

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18

# The Genome of the Human Pathogen *Candida albicans* is Shaped by Mutation and Cryptic Sexual Recombination

Joshua M. Wang<sup>1,2,3</sup>, Richard J. Bennett<sup>\*1</sup>, and Matthew Z. Anderson <sup>\*1,2,3</sup>

**Short title: Natural evolution of *C. albicans* genomes**

<sup>1</sup>Department of Molecular Microbiology and Immunology, Brown University, Providence, RI, 02912, USA.

<sup>2</sup>Department of Microbiology, The Ohio State University, Columbus, OH, 43210, USA.

<sup>3</sup>Department of Microbial Infection and Immunity, The Ohio State University, Columbus, OH, 43210, USA.

\* **correspondence:** Richard\_Bennett@brown.edu; anderson.3196@osu.edu

19 **ABSTRACT**

20           The opportunistic fungal pathogen *Candida albicans* lacks a conventional sexual  
21 program and is thought to evolve, at least primarily, through the clonal acquisition of genetic  
22 changes. Here, we performed an analysis of heterozygous diploid genomes from 21 clinical  
23 isolates to determine the natural evolutionary processes acting on the *C. albicans* genome.  
24 Consistent with a model of inheritance by descent, most single nucleotide polymorphisms  
25 (SNPs) were shared between closely related strains. However, strain-specific SNPs and  
26 insertions/deletions (indels) were distributed non-randomly across the genome. For example,  
27 base substitution rates were higher in the immediate vicinity of indels, and heterozygous regions  
28 of the genome contained significantly more strain-specific polymorphisms than homozygous  
29 regions. Loss of heterozygosity (LOH) events also contributed substantially to genotypic  
30 variation, with most long-tract LOH events extending to the ends of the chromosomes  
31 suggestive of repair via break-induced replication. Importantly, some isolates contained highly  
32 mosaic genomes and failed to cluster closely with other isolates within their assigned clades.  
33 Mosaicism is consistent with strains having experienced inter-clade recombination during their  
34 evolutionary history and a detailed examination of nuclear and mitochondrial genomes revealed  
35 striking examples of recombination. Together, our analyses reveal that both (para)sexual  
36 recombination and mitotic mutational processes drive evolution of this important pathogen in  
37 nature. To further facilitate the study of genome differences we also introduce an online  
38 platform, SNPMap, to examine SNP patterns in sequenced *C. albicans* genomes.

39

40 **AUTHOR SUMMARY**

41           Mutations introduce variation into the genome upon which selection can act. Defining  
42 the nature of these changes is critical for determining species evolution, as well as for  
43 understanding the genetic changes driving important cellular processes such as carcinogenesis.  
44 The fungus *Candida albicans* is a heterozygous diploid species that is both a frequent  
45 commensal organism and a prevalent opportunistic pathogen. Prevailing theory is that *C.*  
46 *albicans* evolves primarily through the gradual build-up of mutations, and a pressing question is  
47 whether sexual or parasexual processes also operate within natural populations. Here, we  
48 determine the evolutionary patterns of genetic change that have accompanied species evolution  
49 in nature by examining genomic differences between clinical isolates. We establish that the *C.*  
50 *albicans* genome evolves by a combination of base-substitution mutations, insertions/deletion  
51 events, and both short-tract and long-tract loss of heterozygosity (LOH) events. These  
52 mutations are unevenly distributed across the genome, and reveal that non-coding regions and  
53 heterozygous regions are evolving more quickly than coding regions and homozygous regions,  
54 respectively. Furthermore, we provide evidence that genetic exchange has occurred between  
55 isolates, establishing that sexual or parasexual processes have transpired in *C. albicans*  
56 populations and contribute to the diversity of both nuclear and mitochondrial genomes.

57

## 58 INTRODUCTION

59 A wide variety of genetic events contribute to the evolution of eukaryotic genomes. In  
60 asexual cells, haploid genomes evolve via the accumulation of point mutations as well as  
61 undergo recombination events that drive DNA expansions/contractions (indels). Heterozygous  
62 diploid genomes also have the capacity to experience loss of heterozygosity (LOH) events, in  
63 which genetic information is lost from one of the two chromosome homologs. In addition, both  
64 haploid and diploid genomes may experience large-scale chromosomal changes such as gross  
65 rearrangements, acquisition of supernumerary chromosomes or other forms of aneuploidy [1, 2].

66 Many eukaryotic species also generate genetic diversity via sexual reproduction. Here,  
67 recombination between individuals provides an efficient mechanism for producing diverse  
68 progeny. Sexual reproduction can therefore promote adaptation to new environments more  
69 rapidly than asexual propagation [3, 4]. However, this comes at a fitness cost due to the  
70 associated energetic requirements and the fact that only 50% of parental alleles are passed on  
71 to single progeny [5-7]. Sex can also be detrimental by breaking up beneficial allelic  
72 combinations [5, 8]. Facultative sexuality, the ability to alternate between sexual and asexual  
73 forms of reproduction, promotes a flexible lifestyle that can accelerate adaptation in response to  
74 environmental pressures [4, 9].

75 Sexual reproduction has been extensively studied in the *Saccharomyces* clade, where  
76 the model yeast *Saccharomyces cerevisiae* divides mitotically but can also undergo mating and  
77 meiosis to generate recombinant progeny. The related *Candida* clade includes some of the  
78 most important human fungal pathogens encountered in the clinic [10, 11], although the  
79 *Saccharomyces* and *Candida* clades diverged from one another ~235 million years ago [12].  
80 The most clinically-relevant *Candida* species is *C. albicans* that, like all *Candida* species, was  
81 originally designated an obligate asexual organism. However, mating of diploid cells has been  
82 observed in the laboratory and produces tetraploid cells that return to the diploid state via a

83 parasexual process of concerted chromosome loss (CCL) [13-16]. Mating requires that *C.*  
84 *albicans* cells undergo a phenotypic transition from the sterile “white” state to the mating-  
85 competent “opaque” state [17]. Conjugation of opaque cells can occur via heterothallic or  
86 homothallic mating [18], and recombination during CCL involves Spo11, a conserved 'meiosis-  
87 specific' factor involved in DNA double-strand break formation across diverse eukaryotes [14,  
88 19].

89         Clinical isolates of *C. albicans* exhibit a largely clonal population structure despite the  
90 potential for recombination via parasexual reproduction [20, 21]. Multilocus sequence typing  
91 (MLST) separates *C. albicans* isolates into 17 clades although previously described  
92 incompatibility between MLST haplotypes and individual mutations suggests that recombination  
93 may act to generate new allelic variants [22]. Analysis of a limited number of haploid  
94 mitochondrial loci also reveals allelic mixtures that suggest recombination may have occurred  
95 within *C. albicans* populations [23, 24]. However, despite these observations, *C. albicans* is still  
96 commonly assumed to be an asexual species that does not undergo mating or recombination in  
97 nature [25, 26]. Prior studies focused on a subset of genomic loci and present conflicting  
98 evidence regarding the role of recombination in shaping *C. albicans* evolution [20, 22-24, 27],  
99 which can now be addressed by a detailed analysis of full genome sequences.

100         In this work, we examined evolutionary patterns in 21 sequenced *C. albicans* isolates  
101 that represent different clades, different sites of infection in the host, and different countries of  
102 origin [28, 29]. Our analyses provides a detailed picture of how mutational events drive  
103 evolution of the diploid *C. albicans* genome. We reveal that mutations preferentially accumulate  
104 in heterozygous regions of the genome, and that emergent SNPs and indels often cluster  
105 together. Moreover, we highlight isolates whose nuclear and mitochondrial genomes appear  
106 highly admixed and therefore display evidence of genetic contributions from multiple clades.  
107 These results establish that the *C. albicans* genome is a dynamic landscape shaped both by

108 local mutations and large-scale rearrangements, and that sexual or parasexual mating has  
109 made a significant contribution to genotypic variation.

## 110 **RESULTS**

111           The availability of whole genome sequencing data for 21 diverse *C. albicans* isolates  
112 [28] provided an opportunity to determine how genetic diversity is generated between strains in  
113 nature. The *C. albicans* diploid genome is ~14 megabases (Mb) and consists of eight  
114 chromosomes encoding ~6100 genes [28, 30, 31]. SNPs occur at a frequency of ~0.3%  
115 between chromosome homologs in the standard laboratory strain SC5314 (i.e., an average of 1  
116 SNP every 330 bp) [28, 32]. Among the 21 isolates, we found SNP frequencies varied from  
117 0.5% between closely related strains within Clade I to 1.1% between strains from different  
118 clades (Table S1). A previous phylogenetic reconstruction using 112,223 SNP positions found  
119 that most strains matched their previously assigned fingerprinting clades and MLST subtypes,  
120 with the exception of P94015 which clustered separately from other Clade I strains (Fig. 1A)  
121 [28]. Strong bootstrap values across the constructed phylogeny of these strains supports a  
122 primarily clonal lifestyle in which most polymorphisms are mutations consistent with inheritance  
123 by descent (Fig. 1B). Accordingly, SNPs and indels fit a nonrandom distribution across the 21  
124 sequenced isolates  $\chi^2$  ((SNPs; 20, N = 302641) = 83118,  $p < 2E-16$ , indels; 20, N = 19581) =  
125 13825,  $p < 2E-16$ , Fig. S1).

126

### 127 **Base substitutions in *C. albicans***

128           LOH events can distort the patterns of SNPs inherited from ancestral strains (Fig. S2).  
129 To help limit these confounding effects, we restricted most analyses to strain-specific SNPs and  
130 indels that are unique to individual strains. Approximately 25% of all SNP positions and 10% of  
131 all indel positions were strain-specific (66,086 and 6,474 events, respectively; Tables S2, S3).  
132 As expected, the number of strain-specific mutations increased with longer branch lengths from  
133 the nearest node in the phylogenetic tree (SNPs,  $r_s = 0.60$ ,  $p = 4.2E-3$ ; indels,  $r_s = 0.42$ ,  $p =$   
134 0.055; Fig. 1B and Fig. S3). Correlation between these metrics of strain identity supports the  
135 use of strain-specific mutations in assessing mutational patterns.

136 In many eukaryotes, base-substitution mutations are biased towards transitions over  
137 transversions, although the cause of this bias is not completely clear [33]. In *C. albicans*, base  
138 substitutions also favored transitions over transversions for both strain-specific SNPs and total  
139 SNPs,  $\chi^2$  ((11, N = 66086) = 18182,  $p < 2E-16$  and (11, N = 302641) = 628000,  $p < 2E-16$ ,  
140 respectively). The ratio of transitions to transversions was 2.21 for strain-specific SNPs and  
141 2.50 for all SNPs (Fig. S4). Both coding and noncoding regions encoded more strain-specific  
142 transitions than transversions, although coding sequences were more biased than noncoding  
143 regions (2.74 versus 1.80, respectively). Base substitutions displayed a 1.39-fold bias towards  
144 introducing A/T instead of G/C for strain-specific SNPs that shrank to 1.03-fold when including  
145 all SNPs. The fact that substitutions favor transitions resulting in A/T suggests that this may  
146 contribute to the overall A/T richness of the *C. albicans* genome [31].

147

#### 148 **Distribution of strain-specific polymorphisms across the *C. albicans* genome**

149 Analysis of the global distribution of strain-specific SNPs revealed a bias against the  
150 accumulation of these mutations within protein-coding genes. Thus, most strain-specific  
151 polymorphisms (33,818 of 66,086 SNPs and 5,502 of 6,474 indels) were present within the  
152 36.7% of the genome representing intergenic regions, suggesting that mutations in coding  
153 sequences are selected against ( $p = 9.71E-16$ ; Fig. 2A,B). As a result, relatively few strain-  
154 specific SNPs were present within ORFs across the twenty-one sequenced strains (Fig. 2C).  
155 We found that 259 genes exhibited significantly greater SNP densities per nucleotide (nt) than  
156 the 0.004 SNPs/nt average for all *C. albicans* ORFs (Fig. 2C, Table S4). SNP densities within  
157 enriched genes were equal to or greater than the intergenic average (0.0066 vs. 0.0063,  
158 respectively). Protein-coding genes within this group lacked any enrichment for gene ontology  
159 (GO) annotations or pathways (Table S4). However, noncoding snoRNAs (small nucleolar  
160 RNAs) were significantly overrepresented among 'faster-evolving genes' by GO term analysis,



161  $\chi^2$  ((2, N = 5) = 15.6, p = 7.90E-5; Fig. 2D). The five snoRNAs identified from GO enrichment  
162 had mutation rates greater than 0.02 SNPs/nt, significantly higher than that of the average rate  
163 of 0.0063 SNPs/nt within intergenic regions. Strain-specific polymorphisms clustered towards  
164 the 5' end of the snoRNAs (Fig. 2E) and could contribute to variation in functional aspects of  
165 protein translation, although this possibility was not explored here.

166 An inspection of strain-specific indels revealed that 3527 (54.5%) were deletions and  
167 2948 (45.5%) were insertions. Indels ranged in size from 1 bp to 10 bp with the majority of  
168 longer events being insertions (Fig. S5). The frequency of both insertions and deletions  
169 decreased as mutations became larger, suggesting that smaller events occur more frequently or  
170 are less detrimental to the cell and therefore are retained more often. The incidence of  $\pm 3$   
171 nucleotide indels (21.9% of the total) was higher than that expected by chance. When indels  
172 were separated into genic or intergenic mutations, intergenic mutations followed a Poisson  
173 distribution centered on 0, whereas genic mutations were vastly overrepresented for  $\pm 3$   
174 nucleotide indels that do not shift the reading frame (Fig S5). Only ~15% of all indels fell within  
175 ORFs (p < 2.2E-16) suggesting that, as with SNPs, indels are selected against within coding  
176 sequences (Fig. 2B).

177 Indels have been commonly associated with specific genomic features such as repetitive  
178 sequences in other species [34, 35]. Across the sequenced *C. albicans* isolates, there was a  
179 total of 19,581 indel sites across the genome. Of these, 465 indel sites (2.37%) were located  
180 within annotated repetitive sequences (long terminal repeats (LTRs), major repeat sequences  
181 (MRSs), and retrotransposons). Total indels are therefore overrepresented within these  
182 repetitive features (two-tailed Brunner-Munzel (BM) test = 5.15, df = 182.05, p = 6.65E-7).  
183 Likewise, strain-specific indels were significantly enriched within repetitive features (47 of 6475;  
184 BM test = 13.004, df=182, p<2.2E-16). Both total and strain-specific SNPs also clustered within  
185 repetitive elements (BM test = 14.98, df = 315.62, p < 2.2E-16 and BM test = 12.26, df = 240.02,

186  $p < 2.2E-16$ , respectively). Thus, mutations within the *C. albicans* genome are enriched within  
187 repetitive regions similar to what has been observed in other species [34, 35].

188 Analysis of the genome-wide distribution of strain-specific SNPs and indels across the  
189 21 genomes revealed that these mutation types showed significant clustering with one another  
190 (Fig. 1C) (Pearson,  $t = 11.64$ ,  $df = 286$ ,  $p < 2.2E-16$ ). Multiple SNPs often occurred within 100  
191 bp of an indel (Fig. 1D) as was confirmed via Sanger sequencing of selected regions (Fig. S6).  
192 Enrichment of SNPs was observed immediately adjacent to indels (within 10 bp) but not within  
193 indels (Wilcoxon test ( $W(1.79E7)$ ),  $p < 2.2E-16$ ; Fig. 1E). Three strains, P60002, P75010 and  
194 P94015, encoded a large proportion of strain-specific mutations reflective of their longer branch  
195 lengths in the phylogenetic tree, which could potentially skew the analysis (Fig. 1A). However,  
196 even after removing these three strains from the analysis and reducing the four major clades to  
197 three representative strains each, we still observed a significant association between SNPs and  
198 indels (Wilcoxon test ( $W(3.04E7)$ ),  $p < 2.2E-16$ , Fig. S7A). This association highlights that  
199 distinct mutagenic events occur in close proximity to one another, and suggests that indel  
200 formation or the associated DNA repair processes may be mutagenic in *C. albicans*.

201 In some species, the introduction of indels can influence the observed mutational bias  
202 towards either transitions or transversions [35, 36]. To address this possibility in *C. albicans*,  
203 the transition:transversion ratio was determined for the ~500 strain-specific SNPs located within  
204 10 bp of strain-specific indels. Although base substitutions still slightly favored transitions, the  
205 1.17 transition:transversion ratio was significantly lower than the genome-wide average ratio of  
206 2.21 ( $p = 5.87E-7$ ). This is consistent with mutations close to indels exhibiting a reduced bias  
207 towards transitions over transversions due to recruitment of error-prone polymerases during  
208 DNA repair [35]. We therefore suggest that a similar mechanism operates in *C. albicans* and  
209 can account for the increased mutation rate adjacent to indels, as well as the local bias in the  
210 transition:transversion ratio.

211

## 212 **Association between LOH recombination events and base-substitution mutations**

213         The previous study by Hirakawa *et al.* identified extensive loss of heterozygosity (LOH)  
214 tracts in the 21 sequenced *C. albicans* isolates [28]. Consequently, LOH breakpoints were  
215 mapped in each isolate and emphasis was placed on the distribution of LOH events around the  
216 mating type-like (*MTL*) locus on Chr5. The current study extends the analysis of LOH patterns  
217 in *C. albicans* genomes by determining if genome-wide patterns of LOH exist, and if there is an  
218 association between LOH tracts and other mutational classes such as base substitutions or  
219 indels.

220         LOH regions were defined in Hirakawa *et al.* using several parameters including  
221 contiguous 5 kilobase (kb) windows with a high frequency of homozygous SNPs (>0.4 events  
222 per kb; see Methods and [28]). Plotting the incidence of LOH for all chromosomes (Chr) in each  
223 of the isolates revealed a striking pattern, whereby the prevalence of LOH increased along each  
224 chromosome arm when progressing from centromere to telomere (Fig. 3A and Fig. S8). In fact,  
225 the overwhelming majority of all long-tract LOH regions (155 out of 170 regions larger than 50  
226 kb) extended to the ends of the corresponding chromosomes (Table S5). This reveals that out  
227 of a total of 336 chromosome arms in the 21 isolates, 155 of these arms show evidence of  
228 having undergone a long-tract LOH event. LOH frequency decreased towards the centromeres  
229 and did not occur across centromeres except during LOH of whole chromosomes (Fig. 3B).  
230 Interestingly, LOH frequencies remained low across the entirety of the right arms of Chr2 and  
231 Chr4 (Fig. S8), suggesting that heterozygosity of loci on these arms may be maintained by  
232 selection. Aneuploidy did not significantly alter the frequency of heterozygous and homozygous  
233 intervals along aneuploid chromosomes relative to euploid chromosomes ( $p = 0.756$ ).

234         Several studies have revealed that mutation rates can be impacted by the underlying  
235 genomic context. For example, accumulation of SNPs was found to be increased in regions  
236 adjacent to indels in diverse eukaryotic species [36], and mutation rates were higher in

237 heterozygotes than in homozygotes during meiosis [37]. We therefore examined mutational  
238 patterns in *C. albicans* genomes that are a mosaic of heterozygous and homozygous regions.  
239 We subdivided *C. albicans* genomes into heterozygous (het) or homozygous (hom) regions  
240 using defined criteria on all SNPs (see Methods), resulting in 468 het and 445 hom regions,  
241 respectively (Table S5). Het regions covered a total of 71.1% of the genome and hom regions  
242 28.9%; het tracts were therefore considerably longer on average than hom tracts (~480,000 bp  
243 vs. 186,000 bp, respectively). Definition of het and hom regions using all SNPs allowed  
244 subsequent examination of the frequencies of strain-specific SNPs within these regions. Strain-  
245 specific SNPs comprise only 3% of all SNPs within these genomes and, therefore, do not  
246 contribute substantially to the designation of het and hom regions. Notably, het regions  
247 contained significantly higher frequencies of strain-specific SNPs than hom regions ( $1.4E-4$  vs.  
248  $7.3E-5$  SNPs/bp, respectively; BM test = -10.6, df = 786.6,  $p < 2E-16$ ; Fig. 3C). Even after the  
249 exclusion of the three outlier strains, P60002, P75010 and P94015, and reducing the four major  
250 clades to three representative strains each, there were still significantly higher frequencies of  
251 strain-specific SNPs within het than within hom regions (BM test = -7.558, df = 377.0,  $p = 3.14E-$   
252 13). Furthermore, all 21 isolates exhibited the same bias towards het regions containing more  
253 strain-specific SNPs than hom regions (two-tailed BM test = -1.11, df = 38.5,  $p = 0.28$ ),  
254 indicating that mutations preferentially accumulate in het over hom regions during natural  
255 evolution of *C. albicans* isolates.

256 We note that heterozygous SNPs may have arisen in hom regions but, in some cases,  
257 been eliminated by a subsequent LOH event. As LOH can occur in one of two possible  
258 directions (due to loss of either homolog A or homolog B), we accounted for mutations  
259 potentially lost via LOH by doubling the number of homozygous, strain-specific SNPs within  
260 hom regions. Even with this adjustment, het regions still contained a greater density of strain-  
261 specific SNPs than hom regions (two-sided BM test = -8.74, df = 717.58,  $p < 2.2E-16$ ). The ~2-  
262 fold greater accumulation of polymorphisms in het over hom regions of the *C. albicans* genome

263 shows parallels with the ~3.5-fold higher mutation rate observed in het vs. hom regions of the  
264 *Arabidopsis* genome during meiosis [37].

265 Sites close to recombination events, including LOH events, have been shown to be  
266 associated with elevated mutation rates in some species [36, 38-40]. To determine if there is an  
267 increased frequency of SNPs in regions proximal to LOH tracts in *C. albicans*, the density of  
268 SNPs at heterozygous-homozygous transition points was investigated. Analysis of the 745  
269 identified transition regions included 1 kb of DNA on either side of the junction points between  
270 het and hom regions (with the latter inferred to represent LOH tracts). The SNP density within  
271 these transition regions was significantly lower than that in the rest of the *C. albicans* genome  
272 (one-sided BM test = -35.415, df = 748.8,  $p < 2.2E-16$ ; Fig. 3D). Furthermore, SNP density was  
273 similar on both het and hom sides of the LOH breakpoint. Thus, base substitutions appear to  
274 accumulate less frequently in regions proximal to het/hom breakpoints in the *C. albicans*  
275 genome. One caveat noted here is that this result may be influenced by difficulty in the  
276 identification of precise breakpoints between het and hom regions of the genome.

277

### 278 **Identity by descent during *C. albicans* evolution**

279 A hallmark of phylogenetic reconstructions in asexual species is the ability to track the  
280 relatedness of isolates based on inherited polymorphisms [22, 41, 42]. Reconstruction often  
281 relies on a maximum parsimony model of 'identity by descent', in which more closely related  
282 strains share a greater percentage of shared polymorphisms (Fig. 4A). *C. albicans* SNP  
283 patterns generally follow identity by descent; these patterns matched the phylogenetic tree as  
284 evidenced by strong bootstrap support at almost all nodes [28], as well as a visual examination  
285 of SNP patterns within specific regions (Fig. 4B,C, and Fig. S9). Strikingly, however, certain  
286 regions of the genome exhibited clear violations of identity by descent (Fig. 4D,E, and Fig. S10).  
287 In heterozygous diploid genomes, these deviations could potentially arise through two  
288 mechanisms: (1) by sexual recombination between genetically distinct isolates, or (2) through

289 multiple, independent LOH events that obfuscate the actual pattern of descent (Fig. S11). In the  
290 latter case, multiple LOH events could cause loss or retention of SNPs through homozygosis of  
291 one chromosome homolog or the other, thereby generating a subset of isolates that appear  
292 “recombinant”, i.e., appear to have intermixed genetic content from two different relatives. Such  
293 a history can sometimes be inferred by a comparison of SNP patterns within the region of  
294 interest in multiple extant strains (Fig. 4A).

295 We note that certain regions appear to have undergone divergent, short-tract LOH  
296 events across the 21 isolates, consistent with these events often occurring during asexual  
297 divisions (Fig. 4D,E, Fig. S10 and Supplementary Material). In line with this, we identified 514  
298 non-overlapping 25 kb windows across the 21 sequenced isolates that do not encode similar  
299 polymorphisms to closely-related strains, and could represent regions that had experienced  
300 LOH. Interestingly, two strains, P60002 and P94015, contained 390 of these regions (75.9%),  
301 although only 214 of the 390 incongruent regions in these two strains (54.9%) overlapped with  
302 LOH tracts (hom regions) in these isolates. In contrast, the majority of the incongruent regions  
303 in all other strains (117 of 124 regions) overlapped with LOH tracts. This suggests that  
304 incongruence in polymorphisms in most strains likely results from divergent LOH events but that  
305 LOH does not obviously explain the majority of incongruent polymorphic patterns observed in  
306 P60002 and P94015.

307

### 308 **Evidence for recombination in natural isolates of *C. albicans***

309 Previous studies have provided conflicting messages regarding recombination in natural  
310 populations of *C. albicans* [20, 22-24, 27], and none have examined whole genome data for  
311 evidence of inter-clade mixing. We therefore examined the 21 sequenced genomes for mixed  
312 evolutionary histories. The similarity of genomic segments from each strain to the overall  
313 phylogenetic tree was compared by analysis of SNP patterns using 25 kb sliding windows. To  
314 aid visualization of SNP patterns we developed a custom interactive tool, SNPMap

315 (<http://snpmap.asc.ohio-state.edu/>), which allows users to map the positions of individual  
316 mutations, mutation types, and het/hom tracts across user-defined regions of the 21 genomes.

317 Our analysis largely focused on two clinical isolates, P60002 and P94015, that have the  
318 weakest bootstrap support within the *C. albicans* phylogeny and that cluster with different  
319 strains when analyzed by MLST or DNA fingerprinting [28], suggesting their genomes may be  
320 recombinant. In support of this, examination of Chr4 in P94015 identified one region with clear  
321 homology to Clade I which was in close proximity to a region highly homologous to Clade SA  
322 (Fig. 5A, Table S6). The region with homology to Clade I (labeled P94015-A in Fig. 5A) shares  
323 a large number of SNPs that are present throughout Clade I but are absent in all other strains  
324 with the exception of P94015. Next to this region, a 1 kb segment (region P94015-B) lacks  
325 clear homology to any of the other sequenced isolates while, adjacent to this, the SNP pattern in  
326 P94015 is virtually identical to that of two Clade SA isolates (region P94015-C).

327 Mating between isolates from different clades would be expected to generate hybrid  
328 DNA regions, with SNPs on one homolog of the recombinant strain matching those in one clade  
329 and SNPs on the other homolog matching those in a second clade. Identifying inherited SNPs  
330 following *C. albicans* mating in nature is complicated by the fact that, with the exception of  
331 SC5314 [32], haplotypes are not available for the 21 *C. albicans* genomes. Despite this,  
332 phasing of heterozygous SNPs for some isolates can be inferred using SNP patterns from  
333 related strain(s) that have undergone LOH for that region (Fig. 5B). The region that  
334 experienced LOH will only retain the SNPs that reside on the same homolog (i.e., those that are  
335 phased). Using this approach, we phased SNPs for chromosomal regions of closely related  
336 strains that are heterozygous in some isolates but homozygous in others. Multiple isolates that  
337 have undergone LOH for both alleles strengthen the confidence of phasing assignments within  
338 a given clade.

339 We applied this approach to a region on Chr2 in P94105 that contains polymorphisms  
340 identical both to those on homolog A of a Clade SA strain (12 of 12 SNPs are identical) and to

341 those on homolog A from a Clade III isolate (15 of 15 SNPs are identical) (Fig. 5B). Both SNP  
342 positions and nucleotide identities are conserved across this hybrid region in P94015 when  
343 compared to the corresponding homologs from Clade SA and Clade III isolates. This region  
344 therefore provides a striking example of P94015 inheriting one homolog from a Clade SA strain  
345 and one homolog from a Clade III strain, and establishes a non-clonal origin for this isolate.  
346 Analysis of additional regions for isolates P94015 and P60002 provides support for the  
347 existence of multiple recombination tracts indicating that mixing has occurred between strains  
348 from different *C. albicans* clades (Fig. S12).

349 To examine global patterns of admixing among the set of 21 isolates, the distribution of  
350 all variant positions in each strain was compared to the consensus pattern for each clade using  
351 sliding 25 kb windows. The SNP patterns of most isolates resembled the consensus pattern for  
352 their assigned clade (98.5% of genomic windows matched their assigned clade), as expected  
353 for a population propagating clonally (Fig. 5C). In contrast, many regions within the P60002 and  
354 P94015 genomes showed homology to multiple different clades, producing highly mosaic  
355 genomes (Fig. 5C). Here, the genomes of P60002 and P94015 matched their assigned clades  
356 for only 58.3% and 76.7% of sliding windows, respectively ( $p = 1.14E-10$ ). The majority of  
357 genomic regions in P94015 aligned with Clade I (genome is mostly red in Fig. 5C), whereas  
358 numerous segments aligned to regions from three other major clades (SA, II, and III). In the  
359 case of P60002, Clade SA regions made up the majority of the genome with a smaller number  
360 of regions matching Clade I or, to a lesser extent, Clade II. In line with this, the branchpoint  
361 leading to P60002 is the least well-supported node in the phylogenetic reconstruction of all 21  
362 isolates [28]. The most parsimonious explanation for these highly mosaic genomes is that they  
363 are the products of mating and recombination between isolates from multiple *C. albicans* clades.

364

365 **Analysis of mitochondrial genomes in *C. albicans* isolates**



366 Haploid mitochondrial genomes provide a more simplified context to search for evidence  
367 of recombination than heterozygous diploid genomes. In *S. cerevisiae*, mitochondrial genomes  
368 are biparentally inherited following mating, and recombination can occur between parental  
369 genomes prior to zygote division [43]. We therefore performed the first comparative analysis of  
370 global SNP patterns in *C. albicans* mitochondrial genomes using sequencing data from the set  
371 of 21 isolates. The mitochondrial genome in *C. albicans* is ~41 kb in size and a total of 1847  
372 SNPs (and 0 indels) were annotated within the 21 isolates, with an average SNP density of 1  
373 polymorphism every 476 bp. The mitochondrial genome was highly heterogeneous including  
374 areas of high SNP density (e.g., ChrM: 15000-20000) and regions devoid of polymorphisms  
375 (e.g., ChrM: 6000-12000) among sequenced isolates. Furthermore, of the 1847 annotated  
376 mitochondrial SNPs, only 39 were strain-specific in the set of 21 sequenced genomes (Table  
377 S2).

378 *C. albicans* mitochondrial genomes generally showed clade-specific SNP patterns that  
379 were again consistent with a clonal population structure, although resolution of SNP patterns  
380 was low due to relatively few clade-defining mitochondrial SNPs (Fig. 6). As with nuclear  
381 genomes, examination of mitochondrial genomes of P60002 and P94015 again showed clear  
382 evidence of inter-clade mixing. For example, the mitochondrial genome of P94015 contained  
383 regions that aligned with mitochondrial segments from both Clade SA and Clade II (Fig. 6A).  
384 Here, there are three polymorphisms that are Clade SA-specific on ChrM: 1-6000 (region  
385 P94015-A), and all three are present in P94015 (Fig. 6A). An additional two of the fifteen  
386 P94015 polymorphisms in this region are specific to this strain (Table S2). The remainder of the  
387 mitochondrial genome in P94015 (region P94015-B) encodes 144 polymorphisms matching the  
388 SNP pattern found in Clade II (with the exception of two strain-specific SNPs). Recombination  
389 between Clades I and II was even clearer in the mitochondrial genome of P60002, as the  
390 majority of this genome was identical to Clade I (P60002 – regions A and C), but a 4-kb  
391 segment (ChrM: 19000-22000) encoded 34 polymorphisms that were identical to Clade II and

392 were entirely absent in Clade I (P60002 – region B, Fig. 6B and Fig. S13). We also found that  
393 maximum parsimony approaches mostly separated the P90145 mitochondrial genome into  
394 Clade II and SA regions and the P60002 mitochondrial genome into Clade I and II regions (Fig.  
395 6D), consistent with visual alignments (Fig. 6A,B). Direct Sanger sequencing of the  
396 mitochondrial genome (regions 5000-5700,18500-19500 and 29000-30000) supported the SNP  
397 designations from whole genome sequencing and therefore establish that recombination has  
398 occurred within the mitochondrial genomes of P60002 and P94015 (Fig. 6E, S14).

399 We further note that the single representative of Clade E in our collection, P75010, also  
400 displays strong evidence of recombination in its mitochondrial genome (Fig. 6C). The first ~12  
401 kb of P75010 (P75010-A) aligns closely with Clade II and encodes all three clade II-specific  
402 SNPs in this region. The following ~7 kb region of P75010 (P75010-B) matches that of Clade I  
403 strains, and is followed by a region with clear homology to Clade SA, encoding 22 of 25 Clade  
404 SA-specific SNPs. Identity to several clades infers that multiple recombination events have  
405 given rise to this complex SNP pattern. Taken together, these results reveal recombination  
406 events have occurred within *C. albicans* mitochondrial genomes and provide clear evidence that  
407 sexual/parasexual processes have occurred during *C. albicans* evolution.

408

409

410

## 411 **DISCUSSION**

412           Our analysis of *C. albicans* genome structure reveals a number of important aspects  
413 concerning mutational patterns during natural evolution of the species. We highlight that (1)  
414 non-coding and heterozygous regions of the genome accumulate more mutations than coding  
415 and homozygous regions, respectively, (2) there is a significant association between the  
416 positions of emergent SNPs and indels, (3) diverse LOH events contribute to genetic  
417 inheritance, including long-tract LOH events that extend to the ends of the chromosomes, (4)  
418 there is evidence for selection acting on natural populations, (5) a subset of strains exhibit  
419 mosaic nuclear and mitochondrial genomes, and (6) analysis of specific chromosomal regions  
420 reveals clear evidence for inter-clade recombination.

### 421 **Mutations driving natural evolution of the *C. albicans* genome**

422           Mutation rates vary across eukaryotic genomes in a context-dependent manner [44, 45].  
423 We found that mutation patterns arising in natural *C. albicans* populations similarly exhibit a  
424 non-random distribution across the genome. Our analysis focused on the distribution of strain-  
425 specific SNPs, as these SNPs are likely to have emerged since these strains diverged from one  
426 another. We found that the location of strain-specific SNPs was biased towards heterozygous  
427 over homozygous regions of the genome. This is consistent with recent studies that showed  
428 that mutations arising during meiosis occurred more frequently within heterozygous than  
429 homozygous regions of eukaryotic genomes, although the mechanism driving this bias is  
430 unknown [46, 47]. Our results extend these findings by indicating preferential accumulation of  
431 mutations within heterozygous regions during mitotic growth in *C. albicans*, suggesting that  
432 elevated mutation rates may be a common feature associated with heterozygosity.

433           Our analysis also reveals that emergent SNPs and indels cluster together within the *C.*  
434 *albicans* genome, with a significant enrichment of SNPs within 10 bp of emergent indels. This is  
435 similar to what has been observed in other eukaryotic species where indels were found to  
436 promote elevated substitution rates close to the founding indel [34-36, 48, 49]. Moreover, the

437 transition:transversion ratio was significantly lower at these mutation sites compared to the  
438 genome average. This is consistent with a model in which the recruitment of error-prone  
439 polymerases results in increased mutation rates during indel formation [35].

440         Given that recombination events can be mutagenic [37, 40, 50], we also examined  
441 whether *de novo* SNPs were more prevalent close to the breakpoints between  
442 heterozygous/homozygous regions. However, we found that the regions flanking the borders of  
443 *C. albicans* LOH tracts encoded fewer mutations than the genome average. Other studies have  
444 similarly found varying associations between recombination crossovers and mutation rates;  
445 some recombination break points exhibited no effect on mutation rates [51], some showed  
446 increased rates [47], and some displayed reduced rates [52], similar to the current analysis in *C.*  
447 *albicans*.

448         The location of strain-specific polymorphisms was biased towards non-coding segments  
449 of the *C. albicans* genome. Overall, 51.2% of strain-specific SNPs and 85.0% of strain-specific  
450 indels resided in non-coding regions, despite these regions representing just 36.7% of the  
451 genome. This indicates that selection is likely limiting both SNP and indel mutations from  
452 accumulating in protein-coding sequences in *C. albicans*. Where mutations did occur within  
453 genes, these were biased towards certain gene classes and mutation types. Genes containing  
454 the highest SNP frequencies encoded snoRNAs, consistent with studies in other eukaryotes  
455 [53]. High expression of these noncoding RNAs makes them particularly vulnerable to  
456 mutations during collisions between DNA and RNA polymerases [54] and these genes are also  
457 generally more accepting of mutations than protein-coding genes [55].

#### 458 **Evolutionary impact of loss of heterozygosity**

459         Loss of heterozygosity events are a frequent occurrence in *C. albicans* genomes [1, 56].  
460 Previous analysis of the 21 *C. albicans* genomes noted that LOH tracts can be highly variable in  
461 size, with several isolates containing chromosomes that had experienced whole chromosome  
462 LOH [28]. Here, we made the observation that the vast majority of very large LOH tracts

463 (defined as LOH tracts >50 kb that were not whole chromosome LOH events) initiated between  
464 the centromere and the telomere and extended all the way to the chromosome ends. Such  
465 LOH events are likely due to break-induced replication, in which one chromosome homolog is  
466 used as a template to repair a double-strand DNA break in the other homolog, although  
467 reciprocal crossover events can also generate long LOH tracts [26, 57, 58]. The breakpoint for  
468 most of these long-tract LOH events varied between individual strains suggesting that  
469 independent LOH events had occurred. Short-tract LOH (<50 kb) was also common with  
470 approximately of half of these events being shared between related strains (Fig. 3, Fig. S8,  
471 Table S5), indicating that homozygosity of certain regions may confer a selective advantage or  
472 are not sufficiently deleterious to be selected against. LOH of specific alleles has been shown  
473 to alter phenotypes that impact a range of *C. albicans* traits from growth rates to virulence to  
474 drug resistance [29, 59, 60].

475 The cumulative effects of all of these mutational forces on *C. albicans* genomes are  
476 accelerated evolution of heterozygous regions relative to homozygous regions. Eventually, the  
477 rapid accumulation of deleterious mutations during clonal growth would be expected to result in  
478 a fitness decline due to Muller's ratchet [61, 62]. The occurrence of LOH may counterbalance  
479 these forces both by culling mutations from the genome and by reducing heterozygosity to slow  
480 evolutionary rates across the genome. LOH events appear at different points during the  
481 evolution of individual strains as most LOH tracts are not shared even between closely-related  
482 isolates. Indeed, accumulation of emergent strain-specific SNPs within more ancestral LOH  
483 tracts demonstrates that these LOH events are not recent occurrences and have further  
484 mutated since their origin.

#### 485 **Evidence for genetic exchange in clinical *C. albicans* isolates**

486 Studies of *C. albicans* population structure point to a largely clonal mode of reproduction,  
487 yet there is also evidence of mixed evolutionary histories indicative of a sexually/parasexually

488 reproducing species [21, 22, 24, 27, 63-65]. Furthermore, genes encoded at the mating locus  
489 show evidence for ongoing selection consistent with a conserved role in regulating  
490 sexual/parasexual reproduction [66]. In this study, we interrogated whole genome data for  
491 evidence of genetic admixture, while noting that LOH events can complicate analysis of  
492 recombination in diploid strains (Fig. 4, S10). We reveal that a subset of isolates contain  
493 mosaic genomes, consistent with these genomes being the products of mating between  
494 different *C. albicans* clades. Both nuclear and mitochondrial genomes of P60002 and P94015  
495 show recombinant genotypes supporting a (para)sexual origin for these strains. Genetic  
496 information from multiple clades contributed to these genomes and recombination tracts varied  
497 in length from a few kb to hundreds of kb (Fig. 5 and 6). The existence of a subset of *C.*  
498 *albicans* strains with mosaic genomes is similar to what has been observed in wild and  
499 domesticated strains of *S. cerevisiae*, where both non-mosaic and recombinant mosaic  
500 genomes have been identified [67]. Analysis of admixed genomes in P60002 and P94015  
501 suggests that recombination events may be relatively ancient, as recombination involves  
502 multiple clades and more recent mutational events have obscured the precise evolutionary  
503 histories of these strains.

504         Disagreement between strain phylogeny by MLST and Ca3 fingerprinting may also be  
505 indicative of recombination in the population. Based on Ca3 fingerprinting, P94015 should  
506 cluster with other Clade I strains and is supported by MLST analysis, which groups P94015 and  
507 other MLST 6 strains closer to MLST 1 / Clade I than other groups [28]. Yet, whole genome  
508 sequencing clusters P94015 closest to Clade SA (and even Clade II and Clade III strains) than  
509 other Clade I strains. This could reflect how recombination can distort the phylogenetic  
510 relationship between strains when based on analysis of a small subset of loci. Analysis of  
511 additional *C. albicans* isolates will help define how prevalent recombination is in the species and  
512 whether these recombination events are ancient or more recent occurrences.

513           Critically, we were able to identify regions in *C. albicans* genomes that exactly matched  
514 the pattern of recombinant SNPs expected from mating events between two extant clades. This  
515 was exemplified by one region in P94015 which consisted of multiple SNPs that exactly  
516 matched those present in Clade I followed, after a short gap, by a run of SNPs that precisely  
517 matched those in Clade SA (Fig. 5B). Recombination events were also clearly evident in  
518 mitochondrial genomes of at least 3 of the 21 isolates examined (P60002, P75010 and  
519 P94015). Taken together, these studies provide the clearest evidence to date that *C. albicans*  
520 populations have been shaped by (para)sexual exchange.

521           In summary, *C. albicans* genomes reveal multiple signatures of the forces that have  
522 shaped genetic diversity within the species. Both short and long LOH events have played a  
523 major role in increasing population diversity, with large tracts extending along the terminal  
524 regions of many chromosome arms that impact hundreds to thousands of polymorphisms.  
525 Base-substitution mutations and indels cluster within heterozygous regions of the genome,  
526 suggestive of faster evolution of these regions, and recombination between isolates has  
527 generated mosaic nuclear and mitochondrial genomes with potentially profound consequences  
528 for adaptation. The diploid heterozygous genome of *C. albicans* is therefore a highly dynamic  
529 platform on which selection can act.

530

## 531 **Materials and Methods**

### 532 **Variant calling, Processing, and Display**

533 Whole genome sequencing, variant identification, and Loss of Heterozygosity (LOH)  
534 windows were identified in a previous study (Hirakawa et al. 2015). Briefly, BWA 0.5.9 [68] read  
535 alignments were filtered with a minimum mapping quality of 30 using SAMtools [69]. To reduce  
536 the incidence of false positive SNPs called near indels, poorly aligned regions were realigned  
537 using GATK RealignerTargetCreator and IndelRealigner (GATK version 1.4-14) [70]. Both prior  
538 SNP variants and variants present in the mitochondrial genomes and indels for both  
539 mitochondrial and nuclear genomes described here used GATK UnifiedGenotyper, and filtered  
540 using the GATK VariantFiltration using hard filters (QD<2.0, MQ<40.0, FS>60.0,  
541 MQRankSum<-12.5, ReadPosRankSum<-8.0). The genome sequences used in this study are  
542 available under BioProject ID PRJNA193498 (<http://www.ncbi.nih.gov/bioproject>). SNP data is  
543 available from dbSNP (<http://www.ncbi.nih.gov/projects/SNP>) under noninclusive ss#  
544 1456786277 to 1457237021. SNPs were called homozygous if greater than 90% of the reads  
545 contained the non-reference nucleotide. This high threshold also reduced miscalling due to  
546 trisomic chromosomes/chromosomal regions.

547 To assess heterozygous and homozygous variant calls, the number of reads at each  
548 variant position was divided by the total number of reads at that position in that strain. On  
549 average, each variant position had 52.17 reads with an interquartile range from 33 to 71 reads.  
550 The distribution of the allelic ratio at each variant is plotted in Figure S15. The mean number of  
551 reads of the allele divided by the total number of reads for heterozygous positions was 0.4499  
552 and for homozygous regions, 0.9889. Thus, a 90% threshold was used for homozygous regions  
553 whereas heterozygous regions span from 10-50%.

554 From this dataset, strain-specific SNPs and indels were parsed into a separate set for  
555 additional analysis. Strain-specific variant features were required to be uniquely called in only  
556 one of the 21 strains at its genomic position. The data sets are available online in a searchable



557 interface, using R shiny for the backend: [snpmapp.asc.ohio-state.edu]. Manual interrogation of  
558 100 variants using the Integrated Genome Viewer [71] confirmed the variant presence and call  
559 quality metrics in all 100. Manual interrogation of SNP-indel pairs in the pileup confirmed 98%  
560 (49 of 50) are valid using the same criteria.

561 To adjust for mutations in homozygous regions that may have occurred but then lost due  
562 to LOH, the following calculation was used for each homozygous region:  $((\text{homozygous}$   
563  $\text{SNPs}/\text{all SNPs})+1)*\text{homozygous SNPs}/\text{length}$ . This effectively doubles the number of  
564 homozygous SNPs in hom regions to account for heterozygous SNPs that are lost due to LOH.  
565 Even with this doubling, significantly more emergent SNPs occurred in het regions than in hom  
566 regions.

#### 567 **Phylogenetic construction of strain relatedness**

568 The methods used to construct the phylogeny of the sequenced strains was previously  
569 described in Hirakawa *et al.* [28]. Briefly, the phylogenetic relationship of strain relatedness was  
570 constructed using all whole genome SNP calls, which totaled 113,339. A distance based tree  
571 was estimated relying on maximum parsimony and a stepwise matrix where homozygous  
572 positions are two steps away compared to one step for heterozygous positions. SNP positions  
573 were resampled in 1,000 bootstrapped samples and each node indicates the bootstrap support.

#### 574 **Determination of LOH**

575 Previously defined heterozygous and homozygous genomic region were used as  
576 detailed in [28]. Briefly, the SNP density of homozygous and heterozygous positions was  
577 calculated across the genome in non-overlapping 5 kb windows for each isolate. The resulting  
578 5 kb windows were managed as follows: 1) Homozygous regions shared between individual  
579 strains and SC5314 were identified and marked as homozygous; 2) a single 5 kb window  
580 adjacent to these regions lacking any polymorphisms was merged into homozygous tracts if  
581 present; 3) contiguous, adjacent windows with a significantly higher frequency of homozygous  
582 SNPs than SC5314 homozygous regions ( $>0.4$  SNPs/kb) were merged allowing one intervening

583 window lacking sufficient polymorphism into homozygous tracts. These regions were defined as  
584 homozygous whereas remaining regions covered by 5 kb windows of the genome were  
585 designated heterozygous and contained significantly more heterozygous SNPs. The borders  
586 between homozygous and heterozygous regions were manually inspected for accuracy.

### 587 **Introgression/Tree Violations**

588 Two independent methods were used to assign clade designations for nuclear and  
589 mitochondrial genomes in non-overlapping 25,000 or 750 bp windows, respectively. The first  
590 process assigned the clade that best resembled the query strain's SNP pattern. However,  
591 Clade I contains fewer SNP calls due to alignment to the SC5314 reference genome (also a  
592 Clade I isolate) that can introduce an artificial bias. Therefore, a dataframe was constructed  
593 where each row represented any SNP contained within any strain in the query window. A SNP  
594 could only be counted as a single row so the identity of that SNP position was recorded as "0"  
595 for the absence of the SNP, "1" for a heterozygous position, and "2" for a homozygous position.  
596 The correlation between the target strain's resulting numeric "SNP" profile and each other  
597 strains was individually calculated. Scores from strains within the same clade were averaged  
598 and taken as the absolute value:  $Clade\ score = abs(mean(cor(query\ strain\ SNP, strain\ X$   
599  $SNP)))$ . The clade with the highest similarity score (and thus greatest proportion of shared  
600 SNPs) was selected as the most similar clade for that window. This process was repeated  
601 across the full genome. As a followup, this process was repeated with removing all strains  
602 within the query strain's clade and the scores recalculated.

603 The second phylogenetic approach used an expression matrix listing the 21 strains  
604 against all possible SNP positions present for each non-overlapping window. For each cell, we  
605 denoted a 0 if the respective SNP was not present in the respective strain, a 1 if one copy of the  
606 SNP was present, and a 2 if both copies were present. A distance matrix using a binary method  
607 (R dist) was constructed from this data and finally resulted in a phylogenetic tree (R hclust). The  
608 appropriate number of K clusters for the phylogenetic output of each window was estimated by

609 traversing through all 21 possible values (1 to 21) incrementally using R cutree. The first k-value  
610 was chosen that allowed for the target strain to cluster with at least two other strains that were  
611 members of the same clade as defined by the current phylogenetic relationship among the  
612 sequenced strains. These criteria effectively eliminate Clade E because it only contains one  
613 strain. Additionally, this approach was assessed for congruence manually across candidate  
614 regions.

### 615 **Sanger sequencing**

616 The association between SNPs and indels identified in the NGS data was tested by  
617 amplification of specific regions that were either shared among a number of strains (P75063) or  
618 strain-specific (P60002). Two regions were PCR amplified using primers  
619 AGTCGGTGATGTCTATAGTG / GCTGTCCTTGGATCATTGAT to amplify Chr7:48017..48664  
620 in P75063 and TTCTGCTGTTGCTGCTGCTA / CTGTCAACTGTCAACCAAAG to amplify  
621 ChrR: 19979596..1998145 in P60002. Amplicons were purified and Sanger sequenced.

622 Mitochondrial SNPs were verified using 3 sets of primers to amplify different regions of  
623 ChrM across the 21 natural isolates. Two isolates were analyzed from each clade including  
624 P60002, P75010, and P94015 strains. Primers TTAGTAGTGTCGGTGTCTTC /  
625 AGAGAGGGTTTTGGTTAGGG were used to PCR amplify ChrM: 4899..6076,  
626 GAATCTCAGAGACTACACGT / GTGGTATACGACGAGGCATT were used for ChrM:  
627 18265..20660, and TGGGAAGTAGAGGCTGAAGA / AGGGGCATTATAAGGAGGAG were  
628 used for ChrM: 28094..29548. PCR amplicons were purified and Sanger sequenced.

### 629 **Statistical Testing**

630 Statistics were performed as Student's T-test unless otherwise indicated. All statistical  
631 tests were performed in the R 3.2.5 programming environment [72].

632

### 633 **ACKNOWLEDGEMENTS**

634           We thank Iuliana Ene and Christophe D'Enfert for comments on the manuscript,  
635   Christina Cuomo for technical advice, and members of the Bennett and Anderson labs for  
636   helpful discussions. This work was supported by National Institutes of Health grants AI081704  
637   and AI122011 (to R.J.B.), a PATH award from the Burroughs Wellcome Fund (to R.J.B.), and by  
638   a Karen T. Romer Undergraduate Teaching and Research Award (to J.M.W.).

639

640

## 641 References

- 642 1. Bennett RJ, Forche A, Berman J. Rapid mechanisms for generating genome diversity: whole  
643 ploidy shifts, aneuploidy, and loss of heterozygosity. *Cold Spring Harb Perspect Med.* 2014;4(10). Epub  
644 2014/08/02. doi: 10.1101/cshperspect.a019604. PubMed PMID: 25081629; PubMed Central PMCID:  
645 PMCPMC4200206.
- 646 2. Dunham MJ, Badrane H, Ferea T, Adams J, Brown PO, Rosenzweig F, et al. Characteristic genome  
647 rearrangements in experimental evolution of *Saccharomyces cerevisiae*. *Proc Natl Acad Sci U S A.*  
648 2002;99(25):16144-9. doi: 10.1073/pnas.242624799. PubMed PMID: 12446845; PubMed Central PMCID:  
649 PMCPMC138579.
- 650 3. Goddard MR. Molecular evolution: Sex accelerates adaptation. *Nature.* 2016;531(7593):176-7.  
651 doi: 10.1038/nature17304. PubMed PMID: 26909572.
- 652 4. McDonald MJ, Rice DP, Desai MM. Sex speeds adaptation by altering the dynamics of molecular  
653 evolution. *Nature.* 2016;531(7593):233-6. doi: 10.1038/nature17143. PubMed PMID: 26909573;  
654 PubMed Central PMCID: PMCPMC4855304.
- 655 5. Lehtonen J, Jennions MD, Kokko H. The many costs of sex. *Trends Ecol Evol.* 2012;27(3):172-8.  
656 Epub 2011/10/25. doi: 10.1016/j.tree.2011.09.016. PubMed PMID: 22019414.
- 657 6. Roze D, Otto SP. Differential selection between the sexes and selection for sex. *Evolution.*  
658 2012;66(2):558-74. doi: 10.1111/j.1558-5646.2011.01459.x. PubMed PMID: 22276548.
- 659 7. Grimberg B, Zeyl C. The effects of sex and mutation rate on adaptation in test tubes and to  
660 mouse hosts by *Saccharomyces cerevisiae*. *Evolution.* 2005;59(2):431-8. PubMed PMID: 15807427.
- 661 8. Nielsen K, Heitman J. Sex and virulence of human pathogenic fungi. *Adv Genet.* 2007;57:143-73.  
662 doi: 10.1016/S0065-2660(06)57004-X. PubMed PMID: 17352904.
- 663 9. Becks L, Agrawal AF. Higher rates of sex evolve in spatially heterogeneous environments.  
664 *Nature.* 2010;468(7320):89-92. doi: 10.1038/nature09449. PubMed PMID: 20944628.
- 665 10. Horn DL, Neofytos D, Anaissie EJ, Fishman JA, Steinbach WJ, Olyaei AJ, et al. Epidemiology and  
666 outcomes of candidemia in 2019 patients: data from the prospective antifungal therapy alliance registry.  
667 *Clin Infect Dis.* 2009;48(12):1695-703. doi: 10.1086/599039. PubMed PMID: 19441981.
- 668 11. Sardi JC, Scorzoni L, Bernardi T, Fusco-Almeida AM, Mendes Giannini MJ. *Candida* species:  
669 current epidemiology, pathogenicity, biofilm formation, natural antifungal products and new  
670 therapeutic options. *J Med Microbiol.* 2013;62(Pt 1):10-24. doi: 10.1099/jmm.0.045054-0. PubMed  
671 PMID: 23180477.
- 672 12. Taylor JW, Berbee ML. Dating divergences in the Fungal Tree of Life: review and new analyses.  
673 *Mycologia.* 2006;98(6):838-49. PubMed PMID: 17486961.
- 674 13. Bennett RJ, Johnson AD. Completion of a parasexual cycle in *Candida albicans* by induced  
675 chromosome loss in tetraploid strains. *EMBO J.* 2003;22(10):2505-15. Epub 2003/05/14. doi:  
676 10.1093/emboj/cdg235. PubMed PMID: 12743044; PubMed Central PMCID: PMCPMC155993.
- 677 14. Forche A, Alby K, Schaefer D, Johnson AD, Berman J, Bennett RJ. The parasexual cycle in *Candida*  
678 *albicans* provides an alternative pathway to meiosis for the formation of recombinant strains. *PLoS Biol.*  
679 2008;6(5):e110. Epub 2008/05/09. doi: 10.1371/journal.pbio.0060110. PubMed PMID: 18462019;  
680 PubMed Central PMCID: PMCPMC2365976.
- 681 15. Magee BB, Magee PT. Induction of mating in *Candida albicans* by construction of MTL $\alpha$  and  
682 MTL $\alpha$  strains. *Science.* 2000;289(5477):310-3. Epub 2000/07/15. PubMed PMID: 10894781.
- 683 16. Hull CM, Raisner RM, Johnson AD. Evidence for mating of the "asexual" yeast *Candida albicans* in  
684 a mammalian host. *Science.* 2000;289(5477):307-10. Epub 2000/07/15. PubMed PMID: 10894780.
- 685 17. Miller MG, Johnson AD. White-opaque switching in *Candida albicans* is controlled by mating-  
686 type locus homeodomain proteins and allows efficient mating. *Cell.* 2002;110(3):293-302. Epub  
687 2002/08/15. PubMed PMID: 12176317.

- 688 18. Alby K, Schaefer D, Bennett RJ. Homothallic and heterothallic mating in the opportunistic  
689 pathogen *Candida albicans*. *Nature*. 2009;460(7257):890-3. doi: 10.1038/nature08252. PubMed PMID:  
690 19675652; PubMed Central PMCID: PMCPMC2866515.
- 691 19. Keeney S. Spo11 and the Formation of DNA Double-Strand Breaks in Meiosis. *Genome Dyn Stab*.  
692 2008;2:81-123. doi: 10.1007/7050\_2007\_026. PubMed PMID: 21927624; PubMed Central PMCID:  
693 PMCPMC3172816.
- 694 20. Tavanti A, Davidson AD, Fordyce MJ, Gow NA, Maiden MC, Odds FC. Population structure and  
695 properties of *Candida albicans*, as determined by multilocus sequence typing. *J Clin Microbiol*.  
696 2005;43(11):5601-13. doi: 10.1128/JCM.43.11.5601-5613.2005. PubMed PMID: 16272493; PubMed  
697 Central PMCID: PMCPMC1287804.
- 698 21. Nebavi F, Ayala FJ, Renaud F, Bertout S, Eholie S, Moussa K, et al. Clonal population structure  
699 and genetic diversity of *Candida albicans* in AIDS patients from Abidjan (Cote d'Ivoire). *Proc Natl Acad Sci*  
700 *U S A*. 2006;103(10):3663-8. doi: 10.1073/pnas.0511328103. PubMed PMID: 16501044; PubMed Central  
701 PMCID: PMCPMC1450139.
- 702 22. Odds FC, Bougnoux ME, Shaw DJ, Bain JM, Davidson AD, Diogo D, et al. Molecular phylogenetics  
703 of *Candida albicans*. *Eukaryot Cell*. 2007;6(6):1041-52. doi: 10.1128/EC.00041-07. PubMed PMID:  
704 17416899; PubMed Central PMCID: PMCPMC1951527.
- 705 23. Anderson JB, Wickens C, Khan M, Cowen LE, Federspiel N, Jones T, et al. Infrequent genetic  
706 exchange and recombination in the mitochondrial genome of *Candida albicans*. *J Bacteriol*.  
707 2001;183(3):865-72. doi: 10.1128/JB.183.3.865-872.2001. PubMed PMID: 11208783; PubMed Central  
708 PMCID: PMCPMC94952.
- 709 24. Jacobsen MD, Rattray AM, Gow NA, Odds FC, Shaw DJ. Mitochondrial haplotypes and  
710 recombination in *Candida albicans*. *Med Mycol*. 2008;46(7):647-54. doi: 10.1080/13693780801986631.  
711 PubMed PMID: 18608923.
- 712 25. Zhang L, Yan L, Jiang J, Wang Y, Jiang Y, Yan T, et al. The structure and retrotransposition  
713 mechanism of LTR-retrotransposons in the asexual yeast *Candida albicans*. *Virulence*. 2014;5(6):655-64.  
714 doi: 10.4161/viru.32180. PubMed PMID: 25101670; PubMed Central PMCID: PMCPMC4139406.
- 715 26. Diogo D, Bouchier C, d'Enfert C, Bougnoux ME. Loss of heterozygosity in commensal isolates of  
716 the asexual diploid yeast *Candida albicans*. *Fungal Genet Biol*. 2009;46(2):159-68. doi:  
717 10.1016/j.fgb.2008.11.005. PubMed PMID: 19059493.
- 718 27. Bougnoux ME, Pujol C, Diogo D, Bouchier C, Soll DR, d'Enfert C. Mating is rare within as well as  
719 between clades of the human pathogen *Candida albicans*. *Fungal Genet Biol*. 2008;45(3):221-31. doi:  
720 10.1016/j.fgb.2007.10.008. PubMed PMID: 18063395; PubMed Central PMCID: PMCPMC2275664.
- 721 28. Hirakawa MP, Martinez DA, Sakthikumar S, Anderson MZ, Berlin A, Gujja S, et al. Genetic and  
722 phenotypic intra-species variation in *Candida albicans*. *Genome Res*. 2015;25(3):413-25. doi:  
723 10.1101/gr.174623.114. PubMed PMID: 25504520; PubMed Central PMCID: PMCPMC4352881.
- 724 29. Wu W, Lockhart SR, Pujol C, Srikantha T, Soll DR. Heterozygosity of genes on the sex  
725 chromosome regulates *Candida albicans* virulence. *Mol Microbiol*. 2007;64(6):1587-604. doi:  
726 10.1111/j.1365-2958.2007.05759.x. PubMed PMID: 17555440.
- 727 30. van het Hoog M, Rast TJ, Martchenko M, Grindle S, Dignard D, Hogues H, et al. Assembly of the  
728 *Candida albicans* genome into sixteen supercontigs aligned on the eight chromosomes. *Genome Biol*.  
729 2007;8(4):R52. doi: 10.1186/gb-2007-8-4-r52. PubMed PMID: 17419877; PubMed Central PMCID:  
730 PMCPMC1896002.
- 731 31. Butler G, Rasmussen MD, Lin MF, Santos MA, Sakthikumar S, Munro CA, et al. Evolution of  
732 pathogenicity and sexual reproduction in eight *Candida* genomes. *Nature*. 2009;459(7247):657-62. Epub  
733 2009/05/26. doi: 10.1038/nature08064. PubMed PMID: 19465905; PubMed Central PMCID:  
734 PMCPMC2834264.

- 735 32. Muzzey D, Schwartz K, Weissman JS, Sherlock G. Assembly of a phased diploid *Candida albicans*  
736 genome facilitates allele-specific measurements and provides a simple model for repeat and indel  
737 structure. *Genome Biol.* 2013;14(9):R97. doi: 10.1186/gb-2013-14-9-r97. PubMed PMID: 24025428;  
738 PubMed Central PMCID: PMC4054093.
- 739 33. Rosenberg MS, Subramanian S, Kumar S. Patterns of transitional mutation biases within and  
740 among mammalian genomes. *Mol Biol Evol.* 2003;20(6):988-93. doi: 10.1093/molbev/msg113. PubMed  
741 PMID: 12716982.
- 742 34. McDonald MJ, Wang WC, Huang HD, Leu JY. Clusters of nucleotide substitutions and  
743 insertion/deletion mutations are associated with repeat sequences. *PLoS Biol.* 2011;9(6):e1000622. doi:  
744 10.1371/journal.pbio.1000622. PubMed PMID: 21697975; PubMed Central PMCID: PMC3114760.
- 745 35. Jovelin R, Cutter AD. Fine-scale signatures of molecular evolution reconcile models of indel-  
746 associated mutation. *Genome Biol Evol.* 2013;5(5):978-86. doi: 10.1093/gbe/evt051. PubMed PMID:  
747 23558593; PubMed Central PMCID: PMC3673634.
- 748 36. Tian D, Wang Q, Zhang P, Araki H, Yang S, Kreitman M, et al. Single-nucleotide mutation rate  
749 increases close to insertions/deletions in eukaryotes. *Nature.* 2008;455(7209):105-8. doi:  
750 10.1038/nature07175. PubMed PMID: 18641631.
- 751 37. Nikpay M, Goel A, Won HH, Hall LM, Willenborg C, Kanoni S, et al. A comprehensive 1,000  
752 Genomes-based genome-wide association meta-analysis of coronary artery disease. *Nat Genet.*  
753 2015;47(10):1121-30. doi: 10.1038/ng.3396. PubMed PMID: 26343387; PubMed Central PMCID:  
754 PMC4589895.
- 755 38. Arbeithuber B, Betancourt AJ, Ebner T, Tiemann-Boege I. Crossovers are associated with  
756 mutation and biased gene conversion at recombination hotspots. *Proc Natl Acad Sci U S A.*  
757 2015;112(7):2109-14. doi: 10.1073/pnas.1416622112. PubMed PMID: 25646453; PubMed Central  
758 PMCID: PMC4343121.
- 759 39. Deem A, Keszthelyi A, Blackgrove T, Vayl A, Coffey B, Mathur R, et al. Break-induced replication  
760 is highly inaccurate. *PLoS Biol.* 2011;9(2):e1000594. doi: 10.1371/journal.pbio.1000594. PubMed PMID:  
761 21347245; PubMed Central PMCID: PMC3039667.
- 762 40. Strathern JN, Shafer BK, McGill CB. DNA synthesis errors associated with double-strand-break  
763 repair. *Genetics.* 1995;140(3):965-72. PubMed PMID: 7672595; PubMed Central PMCID:  
764 PMC1206680.
- 765 41. Carbone I, Kohn L. Inferring process from pattern in fungal population genetics. *Applied*  
766 *mycology and biotechnology.* 2004;4(1):30.
- 767 42. Harvey PH, Pagel MD. *The comparative method in evolutionary biology*: Oxford university press  
768 Oxford; 1991.
- 769 43. Chen XJ, Butow RA. The organization and inheritance of the mitochondrial genome. *Nat Rev*  
770 *Genet.* 2005;6(11):815-25. doi: 10.1038/nrg1708. PubMed PMID: 16304597.
- 771 44. Hodgkinson A, Eyre-Walker A. Variation in the mutation rate across mammalian genomes. *Nat*  
772 *Rev Genet.* 2011;12(11):756-66. doi: 10.1038/nrg3098. PubMed PMID: 21969038.
- 773 45. Yue JX, Li J, Aigrain L, Hallin J, Persson K, Oliver K, et al. Contrasting evolutionary genome  
774 dynamics between domesticated and wild yeasts. *Nat Genet.* 2017;49(6):913-24. doi: 10.1038/ng.3847.  
775 PubMed PMID: 28416820; PubMed Central PMCID: PMC5446901.
- 776 46. Xie Z, Wang L, Wang L, Wang Z, Lu Z, Tian D, et al. Mutation rate analysis via parent-progeny  
777 sequencing of the perennial peach. I. A low rate in woody perennials and a higher mutagenicity in  
778 hybrids. *Proc Biol Sci.* 2016;283(1841). doi: 10.1098/rspb.2016.1016. PubMed PMID: 27798292; PubMed  
779 Central PMCID: PMC5095371.
- 780 47. Yang S, Wang L, Huang J, Zhang X, Yuan Y, Chen JQ, et al. Parent-progeny sequencing indicates  
781 higher mutation rates in heterozygotes. *Nature.* 2015;523(7561):463-7. doi: 10.1038/nature14649.  
782 PubMed PMID: 26176923.

- 783 48. Abyzov A, Li S, Kim DR, Mohiyuddin M, Stutz AM, Parrish NF, et al. Analysis of deletion  
784 breakpoints from 1,092 humans reveals details of mutation mechanisms. *Nat Commun.* 2015;6:7256.  
785 doi: 10.1038/ncomms8256. PubMed PMID: 26028266; PubMed Central PMCID: PMC4451611.
- 786 49. Longman-Jacobsen N, Williamson JF, Dawkins RL, Gaudieri S. In polymorphic genomic regions  
787 indels cluster with nucleotide polymorphism: Quantum Genomics. *Gene.* 2003;312:257-61. PubMed  
788 PMID: 12909362.
- 789 50. Brown TC, Jiricny J. A specific mismatch repair event protects mammalian cells from loss of 5-  
790 methylcytosine. *Cell.* 1987;50(6):945-50. PubMed PMID: 3040266.
- 791 51. Wang L, Zhang Y, Qin C, Tian D, Yang S, Hurst LD. Mutation rate analysis via parent-progeny  
792 sequencing of the perennial peach. II. No evidence for recombination-associated mutation. *Proc Biol Sci.*  
793 2016;283(1841). doi: 10.1098/rspb.2016.1785. PubMed PMID: 27798307; PubMed Central PMCID:  
794 PMC45095386.
- 795 52. Seplyarskiy VB, Logacheva MD, Penin AA, Baranova MA, Leushkin EV, Demidenko NV, et al.  
796 Crossing-over in a hypervariable species preferentially occurs in regions of high local similarity. *Mol Biol*  
797 *Evol.* 2014;31(11):3016-25. doi: 10.1093/molbev/msu242. PubMed PMID: 25135947; PubMed Central  
798 PMCID: PMC4209137.
- 799 53. Bryan Thornclaw JH, Jackie Roger, Henry Gong, Todd Lowe, Russell Corbett-Detig. Transfer RNA  
800 genes experience exceptionally elevated mutation rates. 2017. Epub December 6th, 2017.
- 801 54. Helmrich A, Ballarino M, Tora L. Collisions between replication and transcription complexes  
802 cause common fragile site instability at the longest human genes. *Mol Cell.* 2011;44(6):966-77. Epub  
803 2011/12/27. doi: 10.1016/j.molcel.2011.10.013. PubMed PMID: 22195969.
- 804 55. Mattick JS. The genetic signatures of noncoding RNAs. *PLoS Genet.* 2009;5(4):e1000459. Epub  
805 2009/04/25. doi: 10.1371/journal.pgen.1000459. PubMed PMID: 19390609; PubMed Central PMCID:  
806 PMC2667263.
- 807 56. Selmecki A, Forche A, Berman J. Genomic plasticity of the human fungal pathogen *Candida*  
808 *albicans*. *Eukaryot Cell.* 2010;9(7):991-1008. doi: 10.1128/EC.00060-10. PubMed PMID: 20495058;  
809 PubMed Central PMCID: PMC2901674.
- 810 57. Forche A, Abbey D, Pisithkul T, Weinzierl MA, Ringstrom T, Bruck D, et al. Stress alters rates and  
811 types of loss of heterozygosity in *Candida albicans*. *MBio.* 2011;2(4). doi: 10.1128/mBio.00129-11.  
812 PubMed PMID: 21791579; PubMed Central PMCID: PMC3143845.
- 813 58. Feri A, Loll-Kripplleber R, Commere PH, Maufrais C, Sertour N, Schwartz K, et al. Analysis of  
814 Repair Mechanisms following an Induced Double-Strand Break Uncovers Recessive Deleterious Alleles in  
815 the *Candida albicans* Diploid Genome. *MBio.* 2016;7(5). doi: 10.1128/mBio.01109-16. PubMed PMID:  
816 27729506; PubMed Central PMCID: PMC45061868.
- 817 59. Ford CB, Funt JM, Abbey D, Issi L, Guiducci C, Martinez DA, et al. The evolution of drug  
818 resistance in clinical isolates of *Candida albicans*. *Elife.* 2015;4:e00662. doi: 10.7554/eLife.00662.  
819 PubMed PMID: 25646566; PubMed Central PMCID: PMC4383195.
- 820 60. Hickman MA, Zeng G, Forche A, Hirakawa MP, Abbey D, Harrison BD, et al. The 'obligate diploid'  
821 *Candida albicans* forms mating-competent haploids. *Nature.* 2013;494(7435):55-9. Epub 2013/02/01.  
822 doi: 10.1038/nature11865. PubMed PMID: 23364695; PubMed Central PMCID: PMC3583542.
- 823 61. Muller HJ. The Relation of Recombination to Mutational Advance. *Mutat Res.* 1964;106:2-9.  
824 PubMed PMID: 14195748.
- 825 62. Tucker AE, Ackerman MS, Eads BD, Xu S, Lynch M. Population-genomic insights into the  
826 evolutionary origin and fate of obligately asexual *Daphnia pulex*. *Proc Natl Acad Sci U S A.*  
827 2013;110(39):15740-5. doi: 10.1073/pnas.1313388110. PubMed PMID: 23959868; PubMed Central  
828 PMCID: PMC3785735.
- 829 63. Pujol C, Reynes J, Renaud F, Raymond M, Tibayrenc M, Ayala FJ, et al. The yeast *Candida albicans*  
830 has a clonal mode of reproduction in a population of infected human immunodeficiency virus-positive



- 831 patients. Proc Natl Acad Sci U S A. 1993;90(20):9456-9. PubMed PMID: 8415722; PubMed Central  
832 PMCID: PMC47587.
- 833 64. Schmid J, Cannon, R.D., Holland, B. A futile act? Thoughts on the reproductive biology of  
834 *Candida albicans*. Mycologist. 2004;18(4):158-63.
- 835 65. Odds FC. Molecular phylogenetics and epidemiology of *Candida albicans*. Future Microbiol.  
836 2010;5(1):67-79. doi: 10.2217/fmb.09.113. PubMed PMID: 20020830.
- 837 66. Zhang N, Magee BB, Magee PT, Holland BR, Rodrigues E, Holmes AR, et al. Selective Advantages  
838 of a Parasexual Cycle for the Yeast *Candida albicans*. Genetics. 2015;200(4):1117-32. doi:  
839 10.1534/genetics.115.177170. PubMed PMID: 26063661; PubMed Central PMCID: PMC4574235.
- 840 67. Liti G, Carter DM, Moses AM, Warringer J, Parts L, James SA, et al. Population genomics of  
841 domestic and wild yeasts. Nature. 2009;458(7236):337-41. doi: 10.1038/nature07743. PubMed PMID:  
842 19212322; PubMed Central PMCID: PMC2659681.
- 843 68. Li H, Durbin R. Fast and accurate long-read alignment with Burrows-Wheeler transform.  
844 Bioinformatics. 2010;26(5):589-95. doi: 10.1093/bioinformatics/btp698. PubMed PMID: 20080505;  
845 PubMed Central PMCID: PMC2828108.
- 846 69. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence Alignment/Map  
847 format and SAMtools. Bioinformatics. 2009;25(16):2078-9. doi: 10.1093/bioinformatics/btp352.  
848 PubMed PMID: 19505943; PubMed Central PMCID: PMC2723002.
- 849 70. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al. The Genome  
850 Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. Genome  
851 Res. 2010;20(9):1297-303. doi: 10.1101/gr.107524.110. PubMed PMID: 20644199; PubMed Central  
852 PMCID: PMC2928508.
- 853 71. Robinson JT, Thorvaldsdottir H, Winckler W, Guttman M, Lander ES, Getz G, et al. Integrative  
854 genomics viewer. Nat Biotechnol. 2011;29(1):24-6. doi: 10.1038/nbt.1754. PubMed PMID: 21221095;  
855 PubMed Central PMCID: PMC3346182.
- 856 72. Team RC. R: A language and environment for statistical computing. . R Foudation for Statistical  
857 Computing. 2016.

858

859

860 **FIGURE LEGENDS**

861 **Figure 1. Distribution of polymorphisms among 21 clinical isolates of *C. albicans*.**

862 **A.** Number of heterozygous (purple) and homozygous (orange) strain-specific SNPs  
863 and insertion/deletions (indels) are plotted for each isolate. Clade designations for each  
864 isolate are color coded. **B.** There are two types of sequence variants among the set of  
865 21 isolates; single nucleotide polymorphisms (SNPs) shared by multiple isolates and  
866 SNPs that are specific to individual strains. Variants encoded by multiple strains  
867 suggest origins in a common ancestor, whereas strain-specific polymorphisms likely  
868 arose specifically in individual strain backgrounds. **C.** Relative frequency of strain-  
869 specific SNPs (blue) and strain-specific indels (orange) were plotted across the genome  
870 using 5 kb sliding windows. **D.** Number of strain-specific SNPs within 100 bp of each  
871 strain-specific indel was plotted. The average number of strain-specific SNPs in an  
872 equal number of random 100 bp windows bootstrapped 1000 times is shown (red line).  
873 **E.** Distance to the nearest strain-specific SNP for each strain-specific indel was  
874 measured in a 100 bp window (with non-overlapping 10 bp intervals) surrounding the  
875 indels.

876

877 **Figure 2. Strain-specific SNPs are enriched at intergenic positions and snoRNAs.**

878 The ratio of genic to intergenic strain-specific SNPs (**A**) and indels (**B**) was calculated  
879 for each of the sequenced isolates and falls below the fraction of the genome that  
880 encodes protein-coding genes (red line). **C.** The density of strain-specific SNPs from all  
881 sequenced isolates was measured for all genes in the *C. albicans* genome and plotted  
882 against the average SNP density for all isolates (red line). **D.** The placement of the

883 SNPs was determined by breaking snoRNAs into five equal segments. SNPs were  
884 significantly enriched in the 5' end of the RNAs. \* denotes  $p < 0.01$ . **E.** Genes enriched  
885 for strain-specific SNPs were identified and functional enrichment was assessed using  
886 GO term analysis. snoRNAs were significantly enriched for strain-specific SNPs.

887

888 **Figure 3. Loss of heterozygosity (LOH) events influence the *C. albicans* genome**

889 **landscape.** **A.** Number of strains that show LOH for 50 kb windows across the *C.*  
890 *albicans* genome aligned to their chromosomal position. **B.** Number of chromosomes  
891 that are homozygous for 5 kb windows within 150 kb of the centromere. **C.** SNP  
892 frequency for each of the 445 heterozygous (het) and 468 homozygous (hom) regions  
893 across the 21 sequenced genomes was plotted. Normalized SNP density is calculated  
894 using the number of SNPs within a het or hom region divided by the region length, and  
895 is significantly elevated for het regions (red) relative to hom regions (blue). **D.** One kb  
896 segments flanking the 745 LOH transition points were separated into their respective  
897 het and hom segments and assessed for SNP density. The genome average was  
898 calculated by randomly selecting an equal number of 2 kb windows across the genome  
899 and bootstrapping for 1000 iterations.

900

901 **Figure 4. Mutational patterns following identity by descent and loss of**

902 **heterozygosity are widespread.** **A.** Two patterns of polymorphisms exist within  
903 sequenced genomes. One shows polymorphisms are phylogenetically congruent  
904 indicative of identity by descent, and the other shows polymorphisms that violate the  
905 phylogenetic relationship implicating mechanisms other than direct inheritance.

906 Shading denotes similarity in overall SNP patterns. Heterozygous SNPs are purple and  
907 homozygous SNPs are orange. Polymorphisms for the set of *C. albicans* genomes are  
908 plotted for two loci (**B**, **C**) that display SNP patterns consistent with inheritance by  
909 descent. **D**. The polymorphism pattern is plotted for a locus that does not follow  
910 inheritance by descent. Similar genotypes are color-coded and connected to each  
911 other. **E**. A cartoon depicting LOH of heterozygotes in opposing directions provides the  
912 most parsimonious explanation for the observed SNP patterns.

913

914 **Figure 5. Evidence for recombination in *C. albicans* isolates.** **A**. All SNPs are  
915 shown for a 20 kb region of Chr4 for the 21 sequenced genomes. For each strain, dark  
916 grey bars are heterozygous genomic regions while lighter bars indicate regions that are  
917 mostly homozygous. The SNP pattern in P94015 indicates one region with homology to  
918 Clade I (P94015-A) next to a region without clear homology to any specific clade  
919 (P94015-B) followed by a region with homology to Clade SA (P94015-C). **B**. SNPs for a  
920 2.5 kb region of Chr3 were phased to individual homologs for Clade SA and Clade III by  
921 using LOH of a closely-related strain for reference. One homolog from each clade  
922 matched the exact SNP pattern in a 'hybrid' region present in P94015 (both the position  
923 of the SNP and the actual base substitution matched between P94015 and the Clade  
924 SA or Clade III homolog). **C**. Consensus SNP patterns for each clade were used to  
925 assess genome similarities between all isolates in 25 kb sliding windows. The closest  
926 match for each window was color-coded by clade. The SNP patterns for two strains,  
927 P60002 and P94015, contained regions assigned to multiple clades. dark grey-SA,  
928 blue-Clade III, mustard-Clade II, red-Clade I, light grey-no clade consensus.

929

930 **Figure 6. Mitochondrial genomes in *C. albicans* display recombinant genotypes.**

931 The mitochondrial (mt) genome sequences of 21 clinical isolates of *C. albicans* were  
932 compared. The positions of SNPs that differ from the SC5314 assembly are shown  
933 excluding ChrM:35000-41000 due to the absence of any SNPs in this region. The mt  
934 genomes for P94015 (**A**), P60002 (**B**), and P75010 (**C**) are highlighted to show the  
935 mosaic configuration of SNPs relative to other clades. The mt genome in P94015  
936 contains regions that align with those of Clade II and Clade SA (clade-specific pattern  
937 marked with an asterisk with key positions shaded), whereas the P60002 mt genome  
938 aligns with sequences for Clades I and II. A 6 kb region devoid of SNPs is more lightly  
939 shaded. The three clade SA-specific SNPs in the region that demonstrate alignment of  
940 P94015-A are colored black to aid visual alignment. **D.** The similarity of mt regions from  
941 each isolate was compared by analyzing 2 kb windows from each strain relative to  
942 consensus SNP patterns for each clade. The window was color-coded to designate the  
943 clade with greatest similarity. dark grey-SA, blue-Clade III, mustard-Clade II, red-Clade  
944 I, light grey-no clade consensus. **E.** Two strains from Clade SA and Clade I along with  
945 P60002 and P94015 were Sanger sequenced across three separate 1 kb regions of  
946 their mitochondrial genomes. Chromatograms highlighting variant positions consistent  
947 with recombination between clades producing the SNP patterns present in P60002 and  
948 P94015 are shown along each chromosome.

949

950

## 951 **Supporting Information Legends**

952 **Figure S1. The distribution of polymorphisms is non-random.** The number of  
953 strains encoding each SNP (**A**) or indel (**B**) was determined and the frequency plotted.  
954 A best-fit line (red) was plotted for each distribution.

955 **Figure S2. LOH can obfuscate inheritance by descent patterns of shared**  
956 **polymorphisms.** LOH of opposing alleles in a common ancestor can make it appear  
957 that mutations arose independently despite mutations sharing a common origin. LOH of  
958 homolog A in Clade III produces a different SNP pattern than that of Clade II although  
959 both arose from the same ancestral strain.

960 **Figure S3. Number of polymorphisms correlates with branch length.** The branch  
961 length on the phylogenetic tree to the nearest node was correlated against the number  
962 strain-specific SNPs (**A**) or strain-specific indels (**B**).

963 **Figure S4. Transitions are present more frequently than transversions during**  
964 **strain evolution.** The percentage of SNPs that result in transitions and transversions  
965 was calculated both for all SNPs and for strain-specific SNPs.

966 **Figure S5. Characterization of indels across *C. albicans* isolates.** The number of  
967 strain-specific indels was plotted for either intergenic (**A**) or genic (**B**) mutations based  
968 on the indel size, ranging from 1-10 nucleotides. Blue indicates deletions and yellow  
969 indicates insertions.

970 **Figure S6. Verification of SNP-indel association by Sanger sequencing.** Four  
971 regions that contained SNPs tightly linked to indels were chosen to be assessed by  
972 Sanger sequencing. The genomic DNA from Chr7 in strain P76055 (**A**), Chr5 in GC75

973 (B), Chr7 in P87 (C), and Chr2 in P57072 were Sanger sequenced and encoded the  
974 putative SNPs and indels (colored boxes) as expected from whole genome sequencing.  
975 The reference SC5314 sequence is shown for comparison to the Sanger sequenced  
976 DNA below along with the chromatogram aligned to the polymorphisms identified from  
977 whole genome sequencing. Below each schematic are listed the informative sites.

978 **Figure S7. Associations of strain-specific variants corrected for clade bias.**

979 Strains-specific mutations from 12 strains (Clade I: 12C, L26, P78048; Clade II: P57072,  
980 P76055, P76067; Clade III: P34048, P78042, P57055; Clade SA: P87, GC75, P75063)  
981 were retained to ensure that clade representation did not bias these associations. A.  
982 The relative frequency of strain-specific SNPs (blue) and indels (orange) was plotted  
983 across the genome using 5000 bp sliding windows. Distance to the nearest SNP for  
984 each indel was measured in a 100 bp window (with non-overlapping 5 bp intervals)  
985 surrounding the indels for either strain specific (B) or all (C) variants.

986 **Figure S8. Mapping of heterozygous and homozygous regions of *C. albicans***

987 **genomes.** Schematic showing heterozygous (red) and homozygous (blue) regions of  
988 sequenced *C. albicans* isolates. Note that most long-tract LOH events (>50 kb) start  
989 between the centromere and the telomere and extend to the ends of the chromosomes.  
990 Chromosomes are displayed along the bottom with green circles indicating  
991 centromeres, blue boxes denoting major repeat sequence (MRS) loci, and the orange  
992 box signifying the *MTL* locus.

993 **Figure S9. *C. albicans* SNP patterns generally follow inheritance by descent.**

994 Polymorphisms for the set of *C. albicans* genomes are plotted for two loci that display

995 SNP patterns consistence with inheritance by descent. Heterozygous and homozygous  
996 SNPs are purple and orange, respectively.

997 **Figure S10. LOH alters inheritance patterns of *C. albicans* polymorphisms.**

998 Polymorphism patterns are plotted for two loci (**A, C**) that do not follow inheritance by  
999 descent. Similar genotypes are color-coded and connected to each other. Cartoons  
1000 depicting LOH of heterozygotes in opposing directions to produce the observed SNP  
1001 pattern (**B, D**) provide the most parsimonious explanation for those two loci.  
1002 Heterozygous and homozygous SNPs are purple and orange, respectively.

1003 **Figure S11. Modes of inheritance in *C. albicans*.** Both LOH and mating can impact

1004 SNP patterns in *C. albicans* strains. Analysis of the distribution of heterozygous and  
1005 homozygous SNPs can help differentiate between these two possibilities'. Following  
1006 LOH, all affected SNPs are homozygous and derived strains will therefore contain  
1007 homozygous variants of homolog A or homolog B. LOH can therefore make the  
1008 precursor to these LOH events (Strain B in the figure) appear 'recombinant', as it will  
1009 contain heterozygous SNPs at positions that are homozygous in the lineages that  
1010 underwent LOH. In contrast, strains that are related due to recombination may share  
1011 only heterozygous SNPs between lineages or may share a mix of heterozygous and  
1012 homozygous SNPs, with homozygous positions due to inheritance of the same SNP  
1013 from both parents (see also Figure 3B). These patterns may be further complicated by  
1014 additional short-tract LOH events, as indicated on the right of the figure.

1015 **Figure S12. Evidence for recombination in two *C. albicans* strains.** The SNP  
1016 patterns for two regions highlight recombinant genotypes in strains P60002 (**A**) and  
1017 P94015 (**B**). The DNA segments corresponding to different clades are aligned next to



1018 the appropriate group and labeled according to homologous tracts. Heterozygous and  
1019 homozygous SNPs are purple and orange, respectively. **C.** Consensus SNP patterns  
1020 for each clade were used to assess genome similarities between all isolates in 50 kb  
1021 sliding windows. In this case, all strains from the same clade were removed to force the  
1022 closest match for each window to be assigned by color-coding to the most similar clade.  
1023 The SNP patterns for two strains, P60002 and P94015, stand out with respect to their  
1024 position in the constructed phylogenetic tree and assigned clade. dark grey-SA, blue-  
1025 Clade III, mustard-Clade II, red-Clade I, light grey-no clade consensus.

1026 **Figure S13. Alternative mode of recombination in the P60002 mtDNA genome.**

1027 The mitochondrial genome of P60002 may be the result of recombination between  
1028 clade SA (P60002-B,E), clade I (P60002-A,C,F), and clade II (P60002-D,G), although  
1029 we note that the resolution of SNP patterns cannot distinguish all clades (e.g., regions  
1030 B, E have identical SNP patterns in both clade SA and clade I).

1031 **Figure S14. Verification of mitochondrial SNPs by Sanger sequencing.** Two

1032 strains from Clades I, II, III and SA, as well as P60002, P75010, and P94015 were  
1033 sequenced across three separate 1 kb regions of the mitochondrial genomes to assess  
1034 variant calls from whole genome sequencing. Chromatograms highlighting variant  
1035 positions consistent with recombination between clades producing the SNP patterns  
1036 present in P60002 and P94015 are shown along each chromosome.

1037 **Figure S15. Read depth of variant calls define heterozygous and homozygous**

1038 **positions.** The proportion of reads encoding a variant at each position was calculated  
1039 relative to the total number of mapped reads at that position. Plotting each position for  
1040 all strains produced a distribution that was separated at 90% of reads encoding a

1041 variant allele. The distribution of variant position <90% were defined as heterozygous  
1042 and those >90% as homozygous.

1043 **Figure S16. Positioning of SNPs and indels across *C. albicans* genomes. A.** The  
1044 relative frequency of all SNPs (blue) and indels (orange) was plotted across the genome  
1045 using 5000 bp sliding windows. **B.** Number of SNPs within 100 bp of each strain-  
1046 specific indel was plotted. The average number of SNPs in an equal number of random  
1047 100 bp windows bootstrapped 1000 times is shown (red line). **C.** The distance of the  
1048 nearest SNP to all indels for each strain pair was plotted. (Compare to association of  
1049 strain-specific SNPs and indels in Figure 1C).

1050 **Figure S17. Associations of strain-specific variants with genic content when**  
1051 **corrected for clade bias.** Strains-specific mutations from 12 strains (Clade I: 12C,  
1052 L26, P78048; Clade II: P57072, P76055, P76067; Clade III: P34048, P78042, P57055;  
1053 Clade SA: P87, GC75, P75063) were retained to ensure that clade representation did  
1054 not bias these associations. The ratio of genic to intergenic strain-specific SNPs (**A**) and  
1055 indels (**B**) was calculated for each of the 12 isolates and falls below the fraction of the  
1056 genome that encodes protein-coding genes (red line). **C.** The density of strain-specific  
1057 SNPs from these 12 isolates was measured for all genes in the *C. albicans* genome and  
1058 plotted against the average SNP density for these isolates (red line). **D.** Genes  
1059 enriched for strain-specific SNPs were identified and functional enrichment was  
1060 assessed using GO term analysis. snoRNAs were significantly enriched for strain-  
1061 specific SNPs in these 12 isolates. **E.** The placement of the SNPs was determined by  
1062 breaking snoRNAs into five equal segments. SNPs were significantly enriched in the 5'  
1063 end of the RNAs.

1064 **Table S1. Percentage of the genome encoding SNPs among the sequenced *C.***  
1065 ***albicans* isolates.**

1066 **Table S2. Strain-specific SNPs among sequenced *C. albicans* isolates.**

1067 **Table S3. Strain-specific indels among sequenced *C. albicans* isolates.**

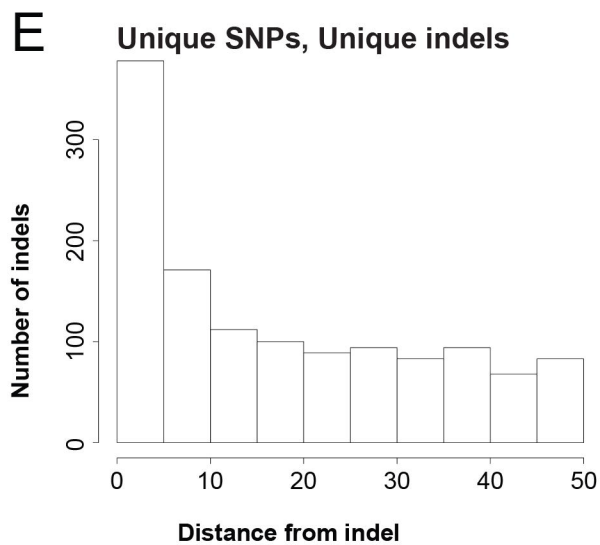
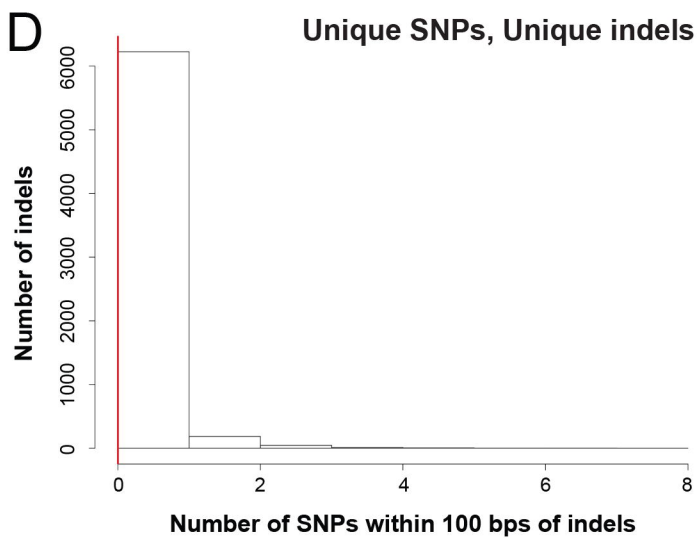
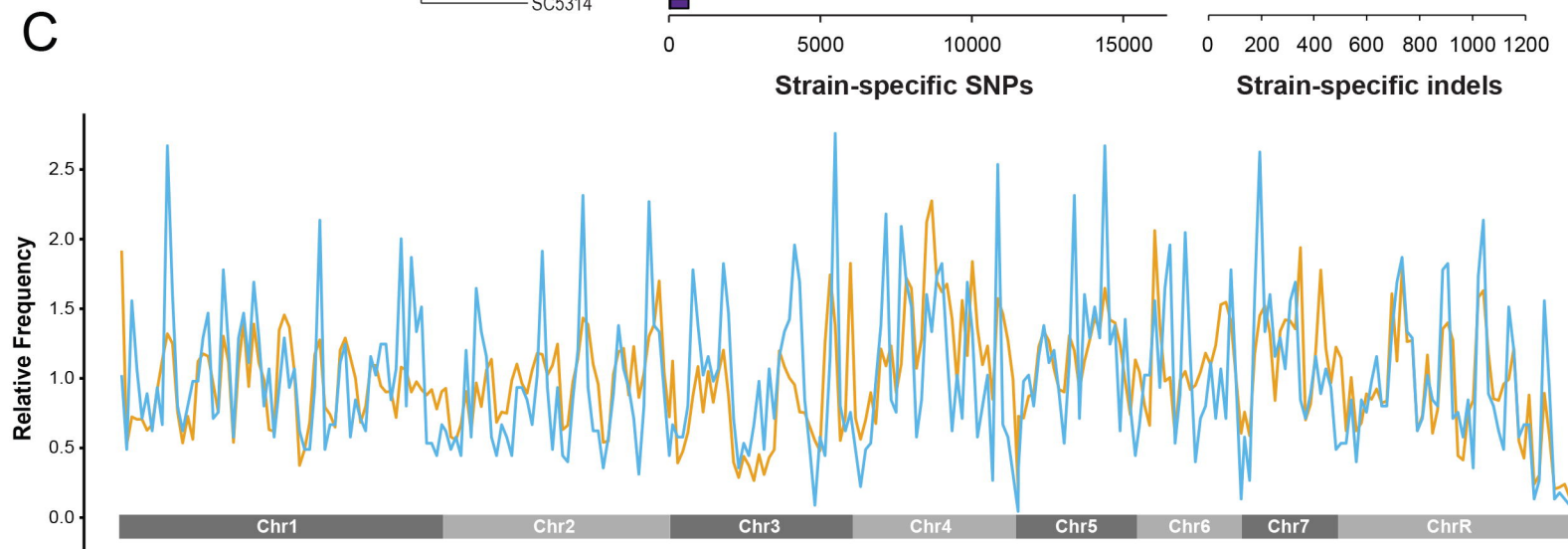
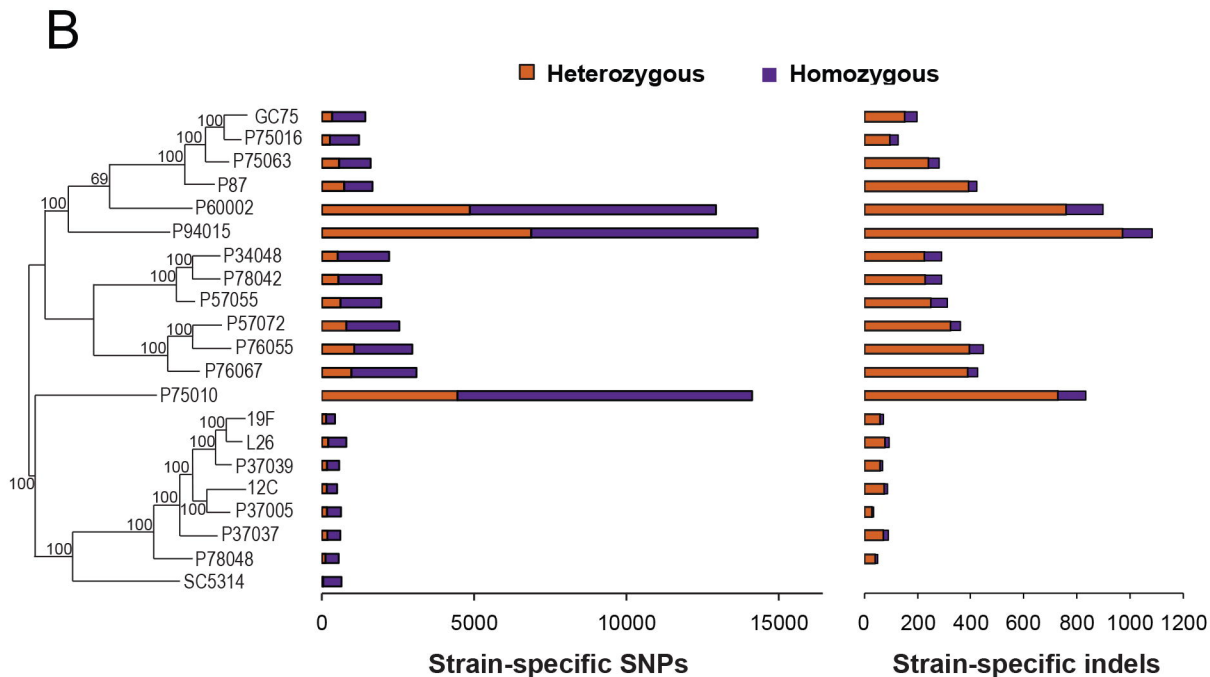
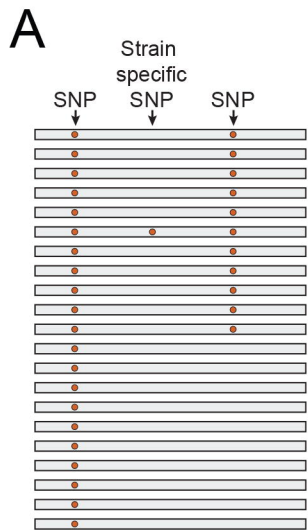
1068 **Table S4. Genes with significant enrichment of strain-specific SNPs.** Gene lists are  
1069 provided at three different SNPs/nt cutoffs.

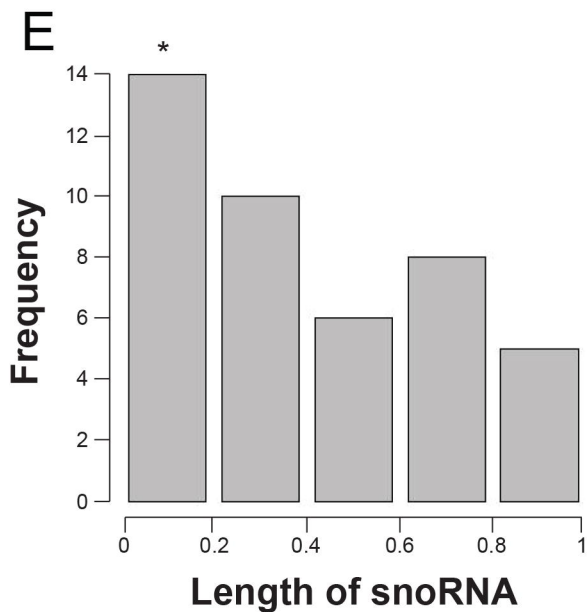
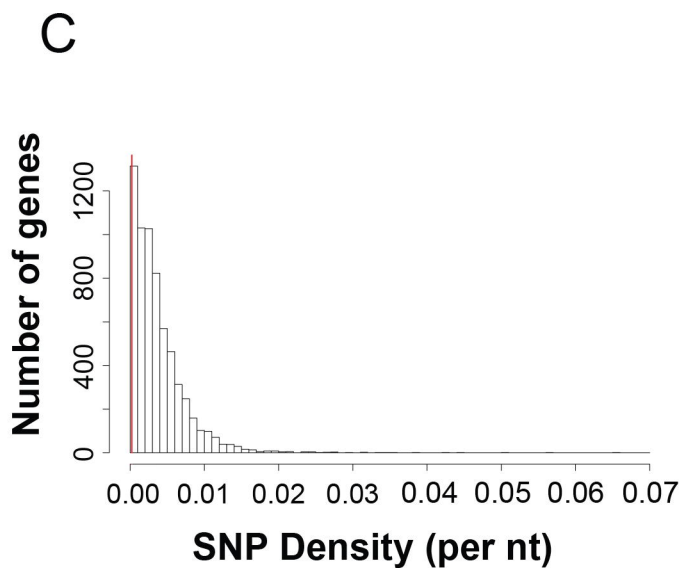
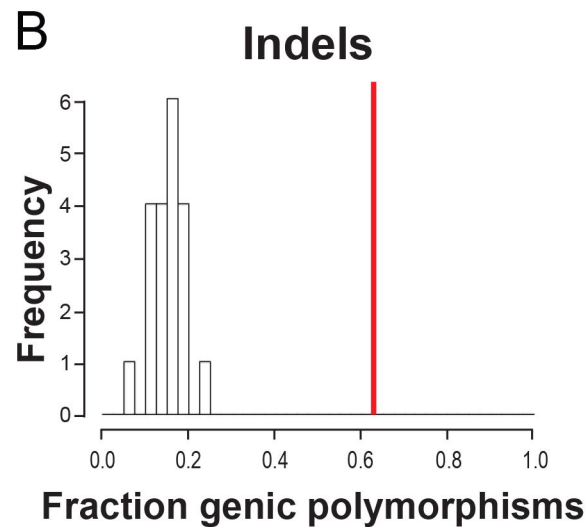
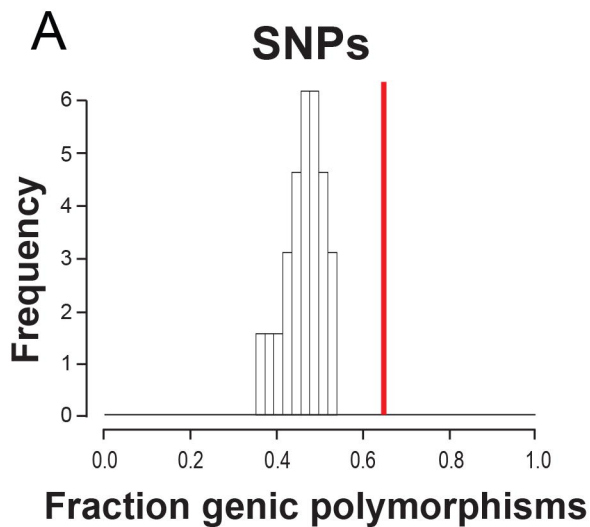
1070 **Table S5. Heterozygous and homozygous regions for the sequenced *C. albicans***  
1071 **isolates.**

1072

1073

1074





<u>Process</u>	<u>Genes included (of 37)</u>	<u>Total possible</u>	<u>Probability (q)</u>
Box C/D snoRNP complex	5	79	5.5e-4

