

1 **Chromosome-level Subgenome-aware *de novo* Assembly of**
2 ***Saccharomyces bayanus* Provides Insight into Genome Divergence**
3 **after Hybridization**

4

5

6 Cory Gardner^{1,2,#}, Junhao Chen^{3,#}, Christina Hadfield², Zhaolian Lu³, David Debruin², Yu Zhan³,

7 Maureen J. Donlin^{2,4}, Zhenguo Lin^{2,3*} and Tae-Hyuk Ahn^{1,2*}

8

9 ¹ Department of Computer Science, Saint Louis University, St. Louis, MO, USA

10 ² Program in Bioinformatics and Computational Biology, Saint Louis University, St. Louis, MO, USA

11 ³ Department of Biology, Saint Louis University, St. Louis, MO, USA

12 ⁴ Department of Biochemistry and Molecular Biology, Saint Louis University, St. Louis, MO, USA

13 # These authors contributed equally to this work

14 * To whom correspondence should be addressed

15 Zhenguo Lin: zhengou.lin@slu.edu

16 Tae-Hyuk Ahn: taehyuk.ahn@slu.edu

17

18 **Keywords:** *de novo* genome assembly, genome annotation, yeast, *Saccharomyces bayanus*, hybrid

19 genome, subgenome-aware assembly, LRS Special Issue

20

21

22 **Abstract**

23 Interspecies hybridization is prevalent in various eukaryotic lineages and plays important roles in
24 phenotypic diversification, adaption, and speciation. To better understand the changes that occurred in the
25 different subgenomes of a hybrid species and how they facilitated adaptation, we completed chromosome-
26 level *de novo* assemblies of all 16 pairs chromosomes for a recently formed hybrid yeast, *Saccharomyces*
27 *bayanus* strain CBS380 (IFO11022), using Nanopore MinION long-read sequencing. Characterization of
28 *S. bayanus* subgenomes and comparative analysis with the genomes of its parent species, *S. uvarum* and *S.*
29 *eubayanus*, provide several new insights into understanding genome evolution after a relatively recent
30 hybridization. For instance, multiple recombination events between the two subgenomes have been
31 observed in each chromosome, followed by loss of heterozygosity (LOH) in most chromosomes in nine
32 chromosome pairs. In addition to maintaining nearly all gene content and synteny from its parental
33 genomes, *S. bayanus* has acquired many genes from other yeast species, primarily through the introgression
34 of *S. cerevisiae*, such as those involved in the maltose metabolism. In addition, the patterns of recombination
35 and LOH suggest an allotetraploid origin of *S. bayanus*. The gene acquisition and rapid LOH in the hybrid
36 genome probably facilitated its adaption to maltose brewing environments and mitigated the maladaptive
37 effect of hybridization.

38

39

40 **Introduction**

41 It has generally been believed that hybridization between closely related species often leads to inviability
42 and sterility, a phenomenon known as hybrid incompatibility. The Dobzhansky-Muller (DM) model, which
43 proposes that it results from negative epistatic interactions between genes with different evolutionary
44 histories, is a well-regarded explanation for hybrid incompatibility (Dobzhansky 1982; Price et al. 2010).
45 Hybrid incompatibility can act as a reproductive isolating barrier contributing to speciation (Coyne and Orr
46 2004). Additionally, reduced fertility in hybrids can result from abnormal chromosome segregation during
47 meiosis if the parental genomes are divergent (Coyne and Orr 2004). Nevertheless, recent studies show that
48 interspecies hybridization is prevalent in major eukaryotic lineages, particularly in angiosperms and yeasts,
49 and it is believed to contribute to adaptation to novel environments (Langdon et al. 2019; Taylor and Larson
50 2019; Gabaldon 2020; Moran et al. 2021; Suvorov et al. 2022). Given that the exchange of genomic content
51 between species is pervasive, it is important to better characterize the impact of hybridization on evolution
52 of hybrid genomes, which will improve our understanding of the genetic basis underlying the adaptation
53 and divergence of species.

54 The *Saccharomyces* budding yeast species involved in fermentation of various products is a group of
55 organisms in which hybrids are most commonly found (Langdon et al. 2019; Gabaldon 2020). The
56 allopolyploid genome of *Saccharomyces cerevisiae* has been extensively studied. The ancestral
57 *Saccharomyces* lineage experienced a whole genome duplication (WGD) about 100 million years ago
58 (Wolfe and Shields 1997; Kellis et al. 2004). New evidence suggests that the WGD in the *Saccharomyces*
59 lineage was caused by interspecies hybridization (Marcet-Houben and Gabaldon 2015). Soon after the
60 WGD, there was a period of rapid losses of duplicate genes and only ~10% of WGD ohnologs survived.
61 The retained WGD duplicates are enriched in genes related to glucose metabolism or rapid growth, such as
62 glycolysis genes (Conant and Wolfe 2007), hexose transporters (Lin and Li 2011), and ribosomal protein
63 genes (Mullis et al. 2019). These studies suggested that the WGD or hybridization event played a significant
64 role in the adaptation of *Saccharomyces* species toward aerobic fermentation (Kellis et al. 2004; Thomson

65 et al. 2005; Conant and Wolfe 2007; Lin and Li 2014) and speciation events (Scannell et al. 2006). These
66 studies improved our understanding of the biological significance of interspecies hybridization in speciation
67 and adaptation.

68 At the genomic level, questions related to what occurred to the genome after a recent allopolyploidy
69 event, such as the earliest genome rearrangements, the mechanisms of gene loss, recombination between
70 subgenomes, and loss of heterozygosity, are not completely understood (Morales and Dujon 2012). The
71 ancient hybridization events, such as the WGD in the ancestral *Saccharomyces* lineage, may not be useful
72 to address these questions as most duplicate genes have been lost. In addition to the ancient hybridization
73 event, recent interspecific hybridization is prevalent in the *Saccharomyces* lineage as they are used to
74 produce fermented beverages (Langdon et al. 2019). The genomes of these recently generated hybrid
75 genomes may serve as ideal systems to study how genomes evolve after hybridization and contributed to
76 adaptation to specific niches. For instance, *S. pastorianus*, which is an interspecies hybrid between *S.*
77 *cerevisiae* and *S. eubayanus*, is widely used for brewing lager style beers under low temperature in Europe
78 (Libkind et al. 2011). Some chromosomes in *S. pastorianus* strains may have 5 copies, suggesting its highly
79 aneuploid nature (van den Broek et al. 2015; Gorter de Vries et al. 2017). The chromosome-level assembly
80 for *S. pastorianus* strain CBS 1483, based on MinION long-read sequencing, enables the assembly and
81 exploration of the unstable subtelomeric regions, which contain industrially-relevant genes such as the
82 MAL genes (Salazar et al. 2019).

83 *Saccharomyces bayanus* is another interspecies hybrid yeast commonly found in industrial brewing
84 environments, but it is viewed as a contaminant in some brewing processes due to the production of
85 undesired byproducts (Rainieri et al. 2003). The taxonomic classification of *S. bayanus* has been a
86 controversial process (Hittinger 2013). Thanks to the discovery of a wild species *S. eubayanus* (Libkind et
87 al. 2011), it is now commonly accepted that *S. bayanus* is a hybrid between *S. uvarum*, and *S. eubayanus*
88 (Perez-Traves et al. 2014; Peris et al. 2014). *S. bayanus* isolates are highly heterogeneous in genetic and
89 metabolic characteristics, probably resulting from many independent hybridization events between *S.*
90 *eubayanus* and *S. uvarum*, creating many different strains (Rainieri et al. 2006; Libkind et al. 2011; Langdon

91 et al. 2019). Genome sequencing using Illumina has been carried out for over 40 *S. bayanus* strains, such
92 as CBS 380, NCAIM 676, FM1309 and NBRC1948 (Libkind et al. 2011; Almeida et al. 2014; Langdon et
93 al. 2019). Mapping Illumina reads to different *Saccharomyces* species showed that the contributions of
94 genome content from *S. uvarum* and *S. eubayanus* are highly variable among *S. bayanus* strains.
95 Specifically, the genome content deriving from *S. uvarum* ranges from 36.6% to 98.8% (Langdon et al.
96 2019). In addition, small introgressed regions from *S. cerevisiae* are present in some *S. bayanus* strains
97 (Nguyen et al. 2011). However, due to the limitation of Illumina short reads, these *S. bayanus* genome
98 assemblies are fragmented.

99 A chromosomal-level subgenome assembly of *S. bayanus* is expected to provide much more detail in
100 the genome evolution following a recent allopolyploidy event. In this study, we sequenced the genome of
101 *S. bayanus* strain CBS 380 (BY20106, IFO11022) using the Nanopore MinION. The strain CBS 380 is the
102 most representative isolate of *S. bayanus*, which has been widely used in many studies (Libkind et al. 2011;
103 Nguyen et al. 2011; Caudy et al. 2013; Perez-Traves et al. 2014). We generated chromosome-level
104 subgenome assemblies based on MinION reads and characterized the evolution of genome structure and
105 gene content. Our results show that *S. uvarum* contributed to over 60% of the hybrid genome. Many
106 chromosomes exhibit mosaic segments of different origins, suggesting multiple recombination events
107 occurred between the two subgenomes after hybridization. Rapid loss of heterozygosity (LOH) following
108 hybridization was also observed, resulting in over 56% of the genome regions becoming homozygous.
109 Introgression from a third species, *S. cerevisiae*, was also detected, contributing to the expansion of maltose
110 metabolism genes in the *S. bayanus* genome. These observations provide detailed examples illustrating how
111 genome evolved immediately after hybridization occurred, improving our understanding of genetic basis
112 of a hybrid species' survival by overcoming hybrid incompatibility.

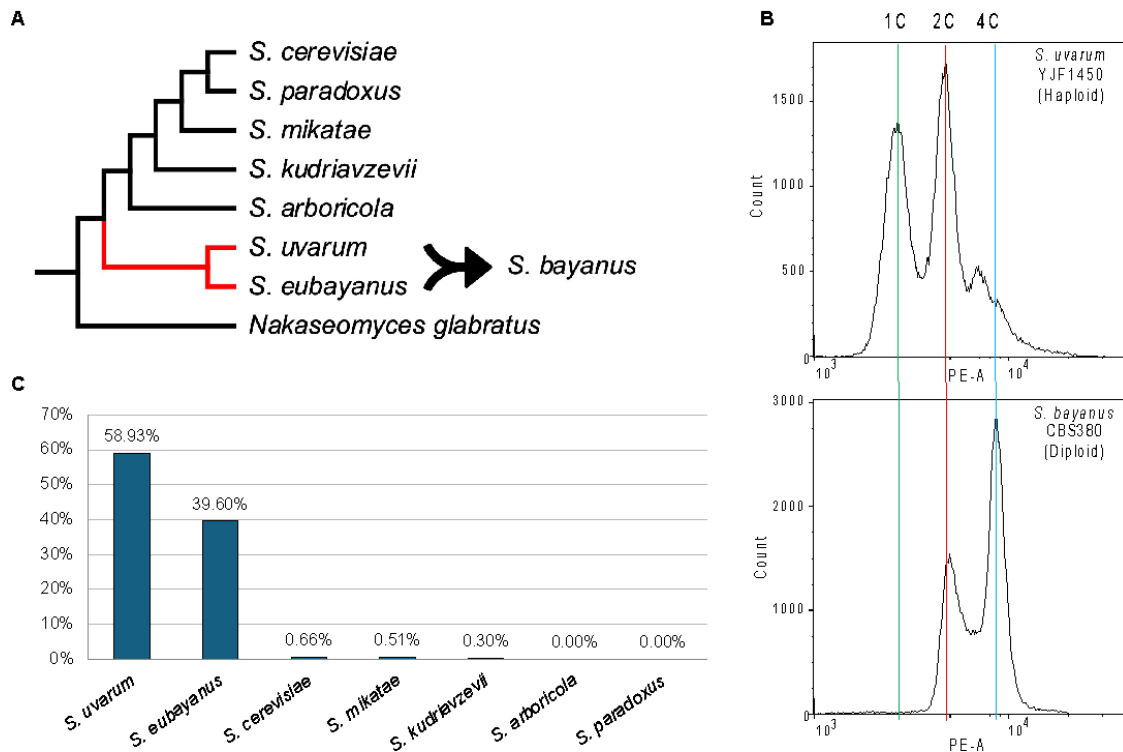
113

114 **Results**

115 **MinION sequencing, ploidy analysis, and parental inference of *S. bayanus* CBS 380 genome**

116 It is well accepted that *S. bayanus* arose from interspecies hybridization between two closely related
117 *Saccharomyces sensu stricto* yeast species *S. uvarum* and *S. eubayanus* (Figure 1A). To confirm the ploidy
118 levels of the *S. bayanus* CBS 380 strain, we assessed its relative genomic DNA content by fluorescence
119 flow cytometry analysis using a haploid yeast strain *S. bayanus* YJF1450 as a control (Figure 1B). Dual
120 peaks of fluorescence were observed in both strains, with the first peak indicating the DNA content of G1
121 phase and the second peak showing DNA content after DNA synthesis (G2/M phase). As shown in Fig 1B,
122 the relative genomic DNA content in G1 phase of *S. bayanus* CBS 380 is similar to the G2/M phase the
123 haploid control *S. bayanus* YJF1450, confirming that two sets of chromosomes are present in *S. bayanus*
124 CBS 380.

125 Sequencing of *S. bayanus* CBS 380 with Oxford Nanopore's MinION yielded 2.2 gigabase pairs
126 (Gb) of data (~170x coverage), with 2.04 Gb passing quality control (Supplemental Figure S1). Among
127 these, 100 reads exceeded 100 kilobase pairs (Kb) with the longest extending to 158,255 base pairs (bp).
128 We hypothesized that most of our reads would map to the suspected parental species, *S. eubayanus* and *S.*
129 *uvarum*, while fewer, if any, reads would map to the other more distantly related species. We used sppIDer
130 (Langdon et al. 2018), which maps sequencing reads to the reference genomes of multiple species of
131 interest, to validate the strain sequenced as a hybrid of *S. eubayanus* and *S. uvarum*, and to determine the
132 relative genetic contribution by each parent. The genomes of *S. uvarum* (CBS 7001), *S. eubayanus* (FM
133 1318), *S. cerevisiae* (BY 4742), *S. mikatae* (IFO 1815), *S. kudriavzevii* (IFO1802), *S. arboricola* (ZP960)
134 and *S. paradoxus* (CBS 432), were used as reference genomes for read mapping. The majority of the reads
135 mapped to *S. eubayanus* and *S. uvarum*, as expected, with 39.6% mapping to *S. eubayanus* and 58.93%
136 mapping to *S. uvarum*. Of the remaining, 5.52% mapped to *S. cerevisiae*, 2.67% to *S. mikatae* and 2.14%
137 to *S. paradoxus* (Figure 1C). This confirms the identity of our sequenced strain as *S. bayanus*, a hybrid of
138 *S. eubayanus* and *S. uvarum*.



139

140 **Figure 1.** *S. bayanus* CBS 380 has a diploid genome, resulting from a hybridization event between *S. uvarum* and *S.*
141 *eubayanus*. (A) Schematic illustration of evolutionary relationships among *S. bayanus* and closely related species. (B)
142 Ploidy analysis by flow cytometry of *S. bayanus*. The top histogram shows cell count of *S. uvarum* YJF1450 and the
143 bottom histogram is CBS 380. The x-axis indicates the amount of DNA that is stained by propidium iodide. The green
144 line shows the DNA amount of the G1 phase of the YJF1450 cells (one copy of haploid genome, 1C). The red line
145 shows the G2 phase of the YJF1450 (two copies of haploid genome, 2C) and the G1 phase of the CBS 380 (one copy
146 of diploid genome). The blue line indicates the G2 phase of the CBS 380 (two copies of diploid genomes, 4C). (C)
147 sppIDer results show that most reads from our *S. bayanus* sequencing are mapped to either *S. eubayanus* or *S. uvarum*,
148 confirming the species we sequenced is a hybrid of *S. eubayanus* and *S. uvarum*.
149

150 **De novo assembly and subgenome phasing**

151 Our genome assembly process examined several tools, detailed in the methods section and in Supplemental
152 Table S1, to address the challenges posed by the diploid nature of the target organism. Among the various
153 tools tested, Flye stood out by producing a collapsed-consensus assembly with the highest quality, as
154 reflected in a 96.6% completeness score according to BUSCO analysis. This high score indicates a
155 successful capture of the genomic features we aimed to assemble.

156 Given the diploid nature of our target organism, we aimed to separately assemble the two
 157 subgenomes, diverging from traditional methods that generate a single, collapsed consensus sequence.
 158 Using the MinION platform's long reads, we produced separate and accurate assemblies for each
 159 subgenome. Among the methods employed, phasing the Flye collapsed-consensus assembly via the
 160 Whatshap pipeline proved the most successful at constructing a high-fidelity diploid genomic
 161 representation of *S. bayanus* CBS 380 (Patterson et al. 2015). Post-assembly correction and polishing
 162 resulted in a robust genomic structure ready for further analysis (Table 1). To address the inherent
 163 complexity of the diploid genome of *S. bayanus* CBS 380, our methodology successfully assembled two
 164 distinct subgenomes, technically designated as haplotype-a and haplotype-b, allowing for a sophisticated
 165 analysis of the dual genome architecture (Table 1 and Supplemental Table S2). These precise subgenome
 166 reconstructions paved the way for gene prediction and other evolutionary insights, with the phased variant
 167 calls described in detail in the Materials and Methods.

168
 169 **Table 1.** Assembly statistics for the *S. bayanus* genome and subgenome grouping. The table details the genomic
 170 assembly metrics for the *S. bayanus* species, including the total genome and two subgenomes, Haplotype-a and
 171 Haplotype-b, along with the mitochondrial genome. As ancestral subgenomes cannot be directly inferred, homologous
 172 chromosomes have been categorized into two hypothetical subgenomes to facilitate analysis.
 173

		Total genome excluding mtDNA	Subgenome grouping		Mitochondrial DNA (mtDNA)
			Haplotype-a	Haplotype-b	
Assembly	Genome Size (bp)	23,484,151	11,829,624	11,654,527	64,655
	# of Sequences	32	16	16	1
	Largest (bp)	1,292,201	1,292,201	1,163,801	64,655
	Smallest (bp)	208,383	217,795	208,383	64,655
	Mean (bp)	733,879	739,352	728,408	64,655
	N50 (bp)	912,922	912,922	919,249	64,655
	GC (%)	40.1	40.1	40.1	16.23

	N Count	300	100	200	0
Annotation	Genes	11,545	5,737	5,808	20
	CDS	12,789	6,318	6,471	8

174

175

176 **Genome annotation**

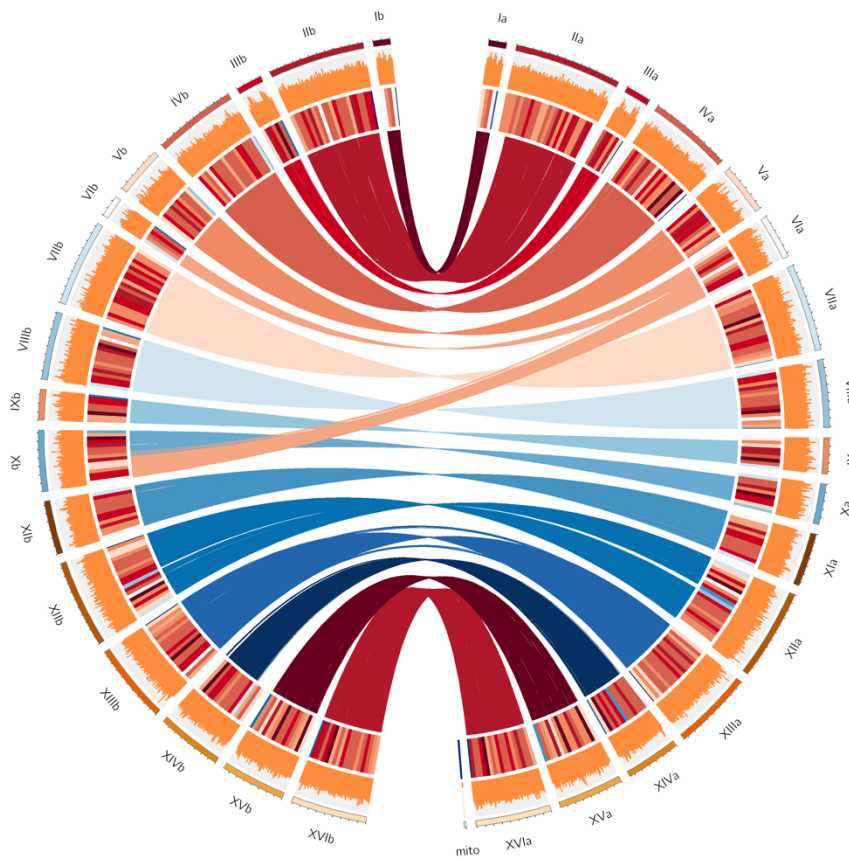
177 Evidence based prediction and annotation of protein-coding genes for each subgenome/haplotype
178 of *S. bayanus* CBS 380 was carried out using the GALBA pipeline (Bruna et al. 2023). The pipeline
179 is perfectly suited to our use case, given its capability to leverage high-quality protein sequences from
180 closely related species. The output revealed a total of 11,547 protein-coding genes identified across both
181 haplotypes, with 5,737 genes in haplotype a and 5,808 genes in haplotype b (Table 1 and Figure 2). The
182 variation in gene count between the two haplotypes is in direct proportion to their chromosomal lengths.

183 The completeness of the genome annotation was assessed by BUSCO based on
184 saccharomycetes_odb10 database, which indicates a high degree of completeness (2095 of 2137,
185 98%). As the BUSCO analysis was based on both haplotype assemblies, most genes are expected
186 to have two copies. As a result, 86.2% of genes (1,843) were classified as duplicates, while only
187 11.8% (252) were identified as unique single-copy genes. Only 13 genes (0.6%) from the
188 saccharomycetes_odb10 gene set were absent from our predicted list. The genome annotation, CDS, and
189 protein sequences are available at <https://github.com/BioHPC/Saccharomyces-bayanus>.

190 Functional annotation of the predicted genes was conducted using EggNog-mapper (Cantalapiedra
191 et al. 2021), which assigned key functional information, such as descriptions of biological functions,
192 orthologous genes in *S. cerevisiae*, Gene Ontology, KEGG pathway and Pfam domains, to 10,985 genes,
193 accounting for 95.1% of the total identified genes (Supplemental Table S3). The combination of functional
194 annotation and BUSCO assessments confirms that our annotation results are comprehensive, providing a
195 solid foundation for our further analysis.

196

197



198

199 **Figure 2.** Circos plot representing the 16 pairs of chromosomes of the *S. bayanus* genome, showing different genomic
200 features across four concentric circles. The outermost circle represents the karyotype of the *S. bayanus* genome for
201 two haplotypes, with the right part representing haplotype-a and the left part representing haplotype-b. The second
202 outermost circle represents the GC content on each chromosome of the genome. The third circle provides information
203 on gene density within the chromosomes. The innermost circle highlights syntenic blocks between haplotypes,
204 illustrating regions of genetic similarity and divergence between haplotype-a and haplotype-b.
205

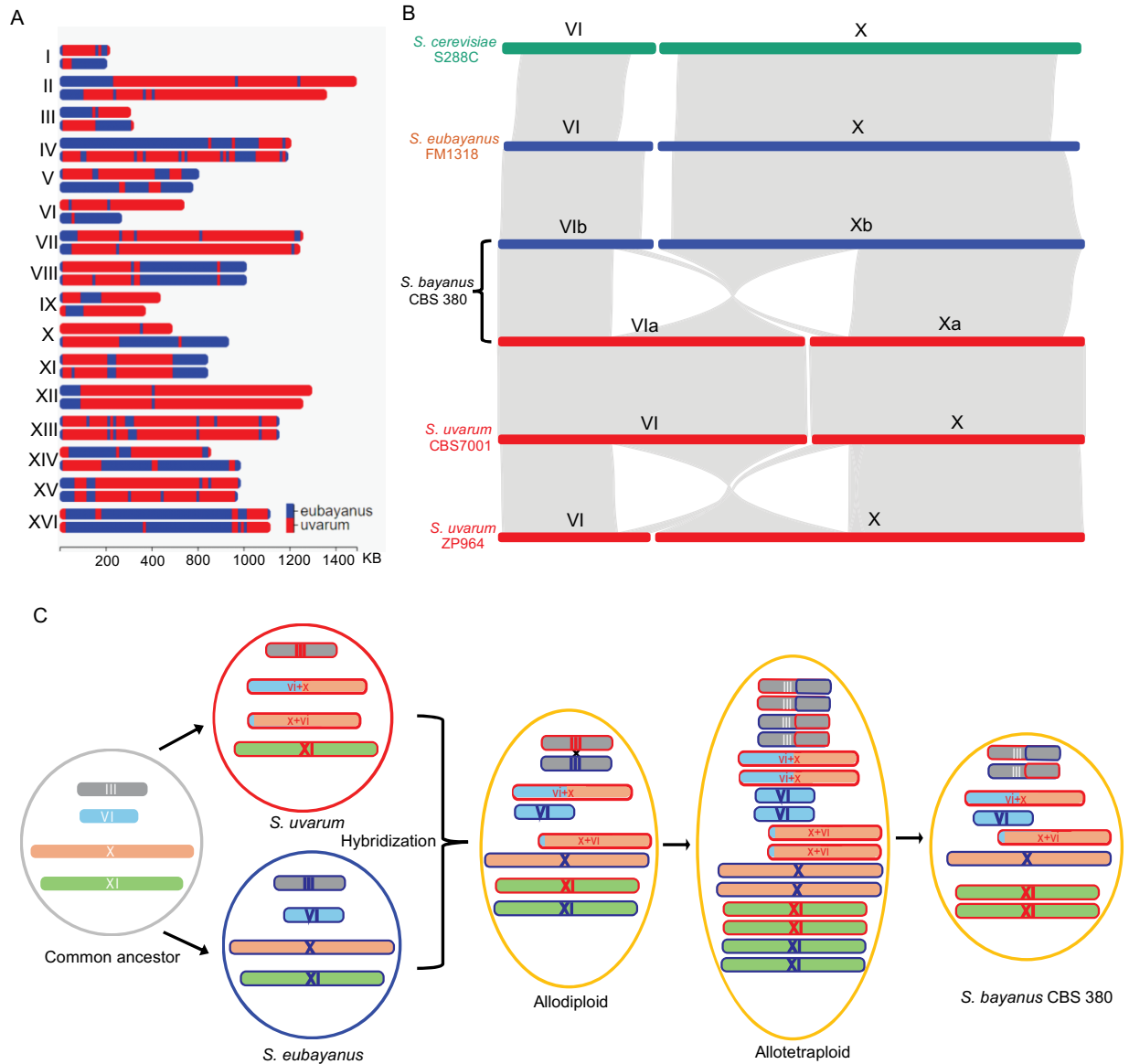
206 **Inference of parental genomic regions**

207 To identify the major genomic events that have occurred in the *S. bayanus* genome since hybridization,
208 including recombination, chromosomal rearrangements, and loss of heterozygosity, we first used two
209 approaches to determine the origin of genomic regions in the hybrid genome. Our first approach is based
210 on BLAST searches of non-overlapping blocks of 5,000 bp for every chromosome against the genomes of
211 *S. eubayanus* and *S. uvarum* (see Methods and Materials). In brief, the origin of each genomic block was

212 determined by its best hit of BLAST search. As illustrated in Figure 3A, each haplotype chromosome
213 contains regions that originated from both *S. uvarum* and *S. eubayanus*, suggesting the recombination
214 between the two orthologous chromosomes of the two subgenomes, creating mosaic chromosomes
215 composed of genomic regions of heterozygous origins. However, the proportions of each subgenome vary
216 substantially across different chromosomes. For instance, segments of *S. eubayanus* origin make up 81%
217 of Chr IVa, whereas they make up only 14% of Chr IVb. In addition, nine of the 16 chromosome pairs
218 have a high degree of homozygosity, meaning that the genomic origin and recombinants are very similar
219 between haplotypes a and b, showing that heterozygosity was quickly lost after hybridization.

220 To confirm the robustness of the results based on the BLAST method, we employed an approach
221 based on the degree of divergence of synonymous sites in protein-coding regions, as it is generally assumed
222 that synonymous mutations are selectively neutral (see Methods and Materials). To summarize, the method
223 first compared rates of synonymous substitution (K_s) between the two alleles of *S. bayanus* and then
224 between orthologous genes from the three species. We then determined if an allele in *S. bayanus* is more
225 similar to its orthologous gene in *S. uvarum* or in *S. bayanus*. We identified a total number of 5,497
226 orthologous groups from the three species using OrthoFinder (Emms and Kelly 2019) (Supplemental
227 Materials). The distribution of K_s between all pairs of alleles in *S. bayanus* shows two distinct peaks
228 (Supplemental Figure S2). The left peak, which consists of lower K_s values, represents sequence divergence
229 between two homozygous alleles (alleles originated from one single parental genome). In contrast, the right
230 peak contains higher K_s values that were obtained from two alleles originated from different parental
231 genomes (heterozygous origins). Consistently, the distribution of K_s values between orthologous genes
232 between the two parental genomes largely overlaps with the right peak of K_s values in *S. bayanus*
233 (Supplemental Figure S2). Next, we calculated K_s values between an allele from *S. bayanus* and its
234 orthologous gene in *S. uvarum* (K_{su}) and *S. eubayanus* (K_{se}) respectively. The origin of each allele was
235 then determined by comparing K_{su} and K_{se} to the two K_s peaks. The results of genomic origin obtained
236 based on our K_s method are highly consistent with the BLAST method (Supplemental Figure 3), supporting

237 the occurrence of multiple recombination events on most chromosomes and rapid loss of heterozygosity
 238 after interspecies hybridization.
 239



240
 241 **Figure 3.** The origin and evolution of *S. bayanus* chromosomes. (A) Origin of genomic regions of each chromosome
 242 in the *S. bayanus* genome based on BLAST searches of non-overlapping 5,000 bp blocks. Genomic Regions that
 243 originated from *S. eubayanus* are shown in blues, while regions inherited from *S. uvarum* are shown in red. (B)
 244 Synteny block of Chromosome VI and X between *S. cerevisiae*, *S. eubayanus*, both *S. bayanus* haplotypes, *S. uvarum*
 245 strain CBS 7001, *S. uvarum* strain ZP964. (C) An evolutionary model of *S. bayanus* chromosomes. For simplification
 246 purposes, only four chromosomes are shown, representing different patterns of chromosome
 247 inheritances. Translocation between Chr VI and X occurred in *S. eubayanus* prior to its hybridization with *S. bayanus*.
 248 Recombination and whole genome duplication occurred in the hybrid *S. bayanus* genome. Subsequent genome

249 reduction by chromosome losses, created some heterozygous chromosomes, such as Chr III, and some homozygous
250 chromosomes, such as Chr XI.
251

252 **A model of allotetraploid origin of *S. bayanus***

253 We found distinct differences in the lengths of chromosomes VI and X between the two subgenomes
254 (haplotypes) of *S. bayanus* CBS380 (Figures 2 and 3A, Supplemental Table S2). Specifically, Chr VIa is
255 ~277k bp longer than Chr VIb (544k vs. 267k), while Chr Xa is ~244k bp shorter than Chr Xb. Our analysis
256 of syntenic regions between the two haplotypes shows that a significant portion of Chr VIa has syntenic
257 regions to Chr Xb. These observations suggest that a translocation occurred between Chr VI and X. Next,
258 we sought to determine whether the translocation was from Chr VI to Chr X or visa versa, and whether the
259 translocations occurred in the parent genomes or after hybridization. The answers to these questions are
260 key to better understanding how hybrid *S. bayanus* arose, and the mechanism by which heterozygosity is
261 rapidly lost on most chromosomes.

262 To investigate the direction and timing of the translocation event, we conducted a detailed
263 examination of all available genomes of *S. uvarum* and *S. eubayanus* strains from the NCBI database. Our
264 results showed that the patterns of chromosome lengths in all *S. eubayanus* strains were consistent with
265 those observed in *S. cerevisiae*, i.e., a short chromosome VI and a long chromosome X (Figure 3B). In
266 contrast, heterogenous lengths of Chr VI and Chr X are observed among *S. uvarum* strains (Figure 3B). For
267 instance, *S. uvarum* ZP964 has similar lengths of Chr VI and Chr X to those of *S. eubayanus* and *S.*
268 *cerevisiae*. In contrast, *S. uvarum* CBS 7001 has a much longer Chr VI, but much shorter Chr X, similar to
269 the haplotype a in *S. bayanus* (Figure 3B). Gene collinearity analysis of the two chromosomes among these
270 species further showed that the translocation occurred only once in the lineage of *S. uvarum* CBS 7001
271 prior to hybridization. *S. bayanus* inherited translocated Chr VI and Chr X from a parental species that is
272 closely related to *S. uvarum* CBS 7001. These results also support that the translocation was generated by
273 exchanging ~270 KB segment at the left end of Chr X with ~30 KB region at the right end of Chr VI (Figure
274 3B).

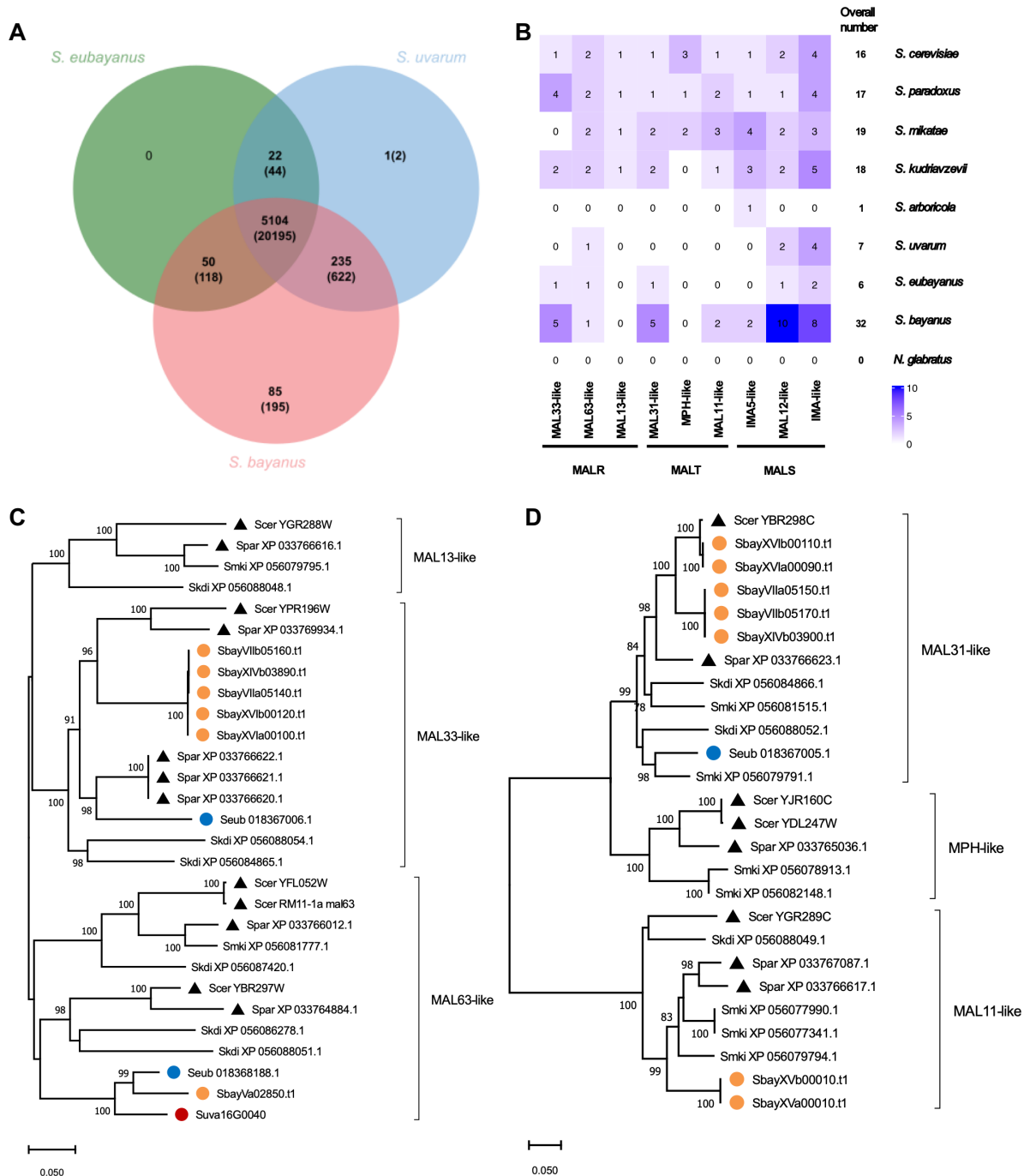
275 Our analysis of origins of genomic regions in *S. bayanus* demonstrates that only seven pairs of
276 chromosomes maintained heterozygous status, and loss of heterozygosity occurred to other chromosome
277 pairs (Figure 3A). Several genetic mechanisms have been proposed to explain the LOH after hybridization,
278 such as whole-genome duplication followed by chromosome loss, duplication or loss of individual
279 chromosomes, and gene conversion (Marcet-Houben and Gabaldon 2015; Wolfe 2015; Wertheimer et al.
280 2016). Duplication and loss of individual chromosomes often result in chromosomal aneuploidies.
281 However, we did not observe obvious chromosomal aneuploidies in *S. bayanus* based on read depth of
282 most, if not all, chromosomes. In addition, the track length of gene conversion is usually limited, which is
283 not supported by our observations that the track length of LOH covers almost entire chromosomes.
284 Furthermore, the locations of recombination events are very similar between haplotypes in most
285 chromosomes, such as Chr VIII, Chr XI, and Chr XII (Figure 3A). Based on these observations, it is mostly
286 parsimonious to propose that the hybrid allopolyploid genome may have undergone duplication without cell
287 division (non-disjunction), resulting in a temporary allotetraploid genome. Subsequence loss of
288 chromosomes may have occurred to the allotetraploid genome, resulting in haploid status (Figure 3C).
289 Therefore, the allotetraploid origin of *S. bayanus* is similar to the evolutionary history of *S. pastorianus*
290 (Dunn and Sherlock 2008; Nakao et al. 2009; Libkind et al. 2011).

291

292 **Evolution of gene content after hybridization and expansion of genes involved in maltose** 293 **metabolism**

294 To better understand the evolution of gene content after hybridization, we carried out further analyses on
295 the 5,497 orthologue groups (OG) in *S. bayanus* and its parental species *S. uvarum* and *S. eubayanus*. A
296 total number of 5,412 unique OGs are present in the two parent genomes. 5,389 of them (99.6%) are also
297 present in *S. bayanus*, suggesting that gene loss in *S. bayanus* is very limited after hybridization (Figure
298 4A). Interestingly, 85 OGs (195 genes) are only present in the genome of *S. bayanus*. Given that ~5% of

299 MinION reads were specifically mapped to *S. cerevisiae*, we speculated that these genes were likely
 300 originated from introgression events of *S. cerevisiae* as proposed in a previous study (Nguyen et al. 2011).



301
 302 **Figure 4.** The evolution of gene content in the hybrid genome of *S. bayanus* CBS 380. (A) A Venn diagram showing
 303 the numbers of shared and species-specific orthologous groups (OGs). (B) Evolutionary changes of the three MAL
 304 gene families in nine OGs in the *Saccharomyces sensu stricto* members and *Nakaseomyces glabratus*. Increased gene
 305 copy numbers were observed in each MAL gene family in *S. bayanus*. (C) A phylogenetic tree of the MALR gene

306 family suggests introgression of MALR genes in *S. bayanus*. (D) A phylogenetic tree of the MALT gene family
307 suggests that all MALT genes in *S. bayanus* were likely acquired by introgression events from other species.
308

309 Among the 5,104 OGs present in all three species, 4,065 OGs (79.6%) exhibit a 1:1:2 ratio, indicating
310 conservation of gene copy numbers in all three species. It is worth noting that 200 of these OGs contain
311 more than two copies of genes in the diploid genome of *S. bayanus*, suggesting expansion of gene copy
312 number in *S. bayanus*. Nine of these expanded OGs are involved in maltose/maltotriose utilization (Figure
313 4B). *S. bayanus* is mainly found in brewing environments and maltose is the most abundant sugar (~60%)
314 in brewer's wort (Magalhaes et al. 2016). Therefore, the expansion of genes involved in maltose or
315 maltotriose utilization (MAL genes) might have facilitated the adaptation of *S. bayanus* to maltose-rich
316 environments.

317 MAL genes were classified into three families based on their functions, including maltose
318 transporter (MALT), enzymes that break down maltose (MALS), and genes that regulate the expression of
319 the pathway (MALR). These genes are often organized into clusters and located near the ends of
320 chromosomes (subtelomeric). To elucidate the origins and expansion of MAL genes in *S. bayanus*, we first
321 identified all MAL genes in *S. bayanus* and eight other *Saccharomyces* species. The total numbers of MAL
322 genes vary substantially among non-hybrid species, ranging from 0 in *Nakaseomyces glabratus*, to 19 in *S.*
323 *mikatae* (Figure 4B). The diploid genome of *S. bayanus* contains a significantly higher number of MAL
324 genes (32 in total) compared to 7 in *S. uvarum* and 6 in *S. eubayanus* (Figure 4B), suggesting a significant
325 expansion in copy numbers of MAL genes in *S. bayanus* after hybridization. This is particularly noticeable
326 that the MALS gene family was increased from 3 and 6 in their parental species genomes to 20 in *S.*
327 *bayanus*.

328 To infer the origin and mechanism of MAL gene expansion in *S. bayanus*, we carried out
329 phylogenetic analyses for each of the three MAL families using their amino acid sequences (Figure 4C-D,
330 Supplemental Figure S4). Based on the tree topology, we found that the majority of the MAL genes in *S.*
331 *bayanus* were not inherited directly from its parental genomes. Instead, these genes seem to be likely

332 acquired through introgression from other species, mostly from *S. cerevisiae*, followed by multiple gene
333 duplication events. For example, 6 copies of MALR genes are present in *S. bayanus*. Only one MALR gene
334 (MAL63-like) is group with *S. eubayanus*, while the other five form a well-supported clade that is closely
335 related to MAL33-like genes in *S. cerevisiae* and *S. paradoxus* (Figure 4C), suggesting multiple rounds of
336 gene duplication to MAL33-like genes after acquisition of MAL33-like genes from ancestral *S. cerevisiae*
337 or *S. paradoxus*. Similarly, none of the seven MALT genes was grouped with either *S. eubayanus* or *S.*
338 *uvarum* (Figure 4D). In the case of MALS genes, a total number of 20 MALS genes are found in *S. bayanus*,
339 and only 8 of them appear to be originated from *S. eubayanus*, and expanded by gene duplication. Similar
340 to other *Saccharomyces* species, most MAL genes in *S. eubayanus* also form clusters and reside in
341 subteleomeric regions (Supplemental Figure S5).

342 We noticed that none of *S. bayanus* MAL genes were inherited from *S. uvarum* (Figure 4C-D,
343 Supplemental Figure S4). It suggests that there was a preferred retention of MAL genes inherited from *S.*
344 *eubayanus* or a preferred loss of MAL genes inherited from *S. uvarum*. Given that *S. uvarum* has contributed
345 over 60% of the genetic makeup of *S. bayanus*, the strong exclusion of *S. uvarum* MALs genes in the *S.*
346 *bayanus* genome were unlikely due to random events. One possibility is that MAL genes from *S. uvarum*
347 might imposed selective disadvantages under maltose-rich brewing environments. Future studies on the
348 growth effects of *S. uvarum* MAL genes may provide new insights into the biased retention of MAL genes.
349

350 **Discussion**

351 We present the first chromosome level subgenome assembly and annotations of the hybrid yeast, *S. bayanus*
352 (CBS 380) which will serve as an excellent reference for future studies of this important yeast and other
353 yeast strains. The assembly was completed using only Oxford Nanopore technology on a single MinION
354 flow cell. Thus, we show the utility of high read depth sequencing, that is available for moderate costs using
355 this technology. We assessed the assemblies from fifteen different *de novo* assembly pipelines, all run on
356 relatively modestly equipped computer workstations, and concluded that the Flye method outperformed the

357 others in producing an assembly with the fewest contigs and high N50 scores. The successful application
358 of the GALBA pipeline allowed for high-fidelity annotation of the two subgenomes, revealing a total of
359 11,547 protein-coding genes and confirming the completeness of our genome assembly with a high BUSCO
360 score. This type of sequencing can be carried out in most laboratories without previous sequencing
361 experience or high-performance computational resources.

362 Through a dual approach involving BLAST searches and synonymous site divergence, we traced
363 recombination events and chromosomal rearrangements that describe the history of *S. bayanus* after
364 hybridization. In particular, the discovery of mosaic chromosomes with heterozygous origins in *S. uvarum*
365 and *S. eubayanus* speaks to the dynamic evolutionary past of this species. Our data suggest a rapid loss of
366 heterozygosity, which can be attributed to multiple genetic mechanisms, including a proposed transient
367 polyploid phase and selective chromosome loss. This model not only parallels the evolutionary history of
368 related species such as *S. pastorianus*, but also provides a plausible explanation for the observed genome
369 organization.

370 Although the majority of genes in *S. bayanus* maintained their copy numbers, large copy number
371 variations were found in the three MAL families. In addition to direct transmission from parent species, our
372 phylogenetic analysis suggested that many of them were acquired from other yeast species, mostly from *S.*
373 *cerevisiae*. Subsequent gene duplication events on MAL genes further increased its copy number. It is
374 reasonable to postulate that the introgression and duplication of MAL genes provide selective advantages
375 in maltose-rich brewing environments.

376 Despite *S. bayanus* inheriting most chromosomal segments from *S. uvarum*, none of the MAL genes
377 in *S. bayanus* were traced to its *S. uvarum* progenitor. It is unlikely due to random loss of *S. uvarum* copies.
378 One possibility is that *S. eubayanus* MALs genes are more efficient or provide more selective advantages
379 than those of *S. uvarum*. Further studies can be performed to examine the functional differences of these
380 MAL genes in maltose metabolism between the two species, which could provide new valuable information
381 for improving industrial brewing using maltose-rich materials.

382

383 **Materials and Methods**

384 **Yeast strain, growth condition, genomic DNA isolation**

385 *Saccharomyces bayanus* CBS 380's cells were grown on YPD medium (1% yeast extract, 2% peptone, and
386 2% glucose) at 30 °C for 16 hours. Extraction of high molecular weight genomic DNA (HMW gDNA)
387 from *S. bayanus* cells was carried out by following a protocol described by Denis et al (Denis E 2018). In
388 brief, *S. bayanus* cell wall was first lysed with Zymolyase (MP Biomedicals). Spheroplasts were then
389 collected and resuspended in SDS buffer with RNase A. Proteins were precipitated and removed with
390 potassium acetate and centrifugation. The supernatants were used to precipitate DNA with isopropanol.
391 DNA pellet was then washed with 70% ethanol and dissolved in TE buffer. The quality and quantity of the
392 extracted DNA were determined using Qubit (Invitrogen). HMW gDNA was sheared into 20kb fragments
393 using g-TUBE (Covaris Inc).

394

395 **Determination of ploidy**

396 We performed a flow cytometry analysis to determine the ploidy of the *S. bayanus* CBS 380 following the
397 protocol (Todd et al. 2018). We also used a haploid *S. uvarum* strain YJF1450 (MAT α ho Δ ::NatMX,
398 derived from CBS 7001, a gift from J. Fay lab at Rochester University) as a control. Briefly, yeast cells
399 were grown to log-phase (OD = 0.3) in YPD medium on a shaker platform at 30 °C by rotation at 225 RPM.
400 Then, cells were fixed in 70% ethanol at 4 °C overnight and then sonicated to separate cells. After RNase
401 A (0.5 mg/ml) treatment for 2 hours, the cells were stained with 25 μ g/ml of propidium iodide at 4 °C
402 overnight. Finally, the stained cells were analyzed using BD Accuri C6 Plus and the data were analyzed in
403 FlowJo v10.8.1.

404

405 **MinION library preparation and sequencing**

406 HMW gDNA were then used to prepare MinION sequencing library using the Nanopore Rapid Sequencing
407 Kit (SQK-RAD004) following the manufacturer's instruction. Briefly, the sample mix was prepared with

408 7.5 µl template DNA (~2µg) and 2.5 µl fragmentation mix and incubated at 30°C for 1 min and then at
409 80 °C for 1min. 1 µl Rapid Adapter was added to the sample mix and incubated for 5 min at room
410 temperature. Priming mix was prepared by adding 30 µl of Flush Tether and Flush Buffer. The priming mix
411 was loaded into the flow cell via the priming port. Sequencing mix was prepared with DNA sample mix
412 and was loaded to the flow cell via the SpotON sample port.

413

414 **Adapter removal**

415 Porechop v0.2.4 (Wick et al. 2017) was used for adapter identification and removal using default thresholds.
416 In all, 179,725 reads had adapters trimmed from their start (15,472,707 bases removed), and 778 reads were
417 split based on middle adapters. (Supplemental Figure 1). A full list of commands and parameters is available
418 in the Supplemental Materials.

419

420 **Genome assembly, post-assembly correction, and genome polishing**

421 Draft collapsed-consensus assemblies were generated using Canu v2.2 (Koren et al. 2017), Flye v2.9
422 (Kolmogorov et al. 2019), Wtdbg2 v2.5 (Ruan and Li 2020), NECAT v0.0.1 (Chen et al. 2021),
423 SMARTdenovo v1.0.0 (Liu et al. 2021), NextDenovo v2.5.0 (NextOmics 2021), Raven v1.8.0 (Vaser et al.
424 2017), and Ra v0.2.1 (Vaser 2019), with both uncorrected and Canu corrected and trimmed reads
425 (Supplemental Table 1). These methods were executed on a general workstation-level computer (36 cores
426 and 128GB memory), demonstrating the feasibility of ONT-based de novo assembly for small genomes in
427 modestly equipped laboratories.

428

429 **Complete subgenome-aware de novo genome assembly**

430 Given the diploid nature of our target organism, we aimed to generate a diploid-level representation of each
431 chromosome. We employed long-read sequencing to facilitate the generation of full-length, phased

432 haplotype *de novo* assemblies, using a suite of assembly tools, as detailed below and in the Supplemental
433 Material.

434 ***Haplotype-aware de novo genome assembly***

435 We experimented with haplotype-aware assembly methods such as Flye (with haplotype preservation
436 enabled) (Kolmogorov et al. 2019), Shasta (Shafin et al. 2020), Phasebook (Luo et al. 2021), and CanuTrio
437 (which organizes reads into haplotype-specific bins before assembly) (Koren et al. 2017). These approaches
438 did not yield high-quality assemblies that were both contiguous and reflective of the expected genome size,
439 leading to their exclusion from analysis.

440 ***Phasing-based diploid genome assembly***

441 To tackle the complexities of *S. bayanus* CBS 380's diploid genome, we undertook a phasing-based
442 assembly strategy, leveraging the long reads generated from Oxford Nanopore's MinION platform. Prior
443 to phasing, the purge_dups pipeline was used to remove haplotype duplication in the primary assemblies
444 (Guan et al. 2020). To construct a phased diploid genome assembly, we first called variants using Claire
445 (Zheng et al. 2022). The variant calls were processed through the WhatsHap pipeline which exploits the
446 connectivity between heterozygous variants within individual reads to generate phased haplotypes
447 (Patterson et al. 2015). To generate a haplotype-specific genomic representation, we used BCFtools
448 'consensus' followed by WhatsHap manual. This allowed us to extract the separate FASTA representations
449 for each haplotype, effectively translating the phased information into a coherent, usable format for further
450 analysis. BUSCO was used to assess the assembly's completion (Simao et al. 2015). A comprehensive list
451 of commands and parameters, along with the phased variant calls, are accessible in the Supplemental
452 Materials, offering a resource for future genetic and evolutionary studies.

453 ***Genome correction and polishing***

454 For assembly correction and polishing, the raw ONT sequencing reads were split via the 'whatshap split'
455 subcommand to segregate the set of unmapped reads according to their haplotypes. This generated two
456 distinct FASTQ files, each corresponding to one of the haplotypes identified within the sample. The
457 assembled contigs were then passed to a series of correction and polishing steps to enhance their accuracy,

458 utilizing Racon (v1.4.3) (Vaser et al. 2017) and Medaka (v1.9.1) (<https://github.com/nanoporetech/medaka>)
459 for error correction and sequence improvement. This correction process was executed separately for each
460 haplotype, utilizing their respective reads. A total of four iterative rounds of correction were iteratively
461 performed with Racon for each haplotype. This cycle involved mapping the haplotype-resolved reads to
462 the assembled contigs using Minimap2 (using the ONT-specific `-x map-ont` option), followed by Racon-
463 based correction to refine assembly quality progressively. After completing the Racon correction cycles, a
464 final round of polishing was conducted using Medaka. This step uses a neural network-based approach to
465 correct consensus sequence errors, further enhancing the accuracy of the assembled haplotypes.`

466

467 **Genome annotation**

468 We employed the GALBA pipeline to annotate protein-coding genes for the assembled nuclear genome
469 (Bruna et al. 2023). Specifically, we used amino acid sequences from *Saccharomyces cerevisiae*,
470 *Saccharomyces uvarum*, and *Saccharomyces eubayanus* as inputs. These protein sequences were aligned
471 to both subgenomes of *S. bayanus* using the Miniprot (Li 2023), followed by gene annotation using
472 AUGUSTUS (Stanke et al. 2006). The output GTF files were processed using AGAT
473 (<https://github.com/NBISweden/AGAT>) for format cleaning and conversion. The completeness of the gene
474 annotation was evaluated using the BUSCO (version: 5.5.0) (Simao et al. 2015), employing the
475 *saccharomycetes_odb10* database for assessment. For functional annotation of predicted genes, we utilized
476 the web version of eggno-mapper (Cantalapiedra et al. 2021) to upload the *S. bayanus* protein files. All
477 other parameters were retained as default settings.

478

479 **Ancestral inference of *S. bayanus* using a BLAST-based approach**

480 To infer the ancestral parentage of the hybrid yeast, we conducted a comparative genomic analysis using a
481 custom script that performs local BLAST (Camacho et al. 2009) homology searches (see Supplemental
482 Materials). The hybrid yeast genome was segmented into consecutive, non-overlapping chunks of 5,000

483 base pairs, which were then individually compared against the genomes of the two parental strains using
484 the BLAST algorithm. This approach allowed for the identification of the closest matching regions between
485 the hybrid and each parent genome, based on the highest bitscore values obtained from the BLAST results.
486 The bitscore, serving as a measure of sequence similarity, was selected as the primary criterion for parental
487 inference. A score threshold of 100 bits was set to distinguish between significant and non-significant
488 matches, thereby facilitating the identification of the most probable ancestral parent for each genomic
489 segment of the hybrid yeast.

490

491 **Inference of gene origin in *S. bayanus* based on similarity of synonymous sites**

492 To delineate the evolutionary lineage of genes within *S. uvarum*, *S. eubayanus*, and *S. bayanus*, we initially
493 conducted an OrthoFinder (Emms and Kelly 2019) analysis on the coding sequence (CDS) datasets of these
494 species to identify orthologous genes. Subsequently, alignments were performed utilizing PRANK
495 (Loytynoja 2021) with codon model. Following alignment, we employed KaKs_Calculator (Zhang 2022)
496 (version 3.0) to compute the synonymous substitution rate (K_s). The derived K_s values were then utilized
497 to categorize gene lineage based on their similarity, providing insights into the genetic heritage of the genes
498 of *S. bayanus*.

499

500 **Comparative Genomic Analysis and Evolutionary Study of the MAL Gene family in *S.*** 501 ***bayanus***

502 To elucidate the evolutionary relationships among *Saccharomyces bayanus* and closely related species, we
503 analyzed coding sequence (CDS) datasets for *S. cerevisiae*, *S. paradoxus*, *S. mikatae*, *S. kudriavzevii*, *S.*
504 *arboricola*, *S. uvarum*, *S. eubayanus*, and *Nakaseomyces glabratus*. Using OrthoFinder (Emms and Kelly
505 2019), we identified orthogroups to enable a comparative genomics study. Adopting the protocols from
506 (Brown et al. 2010) and (Baker et al. 2015), we identified genes belonging to the maltose utilization (MAL)
507 gene families. Sequence alignment was conducted with MAFFT using the L-INS-i strategy (Katoh and

508 Standley 2013). Phylogenetic trees were generated using FastTree v2.1 (Price et al. 2010). The evolutionary
509 history was inferred using the Neighbor-Joining method (Saitou and Nei 1987). The original tree is shown.
510 The percentage of replicate trees in which the associated taxa clustered together in the bootstrap test (1000
511 replicates) are shown next to the branches (Felsenstein 1985). The tree is drawn to scale, with branch lengths
512 in the same units as those of the evolutionary distances used to infer the phylogenetic tree. The evolutionary
513 distances were computed using the Maximum Composite Likelihood method and are in the units of the
514 number of base substitutions per site. These trees were visualized and refined with MEGA11 (Tamura et
515 al. 2021).

516

517 **Acknowledgments**

518 C.G., C.H. D.D. and T.A were supported by the National Science Foundation (NSF) under Grant No.
519 1564894 and Z. L was supported by NSF under Grant No. 1951332. M.J.D receives funding from the
520 National Institute of Allergy and Infectious Diseases of the National Institutes of Health under award
521 number R01AI123407. We thank Dr. Justin Fay for providing *S. uvarum* YJF1450 strain for ploidy analysis.

522

523 **Data Availability**

524 Sequencing and genome assembly data generated in this work have been deposited at the NCBI repository
525 under the BioProject accession PRJNA741321. Annotations and supplemental materials are available at
526 <https://github.com/BioHPC/Saccharomyces-bayanus>.

527

528 **Author Contributions**

529 T.A. and Z.Lin conceived the idea. T.A, Z.Lin, and M.J.D. supervised this study. Z.Lu isolated DNA,
530 prepared libraries and performed Nanopore sequencing. Y.Z. performed flow cytometry. C.G., J.C, C.H.,
531 and D.D. analyzed the data. All authors wrote the manuscript and approved the final version of the
532 manuscript.

533

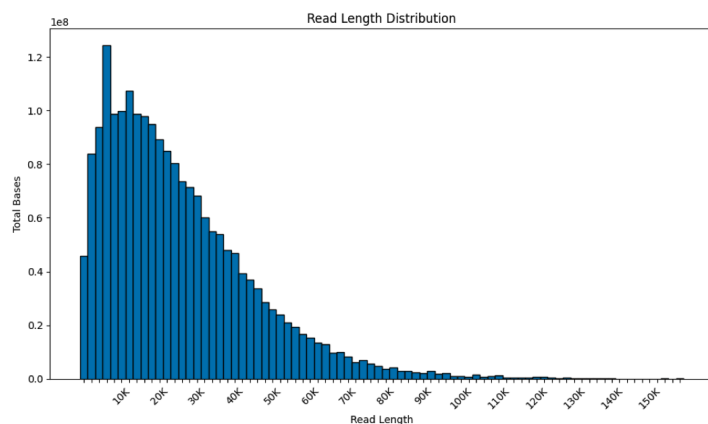
534 References

- 535 Almeida P, Goncalves C, Teixeira S, Libkind D, Bontrager M, Masneuf-Pomarede I, Albertin W, Durrens
536 P, Sherman DJ, Marullo P et al. 2014. A Gondwanan imprint on global diversity and
537 domestication of wine and cider yeast *Saccharomyces uvarum*. *Nat Commun* **5**: 4044.
- 538 Baker E, Wang B, Bellora N, Peris D, Hulfachor AB, Koshalek JA, Adams M, Libkind D, Hittinger CT.
539 2015. The Genome Sequence of *Saccharomyces eubayanus* and the Domestication of Lager-
540 Brewing Yeasts. *Mol Biol Evol* **32**: 2818-2831.
- 541 Brown CA, Murray AW, Verstrepen KJ. 2010. Rapid expansion and functional divergence of
542 subtelomeric gene families in yeasts. *Curr Biol* **20**: 895-903.
- 543 Bruna T, Li H, Guhlin J, Honsel D, Herbold S, Stanke M, Nenasheva N, Ebel M, Gabriel L, Hoff KJ.
544 2023. Galba: genome annotation with miniprot and AUGUSTUS. *BMC Bioinformatics* **24**: 327.
- 545 Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL. 2009. BLAST+:
546 architecture and applications. *BMC Bioinformatics* **10**: 421.
- 547 Cantalapiedra CP, Hernandez-Plaza A, Letunic I, Bork P, Huerta-Cepas J. 2021. eggNOG-mapper v2:
548 Functional Annotation, Orthology Assignments, and Domain Prediction at the Metagenomic
549 Scale. *Mol Biol Evol* **38**: 5825-5829.
- 550 Caudy AA, Guan Y, Jia Y, Hansen C, DeSevo C, Hayes AP, Agee J, Alvarez-Dominguez JR, Arellano H,
551 Barrett D et al. 2013. A new system for comparative functional genomics of *Saccharomyces*
552 yeasts. *Genetics* **195**: 275-287.
- 553 Chen Y, Nie F, Xie SQ, Zheng YF, Dai Q, Bray T, Wang YX, Xing JF, Huang ZJ, Wang DP et al. 2021.
554 Efficient assembly of nanopore reads via highly accurate and intact error correction. *Nat Commun*
555 **12**: 60.
- 556 Conant GC, Wolfe KH. 2007. Increased glycolytic flux as an outcome of whole-genome duplication in
557 yeast. *Mol Syst Biol* **3**: 129.
- 558 Coyne JA, Orr HA. 2004. *Speciation*. Sinauer Associates, Sunderland, Mass.
- 559 Dobzhansky T. 1982. *Genetics and the origin of species*. Columbia University Press, New York.
- 560 Dunn B, Sherlock G. 2008. Reconstruction of the genome origins and evolution of the hybrid lager yeast
561 *Saccharomyces pastorianus*. *Genome Res* **18**: 1610-1623.
- 562 Emms DM, Kelly S. 2019. OrthoFinder: phylogenetic orthology inference for comparative genomics.
563 *Genome Biol* **20**: 238.
- 564 Felsenstein J. 1985. Confidence Limits on Phylogenies: An Approach Using the Bootstrap. *Evolution* **39**:
565 783-791.
- 566 Gabaldon T. 2020. Hybridization and the origin of new yeast lineages. *FEMS Yeast Res* **20**.
- 567 Gorter de Vries AR, Pronk JT, Daran JG. 2017. Industrial Relevance of Chromosomal Copy Number
568 Variation in *Saccharomyces* Yeasts. *Appl Environ Microbiol* **83**.
- 569 Guan D, McCarthy SA, Wood J, Howe K, Wang Y, Durbin R. 2020. Identifying and removing haplotypic
570 duplication in primary genome assemblies. *Bioinformatics* **36**: 2896-2898.
- 571 Hittinger CT. 2013. *Saccharomyces* diversity and evolution: a budding model genus. *Trends Genet* **29**:
572 309-317.
- 573 Katoh K, Standley DM. 2013. MAFFT multiple sequence alignment software version 7: improvements in
574 performance and usability. *Mol Biol Evol* **30**: 772-780.
- 575 Kellis M, Birren BW, Lander ES. 2004. Proof and evolutionary analysis of ancient genome duplication in
576 the yeast *Saccharomyces cerevisiae*. *Nature* **428**: 617-624.
- 577 Kolmogorov M, Yuan J, Lin Y, Pevzner PA. 2019. Assembly of long, error-prone reads using repeat
578 graphs. *Nat Biotechnol* **37**: 540-546.

- 579 Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. 2017. Canu: scalable and
580 accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome Res* **27**:
581 722-736.
- 582 Langdon QK, Peris D, Baker EP, Oplente DA, Nguyen HV, Bond U, Goncalves P, Sampaio JP, Libkind
583 D, Hittinger CT. 2019. Fermentation innovation through complex hybridization of wild and
584 domesticated yeasts. *Nat Ecol Evol* **3**: 1576-1586.
- 585 Li H. 2023. Protein-to-genome alignment with miniprot. *Bioinformatics* **39**.
- 586 Libkind D, Hittinger CT, Valerio E, Goncalves C, Dover J, Johnston M, Goncalves P, Sampaio JP. 2011.
587 Microbe domestication and the identification of the wild genetic stock of lager-brewing yeast.
588 *Proc Natl Acad Sci U S A* **108**: 14539-14544.
- 589 Lin Z, Li W-H. 2014. Comparative Genomics and Evolutionary Genetics of Yeast Carbon Metabolism. In
590 *Molecular Mechanisms in Yeast Carbon Metabolism*, doi:10.1007/978-3-642-55013-3_5 (ed. J
591 Piškur, C Compagno), pp. 97-120. Springer Berlin Heidelberg, Berlin, Heidelberg.
- 592 Lin Z, Li WH. 2011. Expansion of hexose transporter genes was associated with the evolution of aerobic
593 fermentation in yeasts. *Mol Biol Evol* **28**: 131-142.
- 594 Liu H, Wu S, Li A, Ruan J. 2021. SMARTdenovo: a de novo assembler using long noisy reads. *GigaByte*
595 **2021**: gigabyte15.
- 596 Loytynoja A. 2021. Phylogeny-Aware Alignment with PRANK and PAGAN. *Methods Mol Biol* **2231**:
597 17-37.
- 598 Luo X, Kang X, Schonhuth A. 2021. phasebook: haplotype-aware de novo assembly of diploid genomes
599 from long reads. *Genome Biol* **22**: 299.
- 600 Magalhaes F, Vidgren V, Ruohonen L, Gibson B. 2016. Maltose and maltotriose utilisation by group I
601 strains of the hybrid lager yeast *Saccharomyces pastorianus*. *FEMS Yeast Res* **16**.
- 602 Marcet-Houben M, Gabaldon T. 2015. Beyond the Whole-Genome Duplication: Phylogenetic Evidence
603 for an Ancient Interspecies Hybridization in the Baker's Yeast Lineage. *PLoS Biol* **13**: e1002220.
- 604 Morales L, Dujon B. 2012. Evolutionary role of interspecies hybridization and genetic exchanges in
605 yeasts. *Microbiol Mol Biol Rev* **76**: 721-739.
- 606 Moran BM, Payne C, Langdon Q, Powell DL, Brandvain Y, Schumer M. 2021. The genomic
607 consequences of hybridization. *Elife* **10**.
- 608 Nakao Y, Kanamori T, Itoh T, Kodama Y, Rainieri S, Nakamura N, Shimonaga T, Hattori M, Ashikari T.
609 2009. Genome sequence of the lager brewing yeast, an interspecies hybrid. *DNA Res* **16**: 115-129.
- 610 NextOmics. 2021. NextDeNovo.
- 611 Nguyen HV, Legras JL, Neueglise C, Gaillardin C. 2011. Deciphering the hybridisation history leading
612 to the Lager lineage based on the mosaic genomes of *Saccharomyces bayanus* strains NBRC1948
613 and CBS380. *PLoS One* **6**: e25821.
- 614 Patterson M, Marschall T, Pisanti N, van Iersel L, Stougie L, Klau GW, Schonhuth A. 2015. WhatsHap:
615 Weighted Haplotype Assembly for Future-Generation Sequencing Reads. *J Comput Biol* **22**: 498-
616 509.
- 617 Perez-Traves L, Lopes CA, Querol A, Barrio E. 2014. On the complexity of the *Saccharomyces bayanus*
618 taxon: hybridization and potential hybrid speciation. *PLoS One* **9**: e93729.
- 619 Peris D, Sylvester K, Libkind D, Goncalves P, Sampaio JP, Alexander WG, Hittinger CT. 2014.
620 Population structure and reticulate evolution of *Saccharomyces eubayanus* and its lager-brewing
621 hybrids. *Mol Ecol* **23**: 2031-2045.
- 622 Price MN, Dehal PS, Arkin AP. 2010. FastTree 2--approximately maximum-likelihood trees for large
623 alignments. *PLoS One* **5**: e9490.
- 624 Rainieri S, Kodama Y, Kaneko Y, Mikata K, Nakao Y, Ashikari T. 2006. Pure and mixed genetic lines of
625 *Saccharomyces bayanus* and *Saccharomyces pastorianus* and their contribution to the lager
626 brewing strain genome. *Appl Environ Microbiol* **72**: 3968-3974.
- 627 Rainieri S, Zambonelli C, Kaneko Y. 2003. *Saccharomyces sensu stricto*: systematics, genetic diversity
628 and evolution. *J Biosci Bioeng* **96**: 1-9.
- 629 Ruan J, Li H. 2020. Fast and accurate long-read assembly with wtdbg2. *Nat Methods* **17**: 155-158.

- 630 Saitou N, Nei M. 1987. The neighbor-joining method: a new method for reconstructing phylogenetic
631 trees. *Mol Biol Evol* **4**: 406-425.
- 632 Salazar AN, Gorter de Vries AR, van den Broek M, Brouwers N, de la Torre Cortes P, Kuijpers NGA,
633 Daran JG, Abeel T. 2019. Chromosome level assembly and comparative genome analysis
634 confirm lager-brewing yeasts originated from a single hybridization. *BMC Genomics* **20**: 916.
- 635 Scannell DR, Byrne KP, Gordon JL, Wong S, Wolfe KH. 2006. Multiple rounds of speciation associated
636 with reciprocal gene loss in polyploid yeasts. *Nature* **440**: 341-345.
- 637 Shafin K, Pesout T, Lorig-Roach R, Haukness M, Olsen HE, Bosworth C, Armstrong J, Tigyi K, Maurer
638 N, Koren S et al. 2020. Nanopore sequencing and the Shasta toolkit enable efficient de novo
639 assembly of eleven human genomes. *Nat Biotechnol* **38**: 1044-1053.
- 640 Simao FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. 2015. BUSCO: assessing
641 genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics* **31**:
642 3210-3212.
- 643 Stanke M, Keller O, Gunduz I, Hayes A, Waack S, Morgenstern B. 2006. AUGUSTUS: ab initio
644 prediction of alternative transcripts. *Nucleic Acids Res* **34**: W435-439.
- 645 Suvorov A, Kim BY, Wang J, Armstrong EE, Peede D, D'Agostino ERR, Price DK, Waddell P, Lang M,
646 Courtier-Orgogozo V et al. 2022. Widespread introgression across a phylogeny of 155 Drosophila
647 genomes. *Curr Biol* **32**: 111-123 e115.
- 648 Tamura K, Stecher G, Kumar S. 2021. MEGA11: Molecular Evolutionary Genetics Analysis Version 11.
649 *Mol Biol Evol* **38**: 3022-3027.
- 650 Taylor SA, Larson EL. 2019. Insights from genomes into the evolutionary importance and prevalence of
651 hybridization in nature. *Nat Ecol Evol* **3**: 170-177.
- 652 Thomson JM, Gaucher EA, Burgan MF, De Kee DW, Li T, Aris JP, Benner SA. 2005. Resurrecting
653 ancestral alcohol dehydrogenases from yeast. *Nat Genet* **37**: 630-635.
- 654 Todd RT, Braverman AL, Selmecki A. 2018. Flow Cytometry Analysis of Fungal Ploidy. *Curr Protoc*
655 *Microbiol* **50**: e58.
- 656 van den Broek M, Bolat I, Nijkamp JF, Ramos E, Luttik MA, Koopman F, Geertman JM, de Ridder D,
657 Pronk JT, Daran JM. 2015. Chromosomal Copy Number Variation in *Saccharomyces pastorianus*
658 Is Evidence for Extensive Genome Dynamics in Industrial Lager Brewing Strains. *Appl Environ*
659 *Microbiol* **81**: 6253-6267.
- 660 Vaser R, Šikić, M. 2019. Yet another de novo genome assembler. In *International Symposium on Image*
661 *and Signal Processing and Analysis (ISPA)*, pp. 147-151.
- 662 Vaser R, Sovic I, Nagarajan N, Sikic M. 2017. Fast and accurate de novo genome assembly from long
663 uncorrected reads. *Genome Res* **27**: 737-746.
- 664 Wertheimer NB, Stone N, Berman J. 2016. Ploidy dynamics and evolvability in fungi. *Philos Trans R Soc*
665 *Lond B Biol Sci* **371**.
- 666 Wick RR, Judd LM, Gorrie CL, Holt KE. 2017. Completing bacterial genome assemblies with multiplex
667 MinION sequencing. *Microb Genom* **3**: e000132.
- 668 Wolfe KH. 2015. Origin of the Yeast Whole-Genome Duplication. *PLoS Biol* **13**: e1002221.
- 669 Wolfe KH, Shields DC. 1997. Molecular evidence for an ancient duplication of the entire yeast genome.
670 *Nature* **387**: 708-713.
- 671 Zhang Z. 2022. KaKs_Calculator 3.0: Calculating Selective Pressure on Coding and Non-coding
672 Sequences. *Genomics Proteomics Bioinformatics* **20**: 536-540.
- 673 Zheng Z, Li S, Su J, Leung AW, Lam TW, Luo R. 2022. Symphonizing pileup and full-alignment for
674 deep learning-based long-read variant calling. *Nat Comput Sci* **2**: 797-803.
- 675

Supplemental Materials



	Passing basecaller QC	After Porechop
Sequences	192,999	192,837
Bases	2,052,301,802 bp	2,036,617,616 bp
Mean length	10,634 bp	10,561 bp
Median length	6,016 bp	5,949 bp
N50	21,807 bp	21,844 bp
Longest sequence	158,341 bp	158,255 bp

Supplemental Figure S1: Overview of *S. bayanus* CBS 380 genome sequencing results using Oxford Nanopore's MinION. The sequencing effort generated 2.2 Gb of data, achieving approximately 170-fold coverage. After quality control, 2.04 Gb of data were retained. Notably, the sequencing run produced 100 reads surpassing 100 Kb, with the longest read measuring 158,255 bp.

Supplemental Table S1: Comparative analysis of genome assembly tools used for *S. bayanus* CBS 380. This table summarizes the performance of eight different assembly algorithms, including Flye, NextDeNovo, WTDBG, Ra, NECAT, SmartDeNovo, Canu, and their optimized versions (denoted by an asterisk), with an emphasis on addressing diploidy complexities.

Assembler	Flye	NextDeNovo	WTDBG	Ra	NECAT	SmartDeNovo	Canu	SmartDeNovo*	Flye*
Length (bp)	14,635,841	11,870,454	11,811,864	12,017,687	15,428,960	11,578,645	16,115,554	12,391,933	13,538,319
Total seq. #	38	17	21	24	32	16	49	21	62
Max seq.	1,286,589	1,293,072	2,164,800	1,101,538	1,293,344	1,293,262	1,292,830	1,101,177	1,285,351
Min seq.	2,346	122,189	17,265	8,314	59,794	61,853	10,281	59,343	2,074
Med seq len.	301,095	773,390	462,238	528,214	442,886	780,948	213,365	610,360	84,583
N50 (length)	647,014	816,895	813,059	785,083	646,646	919,172	612,488	787,374	785,325
N50 (ctg #)	9	6	5	7	9	6	9	7	7
BUSCO	C:96.6%[S:74.5% %D:22.1%],F:2.3%,M:1.1%	C:96.6%[S:74.5% %D:22.1%],F:2.3%,M:1.1%	C:96.6%[S:74.5% %D:22.1%],F:2.3%,M:1.1%	C:96.6%[S:74.5% %D:22.1%],F:2.3%,M:1.1%	C:96.6%[S:74.5% %D:22.1%],F:2.3%,M:1.1%	C:96.6%[S:74.5% %D:22.1%],F:2.3%,M:1.1%	C:96.6%[S:74.5% %D:22.1%],F:2.3%,M:1.1%	C:96.6%[S:74.5% %D:22.1%],F:2.3%,M:1.1%	C:96.6%[S:74.5% %D:22.1%],F:2.3%,M:1.1%

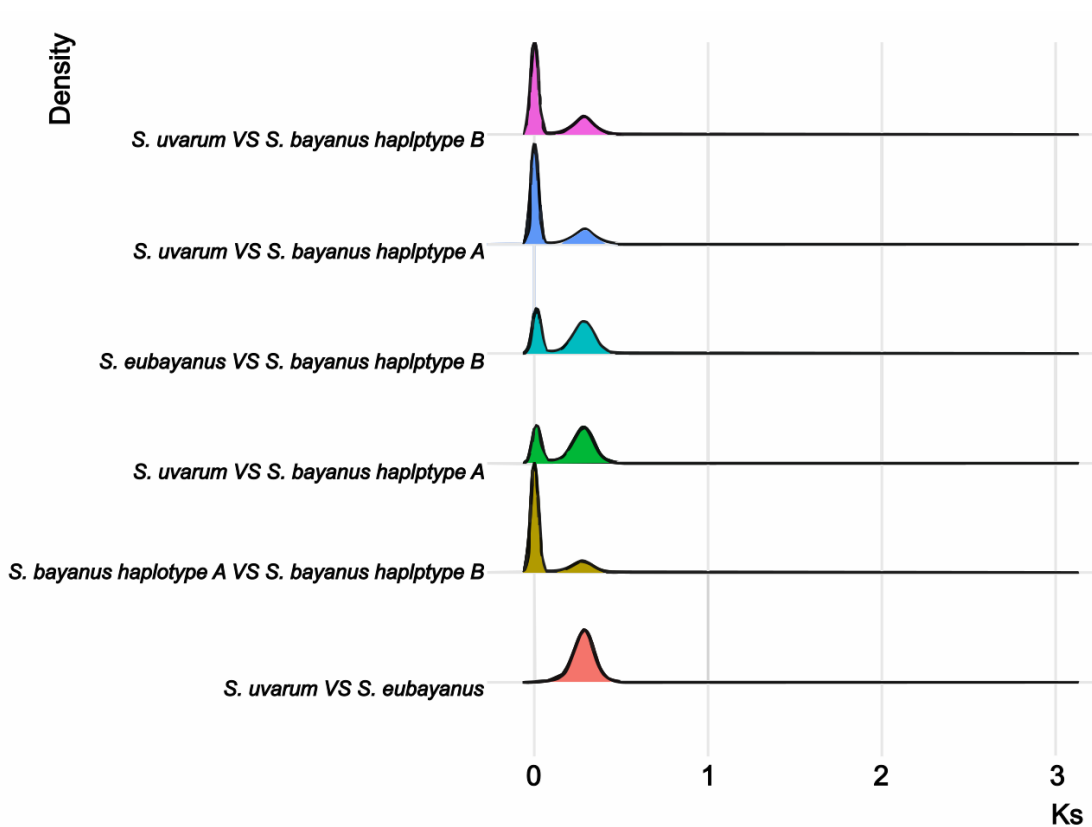
* Ran with Canu-corrected and trimmed reads

Supplemental Table S2: Our methodology successfully assembled two distinct subgenomes, technically designated as haplotype-a and haplotype-b, allowing for a sophisticated analysis of the dual genome architecture. This table shows the different size of each chromosome in subgenomes.

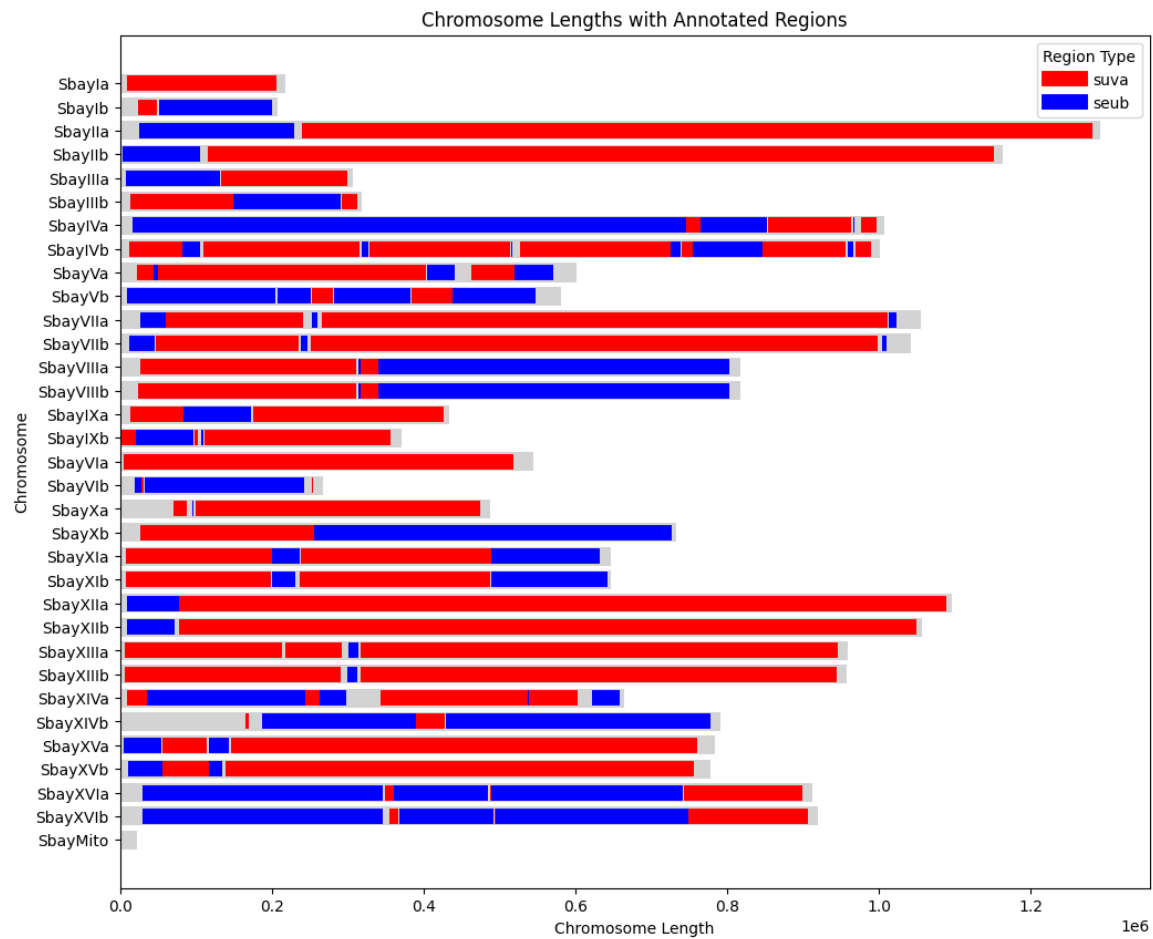
Chr.	Haplotype-a	Haplotype-b
------	-------------	-------------

I	217,795	208,383
II	1,292,201	1,163,801
III	306,526	319,184
IV	1,007,622	1,000,830
IX	433,831	370,629
V	601,693	581,353
VI	544,124	266,821
VII	1,055,917	1,042,311
VIII	817,384	817,171
X	488,404	732,591
XI	646,725	646,334
XII	1,096,630	1,057,074
XIII	959,152	958,029
XIV	664,155	791,919
XV	784,543	778,848
XVI	912,922	919,249
Mt	64,655	

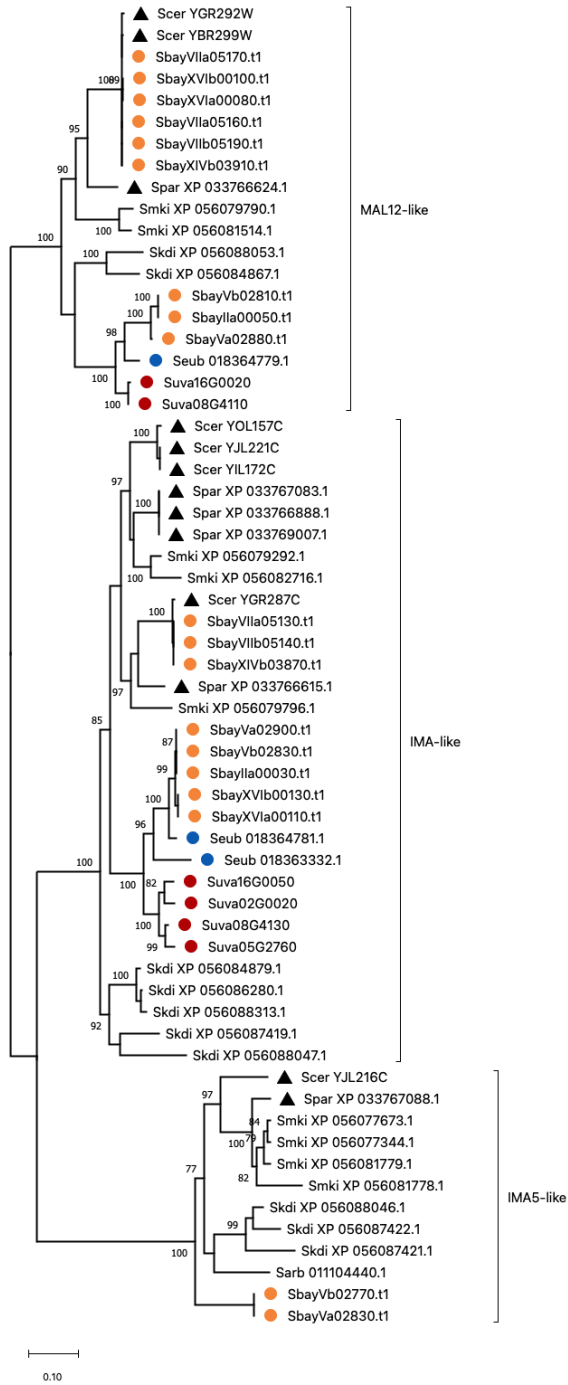
Supplemental Data D1: A list of 5,497 orthologous groups and gene members from the three species are available at <https://github.com/BioHPC/Saccharomyces-bayanus>.



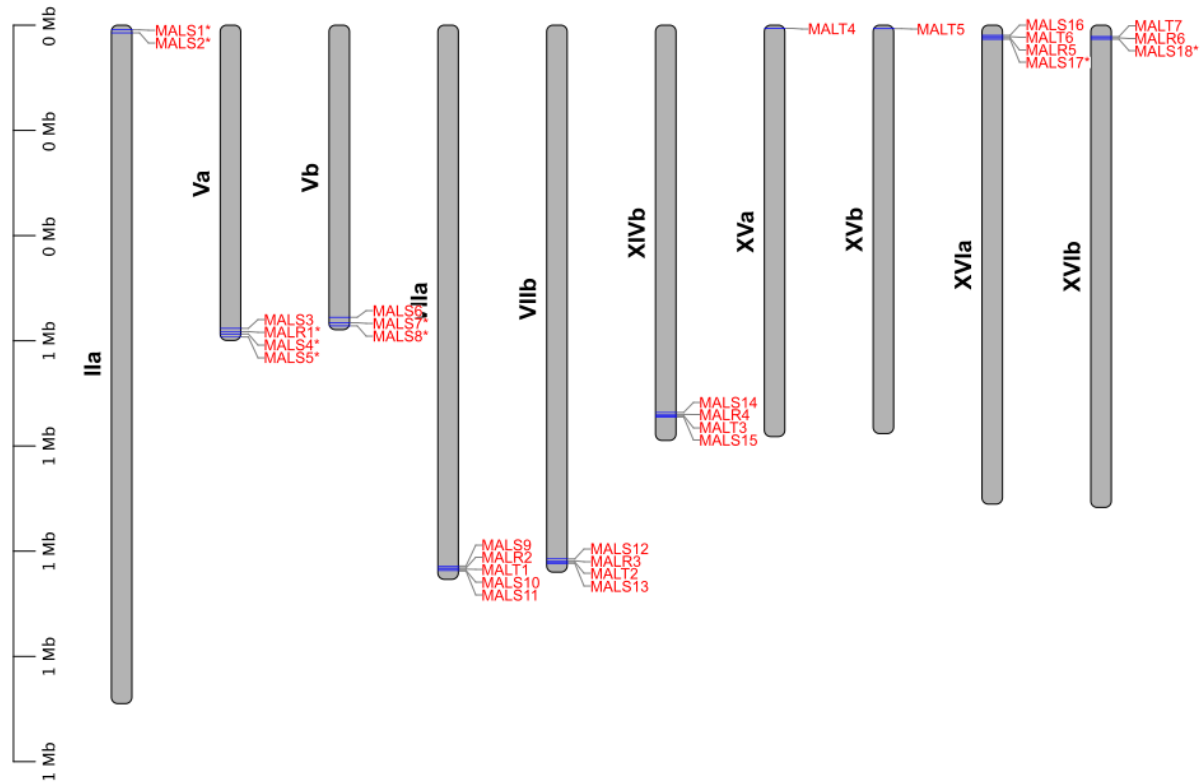
Supplemental Figure S2: Distribution of K_s values.



Supplemental Figure S3: Genomic origin based on gene K_s values.



Supplemental Figure S4: A phylogenetic tree of the MALS gene family.



Supplemental Figure S5: The location of MAL gene families. Gene with star means inherited from *S. eubayanus*.