1

2

**A Global Analysis of Mutations Accompanying Microevolution in the Heterozygous**

**Diploid Pathogen *Candida albicans***

5

Iuliana V. Ene[1], Rhys A. Farrer[2], Matthew P. Hirakawa[1], Kennedy Agwamba[2],

Christina A. Cuomo[2] and Richard J. Bennett[1*]

8

9

[1]Department of Molecular Microbiology and Immunology, Brown University, Providence,

Rhode Island 02912, USA

[2]Broad Institute of MIT and Harvard, Cambridge, Massachusetts 02142, USA

*corresponding author

14

15

16  **Abstract**

17  *Candida albicans* is a heterozygous diploid yeast that is a commensal of the human

18  gastrointestinal (GI) tract and a prevalent opportunistic pathogen.  Here, whole-genome

19  sequencing was performed on multiple *C. albicans* isolates passaged in different niches to

20  characterize the complete spectrum of mutations arising during microevolution.  We reveal that

21  evolution during short time-scales (<600 generations) is driven by both *de novo* base

22  substitutions and short-tract loss of heterozygosity (LOH) events.  In contrast, large-scale

23  chromosomal changes are relatively rare, although chromosome 7 trisomies repeatedly

24  emerged during passaging in one GI colonization model.  Both strain background and

25  chromosomal features affected mutational patterns, with mutation rates being greatly elevated

26  in regions adjacent to emergent LOH tracts.  Mutation rates were also elevated during host

27  infection where genomes showed strong evidence of purifying selection. These results establish

28  the genetic events driving *C. albicans* evolution and that this heterozygous diploid is extensively

29  shaped by purifying selection.

30

31

32    **Introduction**

33         Microbial evolution studies provide insights into genome dynamics and the factors that

34    shape genome evolution (1-5).  Most studies have focused on haploid or homozygous diploid

35    genomes, yet there is an increasing interest in defining genome variation in heterozygous

36    diploid genomes (4, 6-9).  Sexual reproduction often plays an important role in accelerating

37    adaptation in eukaryotic species, as meiotic recombination promotes genome rearrangements

38    and mutation rates are elevated compared to vegetative cells (5, 10, 11).  However, mitotic

39    populations can still generate diversity by a number of mechanisms including *de novo* base

40    substitutions, loss of heterozygosity (LOH) events, and insertion or deletions (indels), as well as

41    larger rearrangements and even the acquisition of chromosomal aneuploidies.  Many of these

42    events have been defined in model eukaryotes and are also recognized as major drivers of

43    somatic mosaicism and cancer development in the human genome (12, 13).

44         In this study, we provide a high-resolution picture of genome microevolution in the

45    heterozygous diploid yeast *Candida albicans*.  *C. albicans* is an opportunistic pathogen

46    responsible for a variety of debilitating mucosal infections and life-threatening systemic

47    infections (14, 15).  A common resident of the human microbiota, *C. albicans* can become

48    pathogenic in individuals that are immunocompromised or undergo prolonged antibiotic use (15-

49    17).  The genome consists of 8 chromosomes with the reference isolate SC5314 containing

50    $\sim$70,000 heterozygous positions representing $\sim$0.5% of the 14.3 Mb genome (18-20).

51    Heterozygous regions of the genome can undergo LOH which is increased in response to

52    stressful environments (21-24).  Isolates can also experience large-scale changes including the

53    acquisition of aneuploid forms (24-27).  *C. albicans* isolates therefore display extensive genomic

54    plasticity due to a variety of events including acquisition of SNPs and indels, LOH events, and

55    changes in gene and chromosome copy number (22, 24, 27-32).

56       Here, we define the complete spectrum of mutations accompanying microevolution in *C.*

57    *albicans*. We performed deep sequencing on multiple clinical isolates to precisely determine the

58    mutations arising during both *in vitro* and *in vivo* passaging. Our experiments reveal that

59    microevolution is driven by widespread, small-scale genetic changes, overwhelmingly

60    represented by *de novo* base substitutions and short-tract LOH events. In contrast, large-scale

61    genomic changes are rare, although both long-tract LOH events and the acquisition of

62    supernumerary chromosomes were observed, with the latter found to be a niche-specific

63    alteration. We also identify hypermutable domains within the genome including repetitive and

64    telomeric regions. Furthermore, we show that DNA recombination events are themselves highly

65    mutagenic and contribute to genomic variation by introducing a large number of *de novo*

66    mutations. Genetic events leading to gains and losses of heterozygosity occurred at similar

67    rates so that global heterozygosity levels were, in most cases, stably maintained throughout

68    microevolution. Finally, we demonstrate that mutational patterns reveal a dominant role for

69    purifying selection, with emergent mutations that alter protein-coding sequences often purged

70    from the genome during infection of the mammalian host.

**Results**

**Microevolution of *C. albicans* diploid genomes**

We selected four clinical isolates of *C. albicans* (SC5314, P78048, P76055 and P57055) for microevolution experiments.  These isolates belong to three major *C. albicans* clades (I, I, II and III, respectively), exhibit normal fitness and morphology, and have heterozygous diploid genomes (with no chromosomal aneuploidies).  Heterozygous positions represent 0.41% to 0.55% of these genomes (Supplementary Tables 1-3) (24).  Strains were passaged both *in vitro* and in three different murine models *in vivo* (Figure 1A).  The latter included two commensal models of gastrointestinal (GI) colonization using either a standard diet (SD) that requires antibiotics for *C. albicans* colonization (33)  or a purified diet (PD) that does not require antibiotics for stable colonization (34).  A model of systemic infection was also utilized in which fungal cells were introduced into the murine tail vein and subsequently recovered from the kidney, the major organ targeted by *C. albicans* (35).  For GI colonization, fungal cells were collected from fecal pellets after 42 days (~227 generations, *n*=2-3).  For systemic infection, five sequential passages were performed in which fungal cells were isolated from infected kidneys three days post infection and used for infection of new hosts (~240 generations, *n*=2).  For comparison, *in vitro* passaging was performed daily under standard laboratory conditions (YPD medium, 30°C) and isolates collected after 80 days (~600 generations, *n*=1-2).  The genomes of evolved isolates were analyzed by Illumina ultra-deep sequencing (average of 185X coverage, 97.7% of SC5314 assembly covered by reads, see Supplementary Table 2 and Materials and Methods).  Using high depth read alignments and stringent variant calling with Haplotype Caller (GATK) and Pilon (36), single nucleotide polymorphisms (SNPs), heterozygous positions and indels were identified for each isolate (Supplementary Table 2), and a number of these were further validated as described below.

96    **Large-scale chromosomal changes acquired during microevolution**

97        Sequence read depth across each of the 28 microevolved genomes revealed that a

98    small subset of isolates underwent changes at the chromosomal level.  Aneuploidy was

99    observed in three out of 28 evolved isolates and in each case involved chromosome (Chr) 7

100   trisomies (Figure 1B and Supplementary Figure 1A).  These aneuploid forms emerged in three

101   different strain backgrounds that were each passaged in the GI SD model, suggesting a fitness

102   benefit may be associated with Chr 7 trisomy under these growth conditions.  One of the three

103   aneuploid isolates also became monosomic for two terminal regions involving the right arm of

104   Chr 2 (1.15 Mbp region) and the left arm of Chr 3 (0.27 Mbp region) (P76055 GI SD isolate C,

105   Figure 1B and Supplementary Figure 1B).  Previous studies have similarly observed aneuploid

106   forms among natural isolates and that chromosome-level changes can arise during passaging

107   or in response to antifungal treatment (24, 27, 29, 31, 32, 37, 38).  An analysis of copy number

108   variation across smaller genomic regions (100-1000 bp windows) is included in Supplementary

109   Material and Supplementary Figures 7-8.

110

111   **Common patterns of microvariation in *C. albicans* genomes**

112       A detailed analysis of microevolved isolates was performed to determine the spectrum of

113   nucleotide changes in each genome.  A total of 564 mutations were identified across the 28

114   microevolved isolates in the four lineages.  All 564 mutations were individually evaluated using

115   IGV (39) and 63 sites were additionally verified using an allele-specific fluorescent PCR

116   technology (KASP genotyping, LGC).  Of these, 55 positions (87%) matched the genotypes

117   from genome sequencing (Supplementary Table 4).  Validated mutations included both SNPs

118   and indels in both genic and intergenic regions (Supplementary Figure 1C).

119       Mutations were subdivided into those leading to gains or losses of heterozygosity (GOH

120   and LOH, respectively).  GOH and LOH were further classified as resulting from either indels or

121    changes in SNPs. For example, *in vitro* passaging of *C. albicans* isolates for 600 generations

122    revealed a total of 31 mutations comprising 6 indels (4 insertions and 2 deletions) and 25 SNPs

123    (17 transitions and 8 transversions) across four isolates. These 31 mutations were the result of

124    8 GOH events (7 *de novo* base substitutions and 1 indel) and 11 LOH events due to

125    recombination. 19 of the mutations occurred in intergenic regions and 12 occurred in coding

126    sequences, of which there were 7 synonymous and 5 nonsynonymous mutations.

127        Microevolution consistently resulted in more SNPs (both GOH and LOH events) than

128    indels, independent of strain background or evolution niche. Thus, an average of 87.2% of

129    mutations involved SNPs and 12.8% involved indels (Figure 1C). For GOH events, the average

130    ratio of base substitutions to indels was 1:0.17, which is much lower than ratios reported for

131    *Saccharomyces cerevisiae* (~1:0.03) (40-42), suggesting that *C. albicans* experiences

132    proportionally higher rates of indels to *de novo* substitutions than *S. cerevisiae*.

133        An average of 41% of all mutations (including both GOH and LOH SNPs and indels)

134    occurred in coding regions and comprised 60.3% synonymous and 39.7% nonsynonymous

135    mutations (Figure 1D, E). Nonsynonymous mutations predicted to disrupt ORF function were

136    rare; only 14 nonsense mutations and 5 readthrough mutations occurred across all evolved

137    isolates, and 17 of the 19 mutations were the direct result of three very large LOH events that

138    occurred in two microevolved lineages (described below and Supplementary Table 5). Base

139    substitutions (GOH SNPs) in evolved isolates were the result of a higher fraction of transitions

140    (54.4%) than transversions, with a *Ts/Tv* ratio of 1.3:1, which is lower than the 2:1 ratio reported

141    for model yeast genomes (41, 43) (Supplementary Figure 1D). These mutational patterns were

142    consistent across microevolved isolates revealing that they are independent of genetic

143    background and the environment in which isolates are passaged (Figure 1C-E and

144    Supplementary Figure 1E-G).

145

146 **Purifying selection shapes the evolution of *C. albicans* genomes**

147      In the absence of bottlenecks, new mutations that have deleterious effects may be

148 purged from the population via purifying selection, and we therefore tested whether mutational

149 patterns in our dataset showed evidence for selection.  If occurring randomly, mutations will

150 accumulate in intergenic and coding regions at frequencies proportional to their representation

151 in the genome (40, 41, 44).  In our experiments, 48% of all *de novo* base substitutions (GOH

152 SNPs) and 61.5% of GOH indels were present in intergenic regions, even though these regions

153 account for only 36.2% of the genome ($P < 0.05$, Figure 2A).  Moreover, none of the indels (0/5)

154 found in coding sequences resulted in frameshifts (i.e., all were a multiple of 3 nucleotides),

155 whereas only 1/8 indels observed in intergenic regions consisted of multiples of 3 nucleotides ($P$

156 $< 0.05$, difference between intergenic and coding indels is significant using a binomial

157 distribution model, Figure 2B).  The fraction of synonymous to nonsynonymous mutations also

158 differed from that expected by chance; ~25% of coding substitutions are expected to be

159 synonymous if mutations occur randomly (40, 44) yet over 48% of base substitutions were

160 synonymous in our dataset ($P < 0.05$, Figure 2C).  This suggests that selection frequently acts

161 to limit the accumulation of mutations that alter the protein-coding sequence.

162      We also estimated the fraction of mutations impacted by selection by examining how

163 many nonsynonymous mutations would be expected during microevolution based on the

164 number of synonymous or intergenic mutations observed.  Assuming an even distribution of

165 mutational events (see Discussion), selection effectively removed an average of 71-79% of the

166 nonsynonymous mutations predicted to occur during microevolution.  Selection coefficients for

167 nonsynonymous mutations in evolved isolates averaged 0.0053 (0.0047 using estimates based

168 on intergenic mutations and 0.0059 using synonymous mutation rates).  Selection coefficients

169 were highest for isolates passaged in the two GI colonization models ($P < 0.0005$, Figure 2D).

8

170    Together, these results establish that *C. albicans* isolates display much higher

171    synonymous:nonsynonymous and intergenic:genic mutation ratios than that expected by

172    chance, implying that purifying selection removes a large fraction of the mutations impacting

173    protein-coding genes during passaging *in vitro* and *in vivo*.

174

175    **Impact of strain background and environment on *C. albicans* mutation rates**

176    Mutation rates were compared between the four clinical isolates and across culture

177    conditions to examine cell-intrinsic and cell-extrinsic factors that impact microevolution.  Strains

178    passaged *in vitro* displayed an average rate of $1.17 \times 10^{-10}$ base substitutions per base pair (bp)

179    per generation (Supplementary Figure 2A).  These *de novo* substitution rates are similar to

180    those reported for asexual populations of *S. cerevisiae* and *Schizosaccharomyces pombe* (40-

181    42, 45) (Supplementary Figure 2B).  Mutations in *C. albicans* cells passaged *in vitro* reflected

182    mutational patterns common to all microevolution experiments, with more frequent changes due

183    to SNPs than to indels, and fewer mutations affecting coding regions than expected by chance

184    (Supplementary Figure 2A).  Isolates passaged *in vitro* displayed an average LOH rate of 1.61 x

185    $10^{-10}$ per bp per generation, resulting from 2.75 LOH events per strain every 600 generations.

186    Mutation rates varied considerably depending on both the genetic background and the

187    environment.  The standard 'laboratory' strain SC5314 displayed the lowest mutation rates (both

188    for GOH and LOH events) as rates in the other three lineages were 1.3 – 5.6-fold higher (Figure

189    3A).  The environment also significantly impacted the mutation frequency; strains grown *in vivo*

190    (either in the GI or in systemic models of infection) showed GOH rates that were 6.7 – 9.6-fold

191    higher than those *in vitro* (Figure 3B).  LOH rates were also higher *in vivo* than *in vitro* (6.8 –

192    12.7-fold; Figure 3B).  Thus, *C. albicans* cells exhibit significantly higher mutation rates when

193    passaged in the host (either in systemic or GI infection models) relative to *in vitro* passaging.

9

194    In contrast to overall GOH and LOH mutations, rates of indel formation did not differ

195    significantly between experiments in different strain backgrounds or in different niches.  The

196    P76055 lineage displayed the lowest indel frequency (0.8 x 10$^{-10}$ per bp per generation)

197    compared to rates that were 3.4 - 4.4-fold higher in the other lineages (Supplementary Figure

198    2C).  Indel rates were also elevated 3.1 – 4.6-fold in the bloodstream and in GI SD infection

199    models relative to *in vitro* passaging, although these differences did not reach significance due

200    to the small number of events (Supplementary Figure 2D).  We note that precise *in vivo*

201    mutation rates are difficult to determine due to approximated generation times (see Materials

202    and Methods, Supplementary Table 3).  Previous studies estimated 0.09 generations/h during

203    systemic infection (31) and 0.14 generations/h during GI colonization (46), rates that are only 2-

204    3-fold lower than *in vitro* growth.  It is therefore possible that these *in vivo* generation times are

205    overestimates and mutation rates *in vivo* could be even higher than those presented here.

206

207    **Genome heterozygosity levels are maintained due to balanced GOH and LOH rates**

208    An important question is how do *C. albicans* strains maintain genome heterozygosity

209    levels despite frequent LOH events and the absence of conventional outcrossing (22, 24, 25,

210    47-49).  Heterozygosity patterns were compared before and after passaging and revealed that

211    LOH rates and GOH rates were often balanced in each experiment (Supplementary Figure 2E).

212    Thus, genome heterozygosity levels in the four strain backgrounds were virtually unchanged

213    following most passaging experiments, with levels within -1.5% to +2.1% of starting

214    heterozygosity levels (Figure 3C).  The exceptions to this pattern were two passaged isolates

215    that experienced very large (>0.27 Mb) LOH events ('LLOH tracts') and therefore exhibited

216    significant decreases (-8.5% and -11.3%) in heterozygous sites across their genomes.

217    Together, these results establish that mitotic recombination events (driving LOH) and GOH

218    events (resulting from both *de novo* base substitutions and indels) often occur at similar

10

219    frequencies and that, in the absence of large LOH events, genome heterozygosity levels can be

220    stably maintained independent of genetic background or environment.

221

222    **Defining hypervariable regions within the *C. albicans* genome**

223            To determine the impact of genomic context on mutation rates, we compared the

224    frequency of mutations arising at a number of chromosomal features including centromeres,

225    terminal chromosome regions, subtelomeric *TLO* genes, *ALS* (agglutinin-like sequence) genes,

226    other glycosylphosphatidylinositol (GPI)-linked genes, and annotated DNA repeat regions (see

227    Supplementary Tables 6-7).  Several of these features have been associated with higher

228    mutation rates in *C. albicans* and other model organisms (4, 50-56).  For each of these features,

229    we examined the frequency of SNP and indel mutations relative to the genome average.  We

230    found that mutation rates were significantly elevated at the ends of chromosomes (6.7 fold

231    increase within the 10 kb terminal regions), as well as in the subtelomeric *TLO* genes (48.8 fold

232    increase) of a subset of lineages (Figure 3D and Supplementary Figure 3A,B).  This is in line

233    with previous observations that *C. albicans* telomeric and subtelomeric regions are highly

234    dynamic (56-58).  In contrast, centromeric regions did not display altered mutation rates relative

235    to the genome average (Supplementary Figure 3C).  This contrasts with the high mutation rates

236    observed at *S. cerevisiae* centromeres (54), however the 3-4.5 kb regional centromeres present

237    in *C. albicans* are much larger than the point centromeres (often <400 bp) found in this model

238    yeast (59, 60).

239            Repeat sequences are common in fungal genomes and have been associated with

240    mobile elements and gene regulation (61, 62).  However, studies in model yeast have tended to

241    exclude analysis of repeat regions to simplify genome analyses (40, 45).  *C. albicans* is unusual

242    among *Candida* species in that it contains 9 large Major Repeat Sequence (MRS) elements that

243    span ~1.7% of the genome and are linked to chromosome translocations and chromosome

244     length polymorphisms (62-66).  We found that microevolved *C. albicans* lineages exhibited

245     significant differences in mutation rates between repeat and non-repeat regions.  For example,

246     both MRS elements and long terminal repeat (LTR) retrotransposons showed mutation rates

247     that were 13.1-fold higher than the genome average (Figure 3D and Supplementary Figure 3D).

248     Importantly, a number of these mutations were validated by KASP analysis confirming that they

249     arose during microevolution (Supplementary Table 4).  Mutation rates in repeat regions were

250     also significantly higher in strains passaged *in vivo* than *in vitro* (Supplementary Figure 3D).

251          Genes encoding GPI-linked cell wall proteins and *ALS* family genes are rich in internal

252     tandem repeats that can vary in number and thereby contribute to allelic diversity and

253     phenotypic variation (67-70).  In line with these observations, we found that both GPI and *ALS*

254     gene families accumulated mutations at much higher rates than the genome average (~8.6-fold

255     increase, Figure 3D and Supplementary Figure 3E).  These results establish that numerous

256     chromosomal features including genomic repeats, telomeric regions, *ALS* genes and GPI-linked

257     cell wall genes undergo evolution faster than the rest of the genome.

258          We also compared the distribution of mutations arising in heterozygous versus

259     homozygous regions of the genome, while noting that LOH events in *C. albicans* can often

260     promote adaptation (30, 71, 72).  Heterozygous and homozygous regions were mapped in each

261     of the four parental strains based on the density of heterozygous positions per 5 kb window

262     (24).  Using this metric, heterozygous regions in the four parental isolates varied between

263     69.5% and 84.2% of the genome (Supplementary Table 8).  Comparison of the frequency of

264     mutations between heterozygous and homozygous regions revealed that these regions

265     accumulated mutations in line with their relative abundance (Figure 3D).  For example, 84.2% of

266     the P78048 genome is represented by heterozygous regions and 80.2% of all mutations in

267     evolved derivatives of this lineage occurred in heterozygous regions.  As LOH events are likely

268     biased by the higher frequency of heterozygous positions in HET regions than in HOM regions,

12

269    we repeated this analysis using only GOH events (base substitutions and indels). We again

270    found that GOH events accumulated in HET regions at similar levels to their proportion in the

271    genome (Supplementary Figure 3F). We therefore did not observe bias in the pattern of *de*

272    *novo* mutations towards either heterozygous or homozygous regions of the genome.

273

274    **Microevolution is punctuated by frequent short-tract LOH events and small indels**

275         A wide variety of LOH events have been described in *C. albicans,* with elevated LOH

276    rates observed during exposure to stress, antifungal treatment, DNA damage, and host passage

277    (22, 31, 32, 37, 73). To provide a global picture of these events during microevolution, we

278    divided LOH events into three categories based on length (in kb) and the number of

279    heterozygous positions affected: (1) microLOH (mLOH) events that involved loss of single

280    heterozygous positions, (2) short-tract LOH (SLOH) events that involved loss of 2 or more

281    heterozygous positions and covered small genomic regions (≤10 kb), and (3) long-tract LOH

282    (LLOH) events that were >10 kb and affected hundreds of heterozygous positions (Figure 4A).

283    The relative frequency of mLOH and SLOH events observed during microevolution was similar

284    across experiments, with the minimum sizes of these events ranging between 1 and 3090 bp

285    ($L_{min}$ size, Figure 4B and Materials and Methods). Thus, isolates underwent an average of

286    52.3% mLOH and 46.4% SLOH events during passaging (Figure 4C,D). Analysis of the

287    average size of LOH tracts revealed that $L_{avg}$ varied between 222 bp (P57055) and 889 bp

288    (SC5314) for the four strain backgrounds and impacted between 1.6 and 3.2 heterozygous

289    positions (when excluding LLOH events, Supplementary Figure 4A,B). In contrast to frequent

290    mLOH and SLOH events, long-tract LOH events occurred in only 2 passaged isolates (P76055

291    and P78048 grown in the GI SD model) and involved tracts of 273-1230 kb that extended to the

292    ends of the chromosomes.

293     We also analyzed the 41 indels that occurred in the 28 evolved isolates (Figure 4E,F),

294     finding that they were represented by similar numbers of insertions (21 events) and deletions

295     (20 events).  In *S. cerevisiae*, indels were biased towards insertions in haploid lines and towards

296     deletions in diploid lines (40, 41).  Indel sizes averaged 3.3 bp for *in vitro*-evolved isolates

297     whereas indels were 2-3-fold larger for *in vivo*-passaged isolates (Figure 4F), although no

298     significant differences in size were found between different lineages or different niches either for

299     LOH events or for indels ($P > 0.05$).  Together, these results provide the first comprehensive

300     analysis of LOH and indel events in *C. albicans* and highlight that short-tract LOH events, most

301     of which involve homozygosis of single heterozygous positions, are the most common LOH

302     event occurring during microevolution.

303

**304     LOH events are overrepresented in repeat regions and telomeric regions**

305     The frequency and distribution of LOH events arising within the *C. albicans* genome

306     were examined for potential relationships with underlying chromosomal features.  All LOH

307     events (196 tracts) were mapped along the genome (Figure 5A) and their frequency determined

308     per 0.2 Mb window (Supplementary Figure 4C).  Mapping the distance of each LOH to the

309     closest chromosomal feature revealed that a high proportion of LOH tracts (21%) arose either

310     within MRS regions or in the 1 kb tracts adjacent to MRS or telomeric regions (Figure 5B and

311     Supplementary Figure 5C).  MRS elements were the main hotspot for these recombination

312     events, with the start sites for 12% of all LOH tracts being located at these elements (Figure

313     5B).  Not all MRS tracts showed the same propensity for recombination; MRS6 and MRS7b

314     displayed the highest LOH frequencies and MRS2, 3, and 5 displayed the lowest LOH

315     frequencies (Supplementary Figure 4C).  In contrast, no LOH events were detected within 1 kb

316     of the centromeres (Figure 5B).  These findings establish that MRS and telomeric regions are

317     hotspots for recombination in *C. albicans* and thereby promote genetic variation.

14

318

**The emergence of new heterozygous SNPs is detected within large LOH tracts**

320        Three very large LOH tracts (0.27-1.23 Mbp) were formed during passaging, and these

321    involved two isolates cultured in the GI (antibiotic-treated) model with the standard mouse diet

322    (SD).  LOH events involved the terminal regions of Chr 2 and Chr 3 in P76055 (GI SD C) and

323    Chr R in P78048 (GI SD B; Figure 5A and Supplementary Figure 5A).  LOH occurred via

324    truncation of chromosome arms in P76055 GI SD C, as the resulting LOH tracts were

325    monosomic (displayed half the read coverage for LOH regions on both Chr 2 and Chr 3,

326    Supplementary Figure 1B).  In contrast, LOH likely involved break-induced replication (BIR) or

327    inter-homolog crossing-over in P78048 GI SD B, as the chromosome was still disomic for the

328    emergent homozygous region (Supplementary Figure 1B).  Analyses revealed that these large

329    events led to LOH of thousands of heterozygous positions (as well as tens of indels) that were

330    present in the parental genomes.  In total, the three LLOH led to the homozygosis of 5,419 sites

331    in coding regions, 2135 of which resulted in nonsynonymous changes, 14 produced nonsense

332    mutations and three resulted in readthrough mutations (Supplementary Figure 5B and

333    Supplementary Table 5).  Interestingly, the LLOH region in P78048 GI SD B showed the

334    reemergence of heterozygous positions due to *de novo* base substitutions (GOH SNPs) within

335    the tract that had undergone LOH.  In fact, a 6.6-fold higher rate of GOH events was detected

336    within this LLOH tract than was evident in the rest of the P78048 GI SD B genome (Figure 5C).

337    The high frequency of GOH mutations within the LOH tract is consistent with a high rate of *de*

338    *novo* base substitutions emerging during BIR, as shown for *S. cerevisiae* where BIR was highly

339    mutagenic (74).

340

**Impact of LOH events on mutational patterns**

342       Chromosomal crossovers can induce *de novo* mutations in DNA regions close to the

343    crossover in a wide variety of species (4, 50, 51, 75, 76). We therefore examined whether the

344    regions flanking emergent LOH tracts in microevolved *C. albicans* isolates showed altered

345    mutation rates relative to the genome average. Strikingly, sites adjacent to LOH tracts

346    appeared highly enriched for mutations; 44 out of 136 GOH events (32 *de novo* base

347    substitutions and 12 indels) were located within 500 bp of emergent LOH tracts (Figure 5D). In

348    fact, 36 of these GOH events were located within just 100 bp of LOH tracts. Thus, 32.4% of all

349    GOH events (26% of base substitutions and 92.3% of indels) were found in regions close to

350    new LOH tracts, even though these regions represent only ~1.3% of the genome. GOH rates in

351    LOH-adjacent regions (defined as 500 bp up/down of LOH tract) were therefore 840-fold higher

352    than in the rest of the genome (Figure 5E), and were significantly higher in both systemic and GI

353    SD infection models than in other environments (Supplementary Figure 5C). The distribution of

354    GOH events adjacent to LOH tracts differed from that in the rest of the genome. For example,

355    indels were highly enriched in LOH-adjacent regions, whereas GOH mutations arising in *TLO*

356    genes, centromeric, and telomeric regions were not closely associated with LOH events

357    (Supplementary Figure 5D).

358       We further note that *de novo* base substitutions in LOH-adjacent regions showed a

359    *Ts/Tv* ratio of 1.13:1 compared to a genome average of 1.27:1. This is consistent with

360    increased transversion rates resulting from translesion polymerases acting to repair DNA

361    lesions at or close to recombination tracts (77, 78). In addition, both homologous recombination

362    and non-homologous end-joining are considered to be error-prone mechanisms that can

363    introduce indels close to the DNA break site (79). Our data now reveal that approximately a

364    third of all GOH indels and substitutions in *C. albicans* arise in regions flanking LOH events and

365    that this is likely due to highly mutagenic DNA repair mechanisms.

366 **Discussion**

367       This study defines the spectrum of mutations that emerge in heterozygous diploid

368 genomes of *C. albicans* during microevolution, including a comparison of mutational patterns

369 during *in vitro* culture with those that occur during infection of a mammalian host.  Numerous

370 studies have established that *C. albicans* exhibits extensive genomic plasticity, from variation at

371 the level of single-nucleotide polymorphisms to changes at the whole-chromosome level (24,

372 25, 30, 32, 64, 80).  However, this work provides the first comprehensive picture of the genetic

373 changes accompanying microevolution in this important pathogen.

374

375 **Global patterns of mutation in *C. albicans*.**  Microevolution resulted in similar mutational

376 patterns regardless of the strain background or culture niche.  In each case, microevolution was

377 driven almost exclusively by multiple, small-scale changes in heterozygous polymorphisms

378 (87.2%) and indels (12.8%).  This reveals that 'micro-scale changes' are by far the most

379 frequent events arising in the *C. albicans* genome.  We establish that *C. albicans* displays an

380 average *de novo* base-substitution rate of $1.17 \times 10^{-10}$ per bp per generation during *in vitro*

381 passaging.  This is the first genome-wide estimate of *C. albicans* mutation rates and is close to

382 those reported for mitotically dividing cells in the model yeast *S. cerevisiae* and *S. pombe* (8, 40,

383 41, 45).  *C. albicans* therefore exhibits a *de novo* substitution rate similar to that of haploid or

384 homozygous diploid yeast genomes.  Critically, we show that microvariation in *C. albicans* is

385 equally driven by LOH events, as these recombination events occur at frequencies ($1.61 \times 10^{-10}$

386 per bp per generation) that are close to those of *de novo* substitution rates and impact a similar

387 number of nucleotide positions.

388

389 **Genome architecture and environmental pressures impact *C. albicans* microevolution.**

390 While overall mutational patterns were similar between microevolution lineages, *de novo*

17

391    substitution and LOH rates varied significantly between different strain backgrounds and

392    environments.  For example, mutation rates varied by up to 5.6-fold between strains, although

393    no obvious genetic differences were found (such as 'mutator' genotypes due to disruptions in

394    DNA repair genes) that could account for these differences (see Supplementary Material).  The

395    niche in which strains were evolved had an even bigger impact on mutation rates; strains

396    passaged *in vivo* showed up to 12.7-fold higher mutation rates than those passaged *in vitro*.  *C.*

397    *albicans* therefore experiences environment- or stress-induced mutagenesis as demonstrated

398    for a number of bacterial, fungal, plant and human studies (81-83).  In support of this, *C.*

399    *albicans* was previously shown to undergo stress-induced LOH events *in vitro* (22), and certain

400    long-tract LOH events were more frequent during bloodstream passage than during *in vitro*

401    culture (31).

402         Our studies also establish that a number of chromosomal features impact *C. albicans*

403    mutation rates.  Mutation rates were higher in repeat regions, telomeric/subtelomeric regions,

404    and in genes encoding GPI-linked cell wall proteins (including *ALS* family genes) than in the rest

405    of the genome (see schematic in Figure 6).  These results are consistent with multiple reports

406    linking higher mutation rates within repetitive and telomere-proximal regions of the *C. albicans*

407    genome (56, 62, 64, 67, 69, 84).

408

409    **Purifying selection acts on emerging mutations.**  Mutations accumulated at significantly

410    higher rates in intergenic regions than in coding regions, and the ratio of synonymous to

411    nonsynonymous mutations was also greater than that expected by chance.  We estimate that

412    71-79% of nonsynonymous mutations were effectively removed from the population by selection

413    based on the number of synonymous substitutions and the number of substitutions observed in

414    intergenic regions.  These results imply that purifying selection frequently acts to remove

415    fitness-reducing mutants from the population.  Natural isolates of *S. cerevisiae* also show a

416    significant bias towards intergenic over genic SNPs (85), and more than a third of

417    nonsynonymous mutations were implicated as being deleterious in one study (86).  The current

418    study provides striking evidence for purifying selection acting broadly on the diploid *C. albicans*

419    genome even over relatively short evolutionary periods.

420

421    **Aneuploid forms frequently arise during passaging in one microevolution niche.**

422    Aneuploid forms have frequently been described in *C. albicans* (22, 24, 25, 87), yet only 3/28

423    microevolved lineages became aneuploid during our studies.  In each case, diploid strains

424    acquired a third copy of Chr 7 and these aneuploidies emerged in three different strain

425    backgrounds during passaging in the GI tract using a standard mouse diet together with

426    antibiotics.  This suggests that being trisomic for Chr 7 provides a significant advantage under

427    these growth conditions.  To our knowledge, this is the first time this trisomy has been

428    associated with passaging of *C. albicans* in the host.  Interestingly, chromosome 7 trisomies did

429    not emerge in isolates passaged in an alternative GI model that did not involve the use of

430    antibiotics.  This suggests that this trisomy does not enhance growth in the GI tract *per se* but

431    could provide a specific advantage to attributes of the GI environment when mice are on the

432    standard diet.

433

434    **MicroLOH events are a major driver of genome dynamics.**  LOH has long been recognized

435    as an important mechanism for introducing diversity into *C. albicans* populations (21, 22, 24, 25,

436    30), although a global analysis of dynamic events had not previously been performed.  Studies

437    have examined the length of LOH tracts in diploid *S. cerevisiae* strains and showed that mitotic

438    tracts are generally longer than meiotic tracts, with the former averaging 2 - 12 kb (88, 89).  A

439    recent study examined genome-wide recombination events in mitotic *S. cerevisiae* cells and

440    showed that LOH tracts range from <100 bp to >100 kb with small LOH tracts (<1 kb) attributed

441    to local gene conversions, although many of the smallest LOH events were excluded from this

442    analysis (8).  The current study now defines the total spectrum of LOH events occurring during

443    *C. albicans* microevolution.  *C. albicans* LOH rates were $1.61 \times 10^{-10}$ per bp per generation (or

444    $4.5 \times 10^{-3}$ per cell division) which are slightly lower than previous estimates of $1.3 - 2.9 \times 10^{-9}$ per

445    bp per generation based on events at three select loci (66, 90).  However, we note that LOH

446    rates varied considerably between different genomic regions, with MRS and telomeric regions

447    representing relative hotspots for LOH.

448         We reveal that the majority of LOH events in *C. albicans* involve very short microLOH

449    tracts (mLOH, estimated $L_{avg}$ size = 368 bp).  Indeed, over half (52%) of all LOH events

450    impacted only a single heterozygous position and, critically, a number of these events were

451    validated by KASP genotyping.  In fact, when examining all genetic changes accrued during

452    microevolution, >30% of these changes were due to LOH at single heterozygous positions

453    revealing that these represent a very high frequency event in the *C. albicans* genome.

454    Consistent with our data, experiments studying the repair of DNA double-strand breaks in *C.*

455    *albicans* showed frequent short-tract LOH events via gene conversion, and only rarely were

456    long-tract LOH events observed due to BIR (Break-Induced Replication) or reciprocal

457    recombination (73).  Similarly, recent analysis of passaged lineages in the water flea *Daphnia*

458    *pulex* also found that short-tract LOH events (median ~221 bp) were prevalent (7).  We

459    therefore suggest that short-tract LOH events represent a common occurrence in heterozygous

460    diploid genomes but will have been missed by studies that lack nucleotide-level resolution.

461         In contrast to frequent microLOH events, large LOH tracts were rarely observed in our

462    experiments and involved only two isolates passaged in the GI SD model.  These LOH events

463    involved long DNA tracts (0.27-1.23 Mb) that extended to the ends of the chromosomes.

464    Previous studies also detected large LOH events in *C. albicans* strains grown both *in vitro* and

465    *in vivo* (21, 22, 24, 25, 73, 91).  We note that while large-scale chromosomal changes are

20

466    relatively rare, these impact a large number of genes and are therefore the most likely to have

467    phenotypic consequences.  Overall, the three large LOH events identified here led to

468    homozygosis of over 10 thousand heterozygous positions, resulting in 2,135 nonsynonymous

469    changes, 14 nonsense mutations and 3 readthrough mutations.  This reveals that large LOH

470    tracts drive extensive genotypic changes but may also be heavily selected against given the

471    large number of positions impacted by such events.

472        We were also surprised to find that global LOH rates were balanced by equivalent rates

473    of *de novo* GOH mutations during microevolution, regardless of genetic background or evolution

474    niche.  Because of this balance, overall heterozygosity levels were stably maintained (± 2%) in

475    the majority of passaging experiments.  This finding sheds light on an important question in *C.*

476    *albicans* – how do strains maintain genome heterozygosity levels in the face of frequent LOH

477    events?  The observation that *de novo* mutation rates often match LOH rates indicates that

478    genomes can frequently remain heterozygous even in the absence of outcrossing events.

479

480    **An association between recombination events and *de novo* mutations**.  We found that

481    there was a striking correlation between *de novo* mutations (both base substitutions and indels)

482    and their proximity to recombination events in *C. albicans*.  Specifically, GOH rates were

483    elevated 840-fold within 500 bp of emergent LOH tracts relative to the genome average.

484    Consequently, a third of all GOH events (substitutions and indels) were in regions flanking new

485    LOH tracts, despite these tracts representing only ~1.3% of the genome.  This phenomenon

486    was observed independent of strain background and was most evident during host passage

487    where LOH rates were higher.  The high rate of mutations adjacent to LOH tracts is likely due to

488    LOH being mutagenic and introducing *de novo* mutations into neighboring regions of the

489    genome.  This is consistent with DNA double-strand break repair processes being both

490    recombinogenic and mutagenic, as has been described in several cell types (51, 74, 92).

21

491     Increased rates of transversions in these mutated regions support the activity of translesion

492     polymerases acting to repair DNA lesions at these sites.  However, this is the first evidence that

493     recombination events in *C. albicans* introduce *de novo* mutations during the repair process.

494     Thus, LOH is a stress-inducible event in *C. albicans (22)* and also introduces additional

495     mutations into the genome, both of which will accelerate adaptation.

496

497     **Concluding remarks**.  This study provides a high-resolution analysis of the spectrum of

498     mutations accumulating in a heterozygous diploid pathogen.  We demonstrate that both cell-

499     intrinsic properties (e.g., strain background, repetitive chromosomal features) and cell-extrinsic

500     factors (e.g., *in vivo* versus *in vitro* passage) impact the frequency and distribution of genetic

501     fluctuations.  Frequent micro-scale changes (predominantly *de novo* substitutions and short-

502     tract LOH events) and occasional larger-scale rearrangements (long-tract LOH or chromosomal

503     aneuploidies) determine genome dynamics.  Furthermore, purifying selection plays a dominant

504     role in dictating which genetic changes are retained during evolution.  Our results provide a

505     detailed picture against which genomic changes in other heterozygous diploid species can be

506     evaluated, and establish the foundation for understanding how *C. albicans* can adapt to a wide

507     variety of distinct host niches.

508

516  of Allergy and Infectious Diseases, National Institutes of Health, Department of Health and

517  Human Services, under Contract No.:HHSN272200900018C and Grant Number U19AI110818

518  to the Broad Institute.

519

520  **Author contributions**.  I.V.E. and R.J.B. planned the experiments, I.V.E. and M.P.H performed

521  the experiments, I.V.E, R.A.F., K.A. and C.A.C. performed the bioinformatics analyses.  I.V.E

522  and R.J.B. drafted the manuscript with contributions from R.A.F. and C.A.C.

523

524

**Materials and Methods**

**Strains and growth conditions.** The *C. albicans* strains used in this study are listed in Supplementary Table 1. Unless otherwise stated, strains were grown at 30°C in YPD medium (93). For *in vitro* evolution experiments, cultures were serially diluted (1/100) every day for 80 days (bottlenecks every 6.8 generations) and cells collected once a week. *In vitro* isolates were collected as a pool and used to prepare genomic DNA.

**Determination of doubling time and generation times.** For determination of doubling times *C. albicans* strains were grown in YPD at 30°C for 18-24 h and cell densities were recorded every 10-15 min in a Biotek Synergy HT plate reader/incubator. Exponential growth intervals were selected for doubling time estimates. Doubling times (D) were calculated using the formula $D = t / \log_2(N_t/N_0)$ with t = duration of growth interval, $N_0$ = number of cells at the start of the selected interval and $N_t$ = number of cells at the end of the selected interval. The four original starting strains displayed similar doubling times (~1.5 h, Supplementary Table 3). The number of generations during *in vitro* microevolution experiments was calculated using the starting and final cell densities as $G = \log_2(N_t/N_0)$ with $N_0$ = number of cells at the start of the culture and $N_t$ = number of cells at the end of the culture (~6.8 generations/day, Supplementary Table 3). The reported doubling and generation times represent the averages of three biological replicates (with 2 technical replicates performed for each biological replicate).

**Murine experiments.** For animal infections, 6-7-week-old female BALB/c mice (~18 g) from Charles River Laboratories were housed together with free access to food and water. For systemic infection, *C. albicans* cells were grown overnight in YPD medium at 30°C, washed in phosphate-buffered saline (PBS), and for each passage 3 mice were infected via the tail-vein with a total inoculum of $6 \times 10^5$ cells in a 200 $\mu$l volume. Mice were monitored for signs of

24

550    infection and their weights, posture and motility scored daily.  3 days post-infection the mice

551    were euthanized and fungal cells isolated from infected kidneys using a PBS solution

552    supplemented with an antibiotic mixture (500 µg/mL penicillin, 500 µg/mL ampicillin, 250 µg/mL

553    streptomycin, 225 µg/mL kanamycin, 125 µg/mL chloramphenicol, and 125 µg/mL doxycycline).

554    The number of colony forming units (CFUs) from the kidneys was determined by plating cells

555    onto YPD medium.  Isolates for subsequent passages were selected by picking single colonies

556    from the mouse showing the highest virulence outcome score.  Virulence outcome scores were

557    determined by assessing kidney fungal burdens and weight changes at 72 h using the formula:

558    outcome score = log (kidney CFU/g) − (0.5 × percentage weight change) (94).  After five 3-day

559    passages, the last isolates were collected, mice were humanely sacrificed and fungal burdens

560    determined from the kidney.  Colonies from the mouse showing the highest outcome score at

561    passage 5 were used to prepare genomic DNA.

562        For GI colonization experiments, two different murine models of commensalism were

563    used.  For the standard diet (SD) model, mice were fed standard rodent chow (FormuLab 5001,

564    PMI Nutrition International) and their water was supplemented with antibiotics (1500 units/mL of

565    penicillin, 2 mg/mL of streptomycin) and 5% glucose for taste (95).  The antibiotic treatment was

566    initiated 4 days prior to infection to reduce the endogenous gastrointestinal microbiota.

567    Alternatively, mice were fed a purified diet (PD) starting 4 days prior to inoculation in which case

568    their water was not supplemented with antibiotics (34).  In both models, mice were orally

569    gavaged with a 20G x 38 mm plastic feeding tube (Instech Laboratories, Inc.) with $10^8$ C.

570    albicans cells in a 500 µl volume and continued with their respective diet and water for 6 weeks.

571    To prevent contamination between independent evolution experiments, each mouse was

572    housed in a separate cage.  Fecal samples were collected weekly and fungal cells isolated

573    using a PBS solution supplemented with antibiotics.  At 42 days, isolates were collected as a

574    pool and used to prepare genomic DNA.  After the last isolate collection, mice were humanely

25

575    sacrificed and fungal burdens determined from the GI organs (stomach, small intestine, colon

576    and caecum).

577

578    **Whole-genome sequencing and variant identification.**  To extract genomic DNA, *C. albicans*

579    isolates (Supplementary Table 1) were grown overnight in YPD at 30°C and DNA isolated from

580    ~$10^9$ cells using a Qiagen Genomic Buffer Set, a Qiagen Genomic-tip 100/G or the MasterPure

581    Yeast DNA Purification kit (Epicentre).  Each isolate was sequenced using Illumina HiSeq 2000

582    generating 101 bp paired reads.  The nuclear genome sequences and General Feature Files

583    (GFF) for *C. albicans* SC5314 reference genome (version A21-s02-m08-r01) were downloaded

584    from http://www.candidagenome.org/.  We randomly down-sampled the paired-end Illumina

585    reads for isolate SC5314 from 13 SRA runs (SRR1106648; SRR1106646; SRR1106647;

586    SRR1106651; SRR1106653; SRR1106654; SRR1106656; SRR1106658; SRR1106664;

587    SRR1106643; SRR1106645; SRR1106649; SRR1106655) to 45,130,695 paired reads (~300X

588    deep, where the range for all isolates is 70X - 547X).  Reads were aligned to the SC5314

589    reference genome assembly using Burrows-Wheeler Aligner (BWA) v0.7.4-r385 mem (96), and

590    converted to sorted BAM format using Samtools v0.1.9 (r783) (97).  The Genome Analysis

591    Toolkit (GATK) (98) v2.7-4-g6f46d11 was used to call both variant and reference bases from the

592    alignments.  Briefly, the Picard tools (http://picard.sourceforge.net/) AddOrReplaceReadGroups,

593    MarkDuplicates, CreateSequenceDictionary and ReorderSam were used to preprocess the

594    alignments.  We used GATK RealignerTargetCreator and IndelRealigner for resolving

595    misaligned reads close to indels on parental-progeny pairs of isolates to avoid discrepancies

596    between isolates.  Next, GATK Haplotype Caller and Pilon (36) (with diploid genotyper ploidy

597    setting) were run with both SNP and INDEL genotype likelihood models (GLM).  We then

598    merged and sorted all the calls from Haplotype Caller, and ran VariantFiltration with the

599    following filters "QD < 2.0, FS > 60.0, MQ < 40.0, MQRankSum < -12.5, ReadPosRankSum < -

26

600    8".  Next, we removed any base that had less than a minimum genotype quality of 50, or a

601    minimum depth of 20.  Finally, we removed any positions that were called by both GLMs (i.e.,

602    incompatible indels and SNPs), any marked as "LowQual" by GATK, nested indels, or sites that

603    did not include a PASS flag.  Similar filtering was performed for Pilon calls, removing low quality

604    sites and setting a minimum depth of 20.  All mutations in evolved isolates were visually

605    inspected using IGV (http://software.broadinstitute.org/software/igv/).  Identical mutations that

606    were present in multiple isolates from the same lineage were removed from the analyses under

607    the likelihood that they had been present in the parental strains.  The final base calls covered

608    >97% of the genome for any given isolate (Supplementary Table 2).  We then categorized every

609    single base between a parent and progeny (summarized in Supplementary Table 2), and

610    annotated those changes using the GFF (VCFannotator, Broad Institute).

611

612    **KASP genotyping.**  To validate sequence variants, genomic DNA was subjected to allele

613    specific PCR (KASP genotyping technology, LGC group), a fluorescent technique which

614    enables testing of SNPs and indels at specific loci.  Primers were designed with 50 bp flanks

615    around the site of interest for each variant allele and genomic DNA from original and

616    microevolved isolates was tested across 80 unique sites (20 for each strain background).  Allele

617    frequencies were calculated for each site and genotyping was assigned by cluster analysis.

618    Sites were selected so that they represent all mutation categories (GOH, LOH, SNPs, indels,

619    transitions, transversions) as well as different regions of the genome (coding, intergenic, MRS,

620    retrotransposons, adhesins and telomeric genes).  Out of the 80 SNPs tested, 63 (78.8%) were

621    successfully genotyped via KASP.  The 63 genotyped sites were then compared with mutations

622    called from genome sequencing data with an 87.3% success rate.

623

624  **Ploidy and copy number variation.**  To examine ploidy variation across the genome, the

625  Illumina read alignment depth was calculated for 100 bp windows across the genome, using

626  BEDTools 2.18 (99), SAMtools 1.3 (97) and the GATK 3.7 Depth of Coverage module.  The

627  read depth was calculated as the number of bases aligned per window divided by the length of

628  the window and normalized to the average depth for each strain and to the GC content, as this

629  can influence both the sequencing chemistry and the alignment quality (100).  The read depth

630  was also normalized per the effective window length by removing any ambiguous sites in the

631  respective window.  The normalized alignment depth for each 100 bp window was then plotted

632  and large scale variations in ploidy (2 fold up or down coverage) were identified.  These include

633  whole chromosome and segmental aneuploidies larger than 0.1 Mbp.  Smaller regions showing

634  read depth variation were designated as copy number variants (CNVs) and their numbers

635  plotted based on the nature of the variation (2 fold up or down coverage).

636

637  **LOH analysis.**  To identify LOH events, all variants were classified based on how each

638  mutation alters heterozygosity at the respective site: losses of heterozygosity (LOH), gains of

639  heterozygosity (GOH), or mutations that do not alter heterozygosity (het neutral).  LOH tracts

640  were defined using each heterozygous site identified to have undergone LOH and the size of

641  the tracts was determined by visual inspection in IGV.  To calculate $L_{min}$ for an LOH event, tracts

642  began at the first converted LOH SNP identified and ended at the final converted LOH SNP with

643  no interruption by a heterozygous position.  For LOH events encompassing a single LOH SNP

644  (i.e., not flanked by a consecutive LOH SNP) $L_{min}$ was 1.  To calculate $L_{max}$ for an LOH event,

645  tracts measure the distance between the nearest upstream, non-converted position and the

646  nearest first downstream, non-converted position of the respective LOH tract.  The average size

647  of LOH tracts ($L_{avg}$) was calculated by averaging the minimum ($L_{min}$) and maximum ($L_{max}$)

648  lengths for each observed event.  LOH tracts were then classified based on genomic size:

28

649    microLOH (mLOH, affecting single heterozygous positions and with an $L_{min}$ of 1 bp), short tract

650    LOH (SLOH, affecting two or more heterozygous positions and with an $L_{min} \leq 10$ kb), and long

651    tract LOH (LLOH, affecting hundreds of heterozygous positions and with an $L_{min} > 10$kb).  The

652    numbers of LOH were then assessed for each lineage, niche of evolution and chromosome.

653    The LOH distribution was examined across all isolates for each lineage and for each niche

654    using the $L_{min}$ genomic size (Figure 4C,D) of individual LOH events or the number of

655    heterozygous positions that were impacted by each LOH (Supplementary Figure 4B).

656

657    **Mutations in different genomic regions or regions adjacent to LOH events.**  We identified

658    mutations in specific regions using the genomic coordinates of these regions - strain specific

659    HET (heterozygous) and HOM (homozygous) regions, HET/HOM junctions, repeat regions

660    (MRS, LTR and genes associated with repeats), centromeres, chromosomal ends, *TLO* genes,

661    and *ALS* and GPI-linked genes.  Genomic coordinates for these chromosomal features were

662    obtained from the Candida Genome Database (http://www.candidagenome.org) and genomic

663    coordinates are provided in Supplementary Tables 6 and 7.  HET and HOM regions were

664    previously defined for the four starting strains (24) and are included in Supplementary Table 8.

665

666    **Statistical analyses.**  Statistical analyses were performed using two-tailed Student's t-tests and

667    by calculating probability values of binomial model distributions using Microsoft Excel 2016

668    (Microsoft) and Prism 6 (GraphPad).  Significance was assigned for *P* values < 0.05, and

669    asterisks denote *P* values that satisfy this condition.

670

671    **Data access.**  The sequence data from this study have been submitted to the NCBI SRA under

672    BioProject ID PRJNA345600 (http://www.ncbi.nlm.nih.gov/bioproject).

673

674    **Ethics Statement.** This study was carried out in strict accordance with the recommendations in

675    the Guide for the Care and Use of Laboratory Animals as defined by the National Institutes of

676    Health (PHS Assurance #A3284-01).  Animal protocols were reviewed and approved by the

677    Institutional Animal Care and Use Committee (IACUC) of Brown University.  All animals were

678    housed in a centralized and AAALAC-accredited research animal facility that is fully staffed with

679    trained husbandry, technical and veterinary personnel.

680

681

**References**

1.      Byrnes EJ, 3rd, Li W, Lewit Y, Ma H, Voelz K, Ren P, et al. Emergence and pathogenicity of highly virulent *Cryptococcus gattii* genotypes in the northwest United States. PLoS pathogens. 2010;6(4):e1000850.

2.      Gire SK, Goba A, Andersen KG, Sealfon RS, Park DJ, Kanneh L, et al. Genomic surveillance elucidates Ebola virus origin and transmission during the 2014 outbreak. Science. 2014;345(6202):1369-72.

3.      Jerison ER, Desai MM. Genomic investigations of evolutionary dynamics and epistasis in microbial evolution experiments. Curr Opin Genet Dev. 2015;35:33-9.

4.      Yang S, Wang L, Huang J, Zhang X, Yuan Y, Chen JQ, et al. Parent-progeny sequencing indicates higher mutation rates in heterozygotes. Nature. 2015;523(7561):463-7.

5.      McDonald MJ, Rice DP, Desai MM. Sex speeds adaptation by altering the dynamics of molecular evolution. Nature. 2016;531(7593):233-6.

6.      Sellis D, Kvitek DJ, Dunn B, Sherlock G, Petrov DA. Heterozygote advantage is a common outcome of adaptation in *Saccharomyces cerevisiae*. Genetics. 2016;203(3):1401-13.

7.      Keith N, Tucker AE, Jackson CE, Sung W, Lucas Lledo JI, Schrider DR, et al. High mutational rates of large-scale duplication and deletion in *Daphnia pulex*. Genome Res. 2016;26(1):60-9.

8.      Dutta A, Lin G, Pankajam AV, Chakraborty P, Bhat N, Steinmetz LM, et al. Genome dynamics of hybrid *Saccharomyces cerevisiae* during vegetative and meiotic divisions. G3 (Bethesda). 2017;7(11):3669-79.

9.      Farrer RA, Voelz K, Henk DA, Johnston SA, Fisher MC, May RC, et al. Microevolutionary traits and comparative population genomics of the emerging pathogenic fungus *Cryptococcus gattii*. Philos Trans R Soc Lond B Biol Sci. 2016;371(1709).

706  10.      Goddard MR, Godfray HC, Burt A. Sex increases the efficacy of natural selection in

707  experimental yeast populations. Nature. 2005;434(7033):636-40.

708  11.      Heitman J. Sexual reproduction and the evolution of microbial pathogens. Current

709  biology : CB. 2006;16(17):R711-25.

710  12.      Ryland GL, Doyle MA, Goode D, Boyle SE, Choong DY, Rowley SM, et al. Loss of

711  heterozygosity: what is it good for? BMC Med Genomics. 2015;8:45.

712  13.      Campbell CD, Chong JX, Malig M, Ko A, Dumont BL, Han L, et al. Estimating the human

713  mutation rate using autozygosity in a founder population. Nat Genet. 2012;44(11):1277-81.

714  14.      Brown GD, Denning DW, Gow NA, Levitz SM, Netea MG, White TC. Hidden killers:

715  human fungal infections. Sci Transl Med. 2012;4(165):165rv13.

716  15.      Kullberg BJ, Arendrup MC. Invasive Candidiasis. N Engl J Med. 2015;373(15):1445-56.

717  16.      Calderone RA, Clancy CJ. *Candida* and Candidiasis. Washington DC: ASM Press; 2011.

718  17.      Pfaller MA, Diekema DJ. Epidemiology of invasive candidiasis: a persistent public health

719  problem. Clinical microbiology reviews. 2007;20(1):133-63.

720  18.      Jones T, Federspiel NA, Chibana H, Dungan J, Kalman S, Magee BB, et al. The diploid

721  genome sequence of *Candida albicans*. Proceedings of the National Academy of Sciences of

722  the United States of America. 2004;101(19):7329-34.

723  19.      Braun BR, van Het Hoog M, d'Enfert C, Martchenko M, Dungan J, Kuo A, et al. A

724  human-curated annotation of the *Candida albicans* genome. PLoS Genet. 2005;1(1):36-57.

725  20.      Muzzey D, Schwartz K, Weissman JS, Sherlock G. Assembly of a phased diploid

726  *Candida albicans* genome facilitates allele-specific measurements and provides a simple model

727  for repeat and indel structure. Genome Biol. 2013;14(9):R97.

728  21.      Diogo D, Bouchier C, d'Enfert C, Bougnoux ME. Loss of heterozygosity in commensal

729  isolates of the asexual diploid yeast *Candida albicans*. Fungal Genet Biol. 2009;46(2):159-68.

730   22.     Forche A, Abbey D, Pisithkul T, Weinzierl MA, Ringstrom T, Bruck D, et al. Stress alters

731   rates and types of loss of heterozygosity in *Candida albicans*. mBio. 2011;

732   2(4):10.1128/mBio.00129-11. Print 2011.

733   23.     Rosenberg SM. Stress-induced loss of heterozygosity in *Candida*: a possible missing

734   link in the ability to evolve. MBio. 2011;2(5).

735   24.     Hirakawa MP, Martinez DA, Sakthikumar S, Anderson MZ, Berlin A, Gujja S, et al.

736   Genetic and phenotypic intra-species variation in *Candida albicans*. Genome Res.

737   2015;25(3):413-25.

738   25.     Forche A. Large-scale chromosomal changes and associated fitness consequences in

739   pathogenic fungi. Curr Fungal Infect Rep. 2014;8(2):163-70.

740   26.     Morrow CA, Fraser JA. Ploidy variation as an adaptive mechanism in human pathogenic

741   fungi. Seminars in cell & developmental biology. 2013;24(4):339-46.

742   27.     Selmecki A, Forche A, Berman J. Aneuploidy and isochromosome formation in drug-

743   resistant *Candida albicans*. Science (New York, NY). 2006;313(5785):367-70.

744   28.     Yang F, Yan TH, Rustchenko E, Gao PH, Wang Y, Yan L, et al. High-frequency genetic

745   contents variations in clinical *Candida albicans* isolates. Biol Pharm Bull. 2011;34(5):624-31.

746   29.     Rustchenko E. Chromosome instability in *Candida albicans*. FEMS Yeast Res.

747   2007;7(1):2-11.

748   30.     Bennett RJ, Forche A, Berman J. Rapid mechanisms for generating genome diversity:

749   whole ploidy shifts, aneuploidy, and loss of heterozygosity. Cold Spring Harb Perspect Med.

750   2014;4(10).

751   31.     Forche A, Magee PT, Selmecki A, Berman J, May G. Evolution in *Candida albicans*

752   populations during a single passage through a mouse host. Genetics. 2009;182(3):799-811.

753   32.     Ford CB, Funt JM, Abbey D, Issi L, Guiducci C, Martinez DA, et al. The evolution of drug

754   resistance in clinical isolates of *Candida albicans*. Elife. 2015;4:e00662.

755   33.    Chen C, Pande K, French SD, Tuch BB, Noble SM. An iron homeostasis regulatory

756   circuit with reciprocal roles in *Candida albicans* commensalism and pathogenesis. Cell host &

757   microbe. 2011;10(2):118-35.

758   34.    Yamaguchi N, Sonoyama K, Kikuchi H, Nagura T, Aritsuka T, Kawabata J. Gastric

759   colonization of *Candida albicans* differs in mice fed commercial and purified diets. J Nutr.

760   2005;135(1):109-15.

761   35.    MacCallum DM, Odds FC. Temporal events in the intravenous challenge model for

762   experimental *Candida albicans* infections in female mice. Mycoses. 2005;48(3):151-61.

763   36.    Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: an

764   integrated tool for comprehensive microbial variant detection and genome assembly

765   improvement. PLoS One. 2014;9(11):e112963.

766   37.    Selmecki A, Gerami-Nejad M, Paulson C, Forche A, Berman J. An isochromosome

767   confers drug resistance *in vivo* by amplification of two genes, *ERG11* and *TAC1*. Molecular

768   microbiology. 2008;68(3):624-41.

769   38.    Arbour M, Epp E, Hogues H, Sellam A, Lacroix C, Rauceo J, et al. Widespread

770   occurrence of chromosomal aneuploidy following the routine production of *Candida albicans*

771   mutants. FEMS yeast research. 2009;9(7):1070-7.

772   39.    Thorvaldsdottir H, Robinson JT, Mesirov JP. Integrative Genomics Viewer (IGV): high-

773   performance genomics data visualization and exploration. Brief Bioinform. 2013;14(2):178-92.

774   40.    Zhu YO, Siegal ML, Hall DW, Petrov DA. Precise estimates of mutation rate and

775   spectrum in yeast. Proc Natl Acad Sci U S A. 2014;111(22):E2310-8.

776   41.    Lynch M, Sung W, Morris K, Coffey N, Landry CR, Dopman EB, et al. A genome-wide

777   view of the spectrum of spontaneous mutations in yeast. Proc Natl Acad Sci U S A.

778   2008;105(27):9272-7.

779   42.     Nishant KT, Wei W, Mancera E, Argueso JL, Schlattl A, Delhomme N, et al. The baker's

780   yeast diploid genome is remarkably stable in vegetative growth and meiosis. PLoS Genet.

781   2010;6(9):e1001109.

782   43.     Keightley PD, Trivedi U, Thomson M, Oliver F, Kumar S, Blaxter ML. Analysis of the

783   genome sequences of three *Drosophila melanogaster* spontaneous mutation accumulation

784   lines. Genome Res. 2009;19(7):1195-201.

785   44.     Graur D. Single-base mutation. Encyclopedia of Life Sciences. 2008.

786   45.     Farlow A, Long H, Arnoux S, Sung W, Doak TG, Nordborg M, et al. The spontaneous

787   mutation rate in the fission yeast *Schizosaccharomyces pombe*. Genetics. 2015;201(2):737-44.

788   46.     Prieto D, Pla J. Distinct stages during colonization of the mouse gastrointestinal tract by

789   *Candida albicans*. Front Microbiol. 2015;6:792.

790   47.     Forche A, Schonian G, Graser Y, Vilgalys R, Mitchell TG. Genetic structure of typical

791   and atypical populations of *Candida albicans* from Africa. Fungal Genet Biol. 1999;28(2):107-25.

792   48.     Graser Y, Volovsek M, Arrington J, Schonian G, Presber W, Mitchell TG, et al. Molecular

793   markers reveal that population structure of the human pathogen *Candida albicans* exhibits both

794   clonality and recombination. Proceedings of the National Academy of Sciences of the United

795   States of America. 1996;93(22):12473-7.

796   49.     Balloux F, Lehmann L, de Meeus T. The population genetics of clonal and partially

797   clonal diploids. Genetics. 2003;164(4):1635-44.

798   50.     Arbeithuber B, Betancourt AJ, Ebner T, Tiemann-Boege I. Crossovers are associated

799   with mutation and biased gene conversion at recombination hotspots. Proc Natl Acad Sci U S A.

800   2015;112(7):2109-14.

801   51.     Malkova A, Haber JE. Mutations arising during repair of chromosome breaks. Annu Rev

802   Genet. 2012;46:455-73.

803  52.  Amos W. Heterozygosity and mutation rate: evidence for an interaction and its

804  implications: the potential for meiotic gene conversions to influence both mutation rate and

805  distribution. Bioessays. 2010;32(1):82-90.

806  53.  Hsueh YP, Idnurm A, Heitman J. Recombination hotspots flank the *Cryptococcus*

807  mating-type locus: implications for the evolution of a fungal sex chromosome. PLoS genetics.

808  2006;2(11):e184.

809  54.  Bensasson D. Evidence for a high mutation rate at rapidly evolving yeast centromeres.

810  BMC Evol Biol. 2011;11:211.

811  55.  Yang Y, Sterling J, Storici F, Resnick MA, Gordenin DA. Hypermutability of damaged

812  single-strand DNA formed at double-strand breaks and uncapped telomeres in yeast

813  *Saccharomyces cerevisiae*. PLoS Genet. 2008;4(11):e1000264.

814  56.  Anderson MZ, Wigen LJ, Burrack LS, Berman J. Real-time evolution of a subtelomeric

815  gene family in *Candida albicans*. Genetics. 2015;200(3):907-19.

816  57.  McEachern MJ, Hicks JB. Unusually large telomeric repeats in the yeast *Candida*

817  *albicans*. Mol Cell Biol. 1993;13(1):551-60.

818  58.  Sadhu C, McEachern MJ, Rustchenko-Bulgac EP, Schmid J, Soll DR, Hicks JB.

819  Telomeric and dispersed repeat sequences in *Candida* yeasts and their use in strain

820  identification. J Bacteriol. 1991;173(2):842-50.

821  59.  Sanyal K, Baum M, Carbon J. Centromeric DNA sequences in the pathogenic yeast

822  *Candida albicans* are all different and unique. Proc Natl Acad Sci U S A. 2004;101(31):11374-9.

823  60.  Mishra PK, Baum M, Carbon J. Centromere size and position in *Candida albicans* are

824  evolutionarily conserved independent of DNA sequence heterogeneity. Mol Genet Genomics.

825  2007;278(4):455-65.

826  61.  Dujon B, Sherman D, Fischer G, Durrens P, Casaregola S, Lafontaine I, et al. Genome

827  evolution in yeasts. Nature. 2004;430(6995):35-44.

36

828    62.    Chibana H, Magee PT. The enigma of the major repeat sequence of *Candida albicans*.

829    Future Microbiol. 2009;4(2):171-9.

830    63.    Chibana H, Beckerman JL, Magee PT. Fine-resolution physical mapping of genomic

831    diversity in *Candida albicans*. Genome Res. 2000;10(12):1865-77.

832    64.    Freire-Beneitez V, Price RJ, Tarrant D, Berman J, Buscaino A. *Candida albicans*

833    repetitive elements display epigenetic diversity and plasticity. Sci Rep. 2016;6:22989.

834    65.    Lephart PR, Chibana H, Magee PT. Effect of the major repeat sequence on

835    chromosome loss in *Candida albicans*. Eukaryot Cell. 2005;4(4):733-41.

836    66.    Lephart PR, Magee PT. Effect of the major repeat sequence on mitotic recombination in

837    *Candida albicans*. Genetics. 2006;174(4):1737-44.

838    67.    Hoyer LL, Green CB, Oh SH, Zhao X. Discovering the secrets of the *Candida albicans*

839    agglutinin-like sequence (*ALS*) gene family--a sticky pursuit. Medical mycology : official

840    publication of the International Society for Human and Animal Mycology. 2008;46(1):1-15.

841    68.    Oh SH, Cheng G, Nuessen JA, Jajko R, Yeater KM, Zhao X, et al. Functional specificity

842    of *Candida albicans* Als3p proteins and clade specificity of *ALS3* alleles discriminated by the

843    number of copies of the tandem repeat sequence in the central domain. Microbiology.

844    2005;151(Pt 3):673-81.

845    69.    Zhao X, Oh SH, Jajko R, Diekema DJ, Pfaller MA, Pujol C, et al. Analysis of *ALS5* and

846    *ALS6* allelic variability in a geographically diverse collection of *Candida albicans* isolates.

847    Fungal Genet Biol. 2007;44(12):1298-309.

848    70.    Zhao X, Pujol C, Soll DR, Hoyer LL. Allelic variation in the contiguous loci encoding

849    *Candida albicans ALS5*, *ALS1* and *ALS9*. Microbiology. 2003;149(Pt 10):2947-60.

850    71.    Dunkel N, Blass J, Rogers PD, Morschhauser J. Mutations in the multi-drug resistance

851    regulator *MRR1*, followed by loss of heterozygosity, are the main cause of *MDR1*

852    overexpression in fluconazole-resistant *Candida albicans* strains. Mol Microbiol.

853    2008;69(4):827-40.

854    72.    Coste A, Turner V, Ischer F, Morschhauser J, Forche A, Selmecki A, et al. A mutation in

855    Tac1p, a transcription factor regulating *CDR1* and *CDR2*, is coupled with loss of heterozygosity

856    at chromosome 5 to mediate antifungal resistance in *Candida albicans*. Genetics.

857    2006;172(4):2139-56.

858    73.    Feri A, Loll-Krippleber R, Commere PH, Maufrais C, Sertour N, Schwartz K, et al.

859    Analysis of repair mechanisms following an induced double-strand break uncovers recessive

860    deleterious alleles in the *Candida albicans* Diploid Genome. MBio. 2016;7(5).

861    74.    Deem A, Keszthelyi A, Blackgrove T, Vayl A, Coffey B, Mathur R, et al. Break-induced

862    replication is highly inaccurate. PLoS Biol. 2011;9(2):e1000594.

863    75.    Hollister JD, Ross-Ibarra J, Gaut BS. Indel-associated mutation rate varies with mating

864    system in flowering plants. Mol Biol Evol. 2010;27(2):409-16.

865    76.    Tian D, Wang Q, Zhang P, Araki H, Yang S, Kreitman M, et al. Single-nucleotide

866    mutation rate increases close to insertions/deletions in eukaryotes. Nature.

867    2008;455(7209):105-8.

868    77.    Mudrak SV, Welz-Voegele C, Jinks-Robertson S. The polymerase eta translesion

869    synthesis DNA polymerase acts independently of the mismatch repair system to limit

870    mutagenesis caused by 7,8-dihydro-8-oxoguanine in yeast. Mol Cell Biol. 2009;29(19):5316-26.

871    78.    Livneh Z. DNA damage control by novel DNA polymerases: translesion replication and

872    mutagenesis. J Biol Chem. 2001;276(28):25639-42.

873    79.    Rodgers K, McVey M. Error-prone repair of DNA double-strand breaks. J Cell Physiol.

874    2016;231(1):15-24.

875    80.    MacCallum DM, Castillo L, Nather K, Munro CA, Brown AJ, Gow NA, et al. Property

876    differences among the four major *Candida albicans* strain clades. Eukaryotic cell.

877    2009;8(3):373-87.

878    81.    Shor E, Fox CA, Broach JR. The yeast environmental stress response regulates

879    mutagenesis induced by proteotoxic stress. PLoS Genet. 2013;9(8):e1003680.

880    82.    Galhardo RS, Hastings PJ, Rosenberg SM. Mutation as a stress response and the

881    regulation of evolvability. Crit Rev Biochem Mol Biol. 2007;42(5):399-435.

882    83.    Jiang C, Mithani A, Belfield EJ, Mott R, Hurst LD, Harberd NP. Environmentally

883    responsive genome-wide accumulation of de novo *Arabidopsis thaliana* mutations and

884    epimutations. Genome Res. 2014;24(11):1821-9.

885    84.    Clutterbuck AJ. Genomic evidence of repeat-induced point mutation (RIP) in filamentous

886    ascomycetes. Fungal Genet Biol. 2011;48(3):306-26.

887    85.    Wohlbach DJ, Rovinskiy N, Lewis JA, Sardi M, Schackwitz WS, Martin JA, et al.

888    Comparative genomics of *Saccharomyces cerevisiae* natural isolates for bioenergy production.

889    Genome Biol Evol. 2014;6(9):2557-66.

890    86.    Doniger SW, Kim HS, Swain D, Corcuera D, Williams M, Yang SP, et al. A catalog of

891    neutral and deleterious polymorphism in yeast. PLoS Genet. 2008;4(8):e1000183.

892    87.    Weil T, Santamaria R, Lee W, Rung J, Tocci N, Abbey D, et al. Adaptive mistranslation

893    accelerates the evolution of fluconazole resistance and induces major genomic and gene

894    expression alterations in *Candida albicans*. mSphere. 2017;2(4).

895    88.    Judd SR, Petes TD. Physical lengths of meiotic and mitotic gene conversion tracts in

896    *Saccharomyces cerevisiae*. Genetics. 1988;118(3):401-10.

897    89.    Mancera E, Bourgon R, Brozzi A, Huber W, Steinmetz LM. High-resolution mapping of

898    meiotic crossovers and non-crossovers in yeast. Nature. 2008;454(7203):479-85.

899    90.    Enloe B, Diamond A, Mitchell AP. A single-transformation gene function test in diploid

900    *Candida albicans*. J Bacteriol. 2000;182(20):5730-6.

901    91.    Moorhouse AJ, Rennison C, Raza M, Lilic D, Gow NA. Clonal strain persistence of

902    *Candida albicans* isolates from chronic mucocutaneous candidiasis patients. PLoS One.

903    2016;11(2):e0145888.

904    92.    Sakofsky CJ, Ayyar S, Malkova A. Break-induced replication and genome stability.

905    Biomolecules. 2012;2(4):483-504.

906    93.    Sherman F. Getting started with yeast. Methods Enzymol. 1991;194:3-21.

907    94.    MacCallum DM, Coste A, Ischer F, Jacobsen MD, Odds FC, Sanglard D. Genetic

908    dissection of azole resistance mechanisms in *Candida albicans* and their validation in a mouse

909    model of disseminated infection. Antimicrobial Agents and Chemotherapy. 2010;54(4):1476-83.

910    95.    Pande K, Chen C, Noble SM. Passage through the mammalian gut triggers a phenotypic

911    switch that promotes *Candida albicans* commensalism. Nature genetics. 2013.

912    96.    Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.

913    arXiv. 2013;arXiv:1303.3997.

914    97.    Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The sequence

915    alignment/map format and SAMtools. Bioinformatics. 2009;25(16):2078-9.

916    98.    McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al. The

917    Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA

918    sequencing data. Genome Res. 2010;20(9):1297-303.

919    99.    Quinlan AR. BEDTools: The swiss-army tool for genome feature analysis. Curr Protoc

920    Bioinformatics. 2014;47:11 2 1-34.

921    100.    Risso D, Schwartz K, Sherlock G, Dudoit S. GC-content normalization for RNA-Seq

922    data. BMC bioinformatics. 2011;12:480.

923    101.    Haran J, Boyle H, Hokamp K, Yeomans T, Liu Z, Church M, et al. Telomeric ORFs

924    (*TLOs*) in *Candida* spp. Encode mediator subunits that regulate distinct virulence traits. PLoS

925    Genet. 2014;10(10):e1004658.

926    102.    Ramsdale M, Selway L, Stead D, Walker J, Yin Z, Nicholls SM, et al. *MNL1* regulates

927    weak acid-induced stress responses of the fungal pathogen *Candida albicans*. Molecular

928    biology of the cell. 2008;19(10):4393-403.

929    103.    Yu EY, Steinberg-Neifach O, Dandjinou AT, Kang F, Morrison AJ, Shen X, et al.

930    Regulation of telomere structure and functions by subunits of the *INO80* chromatin remodeling

931    complex. Mol Cell Biol. 2007;27(16):5639-49.

932    104.    Kathe SD, Barrantes-Reynolds R, Jaruga P, Newton MR, Burrows CJ, Bandaru V, et al.

933    Plant and fungal Fpg homologs are formamidopyrimidine DNA glycosylases but not 8-

934    oxoguanine DNA glycosylases. DNA Repair (Amst). 2009;8(5):643-53.

935    105.    Forche A, Alby K, Schaefer D, Johnson AD, Berman J, Bennett RJ. The parasexual

936    cycle in *Candida albicans* provides an alternative pathway to meiosis for the formation of

937    recombinant strains. PLoS biology. 2008;6(5):e110.

938    106.    Legrand M, Chan CL, Jauert PA, Kirkpatrick DT. Role of DNA mismatch repair and

939    double-strand break repair in genome stability and antifungal drug resistance in *Candida*

940    *albicans*. Eukaryot Cell. 2007;6(12):2194-205.

941

942

943

944 **Figure legends**

945 **Figure 1.** Microevolution of *C. albicans* genomes. (A) Schematic of *in vitro* and *in vivo*

946 microevolution experiments. GI SD, gastointestinal standard diet; GI PD, gastointestinal purified

947 diet. (B) Ploidy variation based on read depth for each evolved isolate and each chromosome,

948 normalized by the average genomic read depth. Isolates with significant ploidy changes

949 (including full and segmental aneuploidies) as well as those with large chromosomal events are

950 marked in red. (C) Distribution of SNPs and indels, intergenic and coding mutations,

951 synonymous and nonsynonymous mutations across microevolved isolates averaged for each

952 lineage. Note that panels include both GOH and LOH events.

953

954 **Figure 2.** Selection shapes microevolution of *C. albicans* genomes. (A) Frequency of observed

955 and expected GOH mutations in intergenic regions (these regions represent 36.2% of the

956 *C. albicans* genome). GOH SNP mutations represent *de novo* base substitutions. (B) Number

957 of GOH indels identified in coding and intergenic regions. Indels are classified based on

958 whether they are multiple of 3 bp in length. (C) Observed and expected fractions of

959 synonymous GOH SNPs in coding regions during microevolution experiments. ~75% of all base

960 substitutions are expected to be nonsynonymous if they occurred randomly and were not

961 subject to selection. (D) Selection coefficients for nonsynonymous GOH SNPs calculated

962 based on the number of observed vs. expected nonsynonymous substitutions. Expected

963 nonsynonymous substitutions were estimated for each microevolution experiment based on

964 observed synonymous substitutions (squares) or observed intergenic substitutions (circles).

965 Only isolates for which both nonsynonymous and synonymous/intergenic substitutions were

966 observed were included in this analysis. For each panel, asterisks indicate significant

967 differences (t-test, $P < 0.05$).

968

969 **Figure 3.** Mutation rates in *C. albicans* are impacted by strain background, environmental

970 niche, and chromosomal features. (A,B) Effect of strain background (A) and evolution niche (B)

971 on GOH and LOH mutation rates. Rates include both SNP and indel mutations. (C)

972 Fluctuations in genome heterozygosity during microevolution relative to starting heterozygosity

973 levels (red line). Major decreases in heterozygosity are only observed for isolates that

974 underwent large-tract LOH events (LLOH; red symbols). (D) Mutation rates in specific regions

975 of the *C. albicans* genome. These include heterozygous (HET) regions, repeat regions (MRS

976 and LTR), Chr END regions (final 10 kb of each chromosome arm), centromeres, genes

977 encoding *ALS* and GPI-linked proteins, and *TLO* genes. Mutations include SNPs and indels

978 resulting from both GOH and LOH. Asterisks indicate significant differences (t-test, $P < 0.05$).

979

980 **Figure 4.** Microevolution is punctuated by frequent short-tract LOH events and small indels.

981 (A) Schematic of different types of LOH events, including microLOH (mLOH, involve loss of

982 single heterozygous positions), short-tract LOH (SLOH, involve loss of two or more

983 heterozygous positions and are <10 kb), and long-tract LOH (LLOH, affect hundreds of

984 heterozygous positions and are >10 kb). (B) Distribution of LOH events showing the $L_{min}$, $L_{avg}$

985 and $L_{max}$ size for each LOH event. (C,D) $L_{min}$ size distribution of LOH events, including

986 microLOH, SLOH and LLOH (shown in red), for each lineage (C) and niche (D). (E,F) Size

987 distribution of indels, including insertions and deletions for each lineage (E) and niche (F).

988

989 **Figure 5.** Relationship between LOH location and different genomic regions. (A) Chromosomal

990 location of all LOH events (using $L_{min}$) with triangles marking the start (red) and end (blue) of

991 each event. Location of centromeres (CEN) and MRS regions are shown. (B) Proximity of LOH

992 events to the closest genomic feature, including MRS regions, telomeres (Chr ENDs or *TLO*

43

993    genes), and centromeres.  Each LOH event is uniquely mapped to the closest of these features

994    on the same chromosome arm.  Distances equal to 0 indicate an LOH start site inside the

995    respective genomic region.  (C) GOH rates (including SNPs and indels) in the duplicated LLOH

996    region versus the rest of the genome in the P78048 GI SD B isolate.  Only GOH SNPs (base

997    substitutions) were observed in the duplicated LLOH region. (D) Number of GOH events (SNPs

998    and indels per 25 bp) observed within 500 bp of LOH tracts in microevolved isolates. (E) GOH

999    rates (including SNPs and indels) in regions adjacent to LOH tracts (within 500 bp) compared to

1000    rates in the rest of the genome.

1001

1002    **Figure 6.**  Schematic illustrating the pattern of mutational events across the *C. albicans*

1003    genome.  Figure highlights how certain chromosomal features are associated with elevated

1004    mutation rates.

1005

**Supplemental Material**

**Supplementary Text 1. Analysis of copy number variation during microevolution**

Read depth analysis revealed copy number variation (CNV) during microevolution. Several genomic regions (average of ~2% of all 100 bp windows) showed a two-fold increase or decrease in coverage relative to the parental strain. Both the number of CNV regions and the nature of the variation differed between strain backgrounds (Supplementary Figure 6A). For example, the SC5314 lineage showed the least CNV, with only ~0.2% of 100 bp windows displaying a two-fold decrease (2X down) relative to the starting isolate. In contrast, ~0.7% of the P57055 windows displayed a two-fold decrease in coverage and ~1.7% of P57055 windows displayed a two-fold increase (2X up) in coverage (Supplementary Figure 6A).

A significant proportion of CNV was associated with specific regions of the genome. Analysis focused on the major repeat sequences (9 MRS elements span 1.7% of the genome and have been linked to chromosome translocations (62-66)), as well as the terminal 5 kb regions of chromosomes and the *TLO* family of subtelomeric genes (56, 62, 64, 66, 101) (see Supplementary Tables 6 and 7 for genomic coordinates). In each of these regions, evolved isolates displayed a higher percentage of windows with variable coverage relative to the rest of the genome (Supplementary Figure 6B-D). Chromosomal ends and *TLO* genes more often displayed a higher number of windows with increased read coverage following passaging (Supplementary Figure 6B,C, and Supplementary Figure 7A,B), whereas MRS regions often displayed reduced coverage (Supplementary Figure 6D). This suggests that terminal regions and *TLO* genes often underwent expansion during microevolution experiments whereas MRS regions more often underwent contraction. A closer inspection revealed that decreased coverage (2X down) was equally observed across MRS subunits, whereas increased coverage (2X up) was most common in the RB2 subunit (Supplementary Figure 6E,F). This result was surprising as variation was expected to be greatest in the tandem array of highly repetitive RPS

45

1031   subunits (62).  While the RPS regions displayed the highest coverage across the MRS

1032   (suggestive of multiple repeats being present), they displayed little CNV during our

1033   microevolution experiments (Supplementary Figure 7C).

1034        We also note that CNV occurred between genes from different lineages, which could

1035   contribute to phenotypic differences between strains.  For example, clade I strains displayed

1036   higher relative coverage of orf19.5474, encoding a protein of unknown function that is induced

1037   by Mnl1 during acid stress (102) (Supplementary Figure 7A).  In contrast, this gene displayed

1038   lower coverage levels in the other 2 lineages.  Similarly, variation in the relative coverage levels

1039   of *TLO* genes was noted between lineages (Supplementary Figure 7B).  These likely represent

1040   variations in gene copy number but local effects of accessibility to DNA isolation and purification

1041   could also contribute to these differences.  Overall, these analyses indicate that extensive copy

1042   number differences exist between different isolates of *C. albicans*.  In addition, particular regions

1043   of the genome are more prone to CNV, including both the terminal chromosome and MRS

1044   regions which fluctuate between 0 and 40 copies, and the *TLO* genes which fluctuate between 0

1045   and 8 copies (Supplementary Figure 7).

1046

1047   **Supplementary Text 2. Impact of mutations in DNA repair and maintenance genes on**

1048   **microevolution.**

1049        To examine whether preexisting mutations could affect differences in mutation rates

1050   between lineages, we screened for mutations in ~60 genes involved in DNA maintenance and

1051   repair (Supplementary Table 9).  This screen identified nonsense and readthrough mutations in

1052   several of the starting strains (Supplementary Table 10).  For example, both clade I strains

1053   SC5314 and P78048 had 2 premature heterozygous stop codon mutations in *ARP8*, a predicted

1054   component of the chromatin-remodeling enzyme complex (103).  These strains also had a

1055   heterozygous readthrough mutation in *FPG1*, a DNA glycosylase involved in the repair of

1056    irradiated DNA (104), whereas P76055 had a homozygous readthrough mutation in *FPG1*

1057    (Supplementary Table 10).  The position of the *FPG1* mutations was identical between strains,

1058    indicating that they shared a common ancestor and that the readthrough mutation occurred

1059    before clade I and II diverged.  In addition, P78048 had a unique heterozygous nonsense

1060    mutation in *SPO11,* a DNA endonuclease implicated in genetic recombination during parasex

1061    (105).  As most of these mutations were heterozygous or present in genes that are not essential

1062    for DNA maintenance under our experimental conditions, it is likely that they did not play a

1063    major role in impacting mutation rates during microevolution.

1064        We also examined whether new mutations emerged which could alter the function of

1065    DNA repair genes during passaging.  Overall, the 28 isolates acquired a total of 128 mutations

1066    in genes with known or predicted roles in DNA maintenance or DNA repair, 120 of which were

1067    directly associated with large chromosomal events such as LLOH tracts (Supplementary Table

1068    11).  The 120 mutations associated with LLOH tracts involved loss of heterozygosity mutations

1069    indicating that they were a direct consequence of the LOH event (Supplementary Table 11).

1070    The remaining mutations included 2 nonsynonymous changes (LOH mutations) in the *RAD50*

1071    gene of the P76055 GI PD isolate, which encodes a double strand repair protein with roles in

1072    stress responses (106).  However, none of the identified mutations were associated with

1073    obviously increased mutation rates in the respective isolates (Figure 3A).  Mutations in *RAD57*

1074    (orf19.4275) in one isolate (P76055 GI SD C) included 4 nonsynonymous missense mutations

1075    (Supplementary Table 11).  These mutations were in the first third of the amino acid sequence

1076    and could disrupt function.  *RAD57* encodes a key protein involved in DNA recombination and

1077    these mutations possibly altered the ability of this isolate to successfully undergo BIR and return

1078    to normal diploid levels, as illustrated by copy number variation analyses (Supplementary Figure

1079    1B).  Therefore, large-scale chromosomal changes via LOH events may have long-term

1080    consequences for genome evolution by disrupting important DNA repair pathways.

1081

1082 **Supplementary Figure legends**

1083 **Supplementary Figure 1.** Large chromosomal events identified during microevolution. (A) Chr

1084 7 trisomies were present in isolates recovered from the GI (SD model) in lineages SC5314,

1085 P78048 and P76055. (B) Segmental aneuploidies for Chr 2 and 3 were present in isolate

1086 P76055 GI SD C (boxed). A third LLOH tract was identified on Chr R of P78048 GI SD B

1087 (boxed). Variations in ploidy were determined by calculating the normalized read depth per 1 kb

1088 window. (C) Mutations verified using KASP assays. Fractions show the number of assays that

1089 failed or verified events identified via Illumina sequencing for different types of mutations. (D)

1090 Transitions/Transversion (*Ts/Tv*) ratios calculated for mutations arising during microevolution

1091 and broken down for mutations identified in intergenic and coding regions. Average *Ts/Tv* ratios

1092 are included below each category. (E) Distribution of SNPs and indels, coding and intergenic

1093 mutations, and synonymous and nonsynonymous mutations across microevolved isolates

1094 averaged for each niche. Note that panels include both GOH and LOH events.

1095

1096 **Supplementary Figure 2.** (A) Mutation rates (per bp per generation) calculated for different

1097 types of mutations following *in vitro* passaging of isolates. Asterisks indicate significant

1098 differences (t-test, $P < 0.05$). (B) Comparison of base-substitution rates in *C. albicans* and

1099 model organisms. Base-substitution rates (GOH SNPs) for *C. albicans* are shown as an

1100 average from *in vitro* evolution experiments performed in four different experiments and three

1101 genetic backgrounds. Expected rates reflect estimates based on observed intergenic base-

1102 substitution mutations. (C,D) Effect of strain background (C) and evolution niche (D) on indel

1103 mutation rates. Indels shown are the result of both GOH and LOH events. (E) GOH and LOH

1104 mutation rates across microevolution experiments. No significant difference was found between

1105 the two groups (t-test, $P < 0.05$).

1106

1107    **Supplementary Figure 3.** Mutation rates in specific regions of the *C. albicans* genome.

1108    Panels show mutation (SNPs and indels due to both GOH and LOH) rates relative to whole

1109    genome rates for (A) Chr END regions (final 10 kb of each Chr arm), (B) *TLO* genes, (C)

1110    centromeres, (D) repeat regions (MRS and LTR), and (E) genes encoding *ALS* family proteins

1111    and GPI-linked proteins. Asterisks indicate significant differences relative to either the SC5314

1112    lineage or *in vitro* passaged isolates ($P < 0.05$). (F) Observed fraction of GOH mutations in

1113    heterozygous (HET) regions for each lineage. Expected values (black) represent the % HET

1114    regions in the parental isolates of each lineage based on the density of heterozygous sites per 5

1115    kb genomic windows.

1116

1117    **Supplementary Figure 4.** Size of LOH events arising during microevolution. (A) Average $L_{min}$,

1118    $L_{avg}$ and $L_{max}$ size of LOH events for each lineage and niche. (B) Number of heterozygous

1119    positions affected by LOH tracts, shown for each lineage and niche. The average number of

1120    heterozygous positions impacted by LOH are included for each lineage and niche. For panels A

1121    and B, the three large LOH events were excluded from the analyses. (C) Density of LOH events

1122    relative to genomic location (per 0.2 Mb windows). Centromeres and MRS regions are included

1123    for reference.

1124

1125    **Supplementary Figure 5.** LLOH events represent three contiguous large-tract LOH regions.

1126    (A) Heterozygosity plots indicating the number of heterozygous positions for each 10 kb window

1127    across the 8 *C. albicans* chromosomes. Large LOH tracts on Chr R, 2 and 3 are boxed in red

1128    and shown relative to corresponding parental isolates. (B) Genomic size of tracts and number

1129    of mutations resulting from LLOH events, including a breakdown for mutations in the coding

1130    region. (C) GOH mutation rates in the 500 bp regions flanking LOH events (upstream and

1131   downstream) relative to whole genome GOH rates.  Asterisks indicate significant differences (t-

1132   test, *P* < 0.05) relative to whole genome rates (red dotted line).  (D) Distribution of GOH events

1133   (LOH-adjacent or not) relative to their position in the genome.  Regions showing enrichment of

1134   one GOH category over another are marked with asterisks (*P* < 0.05, using a binomial

1135   distribution model).

1136

1137   **Supplementary Figure 6.**  Copy number variation (CNV) analysis.  (A) Summary of CNV

1138   windows across the genome showing a two-fold increase (2X up) or decrease (2X down) in

1139   coverage relative to the original strain and averaged for each of the four lineages.  (B-D)

1140   Summary of CNV windows showing a two-fold increase or decrease relative to the original

1141   strain in Chr END regions (terminal 5 kb, B), *TLO* genes (C) and MRS regions (D).  (E)

1142   Schematic representation of a typical MRS region, including RB2, RPS and HOK subregions.

1143   (F) Breakdown of CNV windows aligning to different MRS subregions or MRS-5 (only partial

1144   MRS sequences are present on Chr 5) (63).  For all analyses CNV was determined by

1145   calculating the normalized read depth per 100 bp window.

1146

1147   **Supplementary Figure 7.**  Copy number variation (CNV) patterns in specific genomic regions.

1148   Normalized read depths are shown for Chr END regions (terminal 5 kb at the ends of

1149   chromosomes, A), *TLO* genes (B) and MRS regions (C).  CNV was determined by calculating

1150   the normalized read depth per 100 bp window.  Copy numbers for original (starting) strains are

1151   shown in black lines.

1152

1153   **Supplementary Table legends**

1154   **Supplementary Table 1.**  Strains used in microevolution experiments.

1155     **Supplementary Table 2.** Sequencing and coverage information, including the frequency of

1156     heterozygous sites in each isolate. Sequencing variants were identified relative to the SC5314

1157     genome reference strain.

1158     **Supplementary Table 3.** Doubling times of the four lineages and *in vitro* calculation of

1159     generation times. Calculation of estimated *in vivo* generation times were based on (31) and

1160     (46).

1161     **Supplementary Table 4.** Primers and results of the KASP genotyping assays. Probability

1162     values resulting from testing a binomial distribution model on the different types of mutations is

1163     also included.

1164     **Supplementary Table 5.** List of nonsense and readthrough mutations identified during

1165     microevolution experiments.

1166     **Supplementary Table 6.** Genomic coordinates for centromeres, chromosomal end regions,

1167     *TLO* genes, and *ALS* and GPI-liked genes, according to the *Candida* Genome Database (CGD).

1168     **Supplementary Table 7.** Genomic coordinates for major repeat sequences (MRS), long

1169     terminal repeats (LTR) and genes within repeat regions according to the CGD.

1170     **Supplementary Table 8.** Heterozygous (HET) and homozygous (HOM) regions of the four

1171     starting strains, as defined in (24). Included are also the overall genome heterozygosity levels

1172     based on these maps for the four strains.

1173     **Supplementary Table 9.** Genes with known or predicted roles in DNA maintenance and DNA

1174     repair, including their genomic coordinates and CGD annotation.

1175     **Supplementary Table 10**. Nonsense and stop codon mutations in DNA repair genes identified

1176     in the four starting isolates.

1177     **Supplementary Table 11**. Mutations in DNA repair genes identified across microevolved

1178     isolates. Mutations involving LOH SNPs associated with LLOH events and resulting in

1179     nonsynonymous changes to the protein sequence are highlighted in red.

51

**Figure 1**

# Figure 2

# Figure 3

**A** GOH and LOH mutation rates per lineage

**B** GOH and LOH mutation rates per niche

**C** Genome heterozygosity during microevolution

**D** Mutation rich regions

# Figure 4

# Figure 5

**Figure 6**

# Supplementary Figure 1

# Supplementary Figure 2

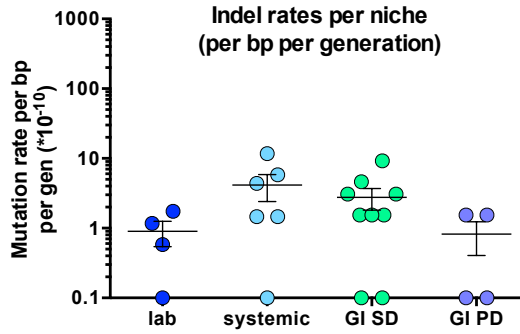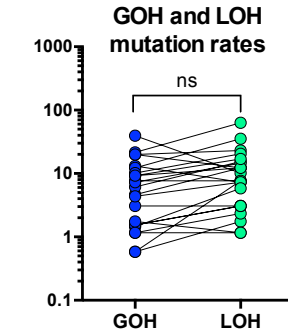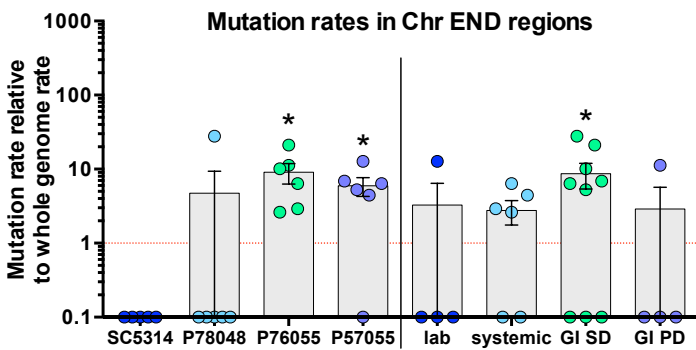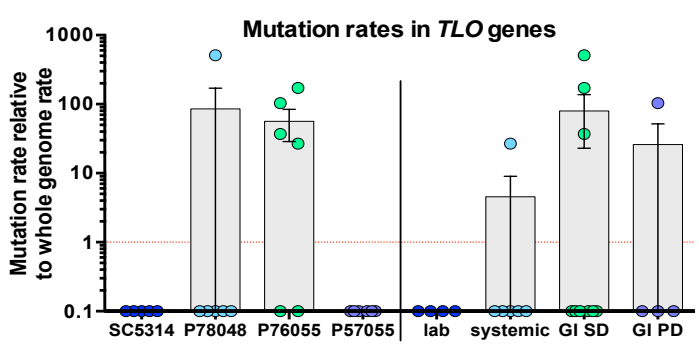# Supplementary Figure 3

# Supplementary Figure 4

**A**

| Average size (bp) | SC5314 | P78048 | P76055 | P57055 | lab | systemic | GI SD | GI PD |
|---|---|---|---|---|---|---|---|---|
| $L_{min}$ | 236.3 | 141.9 | 66.6 | 74.7 | 283.0 | 158.9 | 75.8 | 33.3 |
| $L_{avg}$ | 889.1 | 474.7 | 299.6 | 222.2 | 778.5 | 318.9 | 379.8 | 304.6 |
| $L_{max}$ | 1541.8 | 807.5 | 532.6 | 369.7 | 1273.9 | 478.9 | 683.9 | 575.8 |

**B**



| Number of het positions | SC5314 | P78048 | P76055 | P57055 | lab | systemic | GI SD | GI PD |
|---|---|---|---|---|---|---|---|---|
| average | 3.16 | 2.07 | 2.29 | 1.61 | 2.09 | 3.62 | 1.90 | 1.55 |

**C**

## Supplementary Figure 5

**A**



**B**

| LLOH | P76055 GI SD C | | P78048 GI SD B Chr R |
|---|---|---|---|
| | **Chr 2** | **Chr 3** | |
| size (Mb) | 1.146 | 0.273 | 1.23 |
| BIR | no | no | yes |

**Polymorphisms affected by LLOH**

| | | | |
|---|---|---|---|
| SNPs | 3176 | 1307 | 5812 |
| indels | 50 | 33 | 86 |
| intergenic | 1584 | 649 | 2812 |
| coding | 1642 | 691 | 3086 |

**LLOH resulting mutations in coding regions**

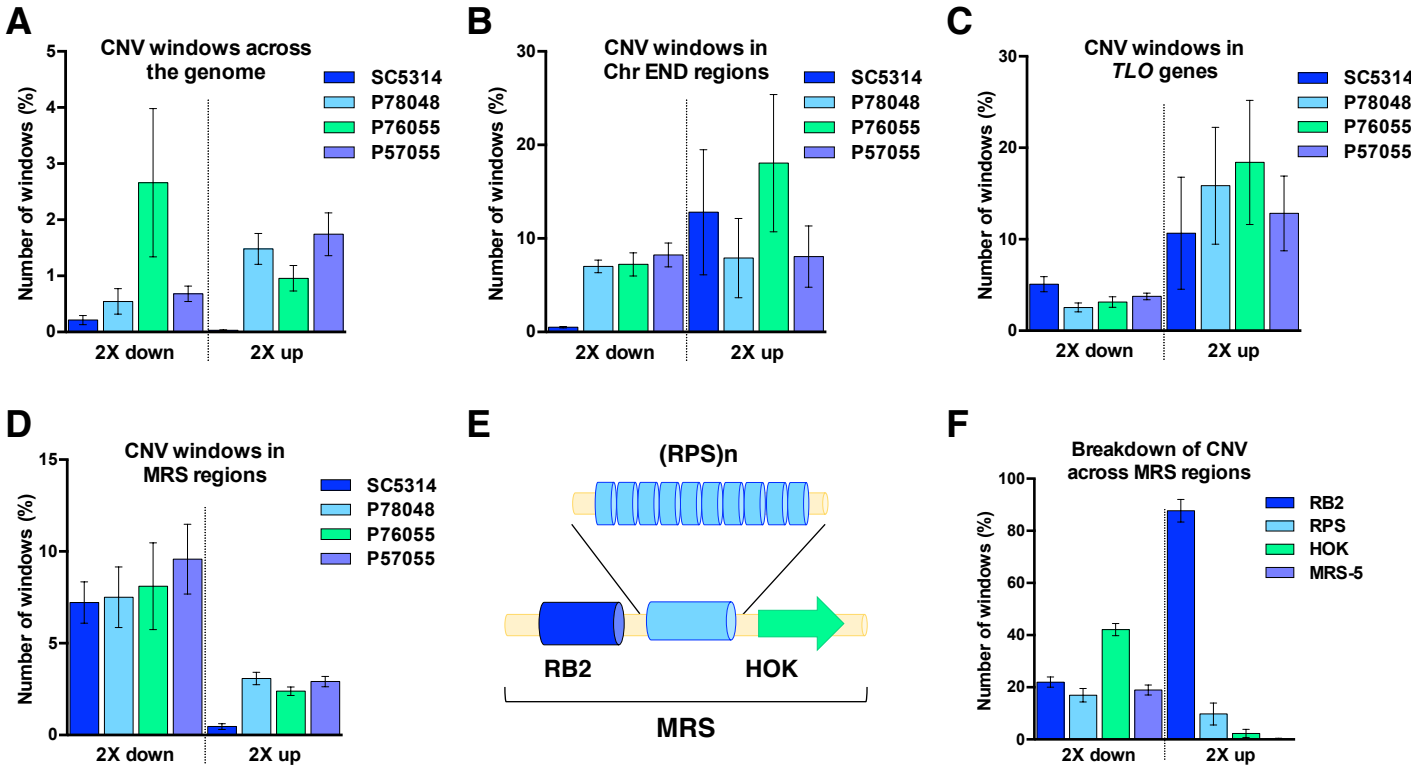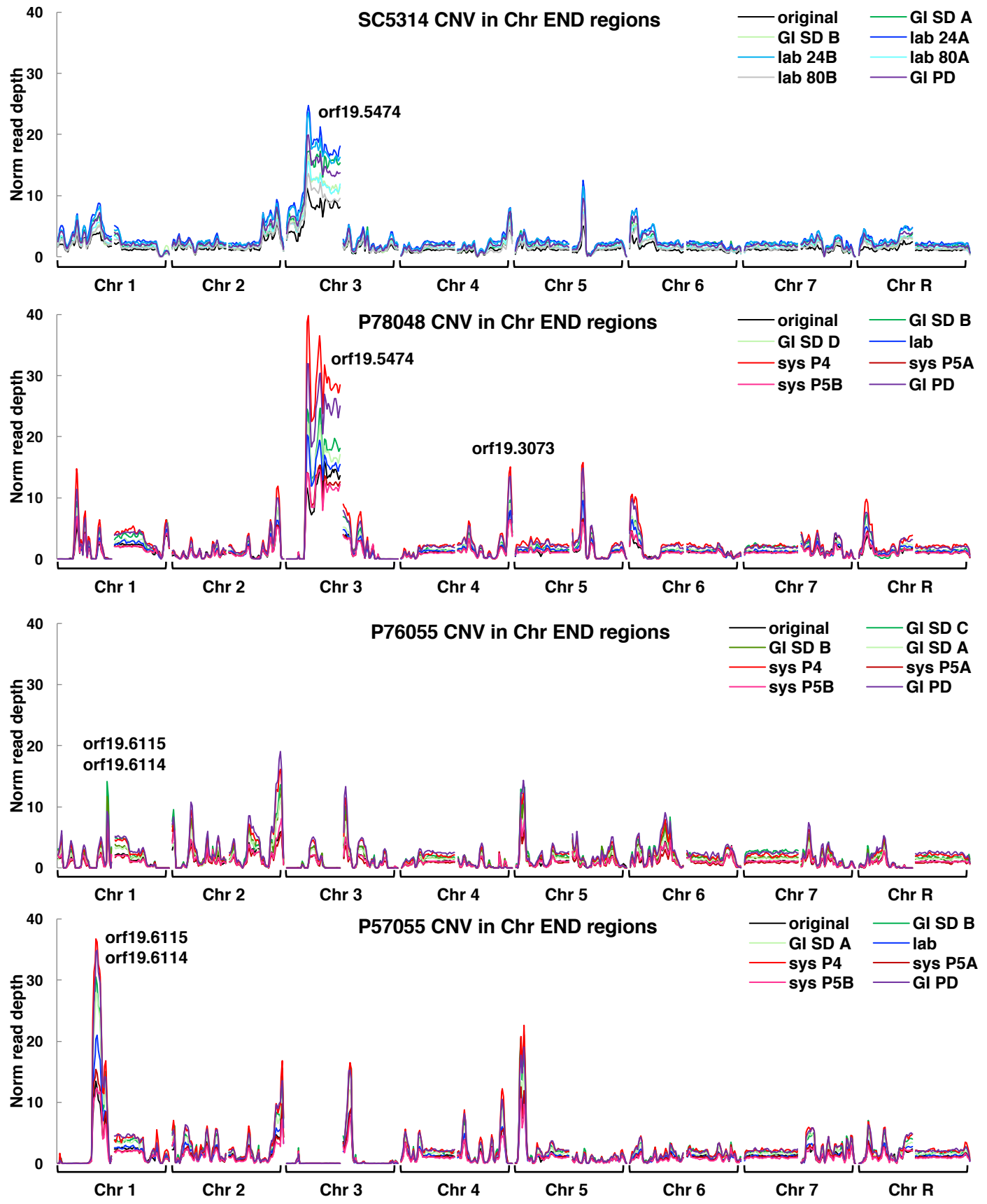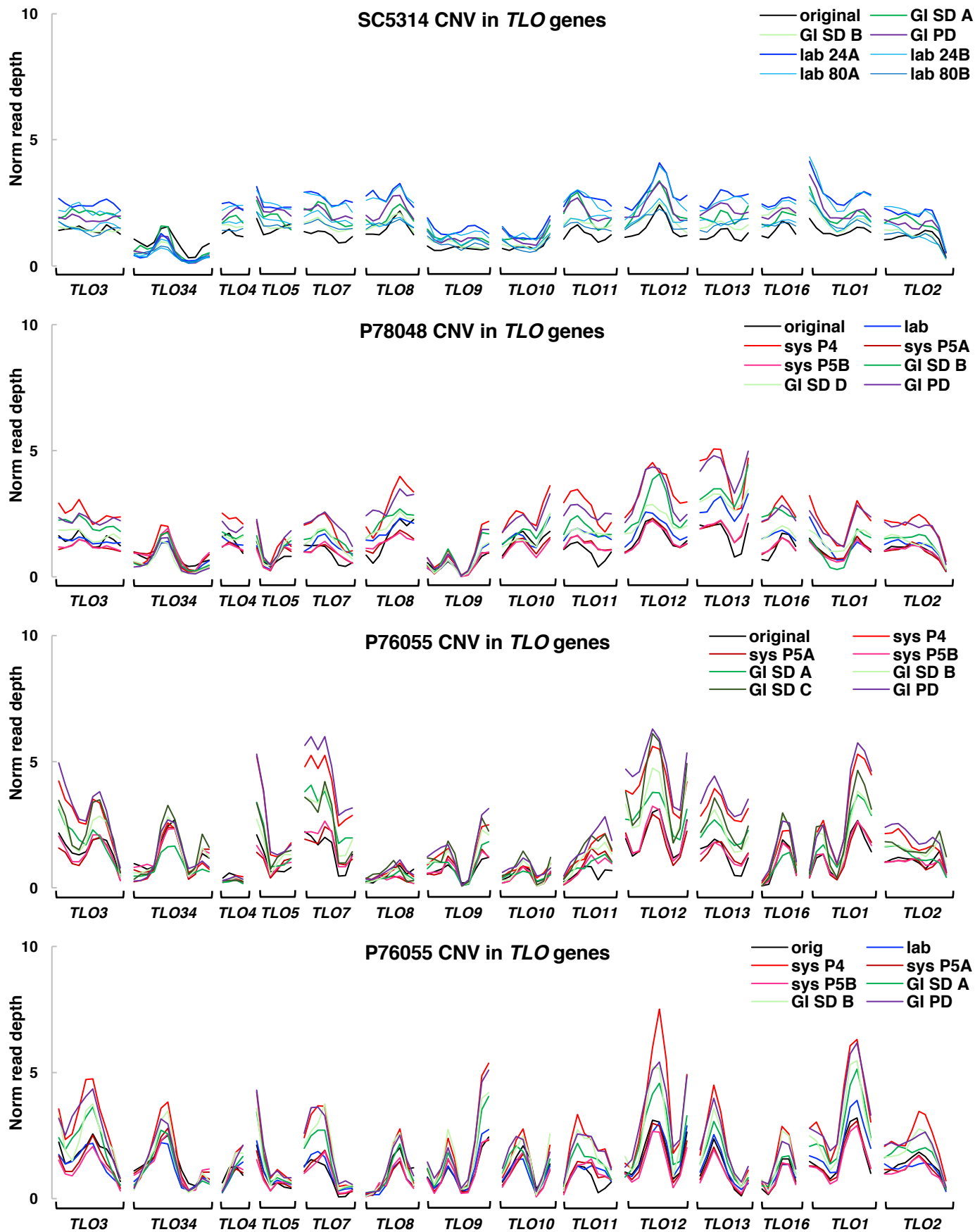| | | | |
|---|---|---|---|
| synonymous | 986 | 374 | 1877 |
| nonsynonymous | 641 | 308 | 1186 |
| readthrough | 1 | 0 | 2 |
| nonsense | 3 | 4 | 7 |
| indels | 11 | 5 | 14 |

**C**



**D**

**Supplementary Figure 6**

# Supplementary Figure 7

## A

**B**

**C**