

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20

Genomic evidence of an early evolutionary divergence event in wild *Saccharomyces cerevisiae*

Devin P Bendixsen*, Noah Gettle, Ciaran Gilchrist, Zebin Zhang, and Rike Stelkens

Division of Population Genetics, Department of Zoology, Stockholm University, Svante Arrheniusväg 18
B, 106 91 Stockholm, Sweden

* Author for Correspondence: Devin P Bendixsen, Department of Zoology Stockholm University;
email: devin.bendixsen@zoologi.su.se

Keywords *Saccharomyces cerevisiae*, yeast, long-read, genome assembly, structural variation, Ty
element

21 **Abstract**

22 Comparative genome analyses have suggested East Asia to be the cradle of the domesticated microbe
23 Brewer's yeast (*Saccharomyces cerevisiae*), used in the food and biotechnology industry worldwide. Here,
24 we provide seven new, high quality long read genomes of non-domesticated yeast strains isolated from
25 primeval forests and other natural environments in China and Taiwan. In a comprehensive analysis of our new
26 genome assemblies, along with other long read *Saccharomycetes* genomes available, we show that the newly
27 sequenced East Asian strains are among the closest living relatives of the ancestors of the global diversity of
28 Brewer's yeast, confirming predictions made from short read genomic data. Three of these strains (termed the
29 East Asian Clade IX Complex here) share a recent ancestry and evolutionary history suggesting an early
30 divergence from other *S. cerevisiae* strains before the larger radiation of the species, and prior to its
31 domestication. Our genomic analyses reveal that the wild East Asian strains contain elevated levels of structural
32 variations. The new genomic resources provided here contribute to our understanding of the natural
33 diversity of *S. cerevisiae*, expand the intraspecific genetic variation found in this this heavily domesticated
34 microbe, and provide a foundation for understanding its origin and global colonization history.

35

36 **Significance statement**

37 Brewer's yeast (*Saccharomyces cerevisiae*) is a domesticated microbe and research model organism with
38 a global distribution, and suspected origin in East Asia. So far only limited genomic resources are available
39 from non-domesticated lineages. This study provides seven new, high quality long read genomes of strains
40 isolated from primeval forests and other natural environments in China and Taiwan. Comparative genomics
41 reveal elevated levels of structural variation in this group, and early phylogenetic branching prior to the global
42 radiation of the species. These new genomic resources expand our understanding of the evolutionary history
43 of Brewer's yeast, and illustrate what the ancestors of this highly successful microbe may have looked like.

44 **Introduction**

45 The history of Brewer's yeast, *Saccharomyces cerevisiae*, is deeply interwoven with that of humanity,
46 having played significant roles in cultural, technological, and societal development for at least 9000 years
47 (McGovern, et al. 2004). While over a hundred years of *S. cerevisiae* research has provided important
48 insights into eukaryotic genomics, evolution and cell physiology, much of its 'wild' ecology as well as its
49 deep human and pre-human evolutionary history have, until recently, largely remained a mystery. Recent
50 broadscale genomic surveys of *S. cerevisiae* and its close relatives, however, are beginning to shed light
51 on important aspects of its population genetic structure, intra- and inter-specific hybridization events, and
52 their interplay in yeast domestication (Duan, et al. 2018; Peter, et al. 2018; Scannell, et al. 2011; Wang, et
53 al. 2012).

54 One of the key results from these broadscale genomic surveys has been increasing evidence for a
55 singular and central radiation event of *S. cerevisiae* from Far East Asia (Wang, et al. 2012; Peter, et al.
56 2018; Duan, et al. 2018). These studies have independently revealed that strains of wild yeast collected in
57 parts of China and Taiwan contain much higher genomic diversity and show greater levels of divergence
58 than all other strains of *S. cerevisiae*. The vast majority of these *S. cerevisiae* genomes, however, have
59 been analyzed using short-read sequencing, resulting in a focus on single nucleotide variants (SNVs).
60 Larger structural variations (SVs), such as inversions, deletions and gene duplications, in addition to
61 repetitive regions such as transposable elements (TE) and telomeres, have gone largely unresolved
62 (Goodwin, et al. 2016). In addition to playing significant roles in yeast adaptation (Payen, et al. 2014;
63 Steenwyk and Rokas 2018; Zhang, et al. 2020), these large structural features can provide increased
64 phylogenetic resolution and key insights about lineage interactions and potential reproductive isolation.
65 SVs have shown to be vital for evolutionary adaptation in many other taxa, supporting the role of
66 inversions in adaptation and speciation, and in the evolution of disease (Merker, et al. 2018;
67 Wellenreuther, et al. 2019).

68 In this study, we generated high quality assemblies of seven of the highly divergent wild East Asian
69 strains and one common laboratory strain (Y55) using both short-reads and PacBio long-reads to better
70 understand the relationships of these strains to the global diversity of *S. cerevisiae*. Analyzing our
71 assemblies in the context of publicly available long-read genomes, we generated a new phylogeny that
72 confirms the place of these East Asian strains at the base of *S. cerevisiae*, and provide further evidence for
73 an out-of-China colonization history of this species. Moreover, we were able to group our sequenced
74 strains belonging to the previously identified CHN IX clade with a Taiwanese strain, both shown in separate
75 studies to be divergent from the rest of *S. cerevisiae*. We show that this combined clade likely has deep
76 roots in mainland China and has had little gene flow with other *S. cerevisiae* strains.

77

78 **Results**

79 *Genome sequencing and assembly*

80 We used whole-genome long-read PacBio sequencing to assemble the genomes of seven highly divergent
81 and one common lab strain of *S. cerevisiae* (**Fig. S1**, average per base genomic coverage = 91.3; average
82 median read length = 3028bp). Initial nuclear and mitochondrial assemblies were highly complete (median
83 # of contigs = 25; median N50 = 821424.5). Final nuclear and mitochondrial assemblies were further
84 resolved to single contigs for each chromosome (**Table S1**, median N50 = 907965.5). Final genome sizes
85 ranged from 11.65 to 11.92 Mbp. Assessment of the completeness of the genome assembly and
86 annotation using BUSCO found that all genomes had similarly high BUSCO scores (C > 96.5%, **Table S2**).

87

88 *Phylogenomics*

89 Our newly constructed consensus species tree, placed six of the newly assembled East Asian strains in a
90 basal position within the *S. cerevisiae* radiation (**Fig. 1**). Three of these strains, EM14S01-3B (Taiwanese)
91 (Peter, et al. 2018), XXYS1.4 and JXXY16.1 (CHN IX) (Duan, et al. 2018), hereon referred to as the East Asian

92 Clade IX Complex, show early divergence from all other *S. cerevisiae* strains. Despite the largely basal
93 placement of our assembled East Asian strains, one strain (BJ4) clustered separately with Y12 and YPS128,
94 strains isolated from Ivory Coast palm wine and Pennsylvanian woodland soil, respectively. The common
95 lab strain, Y55, clustered with two other domesticated strains (DBVPG6044 and SK1) within the West
96 African+ clade. Construction of an Alignment and Assembly-Free (AAF) phylogeny comparing the long-
97 read sequencing data generated in this study and previous short-read data found a high-level of similarity
98 between the two datasets (**Fig. S2-S3**). This analysis also found similar clustering to the consensus species
99 tree, among the East Asian strains and the common lab strain as well as a large amount of divergence of
100 the East Asian Clade IX Complex from the rest of *S. cerevisiae*. However, AAF was unable to resolve the
101 deep early divergence of the East Asian Clade IX Complex from other strains.

102

103 *Structural variation*

104 A comparison of our eight *S. cerevisiae* genomes and previously assembled *Saccharomyces sensu stricto*
105 genomes to the *S. cerevisiae* reference genome (S288C), revealed a high level of collinearity, particularly
106 at larger scales (**Fig. 2, Fig. S4-S10**). We found exceptions to this strict collinearity only in one strain of *S.*
107 *paradoxus* (previously reported (Yue, et al. 2017)) and in the East Asian Clade IX Complex. All three
108 member strains show a ~80kb terminal translocation from chromosome XI to chromosome XII (**Fig. 2A**
109 **inset**). This structural variant in the East Asian Clade IX Complex was further supported by both long- and
110 short-read analyses of alignment coverage (**Fig. S11-S12**). Additional evidence for this unique
111 translocation comes from high short-read coverage of chromosome XII of XXYS1.4, indicating a likely
112 aneuploidy, which extends across the translocated region of chromosome XI (**Fig. S12**). Other notable
113 rearrangements are a large inversion in chromosome X of BJ4 (**Fig. S5**). The common lab strain, Y55,
114 showed a high level of collinearity with only minor deviations from homology (**Fig. S4**).

115 To quantify the extent of smaller structural variations in our genomes, we performed a
116 comprehensive analysis using pairwise comparisons between the 15 *S. cerevisiae* strains with long-read
117 assemblies. We assessed five types of variation: deletions, insertions, duplications, inversions and
118 translocations. This analysis revealed that the wild East Asian strains tend to have higher amounts of total
119 variation (mean=356.5) compared to the other strains (mean=384.7, **Fig. 3A**). The three Clade IX Complex
120 strains (EM14S01-3B, JXXY16.1, XXYS1.4) were among the highest, and in particular strain XXYS1.4 had a
121 significantly higher mean variant count (525.9, **Fig. 4B**). In contrast, the common lab strain Y55 had more
122 moderate levels of total variation. The East Asian Clade IX Complex also had larger numbers of deletions
123 and inversions, and fewer insertions and duplications (**Fig. 4C, Fig. S13-S17**). The Malaysian strain
124 UWOPS03-461.4 had significantly larger numbers of translocations compared to all strains. A closer
125 analysis of the distribution of all structural variations identified in this study along chromosomes revealed
126 areas of elevated variation counts, however we found no strong patterns (**Fig. S18**).

127

128 *Nuclear genome content*

129 In general, our newly assembled long-read genomes were significantly smaller than the currently existing
130 genomes ($t = 2.36$, $df = 10.28$, $p = 0.039$). This difference, however, is largely a result of reduced genome
131 size in members of the East Asian Clade IX Complex ($\bar{x}_{\text{CladeIX}} = 11.72\text{Mbp}$; $\bar{x}_{\text{CladeIX}} = 11.89\text{Mbp}$; $t = 2.93$, $df =$
132 5.25 , $p = 0.031$). These size differences are due to decreases in genic material both in terms of counts ($t =$
133 6.12 , $df = 10.43$, $p < 0.001$) and cumulative gene length ($t = 7.35$, $df = 5.0$, $p < 0.001$), and a relative
134 reduction in non-coding DNA ($t = 4.95$, $df = 6.74$, $p = 0.002$) (**Fig. S19**). Interestingly, these relative
135 reductions in genic material are correlated with increases in identified intronic material, a pattern that is
136 carried throughout all *S. cerevisiae* strains analyzed here ($F = 11.38$; $r^2 = 0.43$; $p = 0.005$).

137

138 *Transposable Element Composition*

139 Transposable elements (TEs) replicate and deteriorate in a way that gives them an evolutionary history
140 that can be unique with regards to their host genomes and can provide hints about past interactions
141 between distinct lineages. To better understand historical relationships between different strains of *S.*
142 *cerevisiae*, we annotated and analyzed all classes of known retrotransposon or *Ty* element in this species.

143 In terms of simple counts, members of the East Asian Clade IX Complex had more *Ty*-associated
144 elements than the rest of the *S. cerevisiae* strains ($t = -6.05$, $df = 6.31$, $p < 0.001$), a result largely based
145 on a disproportionate number of solo long terminal repeats (LTRs) across all classes of *Ty* elements (**Fig.**
146 **4A, Fig. S20-S21**). A similar pattern remained when comparing total length of elements (**Fig. S22**).
147 Although *Ty1/Ty2* LTRs were the most common *Ty* remnant in all strains, the relative frequency of each
148 class of *Ty* element across *S. cerevisiae* strains does not follow the same pattern reported for the reference
149 strain S288C, where $Ty1 > Ty2 > Ty3 > Ty4 > Ty5$. Indeed, *Ty1* elements have often been suggested as being
150 the most prolific TE class in *S. cerevisiae*; however, we did not find any putatively functional *Ty1* elements
151 in 6 of the 15 strains we analyze while finding 30 in the reference strain, S288C, representing a clear outlier
152 at the upper end.

153 As yet, functional *Ty5* elements had only been identified in *S. paradoxus*. “Complete” elements
154 (i.e. elements containing both flanking LTRs and the internal coding region) previously identified in *S.*
155 *cerevisiae* strains are missing a ~2kb portion of the ~5kb internal coding region and are found in very low
156 numbers (1-2 per strain). However, the Clade IX Complex strains show a particularly high abundance of
157 *Ty5*-associated elements (**Fig. 4B**). Further examination revealed six complete *Ty5* elements with fully
158 intact coding regions distributed across two Clade IX Complex strains, EM14S01-3B and JXXY16 (**Fig. S23**).
159 While all “complete” *Ty5* elements that we identified in *S. cerevisiae* outside of the Clade IX Complex are
160 missing the same ~2kb region, only 2/10 Clade IX *Ty5* elements (both in JXXY16.1) are missing this region.
161 Additionally, these elements largely do not share homologous bordering regions. In conclusion, the only

162 putatively functional *Ty5* elements in *S. cerevisiae* are in the Clade IX Complex.

163

164 *Comparative mitochondrial genomics*

165 Overall, the mitochondrial genomes of the *S. cerevisiae* strains showed high-levels of collinearity (**Fig. 5**).

166 Of note, however, is the absence of RPM1, a highly conserved ncRNA component of mitochondrial RNase

167 P in two of the Clade IX Complex strains, JXXY16.1 and XXYS1.4. To further confirm the absence of this

168 gene we aligned the reference RPM1 to the unassembled PacBio reads using BLASTn (Zhang, et al. 2000).

169 We found no full-length alignments of RPM1, a 483 bp gene, in either set of reads; rather the highest

170 scoring alignments (e-value>9e-35) were 149 (JXXY16.1) and 239 bp (XXYS1.4). Similarly, we were unable

171 to find or assemble more than a truncated version of the mitochondrial 21s rRNA in JXXY16.1 None of the

172 strains we sequenced were found to be respiratory incompetents or ρ^- .

173 Previous analyses have suggested that hybridization events can generate discordance between

174 species and mitochondrial phylogenies in yeast (De Chiara, et al. 2020; Peris, et al. 2017). To investigate

175 this, we also included other *Saccharomyces* species with available long read data in our mitochondrial

176 phylogenetic analyses. For the most part, our mitochondrial phylogenies matched our species-level

177 phylogeny with the notable exception of a strain of *S. jurei* (NCYC3947), a recently described European

178 species (Naseeb, et al. 2018) that appears to share mitochondrial ancestry with a subgroup of European

179 strains of *S. paradoxus*. The mitochondrial genomes from this subgroup also contain large structural

180 variations (previously described in Yue et al. 2017) not seen in other strains of *S. paradoxus* and *S.*

181 *cerevisiae*, further supporting their shared ancestry (**Fig. 5**).

182

183 *Intraspecific spore viability*

184 Lastly, we performed intraspecific crosses of each wild East Asian strain with the common lab strain Y55,

185 to assess the level of reproductive isolation. As expected, we found a lower level of viable spores when

186 crossing with a divergent wild strain as compared to self-crossing Y55 (ANOVA $F(7,152) = 9.63$, $p < 0.001$,
187 **Fig. S24**). Most crosses with East Asian strains reduced spore viability by ~50%, while crosses with HN1
188 reduced viability by ~75%.

189

190 **Discussion**

191 Comparative genomic analyses have provided clues about the origin of Brewer's yeast and have suggested
192 an out-of-China origin (Peter, et al. 2018). Here, we provide seven new, high quality long/short-read
193 genomes of highly divergent wild *S. cerevisiae* strains recently isolated in Far East Asia. Phylogenomic
194 analyses of the long-read assemblies agree with previous findings that the wild East Asian strains (CHN,
195 Taiwanese) are basal relative to other *S. cerevisiae* strains (Duan, et al. 2018; Peter, et al. 2018) and, in
196 the case of the CHN IX and Taiwanese clades, show considerable divergence (**Fig. 1**). In addition, we show
197 that the CHN IX clade (represented here by JXXY16.1 and XXYS1.4) and the one strain representing the
198 Taiwanese clade (EM14S01-3B), likely compose a single monophyletic group distinct from not only the
199 other East Asian strains in our study but also all other strains of *S. cerevisiae* sequenced to date.

200 Our comprehensive analysis of structural variations (SVs) further elucidates the evolutionary
201 history and intraspecific diversity of *S. cerevisiae*. Structural variations were identified for each strain pair
202 revealing patterns of genomic divergence. There is a noteworthy pattern of higher amounts of SVs in wild
203 East Asian strains, especially in the three strains within the Clade IX Complex. As a species, *S. cerevisiae*
204 has been shown to accumulate balanced variations at a slower rate compared to *S. paradoxus* (Yue, et al.
205 2017). This is likely due to the different selection histories of these species. Many *S. cerevisiae* strains have
206 long been associated with human activities where domestication, cross-breeding and admixture have
207 resulted in largely mosaic genomes (Hyma and Fay 2013; Liti 2015; Liti, et al. 2009), whereas *S. paradoxus*
208 strains are recently isolated, wild strains. Interestingly, we found that wild East Asian strains accumulated
209 both structural variations at a high rate, more similar to rates normally seen in *S. paradoxus* (**Fig. 3**). It has

210 been suggested that the geographical isolation of some *S. paradoxus* subpopulations may have favored
211 quick fixation of structural rearrangements (Leducq, et al. 2016). We may be witnessing similar patterns
212 in the wild East Asian *S. cerevisiae* strains.

213 In context of the seven previously assembled *S. cerevisiae* long-read genomes and those from
214 *Saccharomyces sensu stricto* species, our results reveal other important aspects of yeast evolutionary
215 genomic history. Not only do the phylogenetic patterns we describe reveal discrete boundaries between
216 certain clade-levels in terms of TEs, indicating that transfer of persisting TEs between deep-rooted clades
217 either through horizontal gene transfer or hybridization is rare (**Fig. 4**). They also give us context for the
218 evolutionary history of these elements in their own right. Interestingly, we found that *Ty5*, a relatively
219 rare retrotransposon with no previously known functional versions in *S. cerevisiae*, has retained
220 functionality in the divergent East Asian Clade IX complex. Additionally, we found that *Ty2*, a TE suggested
221 to be a recent introduction to *S. cerevisiae* via *S. mikatae* (Carr, et al. 2012; Liti, et al. 2005), is also present
222 in the East Asian Clade IX complex. This indicates that this event occurred early in *S. cerevisiae* history,
223 that the donor-donee relationship is reversed, that it happened multiply, or that this element was lost in
224 *S. paradoxus* and other closely related species. With respect to the latter hypothesis, our genomic survey
225 indicates numerous losses of functional different *Ty* elements in various strains suggesting that *Ty*
226 extinction within clades is probably not uncommon and that near complete loss of all traces of extinct
227 elements can occur relatively rapidly (see for example *Ty4* and *Ty5* in **Fig. 4**).

228 In conclusion, we suggest that the divergence of the East Asian Clade IX Complex occurred prior
229 to the genetically close-knit, global radiation of *S. cerevisiae* strains we see today, potentially before their
230 domestication. This begs the question whether there are truly wild *S. cerevisiae* strains outside of Asia at
231 all, especially if the colonization of the rest of the world happened contemporarily with humans. Overall,
232 this study generates new, valuable genomic resources and expands our understanding of the genetic
233 variation and evolutionary history of one of the most important organisms in human history, *S. cerevisiae*.

234 Moreover this set of high-quality genomes, encompassing both domesticated and wild populations from
235 different ecological backgrounds, provides an important resource for future explorations into the
236 dynamics that govern eukaryotic genome evolution.

237

238 **Methods**

239 *Yeast strain origins*

240 We selected eight *Saccharomyces cerevisiae* strains for long-read sequencing and genome assembly
241 (**Table 1**). Seven of these strains originate from East Asia. Six strains were isolated in China (Duan, et al.
242 2018; Wang, et al. 2012) from a variety of ecological niches and one in Taiwan (Peter, et al. 2018). The six
243 Chinese strains cover many of the lineages (CHN I, II, IV, VI, and IX) previously shown to be highly divergent
244 from other *S. cerevisiae* strains based on short-read sequencing. The final strain (Y55) is a common
245 laboratory strain isolated in France with a known mosaic genomic background originating in West Africa.
246 To place our analyses in context, we also included currently publicly available *Saccharomyces sensu stricto*
247 long read genome assemblies as well as assemblies from *Torulaspota delbrueckii* and *Kluyveromyces lactis*
248 (**Table S3**).

249

250 *DNA preparation and long-read sequencing*

251 Before sequencing, strains were sporulated and tetrads were dissected to allow for autodiploidization,
252 making strains homozygous across all loci. Strains were incubated at 30°C in 5ml YEPD (1% yeast extract,
253 2% peptone, 2% dextrose) in a shaking incubator for 24h before we harvested cells by centrifugation. We
254 extracted genomic DNA using NucleoSpin Microbial DNA extraction kit according to the manufacturer's
255 instructions (Macherey-Nagel). Genomic DNA for strain Y55 was extracted independently using the
256 QIAGEN Blood & Culture DNA Midi Kit. Samples were sequenced on PacBio Sequel and Sequel II platforms
257 at the NGI/Uppsala Genome Center (Science for Life Laboratory, Sweden) and the University of Minnesota

258 Sequencing Center (USA). In addition to these PacBio data, we also used publicly available paired-end
259 Illumina sequence data previously generated for each strain (**Table S1**).

260

261 *Genome assembly and annotation*

262 Nuclear contigs were assembled with Flye v2.8.1 (**Fig. S1**, default settings, est. genome size = 12.4Mbp)
263 (Kolmogorov, et al. 2019). We used short read sequences for each strain to error-correct the long reads
264 using FMLRC v2 (Wang, et al. 2018). Corrected long reads and the short reads were subsequently used to
265 polish the Flye assemblies using Racon v1.4.13 (Vaser, et al. 2017) and POLCA v3.4.2 (Zimin and Salzberg
266 2019) respectively. We further scaffolded the contigs based on the reference S288C genome
267 (GCA_000146045.2) using RaGOO v1.1 (Alonge, et al. 2019) and filled any gaps this generated using
268 multiple iterations of LR Gapcloser v1 (Xu, et al. 2019) and Gapcloser (Luo, et al. 2012). To account for any
269 errors introduced by using long reads to fill gaps, we further polished each assembly once more using
270 Racon v1.4.13 and POLCA v3.4.2. Mitochondrial assemblies were largely assembled using Flye without the
271 assumption of even coverage (--metagenomic) using all long reads as input. JXXY16.1 and Y55
272 mitochondrial genomes were assembled using Flye v2.8.1 with default settings. Mitochondrial contigs
273 were extracted by mapping the Flye output to the reference mitochondrial genome using Nucmer
274 (Delcher, et al. 2002). These assemblies were polished and scaffolded following the same process as that
275 of the nuclear assemblies. Completeness of the final genome assemblies was assessed using BUSCO v4.0.5
276 (Simão, et al. 2015; Waterhouse, et al. 2018).

277 We annotated nuclear genes, mitochondrial genes, centromeres, transposable elements, core X
278 elements and Y-prime elements using modified versions of the pipelines within the LRSDAY package (Yue
279 and Liti 2018). In addition to our eight newly assembled genomes, we also used the same method to
280 annotate the previously published long read assemblies (**Table S3**). Nuclear genes orthologous to
281 annotated genes in the *S. cerevisiae* S288C reference genome were identified using Proteinortho v6.0.24

282 (Lechner, et al. 2011). Genes for which no orthologous protein was found in the reference were clustered
283 based on orthology to each other.

284 To further characterize *Ty* elements, we determined potential element viability by translating
285 coding regions of full elements based on reading frames identified for each element in S288C. Elements
286 containing premature stop codons or extensive frameshifts were categorized as putatively being
287 reproductively inviable (loss-of-function or LOF). Additionally, we created gene trees for whole elements
288 of each *Ty* class using MAFFT v7.471 alignments (default settings) with PhyML v3.0 (substitution model =
289 HKY85; bootstrap = 100; tree searching using SRT and NNI; conducted in Unipro UGENE v36.0). To
290 determine the likelihood of closely related elements within a given strain resulting from transposition or
291 segmental genome duplication, we mapped the 10000bp regions containing each element to related
292 intra-strain elements.

293

294 *Phylogenomic analysis*

295 To place our eight assembled genomes within the context of other *Saccharomyces* strains, we employed
296 both a consensus gene tree and Assembly and Alignment-Free (AAF) approaches to phylogenetic tree
297 construction. For consensus species trees, we used OrthoFinder v2.4.0 (**Fig. 1**) in addition to a standard
298 gene tree approach. For the latter, we aligned all orthologous genes found in at least 5 strains (5847 genes)
299 using MUSCLE v3.8.31 (Edgar 2004) and performed maximum-likelihood single-tree inference for each
300 locus using RAxML-NG v1.01 (Kozlov et al. 2019) with a discrete GAMMA model of rate heterogeneity. We
301 used Astral-III v5.7.4 (Zhang et al. 2018) with these gene trees to generate a consensus species tree.

302 AAF v20171001 (Fan, et al. 2015) was used with a k-mer size of 20 nucleotides and a threshold
303 frequency of 7 for each k-mer to be included in the analysis. AAF was used to compare the long-read and
304 short-read sequencing data for the 25 *Saccharomyces* strains (**Fig. S2, Table S3**). Short-read sequencing
305 data for *K. lactis* and *T. delbrueckii* were included as outgroups.

306 To generate the mitochondrial phylogeny, we reoriented the start of each assembly based on the
307 position of the tRNA gene, trnP(ugg), then aligned these assemblies to each other using Mugsy v.1.2.3
308 (Angiuoli and Salzberg 2011). This multiple sequence alignment was then used to create a maximum
309 likelihood tree using IQ-TREE v2.0.5 (options: -m TPM2u+F+R3 -B 1000 -bnni -alrt 1000) (Hoang, et al.
310 2018; Nguyen, et al. 2015). The model was determined using the ModelFinder component of IQ-Tree
311 (Kalyaanamoorthy, et al. 2017).

312

313 *Structural variation detection*

314 To identify the structural variations (SVs) between strains within *S. cerevisiae*, we performed exhaustive
315 pairwise comparisons between the 15 strains with long-read assemblies (210 comparisons). We focused
316 on five types of SV: deletions, insertions, tandem duplications, inversions and translocations. The SVs were
317 detected using MUM&Co (O'Donnell and Fischer 2020), which utilizes MUMmer v3.23 (Marçais, et al.
318 2018) to perform whole-genome alignments and detect SVs \geq 50bps.

319

320 *Spore viability assay*

321 To assess the level of reproductive isolation between the divergent East Asian strains and modern *S.*
322 *cerevisiae*, we crossed all strains with Y55 (α ; ho; leu2 Δ ::HygMX) (and Y55 to itself) and assessed the spore
323 viability of each cross. We sporulated each strain by incubating them in liquid sporulation medium (KAC;
324 2% potassium acetate) for 3 days at 23C. These cultures were then incubated with 10 μ l zymolyase (100
325 U/ml) at 37C for 30min before being plated on YEPD (2.5% agar) in equal mixture with cultures of Y55 (α ;
326 ho; leu2 Δ ::HygMX) and grown for 48h at 30C. This culture was streaked on YEPD + hygromycin and replica
327 plated to minimal media. A single colony was selected from each cross and grown up in liquid YEPD
328 overnight, spun down, put in KAC, and incubated at room temperature with shaking for four days to
329 induce sporulation. The resulting tetrads were treated with zymolyase for 30min at room temperature.

330 500µl of sterile water were added before spores were dissected out of the tetrads onto YPD plates, using
331 a Singer MSM 400 micro-manipulator. We dissected 20 tetrads yielding 80 spores per cross. Plates were
332 incubated at 30C and colonies were counted after 72h, indicating viable spores that were able to
333 germinate.

334 Respiratory competence was determined by plating strains of yeast on rich media containing
335 nonfermentable glycerol as the sole carbon source (1% yeast extract, 2% peptone, 2% glycerol).

336

337 ***Data Availability Statement***

338 The data underlying this article are available in the European Nucleotide Archive and can be accessed with
339 accession number **PRJEB38713**.

340

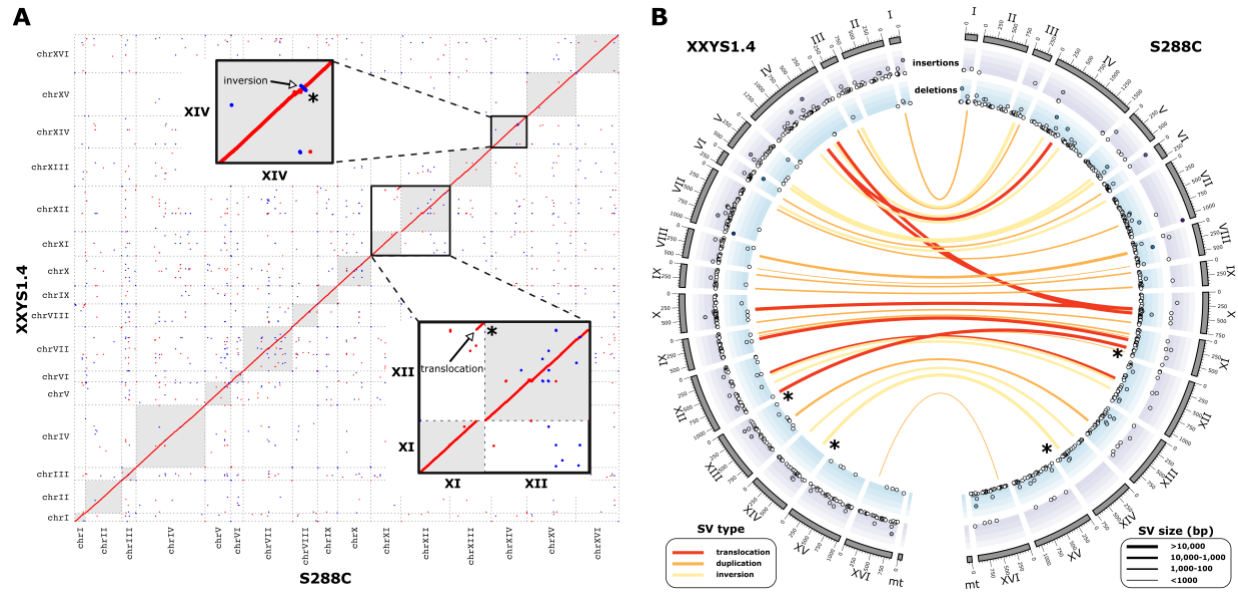
341 ***Acknowledgments***

342 This work was supported by grants from Stockholm University (Science for Life Laboratory sequencing
343 grant SU FV-2.2.2-1843-17 to RS), the Swedish Research Council (2017-04963 to RS), the Knut and Alice
344 Wallenberg Foundation (2017.0163 to RS), the Carl Trygger foundation (CTS 17: 431 to ZZ), the Wenner-
345 Gren Foundations (UPD2018-0196, UPD2019-0110 to DPB), and The University of Minnesota Department
346 of Ecology, Evolution, and Behavior. We acknowledge the support of the National Genomics Infrastructure
347 (NGI)/Uppsala Genome Center and UPPMAX for providing assistance in massive parallel sequencing and
348 computational infrastructure. Work performed at NGI/Uppsala Genome Center (project SNIC 2019/8-23)
349 has been funded by RFI/VR and Science for Life Laboratory, Sweden. We would like to thank Feng-Yan Bai
350 and Gianni Liti for donating strains, Jia-Xing Yue for advice on the LRSDAY pipeline, and Gianni Liti, Samuel
351 O'Donnell and Chris Wheat for discussion.

352 **Table 1: Descriptions of the *Saccharomyces cerevisiae* strains sequenced in this study.**

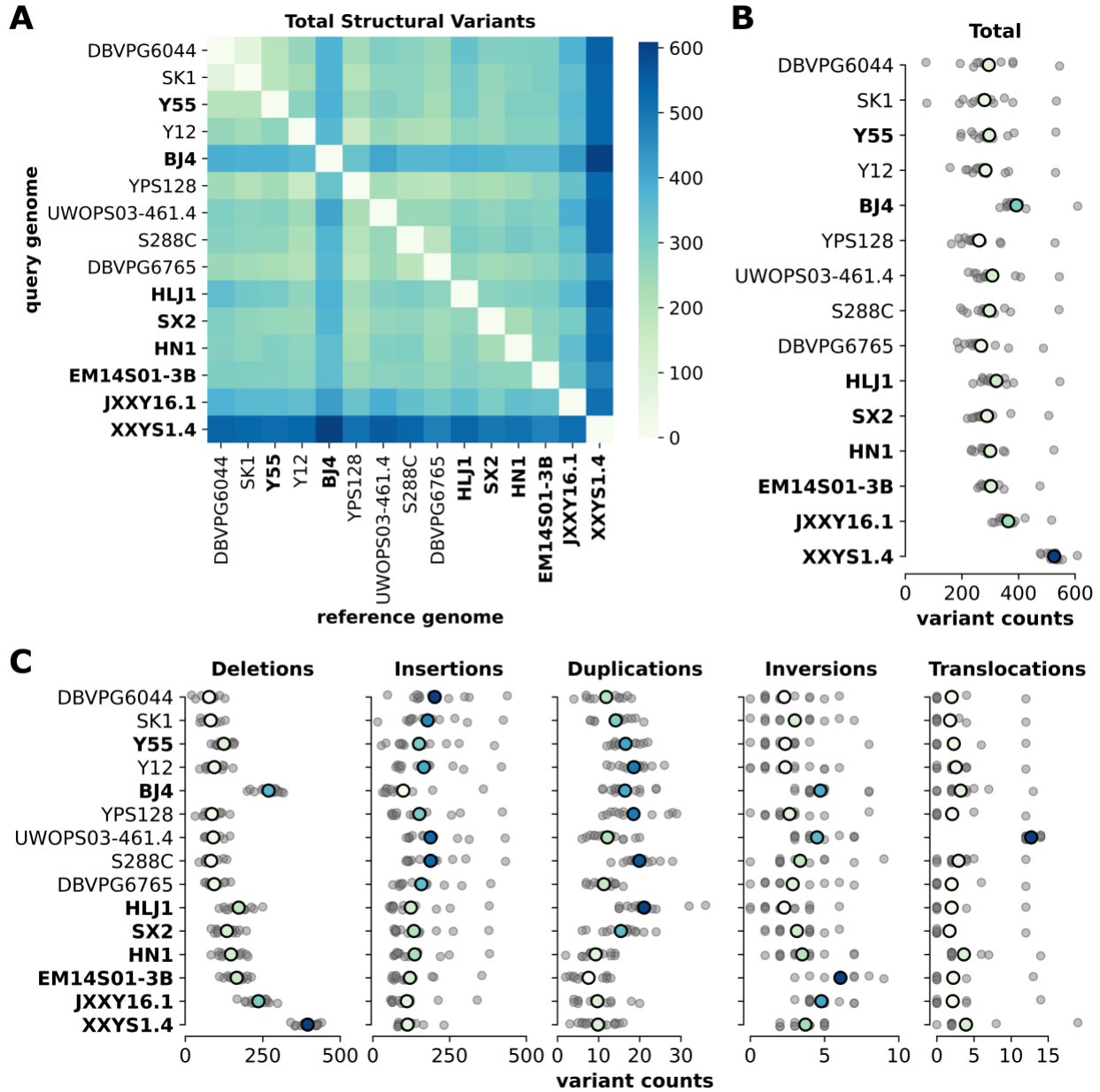
Lineage	Strain	Species	Source	Geographic Location
CHN I	HN1	<i>S. cerevisiae</i>	rotten wood, primeval forest	Diaoluo Mountain, Hainan, China
CHN II	SX2	<i>S. cerevisiae</i>	bark of a Fagaceae tree, primeval forest	Qinling Mountain, Shaanxi, China
CHN IV	HLJ1	<i>S. cerevisiae</i>	bark of <i>Quercus mongolica</i> , secondary forest	Jingbo lake, Heilongjiang, China
CHN VI	BJ4	<i>S. cerevisiae</i>	intestine of a butterfly, park	Haidian, Beijing, China
CHN IX	JXXY16.1	<i>S. cerevisiae</i>	bark, primeval forest	Xiangxiyuan, Hubei Province, China
CHN IX	XXYS1.4	<i>S. cerevisiae</i>	bark, primeval forest	Xiangxiyuan, Hubei Province, China
Taiwanese	EM14S01-3B	<i>S. cerevisiae</i>	soil	Taiwan
West African+	Y55	<i>S. cerevisiae</i>	wine grapes	France

353

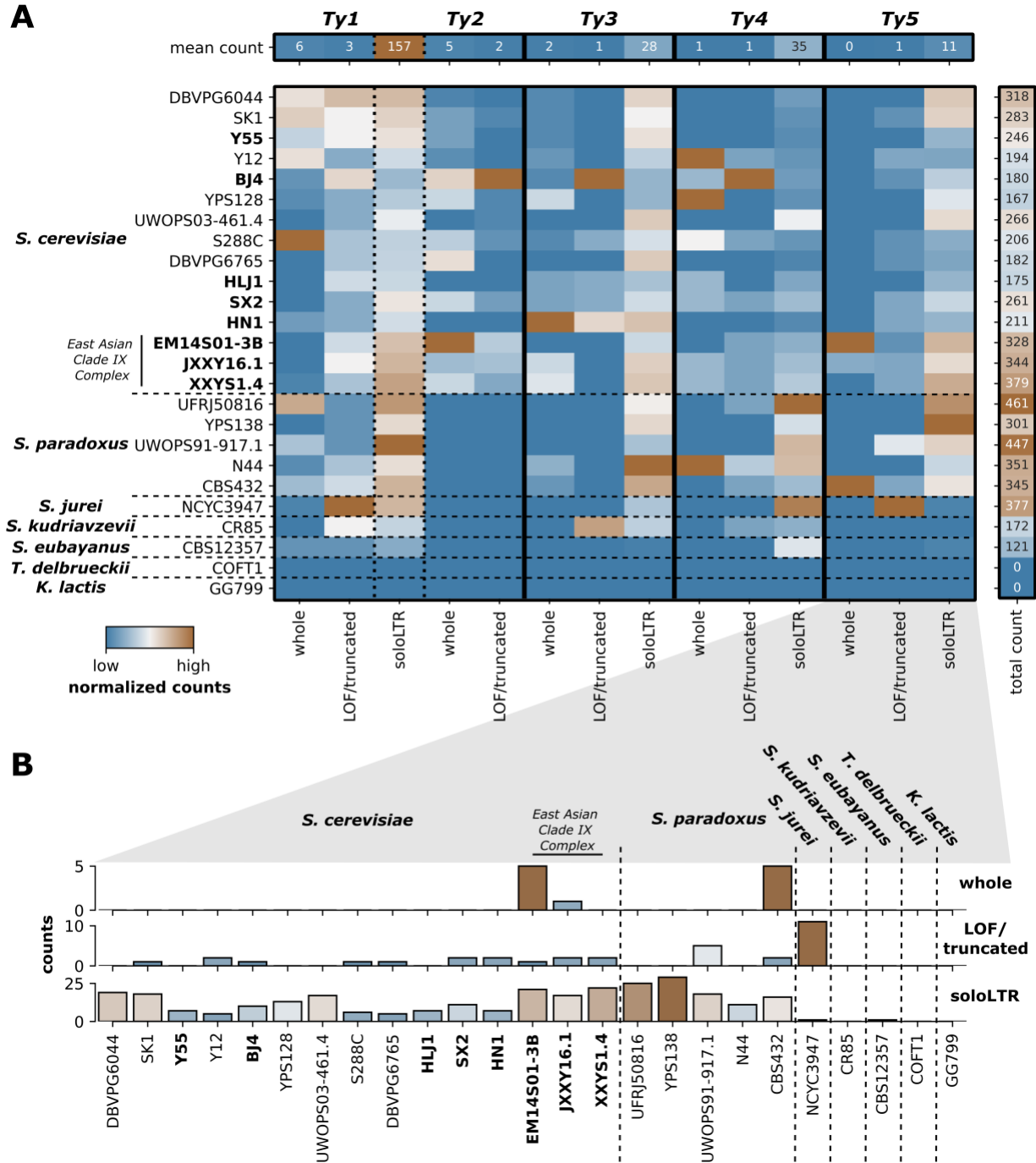


363

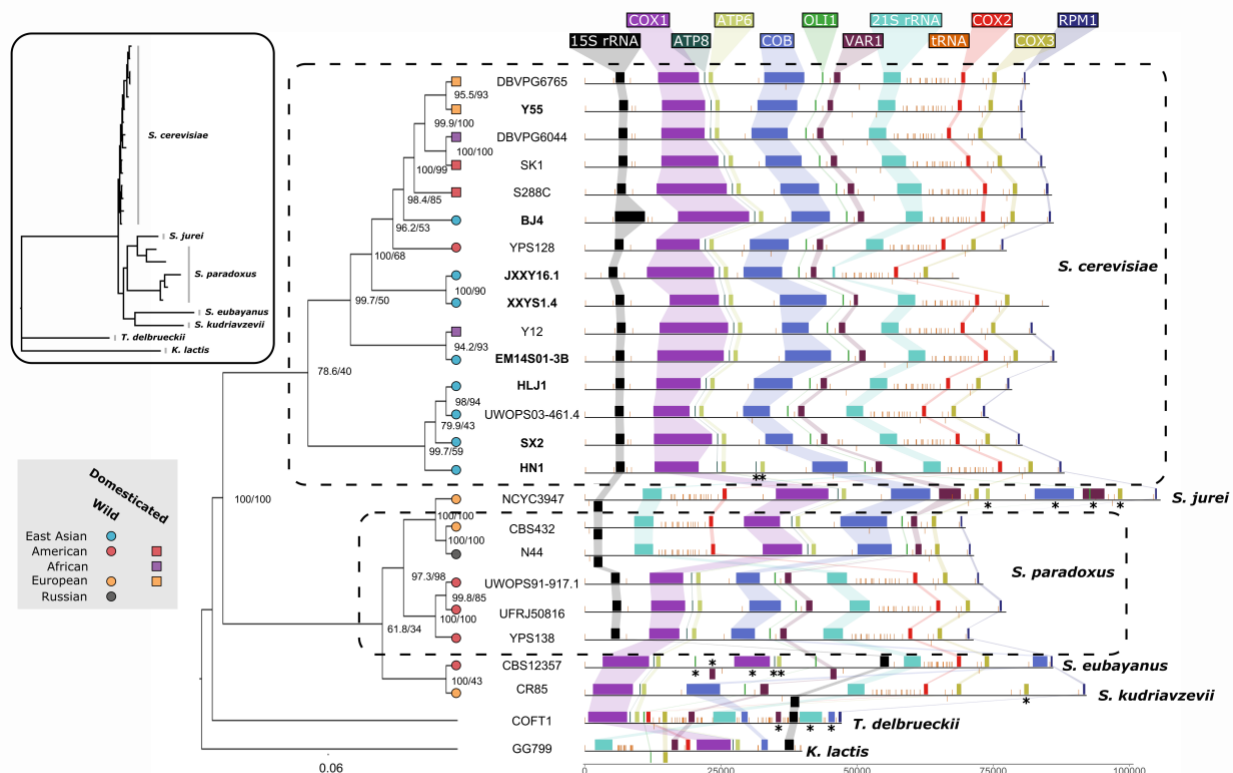
364 **Figure 2: *Saccharomyces cerevisiae* long-read PacBio genome assemblies** (A) Genome comparison of the
365 reference strain, S288C, and a member of the East Asian Clade IX Complex: XXYS1.4. Sequence homology
366 within the dot plots is indicated by red dots for forward matches and blue dots for reverse matches. Insets
367 depict examples of deviations from homology: 1) a large translocation between XI and chromosome XII
368 found conserved in JXXY16.1, XXYS1.4 and EM14SO1-3B and 2) a large inversion in chromosome XIV. (B)
369 CIRCOS plot showing the detected structural variations between reference strain, S288C and XXYS1.4.
370 Translocations (red), duplications (orange) and inversions (yellow) are depicted as links between the two
371 genomes. The width of the link reflects the relative size of the variation (bp). The translocation and
372 inversion depicted in panel A are highlighted with asterisks. Insertions (blue) and deletions (purple) are
373 depicted in the outer tracks. Deletion and insertion size increase towards the outside. Chromosome size
374 is shown on the outside in 1kb units.



375
 376 **Figure 3: Structural variations within *Saccharomyces cerevisiae*.** (A) Pairwise comparisons among all
 377 *Saccharomyces cerevisiae* genome assemblies with the total number of variations. Order of genome
 378 assemblies is consistent with the species tree (Fig. 1). New long-read genome assemblies presented in this
 379 study are bold. (B) The range of total structural variation counts found for each genome serving as
 380 reference genome. Grey dots indicate each pairwise genome comparison. Colored dots indicate the mean
 381 and are colored on a relative scale. (C) The range of structural variation counts for each type of variation.
 382 Grey dots indicate each pairwise genome comparison. Colored dots indicate the mean and are colored on
 383 a relative scale. Corresponding heatmaps for pairwise comparison are shown in Fig. S13 to S17.



384
 385 **Figure 4: Transposable element composition in *Saccharomyces*** (A) Transposable element composition
 386 in total count subdivided by Ty classification for *Saccharomyces sensu stricto* strains. For visual
 387 comparison, each column is normalized (x/x_{max}) for that specific element. For raw values see **Fig. S21**. LTR
 388 = long terminal repeat components of Ty elements without replicative machinery. (B) A closer look at Ty5
 389 elements across *Saccharomyces*.



390
 391 **Figure 5: Mitochondrial phylogenetics and genomic arrangements.** Phylogenetic tree based on
 392 mitochondrial genomic content. Internal branches are labeled with bootstrap support. *Saccharomyces*
 393 strains are colored according to their location of origin and branch tip shape indicates whether it is a
 394 domesticated (square) or wild (circle) strain. New long-read genome assemblies presented in this study
 395 are indicated in bold. Major genomic elements found on mitochondria are shown and colored according
 396 to guide elements at the top. Inverted elements appear on the underside of the line. Duplicated elements
 397 are indicated with asterisk. Inset depicts untransformed phylogenetic tree with species labeled.

398 **References**

- 399 Alonge M, et al. 2019. RaGOO: fast and accurate reference-guided scaffolding of draft genomes.
400 *Genome biology* 20: 1-17.
- 401 Angiuoli SV, Salzberg SL 2011. Mugsy: fast multiple alignment of closely related whole genomes.
402 *Bioinformatics* 27: 334-342.
- 403 Carr M, Bensasson D, Bergman CM 2012. Evolutionary genomics of transposable elements in
404 *Saccharomyces cerevisiae*. *PloS one* 7.
- 405 De Chiara M, et al. 2020. Discordant evolution of mitochondrial and nuclear yeast genomes at
406 population level. *BMC biology* 18: 1-15.
- 407 Delcher AL, Phillippy A, Carlton J, Salzberg SL 2002. Fast algorithms for large-scale genome
408 alignment and comparison. *Nucleic acids research* 30: 2478-2483.
- 409 Duan S-F, et al. 2018. The origin and adaptive evolution of domesticated populations of yeast
410 from Far East Asia. *Nature communications* 9: 1-13.
- 411 Fan H, Ives AR, Surget-Groba Y, Cannon CH 2015. An assembly and alignment-free method of
412 phylogeny reconstruction from next-generation sequencing data. *BMC genomics* 16: 522.
- 413 Goodwin S, McPherson JD, McCombie WR 2016. Coming of age: ten years of next-generation
414 sequencing technologies. *Nature Reviews Genetics* 17: 333.
- 415 Hoang DT, Chernomor O, Von Haeseler A, Minh BQ, Vinh LS 2018. UFBoot2: improving the
416 ultrafast bootstrap approximation. *Molecular biology and evolution* 35: 518-522.
- 417 Hyma KE, Fay JC 2013. Mixing of vineyard and oak-tree ecotypes of *Saccharomyces cerevisiae* in
418 North American vineyards. *Molecular Ecology* 22: 2917-2930.
- 419 Kalyaanamoorthy S, Minh BQ, Wong TK, von Haeseler A, Jermiin LS 2017. ModelFinder: fast
420 model selection for accurate phylogenetic estimates. *Nature methods* 14: 587.
- 421 Kolmogorov M, Yuan J, Lin Y, Pevzner PA 2019. Assembly of long, error-prone reads using
422 repeat graphs. *Nature biotechnology* 37: 540-546.
- 423 Lechner M, et al. 2011. Proteinortho: detection of (co-) orthologs in large-scale analysis. *BMC*
424 *bioinformatics* 12: 124.
- 425 Leducq JB, et al. 2016. Speciation driven by hybridization and chromosomal plasticity in a wild
426 yeast. *Nat Microbiol* 1: 15003. doi: 10.1038/nmicrobiol.2015.3
- 427 Liti G 2015. The fascinating and secret wild life of the budding yeast *S. cerevisiae*. *Elife* 4:e0585.
- 428 Liti G, et al. 2009. Population genomics of domestic and wild yeasts. *Nature* 458: 337-341.
- 429 Liti G, Peruffo A, James SA, Roberts IN, Louis EJ 2005. Inferences of evolutionary relationships
430 from a population survey of LTR-retrotransposons and telomeric-associated sequences in the
431 *Saccharomyces sensu stricto* complex. *Yeast* 22: 177-192.
- 432 Luo R, et al. 2012. SOAPdenovo2: an empirically improved memory-efficient short-read de novo
433 assembler. *GigaScience* 1: 2047-2217X-2041-2018.
- 434 Marçais G, et al. 2018. MUMmer4: A fast and versatile genome alignment system. *PLoS*
435 *computational biology* 14: e1005944.
- 436 McGovern PE, et al. 2004. Fermented beverages of pre-and proto-historic China. *Proceedings of*
437 *the National Academy of Sciences* 101: 17593-17598.
- 438 Merker JD, et al. 2018. Long-read genome sequencing identifies causal structural variation in a
439 Mendelian disease. *Genetics in Medicine* 20: 159-163.

- 440 Naseeb S, et al. 2018. Whole genome sequencing, de novo assembly and phenotypic profiling
441 for the new budding yeast species *Saccharomyces jurei*. *G3: Genes, Genomes, Genetics* 8: 2967-
442 2977.
- 443 Nguyen L-T, Schmidt HA, Von Haeseler A, Minh BQ 2015. IQ-TREE: a fast and effective stochastic
444 algorithm for estimating maximum-likelihood phylogenies. *Molecular biology and evolution* 32:
445 268-274.
- 446 O'Donnell S, Fischer G 2020. MUM&Co: accurate detection of all SV types through whole-
447 genome alignment. *Bioinformatics* 36: 3242-3243.
- 448 Payen C, et al. 2014. The dynamics of diverse segmental amplifications in populations of
449 *Saccharomyces cerevisiae* adapting to strong selection. *G3: Genes, Genomes, Genetics* 4: 399-
450 409.
- 451 Peris D, et al. 2017. Mitochondrial introgression suggests extensive ancestral hybridization
452 events among *Saccharomyces* species. *Molecular phylogenetics and evolution* 108: 49-60.
- 453 Peter J, et al. 2018. Genome evolution across 1,011 *Saccharomyces cerevisiae* isolates. *Nature*
454 556: 339.
- 455 Scannell DR, et al. 2011. The Awesome Power of Yeast Evolutionary Genetics: New Genome
456 Sequences and Strain Resources for the *Saccharomyces sensu stricto* Genus. *G3 (Bethesda)* 1:
457 11-25. doi: 10.1534/g3.111.000273
- 458 Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM 2015. BUSCO: assessing
459 genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics* 31:
460 3210-3212.
- 461 Steenwyk JL, Rokas A 2018. Copy number variation in fungi and its implications for wine yeast
462 genetic diversity and adaptation. *Frontiers in microbiology* 9: 288.
- 463 Vaser R, Sović I, Nagarajan N, Šikić M 2017. Fast and accurate de novo genome assembly from
464 long uncorrected reads. *Genome research* 27: 737-746.
- 465 Wang JR, Holt J, McMillan L, Jones CD 2018. FMLRC: Hybrid long read error correction using an
466 FM-index. *BMC bioinformatics* 19: 50.
- 467 Wang QM, Liu WQ, Liti G, Wang SA, Bai FY 2012. Surprisingly diverged populations of
468 *Saccharomyces cerevisiae* in natural environments remote from human activity. *Molecular*
469 *Ecology* 21: 5404-5417.
- 470 Waterhouse RM, et al. 2018. BUSCO applications from quality assessments to gene prediction
471 and phylogenomics. *Molecular biology and evolution* 35: 543-548.
- 472 Wellenreuther M, Mérot C, Berdan E, Bernatchez L 2019. Going beyond SNPs: the role of
473 structural genomic variants in adaptive evolution and species diversification. *Molecular Ecology*
474 28: 1203-1209.
- 475 Xu G-C, et al. 2019. LR_Gapcloser: a tiling path-based gap closer that uses long reads to
476 complete genome assembly. *GigaScience* 8: giy157.
- 477 Yue J-X, et al. 2017. Contrasting evolutionary genome dynamics between domesticated and
478 wild yeasts. *Nature genetics* 49: 913.
- 479 Yue J-X, Liti G 2018. Long-read sequencing data analysis for yeasts. *Nature protocols* 13: 1213.
- 480 Zhang Z, et al. 2020. Recombining your way out of trouble: The genetic architecture of hybrid
481 fitness under environmental stress. *Molecular biology and evolution* 37: 167-182.
- 482 Zhang Z, Schwartz S, Wagner L, Miller W 2000. A greedy algorithm for aligning DNA sequences.
483 *Journal of Computational biology* 7: 203-214.

484 Zimin AV, Salzberg SL 2019. The genome polishing tool POLCA makes fast and accurate
485 corrections in genome assemblies. bioRxiv.
486