

1 SUPPORTING INFORMATION

2 SI Results

3 Base substitutions in *C. albicans*

4 As expected, the number of strain-specific mutations increased with longer branch
5 lengths from the nearest node in the phylogenetic tree (SNPs, $r_s = 0.60$, $p = 4.2E-3$; indels, $r_s =$
6 0.42 , $p = 0.055$; Fig. 1B and Fig. S13). Correlation between these metrics of strain identity
7 supports the use of strain-specific mutations in assessing mutational patterns.

8 In many eukaryotes, base-substitution mutations are biased towards transitions over
9 transversions, although the cause of this bias is not completely clear [33]. In *C. albicans*, base
10 substitutions also favored transitions over transversions for both strain-specific SNPs and total
11 SNPs, χ^2 ((11, N = 66086) = 18182, $p < 2E-16$ and (11, N = 302641) = 628000, $p < 2E-16$,
12 respectively). The ratio of transitions to transversions was 2.21 for strain-specific SNPs and
13 2.50 for all SNPs (Fig. S14). Both coding and noncoding regions encoded more strain-specific
14 transitions than transversions, although coding sequences were more biased than noncoding
15 regions (2.74 versus 1.80, respectively). Base substitutions displayed a 1.39-fold bias towards
16 introducing A/T instead of G/C for strain-specific SNPs that shrank to 1.03-fold when including
17 all SNPs. The fact that substitutions favor transitions resulting in A/T suggests that this may
18 contribute to the overall A/T richness of the *C. albicans* genome [31].

19 Analysis of the global distribution of strain-specific SNPs revealed a bias against the
20 accumulation of these mutations within protein-coding genes. Thus, most strain-specific
21 polymorphisms (33,818 of 66,086 SNPs and 5,502 of 6,474 indels) were present within the
22 36.7% of the genome representing intergenic regions, suggesting that mutations in coding
23 sequences are selected against ($p = 9.71E-16$; Fig. S15A,B). As a result, relatively few strain-
24 specific SNPs were present within ORFs across the twenty-one sequenced strains (Fig. S15C).

25 We found that 259 genes exhibited significantly greater SNP densities per nucleotide (nt) than
26 the 0.004 SNPs/nt average for all *C. albicans* ORFs (Fig. S15C, Table S6). SNP densities
27 within enriched genes were equal to or greater than the intergenic average (0.0066 vs. 0.0063,
28 respectively). Protein-coding genes within this group lacked any enrichment for gene ontology
29 (GO) annotations or pathways (Table S6). However, noncoding snoRNAs (small nucleolar
30 RNAs) were significantly overrepresented among 'faster-evolving genes' by GO term analysis,
31 χ^2 ((2, N = 5) = 15.6, p = 7.90E-5; Fig. S15D). The five snoRNAs identified from GO enrichment
32 had mutation rates greater than 0.02 SNPs/nt, significantly higher than that of the average rate
33 of 0.0063 SNPs/nt within intergenic regions. Strain-specific polymorphisms clustered towards
34 the 5' end of the snoRNAs (Fig. S15E) and could contribute to variation in functional aspects of
35 protein translation, although this possibility was not explored here.

36 **Indels in *C. albicans***

37 An inspection of strain-specific indels revealed that 3527 (54.5%) were deletions and
38 2948 (45.5%) were insertions. Indels ranged in size from 1 bp to 10 bp with the majority of
39 longer events being insertions (Fig. S16). The frequency of both insertions and deletions
40 decreased as mutations became larger, suggesting that smaller events occur more frequently or
41 are less detrimental to the cell and therefore are retained more often. The incidence of ± 3
42 nucleotide indels (21.9% of the total) was higher than that expected by chance. When indels
43 were separated into genic or intergenic mutations, intergenic mutations followed a Poisson
44 distribution centered on 0, whereas genic mutations were vastly overrepresented for ± 3
45 nucleotide indels that do not shift the reading frame (Fig. S16). Only ~15% of all indels fell
46 within ORFs (p < 2.2E-16) suggesting that, as with SNPs, indels are selected against within
47 coding sequences (Fig. S15B).

48 Indels have been commonly associated with specific genomic features such as repetitive
49 sequences in other species [34, 35]. Across the sequenced *C. albicans* isolates, there was a

50 total of 19,581 indel sites across the genome. Of these, 465 indel sites (2.37%) were located
51 within annotated repetitive sequences (long terminal repeats (LTRs), major repeat sequences
52 (MRSs), and retrotransposons). Total indels are therefore overrepresented within these
53 repetitive features (two-tailed Brunner-Munzel (BM) test = 5.15, df = 182.05, p = 6.65E-7).
54 Likewise, strain-specific indels were significantly enriched within repetitive features (47 of 6475;
55 BM test = 13.004, df=182, p<2.2E-16). Both total and strain-specific SNPs also clustered within
56 repetitive elements (BM test = 14.98, df = 315.62, p < 2.2E-16 and BM test = 12.26, df = 240.02,
57 p < 2.2E-16, respectively). Thus, mutations within the *C. albicans* genome are enriched within
58 repetitive regions similar to what has been observed in other species [34, 35].

59 **Association of SNPs and indels**

60 Similar to our analysis of strain-specific SNPs and indels described in the main text, an
61 examination of all annotated SNPs and indels found them to be associated with each other in *C.*
62 *albicans* genomes (Fig. S17A). Again, all SNPs were significantly enriched within a 100 bp
63 window surrounding each indel (Wilcoxon test (W(1.66E11)), p<2.2E-16, Fig. S17B). These
64 polymorphisms clustered within the proximal 10 bp to each indel but did not overlap the indel
65 position (Fig. S17C), similar to strain-specific mutations. Thus, SNPS and indels display the
66 same tight association with one another regardless of whether these are strain-specific or total
67 polymorphisms.

68 To assess the accuracy of regions containing indels and adjacent SNPs, 21 regions
69 were Sanger sequenced across the indel-SNP associations. Sanger sequencing confirmed 17
70 of 21 associations (81%) identified from variant calling of the assembled 21 genomes. These
71 validated regions contained identical variants at the predicted position and the calculated allele
72 frequencies as that produced from genome assembly (Fig. S4). Importantly, the interrogated
73 regions represent highly repetitive stretches that are more likely to be called incorrectly by whole

74 genome analysis. Thus, the high accuracy of variant calls among the Sanger sequenced
75 positions likely reflects that most variant positions are accurately called.

76 In some species, the introduction of indels can influence the observed mutational bias
77 towards either transitions or transversions [35, 36]. To address this possibility in *C. albicans*,
78 the transition:transversion ratio was determined for the ~500 strain-specific SNPs located within
79 10 bp of strain-specific indels. Although base substitutions still slightly favored transitions, the
80 1.17 transition:transversion ratio was significantly lower than the genome-wide average ratio of
81 2.21 ($p = 5.87E-7$). This is consistent with mutations close to indels exhibiting a reduced bias
82 towards transitions over transversions due to recruitment of error-prone polymerases during
83 DNA repair [35]. We therefore suggest that a similar mechanism operates in *C. albicans* and
84 can account for the increased mutation rate adjacent to indels, as well as the local bias in the
85 transition:transversion ratio.

86 **Loss of Heterozygosity tracts in *C. albicans***

87 Mapping the LOH regions among the sequenced *C. albicans* isolates revealed that out of a total
88 of 336 chromosome arms in the 21 isolates, 155 of these arms show evidence of having
89 undergone a long-tract LOH event. Thus, almost half of the chromosomes contained an LOH of
90 greater than 50 kb. LOH frequency decreased towards the centromeres and did not occur
91 across centromeres except during LOH of whole chromosomes (Fig. S18A).

92 **Association of mutations with LOH**

93 The previous study by Hirakawa *et al.* identified extensive loss of heterozygosity (LOH)
94 tracts in the 21 sequenced *C. albicans* isolates [28]. Consequently, LOH breakpoints were
95 mapped in each isolate and emphasis was placed on the distribution of LOH events around the
96 mating type-like (*MTL*) locus on Chr5. The current study extends the analysis of LOH patterns
97 in *C. albicans* genomes by determining if genome-wide patterns of LOH exist, and if there is an

98 association between LOH tracts and other mutational classes such as base substitutions or
99 indels. Importantly, aneuploidy did not significantly alter the frequency of heterozygous and
100 homozygous intervals along aneuploid chromosomes relative to euploid chromosomes ($p =$
101 0.756).

102 We note that heterozygous SNPs may have arisen in hom regions but, in some cases,
103 been eliminated by a subsequent LOH event. As LOH can occur in one of two possible
104 directions (due to loss of either homolog A or homolog B), we accounted for mutations
105 potentially lost via LOH by doubling the number of homozygous, strain-specific SNPs within
106 hom regions. Even with this adjustment, het regions still contained a greater density of strain-
107 specific SNPs than hom regions (two-sided BM test = -8.74, $df = 717.58$, $p < 2.2E-16$). The ~2-
108 fold greater accumulation of polymorphisms in het over hom regions of the *C. albicans* genome
109 shows parallels with the ~3.5-fold higher mutation rate observed in het vs. hom regions of the
110 *Arabidopsis* genome during meiosis [37].

111 Sites close to recombination events, including LOH events, have been shown to be
112 associated with elevated mutation rates in some species [36, 38-40]. To determine if there is an
113 increased frequency of SNPs in regions proximal to LOH tracts in *C. albicans*, the density of
114 SNPs at heterozygous-homozygous transition points was investigated. Analysis of the 745
115 identified transition regions included 1 kb of DNA on either side of the junction points between
116 het and hom regions (with the latter inferred to represent LOH tracts). The SNP density within
117 these transition regions was significantly lower than that in the rest of the *C. albicans* genome
118 (one-sided BM test = -35.415, $df = 748.8$, $p < 2.2E-16$; Fig. S18B). Furthermore, SNP density
119 was similar on both het and hom sides of the LOH breakpoint. Thus, base substitutions appear
120 to accumulate less frequently in regions proximal to het/hom breakpoints in the *C. albicans*
121 genome. One caveat noted here is that this result may be influenced by difficulty in the
122 identification of precise breakpoints between het and hom regions of the genome.

123 When evaluating all polymorphisms, heterozygous regions of the genome contained
124 significantly higher frequencies of SNPs relative to hom regions (1.4E-4 vs. 7.3E-5 SNPs/bp,
125 respectively; two-tailed Brunner-Munzel (BM) test = -10.6, df=786.6, $p < 2E-16$). Therefore, as
126 with the elevated frequency of strain-specific SNPs in heterozygous regions, we find higher
127 frequencies of all SNPs in het than hom regions, indicating that this pattern has been retained
128 during species evolution.

129 The association of indels with heterozygous and homozygous regions of the genome
130 was more complex. Strain-specific indels were more commonly found in hom regions than in
131 het regions (BM test = -3.22, df=878.4, $p = 0.0013$). In contrast, analysis of all indels displayed a
132 highly significant enrichment in het regions compared to hom regions (BM test = -5.78,
133 df=729.8, $p = 1.08E-8$). Thus, although strain-specific indels are more commonly found in hom
134 regions, het regions encode significantly more indels overall, suggesting that a simple
135 relationship between indel formation/retention and het/hom regions does not exist.

136 LOH regions were defined by 5 kb windows encoding >0.4 events per kb. We note here
137 that this classification could misclassify some 'hom' regions as 'het' regions and vice versa.
138 However, misclassification of a subset of these regions would likely result in our analysis
139 underestimating any divergent mutational trends between het and hom regions.

140 **Impact of phylogenetic branch length on evolutionary trends**

141 Phylogenetic reconstruction produced longer individual branch lengths for three strains,
142 P60002, P75010, and P94015, compared to the other 18 strains in this study (Fig. 1A). To
143 investigate if the described evolutionary associations between mutations are biased by the
144 inclusion of these three 'outlier' strains, or by having unequal numbers of strains representing
145 each clade, three strains were picked from each of the four major clades (I: 12C, L26, P78048;
146 II: P57072, P76055, P76067; III: P34048, P78042, P57055; SA: P87, GC75, P75063) and the

147 mutational trends reassessed. In total, these 12 strains accounted for 53% of all strain-specific
148 SNPs and 25% of all strain-specific indels, leaving 31,090 strain-specific SNPs and 4,947 strain-
149 specific indels, respectively.

150 The association between SNP and indel polymorphisms was still observed when
151 analyzing the more limited set of 12 strains, as we again found a strong association between the
152 positions of strain-specific SNPs and indels (Wilcoxon test ($W(3.04E7)$), $p < 2.2E-16$, Fig. S5A,B).
153 Similarly, analysis of the association between all SNPs and indels revealed a tight association
154 across the 12 strains (Wilcoxon test ($W(2.12E7)$), $p < 2.2E-16$, Fig. S5C). Thus, the association
155 between SNPs and indels across the genome was robust across *C. albicans* isolates.

156 We next examined whether the overrepresentation of polymorphisms accumulating in
157 heterozygous regions of the genome was still observed within the set of 12 strains. Strain-
158 specific polymorphisms were still more prevalent in het than hom regions of the genome when
159 analyzing polymorphisms from this set (SNPs; BM test = -7.558 , $df=377.0$, $p < 3.14E-13$, indels;
160 BM test = -2.35 , $df=379.7$, $p=1.97E-2$, Fig; S19A,B). This association also existed among all
161 polymorphisms, as polymorphisms were overrepresented in heterozygous regions of the
162 genome in the 12 strains (SNPs; BM test = -14.80 , $df=422.6$, $p < 2.2E-16$, indels; BM test = -
163 3.188 , $df=321.8$, $p=1.57E-3$).

164 Finally, we identified genes enriched for strain-specific base substitutions among this set
165 of 12 strains. A total of 15 genes encoded greater than 0.02 SNPs/nt (Figure S19C). As
166 observed for all strain-specific SNPs, snoRNAs were greatly overrepresented among genes for
167 strain-specific polymorphisms (Fig. S19D) and these mutations clustered towards the 5' end of
168 the snoRNA sequence (Fig. S19E). In contrast, no GO terms were significantly enriched among
169 genes with high numbers of all base substitutions among the 12 strains. Therefore, analysis of
170 the more limited set of 12 strains representing all 4 major clades did not significantly alter the
171 association of mutation type, the accumulation of SNPs in heterozygous regions, or the class of

172 functional RNAs most frequently containing emergent base substitutions despite the loss of the
173 majority of strain-specific polymorphisms relative to the set of 21 isolates.

174 **Genome evolution processes contributing to SNP incongruence**

175 The plasticity of the *C. albicans* genome allows it to adapt quickly through a combination
176 of mechanisms that potentially include both LOH and mating (4-6). The ability to decipher
177 between these two possibilities is challenging in heterozygous diploid genomes where different
178 patterns of LOH can produce genomes that appear 'recombinant' despite a clonal origin. For
179 example, consider a strain encoding heterozygous SNPs along both chromosome homologs.
180 Homozygosis of one homolog or the other will generate strains that have different SNP patterns
181 relative to each another, although both will share SNPs with the parental strain that did not
182 undergo homozygosis (Fig. S9). In isolation, the two homozygosed strains would be viewed as
183 unrelated to one another because of their completely different SNP patterns despite sharing a
184 common origin. If the parental strain was also present, it would appear to be the 'recombinant'
185 product of mating between the two homozygosed strains (Fig. S9). Additionally, short-tract LOH
186 events could further complicate analysis due to shuffling of the polymorphic patterns between
187 related strains (Fig. S9).

188 Careful analysis of the 21 genome sequences found examples where 'identity by
189 descent' patterns have been broken by different LOH events occurring in different lineages.
190 One example involves *orf19.4819* which is heterozygous for four SNPs in Clade II (Fig. 2D,E).
191 Most Clade I strains have retained only one of the four SNPs, whereas Clade III isolates and the
192 Clade E isolate have retained the other three SNPs. Given the position of these strains in the
193 phylogenetic tree, the most parsimonious explanation is that Clade II has retained the ancestral
194 pattern of heterozygous SNPs for *orf19.4819*, and that the other lineages underwent loss of one
195 homolog or the other generating the extant patterns shown. A second example is shown for
196 SNPs in *orf19.6605* in which it appears that extant strains lost one homolog or the other,

197 followed by additional mutagenic events (Fig. S8A,B). These patterns are mostly easily
198 recognized when all three SNP configurations are present in a population: one isolate with the
199 parental configuration of heterozygous SNPs, one isolate in which LOH has homozygosed the
200 SNPs for one allele, and another isolate in which SNPs have homozygosed for the other allele
201 (Fig. 2D,E).

202 Mating between strains is a second key mechanism that can distort mutation patterns
203 from that expected by 'identity by descent'. A possible indication of mating is a sudden break of
204 homology along a chromosome so that SNPs switch from matching one clade to matching a
205 second clade (Fig. 3A,B,D). We also found evidence of mating by identifying strains that have
206 inherited one chromosome homolog from two different extant strains (Fig. 3C, S9). The
207 genotype of the resulting strain will be a composite of polymorphisms from each parent,
208 although additional recombination or mutation events can make such regions challenging to
209 identify. Inference of these events also requires phasing of SNPs for each homolog in both
210 parental strains, as we performed for the region on Chr3 in Fig. 3C. In this case, we reason that
211 it is difficult to envisage how this pattern of SNPs could be generated by LOH events, as the
212 SNPs reflect the expected pattern of heterozygous and homozygous positions when the two
213 parental homologs are combined. Furthermore, we emphasize that the SNPs involved are
214 identical (both in their positions and in the actual nucleotide substitutions) in the highlighted
215 regions. The most parsimonious explanation is that this section of the P94015 genome
216 therefore arose via a hybridization event between a Clade III strain and a Clade SA strain.

217

218 **REFERENCES**

- 219 1. Rosenberg MS, Subramanian S, & Kumar S (2003) Patterns of transitional mutation
220 biases within and among mammalian genomes. *Mol Biol Evol* 20(6):988-993.
- 221 2. Butler G, *et al.* (2009) Evolution of pathogenicity and sexual reproduction in eight
222 *Candida* genomes. *Nature* 459(7247):657-662.
- 223 3. Jovelin R & Cutter AD (2013) Fine-scale signatures of molecular evolution reconcile
224 models of indel-associated mutation. *Genome Biol Evol* 5(5):978-986.
- 225 4. Zhang N, *et al.* (2015) Selective advantages of a parasexual cycle for the yeast *Candida*
226 *albicans*. *Genetics* 200(4):1117-1132.
- 227 5. Bennett RJ, Forche A, & Berman J (2014) Rapid mechanisms for generating genome
228 diversity: whole ploidy shifts, aneuploidy, and loss of heterozygosity. *Cold Spring Harb*
229 *Perspect Med* 4(10).
- 230 6. Ciudad T, Hickman M, Bellido A, Berman J, & Larriba G (2016) Phenotypic
231 consequences of a spontaneous loss of heterozygosity in a common laboratory strain of
232 *Candida albicans*. *Genetics* 203(3):1161-1176.

233

234