# A new lineage of *Cryptococcus gattii* (VGV) discovered in the Central Zambezi Miombo Woodlands

Rhys A. Farrer[1,2,3,4], Miwha Chang[5], Michael J. Davis[5], Lucy van Dorp[3], Dong-Hoon Yang[5], Terrance Shea[4], Thomas R. Sewell[6], Wieland Meyer[7], Francois Balloux[3], Hannah M. Edwards[6], Duncan Chanda[8], Geoffrey Kwenda[9], Mathieu Vanhove[6], Yun C. Chang[5], Christina A. Cuomo[4], Matthew C. Fisher[6], Kyung J. Kwon-Chung[5]

[1]Medical Research Council Centre for Medical Mycology, University of Exeter, Exeter, EX4 4PY, UK.
[2]Medical Research Council Centre for Medical Mycology, University of Aberdeen, Aberdeen, AB242TN, UK.
[3]UCL Genetics Institute, University College London, Gower Street, London, WC1E 6BT, UK.
[4]Broad Institute of MIT and Harvard, Cambridge, Massachusetts, United States of America.
[5]Molecular Microbiology Section, Laboratory of Clinical Immunology and Microbiology, National Institute of Allergy and Infectious Diseases, National Institute of Health, Bethesda, Maryland, United States of America.
[6]MRC Centre for Global Infectious Disease Analysis, Imperial College London, London, United Kingdom
[7]Molecular Mycology Research Laboratory, Centre for Infectious Diseases and Microbiology, Faculty of Medicine and Health, Sydney Medical School, Westmead Clinical School, Marie Bashir Institute for Infectious Diseases and Biosecurity, The University of Sydney, Westmead Hospital (Research and Education Network), Westmead Institute for Medical Research, Sydney, NSW, Australia
[8]Adult centre of Excellence, University Teaching Hospital, Lusaka, Zambia
[9]Department of Biomedical Sciences, School of Health Sciences, University of Zambia, Lusaka, Zambia

**Corresponding author:** Kyung J. Kwon-Chung (jkchung@niaid.nih.gov)

## Abstract

We discovered a new lineage of the globally important fungal pathogen *Cryptococcus gattii*, based on six isolates collected from three locations spanning the Central Miombo Woodlands of Zambia, Africa. All isolates were from environments (middens and tree holes) that are associated with a small mammal, the African hyrax. Phylogenetic and population genetic analyses confirmed that these isolates form a distinct, deeply divergent lineage, which we name VGV. VGV comprises two subclades (A and B) that are capable of causing mild lung infection with negligible neurotropism in mice. Comparing the VGV genome to previously identified lineages of *C. gattii* revealed a unique suite of genes together with gene loss and inversion events. However, standard *URA5* RFLP analysis could not distinguish between VGV and VGIV isolates. We therefore developed a new *URA5* RFLP method that can reliably distinguish the newly described lineage. Our work highlights how sampling understudied ecological regions alongside genomic and functional characterisation can broaden our understanding of the evolution and ecology of major global pathogens*.

## Importance

*Cryptococcus gattii* is an environmental pathogen that causes severe systemic infection in immunocompetent individuals more often than in immunocompromised humans. Over the past two decades, researchers have shown *C. gattii* falls within four genetically distinct major lineages. By combining field work from an understudied ecological region (the Central Miombo Woodlands of Zambia, Africa), genome sequencing and assemblies, phylogenetic and population genetic analyses, and phenotypic characterization (morphology, histopathological, drug-sensitivity, survival experiments) we discovered a hither to unknown lineage which we name VGV (variety *gattii* five). The discovery of a new lineage from an under studied

63     ecological region has far-reaching implications for the study and understanding of

64     fungal pathogens and diseases they cause.

65     **Introduction**

66

67     Cryptococcosis is a severe fungal infection responsible for high levels of mortality

68     and morbidity worldwide(1). The etiological agents are two species complexes of the

69     genus *Cryptococcus*: *C. neoformans* and *C. gattii.* Whilst, the first described cases of

70     clinical cryptococcosis due to these two distinct species complexes were reported in

71     the mid-1890s under the names *Saccharomyces hominis*(2) and *S. subcutaneous*

72     *tumefacience*(3) respectively, clinical *Cryptococcus* isolates have been taxonomically

73     treated as a single species (*C. neoformans*) for more than 100 years(4).

74     Heterogeneity among cryptococcosis-causing yeast isolates became increasingly

75     apparent from the middle of the $20^{th}$ century onward, and led to the recognition of

76     four serotypes (A, B, C, D) based on the antigenic determinant of capsular

77     polysaccharide(5, 6). Subsequent discovery of two distinct sexual cycles produced

78     by the isolates of A/D vs. B/C serotypes(7, 8) and phylogenetic analysis using

79     various gene sequences(9–11) confirmed these complexes to be genetically

80     divergent enough to be considered as separate species. Thus, in 2002, the isolates

81     of serotype B/C were formally classified as *C. gattii*(12) while *C. neoformans*

82     includes all serotype A/D strains(13).

83

84     Over the past two decades, population structure analysis of the two species using

85     molecular typing methods such as PCR fingerprinting(14), AFLP analysis(15) and

86     multi-locus sequencing(16) has demonstrated that both species contain genetically

87     diverse lineages that qualify them to be considered as two species complexes, which

88     have been further subdivided into numerous molecular types(14, 17). To date, four

89     major lineages are recognised for *C. gattii*, which are denoted VGI/AFLP4,

90     VGII/AFLP6, VGIII/AFLP5, and VGIV/AFLP7. Recently a fifth genotype was

91     described on the basis of a single strain but with several different designations

92     including Clade B (based on MLST), VGIIIc/VGIV, and *C. decagattii*(17, 18). It has

93     been proposed to elevate these five lineages to separate species(17). However,

94     such taxonomic treatment is currently controversial mainly due to the lack of clear

95     biological differences between the lineages and no clear consensus on the limits and

96    numbers of the putative species boundaries. As such, the various *C. gattii* lineages

97    are collectively considered as the '*C. gattii* species complex' (18).

98

99    In this paper, we describe the discovery of a new lineage/molecular type within the

100    *C. gattii* species complex, which we designate as VGV. The six VGV isolates were

101    identified among 32 *C. gattii* isolates recovered from soil, animal dung and tree bark

102    samples collected in Zambia by Vanhove *et al* in 2013 (19). In this paper, we

103    characterize genomic and phenotypic features of the VGV molecular type.

104    Additionally, we present a new improved genome assembly and gene-sets for *C.*

105    *decagattii* (17) which we confirmed for the first time to be a separate lineage and

106    therefore name as VGVI for consistency with the other lineages.

107

108    **Results**

109

110    **Comparative and population genomics for the six lineages of *C. gattii***

111

112    We discovered a new lineage of *Cryptococcus gattii* from a panel of 32 (out of 55)

113    genome sequenced isolates recovered from Southern tree hyrax (*Dendrohyrax*

114    *arboreus*) middens, midden soil or tree holes from the Central Zambezian Miombo

115    Woodland ecoregion, a densely forested ecoregion that covers much of Central and

116    East Africa (**Fig. 1, Table 1, Table S1**). Isolates from the new lineage, which we

117    have named VGV, were collected from a 430 km span of northern Zambia including

118    the Mupata Hills (Luanshya, Copperbelt Province), Mutinondo wilderness area and

119    Kapishya (Mpika, Northern Province), suggesting that the lineage has a broad

120    regional distribution across this ecoregion (**Fig. 1a**). All VGV isolates were identified

121    as serotype B, which also encompasses strains from VGI, VGII, the VGIIIa

122    subgroup, and rare isolates among VGIV.

123

124    Phylogenetic analyses demonstrates that VGV, VGVI (*C. decagattii*), and the four

125    previously described lineages, are genetically distinct from each other (**Fig. 2**).

126    Indeed, none of the *C. gattii* lineages appear to be the result of hybridisation based

127    on the distribution of private alleles (**Fig. 3a-b**), maximum likelihood phylogenetic

128    reconstruction (**Fig. 2a**), $F_{ST}$ (**Fig. S4**) or NeighborNet Network (**Fig. 2b**). Additional

129    population genetics analyses confirmed low levels of genetic exchange between the

130    six well resolved *C. gattii* lineages. For example, Principal Component Analysis

131    (PCA) resolved distinct grouping for the lineages, with the first component (PC1)

132    separating VGII from all other lineages, forming distinct clusters for VGIII and VGVI

133    on PC2 (**Fig. 2c**). The projection of PC3 and PC4 further allows identication of

134    distinct tight clusters for the VGI, VGIV and VGV lineages (**Fig. 2d**).

135

136    The new VGV lineage is represented by six isolates falling within two distinct

137    subclades (A and B). Clade A comprises three VGV isolates (MF5, MF13, MF54)

138    that were recovered from soil and animal dung sampled in hyrax middens, from

139    which we also identified VGI and VGII isolates (**Fig. 1a**). Clade B comprises a further

140    three VGV isolates: two that were recovered approximately 345 km away from Clade

141    A (MF34 and MF51), and a third (MF56) that was recovered approximately 430 km

142    away from the other Clade B isolates. Clade B isolates were recovered from both a

143    tree hole and also hyrax middens, showing that the lineage can occupy both tree and

144    dung, environments that are both associated with hyrax activity. The fact that Clade

145    A and B were found in different geographic locations might reflect a degree of spatial

146    genetic structure within VGV. All the VGV isolates were located in regions of granite

147    and acidic kopjes/hills that are found occurring patchily across this ecoregion.

148

149    *C. gattii* VGV is highly diverged from all previously recognised *C. gattii* lineages.

150    VGV isolates differ from VGII (reference isolate R265) by ~0.75 million SNPs on

151    average (44 SNPs/Kb), and are thus similarly distant from VGII as the other lineages

152    (**Table S1**). The analysis of the relative proportion of shared private alleles for the *C.*

153    *gattii* lineages (**Fig. 3a-b**), indicates VGII shared the fewest alleles with any of the

154    other lineages, reflecting its greater divergence (<4.6 Kb total; <0.2 SNPs per Kb;

155    **Fig. 2, Fig. 3a**). The newly discovered VGV shared fewer alleles with VGVI (0.24 per

156    Kb) and VGIII (0.29 per Kb), than with VGI (1.18 per Kb) and VGIV (2.53 per Kb).

157    The lineages that shared the most private alleles were VGVI and its closest relative

158    VGIII (92 Kb total; 5.37 per Kb), which account for an average of 12% of all SNPs

159    (based on alignments to VGII) found in isolates from those lineages.

160

161    Nearly one in ten nucleotides in the *C. gattii* genome has an alternative allele across

162    the six lineages ($1.55 \times 10^6$ sites; 9.01% of the *C. gattii* genome). Indeed, >180 Kb of

163  these unique/private alleles were identified for each lineage, including for VGV which

164  had 220 Kb private alleles (12.75 per Kb) (**Fig. 3b**). VGI is the most distinct in terms

165  of both the highest count of private alleles (378 kb / 21.93 SNPs per Kb) and its

166  nucleotide diversity (π) (**Fig. 3c**), which is reflected in the three distinct subclades of

167  VGI isolates in the whole genome phylogeny (**Fig. 2a-b**). Conversely, the three VGVI

168  isolates are thought to be derived from a single isolate recovered from a patient in

169  Mexico and subsequently distributed to different labs where they have been

170  renamed and sequenced(14, 17, 20, 21). Its few clonal differences are illustrated by

171  its low nucleotide diversity (π) (**Fig. 3c**).

172

173  Unsupervised model-based clustering identified highly structured ancestry

174  components enriched in each of the lineages. The clustering solution with the lowest

175  cross-validation error (K=9) grouped the VGV isolates into a single genetically

176  homogenous group (**Fig3d-e, Fig. S1**) whilst identifying four unique components

177  within the VGII lineage. Of these, subclades VGIIx  and VGIIb share small

178  proportions of ancestry with other defined VGII subclades. For example, VGIIb is

179  inferred to share ancestry with other VGII subclades (isolates Ram5 and B8554) and

180  other lineages (B7394 has alleles from VGIV, and B7735 has alleles from VGV).

181  Conversely, none of the isolates in VGIIa and VGIIc have demonstrable admixture

182  with other subclades or lineages, both being formed by single unique ancestry

183  components. VGIII isolate B8212 (a clinical isolate from Oregon, USA in 2007(22)) is

184  also modelled as sharing ancestry with VGVI.

185

186  Finer-scale clustering was performed by considering patterns of genome-wide

187  haplotype sharing in fineSTRUCTURE(23). Here, VGV isolates forming a separate

188  cluster with greater haplotype similarity to isolates from VGI, VGIII and VGIV than

189  VGII (**Fig. S2**). While haplotype sharing patterns were overwhelmingly in accordance

190  with each lineage being genetically distinct, a notable exception was VGIII isolate

191  B8212 that shares haplotypes with VGIV and VGVI (also in accordance with model-

192  based clustering), perhaps owing to a small amount of genetic exchange with one or

193  both of those lineages. As also observed using ADMIXTURE based clustering (**Fig.**

194  **3e**), two isolates from VGII, B7394 and B7735, were also genetically distinct and

195  were assigned to their own cluster which was most closely related to isolates from

196  subclade VGIIb.

197

198    All six VGV isolates were haploid with no evidence for aneuploidy based on allele-

199    frequencies and depth of coverage (**Fig. S3**). However, we did observe copy number

200    variation (CNV) between the three VGVI isolates derived from a single clinical isolate

201    from Mexico(14, 17, 20, 21). Specifically, isolate CBS11687 acquired a ~200Kb

202    duplication of supercontig (sc) 5 (position 1,040,000 through the end of the

203    supercontig). Separately, isolate WM1804 had a smaller 40kb duplication on sc21

204    (positions 150,000 – 190,000). Isolate WM1802 had neither CNV. In terms of base

205    changes, the three VGVI isolates (WM1802, WM1804 and CBS11687) differed by

206    only 419 SNPs, with the fewest found between WM1804 and CBS11687 ($n$=126)

207    and the most found between WM1802 and CBS11687 ($n$=315). These genetic

208    differences may have occurred as a result of micro-evolution during or following

209    passaging or cryo-preserving, although large CNVs are common in *C. gattii*(24, 25).

210    All of the newly isolated VGI ($n$=7) and VGVA ($n$=3) samples from Zambia had a

211    small <10 kb duplication within supercontig 6 of the R265 genome (position 400 kb

212    to 410 kb). This genomic region encodes a single 87aa protein that is conserved in

213    *C. neoformans* and *C. gattii*, but has no functional annotation (PFAM, GO-terms,

214    KEGG-EC, TMHMM or SigP4).

215

216    The results from our phylogenetic and population genetic analyses are in line with

217    previous work(26), indicating that lineages within the *C. gatti* species complex have

218    remained largely genetically isolated since their divergence. Pairwise-lineage

219    calculations of θ, Weir's formulation of Wright's fixation index ($F_{ST}$) suggest very low

220    levels of genetic exchange between each of the lineages (**Fig. S4**) which is also

221    reflected in analyses of genetic structure (**Fig. 2-3, Fig. S1-S2**). Both depth of

222    coverage plots and $F_{ST}$ non-overlapping sliding 10 Kb window plots across the

223    mating type locus (*MAT*) at the start of supercontig 18 demonstrate that all VGIV and

224    VGV isolates included in this study are *MAT*α (the reference genome of R265 is

225    *MAT*α; high depth of coverage and θ > 0.98 across the *MAT* loci). In contrast, for

226    VGI, VGII, VGIII and VGVI *MAT***a** isolates were included in our panel.

227

228    **Genome assembly and analysis of VGV reveals _C. gattii_ lineage-specific**

229    **differences**

230

231  We assembled and annotated a near complete genome for the newly discovered

232  lineage *C. gattii* VGV (isolate MF34) using both Oxford Nanopore and Illumina

233  sequencing reads. The resultant assembly consisted of 15 contigs corresponding to

234  the 14 chromosomes; the single break in one chromosome corresponds to the

235  ribosomal (rDNA) region. Other than under-representing rDNA genes, this assembly

236  provides a complete representation of the genome, with telomeric repeats

237  (TTAGGG) present at 28 contigs ends. Gene annotation revealed 6,322 predicted

238  protein coding genes, which is similar to the seven other representative *C. gattii*

239  isolates with publicly-available complete genomes(26, 27) representing the four

240  previously known major lineages (ranging from 6,092 to 6,763), as well as *C.*

241  *neoformans* H99(28) ($n$=6,962) (**Fig. S5-S6**).

242

243  To establish the evolution of protein coding genes in *C. gattii*, we compared the gene

244  content for two representative annotated genomes per lineage where possible (no

245  second annotated reference genomes were available for VGIV, VGV and VGVI),

246  identifying 4,565 single copy core orthologs that are shared amongst the five

247  lineages of *C. gattii* and *C. neoformans* (~74% of *Cryptococcus* genes) (**Table 2**).

248  For VGVI, we sequenced and assembled the WM1802 isolate obtaining a similar

249  genome length (17.42 Mb) and protein coding gene count ($n$=6,092). For VGII, we

250  included the updated VGII R265 PacBio assembly in our panel of genomes(29)

251  (**Table 2**). Orthology detection between just the two R265 assemblies identified 91%

252  of genes in 1:1 orthology ($n$=5,642), ~4% of genes unique to the new assembly

253  ($n$=252) and ~6% of genes in paralogous clusters ($n$=364). The previous VGII R265

254  assembly had 635 genes that were not called in the new assembly, likely a

255  difference in the annotation protocol. Analysis of Core Eukaryotic Genes (CEGs) and

256  BUSCO revealed a high completeness of gene-sets, and an increased completeness

257  in the new annotation (**Fig. S6**). Furthermore, all assemblies generated using long

258  read sequencing technology assemble into 14 scaffolds/supercontigs, suggesting all

259  *Cryptococcus* lineages/species have conserved numbers of chromosomes.

260

261  Ortholog amino acid differences within and between lineages were consistent with

262  results from our phylogenetic and population genetic analyses. VGV MF34 had the

263  highest amino acid sequence similarity to VGIV IND107 (53,000 amino-acid

264  differences = 97.55% similarity), which is observed in both alignment-based and

265  ortholog-based phylogenies (**Fig. 2, Fig. 4**). The most similar inter-lineage orthologs

266  were between VGIII and VGVI (49,500 predicted amino-acid differences = 97.71%

267  similarity) (**Table S1**), while the most distinct pairwise comparisons were between *C.*

268  *gattii* and *C. neoformans* (between 205,000 and 218,000 amino acid changes; ~90%

269  protein similarity).

270

271  Overall, synteny is conserved within *C. gattii*(26), though with notable differences

272  between some lineages. For example, VGV has a single 171 Kb inversion on

273  supercontig 7 (positions 544,906-716,249) compared with the middle of VGIV

274  IND107 supercontig 7 and the middle of VGIII CA1280 supercontig 5 (**Fig. 4**). VGVI

275  also has some syntenic differences between its closest relative VGIII (**Fig. 2, Fig. 3,**

276  **Fig. 4**). For example, approximately half of VGVI supercontig 5 is syntenic for the

277  start of VGIII (CA1873) supercontig 16, while the second half of VGVI supercontig 5

278  is syntenic for a middle region of VGIII supercontig 1, indicative of a chromosomal

279  translocation. Further improvements and additional genome assemblies should

280  establish the full number and genetic impact of lineage-specific genomic

281  rearrangements.

282

283  Lineage specific genes and multi-lineage specific genes (found in two or more

284  lineages) were identified in each of the lineages (**Fig. S7, Fig. 3f**). Many of these

285  lineage-specific genes (223/605; 37%) were previously identified from a panel of

286  genome assemblies without the addition of VGV and VGVI(26). A further 53/605

287  (9%) of newly detected lineage-specific genes were previously categorised as multi-

288  lineage-specific genes. Lineage specific genes in newly sequenced lineages (VGV

289  and VGVI) include 74 genes that were unique to VGV and 49 genes that were

290  uniquely absent in VGV. Genes unique to VGV include two sugar transporters

291  (D1P53_002216, D1P53_002944) an alcohol dehydrogenase (D1P53_004471), and

292  an aldehyde dehydrogenase (D1P53_006242). Conversely, eight transmembrane

293  proteins and a single uncharacterised secreted protein were uniquely missing in

294  VGV. All of the genes involved in the ergosterol biosynthesis pathway were present

295  in single-copy in VGV. The VGVI WM1802 genome encodes 80 genes that are

296  unique and 127 genes that are uniquely absent. Among the unique genes in

297  WM1802, 14 are predicted to be involved in transport and include three

298  monosaccharide transporters, one hexose transporter, one cadmium ion transporter

299  and one monocarboxylic acid transporter.

300

301  Predictably, *C. neoformans* VNI H99 had the greatest number of lineage-

302  specific/absent genes, with 578 unique genes and 28 absent genes. These included

303  47 genes predicted to be transmembrane proteins (including four sugar transporters,

304  five MFS transporters, and a caffeine resistance transporter), and 34 secreted

305  proteins. Fewer genes were uniquely absent in *C. neoformans*, which included the

306  eukaryotic translation initiation factor 3 subunit B, an ACC oxidase, a copper amine,

307  an allantoate permease of the major facilitator, and a 3-hydroxyacyl-dehydrogenase

308  with oxidoreductase activity. Full details of all lineage-specific genes are provided in

309  **Table S1**.

310

311  **Phenotypic characteristics of VGV**

312

313  All six isolates that belonged to the new VGV lineage based on whole genome

314  sequencing (**Table 1**) were first incorrectly identified as VGIV based on the *URA5*

315  restriction fragment length polymorphism (RFLP) banding pattern. Unlike most of

316  VGIV, all six VGV isolates were serotype B. These were further characterised as

317  *MAT*α, melanin and urease positive (**Fig. S8**) and grew well at 37°C (**Fig. 5A**).

318  However, VGV strains grew slightly more slowly on CGB agar (**Fig. S8**) than the

319  serotype C VGIV strain (MF46) isolated from the same environment in the Central

320  Miombo Woodland (**Table 1**). Hence, the positive blue/green colour development on

321  CGB agar took longer in VGV than the control strains (**Fig. S8**). Since VGV is

322  genetically closest to VGIV (**Fig. 2a-b**, **Fig. S2**), two Serotype C VGIV strains were

323  used as control isolates for further phenotypic comparisons (WM779 isolated from a

324  cheetah in South Africa(16), and MF46 isolated from Miombo tree bark in Zambia

325  near Hyrax middens) (**Table 1**).

326

327  The size and morphology of VGV yeast cells were typical for *Cryptococcus* and

328  indistinguishable from the control strains (**Fig. 5B**). Two distinct patterns of capsule

329  formation were found among the six VGV isolates grown in YEPD broth (**Fig. 5B**).

330  The isolates recovered from soil, Clade A (MF5, MF13, MF54), produced thinner

331   capsule ($\leq$ 1$\mu$m) compared to those recovered from tree bark, Clade B (MF34,

332   MF51, MF56), which produced thick (2-4 $\mu$m) capsules similar to the VGIV control

333   strains.

334

335   The VGV isolates and the control strains of VGIV manifested unusually high

336   resistance toward fluconazole (FLC), particularly given they were sampled from an

337   environmental niche. The three isolates of VGV clade A were more resistant to FLC,

338   with minimal inhibitory concentration (MIC) $\geq$ 128 $\mu$g/ml than the three isolates in

339   clade B which showed MICs of 24-64 $\mu$g/ml. All six VGV isolates showed MIC of

340   0.0625 $\mu$g/ml for 5-fluorocytocine (5-FC) similar to WM779. The MIC of MF46 for 5-

341   FC was unusually high, 4$\mu$g/ml. The VGV MIC of amphotericin B ranged between

342   0.5 to 1$\mu$g/ml, higher than the control strains which had MIC below 0.5 $\mu$g/ml (**Fig.**

343   **5C**).

344

345   To explore the relative pathogenicity amongst VGV subclades we selected two

346   isolates from Clade A (MF5, MF13) and two isolates from Clade B (MF34, MF51) for

347   inoculation in mice models. Mice infected by all four isolates survived for 70 days

348   while all the mice infected by WM779, a virulent serotype C control isolate,

349   succumbed to infection within 30 days post infection (**Fig. 5D**). The VGIV

350   environmental isolate MF46 (serotype C) caused no death in the mouse model.

351   Fungal loads in the lungs of VGV infected mice were substantially lower than that of

352   WM779 and slightly lower than those infected by MF46. Brain fungal loads of mice

353   infected with the VGV strains were also low to negligible. The control isolates of

354   VGIV showed little neurotropism (**Fig. 5E**). Histopathological analysis of the lungs

355   demonstrated significant pathology in WM779 infected mice, with yeast found

356   throughout the lung together with notable disruption of lung tissue. In many locations

357   extensive leukocyte recruitment was evident in areas of concentrated infection (**Fig.**

358   **6**).

359

360   Histopathological analysis of the VGV isolates displayed substantially lower

361   pulmonary yeast levels. That said, mice infected by MF46, MF34 and MF51 had

362   higher levels of yeast than those infected by MF5 and MF13 in both CFU (**Fig. 5E**)

363   and histopathological analysis (**Fig. 6**) in which MF51 was shown to represent Clade

11

364   B. The lung histopathology of the mice infected by MF34 was similar to that of MF51

365   (data not shown). Notable for its absence, leukocyte infiltration was mostly low or

366   absent from sites of VGV infection. MF13 showed some leukocyte infiltration to a

367   subset of infectious foci (**Fig. 6**).

368

369   **<u>Identification of VGV by *URA5* RFLP</u>**

370

371   The patterns of *URA5* RFLP, obtained by double digestion with *Sau*961 and *Hha*1,

372   has been widely used to identify the lineage/molecular type in both *C. neoformans*

373   and *C. gattii* species complexes(14). The *URA5* RFLP of Clade A isolates obtained

374   by *Sau*96I/*Hha*I digest showed identical pattern with that of VGIV (**Fig. S9A**). Those

375   of Clade B, however, produced an additional 1.3 kb amplicon which was absent from

376   Clade A or any other VG molecular type reference isolates. This 1.3 kb amplicon

377   was present even in uncut DNA of Clade B isolates (**Fig. S9A-B**), but its nature is

378   not known at this juncture. Since the *URA5* RFLP patterns of both VGIV and VGV

379   were not conclusively different, we compared the *URA5* gene sequence of the VGV

380   isolate MF34 to that of the VGIV reference strain WM779 to identify possible

381   restriction enzymes that can clearly distinguish the two lineages. This led to us

382   identifying two highly discriminatory restriction enzymes; *Stu*I and *Ear*I. The expected

383   sizes (bp) of the *URA5* gene fragment resulting from *Stu*I digestion are: 221 bp, 237

384   bp and 322 bp in VGIV and 237 bp and 543 bp in VGV. The *Ear*I digestions

385   produced 247 and 501bp fragments in VGIV and 247 and 300 bp fragments in VGV.

386   We compared the *URA5* RFLP of 17 VGIV isolates (**Table S1**) with the 6 VGV

387   isolates by *Stu1* or *Ear*I digestion and the results are shown in **Fig. S9C-D**.

388

389   **<u>Type strain of VGV</u>**

390

391   We have designated MF34 (Clade B isolate) as the type strain of VGV which was

392   isolated from a tree hole located in Mutinondo (latitude -12.45, longitude 31.29),

393   Central Zambezian Miombo Woodlands (**Table 1**). Its genome has been assembled

394   and annotated to near completion (15 scaffolds with N50=1.3Mb and telomeric

395   repeats at 28 contigs ends). MF34 is serotype B and *MAT*$\alpha$ and causes mild

396   pneumonia in C57BL/6 mice with neglible neurotropism. The genome assembly has

397 been submitted to NCBI under the project accession PRJNA487802 and the culture
398 has been deposited at the American Type Culture Collection (accession number
399 pending).
400
401 **Discussion**
402
403 Over the past decade, increased sampling world-wide alongside whole-genome
404 sequencing (WGS) methods have uncovered a greater genetic diversity of important
405 pathogens including the *C. neoformans* and *C. gattii* species complexes. For
406 example, sampling from Botswana revealed the existence of the *C. neoformans* VNB
407 lineage (30), which itself has recently been shown to be deeply split into two
408 genetically isolated lineages, VNBI and VNBII (31). Thus far, VGVI is the only
409 lineage that exists as a single genotype since the three isolates previously
410 designated as *C. decagattii* appear to have been isolated from the same patient (21).
411 Each of the previously identified lineages of *Cryptococcus* have recently been
412 designated as separate taxonomic species based on phylogenetic species
413 recognition criteria (17). While we agree that *Cryptococcus* contains a number of
414 genetically diverse and monophylectic clades that may be viewed as species under
415 an Evolutionary Species Concept (32), we have previously argued that it is
416 premature to give each clade a separate taxonomic name at this juncture (18, 33).
417 One notable concern raised by Kwon-Chung *et al.* (18) was that the proposed seven-
418 species taxonomy (33) was likely to be unstable due to incomplete knowledge of the
419 true extent of *Cryptococcus* diversity worldwide. Our discovery of *C. gattii* VGV from
420 the Miombo woodlands of Zambia clearly shows that we have not yet achieved a full
421 understanding of the global biodiversity of *Cryptococcus*, and that further exploration
422 will likely yield more phylogenetic species. Until we have a more accurate consensus
423 on the true numbers of *Cryptococcus* lineages, we propose that the names 'VN' and
424 'VG' serve as a practical 'zip-code' within *C. neoformans* and *C. gattii*, offering a
425 convenient way to describe newly discovered lineages or recombinants without
426 introducing unwanted nomenclatural instability and confusion.
427
428 Our discovery of *C. gattii* VGV from hyrax-associated environments suggests an
429 association with these mammals. Hyrax are small herbivores that are most closely

13

430    related to elephants (Proboscidea) and sea cows (Sirenia), and are characterised by

431    the behaviour of defecating in communal latrines, usually located in crevices in rocky

432    kopjes, over many generations (34). These locations are often sheltered in rocky

433    caves and droppings are likely to accumulate for upwards of 50,000 years, in some

434    cases forming a stable paleoenvironmental hotspot of urea-rich nitrogenous material

435    (35). *Cryptococcus* has a pronounced trophism for urea as a nutritive substrate, and

436    pigeon guano is known to support prolific growth of *C. neoformans* and (to a lesser

437    extent) *C. gattii* (36). Our finding that hyrax middens are hotspots of *Cryptococcus*

438    diversity suggests that their ecological stability in landscapes that are low in nitrogen

439    availability may lead to them being important arenas for the evolution of

440    *Cryptococcus*, and will likely be fertile ground for further discovery of diversity within

441    this genus.

442

443    Fungal association with small mammals may suggest adaptations that confer

444    pathogenicity, known as the 'endozoan, small-mammal reservoir hypothesis' (37),

445    and deserves to be explored further following our findings of an association of

446    *Cryptococcus* with hyrax. Accordingly, alongside further study of potential

447    mammalian reservoirs, the search for VGV clinical isolates is also needed in order to

448    understand the true virulence potential of VGV and whether it can spillover into

449    humans. Murine models have shown that environmental isolates are less virulent

450    than clinical isolates of the same molecular type in both the *C. gattii* and *C.*

451    *neoformans* species complexes suggesting that polymorphic virulence factors

452    exist(38, 39). However, despite its large genetic distance from all other lineages, the

453    new VGV lineage is not clearly distinguishable from others by existing methods such

454    as serotyping or the routinely used *Sau*961 and *Hha*1 digested *URA5* RFLP

455    analysis(14). Thus, it is possible that previous isolates belonging to VGV may have

456    been misidentified using non-WGS methods. The most likely candidates for the

457    search of clinical VGV are VGIV serotype B isolates recovered from patients.

458    Geographically, the most likely place to find the VGV clinical isolates appear to be in

459    sub-Saharan Africa since the current panel of isolates were found in the Zambian

460    environment within an ecoregion that includes Tanzania, Burundi, Democratic

461    Republic of the Congo, Angola and Malawi.

462

463

464

465 Previous work has shown that most isolates of the *C. gattii* species complex

466 generally cause pulmonary infection in a murine model with low neurotropism(20, 40,

467 41). The four VGV isolates tested here were less neurotropic than the VGIV isolate

468 MF46 that was collected from the same Zambian environment, and all the examined

469 Zambian environmental isolates were significantly less virulent than a VGIV control

470 strain, WM779. It remains to be shown if the differences in neurotropism are due to

471 lineage-specific genes, or alleles in VGV. As previous work has shown in *C.*

472 *neoformans* (42), capsule size difference manifested by Clade A and B *in vitro* was

473 unrelated to virulence in mice

474

475 Although serotypes have not yet been conclusively linked to virulence in

476 *Cryptococcus*, they remain important for strain identification. The majority of *C. gattii*

477 tested to date are serotype B, except for a subset of VGIII and the majority of VGIV

478 isolates which are serotype C. The six VGV isolates are also all serotype B - but due

479 to the slower growth rate on CGB agar, the CGB reaction was weaker than other

480 isolates of serotype B or serotype C. It took 24 hours longer for VGV compared to

481 other VG isolates (VGI-VGIV) to turn the medium dark blue. As the six VGV isolates

482 are all serotype B whilst the majority of VGIV isolates (their most closely related

483 lineage) reported thus far have been serotype C, it is possible that VGV may also

484 occur in serotype C. Additional environmental sampling of VGV is therefore

485 necessary to establish the dominant serotype, since the current sample size of six is

486 insufficient to make definitive conclusions.

487

488 Surprisingly, five of the six VGV isolates and the two control VGIV isolates were

489 highly resistant to fluconazole (MIC of >64$\mu$g/ml), a commonly used anti-fungal drug.

490 The three isolates of the VGV Clade A were more resistant to FLC than those of the

491 VGV Clade B. Although the *C. gattii* species complex was previously known to be on

492 average more resistant to FLC than *C. neoformans*(43), such high resistance to FLC

493 in environmental isolates is notable and has not yet been reported(44). All of the

494 VGV isolates had identical nucleotide sequences for *ERG11* and *AFR1*,

495 demonstrating the resistant isolates are not a result of genetic differences in the

496 target or transporter of FLC. However, innate fungal resistance to FLC can be due to

15

497   multiple factors besides the *ERG11* gene or efflux pumps and the mechanism(s) of

498   FLC resistance in VGV remain a subject for future investigation. Why environmental

499   VGV isolates should have such high resistance to azoles is unclear as it is unlikely

500   that they have come into contact with agrichemicals owing to the relatively pristine

501   environments from which they were recovered. More likely, fluconazole resistance is

502   a pleiotrophic effect that has evolved as a consequence of exposure by xenobiotics

503   other than azoles. Further investigations into the evolution of FLC resistance in VGV

504   may take on additional importance as clinical cases due to VGV are a distinct

505   possibility in the Sub-Saharan regions where 12% of the Zambian population are

506   living with the HIV virus(45).

507

508   In this paper, we present a near complete genome assembly for the VGV type strain,

509   MF34. The MF34 genome allowed us to conclusively establish that VGV is a

510   separate and distinct lineage of *C. gattii* from any previously identified, and not the

511   result of hybridisation, as has been seen for other divergent isolates(31). Indeed,

512   while both *Cryptococcus* species complexes appear to have a conserved

513   chromosome number of 14 based on the current panels of assembled and annotated

514   genomes available, intra- and inter-chromosomal rearrangements as well as large

515   CNV's appear to be common. This chromosomal variation may provide the genetic

516   basis for phenotypic variation and may act as a genetic barrier to recombination

517   between more divergent isolates such as those from separate lineages. At the

518   within-lineage level, there are also a number of unique and uniquely lost "lineage-

519   specific genes', which may contribute to phenotypic differences between lineages.

520   However, it should be noted that many of the main phenotypes routinely measured,

521   including virulence in animal models, growth rates, and ability to cause pulmonary

522   versus CNS infections, appear to vary as much within as between lineages.

523

524   One line of future inquiry towards explaining this phenotypic diversity may come from

525   the characterisation of further transcriptional differences. For example, VGII

526   upregulates many of the ergosterol genes during co-incubation with bone-marrow

527   derived macrophages(46) and it will be important to determine whether other traits

528   exist which differentiate the lineages of *Cryptococcus.* Further, it will be important to

529   examine whether similar lineage-specific differences underpin VGV's increased FLC

530   resistance, and whether clinically-relevant traits such as drug resistance are linked to

16

531 the environment within which these isolates have evolved. Ultimately, our study

532 testifies to the deep reservoir of diversity that exists within *Cryptococcus* which,

533 despite decades of research into this genus, still harbours abundant surprises.

534

546

## Data access

548 The raw sequence and genome assembly of VGV MF34 is available in NCBI under

549 BioProject PRJNA487802.

550

551 **Author contributions**

552 RAF, LvD, and TRS performed the genomic analyses.

553 RAF, MC, CAC, MCF, JK-C, HE, LvD, and FB wrote the manuscript.

554 MC, MJD, and DHY performed the phenotypic assays.

555 JK-C, YCC, WM and CAC sequenced the isolates.

556 TS and CAC assembled the VGV genome

557 MV, DC, GK, and MCF conducted the field work.

558

## Figure and Table Legends

560 **Fig. 1. Environmental sampling of *C. gattii* VGV in Zambia. A**) Location of *C.*

561 *gattii* VGV isolates across Central Zambian Miombo Woodlands. Isolates MF5, MF13

562 and 54 (Mupata Hills, Copperbelt province) and MF56 (Kapishya, Northern Province)

563 were found in or near to Hyrax middens created by Southern tree hyrax

564 (*Dendrohyrax arboreus*). MF34 and MF51 were found from sampling Miombo tree

565 holes in the Mutinondo wilderness area, Northern Province. **B)** A tree hyrax feeding

566 on leaves. **C)** Sampling from hyrax middens at Kapishya from which MF56 was

567 isolated.

568

569 **Fig. 2. Whole genome analysis supports VGV as a distinct lineage. a)** Maximum

570 likelihood (RAxML) phylogeny of 101 *C. gattii* genomes generated over all non-

571 ambiguous sites with at least a SNP in ≥1 isolate (1,518,323 sites, or 8.7% of the

572 total genome). Isolate names are coloured according to lineage (dark blue = VGI,

573 green = VGII, purple = VGIII, light blue = VGIV, orange = VGV, and red = VGVI).

574 Asterisks indicate 100% bootstrap support at each node after 1,000 tree-building

575 replicates. **b)** SplitsTree NeighborNet Network**. c-d)** Principal Component Analysis

576 (PCA) of genomic variant sites, showing separation of isolates into lineages (isolates

577 plotted with 2% random noise (jitter) for clarity of individual points). The first four

578 Principal Components (PCs) account for 71.26% of the total genetic variation.

579

580 **Fig. 3. Population genetic analyses for each of the *C. gattii* lineages (based on**

581 **101 genomes). a)** Shared alleles (SNPs per Kb) between each of the *C. gattii*

582 lineages. **b)** Private alleles (SNPs per Kb) between each of the *C. gattii* lineages. **c)**

583 Nucleotide diversity (π) within each lineage against the number of isolates

584 representing each lineage. **d)** Admixture K optimisation based on Cross Validation

585 Error. **e)** Unsupervised ADMIXTURE clustering analysis of all isolates at K=9. **f)**

586 Lineage specific gene and lineage specific gene-loss counts. The tree topology is

587 based on the core-ortholog RAxML tree setting equal branch lengths, and the

588 number of multi-lineage-specific gene gains and losses are shown above internal

589 nodes.

590

591 **Fig. 4. A phylogenetic tree for ten *Cryptococcus* genomes belonging to the six**

592 ***C. gattii* lineages and one *C. neoformans* lineage outgroup alongside their**

593 **genome synteny.** The phylogenetic tree was constructed in RAxML with branch

594 lengths indicating the mean number of nucleotide substitutions per site. To the right

595 is a synteny plot, visualizing regions that span three or more orthologs between any

596 two species as a connected grey line. Supercontig numbers are shown above each

597   genome axis if longer than 400 kb, where + represents the forward orientation and –
598   represents the negative orientation.
599
600   **Fig.5. Phenotypic characteristics. a)** Growth of six VGV isolates and two VGIV
601   control isolates at 30°C and 37°C on YEPD agar. **b)** India ink staining of VGV cells
602   grown in YEPD broth for 24 hours at 30°C. The isolates of Clade A (MF5, 13, 53)
603   produce thinner capsules than the isolates of Clade B cells which produce a similar
604   size capsule compared to the two VGIV control isolates. Bar = 5μm. **c)** MIC of VGV
605   isolates for FLC, 5-FC and amphotericin B. All tested isolates had high MIC for FLC
606   ranging from 20 to 256 μg/ml. The MIC for 5-FC was low except for MF46, an
607   environmental isolate of VGIV which showed average 4μg/ml. All VGV isolates
608   showed higher MIC for amphotericin B than VGIV controls. **d)** Survival curve of mice
609   infected by four VGV isolates (intrapharyngeal aspiration of 5,000 cell/mouse) and
610   two VGIV control isolates. Only the mice infected with VGIV isolate WM779
611   succumbed to infection. **e)** Lung and brain fungal loads. VGV isolates grew
612   moderately in lungs but the CFU in the brains were negligible.
613
614   **Fig. 6. Histopathology of the lung infected by VGV isolates.** Sections of the
615   mouse lungs infected by three different VGV isolates and two VGIV control strains
616   stained by Alcian blue, Periodic acid-Schiff stain and counterstained with
617   haematoxylin (left and middle columns) or with standard haematoxylin (right column).
618   Note that Alcian blue stains cryptococcal cells blue. Images in the left columns were
619   acquired using a 2.5X objective. Images in the middle and right columns are higher
620   magnification (10X) images of the area indicated by the yellow boxes in the 2.5X
621   images.
622
623   **Table 1.** Environmental isolates of VGIV and VGV from the Central Zambezian
624   Miombo Woodlands.
625
626   **Table 2.** The genome assemblies used for phylogenetic analysis and orthology
627   detection. *indicates newly described genome assemblies for this paper. All others
628   have been described previously(26).
629

19

## Supplemental Figure and Table Legends

**Fig S1** Admixture analysis of 101 isolates based on K=2 through to K=15. The lowest Cross Validation Error was found at K=9.

**Fig. S2** Chromopainter's inferred proportion of genome-wide DNA that each strain shares with every other based on pairwise haplotype matching profiles. The tree at the top provides fineSTRUCTURE's inferred hierarchical merging of clusters based on these profiles. Tick marks at the bottom are coloured according to lineage. VGV form their own cluster, tending to match more haplotypes genome-wide with isolates from VGI, VGIII and VGIV compared to VGII.

**Fig. S3** Normalised read depth across 10 kb sliding windows along each supercontig relative to the R265 reference sequence. Aneuploid regions were only identified in two of the three *C. decagattii* (VGVI) isolates, which are highlighted by red circles. The lower depth of coverage across the start of supercontig 18 indicates a *MAT***a** isolate, compared with the R265 *MAT*α.

**Fig. S4.** Genome-wide variation in θ, Weir's formulation of Wright's fixation index (FST), on pairwise comparisons in each lineage. For comparison of isolates between each VG group, θ was calculated across window lengths of 10 kb. The lower $F_{ST}$ at the start of supercontig 18 shows the location of the *MAT*α locus. Below the non-overlapping windows, mean pairwise $F_{ST}$ values from all nuclear supercontigs are shown.

**Fig. S5.** The numbers of protein coding genes, rRNA, tRNA, genes with PFAMs, KEGG-ECs, GO-terms, predicted secreted genes (SignalP4) and transmembrane genes (TMHMM) for each assembly described in this paper.

**Fig. S6.** Coverage of the 248 Core Eukaryotic Genes (CEGs) by the *C. gattii* and *C. neoformans* gene-sets described in this paper.

**Fig. S7. Synima/Orthofinder identified orthogroups (a)** All orthogroups (representing every gene in each lineage) were grouped into a variety of categories including Orthologs 1:1 and Orthologs 1:>1 – which are orthogroups with 1:1 orthology in all isolates except for within a single isolate of the lineage that includes paralogs. "Orthologs divergent" and "Paralogs divergent" are genes that were previously lineage or isolate specific, but BLASTn revealed them to be unique orthogroups that, when joined to the database search genes, would make a 1:1 ortholog or Paralog conserved, respectively. Paralogs Lineage-specific (L.S.) and Paralogs Strain-specific (S.S.) are those genes that are represented by a single gene in all the other isolates in the category Orthologs 1:1 divergent. Paralogs miscellaneous (misc.) are all other Orthogroups including paralogous clusters. Genes absent in one lineage and lineage specific genes are broken down into further categories in **panel b**. Absent in one strain and present in one strain are self-explanatory, while Miscellaneous (Misc.) contains all remaining Orthogroups (such as genes found in multiple gene-sets but not lineage-specific or an ortholog. **b)** Details the number of genes in each lineage that are either absent in another lineage or specific to just that lineage.

**Fig. S8.** CGB reaction, melanin and urease production by VGV isolates. The melanin and urease production by VGV isolates were similar to the other VG type isolates but the CGB reaction took a longer time due to the slower growth on CGB agar.

**Fig. S9.** Patterns of *URA5* RFLP of VGIV and VGV isolates. **A)** The banding patterns of *URA5* uncut or *Sau*96I/*Hha*I digests are identical between VGIV and VGV Clade A isolates. However, Clade B isolates of VGV show 1.3 kb amplicon (red arrow) both in uncut as well as in *Sau*96I/*Hha*I digested DNA. This 1.3 kb amplicon is also absent in other molecular types. **B)** The uncut as well as *Stu*I and *Ear*I digests all show the 1.3 kb amplicon only in the Clade B isolates. **C)** The RFLP patterns of the 6 VGV isolates and 17 VGIV isolates digested by *Stu*I and **D)** *Ear*I showing clear difference between the two molecular types.

**Table S1. (Tab 1)** Details of all isolates used in this study. Isolates include those newly sequenced, and those presented in previous papers (with select citations included). Details of the alignments to R265 are given, along with the variants called.

696 **(Tab 2) Pairwise comparisons of amino acid differences found among 1:1 core**

697 **orthologs between all lineages of *C. gattii* and *C. neoformans* VNI H99**. Each

698 ortholog orthogroup has been aligned using MUSCLE and concatenated into a

699 contiguous sequence used for phylogenetic reconstruction. **(Tab 3) Lineage**

700 **specific genes for VGV, VGVI and VNI.** Details of all lineage-specific genes and

701 genes uniquely absent in each of the lineages. Columns include unique orthogroup

702 number, gene ID, GO-term annotation, length of gene, PFAMs, GO-terms, SigP4

703 predictions and TMHMM predictions. Genes that are uniquely absent in a lineage are

704 represented by a separate lineage (VNI absent represented by VGV genes, VGVI

705 (Cd) absent represented by VNI genes, VGV absent represented by VNI genes).

706 **(Tab 4) Isolates of VGIV used to distinguish from VGV by *URA5* RFLP.**

707

708 ## Methods

709

710 **Library preparation and sequencing of Zambian isolates**

711 Environmental sampling took place in January and September of 2013. Samples

712 were collected using "Transwab" Amies swabs (MWETM – MW170) and sterilized

713 30-mL screw- capped glass bottles. Amies liquid transport swabs were taken from

714 tree bark ($n$=20), soil ($n$=19) and cracks in granite kopjes or droppings from rock

715 Hyrax middens ($n$=16). Samples were collected and processed according to

716 previously established protocols(47, 48), and the samples were kept at 4°C before

717 being processed on niger seed agar. All samples were collected under license from

718 the Zambian Wildlife Authority (ZAWA).

719

720 Single colonies purified from the original isolation media were maintained

721 cryopreserved at -80°C at Imperial College in London since 2013. The isolates were

722 revived on YPD agar (Yeast extract 1%, Peptone 1%, glucose 2%) and incubated at

723 30°C before use. Genomic DNA was isolated with CTAB extraction method as

724 described previously with modification(49). Paired-end libraries (150 bp) were

725 prepared and sequenced using the Illumina HiSeq 4000 platform by Novogene

726 (Davis, CA). Two Oxford Nanopore libraries of isolate MF34 were constructed from

727 genomic DNA using the 1D library construction kit (SQK-LSK109). A total of 243,660

728 reads with an N50 of 9,827 were generated on a FLO-MIN106 flow cell using a

729 Minion. Reads were base called using Albacore v2.3.1. This resulted in 923,997,900

730     total bases (~46X coverage). Raw sequence data was submitted to the NCBI

731     Sequence Read Archive under BioProject ID PRJNA476154 (all *C. gattii* non-VGV

732     isolates) and PRJNA480403 (all *C. gattii* VGV isolates).

733

734     **Genome assembly and annotation**

735     For Isolate MF34, a hybrid assembly of Oxford Nanopore long-reads and Illumina

736     short reads was generated. An initial assembly of the Oxford reads was generated

737     using Canu v1.5(50) with parameter genomeSize=20,000,000. The assembly was

738     inspected for the presence of telomeric repeat (TTAGGG) at contig ends; for two

739     contig ends missing telomeric repeat, contigs were extended by aligning

740     unassembled Canu contigs to these ends using NUCmer v3.1(51). Base called

741     reads were then aligned to the contigs with BWA mem(52) with flag "-x ont2d", and

742     the alignments used for polishing with Nanopolish(53). Two rounds of Pilon

743     v1.13(54) correction were performed using Illumina BWA read alignments(52).

744     Paired Illumina sequences of *C. decagattii* (VGVI) were assembled and scaffolded

745     using SPAdes v3.1.1(55) with *k*-mer lengths (21, 33, 55 and 77). An assembly

746     statistics summary for the assembly is provided in **Table 2**. Reads were aligned back

747     to the assembly with BWA v0.7.4-r385 mem(52), and Pilon v1.12(54) was further

748     used to improve the assembly. Scaffolds smaller than 1Kb were removed. The

749     genome assembly has been submitted to NCBI under the project accession

750     PRJNA487802.

751

752     The *C. gattii* VGV MF34 and VGVI WM1802 genomes were annotated using

753     Genemark(56), BLASTx against SwissProt(57) and KEGG(58), and HMMER

754     hmmscan(59) against PFAM(60). We ran tRNAscan(61) and RNAmmer(62) to

755     identify non-protein coding genes. Gene predictions were checked for a variety of

756     issues, including overlap with non-coding genes, overlap with coding genes, and the

757     presence of in-frame stops. Genes were named according to evidence from BLASTx

758     and HMMER following order of precedence: (1) SwissProt(57), (2) TIGRfam(63), and

759     (3) KEGG(58), where BLASTx hits must meet the 70% identity and 70% overlap

760     criteria to be considered a good hit and for the name to be applied. Otherwise, genes

761     were named as hypothetical proteins.

762

763    Genes were functionally annotated by assigning PFAM domains(60), GO terms,

764    KEGG assignment and ortholog mapping to genes of known function. HMMER3(59)

765    was used to identify PFAM (release 27) domains, and BLASTx used against the

766    KEGG v65 database(58) (e-value<$1 \times 10^{-10}$). GO terms were assigned using

767    Blast2GO version2.3.5(64), with a minimum e-value of $1 \times 10^{-10}$. SignalP 4.0(65)

768    and TMHMM(66) were used to identify secreted proteins and trans-membrane

769    proteins respectively (**Fig. S5**). Gene sets were aligned to the 248 Core Eukaryotic

770    Genes (CEGs) and BUSCO basidiomycota_odb9 set to evaluate completeness (**Fig.**

771    **S6**).

772

773    **Read alignment and variant identification**

774    The 36 newly sequenced isolates from this study were compared to an additional 65

775    isolates that were sequenced and described in previous studies(20, 26, 38, 67, 68).

776    These additional isolates were obtained from the NCBI Sequence read archive

777    (SRA) and converted from SRA format to FASTQ using SRAtoolkit version 2.3.3–4.

778    Illumina reads were aligned to the *C. gattii* VGII R265 reference genome assembly

779    using Burrows-Wheeler Aligner (BWA) v0.7.4-r385 mem(52) with default parameters

780    and converted to sorted BAM format using SAMtools v0.1.9 (r783)(69).

781

782    Genome Analysis Toolkit (GATK) v2.7-4-g6f46d11(70) was used to call both variant

783    and reference nucleotides from the 101 alignments (as previously described(24)).

784    Briefly, the Picard tools AddOrReplaceReadGroups, MarkDuplicates,

785    CreateSequenceDictionary, and ReorderSam were used to preprocess the

786    alignments (http://broadinstitute.github.io/picard/). GATK RealignerTarget-Creator

787    and IndelRealigner were then used to resolve misaligned reads close to indels. Next,

788    GATK Unified Genotyper (with the haploid Genotyper ploidy setting) was run with

789    both SNP and indel genotype likelihood models (GLM). We also ran Base

790    Recalibrator and PrintReads for base quality score recalibration on those initial sites

791    for GLM SNP. We then recalled variants with Unified Genotyper with the parameter

792    "—output_mode EMIT_ALL_SITES." We merged and sorted all of the calls and then

793    ran Variant Filtration with the parameters "QD < 2.0, FS > 60.0, MQ < 40.0." Next,

794    we removed any base that had less than a minimum genotype quality of 50, a

795    minimum percent alternate allele (AD) of 80%, or a minimum depth of 10. Finally, we

796    removed any positions that were called by both GLMs (i.e., incompatible indels and

797    SNPs), any marked as "LowQual" by GATK, any nested indels, and any sites that did
798    not include a PASS flag.

799

800    **Phylogenetic and population genetic analysis**
801    The variants identified from the 101 alignments were filtered for positions that were
802    homozygous (reference or SNP) and polymorphic in one or more isolate (**Fig. 2**),
803    resulting in an alignment of 1,517,353 nuclear sites and 970 mitochondrial sites. A
804    FASTA file of these positions was created and converted into PHYLIP format, and a
805    phylogenetic tree was generated using RAxML v7.7.8(71) with 1,000 bootstrap
806    replications. RAxML was run with the generalized time-reversible (GTR) and
807    category (CAT) rate approximation with final evaluation of the tree using GTR plus
808    gamma-distributed rates. The same sites were analysed using the NeighborNet
809    Network of SplitsTree v4.14.6(72).

810

811    A multi sample VCF of all 101 genomes was made with VCFtools(73) and converted
812    to ped and map file formats for use in PLINK v1.90(74). Unsupervised
813    ADMIXTURE(75) was run on a moderately Linkage Disequilibrium (LD) pruned
814    alignment for values of K between 1-15. A value of K=9 provided the lowest cross-
815    validation error (**Fig3d-e, Fig. S1**). To explore finer-patterns of population structure
816    amongst our sampled lineages we applied a technique designed to characterise
817    patterns of haplotype sharing between a panel of "donor" and "recipient" haplotypes
818    within a recombining population. We ran Chromopainter v2(23) to infer, at each
819    position in a recipient isolate's genome, which donor they are most closely related to
820    ancestrally relative to all others in the dataset. To do this, we assumed a uniform
821    recombination rate of 1.5 morgans/megabase based on the genome wide
822    recombination rate previously estimated in *C. neoformans*(76) and with
823    Chromopainter's switch and mutation rate parameters estimated using 10 runs of
824    Expectation-Maximisation (-n 190.29, -M 0.0011). We then ran Chromopainter in
825    linked mode using the haploid switch (-j) under an "all-versus-all" framework, painting
826    all samples using all others to produce a pair-wise coancestry matrix describing the
827    amount of DNA each isolate matches to every other under the copying model.

828

829    Haplotype based clustering was then implemented in fineSTRUCTURE(23) with an
830    estimated normalization parameter of *c=0.51*, sampling cluster assignments every

25

831    10,000 iterations for $1\times10^6$ MCMC iterations after $1\times10^6$ initial burn-in steps. We then

832    performed an additional $1\times10^5$ hill-climbing iterations beginning with the MCMC

833    sample with the highest posterior probability. This classified our data into 34 clusters

834    (**Fig. S2**).

835

836    For the *C. neoformans* VNI H99 rooted *C. gattii* tree, we identified 1:1 orthologs

837    among each of the nine isolates with Orthofinder v2.1.2(77) using the Synima

838    pipeline(78). We aligned orthologs with MUSCLE v3.8.31(79), extracted the CDS

839    sequences in a codon context, and trimmed to the smallest contiguous sequence,

840    and then concatenated alignments. In total, we aligned 2.16 Mb of transcripts for

841    each genome. Prottest v3.4(80) was used to determine the best-fitting amino acid

842    transition model (JTT) according to Bayesian information criterion. The final tree was

843    produced using RAxML v7.7.8(71) using the CAT rate approximation and WAG

844    amino acid replacement matrix with 1,000 bootstrap replicates. Synima(78) was

845    used to visualise synteny between each of the genomes. The same pipeline was

846    used to compare the previous and updated R265 genomes.

847

848    **Phenotypic analysis**

849    To determine the growth rate of *C. gattii* VGV, cells of all six VGV isolates were

850    inoculated in YEPD broth and incubated at 30°C on a shaker (200rpm) for 18h. Cells

851    were washed with sterile PBS and $2\times10^5$ cells/ml were resuspended in PBS. Three

852    microliter aliquots of 10-fold serial dilutions were spotted onto YEPD agar and

853    incubated at 30°C and 37°C. For biological confirmation of the species, isolates were

854    inoculated on Canavanine glycine bromothymol blue (CGB) agar(81) for species

855    specific CGB reaction and Christensen's urea agar (Sigma) and norepinephrine

856    agar(82) for urease and melanin production respectively and incubated at 30°C for

857    48 hours. India ink mount of the cells grown on YEPD broth for 24h at 30°C were

858    used for microscopic observation of the cell and polysaccharide capsule size. The

859    reference strains used were WM148 or H99 (serotype A, VNI), WM626 (serotype A,

860    VNII), WM179 (serotype B, VGI), WM178, R265 and R272 (serotype B, VGII),

861    WM161 (serotype B, VGIII), and WM779 (serotype C, VGIV)(14). Mating type of

862    each isolate was determined by PCR using primers specific to the *STE12*α and

863    *STE20*a(83).

864

26

**Determination of MIC for antifungal antibiotics**

MICs for fluconazole (FLC), 5-fluorocytosine (5FC), and Amphotericin B were determined using Etest strips according to the Etest technical guide (AB Biodisk, Solna, Sweden), with slight modification. Fungal cells were grown in 5 ml of YEPD at 30°C for 18 hours. Harvested cells were diluted in sterile saline to an optical density of 0.05 at 600 nm (OD600) and plated on yeast nitrogen base (YNB) agar plates. Etest strips were placed at the center of the plates and incubated at 30°C for 72 hours. The susceptibility endpoint was read at the first growth inhibition ellipse. The concentration ranges tested were: FLC, 0.016 to 256 µg/ml; both 5-FC and Amphotericin B, 0.002 to 32 µg/ml.

*URA5* **gene RFLP**

The *URA5* gene of each isolate was amplified from genomic DNA by PCR to identify the molecular type using 50 ng of two primers: URA5 (5′-ATGTCCTCCCAAGCCCTCGACTCCG-3′) and SJ01 (5′-TTAAG ACCTCTGAACACCGTACTC-3′). Reactions were carried out in a total volume of 50 µL as previously described(14). PCR was performed for 40 cycles at 94°C for 2 min initial denaturation, 30 s of denaturation at 94°C, 30 s annealing at 55°C, and 2 min extension at 72°C. The reactions were completed by a final extension step for 10 min at 72°C. PCR products were analysed by 1% agarose gel electrophoresis and 5 µL of PCR products were double digested with *Sau*96I (10 U/µL) and *Hha*I (20 U/µL) for 3 h at 37°C. Then, digested samples were separated by 3% agarose gel electrophoresis at 80V for 5 h. The RFLP patterns of *URA5* gene were analysed using well-defined bands in the gel images by comparing them with the patterns obtained from the standard reference strains.

**Restriction enzyme analysis of the *URA5* gene to distinguish VGV from VGIV**

We found the *URA5* RFLP banding patterns(14) of VGV and VGIV are not clearly distinguishable although *Sau*96I and *Hha*I (**Fig. S9A**). We compared the DNA sequences of the *URA5* gene from MF34 (VGV) and WM779 (VGIV) and found two restriction enzymes, *Stu*I and *Ear*I that can be used to distinguish the two molecular types based on *URA5* RFLP. Three microliters of *URA5* PCR products were digested with *Stu*I (10 U/µl) or *Ear*I (20U/µl) (New England BioLabs Inc) at 37°C for

27

898    4h and restriction fragments were separated by electrophoresis in 3% agarose Tris-

899    acetate-EDTA (TAE) gels at 80V for 5h. Standard reference strains for molecular

900    typing were used as controls.

901

902    **Virulence in mice**

903    The virulence of four VGV isolates, two from Clade A and two from Clade B, was

904    assessed using seven to eight weeks old female C57BL/6 mice (Taconic Farms).

905    Isolates to be tested in mice were inoculated in YEPD broth and incubated overnight,

906    washed twice and diluted to $2.5X10^5$ cells/ml in PBS. Mice (14 mice per isolate) were

907    inoculated with $20\mu l$ of cell suspension ($5x10^3$/mouse) by pharyngeal aspiration.

908    Eight mice for each isolate were used for the survival rate and six mice each were

909    used for the analysis of fungal burden and histopathology at the indicated time

910    points. Mice were monitored twice per day and differences in survival were

911    determined using GraphPad Prism, version 7 (GraphPad Software, San Diego, CA).

912

913    To assess the organ fungal burden, lungs and brains of four mice from each infected

914    group were inspected. The mice infected with WM779 started to die on day 25 post

915    infection and the lungs and brains were harvested immediately from the dead mice

916    on day 25. Mice infected with other isolates were euthanized on day 60 and organs

917    were harvested.  Harvested lungs and brains were homogenized in 7ml and 2ml

918    sterile water respectively and 5 µl aliquots of 10-fold serial dilutions were plated on

919    YEPD agar and incubated at $30^oC$ for 48h. Fungal colonies were counted and the

920    tissue fungal burden was analysed using GraphPad Prism, version 7 (GraphPad

921    Software, San Diego, CA).

922

923    **Histopathological analysis**

924    For histopathological analysis, organs of infected mice from each group were

925    harvested at 10 and 20 days post inoculation and fixed in 3.7% buffered formalin and

926    embedded in paraffin. Sections were stained with hematoxylin and eosine (H&E) or

927    Alcian blue/periodic acid-Schiff (AB/PAS) at the Histoserv Inc.

928

929    **Ethics statement**

930

The Institutional Animal Care and Use Committee of the National Institute of Allergy and Infectious Diseases approved all animal studies (#LCIM-5E). Studies were performed in accordance with recommendations of the Guide for the Care and Use of Laboratory Animals of the National Institutes of Health.

## References

1. Rajasingham R, Smith RM, Park BJ, Jarvis JN, Govender NP, Chiller TM, Denning DW, Loyse A, Boulware DR. 2017. Global burden of disease of HIV-associated cryptococcal meningitis: an updated analysis. Lancet Infect Dis 17:873–881.

2. Busse O. 1894. Uber parasitare Zelleinschlusse und ihre Zuchtung. Cent Bakt Parasit 16.

3. Barnett JA. 2010. A history of research on yeasts 14: medical yeasts part 2, *Cryptococcus neoformans*. Yeast Chichester Engl 27:875–904.

4. Kwon-Chung KJ, Fraser JA, Doering TL, Wang Z, Janbon G, Idnurm A, Bahn Y-S. 2014. *Cryptococcus neoformans* and *Cryptococcus gattii*, the etiologic agents of cryptococcosis. Cold Spring Harb Perspect Med 4:a019760.

5. Evans EE. 1950. The Antigenic composition of *Cryptococcus neoformans*: I. A serologic classification by means of the capsular and agglutination reactions. J Immunol 64:423–430.

6. Wilson DE, Bennett JE, Bailey JW. 1968. Serologic grouping of *Cryptococcus neoformans*. Proc Soc Exp Biol Med 127:820–823.

7. Kwon-Chung KJ. 1975. A new genus, filobasidiella, the perfect state of *Cryptococcus neoformans*. Mycologia 67:1197–1200.

954   8.  Kwon-Chung KJ. 1976. A new species of Filobasidiella, the sexual state of *Cryptococcus*

955        *neoformans* B and C serotypes. Mycologia 68:943–946.

956   9.  Franzot SP, Salkin IF, Casadevall A. 1999. *Cryptococcus neoformans var. grubii*:

957        separate varietal status for *Cryptococcus neoformans* serotype A isolates. J Clin

958        Microbiol 37:838–840.

959   10. Xu J, Vilgalys R, Mitchell TG. 2000. Multiple gene genealogies reveal recent dispersion

960        and hybridization in the human pathogenic fungus *Cryptococcus neoformans*. Mol Ecol

961        9:1471–1481.

962   11. Diaz MR, Boekhout T, Kiesling T, Fell JW. 2005. Comparative analysis of the intergenic

963        spacer regions and population structure of the species complex of the pathogenic yeast

964        *Cryptococcus neoformans*. FEMS Yeast Res 5:1129–1140.

965   12. Kwon-Chung KJ, Boekhout T, Fell JW, Diaz M. 2002. (1557) Proposal to Conserve the

966        Name *Cryptococcus gattii* against *C. hondurianus* and *C. bacillisporus* (Basidiomycota,

967        Hymenomycetes, Tremellomycetidae). Taxon 51:804–806.

968   13. Kwon-Chung KJ, Varma A. 2006. Do major species concepts support one, two or more

969        species within *Cryptococcus neoformans*? FEMS Yeast Res 6:574–587.

970   14. Meyer W, Castañeda A, Jackson S, Huynh M, Castañeda E, IberoAmerican Cryptococcal

971        Study Group. 2003. Molecular typing of IberoAmerican *Cryptococcus neoformans*

972        isolates. Emerg Infect Dis 9:189–195.

973   15. Boekhout T, Theelen B, Diaz M, Fell JW, Hop WC, Abeln EC, Dromer F, Meyer W.

974        2001. Hybrid genotypes in the pathogenic yeast *Cryptococcus neoformans*. Microbiol

975        Read Engl 147:891–907.

976  16. Meyer W, Aanensen DM, Boekhout T, Cogliati M, Diaz MR, Esposto MC, Fisher M,

977      Gilgado F, Hagen F, Kaocharoen S, Litvintseva AP, Mitchell TG, Simwami SP, Trilles

978      L, Viviani MA, Kwon-Chung J. 2009. Consensus multi-locus sequence typing scheme

979      for *Cryptococcus neoformans* and *Cryptococcus gattii*. Med Mycol 47:561–570.

980  17. Hagen F, Khayhan K, Theelen B, Kolecka A, Polacheck I, Sionov E, Falk R, Parnmen S,

981      Lumbsch HT, Boekhout T. 2015. Recognition of seven species in the *Cryptococcus*

982      *gattii/Cryptococcus neoformans* species complex. Fungal Genet Biol FG B 78:16–48.

983  18. Kwon-Chung KJ, Bennett JE, Wickes BL, Meyer W, Cuomo CA, Wollenburg KR,

984      Bicanic TA, Castañeda E, Chang YC, Chen J, Cogliati M, Dromer F, Ellis D, Filler SG,

985      Fisher MC, Harrison TS, Holland SM, Kohno S, Kronstad JW, Lazera M, Levitz SM,

986      Lionakis MS, May RC, Ngamskulrongroj P, Pappas PG, Perfect JR, Rickerts V, Sorrell

987      TC, Walsh TJ, Williamson PR, Xu J, Zelazny AM, Casadevall A. 2017. The case for

988      adopting the "Species Complex" nomenclature for the etiologic agents of

989      cryptococcosis. mSphere 2.

990  19. Vanhove M, Beale MA, Rhodes J, Chanda D, Lakhi S, Kwenda G, Molloy S,

991      Karunaharan N, Stone N, Harrison TS, Bicanic T, Fisher MC. 2017. Genomic

992      epidemiology of Cryptococcus yeasts identifies adaptation to environmental niches

993      underpinning infection across an African HIV/AIDS cohort. Mol Ecol 26:1991–2005.

994  20. Firacative C, Roe CC, Malik R, Ferreira-Paim K, Escandón P, Sykes JE, Castañón-

995      Olivares LR, Contreras-Peres C, Samayoa B, Sorrell TC, Castañeda E, Lockhart SR,

996      Engelthaler DM, Meyer W. 2016. MLST and whole-genome-based population analysis

997      of *Cryptococcus gattii* VGIII links clinical, veterinary and environmental strains, and

998      reveals divergent serotype specific sub-populations and distant ancestors. PLoS Negl

999      Trop Dis 10:e0004861.

1000  21. Hagen F, Illnait-Zaragozí M-T, Meis JF, Chew WHM, Curfs-Breuker I, Mouton JW,

1001      Hoepelman AIM, Spanjaard L, Verweij PE, Kampinga GA, Kuijper EJ, Boekhout T,

1002      Klaassen CHW. 2012. Extensive genetic diversity within the dutch clinical

1003      *Cryptococcus neoformans* population. J Clin Microbiol 50:1918–1926.

1004  22. Gillece JD, Schupp JM, Balajee SA, Harris J, Pearson T, Yan Y, Keim P, DeBess E,

1005      Marsden-Haug N, Wohrle R, Engelthaler DM, Lockhart SR. 2011. Whole genome

1006      sequence analysis of *Cryptococcus gattii* from the Pacific Northwest reveals unexpected

1007      diversity. PloS One 6:e28550.

1008  23. Lawson DJ, Hellenthal G, Myers S, Falush D. 2012. Inference of population structure

1009      using dense haplotype data. PLOS Genet 8:e1002453.

1010  24. Chen Y, Farrer RA, Giamberardino C, Sakthikumar S, Jones A, Yang T, Tenor JL,

1011      Wagih O, Wyk MV, Govender NP, Mitchell TG, Litvintseva AP, Cuomo CA, Perfect

1012      JR. 2017. Microevolution of serial clinical isolates of *Cryptococcus neoformans var.*

1013      *grubii* and *C. gattii*. mBio 8:e00166-17.

1014  25. Steenwyk JL, Soghigian JS, Perfect JR, Gibbons JG. 2016. Copy number variation

1015      contributes to cryptic genetic variation in outbreak lineages of *Cryptococcus gattii* from

1016      the North American Pacific Northwest. BMC Genomics 17:700.

1017  26. Farrer RA, Desjardins CA, Sakthikumar S, Gujja S, Saif S, Zeng Q, Chen Y, Voelz K,

1018      Heitman J, May RC, Fisher MC, Cuomo CA. 2015. Genome evolution and innovation

1019      across the four major lineages of *Cryptococcus gattii*. mBio 6:e00868-00815.

1020  27. D'Souza CA, Kronstad JW, Taylor G, Warren R, Yuen M, Hu G, Jung WH, Sham A,

1021      Kidd SE, Tangen K, Lee N, Zeilmaker T, Sawkins J, McVicker G, Shah S, Gnerre S,

1022      Griggs A, Zeng Q, Bartlett K, Li W, Wang X, Heitman J, Stajich JE, Fraser JA, Meyer

1023       W, Carter D, Schein J, Krzywinski M, Kwon-Chung KJ, Varma A, Wang J, Brunham

1024       R, Fyfe M, Ouellette BFF, Siddiqui A, Marra M, Jones S, Holt R, Birren BW, Galagan

1025       JE, Cuomo CA. 2011. Genome variation in *Cryptococcus gattii*, an emerging pathogen

1026       of immunocompetent hosts. mBio 2:e00342-00310.

1027  28. Janbon G, Ormerod KL, Paulet D, Iii EJB, Yadav V, Chatterjee G, Mullapudi N, Hon C-

1028       C, Billmyre RB, Brunel F, Bahn Y-S, Chen W, Chen Y, Chow EWL, Coppée J-Y,

1029       Floyd-Averette A, Gaillardin C, Gerik KJ, Goldberg J, Gonzalez-Hilarion S, Gujja S,

1030       Hamlin JL, Hsueh Y-P, Ianiri G, Jones S, Kodira CD, Kozubowski L, Lam W, Marra

1031       M, Mesner LD, Mieczkowski PA, Moyrand F, Nielsen K, Proux C, Rossignol T, Schein

1032       JE, Sun S, Wollschlaeger C, Wood IA, Zeng Q, Neuvéglise C, Newlon CS, Perfect JR,

1033       Lodge JK, Idnurm A, Stajich JE, Kronstad JW, Sanyal K, Heitman J, Fraser JA, Cuomo

1034       CA, Dietrich FS. 2014. Analysis of the genome and transcriptome of *Cryptococcus*

1035       *neoformans var. grubii* reveals complex RNA expression and microevolution leading to

1036       virulence attenuation. PLOS Genet 10:e1004261.

1037  29. Yadav V, Sun S, Billmyre RB, Thimmappa BC, Shea T, Lintner R, Bakkeren G, Cuomo

1038       CA, Heitman J, Sanyal K. 2018. RNAi is a critical determinant of centromere evolution

1039       in closely related fungi. Proc Natl Acad Sci U S A 115:3108–3113.

1040  30. Litvintseva AP, Thakur R, Vilgalys R, Mitchell TG. 2006. Multilocus sequence typing

1041       reveals three genetic subpopulations of *Cryptococcus neoformans var. grubii* (serotype

1042       A), including a unique population in Botswana. Genetics 172:2223–2238.

1043  31. Desjardins CA, Giamberardino C, Sykes SM, Yu C-H, Tenor JL, Chen Y, Yang T, Jones

1044       AM, Sun S, Haverkamp MR, Heitman J, Litvintseva AP, Perfect JR, Cuomo CA. 2017.

1045       Population genomics and the evolution of virulence in the fungal pathogen

1046       *Cryptococcus neoformans*. Genome Res 27:1207–1219.

1047     32. Wiley EO. 1978. The evolutionary species concept reconsidered. Syst Zool 27:17–26.

1048     33. Hagen F, Lumbsch HT, Arsenijevic VA, Badali H, Bertout S, Billmyre RB, Bragulat

1049        MR, Cabañes FJ, Carbia M, Chakrabarti A, Chaturvedi S, Chaturvedi V, Chen M,

1050        Chowdhary A, Colom M-F, Cornely OA, Crous PW, Cuétara MS, Diaz MR, Espinel-

1051        Ingroff A, Fakhim H, Falk R, Fang W, Herkert PF, Rodríguez CF, Fraser JA, Gené J,

1052        Guarro J, Idnurm A, Illnait-Zaragozi M-T, Khan Z, Khayhan K, Kolecka A, Kurtzman

1053        CP, Lagrou K, Liao W, Linares C, Meis JF, Nielsen K, Nyazika TK, Pan W,

1054        Pekmezovic M, Polacheck I, Posteraro B, Telles F de Q, Romeo O, Sánchez M,

1055        Sampaio A, Sanguinetti M, Sriburee P, Sugita T, Taj-Aldeen SJ, Takashima M, Taylor

1056        JW, Theelen B, Tomazin R, Verweij PE, Wahyuningsih R, Wang P, Boekhout T. 2017.

1057        Importance of resolving fungal nomenclature: the case of multiple pathogenic species in

1058        the *Cryptococcus* genus. mSphere 2:e00238-17.

1059     34. Scott L. 1990. Hyrax (Procaviidae) and dassie rat (Petromuridae) middens in

1060        palaeoenvironmental studies in Africa, p. 408–427. *In* Packrat Middens: The Last

1061        40,000 Years of Biotic Change. University of Arizona Press, Tucson.

1062     35. Chase BM, Scott L, Meadows ME, Gil-Romera G, Boom A, Carr AS, Reimer PJ, Truc L,

1063        Valsecchi V, Quick LJ. 2012. Rock hyrax middens: A palaeoenvironmental archive for

1064        southern African drylands. Quat Sci Rev 56:107–125.

1065     36. Nielsen K, Obaldia ALD, Heitman J. 2007. *Cryptococcus neoformans* mates on Pigeon

1066        guano: Implications for the realized ecological niche and globalization. Eukaryot Cell

1067        6:949–959.

1068     37. Taylor JW, Barker BM. 2019. The endozoan, small-mammal reservoir hypothesis and the

1069        life cycle of Coccidioides species. Med Mycol 57:S16–S20.

34

1070    38. Springer DJ, Billmyre RB, Filler EE, Voelz K, Pursall R, Mieczkowski PA, Larsen RA,

1071        Dietrich FS, May RC, Filler SG, Heitman J. 2014. *Cryptococcus gattii* VGIII isolates

1072        causing infections in HIV/AIDS patients in Southern California: identification of the

1073        local environmental source as Arboreal. PLoS Pathog

1074        10:https://doi.org/10.1371/journal.ppat.1004285.

1075    39. Litvintseva AP, Mitchell TG. 2009. Most environmental isolates of *Cryptococcus*

1076        *neoformans var. grubii* (serotype A) are not lethal for mice. Infect Immun 77:3188–

1077        3195.

1078    40. Ngamskulrungroj P, Chang Y, Sionov E, Kwon-Chung KJ. 2012. The primary target

1079        organ of *Cryptococcus gattii* is different from that of *Cryptococcus neoformans* in a

1080        murine model. mBio 3:e00103-12.

1081    41. Davis MJ, Moyer S, Hoke ES, Sionov E, Mayer-Barber KD, Barber DL, Cai H, Jenkins

1082        L, Walter PJ, Chang YC, Kwon-Chung KJ. 2019. Pulmonary iron limitation induced by

1083        exogenous type I IFN protects mice from *Cryptococcus gattii* independently of T Cells.

1084        mBio 10:e00799-19.

1085    42. Dykstra MA, Friedman L, J W Murphy. 1977. Capsule size of *Cryptococcus neoformans*:

1086        control and relationship to virulence. Infect Immun 16:129–135.

1087    43. Gomez-Lopez A, Zaragoza O, Dos Anjos Martins M, Melhem MC, Rodriguez-Tudela

1088        JL, Cuenca-Estrella M. 2008. *In vitro* susceptibility of *Cryptococcus gattii* clinical

1089        isolates. Clin Microbiol Infect Off Publ Eur Soc Clin Microbiol Infect Dis 14:727–730.

1090    44. Khan ZU, Randhawa HS, Kowshik T, Chowdhary A, Chandy R. 2007. Antifungal

1091        susceptibility of *Cryptococcus neoformans* and *Cryptococcus gattii* isolates from

1092      decayed wood of trunk hollows of *Ficus religiosa* and *Syzygium cumini* trees in north-

1093      western India. J Antimicrob Chemother 60:312–316.

1094   45.  https://www.unaids.org/en/regionscountries/countries/zambia.

1095   46. Farrer RA, Ford CB, Rhodes J, Delorey T, May RC, Fisher MC, Cloutman-Green E,

1096      Balloux F, Cuomo CA. 2018. Transcriptional heterogeneity of *Cryptococcus gattii* VGII

1097      compared with non-VGII lineages underpins key pathogenicity pathways. mSphere 3.

1098   47. Litvintseva AP, Carbone I, Rossouw J, Thakur R, Govender NP, Mitchell TG. 2011.

1099      Evidence that the human pathogenic fungus *Cryptococcus neoformans var. grubii* may

1100      have evolved in Africa. PloS One 6:e19688.

1101   48. Randhawa HS, Kowshik T, Khan ZU. 2005. Efficacy of swabbing versus a conventional

1102      technique for isolation of *Cryptococcus neoformans* from decayed wood in tree trunk

1103      hollows. Med Mycol 43:67–71.

1104   49. Fujimura H, Sakuma Y. 1993. Simplified isolation of chromosomal and plasmid DNA

1105      from yeasts. BioTechniques 14:538–540.

1106   50. Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. 2017. Canu:

1107      scalable and accurate long-read assembly via adaptive k-mer weighting and repeat

1108      separation. Genome Res 27:722–736.

1109   51. Delcher AL, Kasif S, Fleischmann RD, Peterson J, White O, Salzberg SL. 1999.

1110      Alignment of whole genomes. Nucleic Acids Res 27:2369–2376.

1111   52. Li H. 2013. Aligning sequence reads, clone sequences and assembly contigs with BWA-

1112      MEM. ArXiv13033997 Q-Bio.

1113    53. Loman NJ, Quick J, Simpson JT. 2015. A complete bacterial genome assembled *de novo*

1114        using only nanopore sequencing data. Nat Methods 12:733–735.

1115    54. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, Cuomo CA, Zeng

1116        Q, Wortman J, Young SK, Earl AM. 2014. Pilon: an integrated tool for comprehensive

1117        microbial variant detection and genome assembly improvement. PloS One 9:e112963.

1118    55. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin VM,

1119        Nikolenko SI, Pham S, Prjibelski AD, Pyshkin AV, Sirotkin AV, Vyahhi N, Tesler G,

1120        Alekseyev MA, Pevzner PA. 2012. SPAdes: a new genome assembly algorithm and its

1121        applications to single-cell sequencing. J Comput Biol J Comput Mol Cell Biol 19:455–

1122        477.

1123    56. Lukashin AV, Borodovsky M. 1998. GeneMark.hmm: new solutions for gene finding.

1124        Nucleic Acids Res 26:1107–1115.

1125    57. Bairoch A, Apweiler R. 2000. The SWISS-PROT protein sequence database and its

1126        supplement TrEMBL in 2000. Nucleic Acids Res 28:45–48.

1127    58. Kanehisa M, Goto S. 2000. KEGG: kyoto encyclopedia of genes and genomes. Nucleic

1128        Acids Res 28:27–30.

1129    59. Finn RD, Clements J, Eddy SR. 2011. HMMER web server: interactive sequence

1130        similarity searching. Nucleic Acids Res 39:W29–W37.

1131    60. Finn RD, Bateman A, Clements J, Coggill P, Eberhardt RY, Eddy SR, Heger A,

1132        Hetherington K, Holm L, Mistry J, Sonnhammer ELL, Tate J, Punta M. 2014. Pfam: the

1133        protein families database. Nucleic Acids Res 42:D222–D230.

1134    61. Lowe TM, Eddy SR. 1997. tRNAscan-SE: a program for improved detection of transfer

1135        RNA genes in genomic sequence. Nucleic Acids Res 25:955–964.

1136    62. Lagesen K, Hallin P, Rødland EA, Staerfeldt H-H, Rognes T, Ussery DW. 2007.

1137        RNAmmer: consistent and rapid annotation of ribosomal RNA genes. Nucleic Acids

1138        Res 35:3100–3108.

1139    63. Haft DH, Selengut JD, White O. 2003. The TIGRFAMs database of protein families.

1140        Nucleic Acids Res 31:371–373.

1141    64. Conesa A, Götz S, García-Gómez JM, Terol J, Talón M, Robles M. 2005. Blast2GO: a

1142        universal tool for annotation, visualization and analysis in functional genomics research.

1143        Bioinforma Oxf Engl 21:3674–3676.

1144    65. Petersen TN, Brunak S, von Heijne G, Nielsen H. 2011. SignalP 4.0: discriminating

1145        signal peptides from transmembrane regions. Nat Methods 8:785–786.

1146    66. Krogh A, Larsson B, von Heijne G, Sonnhammer EL. 2001. Predicting transmembrane

1147        protein topology with a Hidden Markov Model: application to complete genomes. J Mol

1148        Biol 305:567–580.

1149    67. Farrer RA, Voelz K, Henk DA, Johnston SA, Fisher MC, May RC, Cuomo CA. 2016.

1150        Microevolutionary traits and comparative population genomics of the emerging

1151        pathogenic fungus *Cryptococcus gattii*. Phil Trans R Soc B 371:20160021.

1152    68. Engelthaler DM, Hicks ND, Gillece JD, Roe CC, Schupp JM, Driebe EM, Gilgado F,

1153        Carriconde F, Trilles L, Firacative C, Ngamskulrungroj P, Castañeda E, Lazera M dos

1154        S, Melhem MSC, Pérez-Bercoff Å, Huttley G, Sorrell TC, Voelz K, May RC, Fisher

1155        MC, Thompson GR, Lockhart SR, Keim P, Meyer W. 2014. *Cryptococcus gattii* in

1156        North American Pacific Northwest: whole-population genome analysis provides

1157        insights into species evolution and dispersal. mBio 5:e01464-14.

1158    69. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G,

1159        Durbin R, 1000 Genome Project Data Processing Subgroup. 2009. The Sequence

1160        Alignment/Map format and SAMtools. Bioinforma Oxf Engl 25:2078–2079.

1161    70. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K,

1162        Altshuler D, Gabriel S, Daly M, DePristo MA. 2010. The Genome Analysis Toolkit: a

1163        MapReduce framework for analyzing next-generation DNA sequencing data. Genome

1164        Res 20:1297–1303.

1165    71. Stamatakis A. 2006. RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses

1166        with thousands of taxa and mixed models. Bioinformatics 22:2688–2690.

1167    72. Huson DH. 1998. SplitsTree: analyzing and visualizing evolutionary data. Bioinforma

1168        Oxf Engl 14:68–73.

1169    73. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, Handsaker RE,

1170        Lunter G, Marth GT, Sherry ST, McVean G, Durbin R. 2011. The variant call format

1171        and VCFtools. Bioinformatics 27:2156–2158.

1172    74. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, Maller J, Sklar

1173        P, de Bakker PIW, Daly MJ, Sham PC. 2007. PLINK: a tool set for whole-genome

1174        association and population-based linkage analyses. Am J Hum Genet 81:559–575.

1175    75. Alexander DH, Novembre J, Lange K. 2009. Fast model-based estimation of ancestry in

1176        unrelated individuals. Genome Res 19:1655–1664.

1177    76. Roth C, Sun S, Billmyre RB, Heitman J, Magwene PM. 2018. A high-resolution map of

1178        meiotic recombination in *Cryptococcus deneoformans* demonstrates decreased

1179        recombination in unisexual reproduction. Genetics 209:567–578.

1180    77. Emms DM, Kelly S. 2015. OrthoFinder: solving fundamental biases in whole genome

1181        comparisons dramatically improves orthogroup inference accuracy. Genome Biol

1182        16:157.

1183    78. Farrer RA. 2017. Synima: a synteny imaging tool for annotated genome assemblies.

1184        BMC Bioinformatics 18:507.

1185    79. Edgar RC. 2004. MUSCLE: a multiple sequence alignment method with reduced time

1186        and space complexity. BMC Bioinformatics 5:113.

1187    80. Darriba D, Taboada GL, Doallo R, Posada D. 2011. ProtTest 3: fast selection of best-fit

1188        models of protein evolution. Bioinformatics 27:1164–1165.

1189    81. Kwon-Chung KJ, Polacheck I, Bennett JE. 1982. Improved diagnostic medium for

1190        separation of *Cryptococcus neoformans var. neoformans* (serotypes A and D) and

1191        *Cryptococcus neoformans var. gattii* (serotypes B and C). J Clin Microbiol 15:535–537.

1192    82. Kwon-Chung KJ, Polacheck I, Popkin TJ. 1982. Melanin-lacking mutants of

1193        *Cryptococcus neoformans* and their virulence for mice. J Bacteriol 150:1414–1421.

1194    83. Halliday CL, Bui T, Krockenberger M, Malik R, Ellis DH, Carter DA. 1999. Presence of

1195        alpha and a mating types in environmental and clinical collections of *Cryptococcus*

1196        *neoformans var. gattii* strains from Australia. J Clin Microbiol 37:2920–2926.

1197