
601 communities. *PLoS Computational Biology.* 2005.
- 602 3. Schulz F, Alteio L, Goudeau D, Ryan EM, Yu FB, Malmstrom RR, et al. Hidden diversity of soil giant
603 viruses. *Nat Commun.* 2018;
- 604 4. Guthrie L, Gupta S, Daily J, Kelly L. Human microbiome signatures of differential colorectal cancer
605 drug metabolism. *npj Biofilms Microbiomes.* 2017;
- 606 5. Wirbel J, Pyl PT, Kartal E, Zych K, Kashani A, Milanese A, et al. Meta-analysis of fecal
607 metagenomes reveals global microbial signatures that are specific for colorectal cancer. *Nat Med.*
608 2019;
- 609 6. Quince C, Walker AW, Simpson JT, Loman NJ, Segata N. Shotgun metagenomics, from sampling to
610 analysis. *Nature Biotechnology.* 2017.
- 611 7. Brown BL, Watson M, Minot SS, Rivera MC, Franklin RB. MinIONTM nanopore sequencing of
612 environmental metagenomes: A synthetic approach. *Gigascience.* 2017;
- 613 8. Nicholls SM, Quick JC, Tang S, Loman NJ. Ultra-deep, long-read nanopore sequencing of mock
614 microbial community standards. *Gigascience.* 2019;
- 615 9. Fu S, Wang A, Au KF. A comparative evaluation of hybrid error correction methods for error-
616 prone long reads. *Genome Biol.* 2019;
- 617 10. Payne A, Holmes N, Rakyan V, Loose M. Bulkvis: A graphical viewer for Oxford nanopore bulk
618 FAST5 files. *Bioinformatics.* 2019;
- 619 11. Charalampous T, Kay GL, Richardson H, Aydin A, Baldan R, Jeanes C, et al. Nanopore
620 metagenomics enables rapid clinical diagnosis of bacterial lower respiratory infection. *Nat*
621 *Biotechnol.* 2019;
- 622 12. Kafetzopoulou LE, Pullan ST, Lemey P, Suchard MA, Ehichioya DU, Pahlmann M, et al.
623 Metagenomic sequencing at the epicenter of the Nigeria 2018 Lassa fever outbreak. *Science*
624 (80-). 2019;
- 625 13. Chan JFW, Yuan S, Kok KH, To KKW, Chu H, Yang J, et al. A familial cluster of pneumonia
626 associated with the 2019 novel coronavirus indicating person-to-person transmission: a study of
627 a family cluster. *Lancet.* 2020;
- 628 14. Greninger AL, Naccache SN, Federman S, Yu G, Mbala P, Bres V, et al. Rapid metagenomic
629 identification of viral pathogens in clinical samples by real-time nanopore sequencing analysis.
630 *Genome Med.* 2015;
- 631 15. Yang C, Chu J, Warren RL, Birol I. NanoSim: Nanopore sequence read simulator based on
632 statistical characterization. Vol. 6, *GigaScience.* 2017.
- 633 16. Hafezqorani S, Yang C, Lo T, Nip KM, Warren RL, Birol I. Trans-NanoSim characterizes and
634 simulates nanopore RNA-sequencing data. *Gigascience.* 2020;

- 635 17. Buck D, Weirather JL, de Cesare M, Wang Y, Piazza P, Sebastiano V, et al. Comprehensive
636 comparison of Pacific Biosciences and Oxford Nanopore Technologies and their applications to
637 transcriptome analysis. *F1000Research*. 2017;
- 638 18. Wood DE, Salzberg SL. Kraken: Ultrafast metagenomic sequence classification using exact
639 alignments. *Genome Biol.* 2014;
- 640 19. Lu J, Breitwieser FP, Thielen P, Salzberg SL. Bracken: Estimating species abundance in
641 metagenomics data. *PeerJ Comput Sci.* 2017;
- 642 20. White R, Pellefigues C, Ronchese F, Lamiable O, Eccles D. Investigation of chimeric reads using
643 the MinION. *F1000Research*. 2017;
- 644 21. Martin S, Leggett RM. Alvis: a tool for contig and read ALignment VISualisation and chimera
645 detection. *BMC Bioinformatics*. 2021;
- 646 22. Marijon P, Chikhi R, Varré JS. Yacrd and fpa: Upstream tools for long-read genome assembly.
647 *Bioinformatics*. 2020;
- 648 23. Xu Y, Lewandowski K, Lumley S, Pullan S, Vipond R, Carroll M, et al. Detection of viral pathogens
649 with multiplex nanopore MinION sequencing: Be careful with cross-Talk. *Front Microbiol.* 2018;
- 650 24. Tvedte ES, Gasser M, Sparklin BC, Michalski J, Hjelmen CE, Johnston JS, et al. Comparison of long-
651 read sequencing technologies in interrogating bacteria and fly genomes. *G3 Genes|Genomes|Genetics*. 2021;
- 652 25. Wick RR, Judd LM, Holt KE. Deepbinner: Demultiplexing barcoded Oxford Nanopore reads with
653 deep convolutional neural networks. *PLoS Comput Biol.* 2018;
- 654 26. Dilthey AT, Jain C, Koren S, Phillippy AM. Strain-level metagenomic assignment and compositional
655 estimation for long reads with MetaMaps. *Nat Commun.* 2019;
- 656 27. Jia B, Xuan L, Cai K, Hu Z, Ma L, Wei C. NeSSM: A Next-Generation Sequencing Simulator for
657 Metagenomics. *PLoS One.* 2013;
- 658 28. Fritz A, Hofmann P, Majda S, Dahms E, Dröge J, Fiedler J, et al. CAMISIM: Simulating
659 metagenomes and microbial communities. *Microbiome.* 2019;
- 660 29. O'Leary NA, Wright MW, Brister JR, Ciufo S, Haddad D, McVeigh R, et al. Reference sequence
661 (RefSeq) database at NCBI: Current status, taxonomic expansion, and functional annotation.
662 *Nucleic Acids Res.* 2016;
- 663 30. Howe KL, Achuthan P, Allen J, Allen J, Alvarez-Jarreta J, Ridwan Amode M, et al. Ensembl 2021.
664 *Nucleic Acids Res.* 2021;
- 665 31. Kolmogorov M, Bickhart DM, Behsaz B, Gurevich A, Rayko M, Shin SB, et al. metaFlye: scalable
666 long-read metagenome assembly using repeat graphs. *Nat Methods.* 2020;
- 667 32. Patro R, Duggal G, Love MI, Irizarry RA, Kingsford C. Salmon provides fast and bias-aware
668 quantification of transcript expression. *Nat Methods.* 2017;
- 669 33. Martin S, Heavens D, Lan Y, Horsfield S, Clark MD, Leggett RM. Nanopore adaptive sampling: a
670 tool for enrichment of low abundance species in metagenomic samples. *bioRxiv*. 2021;
- 671 34. Proctor LM, Creasy HH, Fettweis JM, Lloyd-Price J, Mahurkar A, Zhou W, et al. The Integrative
672 Human Microbiome Project. *Nature.* 2019;569(7758).

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [Supplementarysubmit.pdf](#)