1 **Inferring species compositions of complex fungal communities from long- and short-read**

2 **sequence data**

3 **Yiheng Hu[1*†], Laszlo Irinyi[2,3,4], Minh Thuy Vi Hoang[2,3,4], Tavish Eenjes[1], Abigail Graetz[1], Eric**

4 **Stone[1,5], Wieland Meyer[2,3,4,6], Benjamin Schwessinger[1*], John P. Rathjen[1*]**

5 [1] Research School of Biology, Australian National University, Canberra, ACT, Australia

6 [2] Molecular Mycology Research Laboratory, Centre for Infectious Diseases and Microbiology, Faculty

7 of Medicine and Health, Sydney Medical School, Westmead Clinical School, The University of Sydney,

8 Sydney, NSW, Australia.

9 [3] Marie Bashir Institute for Infectious Diseases and Biosecurity, The University of Sydney, Sydney,

10 NSW, Australia.

11 [4] Westmead Institute for Medical Research, Westmead, NSW Australia.

12 [5] ANU-CSIRO Centre for Genomics, Metabolomics and Bioinformatics, Canberra, ACT, Australia

13 [6] Westmead Hospital (Research and Education Network), Westmead, NSW, Australia.

14

15 * Correspondence: yiheng.hu@anu.edu.au; benjamin.schwessinger@anu.edu.au;

16 john.rathjen@anu.edu.au

17 † Present address: Department of Microbial Interactions, IMIT/ZMBP, University of Tübingen,

18 Tübingen, Germany

19

20

21

22    **Abstract**

23    Background:

24    The kingdom fungi is crucial for life on earth and is highly diverse. Yet fungi are challenging to

25    characterize. They can be difficult to culture and may be morphologically indistinct in culture.

26    They can have complex genomes of over 1 Gb in size and are still underrepresented in whole

27    genome sequence databases. Overall their description and analysis lags far behind other

28    microbes such as bacteria. At the same time, classification of species via high throughput

29    sequencing without prior purification is increasingly becoming the norm for pathogen

30    detection, microbiome studies, and environmental monitoring. However, standardized

31    procedures for characterizing unknown fungi from complex sequencing data have not yet

32    been established.

33    Results:

34    We compared different metagenomics sequencing and analysis strategies for the

35    identification of fungal species. Using two fungal mock communities of 44 phylogenetically

36    diverse species, we compared species classification and community composition analysis

37    pipelines using shotgun metagenomics and amplicon sequencing data generated from both

38    short and long read sequencing technologies. We show that regardless of the sequencing

39    methodology used, the highest accuracy of species identification was achieved by sequence

40    alignment against a fungi-specific database. During the assessment of classification

41    algorithms, we found that applying cut-offs to the query coverage of each read or contig

42    significantly improved the classification accuracy and community composition analysis

43    without significant data loss.

44    Conclusion:

45    Overall, our study expands the toolkit for identifying fungi by improving sequence-based

46    fungal classification, and provides a practical guide for the design of metagenomics analyses.

47

48    **Introduction**

49    Fungi are ubiquitous yet their presence and impact are often overlooked. It has been

50    estimated that 2.2-3.8 million species inhabit planet earth [1] but only about 4% of these are

51    catalogued [2]. Mora *et al*. estimated that there are 7.8 million and 298,0000 animal and

52    plants species on earth with 12.3% and 72.4% of these characterised scientifically,

53    respectively [3], which points towards a more central role in cultural awareness. In contrast,

54    fungi are introduced to our consciousness via a brief mention in high school textbooks, or as

55    largely side subjects in botany and microbiology courses at university [4,5]. Fungi play diverse

56    roles throughout evolution and are particularly active in mediating the breakdown and uptake

57    of nutrients. They constitute a major disease load to humans, causing millions of deaths per

58    year, and wreak devastating crop losses via a constant toll of disease and epidemics and are

59    an existential threat to many frog species [6,7]. On the other hand, fungi are or are used to

60    manufacture delicious foods and beverages, and have saved countless lives via antibiotic

61    production [8,9]. Therefore, a recent call was made to expand fungal research and improve

62    our awareness of this special kingdom [10].

63    To progress our understanding of fungal biology we need to be able to classify more species

64    more precisely. Fungi have been an independent kingdom since 1969 [11] with addition of

65    further phyla in early 2000 [12–16]. Historically, its taxonomy was based on morphological

66  and reproductive traits but this has been surpassed by DNA-based classification which

67  revolutionized mycology, not only refining the conventional taxonomic tree [17,18] but also

68  standardizing the identification of new species. In the absence of whole genome data, DNA-

69  based classification primarily exploits the internal transcribed spacer (ITS) within the

70  ribosomal RNA genes as a highly polymorphic marker to distinguish species. It is easily

71  amplified and sequenced due to highly conserved flanking sequences and contains a high

72  degree of variation between even closely related species. Although a mature pipeline

73  comprising ITS amplification, IIllumina sequencing and data analysis has been established[19],

74  several studies reported biases from the sequencing technology used and from unevenly

75  amplified fungal marker regions [20–22]. Recently, novel strategies exploiting long-range

76  amplification and long-read sequencing have been developed to improve these classifications

77  [23,24]. In addition, whole genome shotgun sequencing and rapidly expanding genome

78  databases allow mapping of newly generated DNA sequences directly to the database. This

79  strategy allows exploitation of genetic variation throughout the genome and abandonment

80  of the marker gene amplification step, which increases classification accuracy and reduces

81  the biases from the estimation of relative abundance [25].

82  Although advanced sequencing methods allow novel strategies for fungal identification

83  particularly from mixed samples, new demands are placed on data analysis pipelines to

84  improve the accuracy of fungal classification. Different algorithms have been developed to

85  classify DNA sequences at distinct taxonomic ranks based on sequence databases with

86  taxonomic information [26–30]. For example, alignment algorithms such as Basic Local

87  Alignment Search Tool (BLAST) [27] detect matches of each sequence to subjects of the target

88  database along with the taxonomic information assigned to each entry. Alternatively,

89  sequence features represented by short unique subsequences named k-mers can be derived

90    from sequence data and mapped to databases to identify taxa with the highest number of

91    cross-mapping k-mers[28]. Several studies have critically assessed algorithms for species

92    classification on simulated datasets or bacterial community datasets [31–33], but

93    comparisons of sequencing strategies for complex fungal communities alignment using real

94    data and different identification pipelines are extremely rare. In addition to search algorithms,

95    the choice of database also influences classifications dramatically, but only a few studies have

96    researched their impact [34–36]. Therefore, more comprehensive benchmarking of both

97    classification algorithms and databases are needed to optimise identification pipelines.

98    Here, we assessed different combinations of algorithms and databases during processing of

99    both short- and long-read sequencing data for the identification of taxa from complex mock

100   fungal communities. We identified key factors that influence the accuracy of classifications,

101   both for mock community datasets and public datasets. Optimisation of these methods also

102   lead to more accurate community composition analysis. Our results provide guidelines for the

103   design of sequence-based community analysis for fungal species.

104

105   **Results**

106   **Construction of mock fungal community datasets**

107   We constructed two mock communities from the same set of 44 fungal species

108   (Supplementary Table S1). Most of these are human-associated pathogenic yeasts while some

109   are basidiomycete pathogens. One community comprised pooled DNA (PD) from each species

110   and the second was composed of DNA extracted from equal quantities of fungal biomass (PB)

111   of each species that were mixed together prior to extraction. We generated four sequence

112    datasets for each community using Illumina and nanopore technologies, sequencing both

113    shotgun metagenomes and targeted amplicons respectively. The data derived from each

114    strategy are summarized in Table 1.

115

116    **Alignment algorithm against a specific fungal database resulted in the most accurate fungal**

117    **classifications**

118    We compared different analysis strategies for each shotgun dataset. For nanopore datasets,

119    we directly used the quality-controlled reads for classification. For Illumina data, we quality

120    filtered all reads and assembled them into contigs before classification to maximize the

121    classification accuracy. We performed both alignment and k-mer based classifications on

122    these data using BLAST and Kraken2 [27,29] using a 'winner-takes-all' strategy in which the

123    top hit was taken as the identity of the query sequence. For each algorithm, we compared the

124    use of two reference databases: the non-redundant NCBI nucleotide database (nt) [37] and

125    the RefSeq fungi database (RFD) [38] which only contains curated fungal genomes. We first

126    assessed the performance of each alignment tool on both databases for each data input. We

127    compared the concordance in the results of each pipeline at the genus level. We define

128    concordance as the percentage of fungal genera identified by both analyses in a pairwise

129    comparison (Figure 1A). The concordance between analyses on each dataset varied between

130    69% and 86% and generally, Illumina data resulted in higher concordance than did nanopore

131    data.

132    We then aimed to identify the combination of algorithm and database that yielded the most

133    accurate species identification. We used classified proportion and precision to evaluate each

134    classification, where $\text{Classified Proportion} = \frac{\text{\# total basepairs classified.}}{\text{\# total basepairs of input reads}}$, and

135    $\text{Precision} = \frac{\text{\# total basepairs classified correctly}}{\text{\# total basepairs classified}}$ .

136    The number of total basepairs is calculated as total read length for nanopore reads and total

137    coverage of Illumina reads to each contig [32,33]. We plotted the precision and classified

138    proportion for each pipeline and found three regular patterns (Figures 1B and 1C): First, for

139    each dataset, BLAST resulted in higher precision but lower classification proportion by

140    comparison to Kraken2. Second, Illumina contigs returned higher classification proportion

141    and precision than nanopore reads. Third, classification against the RFD database yielded

142    higher precision than those against the nt database. In summary, BLAST alignments against

143    the RFD database yielded the best classification strategy.

144

145    **Applying cut-offs to query coverage improves classification accuracy on shotgun**

146    **metagenomics datasets**

147    We next aimed to improve our classification scheme by filtering the BLAST search results. We

148    reasoned that restricting alignment metrics would reduce the number of false classifications.

149    To investigate changes in classification accuracy after restricting BLAST output parameters,

150    we first BLASTed shotgun metagenomics reads against the RFD database without applying

151    any filter, then applied progressive cut-offs on different parameters of the BLAST results. We

152    evaluated changes in the results based on the metrics precision, remaining rate and

153    completeness. Precision is described above and estimates the accuracy of the classification;

154    remaining rate captures the percentage of the input data remaining after the application of

155    each cut-off; and completeness is the number of taxa captured relative to the total number

156    of taxa within the mock community. We initially applied cut-offs on query length and two

157    alignment metrics; E-value - the number of expected hits of similar quality that could be found

158    by chance alone; and pident – the percentage of identical matches within the region of

159    alignment between query and subject. As shown in Figure 2A, applying progressive cut-offs

160    to query length did not improve the precision, whilst both completeness and remaining rate

161    diminished dramatically from very small cut-off values. Cut-offs applied to alignment E-values

162    removed <20% of the BLAST results, whereas precision showed minor improvement,

163    especially on nanopore datasets (Figure 2B). For Illumina data, applying cut-offs to the E-value

164    increased the precision by around 2% but at the cost of diminished completeness. E-value

165    cut-offs performed better on nanopore datasets, improving precision by 3% (PD) or 4% (PB)

166    with non-identification of only a single genus from the mock community, at $10^{-250}$ or almost

167    $10^{-400}$ respectively. Progressive cut-offs on pident yielded the best results of all three filters.

168    For Illumina data, precision was improved by up to 8% for PB data, and completeness

169    remained at 100% in almost all cases (Figure 2C). For nanopore datasets, pident cut-offs

170    improved the precision by up to ~3% before sharp decreases, with a concurrent filtering of

171    ~60% BLAST result as shown by the remaining rate. Given the characteristically high error rate

172    of nanopore reads, we also applied cut-offs on quality scores to these data. Cut-offs applied

173    to Phred scores did not alter the precision, while a significant proportion of the dataset was

174    lost through filtering (Supplementary Figure S1). Overall, our results suggest that applying

175    each filter to BLAST results performs well on either Illumina or nanopore data but not both,

176    and that cut-offs based on query length or quality scores did not affect the precision

177    significantly.

178    Given the results above, we investigated how the alignment parameters were calculated and

179    explored other variables to improve the classifications. The BLAST E-value is calculated as E =

180    $mn2^{-S}$ in which S is the bits score derived from the number of gaps and mismatches in the

181    alignment, and m and n are the query length and database total length respectively [39].

182    Therefore, the E-value is influenced exponentially by the alignment quality. We next

183    investigated query coverage, a metric based on how much of the query sequence aligned to

184    the subject. We calculated the query coverage as the number of identical matches divided by

185    the read or contig length, and applied progressive cut-offs on this parameter for each

186    dataset/algorithm analysis. As shown in Figure 2D, applying cut-offs on query coverage

187    improved the precision of all four analyses significantly, and did not cause losses of

188    completeness at smaller cut-off values. For example, at a 10% cut-off on query coverage, the

189    precision of all four analyses was 98-99% while the completeness remained at 100% and the

190    removed BLAST results ranged from 10-25%. This result not only supported our hypothesis

191    that the total length of the alignment matters as much as the alignment quality, but also

192    suggested a novel approach to improve the accuracy of fungal classification.

193

194    **Improving taxa identification from published metagenomics datasets using query coverage**

195    **as a filtering parameter**

196    After improving classifications by applying cut-offs to the query coverage on the mock

197    community datasets, we extended this strategy to try to improve the classification of

198    published shotgun metagenomics datasets. We re-analysed ten nanopore and six Illumina

199    shotgun metagenomics datasets [40–43]. These included host-associated fungal samples

200    (nanopore) and host-depleted microbiome data (Illumina). Since the environmental datasets

201    contain unknown species, we followed the concept of classification precision. We calculated

202    the percentage of the dataset that was classified into taxa known to be included in the sample.

203    For example, in re-analysing human clinical samples [42], we included the pathogen

204    (*Pneumocystis*) and the host human (*Homo*) as the true taxa, and calculated the total

205    proportion of query sequences classified to these taxa before and after applying cut-offs on

206    query coverage. Table 2 shows the improvement in taxonomic classification from the

207    published datasets after applying query coverage cut-offs. We initially applied a 20% cut-off

208    on the query coverage for all analyses, but the data loss in most cases was too high. Therefore,

209    we applied query cut-offs that filtered around 20% of the blast result based on our results

210    from the mock fungal community datasets (Figure 2D).

211    For all Illumina datasets, we downloaded the quality-controlled sequences and re-analysed

212    them using the assembly and BLAST pipeline described above against the NCBI nt database.

213    For the nanopore human datasets [42], we used the BLAST results taken directly from the

214    original articles for analysis. For the infected wheat datasets [41], we downloaded the

215    sequences and re-analysed them against the RefSeq fungal database. The precision increased

216    for nearly all datasets after applying query coverage cut-offs (Table 2). For the Illumina

217    microbiome datasets, we first assessed the change of proportions in fungal taxa after applying

218    cut-offs on query coverages using the species lists identified by Donovan *et al.* [44] as

219    confirmed taxa. We observed only marginal increase in percentages for the confirmed fungal

220    communities, due to their low total proportions in the original samples. We then calculated

221    the improvement in precision for the bacterial communities. The Illumina datasets were

222    generated from swine and mouse gut microbiome samples, so we assessed the change in

223    proportions of their core bacterial genera (a group of bacteria commonly present in swine

224    and mouse guts [45,46]). The percentages of confirmed core bacterial genera improved by up

225    to 5.7% after applying cut-offs on query coverage (Table 2). In addition, in the nanopore

226    human datasets, the total percentage of reads classified as *Homo* in the three healthy

227    individual samples were improved by applying cut-offs to query coverage. These results

228    indicated that this strategy may be broadly applicable not only to fungal species, but also to

229    the classification of other eukaryotes and bacteria. One Illumina dataset (d1) and one

230    nanopore dataset (a5) showed decreased percentages of confirmed taxa after applying query

231    coverage cut-offs, which might be because the core microbiome species are not representing

232    the species identified in the Illumina sample, or due to the low coverage and high error rate

233    of nanopore data.

234

235    **Benchmarking classification pipelines for amplicon datasets identified advantages of each**

236    **strategy**

237    We next assessed different strategies for the classification of ITS amplicon datasets. We

238    amplified the ITS region from both mock communities using two different primer pairs and

239    three technical replicates for each sample. Taking advantage of nanopore technology, we

240    performed long-amplicon sequencing of a roughly 3 kb ribosomal RNA gene region covering

241    part of the 28S subunit, ITS1, 5.8S subunit, ITS2 and part of the 18S subunit [19]. For Illumina

242    sequencing we used the well-established ITS1F-ITS2 amplicon of about 300 bp in length [47].

243    Similar to the analysis of the shotgun datasets, we applied both k-mer and alignment-based

244    approaches to the classification of nanopore amplicon data. We used the pair-wise alignment

245    algorithm minimap2 as the alignment algorithm instead of BLAST due to its speed and

246    efficiency. We tested four different databases for classification of long amplicons; the NCBI

247    18S and 28S databases, and two ITS databases from NCBI and UNITE, respectively [38,48].

248    Overall, we found that the k-mer algorithm returned much higher classification proportion

249    than alignment for each nanopore dataset, but the highest precision (~97%) were achieved

250    by combining the minimap2 alignment algorithm with the NCBI ITS database (Figure 3A). For

251    Illumina amplicon datasets, we applied the QIIME2 pipeline which is one of the most widely

252    used strategies for ITS classification and community composition analysis[49]. The QIIME2

253    pipeline groups similar Illumina amplicons into sequence features before classification to

254    reduce the demand on computational resources [50]. Since all individual Illumina reads are

255    grouped into sequence features and all the sequence features are classified, the classification

256    proportion of the Illumina amplicon datasets are 100%. We plotted precision rates from the

257    QIIME2 analysis of both the PD and PB samples with their means (Figure 3B). The mean

258    precision from either Illumina dataset were lower than that from k-mer analysis of the

259    respective nanopore datasets.

260    Although the precision from the amplicon datasets were higher than that from shotgun

261    datasets, the ITS classification did not identify all genera within the mock community, as

262    shown by our completeness analysis (Figure 3C). The nanopore amplicons identified 68% (PD)

263    and 63% (PB) of the total genera in the mock community, whereas the Illumina amplicon

264    datasets covered only 25% and 41% of the genera respectively. We suspect that the low

265    completeness from ITS classifications was due partially to the low quality of this particular

266    dataset (Table 1) and partially due to non-uniform amplification from the different primer

267    pairs. However, there were fewer nanopore amplicon reads than in the Illumina amplicon

268    datasets and the completeness from the nanopore data was higher (Figure 3C). This supports

269    the argument that long amplicons identify a wider range of species and are more accurate in

270    species classification than short amplicons [51,52].

271

272    **Cut-offs on query coverage also improve community composition analysis**

273     We next analysed community compositions using the most accurate classification method for

274     each dataset. Community composition refers to the identity and relative abundances of all

275     taxa in a community. Given the observation that use of a restricted database resulted in

276     higher classification precision, we constructed a database containing only the genomes from

277     species within the mock community and aligned all of the data to the mock community

278     database using BLAST. This forces the precision to 100% as any classification will belong to a

279     species from the mock community. We then BLASTed each dataset against this database and

280     calculated the relative abundance of each genus. We defined this as the 'gold standard' for

281     community composition analysis of the mock fungal community (Figure 4A). We then

282     compared the community composition determined from each combination of algorithms and

283     databases with the gold standard for each dataset, and measured their differences using

284     three statistical distance tests: Bhattacharyya distance, relative Euclidean distance and

285     relative entropy [53–55]. Consistently, BLASTing sequences against the RFD database

286     produced community compositions with the highest similarity to the gold standard analysis

287     (Figure 4B).

288     To assess whether query coverage cut-offs also improved the community composition

289     analysis of shotgun metagenomics data, we plotted the changes in statistical distance after

290     progressive application of query coverage cut-offs (Figure 4). After applying cut-offs on the

291     query coverage, the community composition improved in all cases especially for lower cut-off

292     values. The community compositions from PB-Illumina datasets improved and turned out to

293     be the most similar to the gold standard at query-coverage cut-offs greater than 3 - 4%, which

294     is consistent with the changes in precision rate shown in Figure 2D. Overall, our results

295     illustrated that applying cut-offs on query coverage did not only improve the classification

296     accuracy, but also the community composition analysis.

297

**Discussion:**

299 Here we investigated the taxonomic classification from sequencing data, one of the key steps

300 in all metagenomic workflows, with a particular focus on fungi. After assessing various

301 combinations of algorithms and databases following different sequencing strategies, we

302 found that the combination of BLAST with the specific RFD database always resulted in the

303 most precise classification for all mock fungal community datasets. These classifications were

304 further improved when applying cut-offs on query coverage including positive flow on effects

305 on downstream community composition analysis from shotgun metagenomics datasets.

306 Despite that sampling and DNA extraction substantially influence the outcome of species

307 classifications [56–58], choosing an appropriate sequencing strategy is the primary step

308 towards accurately profiling a sample. For shotgun datasets, our results suggested that both

309 short and long shotgun datasets have comparable accuracy and both higher than the

310 amplicon datasets. However, Illumina shotgun datasets require additional steps to assemble

311 reads into contigs before querying them against a database, and to map all reads back to the

312 assembly to quantify the coverage. These processes are necessary to achieve accurate

313 classification from longer contigs [59], but result in a longer sequence-to-result turnover than

314 the long read shotgun data. In the analysis of the amplicon data, long range amplicons

315 performs better in the classification accuracy and completeness compare to the short ITS data,

316 consistent with other studies [51,52]. Comparing to the results from shotgun datasets, the

317 overall completeness from the result of amplicon datasets is much lower. We think that is

318 because we used much less amplicon data for benchmarking classification pipelines, and the

319     incomplete database which do not contain all taxa present in our mock community. Overall,

320     the long read shotgun datasets returned the most accurate fungal classification.

321     Next, our data supported that alignment algorithm (BLAST) outperform the k-mer based

322     approach (kraken2) in the accuracy of classification [32,60], and also compared progressive

323     cut-offs to major alignment parameters for shotgun metagenomics data. We found that

324     applying read length or read quality cut-offs did not improve the precision of the classification

325     for all shotgun datasets. This observation is different with the previous study based on

326     simulated data, which claimed that the long reads improves the accuracy of classification [60].

327     Cut-offs on pident slightly improved the classification accuracy for illumina datasets, but the

328     error-prone nature of the nanopore data (~10% error rate) is also reflected in the result, as it

329     causes the breakdown of precision when pident cut-offs reach 90% (Figure 2C).

330     We found that query coverage cut-off that filter out 20% blast result worked best. Unlike the

331     E-value weighing the gaps and mismatch as the major factor effecting alignment quality, the

332     query coverage weighs the query length as well as the number of identical matches in the

333     assessment of the alignment quality. In this case, we can eliminate more spurious alignments

334     that are due to a small proportion of reads with high fidelity to the reference, which are

335     commonly found in reads containing conserved genes and repeated sequences. Interestingly,

336     to reach the same 20% filtering threshold, we set up higher cut-offs on the query coverage

337     (10 -20%) in mock community datasets than the real environmental datasets, including few

338     extremely low thresholds of query coverage in the Illumina shotgun datasets. We compared

339     other studies that use simulated data to generate metagenomics contigs for classification,

340     and found that they used 90% query coverage cut-offs as the parameter[62–64]. Together

341     with the different result of read length and read quality cut-offs, this observations highlighted

342    the difference between the use of real environmental data and the simulated data in

343    benchmarking studies, especially for the classification of complex microbial communities.

344    PD and PB samples showed slightly different results in comparing statistical distance with the

345    gold standard. After applying cut-offs on query coverage, both Bhattacharyya distance and

346    Euclidean distance between the best practice and the gold standard classification only

347    showed marginal decrease in PD samples, and slowly reversed as the cut-offs increase. We

348    think that is because about 1/3 of reads were classified as *Candida* in the pooled DNA sample,

349    so the difference on the relative abundance of one *Candida* genus between the gold standard

350    community composition and the best practice is much higher and much more influential to

351    the final distance than that from other genera.

352    Following the importance of the alignment quantity represented by the query coverage, the

353    next question is, how to bring the low quality but high quantity alignment into consideration?

354    Therefore, the winner-takes-all selection strategy itself can be re-designed, as the highly

355    conserved genome regions from different species generate highly close alignment scores

356    between the best alignment and other top alignments. In this case, a weighing statistics and

357    the relative probability for multiple top taxonomic assignments can be explored and

358    introduced to replace the best-hit-takes-all strategy. This will be particularly useful in

359    connection with the rapid expansion of the fungal genome databases.

360    Next to the right classification tool, chosen the appropriate database significantly influences

361    analysis outcomes [33,34]. Based on our observation, we suggest that 'prior knowledge'

362    about the dataset should guide the choice of the appropriate database as this will improve

363    the accuracy of taxonomic classifications. For example, our results suggested that the

364    restricted database resulted in more accurate fungal classifications for shotgun

365    metagenomics datasets. This strategy might be appropriate if queries are initial binned into

366    kingdoms before a more in-depth analysis with kingdom specific databases. Also, Kaehler *et al.* [65] incorporated environment-specific taxonomic abundance information into the

368    analysis of amplicon datasets and showed that these improve classification accuracy. Similar

369    approaches can be applied to metagenomic datasets. In addition, machine learning strategies

370    become increasingly popular for analysing genomic data. Here taxonomic classifiers could be

371    trained on existing labelled sequence datasets before being applied to communities with

372    similar composition to the training datasets or to identify target species from complex

373    communities [60,66].

374

**Conclusion**

376    In this study, we perform an in-depth analysis on how different sequencing strategies,

377    classification algorithms and databases impact fungal classifications using complex real-life

378    mock community sequencing datasets. We find that alignment algorithm (BLAST) with

379    targeted fungal database (RFD) achieve the best classification accuracy and community

380    composition estimates. These can be further improved by applying cut-offs on query coverage.

381    Taken together, the findings from our benchmarking workflows have important implications

382    for mycology studies for multiple stages of metagenomics analysis, and provided a guide to

383    other researchers aiming to study fungal metagenomics.

384

**Methods**

**Code availability**

387 All detailed commands and scripts used in each step were summarized in

388 https://github.com/Yiheng323/Benchmarking-taxonomic-classification-strategies-using-

389 mock-fungal-communities.

390 **Fungal harvesting, DNA extraction and construction of mock communities**

391 Selected fungal strains were cultured onto Sabouraud dextrose agar and incubated for 48

392 hours at 27ºC.

393 For the species in the PD community, an inoculating loop full of fungal cells were scraped into

394 a 1.5 mL microfuge tube and crushed with a pestle and liquid nitrogen. Genomic DNA was

395 then extracted using the Zymo Research *Quick*-DNA Fungal/Bacterial Miniprep Kit (cat. no.

396 D6005 Zymo Research, Irvine, CA, USA). First, BashingBead$^{TM}$ Buffer was added to the crushed

397 fungal cells and vortexed. The mixture was then filtered through a Zymo-SpinTM III-F Column

398 and the filtrate was combined with Genomic Lysis Buffer. The mixture was filtered through a

399 Zymo-Spin$^{TM}$ IICR Column and washed with DNA Pre-Wash buffer and g-DNA Wash Buffer.

400 The DNA was eluted in nuclease free water. DNA concentration was measured using the

401 DeNovix dsDNA Broad Range Kit (DeNovix, Wilmington, DE, USA) and 250 ng of DNA from

402 each strain were then pooled together.

403 For the PB community, two inoculating loops of fungi of each species in teg mock community

404 were scraped into a ceramic mortar. Liquid nitrogen was then poured into the mortar and the

405 fungal mixture was crushed into a fine powder. DNA was then extracted using the Qiagen

406 DNeasy PowerMax Soil Kit (cat. no. 12988-10 Qiagen, Hilden, Germany). PowerBead Solution

407 and Solution C1 were added to the crushed fungal community, vortexed and centrifuged. The

408 supernatant was then added to Solution C2, mixed and centrifuged, which was then repeated

409 with Solution C3. The resulting supernatant was combined with Solution C4 and centrifuged

410    through a column. The column was then washed twice with Solution C5. Final DNA was eluted

411    in nuclease free water and the concentration measured using the DeNovix dsDNA Broad

412    Range Kit.

413    **Library preparation and sequencing**

414    The ITS1 regions of the rRNA gene were amplified with the universal fungal primers, ITS1F

415    (CTTGGTCATTTAGAGGAAGTAA) and ITS2 (GCTGCGTTCTTCATCGATGC)[47]. Sequencing of

416    PCR amplicons was conducted with MiSeq® System of Illumina (Illumina, San Diego, CA, USA)

417    by the Australian Genome Research Facility. The Illumina bcl2fastq 2.18.0.12 pipeline was

418    used to generate the sequence data. Pair-ends reads 2 × 300bp were generated up to 0.15

419    GB per sample for amplicon data. The Illumina amplicon data are then directly imported into

420    QIIME2 for analysis. For shotgun Illumina datasets, we employed the same sequencing

421    pipeline as the amplicon data, with MiSeq® and bcl2fastq 2.18.0.12 pipeline for the 2 x 300bp

422    paired end reads. Raw shotgun Illumina reads were trimmed adapters with Trimomatic [67].

423    Quality controlled, paired end reads were merged and assembled to metagenomics contigs

424    using IDBA_UD [68], which is more suitable for datasets with uneven sequencing depths of

425    each species. After assembly, raw reads were mapped back to the contigs using bwa-mem

426    [69], and the bam files were generated and sorted from sam files using samtools [70].

427    Bedtools [71] was used for generating coverage for each contig, and we used python numpy

428    and pandas module to calculate the average coverage for each contig.

429    For Nanopore sequencing of both shotgun and amplicon sequencing, we used Ligation

430    Sequencing 1D SQK-LSK108 and Native Barcoding Expansion (PCR-free) EXP-NBD103 Kits from

431    ONT (UK), as adapted by Hu and Schwessinger [72], which was adapted from the

432    manufacturer's instructions with the omission of DNA fragmentation and DNA repair. DNA

433  was first cleaned up using a 1× volume of Agencourt AMPure XP beads (cat. no. A63881,

434  Beckman Coulter, Indianapolis, IN, USA) following manufacturer's instructions. We then

435  eluted the beads binded DNA in 51 µl nuclease free water and quantified using NanoDrop®

436  and Quibit™ Fluorometer (Thermo Fisher Scientific, Waltham, MA, USA). DNA was end-

437  repaired (NEBNext Ultra II End-Repair/dA-tailing Module, cat. No. E7546), 1x volume beads

438  cleaned (AMPure XP beads) and eluted in 31 µl nuclease free water. Barcoding reaction was

439  performed by adding 2 µl of each native barcode and 20 µl NEB Blunt/TA Master Mix (cat. No.

440  M0367) into 18 µl DNA, mixing gently and incubating at room temperature for 10 minutes. A

441  1× volume (40 µl) Agencourt AMPure XP clean-up was then performed and the DNA was

442  eluted in 15 µl nuclease free water. Ligation was then performed by adding 20 µl Barcode

443  Adapter Mix (EXP-NBD103 Native Barcoding Expansion Kit, ONT, UK), 20 µl NEBNext Quick

444  Ligation Reaction Buffer, and Quick T4 DNA Ligase (cat. No. E6056) to the 50 µl pooled

445  equimolar barcoded DNA, mixing gently and incubating at room temperature for 10 minutes.

446  The adapter-ligated DNA was cleaned-up by adding a 0.4× volume (40 µl) of Agencourt

447  AMPure XP beads, incubating for 5 minutes at room temperature and resuspending the pellet

448  twice in 140 µl ABB provided in the SQK-LSK108 kit. The purified-ligated DNA was

449  resuspended by adding 15 µl ELB provided in the SQK-LSK108 kit and resuspending the beads.

450  The beads were pelleted again and the supernatant sequencing library was transferred to a

451  new 0.5 ml DNA LoBind tube (Eppendorf, Germany). Nanopore sequencing was carried out

452  by MinION MK1b device using R9.4.1 Flowcells. Raw fast5 files are barcode demultiplexed by

453  deepbiner (ONT), then basecalled by Guppy (v3.6.0, ONT, UK). Quality passed reads in fastq

454  files were trimmed adapters and barcodes using qcat (ONT, UK). For the long amplicon data,

455  we filtered out reads less than 2000 base pairs. All sequencing data was submitted to NCBI

456  Short Read Archive (SRA) under the bioproject PRJNA725368 including eight accessions:

457    SRX10705648, SRX10705649, SRX10705650, SRX10705651, SRX10705695, SRX10705696 and

458    SRX10705697.

459    **Genome assembly**

460    While generating the reference genome database, we found that there were no reference

461    genomes for *Candida rugosa*, *Candida mesorugosa* and *Cryptococcus magnus*, so we

462    performed nanopore sequencing on pure DNA from each species and assembled their draft

463    genomes. These assemblies were of sufficient contiguity and quality (Supplementary Table

464    S2), so we added the new draft genomes into the reference database.

465    The nanopore data of *Candida rugosa*, *Candida mesorugosa* and *Cryptococcus magnus* was

466    generated individually using Ligation Sequencing 1D SQK-LSK108 kit alone, and from

467    independent flowcells. Data from each flowcell was basecalled and quality filtered using the

468    same pipeline as described above. We got roughly 40X coverage for *Candida rugosa* and

469    *Candida mesorugosa*, and 20X coverage for *Cryptococcus magnus.* Draft genomes were

470    assembled with Flye [73] using default parameters and an estimated genome size of 20Mb.

471    After assembly, the contigs were polished ten times with Racon [74] using nanopore reads,

472    followed by one polishing with Medaka (ONT). Polished assembly was assessed completeness

473    using BUSCO [75]. The assembly statistics were reported from Flye.

474    **Database constructions**

475    For shotgun metagenomics analysis, we used three BLAST database and three kraken

476    databases. Two databases (nt and RFD) are from the same NCBI source, downloaded in May

477    2019. BLAST and kraken2 nt databases were downloaded using the updateblastdb.pl script

478    from BLAST+ package[76] and the kraken2 program [29], respectively. The fasta files of RefSeq

479   fungal database was downloaded from the NCBI and converted to BLAST database using

480   makblastdb command from the BLAST+ package[76], and was added to the kraken2 database

481   library using kraken2 command [29]. We also build the standard kraken2 database for

482   masking the contaminated regions within the fungal genomes using kraken2 command [29].

483   To generate the mock community database with only the species from the mock community,

484   we downloaded the genomes of all species in the mock community from the NCBI according

485   to their accessions (Supplementary Table S1), and concatenated them with the three newly

486   assembled genomes of *Candida rugosa, Candida mesorugosa* and *Cryotococcus magnus.*

487   Following the previous pipeline [77], we then performed a kraken2 search to identify the

488   potential contaminated regions in the concatenated fasta, and masked those regions using

489   bedtools [71]. We also masked the low complexity regions using the dustmasker from BLAST+

490   package [76]. To enable new genomes to be indexed by blastn, we updated the taxonomic

491   map file by adding the fasta headers of the three new genomes and manually assigned their

492   taxonomic ID in the file. Lastly, we used the makeblastdb program to construct the mock

493   community database.

494   For amplicon data analysis, we used two versions of fungal ITS database from the NCBI and

495   UNITE, plus the fungal 18S, 28S database from the NCBI. All of them are downloaded as fasta

496   format in February 2020 and added to the kraken2 database library using kraken2 command

497   [29].

498   **Data analysis**

499   For Shotgun metagenomics datasets, we first used blastn (version 2.10.1) and kraken2

500   (version 2.0.8) to assign the NCBI taxonomic ID for each Illumina contig or Nanopore read.

501   During the classification, we found one contamination species *Purpureocillium lilacinum*

502    always present in all samples with a significant abundance (10-20%). Therefore, we added this

503    species into the true species list. The best hit from BLAST or species with the highest k-mer

504    counts for each read and/or contig was retained for further analysis. After classification, we

505    used python pandas module to merge information from different output files, and used ete3

506    module [78] to assign taxonomic information to each read or contigs. The relative abundance

507    of each classification were calculated based on the total length of Nanopore reads of total

508    coverage of Illumina contigs. We used python numpy and math module for all statistical

509    analysis.

510    For amplicon datasets, we sequenced each sample with three technical replicates. The

511    classification workflow was different for datasets with different sequencing technologies. We

512    only used QIME2 workflow plus the UNITE database for the Illumina amplicon data, since it is

513    the only widely used method for classification. The paired end reads were denoised using the

514    DADA2[79] plugin and assigned taxonomic information using the q2-feature-classifier [80]

515    plugin. The QIME2 classifier was trained by the database sequence before classification. The

516    classification output .qzv files were visualized by the QIME2 view website

517    (https://view.qiime2.org/) and the feature-frequency csv file was extracted from the website.

518    We then used python numpy and math module for the mathematical analysis and used

519    seaborn module for generating figures.

520    For nanopore amplicon datasets, we used kraken2 as the k-mer based algorithm and

521    minimap2 as the alignment based algorithm. The kraken2 command is the same as the

522    kraken2 analysis for the shotgun metagenomics datasets, only using different databases. For

523    the minimap2 analysis, we extracted the accessions of the best hits from the output files, and

524    searched their corresponding taxonomic ID from the NCBI taxonomic map (downloaded from

525    https://ftp.ncbi.nih.gov/pub/taxonomy/accession2taxid/nucl_wgs.accession2taxid.gz,    in

526    June 2020) using python pandas module. We then merge information from different output

527    files, and used ete3 module again to assign taxonomic information to each read.

528

### Declarations

529    **Declarations**

530    **Ethics approval and consent to participate**

531    Not applicable.

532    **Consent for publication**

533    Not applicable.

534    **Data Availability**

535    All sequencing data was submitted to NCBI Short Read Archive (SRA) under the BioProject

536    PRJNA725368 including eight accessions: SRX10705648, SRX10705649, SRX10705650,

537    SRX10705651, SRX10705695, SRX10705696 and SRX10705697.

538    **Competing interests**

539    The authors declare that they have no competing interests.

540    **Authors' contributions**

541    WM, ES, BS and JPR conceived the study and designed experiments. YH, LI and WTVH

542    prepared the samples and generated sequencing data. YH, TE and AG performed the

543    bioinformatics analysis. ES provided feedback on statistical analysis. All authors contributed

544    to data analysis and manuscript writing. All authors read and approved the final manuscript.

555

556    **Figure Legends**

557    **Table 1.** The characteristics for each dataset.

558    **Table 2.** Assignment of published sequence data to genera after application of cut-offs to

559    query coverage.

560    **Figure 1.** Analysis of shotgun metagenomics data. (A) Swarmplot showing the concordance in

561    genus identification after varying either the alignment algorithm or querying different

562    databases on different data inputs. nt = NCBI nucleotide database; RFD = RefSeq Fungi

563    database; data inputs are indicated below the line (PD = pooled data; PB = pooled biomass);

564    (B) Identification of fungal genera from PD samples. The classified proportion and precision

565    were derived from different combinations of search algorithms and databases as indicated

566    (box); (C) Identification of fungal genera from PB samples. The classification proportion and

567   precision were derived from the different combinations of search algorithms and databases

568   as indicated.

569   **Figure 2.** Dynamics in precision, completeness and remaining rate after applying progressive

570   cut-offs on BLAST alignment metrics. (A) Cut-offs applied to query length. (B) Cut-offs applied

571   to alignment E-values. (C) Cut-offs applied to the percentage of identical matches. (D) Cut-

572   offs applied to query coverage.

573   **Figure 3.** Benchmarking of amplicon datasets. (A) Scatter plot represented genus level

574   classification proportion and precision for nanopore amplicon data. (B) Genus level precision

575   of Illumina amplicon data. Classification proportion of Illumina data were 100% due to the

576   nature of the QIIME2 pipeline (based on the UNITE ITS database). (C) Genus level

577   completeness of both nanopore and Illumina amplicon datasets. The nanopore results are

578   from minimap2 algorithm and uniteITS database.

579   **Figure 4.** Improving community composition analysis by applying query coverage cut-offs. (A)

580   Experimental flowchart for analysing community compositions. (B) Statistical similarity

581   measures between gold standard community composition and each combination of

582   algorithms and databases. Lower values correspond to greater similarity between the

583   samples and the gold standard. (C) Change in Bhattacharyya distance after applying cut-offs

584   to query coverage for each dataset as indicated. The query coverage gap between each dot

585   point is 0.5%. (D) Change in relative Euclidean distance after applying cut-offs to query

586   coverage for each dataset. The gap between each dot point is 0.5%. (E) Change in relative

587   entropy after applying cut-offs on query coverage for each dataset. The gap between each

588   dot point is 0.5%.

589   **Supplementary Table S1.** Metadata of the mock fungal community

590    **Supplementary Table S2.** Assembly statistics of the draft genomes of *Candida rugosa*,

591    *Candida mesorugosa* and *Cryptococcus magnus* in the mock fungal community.

592    **Supplementary Figure S1.** Change of alignment metrics after applying cut-offs on Phred score.

593

594    **Reference**

595    1. Hawksworth DL, Lücking R. Fungal Diversity Revisited: 2.2 to 3.8 Million Species. Fungal
596    Kingd. 2017;79–95.

597    2. Cheek M, Lughadha EN, Kirk P, Lindon H, Carretero J, Looney B, et al. New scientific
598    discoveries: Plants and fungi. PLANTS PEOPLE PLANET. 2020;2:371–88.

599    3. Mora C, Tittensor DP, Adl S, Simpson AGB, Worm B. How Many Species Are There on Earth
600    and in the Ocean? PLOS Biol. Public Library of Science; 2011;9:e1001127.

601    4. Freimoser F. Start teaching mycology! [Internet]. Nat. Res. Microbiol. Community. 2017
602    [cited         2021         Jan         6].         Available         from:
603    https://naturemicrobiologycommunity.nature.com/posts/20287-start-teaching-mycology

604    5. Editorial. Stop neglecting fungi. Nat Microbiol. Nature Publishing Group; 2017;2:1–2.

605    6. Fisher MC, Gurr SJ, Cuomo CA, Blehert DS, Jin H, Stukenbrock EH, et al. Threats Posed by
606    the Fungal Kingdom to Humans, Wildlife, and Agriculture. mBio [Internet]. American Society
607    for   Microbiology;   2020   [cited   2021   Jan   6];11.   Available   from:
608    https://mbio.asm.org/content/11/3/e00449-20

609    7. Scheele BC, Pasmans F, Skerratt LF, Berger L, Martel A, Beukema W, et al. Amphibian fungal
610    panzootic causes catastrophic and ongoing loss of biodiversity. Science. American Association
611    for the Advancement of Science; 2019;363:1459–63.

612    8. Naranjo-Ortiz MA, Gabaldón T. Fungal evolution: diversity, taxonomy and phylogeny of the
613    Fungi. Biol Rev. 2019;94:2101–37.

614    9. Valverde ME, Hernández-Pérez T, Paredes-López O. Edible Mushrooms: Improving Human
615    Health and Promoting Quality Life [Internet]. Int. J. Microbiol. Hindawi; 2015 [cited 2020 Oct
616    19]. p. e376387. Available from: https://www.hindawi.com/journals/ijmicro/2015/376387/

617    10. Kong HH, Segre JA. Cultivating fungal research. Science. American Association for the
618    Advancement of Science; 2020;368:365–6.

619    11. Whittaker RH. New Concepts of Kingdoms of Organisms. Science. American Association
620    for the Advancement of Science; 1969;163:150–60.

621    12. James TY, Kauff F, Schoch CL, Matheny PB, Hofstetter V, Cox CJ, et al. Reconstructing the
622    early evolution of Fungi using a six-gene phylogeny. Nature. Nature Publishing Group;
623    2006;443:818–22.

624    13. White MM, James TY, O'Donnell K, Cafaro MJ, Tanabe Y, Sugiyama J. Phylogeny of the
625    Zygomycota based on nuclear ribosomal sequence data. Mycologia. Taylor & Francis;
626    2006;98:872–84.

627    14. James TY, Letcher PM, Longcore JE, Mozley-Standridge SE, Porter D, Powell MJ, et al. A
628    molecular phylogeny of the flagellated fungi (Chytridiomycota) and description of a new
629    phylum (Blastocladiomycota). Mycologia. Taylor & Francis; 2006;98:860–71.

630    15. Fischer WM, Palmer JD. Evidence from small-subunit ribosomal RNA sequences for a
631    fungal origin of Microsporidia. Mol Phylogenet Evol. 2005;36:606–22.

632    16. Keeling PJ, Luker MA, Palmer JD. Evidence from Beta-Tubulin Phylogeny that
633    Microsporidia Evolved from Within the Fungi. Mol Biol Evol. 2000;17:23–31.

634    17. Jones MDM, Forn I, Gadelha C, Egan MJ, Bass D, Massana R, et al. Discovery of novel
635    intermediate forms redefines the fungal tree of life. Nature. Nature Publishing Group;
636    2011;474:200–3.

637    18. Adl SM, Simpson AGB, Lane CE, Lukeš J, Bass D, Bowser SS, et al. The Revised Classification
638    of Eukaryotes. J Eukaryot Microbiol. 2012;59:429–514.

639    19. Raja HA, Miller AN, Pearce CJ, Oberlies NH. Fungal Identification Using Molecular Tools: A
640    Primer for the Natural Products Research Community. J Nat Prod. 2017;80:756–70.

641    20. Schirmer M, Ijaz UZ, D'Amore R, Hall N, Sloan WT, Quince C. Insight into biases and
642    sequencing errors for amplicon sequencing with the Illumina MiSeq platform. Nucleic Acids
643    Res. 2015;43:e37.

644    21. Filippis FD, Laiola M, Blaiotta G, Ercolini D. Different Amplicon Targets for Sequencing-
645    Based Studies of Fungal Diversity. Appl Environ Microbiol [Internet]. American Society for
646    Microbiology;    2017    [cited    2021    Feb    15];83.    Available    from:
647    https://aem.asm.org/content/83/17/e00905-17

648    22. Frau A, Kenny JG, Lenzi L, Campbell BJ, Ijaz UZ, Duckworth CA, et al. DNA extraction and
649    amplicon production strategies deeply inf luence the outcome of gut mycobiome studies. Sci
650    Rep. Nature Publishing Group; 2019;9:9328.

651    23. Heeger F, Bourne EC, Baschien C, Yurkov A, Bunk B, Spröer C, et al. Long-read DNA
652    metabarcoding of ribosomal RNA in the analysis of fungi from aquatic environments. Mol Ecol
653    Resour. 2018;18:1500–14.

654    24. D'Andreano S, Cuscó A, Francino O. Rapid and real-time identification of fungi up to the
655    species level with long amplicon Nanopore sequencing from clinical samples. bioRxiv. Cold
656    Spring Harbor Laboratory; 2020;2020.02.06.936708.

657    25. Quince C, Walker AW, Simpson JT, Loman NJ, Segata N. Shotgun metagenomics, from
658    sampling to analysis. Nat Biotechnol. 2017;35:833–44.

659    26. Marcelino VR, Clausen PTLC, Buchmann JP, Wille M, Iredell JR, Meyer W, et al. CCMetagen:
660    comprehensive and accurate identification of eukaryotes and prokaryotes in metagenomic
661    data. Genome Biol. 2020;21:103.

662    27. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool. J
663    Mol Biol. 1990;215:403–10.

664    28. Wood DE, Salzberg SL. Kraken: ultrafast metagenomic sequence classification using exact
665    alignments. Genome Biol. 2014;15:R46.

666    29. Wood DE, Lu J, Langmead B. Improved metagenomic analysis with Kraken 2. Genome Biol.
667    2019;20:257.

668    30. Beghini F, McIver LJ, Blanco-Míguez A, Dubois L, Asnicar F, Maharjan S, et al. Integrating
669    taxonomic, functional, and strain-level profiling of diverse microbial communities with
670    bioBakery 3. bioRxiv. Cold Spring Harbor Laboratory; 2020;2020.11.19.388223.

671    31. Zielezinski A, Girgis HZ, Bernard G, Leimeister C-A, Tang K, Dencker T, et al. Benchmarking
672    of alignment-free sequence comparison methods. Genome Biol. 2019;20:144.

673    32. McIntyre ABR, Ounit R, Afshinnekoo E, Prill RJ, Hénaff E, Alexander N, et al. Comprehensive
674    benchmarking and ensemble approaches for metagenomic classifiers. Genome Biol.
675    2017;18:182.

676    33. Ye SH, Siddle KJ, Park DJ, Sabeti PC. Benchmarking Metagenomics Tools for Taxonomic
677    Classification. Cell. 2019;178:779–94.

678    34. Nasko DJ, Koren S, Phillippy AM, Treangen TJ. RefSeq database growth influences the
679    accuracy of k-mer-based lowest common ancestor species identification. Genome Biol.
680    2018;19:165.

681    35. R. Marcelino V, Holmes EC, Sorrell TC. The use of taxon-specific reference databases
682    compromises metagenomic classification. BMC Genomics. 2020;21:184.

683    36. Heeger F, Wurzbacher C, Bourne EC, Mazzoni CJ, Monaghan MT. Combining the 5.8S and
684    ITS2 to improve classification of fungi. Methods Ecol Evol. 2019;10:1702–11.

685    37. NCBI Resource Coordinators. Database resources of the National Center for Biotechnology
686    Information. Nucleic Acids Res. 2018;46:D8–13.

687    38. O'Leary NA, Wright MW, Brister JR, Ciufo S, Haddad D, McVeigh R, et al. Reference
688    sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and functional
689    annotation. Nucleic Acids Res. 2016;44:D733–45.

690    39. Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W, et al. Gapped BLAST and
691    PSI-BLAST: a new generation of protein database search programs. Nucleic Acids Res.
692    1997;25:3389–402.

693    40. Xiao L, Feng Q, Liang S, Sonne SB, Xia Z, Qiu X, et al. A catalog of the mouse gut
694    metagenome. Nat Biotechnol. Nature Publishing Group; 2015;33:1103–8.

695    41. Hu Y, Green GS, Milgate AW, Stone EA, Rathjen JP, Schwessinger B. Pathogen Detection
696    and Microbiome Analysis of Infected Wheat Using a Portable DNA Sequencer. Phytobiomes J.
697    Scientific Societies; 2019;3:92–101.

698    42. Irinyi L, Hu Y, Hoang MTV, Pasic L, Halliday C, Jayawardena M, et al. Long-read sequencing
699    based clinical metagenomics for the detection and confirmation of Pneumocystis jirovecii
700    directly from clinical specimens: A paradigm shift in mycological diagnostics. Med Mycol.
701    Oxford Academic; 2020;58:650–60.

702    43. Xiao L, Estellé J, Kiilerich P, Ramayo-Caldas Y, Xia Z, Feng Q, et al. A reference gene
703    catalogue of the pig gut microbiome. Nat Microbiol. Nature Publishing Group; 2016;1:1–6.

704    44. Donovan PD, Gonzalez G, Higgins DG, Butler G, Ito K. Identification of fungi in shotgun
705    metagenomics datasets. PLOS ONE. Public Library of Science; 2018;13:e0192898.

706    45. Holman DB, Brunelle BW, Trachsel J, Allen HK. Meta-analysis To Define a Core Microbiota
707    in the Swine Gut. mSystems. 2017;2.

708    46. Wang J, Lang T, Shen J, Dai J, Tian L, Wang X. Core Gut Bacteria Analysis of Healthy Mice.
709    Front Microbiol [Internet]. Frontiers; 2019 [cited 2021 Jan 6];10. Available from:
710    https://www.frontiersin.org/articles/10.3389/fmicb.2019.00887/full

711    47. White TJ, Bruns T, Lee S, Taylor J. AMPLIFICATION AND DIRECT SEQUENCING OF FUNGAL
712    RIBOSOMAL RNA GENES FOR PHYLOGENETICS. PCR Protoc [Internet]. Elsevier; 1990 [cited
713    2021    Mar    19].    p.    315–22.    Available    from:
714    https://linkinghub.elsevier.com/retrieve/pii/B9780123721808500421

715    48. Nilsson RH, Larsson K-H, Taylor AFS, Bengtsson-Palme J, Jeppesen TS, Schigel D, et al. The
716    UNITE database for molecular identification of fungi: handling dark taxa and parallel
717    taxonomic classifications. Nucleic Acids Res. Oxford Academic; 2019;47:D259–64.

718    49. Bharti R, Grimm DG. Current challenges and best-practice protocols for microbiome
719    analysis. Brief Bioinform [Internet]. 2019 [cited 2021 Jan 7]; Available from:
720    https://doi.org/10.1093/bib/bbz155

721    50. Bolyen E, Rideout JR, Dillon MR, Bokulich NA, Abnet CC, Al-Ghalith GA, et al. Reproducible,
722    interactive, scalable and extensible microbiome data science using QIIME 2. Nat Biotechnol.
723    Nature Publishing Group; 2019;37:852–7.

724    51. Johnson JS, Spakowicz DJ, Hong B-Y, Petersen LM, Demkowicz P, Chen L, et al. Evaluation
725    of 16S rRNA gene sequencing for species and strain-level microbiome analysis. Nat Commun.
726    Nature Publishing Group; 2019;10:5029.

727    52. Krehenwinkel H, Pomerantz A, Henderson JB, Kennedy SR, Lim JY, Swamy V, et al.
728    Nanopore sequencing of long ribosomal DNA amplicons enables portable and simple
729    biodiversity assessments with high phylogenetic resolution across broad taxonomic scale.

730    GigaScience [Internet]. 2019 [cited 2021 Jan 7];8. Available from:
731    https://doi.org/10.1093/gigascience/giz006

732    53. Dokmanic I, Parhizkar R, Ranieri J, Vetterli M. Euclidean Distance Matrices: Essential
733    Theory, Algorithms and Applications. 2015 [cited 2021 Jan 7]; Available from:
734    https://arxiv.org/abs/1502.07541v2

735    54. Aherne FJ, Thacker NA, Rockett PI. The Bhattacharyya metric as an absolute similarity
736    measure for frequency coded data. Kybernetika. 1998;34:363–8.

737    55. MacKay DJC, Kay DJCM. Information Theory, Inference and Learning Algorithms.
738    Cambridge University Press; 2003.

739    56. Henderson G, Cox F, Kittelmann S, Miri VH, Zethof M, Noel SJ, et al. Effect of DNA
740    Extraction Methods and Sampling Techniques on the Apparent Structure of Cow and Sheep
741    Rumen Microbial Communities. PLOS ONE. Public Library of Science; 2013;8:e74787.

742    57. Davis A, Kohler C, Alsallaq R, Hayden R, Maron G, Margolis E. Improved yield and accuracy
743    for DNA extraction in microbiome studies with variation in microbial biomass. BioTechniques.
744    Future Science; 2019;66:285–9.

745    58. Douglas CA, Ivey KL, Papanicolas LE, Best KP, Muhlhausler BS, Rogers GB. DNA extraction
746    approaches substantially influence the assessment of the human breast milk microbiome. Sci
747    Rep. Nature Publishing Group; 2020;10:123.

748    59. Tamames J, Cobo-Simón M, Puente-Sánchez F. Assessing the performance of different
749    approaches for functional and taxonomic annotation of metagenomes. BMC Genomics.
750    2019;20:960.

751    60. Liang Q, Bible PW, Liu Y, Zou B, Wei L. DeepMicrobes: taxonomic classification for
752    metagenomics with deep learning. NAR Genomics Bioinforma [Internet]. 2020 [cited 2021 Jan
753    6];2. Available from: https://doi.org/10.1093/nargab/lqaa009

754    61. Pearman WS, Freed NE, Silander OK. Testing the advantages and disadvantages of short-
755    and long- read eukaryotic metagenomics using simulated reads. BMC Bioinformatics.
756    2020;21:220.

757    62. Mallawaarachchi V, Wickramarachchi A, Lin Y. GraphBin: refined binning of metagenomic
758    contigs using assembly graphs. Bioinformatics. 2020;36:3307–13.

759    63. Alneberg J, Bjarnason BS, de Bruijn I, Schirmer M, Quick J, Ijaz UZ, et al. Binning
760    metagenomic contigs by coverage and composition. Nat Methods. Nature Publishing Group;
761    2014;11:1144–6.

762    64. Wickramarachchi A, Mallawaarachchi V, Rajan V, Lin Y. MetaBCC-LR: metagenomics
763    binning by coverage and composition for long reads. Bioinformatics. 2020;36:i3–11.

764  65. Kaehler BD, Bokulich NA, McDonald D, Knight R, Caporaso JG, Huttley GA. Species
765  abundance information improves sequence taxonomy classification accuracy. Nat Commun.
766  Nature Publishing Group; 2019;10:4643.

767  66. Bokulich NA, Dillon MR, Bolyen E, Kaehler BD, Huttley GA, Caporaso JG. q2-sample-
768  classifier: machine-learning tools for microbiome classification and regression. J Open Res
769  Softw [Internet]. 2018 [cited 2021 Jan 6];3. Available from:
770  https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6759219/

771  67. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence
772  data. Bioinformatics. 2014;30:2114–20.

773  68. Peng Y, Leung HCM, Yiu SM, Chin FYL. IDBA-UD: a de novo assembler for single-cell and
774  metagenomic sequencing data with highly uneven depth. Bioinformatics. Oxford Academic;
775  2012;28:1420–8.

776  69. Li H, Durbin R. Fast and accurate short read alignment with Burrows–Wheeler transform.
777  Bioinformatics. 2009;25:1754–60.

778  70. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence
779  Alignment/Map format and SAMtools. Bioinformatics. 2009;25:2078–9.

780  71. Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic features.
781  Bioinformatics. Oxford Academic; 2010;26:841–2.

782  72. Hu Y, Schwessinger B. Amplicon sequencing using MinION optimized from 1D native
783  barcoding genomic DNA [Internet]. protocols.io. 2018 [cited 2018 Sep 27]. Available from:
784  https://www.protocols.io/view/amplicon-sequencing-using-minion-optimized-from-1d-
785  mhkc34w

786  73. Kolmogorov M, Yuan J, Lin Y, Pevzner PA. Assembly of long, error-prone reads using repeat
787  graphs. Nat Biotechnol. Nature Publishing Group; 2019;37:540–6.

788  74. Vaser R, Sović I, Nagarajan N, Šikić M. Fast and accurate de novo genome assembly from
789  long uncorrected reads. Genome Res. 2017;27:737–46.

790  75. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO: assessing
791  genome assembly and annotation completeness with single-copy orthologs. Bioinformatics.
792  Oxford Academic; 2015;31:3210–2.

793  76. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+:
794  architecture and applications. BMC Bioinformatics. 2009;10:421.

795  77. Lu J, Salzberg SL. Removing contaminants from databases of draft genomes. PLOS Comput
796  Biol. 2018;14:e1006277.

797  78. Huerta-Cepas J, Serra F, Bork P. ETE 3: Reconstruction, Analysis, and Visualization of
798  Phylogenomic Data. Mol Biol Evol. 2016;33:1635–8.

799   79. Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJA, Holmes SP. DADA2: High-
800   resolution sample inference from Illumina amplicon data. Nat Methods. Nature Publishing
801   Group; 2016;13:581–3.

802   80. Bokulich NA, Kaehler BD, Rideout JR, Dillon M, Bolyen E, Knight R, et al. Optimizing
803   taxonomic classification of marker-gene amplicon sequences with QIIME 2's q2-feature-
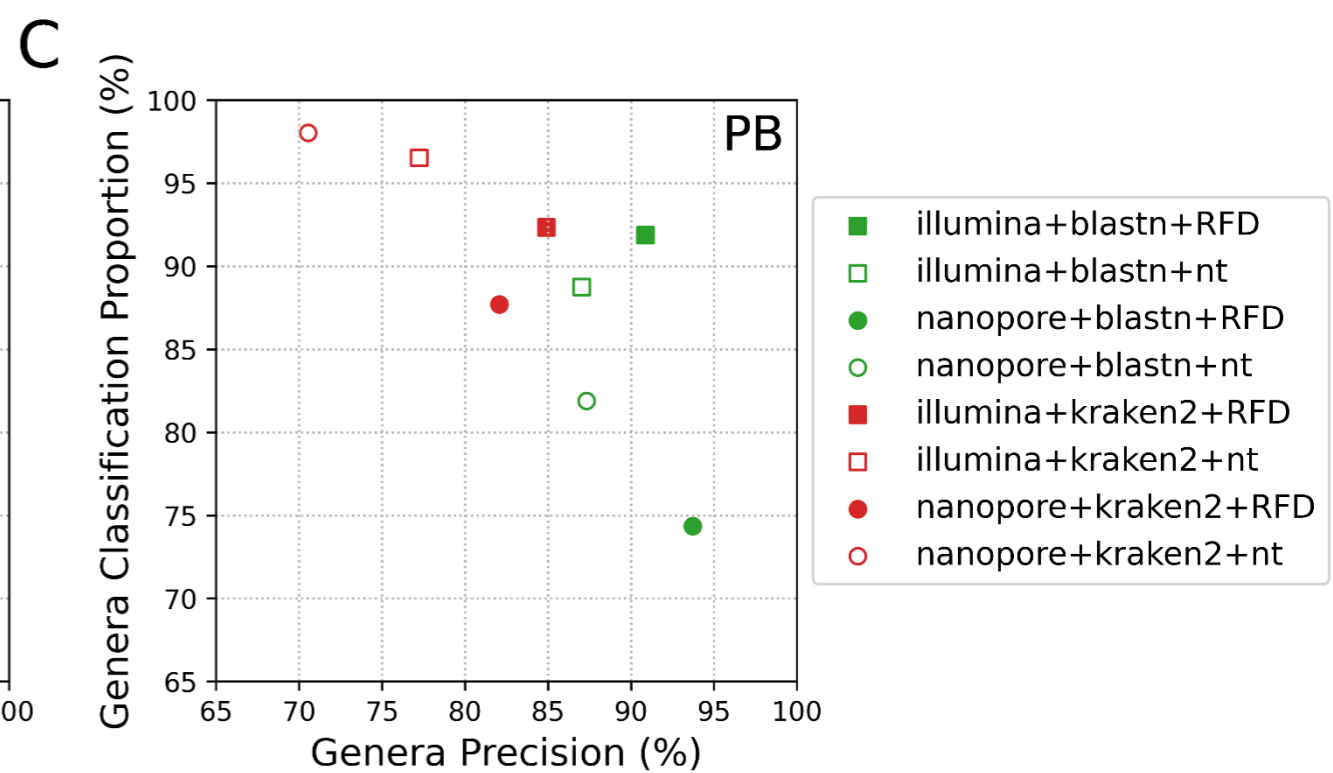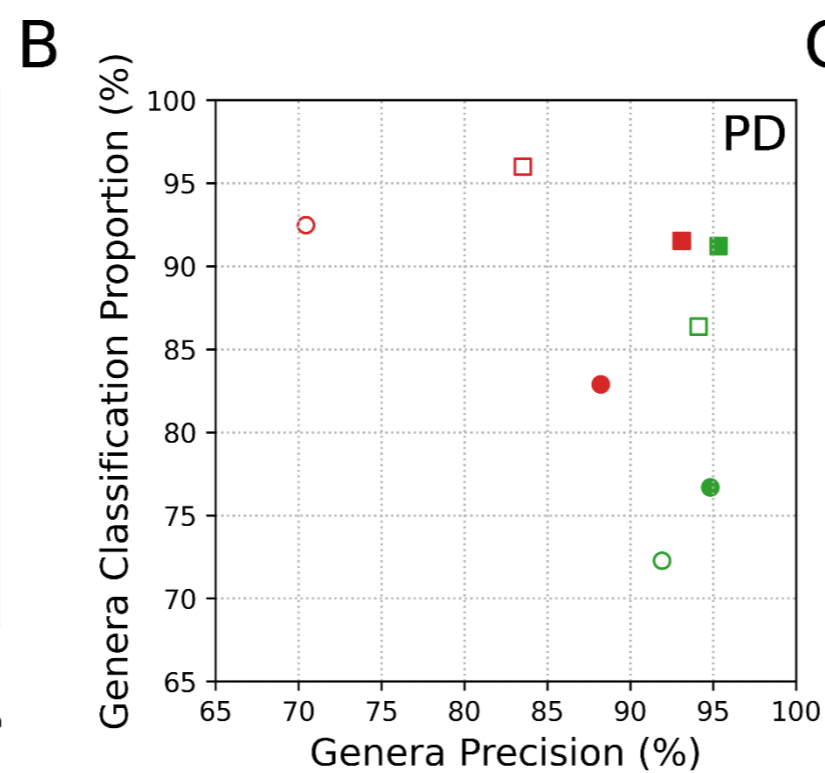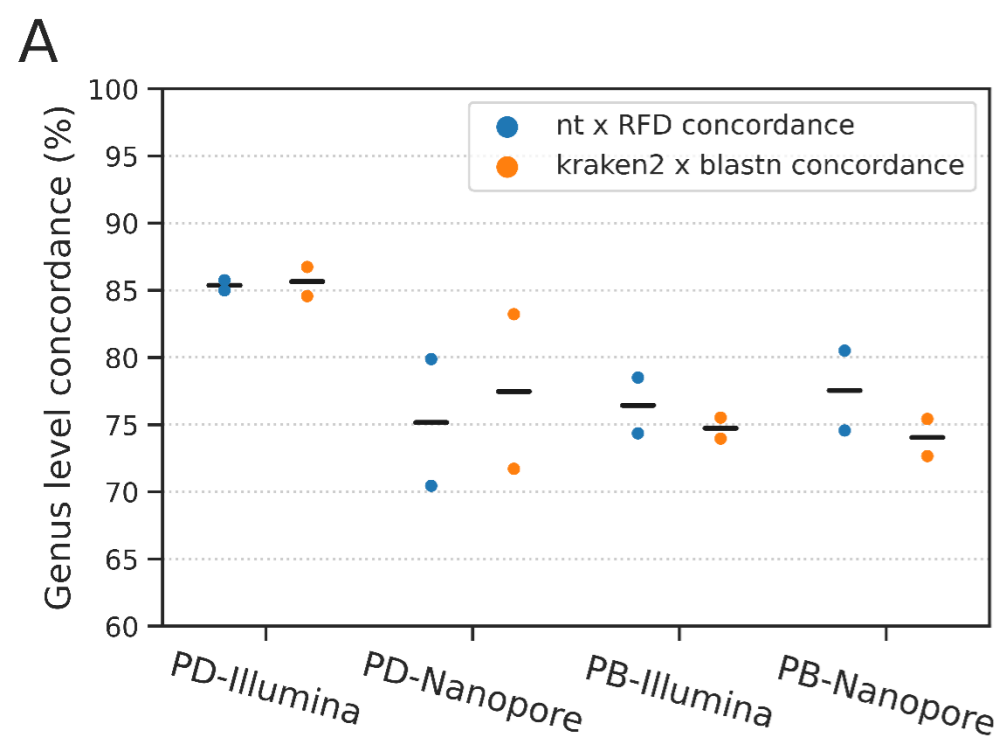804   classifier plugin. Microbiome. 2018;6:90.

805

Table 1. The characteristics for each dataset.

| Sample | Sequencing Tech | Sequencing Strategy | # Basepairs | # reads | # Assembled contigs | # Mapped basepairs (Gb) |
|---|---|---|---|---|---|---|
| PD | Illumina | Shotgun | 3.91 Gb | 14525058 | 338823 | 3.69 |
| | | Amplicon | 66.9/95.8/106.4 Mb[a] | 39374/9614/10236[b] | N/A | N/A |
| | Nanopore | Shotgun | 1.96 Gb | 1273484 | N/A | N/A |
| | | Amplicon | 71.5/72.5/86.5 Mb | 26212/ 26680/ 31826[b] | N/A | N/A |
| PB | Illumina | Shotgun | 3.67 Gb | 13623120 | 345009 | 3.44 |
| | | Amplicon | 55.7/38.1/71.9 Mb[a] | 23613/13828/27093[b] | N/A | N/A |
| | Nanopore | Shotgun | 3.78 Gb | 1043343 | N/A | N/A |
| | | Amplicon | 54.5/49.4/42.0 Mb | 20163/ 18273/ 15502[b] | N/A | N/A |

[a] The total basepairs of each technical replicate were calculated before importing into QIIME2 pipeline.
[b] Number of nanopore reads or paired-end Illumina reads for technical replicate 1/replicate 2/replicate 3 after quality control.

Table 2. Assignment of published sequence data to genera after application of cut-offs to query coverage.

| Sample ID | Sample description | Sequencing tech | Cut-offs on query coverage (%) | Filtered results (%) | Percentage of confirmed genera BEFORE applying cut-offs (%) | Percentage of confirmed genera AFTER applying cut-offs (%) |
|---|---|---|---|---|---|---|
| a1 | Human sputum samples[42] | Nanopore | 59 | 20.2 | 85.9 | 86.5 |
| a2 | | | 53.2 | 20.1 | 97.9 | 98.5 |
| a3 | | | 54 | 20.5 | 96.5 | 97.4 |
| a4 | | | 45.5 | 20.1 | 16.2 | 19.8 |
| a5 | | | 58.5 | 20 | 71.1 | 66.9 |
| a6 | | | 50.4 | 20.1 | 93.6 | 94.7 |
| b1 | Field infected wheat samples[41] | | 5 | 20 | 60.4 | 75.1 |
| b2 | | | 0.77 | 19.9 | 34.8 | 43 |
| b3 | | | 12 | 19.7 | 67 | 82 |
| b4 | | | 0.61 | 20 | 5.8 | 6.2 |
| c1 | Pig gut microbiome samples[43] | Illumina | 2.4 | 20.1 | 32 | 35.4 |
| c2 | | | 3.3 | 20.2 | 34.2 | 36.6 |
| c3 | | | 2.6 | 20.2 | 35.2 | 38.3 |
| d1 | Mouse gut microbiome samples[40] | | 3.4 | 19.8 | 29.1 | 24.3 |
| d2 | | | 14 | 20.1 | 63.7 | 69.4 |
| d3 | | | 4.5 | 20.2 | 38.6 | 42.3 |