1   # Metatranscriptomics as a tool to identify fungal species

2   # and subspecies in mixed communities

3

4

5   Vanesa R. Marcelino[1,2], Laszlo Irinyi[1,2,3], John-Sebastian Eden[1,2], Wieland Meyer[1,2,3], Edward

6   C. Holmes[1,4], Tania C. Sorrell[1,2]

7

8

9   [1]Marie Bashir Institute for Infectious Diseases and Biosecurity and Faculty of Medicine and

10  Health, Sydney Medical School, Westmead Clinical School, The University of Sydney,

11  Sydney, NSW 2006, Australia.

12

13  [2]Molecular Mycology Research Laboratory, Centre for Infectious Diseases and

14  Microbiology, Westmead Institute for Medical Research, Westmead, NSW 2145, Australia.

15

16  [3]Westmead Hospital (Research and Education Network), Westmead, NSW 2145, Australia.

17

18  [4]School of Life & Environmental Sciences, Charles Perkins Centre, The University of

19  Sydney, Sydney, NSW 2006, Australia.

20

## Abstract

High-throughput sequencing (HTS) enables the generation of large amounts of genome sequence data at a reasonable cost. Organisms in mixed microbial communities can now be sequenced and identified in a culture-independent way, usually using amplicon sequencing of a DNA barcode. Bulk RNA-seq (metatranscriptomics) has several advantages over DNA-based amplicon sequencing: it is less susceptible to amplification biases, it captures only living organisms, and it enables a larger set of genes to be used for taxonomic identification. Using a defined mock community comprised of 17 fungal isolates, we evaluated whether metatranscriptomics can accurately identify fungal species and subspecies in mixed communities. Overall, 72.9% of the RNA transcripts were classified, from which the vast majority (99.5%) were correctly identified at the species-level. Of the 15 species sequenced, 13 were retrieved and identified correctly. We also detected strain-level variation within the *Cryptococcus* species complexes: 99.3% of transcripts assigned to *Cryptococcus* were classified as one of the four strains used in the mock community. Laboratory contaminants and/or misclassifications were diverse but represented only 0.44% of the transcripts. Hence, these results show that it is possible to obtain accurate species- and strain-level fungal identification from metatranscriptome data as long as taxa identified at low abundance are discarded to avoid false-positives derived from contamination or misclassifications. This study therefore establishes a base-line for the application of metatranscriptomics in clinical mycology and ecological studies.

## Introduction

41

42 Microscopic fungal species, such as yeasts and some filamentous fungi, do not live in

43 isolation, they are most commonly found within mixed microbial communities inhabiting

44 soil, water systems, plants and animal hosts. Assessing the diversity of fungi in mixed

45 communities is important because different fungal taxa may exhibit distinctive phenotypes,

46 and consequently may have different pathogenicity or functional roles. For example, in the

47 rhizosphere, changes in fungal community composition have been associated with shifts in

48 nutrient cycling (Hannula *et al.* 2017). Humans also harbor, or are exposed to, a diverse

49 fungal community that provides a source of opportunistic pathogens (Bandara *et al.* 2019;

50 Huffnagle & Noverr 2013; Seed 2014). Although it is typically assumed that invasive fungal

51 infections are caused by a single strain, multiple *Candida* strains have been observed

52 during the course of a single episode of infection (Soll *et al.* 1988). Furthermore, nearly 20%

53 of patients with cryptococcosis are infected with multiple strains, with different phenotypes

54 and virulence traits (Desnos-Ollivier *et al.* 2015; Desnos-Ollivier *et al.* 2010). Strain-level

55 fungal diversity may influence therapeutic responsiveness and needs further investigation.

56 Despite its importance, fungal taxonomic diversity is poorly characterized. From

57 over two million fungal species estimated to exist, less than 8% have been described

58 (Hawksworth & Lucking 2017). Even well-known fungal species are often overlooked during

59 routine diagnostic procedures, surveillance and biodiversity surveys (Brown *et al.* 2012;

60 Enaud *et al.* 2018; Yahr *et al.* 2016). This is in part due to challenges in the detection and

61 classification of these organisms, especially microscopic and cryptic species, for example,

62 the etiologic agents of cryptococcosis. Currently, two species complexes are recognized:

63 *Cryptococcus neoformans* and *Cryptococcus gattii* (Kwon-Chung *et al.* 2002). Seven major

64 haploid lineages are found within these two species complexes (*C. neoformans* species

65 complex: VNI, VNII, VNIV, and *C. gattii* species complex: VGI, VGII, VGIII and VGIV) and

66 their recognition as distinct biological species has been debated (Hagen *et al.* 2015; Kwon-

67 Chung *et al.* 2017; Ngamskulrungroj *et al.* 2009). Being able to distinguish closely-related

68 lineages is important because their phenotype, virulence and ecophysiology can vary

69 substantially. For example, the JEC21 and B-3501 strains of *C. neoformans* var.

70 *neoformans* (VNIV) are 99.5% identical at the genomic sequence level but differ

71 substantially in thermotolerance and virulence (Loftus *et al.* 2005). Likewise, different

72 virulence and antifungal tolerance traits were observed within lineages of *C. gattii* VGIII

73 (Firacative *et al.* 2016).

74 The introduction of high-throughput sequencing (HTS) marked a new era in

75 mycological research, where the vast diversity of fungi can be studied without the need for

76    culture (Nilsson *et al.* 2019). To date, amplicon sequencing of marker genes

77    (metabarcoding) has been the most popular HTS method used to identify fungal species in

78    mixed communities. Despite its indisputable utility, metabarcoding surveys are affected by

79    PCR amplification biases, and even abundant species can go undetected due to primer

80    mismatch (Marcelino & Verbruggen 2016; Nilsson *et al.* 2019; Tedersoo *et al.* 2015). In

81    addition, DNA fragments from dead organisms inflate biodiversity estimates in

82    metabarcoding surveys (Carini *et al.* 2016). Stool samples, for instance, naturally contain

83    food-derived DNA, which cannot be distinguished from the genetic material of the resident

84    gut microbiota when using DNA-based methods. These challenges can be circumvented by

85    directly sequencing actively transcribed genes, via RNA-Seq, hence avoiding the

86    amplification step, and obtaining an unbiased characterization of the living microbial

87    community. Metatranscriptomics has been used to identify RNA viruses in a range of

88    animal samples (Shi *et al.* 2016; Shi *et al.* 2017; Wille *et al.* 2018; Zhang *et al.* 2018) and to

89    characterize the functional profile of microbial communities (Bashiardes *et al.* 2016; Kuske

90    *et al.* 2015). Studies applying metatranscriptomics to mycorrhizal communities have

91    provided valuable insights into the functional roles of fungi in these symbiotic systems

92    (Gonzalez *et al.* 2018; Liao *et al.* 2014). However, links between functional and species-level

93    taxonomy have been sought infrequently, likely because fungal identification from

94    metatranscriptome data is considered unreliable below phylum level (Nilsson *et al.* 2019).

95    Critically, it is currently unknown whether metatranscriptomics can accurately identify fungi

96    at the species and subspecies level within a mixed community. This information is

97    fundamental to the investigation of the potential and utility of metatranscriptomics in

98    diagnostics and ecological studies.

99          Herein, we evaluated the utility of metatranscriptomics as a tool for the

100   simultaneous identification of fungal species, using a defined mock community containing

101   15 fungal species. In addition, we investigated whether strains belonging to sister species,

102   such as the *C. neoformans* and *C. gattii* species complexes could be identified correctly

103   using metatranscriptomics. Rather than focusing on marker genes, we sought to classify

104   fungal species using the information from all expressed genes, using the totality of NCBI's

105   nucleotide collection as a reference database. This study paves the way to apply state-of-

106   the art techniques in fungal biodiversity surveys and clinical diagnostics.

107

## Methods

109   A defined fungal community was constructed from 17 isolates, including 15 fungal species

110   and three strains of the *C. neoformans* species complex in addition to one strain of *C. gattii*

111    (Table 1). Fungal strains were obtained from the Westmead Mycology Culture Collection

112    and cultured on Sabouraud agar at 27°C for 72 hours. A sweep of colonies was made with

113    a disposable inoculating loop and dispersed in PBS. Fungal cells were quantified in a

114    Neubauer chamber and their concentration adjusted such that the fungal mixture contained

115    equal concentrations of each species ($10^8$ cells/mL). RNA was isolated with the RNeasy

116    Plus kit (Qiagen), following the manufacture's protocol, with an initial freeze-thaw step in

117    liquid nitrogen to disrupt fungal cells. The quantity and quality of the RNA extract was

118    determined with the Nanodrop Spectrophotometer (Thermo Scientific) and the Agilent 2200

119    TapeStation. As some residual DNA was detected, the RNA extract was further treated with

120    DNase I (Qiagen). Ribosomal depletion (Ribo-Zero Gold technology), library preparation and

121    sequencing (Illumina HiSeq HT, 125bp Paired End) were performed by the Australian

122    Genomics Research Facility. The raw sequence data were deposited in the NCBI Short

123    Read Archive (accession PRJNA521097).

124        Sequence reads containing more than five ambiguous positions or with average

125    quality scores $\leq 25$ were filtered from the dataset using prinseq-lite v.0.20.4 (Schmieder &

126    Edwards 2011) with the options -ns_max_n 5 -min_qual_mean 25 -out_format 3. Assembly

127    of sequence reads into contigs was performed with Trinity v.2.5.1 (Grabherr *et al.* 2011).

128    Contigs were mapped to the NCBI nucleotide collection using KMA (Clausen *et al.* 2018), a

129    novel approach that has proven to be more accurate than other mapping software. Prior to

130    mapping, NCBI's taxonomic identifier codes (taxids) were appended to each sequence

131    record in the nucleotide collection, and the reference database was indexed using KMA's

132    options -NI -Sparse TG. Contigs were then mapped to the indexed database with the

133    options -mem_mode -and -apm f. Matches to the reference database with low support (*i.e.*

134    coverage < 20 and depth < 0.05) were excluded from the analyses. The species-level

135    taxonomic classifications were based on NCBI's taxonomy identifiers (taxids) to minimize

136    the issue of changing species nomenclature (Federhen 2012). Subspecies-level

137    classifications within the *Cryptococcus neoformans* and *C. gattii* species complexes were

138    examined manually.

139        Abundance was estimated at the level of sequence reads and transcripts. For read-

140    level abundances, sequence reads were mapped to transcripts using Bowtie2 (Langmead &

141    Salzberg 2012) and quantified in Transcripts Per Million (TPM) with RSEM (Li & Dewey

142    2011), using the Trinity pipeline. For transcript-level abundances, the depth values

143    estimated within KMA were used, which is the total number of nucleotides (in each contig)

144    covering the reference sequence divided by the length of the reference sequence. The

145    number and length of assembled contigs for each taxon is likely a better proxy for species

146    abundance than read-level abundances (which are subject to gene expression), and

147    therefore were used for graphic representation and analyses. For simplicity, we refer as

148    'abundance' the transcript-level abundance, unless otherwise stated.

149        It would be logical to expect that species with larger and gene-rich genomes would

150    express a greater number of transcripts. To test for this potential correlation, genome sizes

151    and the estimated number of proteins were obtained from the Fungal Genome Size

152    Database (Kullman *et al.* 2005), Loftus et al. (2005), Muñoz et al. (2018) and NCBI's Genome

153    database (Supplementary table S1). The correlation coefficients between genome size,

154    number of proteins and abundance of transcripts were estimated using Person's correlation

155    and visualized using the R package *ggpubr* (Kassambara 2017).

156

157

## Results

159    RNA sequencing yielded a total of 26,558,491 paired end reads, of which 98.3% passed

160    quality control. Overall, 277,404 contigs (transcripts) were obtained, from which 202,219

161    (72.9%) were classified. The majority of the sequence reads (80.2%) mapped to a classified

162    contig. Of the 15 fungal species sequenced, 13 were retrieved and correctly classified at

163    the species level (Figure 1, Table 2, Supplementary table S2). The two false-negatives were

164    *Debaryomyces hansenii* and *Schizosaccharomyces pombe*; these may have been

165    misclassified as another fungus or were lost due to cell pooling inaccuracy and/or RNA

166    extraction biases. A small proportion of bacterial transcripts (0.03%) and other eukaryotic

167    microbes (0.4%, including 31 fungi that were not present in the mock community) was also

168    observed (Table 2, Supplementary table S2), which likely represent laboratory contaminants

169    and misclassifications (see discussion). However, these were present at a consistently

170    lower frequency than true members of the mock community, with the most common –

171    *Candida glycerinogenes* – only present in 0.08% of the transcripts. Some of the transcripts

172    were assigned to entries in GenBank that do not have a species-level classification (*e.g.*

173    *Candida* sp. and *Pichia* sp.). These assignments were considered misclassifications here,

174    although it is possible that the species in our mock community are the correct species-level

175    identity of these GenBank sequences.

176        Overall, the commonest species detected was *C. neoformans*, which was to be

177    expected as it comprised three strains in the mock community and therefore was three

178    times more abundant than other fungal species. Transcripts belonging to *Candida tropicalis*

179    and *Pichia kudriavzevii* (former *Candida krusei*) – were also common (19.2% and 18.8%,

180    respectively), while *C. albicans*, *C. orthopsilosis* and *C. glabrata* (other causes of

6

181 candidaemia in humans) were detected at lower abundance (2.0 – 2.9%). There was no

182 relationship evident between abundance of transcripts and phylogenetic relatedness.

183 Genomes with low GC content can be overrepresented in metagenomic sequencing

184 (Shakya *et al.* 2013). Conversely, some of the species detected here in high abundance

185 (*Cryptococcus neoformans* and *Clavispora lusitaniae*) have a higher GC content than most

186 other fungal species (Dujon 2010), suggesting that GC bias is unlikely to affect our results.

187 No correlation between abundance of transcripts and genome size or estimated number of

188 proteins was observed ($p > 0.05$, Supplementary figure 1).

189       Molecular type and strain-level variation within the *Cryptococcus neoformans* and *C.*

190 *gattii* species complexes was also detected, with contigs matching to *C. gattii* VGI WM 276,

191 *C. neoformans* var. *grubii* VNI H99 and *C. neoformans* var. *neoformans* VNIV strains B-

192 3501A and JEC21 (Figure 2, Supplementary table S3). A proportion of the transcripts (1.6%)

193 matched with equal probability scores to both strains of *C. neoformans* var. *neoformans* (B-

194 3501A and JEC21, Supplementary tables S2 and S3). From the transcripts classified as

195 *Cryptococcus* spp*.,* 99.3% were classified as one of the four *Cryptococcus* strains (or both

196 B-3501A and JEC21) used in the mock community. It is possible that misclassifications

197 occurred within the strains analyzed. For example, transcripts originally from JEC21 might

198 have been classified as B-3501A and *vice versa*. As it is not possible to know from which

199 strain the transcripts originated, these possible misclassifications would be undetected.

200

201

## Discussion

203 Our metatranscriptomics approach yielded taxonomic identification of fungi from a defined

204 mock community with high success, while false-positives were detected at far lower

205 abundance. These results indicate that it is possible to obtain accurate species- and strain-

206 level identifications for fungi from metatranscriptome data, as long as taxa identified at low

207 abundance are removed from the analyses to avoid false-positives derived from

208 contamination or misclassifications.

209       Taxonomic classification at species and strain levels using metabarcoding and

210 metagenomic data has been considered inaccurate (Nilsson *et al.* 2019; Sczyrba *et al.*

211 2017; Yamamoto *et al.* 2014), raising the question of how our metatranscriptomics

212 approach succeeded in identifying closely-related fungal strains. Metabarcoding relies on a

213 single marker gene (Banchi *et al.* 2018; e.g. McGuire *et al.* 2013; Schmidt *et al.* 2013), which

214 does not contain sufficient phylogenetic information to differentiate some closely related

215 fungal lineages (Balasundaram *et al.* 2015; Nilsson *et al.* 2008). Metatranscriptomics, on the

7

216    other hand, yields data on all expressed coding sequences. Classifications derived from

217    metagenomes are likely to be equally accurate as the ones obtained from

218    metatranscriptomes, except that dead organisms might also be sequenced. Additionally,

219    we used a new alignment method that is both highly accurate and fast (Clausen *et al.* 2018),

220    allowing us to map sequences against the complete NCBI nucleotide collection. Complete

221    genomes of all fungal isolates used here are available in NCBI, and it is likely that the

222    accuracy of identifications is reduced for poorly-documented microorganisms. However, it

223    is possible to extract informative genes from metatranscriptome data and subsequently

224    perform phylogenomic analyses to identify rare and novel taxa (e.g. Shi *et al.* 2017; Wille *et*

225    *al.* 2018; Zhang *et al.* 2018). Besides being highly accurate, metatranscriptomics is less

226    susceptible to amplification bias, no information about the community members is required

227    *a priori*, and it only detects functionally active members of a microbial community. These

228    advantages make metatranscriptomics a promising tool in biodiversity surveys, functional

229    assessments of microbial communities, pathogen detection and biosecurity surveillance

230    (e.g. Kuske *et al.* 2015; Shi *et al.* 2016; Wille *et al.* 2018).

231         Even though false-positives were present at low abundance, they pose a challenge

232    in the interpretation of metatranscriptomic and metagenomic data. False-positives generally

233    result from spurious classifications and laboratory contaminants, which may be common in

234    laboratory reagents (Salter *et al.* 2014). However, metatranscriptomics is less sensitive to

235    laboratory contamination than DNA-based metagenomics or metabarcoding, as only living

236    microorganisms are sequenced. Nevertheless, contamination can occur at all stages of the

237    library preparation and is routinely observed in RNA-Seq studies (Quince *et al.* 2017; Strong

238    *et al.* 2014). Misclassifications occur because some genome regions are very similar (or

239    identical) across closely-related species and cannot be differentiated. Errors in reference

240    databases can also result in misclassifications. Sequences attributed to incorrectly-

241    classified species are not uncommon in GenBank and result in downstream classification

242    errors (Li *et al.* 2018). It is also not unusual to find bacterial regions misassembled into

243    eukaryotic genomes (e.g. Koutsovoulos *et al.* 2016), which can result in sequences from

244    common laboratory contaminants being classified as a eukaryote. Filtering out organisms

245    found in low abundance is an option to reduce the incidence of false-positives in

246    downstream analyses. In this study, filtering organisms for which the abundance of

247    transcripts is lower than 0.1% would eliminate false-positives, at the cost of excluding one

248    true-positive from the analyses (Table 2). The application of this abundance-filtering step

249    might not be feasible when sequencing depth (per microbial species) is limited. Species

250  present in low abundances will be represented by a small number of transcripts and so are

251  more likely to be misclassified or undetected.

252      The abundance of transcripts and sequence reads can vary according to genome

253  size, number of coding sequences and gene expression. Therefore, the abundance

254  disparity across species observed here is unsurprising. Interestingly, we found no

255  correlation between the abundance of transcripts and genome size or number of genes

256  (Supplementary figure 1). Imprecise estimates of cell abundance and RNA extraction biases

257  could also have influenced abundance estimates, and might be the reason why two species

258  in the mock community (*D. hansenii* and *S. pombe*) were not detected in the analyses.

259  Metabarcoding studies have suggested that performing DNA extraction in triplicate

260  minimizes biases for bacteria, but it had no effect in fungal communities (Feinstein *et al.*

261  2009). To our knowledge, the effect of RNA extraction bias in metatranscriptomics has yet

262  to be studied. As metagenomics surveys are not affected by gene expression, they might

263  be more appropriate for studies where it is important to quantify species abundance.

264      Although fungal species and their genes can be confidently identified, it remains

265  challenging to link some genes with particular species using metatranscriptomics. A large

266  portion of fungal genomes are highly similar among species, making it difficult, if not

267  impossible, to infer which species in the community are expressing which genes. Recently,

268  a method was developed to perform species-level functional profiling of metagenome data

269  (Franzosa *et al.* 2018). This method, however, relies on a reference database of complete

270  genomes that currently contains few fungal representatives, limiting its application in fungal

271  metagenomics. Contrary to metatranscriptomics, metagenomics yields coding and non-

272  coding sequences, which can facilitate linking genes to species if sequencing depth is large

273  enough to assemble large parts of fungal genomes (e.g. Olm *et al.* 2019).

274      In sum, we show that metatranscriptomics is a useful approach to identify fungal

275  species and subspecies in mixed samples. The major advantages of metatranscriptomics

276  over other HTS technologies include the selective sequencing of living organisms and the

277  ability to detect a wide range of microorganisms in one step, which has multiple

278  applications in biological research, surveillance and diagnosis. There is an increasing

279  literature reporting that virulence and antimicrobial tolerance traits vary within species

280  (Firacative *et al.* 2016; Rizzetto *et al.* 2013; Schauwvlieghe *et al.* 2017; Strope *et al.* 2015)

281  and that multiple strains or species can co-infect a host (Desnos-Ollivier *et al.* 2010; Gupta

282  *et al.* 2014; Seki *et al.* 2014; Soll *et al.* 1988; Tati *et al.* 2016). The high discriminatory power

283  obtained for closely-related lineages of *Cryptococcus* provides a good example of where

284  metatranscriptomics would be valuable in precision medicine, where therapy practices are

285 defined according to strain-specific pathogenicity and drug susceptibility traits. However, it

286 must be acknowledged that metatranscriptomics also has limitations that are common to

287 high-throughput sequencing methods, as it is susceptible to DNA/RNA extraction biases,

288 contamination and misclassifications. These limitations can be significantly minimized if

289 appropriate controls are in place (*e.g.* abundance filtering before statistical analyses).

290 Besides its application to identify well-known fungal species, metatranscriptomics can help

291 to identify novel functional roles of fungi (e.g. Gonzalez *et al.* 2018; Liao *et al.* 2014) and

292 novel species when used within a phylogenomic framework.

293

294

## Acknowledgments

303

304

## References

306 Balasundaram SV, Engh IB, Skrede I, Kauserud H (2015) How many DNA markers are

307         needed to reveal cryptic fungal species? *Fungal Biol* **119**: 940-945.

308 Banchi E, Ametrano CG, Stankovic D*, et al.* (2018) DNA metabarcoding uncovers fungal

309         diversity of mixed airborne samples in Italy. *PLoS One* **13**: e0194489.

310 Bandara H, Panduwawala CP, Samaranayake LP (2019) Biodiversity of the human oral

311         mycobiome in health and disease. *Oral Dis* **25**: 363-371.

312 Bashiardes S, Zilberman-Schapira G, Elinav E (2016) Use of Metatranscriptomics in

313         Microbiome Research. *Bioinform Biol Insights* **10**: 19-25.

314 Brown GD, Denning DW, Gow NA*, et al.* (2012) Hidden killers: human fungal infections. *Sci*

315         *Transl Med* **4**: 165rv113.

316 Carini P, Marsden PJ, Leff JW*, et al.* (2016) Relic DNA is abundant in soil and obscures

317         estimates of soil microbial diversity. *Nat Microbiol* **2**: 16242.

318  Clausen P, Aarestrup FM, Lund O (2018) Rapid and precise alignment of raw reads against
319      redundant databases with KMA. *BMC Bioinformatics* **19**: 307.
320  Desnos-Ollivier M, Patel S, Raoux-Barbot D*, et al.* (2015) Cryptococcosis Serotypes Impact
321      Outcome and Provide Evidence of *Cryptococcus neoformans* Speciation. *MBio* **6**:
322      e00311.
323  Desnos-Ollivier M, Patel S, Spaulding AR*, et al.* (2010) Mixed infections and *In Vivo*
324      evolution in the human fungal pathogen *Cryptococcus neoformans*. *MBio* **1**.
325  Dujon B (2010) Yeast evolutionary genomics. *Nature Reviews: Genetics* **11**: 512-524.
326  Enaud R, Vandenborght LE, Coron N*, et al.* (2018) The Mycobiome: A Neglected
327      Component in the Microbiota-Gut-Brain Axis. *Microorganisms* **6**.
328  Federhen S (2012) The NCBI Taxonomy database. *Nucleic Acids Research* **40**: D136-143.
329  Feinstein LM, Sul WJ, Blackwood CB (2009) Assessment of bias associated with incomplete
330      extraction of microbial DNA from soil. *Applied and Environmental Microbiology* **75**:
331      5428-5433.
332  Firacative C, Roe CC, Malik R*, et al.* (2016) MLST and Whole-Genome-Based Population
333      Analysis of *Cryptococcus gattii* VGIII Links Clinical, Veterinary and Environmental
334      Strains, and Reveals Divergent Serotype Specific Sub-populations and Distant
335      Ancestors. *PLoS Negl Trop Dis* **10**: e0004861.
336  Franzosa EA, McIver LJ, Rahnavard G*, et al.* (2018) Species-level functional profiling of
337      metagenomes and metatranscriptomes. *Nature Methods* **15**: 962-968.
338  Gonzalez E, Pitre FE, Page AP*, et al.* (2018) Trees, fungi and bacteria: tripartite
339      metatranscriptomics of a root microbiome responding to soil contamination.
340      *Microbiome* **6**: 53.
341  Grabherr MG, Haas BJ, Yassour M*, et al.* (2011) Full-length transcriptome assembly from
342      RNA-Seq data without a reference genome. *Nature Biotechnology* **29**: 644-652.
343  Gupta V, Rajagopalan N, Patil M, Shivaprasad C (2014) Aspergillus and mucormycosis
344      presenting with normal chest X-ray in an immunocompromised host. *BMJ Case Rep*
345      **2014**.
346  Hagen F, Khayhan K, Theelen B*, et al.* (2015) Recognition of seven species in the
347      *Cryptococcus gattii/Cryptococcus neoformans* species complex. *Fungal Genetics*
348      *and Biology* **78**: 16-48.
349  Hannula SE, Morrien E, de Hollander M*, et al.* (2017) Shifts in rhizosphere fungal community
350      during secondary succession following abandonment from agriculture. *ISME Journal*
351      **11**: 2294-2304.

352  Hawksworth DL, Lucking R (2017) Fungal Diversity Revisited: 2.2 to 3.8 Million Species.
353      *Microbiol Spectr* **5**.
354  Huffnagle GB, Noverr MC (2013) The emerging world of the fungal microbiome. *Trends in*
355      *Microbiology* **21**: 334-341.
356  Kassambara A (2017) ggpubr:"ggplot2" based publication ready plots. R package version
357      0.1. 6.
358  Koutsovoulos G, Kumar S, Laetsch DR*, et al.* (2016) No evidence for extensive horizontal
359      gene transfer in the genome of the tardigrade Hypsibius dujardini. *Proc Natl Acad Sci*
360      *U S A* **113**: 5053-5058.
361  Kullman B, Tamm H, Kullman K (2005) *Fungal Genome Size Database.*
362      http://www.zbi.ee/fungal-genomesize/
363  Kuske CR, Hesse CN, Challacombe JF*, et al.* (2015) Prospects and challenges for fungal
364      metatranscriptomics of complex communities. *Fungal Ecology* **14**: 133-137.
365  Kwon-Chung KJ, Bennett JE, Wickes BL*, et al.* (2017) The Case for Adopting the "Species
366      Complex" Nomenclature for the Etiologic Agents of Cryptococcosis. *mSphere* **2**.
367  Kwon-Chung KJ, Boekhout T, Fell JW, Diaz M (2002) Proposal to conserve the name
368      *Cryptococcus gattii* against *C. hondurianus* and *C. bacillisporus* (Basidiomycota,
369      Hymenomycetes, Tremellomycetidae). *Taxon* **51**: 804-806.
370  Langmead B, Salzberg SL (2012) Fast gapped-read alignment with Bowtie 2. *Nature*
371      *Methods* **9**: 357-359.
372  Li B, Dewey CN (2011) RSEM: accurate transcript quantification from RNA-Seq data with or
373      without a reference genome. *BMC Bioinformatics* **12**: 323.
374  Li X, Shen X, Chen X*, et al.* (2018) Detection of Potential Problematic Cytb Gene
375      Sequences of Fishes in GenBank. *Front Genet* **9**: 30.
376  Liao HL, Chen Y, Bruns TD*, et al.* (2014) Metatranscriptomic analysis of ectomycorrhizal
377      roots reveals genes associated with *Piloderma-Pinus* symbiosis: improved
378      methodologies for assessing gene expression in situ. *Environmental Microbiology* **16**:
379      3730-3742.
380  Loftus BJ, Fung E, Roncaglia P*, et al.* (2005) The genome of the basidiomycetous yeast and
381      human pathogen Cryptococcus neoformans. *Science* **307**: 1321-1324.
382  Marcelino VR, Verbruggen H (2016) Multi-marker metabarcoding of coral skeletons reveals
383      a rich microbiome and diverse evolutionary origins of endolithic algae. *Scientific*
384      *Reports* **6**: 31508.
385  McGuire KL, Payne SG, Palmer MI*, et al.* (2013) Digging the New York City Skyline: soil
386      fungal communities in green roofs and city parks. *PLoS One* **8**: e58020.

387    Munoz JF, Gade L, Chow NA, *et al.* (2018) Genomic insights into multidrug-resistance,
388        mating and virulence in Candida auris and related emerging species. *Nature*
389        *Communications* **9**: 5346.
390    Ngamskulrungroj P, Gilgado F, Faganello J, *et al.* (2009) Genetic diversity of the
391        *Cryptococcus* species complex suggests that *Cryptococcus gattii* deserves to have
392        varieties. *PLoS One* **4**: e5862.
393    Nilsson RH, Anslan S, Bahram M, *et al.* (2019) Mycobiome diversity: high-throughput
394        sequencing and identification of fungi. *Nature Reviews: Microbiology* **17**: 95-109.
395    Nilsson RH, Kristiansson E, Ryberg M, Hallenberg N, Larsson K-H (2008) Intraspecific ITS
396        Variability in the Kingdom Fungi as Expressed in the International Sequence
397        Databases and Its Implications for Molecular Species Identification. *Evolutionary*
398        *Bioinformatics* **4**.
399    Olm MR, West PT, Brooks B, *et al.* (2019) Genome-resolved metagenomics of eukaryotic
400        populations during early colonization of premature infants and in hospital rooms.
401        *Microbiome* **7**: 26.
402    Quince C, Walker AW, Simpson JT, Loman NJ, Segata N (2017) Shotgun metagenomics,
403        from sampling to analysis. *Nature Biotechnology* **35**: 833-844.
404    Rizzetto L, Giovannini G, Bromley M, *et al.* (2013) Strain dependent variation of immune
405        responses to *A. fumigatus*: definition of pathogenic species. *PLoS One* **8**: e56651.
406    Salter SJ, Cox MJ, Turek EM, *et al.* (2014) Reagent and laboratory contamination can
407        critically impact sequence-based microbiome analyses. *BMC Biology* **12**: 87.
408    Schauwvlieghe A, Vonk AG, Buddingh EP, *et al.* (2017) Detection of azole-susceptible and
409        azole-resistant *Aspergillus* coinfection by cyp51A PCR amplicon melting curve
410        analysis. *Journal of Antimicrobial Chemotherapy* **72**: 3047-3050.
411    Schmidt P-A, Bálint M, Greshake B, *et al.* (2013) Illumina metabarcoding of a soil fungal
412        community. *Soil Biology and Biochemistry* **65**: 128-132.
413    Schmieder R, Edwards R (2011) Quality control and preprocessing of metagenomic
414        datasets. *Bioinformatics* **27**: 863-864.
415    Sczyrba A, Hofmann P, Belmann P, *et al.* (2017) Critical Assessment of Metagenome
416        Interpretation-a benchmark of metagenomics software. *Nature Methods* **14**: 1063-
417        1071.
418    Seed PC (2014) The human mycobiome. *Cold Spring Harb Perspect Med* **5**: a019810.
419    Seki M, Ohno H, Gotoh K, *et al.* (2014) Allergic bronchopulmonary mycosis due to co-
420        infection with *Aspergillus fumigatus* and *Schizophyllum commune*. *IDCases* **1**: 5-8.

421 Shakya M, Quince C, Campbell JH*, et al.* (2013) Comparative metagenomic and rRNA
422   microbial diversity characterization using archaeal and bacterial synthetic
423   communities. *Environmental Microbiology* **15**: 1882-1899.

424 Shi M, Lin XD, Tian JH*, et al.* (2016) Redefining the invertebrate RNA virosphere. *Nature*.

425 Shi M, Neville P, Nicholson J*, et al.* (2017) High-Resolution Metatranscriptomics Reveals the
426   Ecological Dynamics of Mosquito-Associated RNA Viruses in Western Australia.
427   *Journal of Virology* **91**.

428 Soll DR, Staebell M, Langtimm C*, et al.* (1988) Multiple *Candida* strains in the course of a
429   single systemic infection. *Journal of Clinical Microbiology* **26**: 1448-1459.

430 Strong MJ, Xu G, Morici L*, et al.* (2014) Microbial contamination in next generation
431   sequencing: implications for sequence-based analysis of clinical samples. *PLoS*
432   *Pathog* **10**: e1004437.

433 Strope PK, Skelly DA, Kozmin SG*, et al.* (2015) The 100-genomes strains, an *S. cerevisiae*
434   resource that illuminates its natural phenotypic and genotypic variation and
435   emergence as an opportunistic pathogen. *Genome Research* **25**: 762-774.

436 Tati S, Davidow P, McCall A*, et al.* (2016) *Candida glabrata* Binding to *Candida albicans*
437   Hyphae Enables Its Development in Oropharyngeal Candidiasis. *PLoS Pathog* **12**:
438   e1005522.

439 Tedersoo L, Anslan S, Bahram M*, et al.* (2015) Shotgun metagenomes and multiple primer
440   pair-barcode combinations of amplicons reveal biases in metabarcoding analyses of
441   fungi. *MycoKeys* **10**: 1-43.

442 Wille M, Eden JS, Shi M*, et al.* (2018) Virus-virus interactions and host ecology are
443   associated with RNA virome structure in wild birds. *Molecular Ecology* **27**: 5263-
444   5278.

445 Yahr R, Schoch CL, Dentinger BT (2016) Scaling up discovery of hidden diversity in fungi:
446   impacts of barcoding approaches. *Philos Trans R Soc Lond B Biol Sci* **371**.

447 Yamamoto N, Dannemiller KC, Bibby K, Peccia J (2014) Identification accuracy and diversity
448   reproducibility associated with internal transcribed spacer-based fungal taxonomic
449   library preparation. *Environmental Microbiology* **16**: 2764-2776.

450 Zhang YZ, Shi M, Holmes EC (2018) Using Metagenomics to Characterize an Expanding
451   Virosphere. *Cell* **172**: 1168-1172.

452

453

454

455

456 **Tables:**

457

458 **Table 1.** Species and strains used to construct a mock fungal community for

459 metatranscriptome sequencing.

| Fungal species | Strain number |
|---|---|
| *Candida albicans* | WM 229 |
| *Candida auris* | WM 17.110 |
| *Candida glabrata* | WM 13.101 |
| *Candida dubliniensis* | WM 606 |
| *Candida orthopsilosis* | WM 03.136 |
| *Candida tropicalis* | WM 17.08 |
| *Clavispora lusitaniae* (former *Candida lusitaniae*) | WM 14.04 |
| *Cryptococcus gattii* (VGI) | WM 276 |
| *Cryptococcus neoformans* var. *grubii* (VNI) | H99 GC (H99) |
| *Cryptococcus neoformans* var. *neoformans* (VNIV) | WM 01.133 (B-3501A) |
| *Cryptococcus neoformans* var. *neoformans* (VNIV) | WM 01.127 (JEC21) |
| *Debaryomyces hansenii* | WM 36 |
| *Pichia kudriavzevii* (former *Candida krusei*) | WM 14 |
| *Pichia membranifaciens* | WM 46 |
| *Saccharomyces cerevisiae* | WM 318 |
| *Schizosaccharomyces pombe* | WM 72 |
| *Yarrowia lipolytica* | WM 63 |

460

461

462

463

464 **Table 2.** Abundance of reads (TPM) and abundance of transcripts (Depth) per fungal

465 species detected with metatranscriptomics. True members of the mock community – at

466 species level – are shown in bold.

| Species* | TPM (read-level) | Depth (transcript-level) | Relative abundance (transcript-level %) |
|---|---|---|---|
| **Cryptococcus neoformans** | **149692.28** | **11049.04** | **22.464** |
| **Candida tropicalis** | **142496.47** | **9424.30** | **19.161** |
| **Pichia kudriavzevii** | **62133.74** | **9234.05** | **18.774** |
| **Clavispora lusitaniae** | **57317.81** | **7107.80** | **14.451** |
| **Candida auris** | **13402.41** | **3354.39** | **6.820** |
| **Candida dubliniensis** | **52027.94** | **1706.41** | **3.469** |
| **Pichia membranifaciens** | **10860.75** | **1498.26** | **3.046** |
| **Candida albicans** | **42948.69** | **1441.56** | **2.931** |
| **Yarrowia lipolytica** | **52376.03** | **1384.29** | **2.814** |

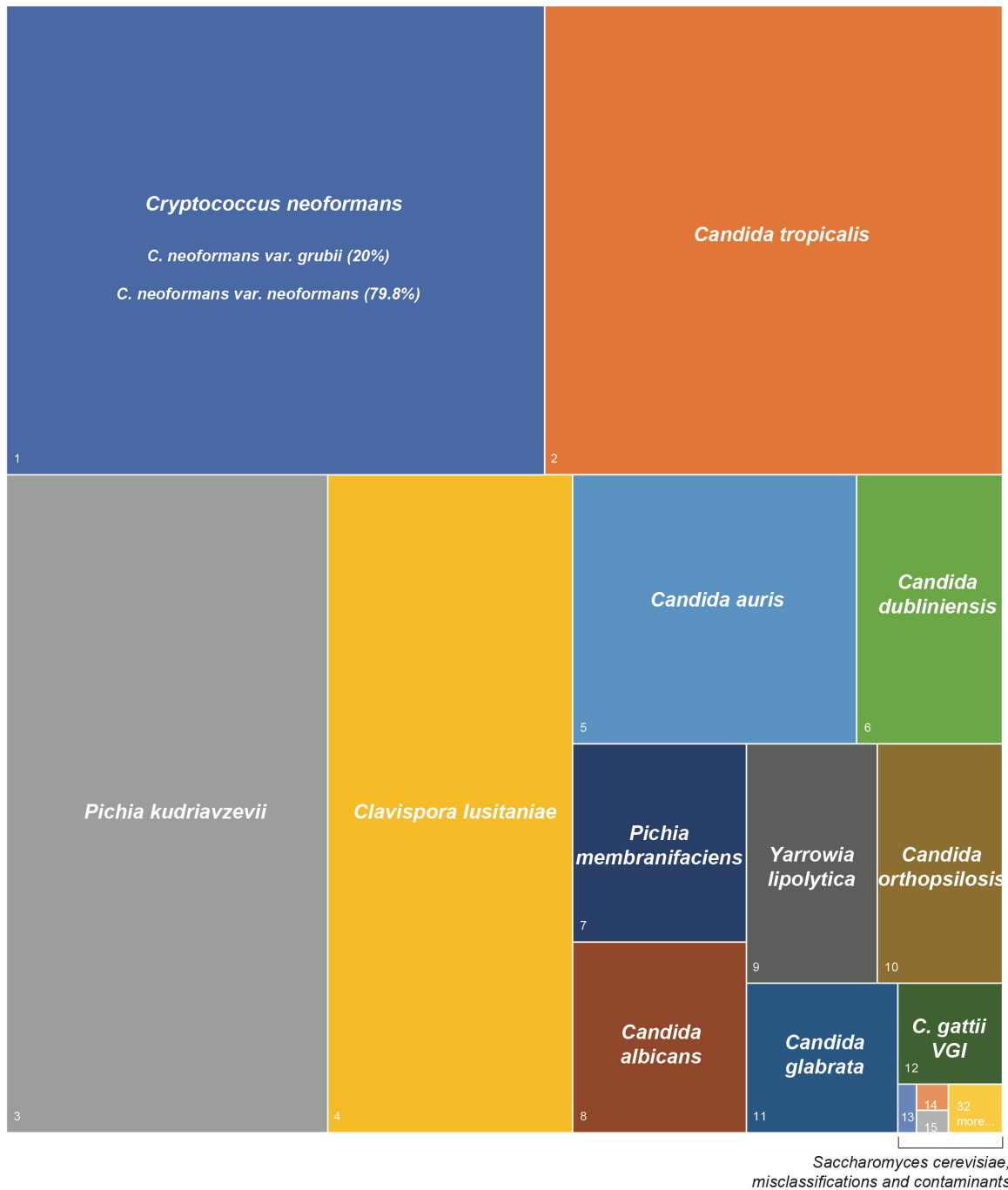| | | | |
|---|---|---|---|
| *Candida orthopsilosis* | **44531.25** | **1308.09** | **2.660** |
| *Candida glabrata* | **50404.11** | **992.55** | **2.018** |
| *Cryptococcus gattii* VGI | **9350.81** | **463.71** | **0.943** |
| *Candida glycerinogenes* | 2345.13 | 41.51 | 0.084 |
| *Nakaseomyces delphensis* | 12821.76 | 35.45 | 0.072 |
| *Candida parapsilosis* | 15989.78 | 31.11 | 0.063 |
| *Candida nivariensis* | 1411.36 | 12.72 | 0.026 |
| *Kluyveromyces marxianus* | 47.70 | 8.90 | 0.018 |
| *Torulaspora delbrueckii* | 15.01 | 6.00 | 0.012 |
| *Kluyveromyces lactis* | 9.24 | 3.93 | 0.008 |
| *Saccharomyces cerevisiae* | **22.84** | **3.77** | **0.008** |
| *Eremothecium sinecaudum* | 25.93 | 3.67 | 0.008 |
| *Pichia cecembensis* | 751.23 | 3.60 | 0.007 |
| *Lodderomyces elongisporus* | 34.30 | 3.59 | 0.007 |
| Uncultured *Candida* | 885.35 | 3.13 | 0.006 |
| *Eremothecium gossypii* | 10.13 | 3.04 | 0.006 |
| *Naumovozyma dairenensis* | 17.37 | 2.80 | 0.006 |
| *Suhomyces tanzawaensis* | 30.87 | 2.25 | 0.005 |
| Dipodascaceae sp. LM136 | 24286.11 | 2.16 | 0.004 |
| *Cyberlindnera jadinii* | 12.32 | 2.04 | 0.004 |
| *Metschnikowia bicuspidata* | 16.02 | 1.29 | 0.003 |
| *Brettanomyces naardenensis* | 96.13 | 1.19 | 0.002 |
| *Pichia norvegensis* | 783.06 | 1.11 | 0.002 |
| *Debaryomyces fabryi* | 22.87 | 0.92 | 0.002 |
| *Candida neerlandica* | 487.65 | 0.69 | 0.001 |
| *Melanotaenium endogenum* | 262.12 | 0.59 | 0.001 |
| *Pichia kluyveri* | 51814.00 | 0.55 | 0.001 |
| *Candida pseudohaemulonis* | 560.14 | 0.49 | 0.001 |
| *Candida* sp. (in: *Saccharomycetales*) | 330.15 | 0.49 | 0.001 |
| *Pichia* sp. 2 TMS-2011 | 0.00 | 0.45 | 0.001 |
| *Cryptococcus neoformans* AD hybrid | 0.00 | 0.44 | 0.001 |
| Saccharomycetales sp. LM594 | 2.60 | 0.30 | 0.001 |
| *Naumovozyma castellii* | 4.56 | 0.29 | 0.001 |
| *Saccharomyces pastorianus* | 0.81 | 0.26 | 0.001 |
| *Cryptococcus gattii* VGIII | 0.48 | 0.21 | 0.000 |
| Other Eukaryotes | 936.75 | 18.83 | 0.038 |
| Bacteria | 423.36 | 15.79 | 0.032 |
| Unclassified | 198000.57 | 8.00 | 0.016 |
| **TOTAL** | **1000000** | **49186** | **100** |

467     * Species defined according to the NCBI taxonomy database. Strain numbers may indicate vouchers

468     rather than genetically different lineages.

469

470

471    **Figures**

472



*Saccharomyces cerevisiae,
misclassifications and contaminants*

473

474    **Figure 1.** Relative abundance of transcripts assigned to microbial species recovered in the

475    metatranscriptome of a mock community. See Table 2 for a full list of species and more

476    details about their abundance.

477

478

479

480

481
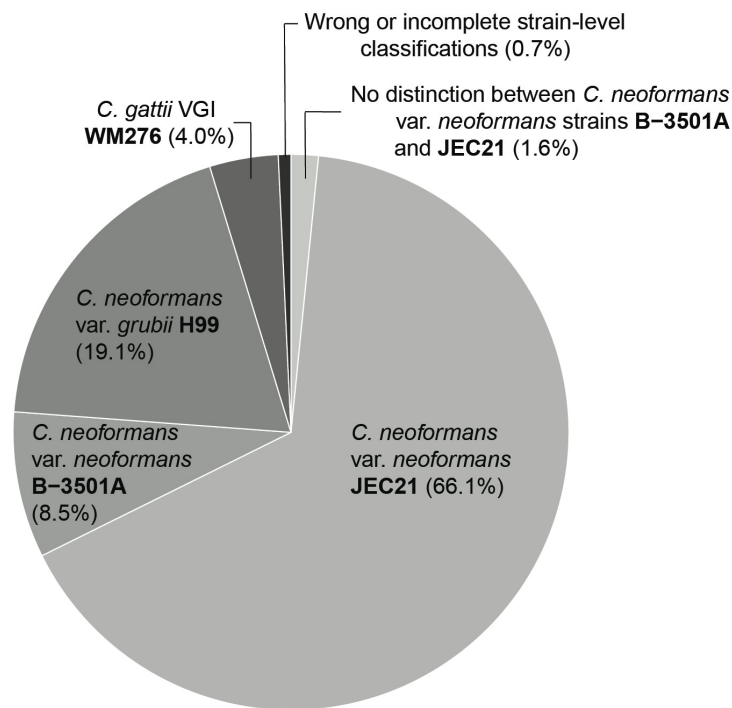
482

483

484

485

486

487

488

489



490 **Figure 2.** Strain-level classifications of taxa within the *Cryptococcus neoformans* and *C.*

491 *gattii* species complexes.

492

493

494

495

496

497

498

499

500

501

502

503