# Genome sequence assembly and annotation of *MATA* and *MATB* strains of *Yarrowia lipolytica*

Narges Zali[1,2], Osama El Damerash[1], Kapeel Chougule[1], Zhenyuan Lu [1],Doreen Ware[1,3] and Bruce Stillman[1]

[1] Cold Spring Harbor Laboratory, 1 Bungtown Road, Cold Spring Harbor, NY 11724, USA

[2] Graduate Program in Genetics, Stony Brook University, Stony Brook, NY 11794, USA

[3] USDA-ARS Robert W. Holley Center for Agriculture and Health, Ithaca, NY, 14853, USA

Correspondence: Bruce Stillman: email stillman@cshl.edu, phone +1 (516) 367-8383

**ABSTRACT**

Yeast is commonly utilized in molecular and cell biology research, and *Yarrowia lipolytica* is favored by bio-engineers due to its ability to produce copious amounts of lipids, chemicals, and enzymes for industrial applications. *Y. lipolytica* is a dimorphic yeast that can proliferate in aerobic and hydrophobic environments conducive to industrial use. However, there is limited knowledge about the basic molecular biology of this yeast, including how the genome is duplicated and how gene silencing occurs. Genome sequences of *Y. lipolytica* strains have offered insights into this yeast species and have facilitated the development of new industrial applications. Although previous studies have reported the genome sequence of a few *Y. lipolytica* strains, it is of value to have more precise sequences and annotation, particularly for studies of the biology of this yeast. To further study and characterize the molecular biology of this microorganism, a high-quality reference genome assembly and annotation has been produced for two related *Y. lipolytica* strains of the opposite mating type, *MATA* (E122) and *MATB* (22301-5). The combination of short-read and long-read sequencing of genome DNA and short-read and long-read sequencing of transcript cDNAs allowed the genome assembly and a comparison with a distantly related *Yarrowia* strain.

**INTRODUCTION**

*Y. lipolytica* is an ascomycete yeast belonging to the class *Saccharomycetes* that proliferates in hydrophobic environments rich in lipids and proteins, in part due to its remarkable lipolytic and proteolytic abilities (1). *Yarrowia* therefore proliferates in environments rich in lipids and proteins, such as meat and dairy products, particularly fermented ones like cheeses and dry meats, as well as sewage or oil-polluted waters (2–4). *Yarrowia* is significantly different from other hemiascomycetous yeasts in terms of its genomic features. For instance, it is a heterothallic yeast with two distinct mating types, *MATA* and *MATB,* and most natural isolates of this yeast are predominantly haploid (5). Additionally, the G/C content in *Y. lipolytica* is notably high, averaging at 49% and reaching nearly 53% in genes, compared to other yeasts, particularly compared to the commonly used *Saccharomyces cerevisiae* with a genome containing 38.3% G/C (6). *Y. lipolytica* and *S. cerevisiae* are estimated to have a common ancestor that existed 300 million years ago (7), making this pair of yeasts attractive for comparing the evolution of fundamental biological processes such as genome function and DNA replication. *Yarrowia* also has an unusually high number of intron-containing genes compared to *S. cerevisiae* and its related

41 species (8). The number of genes in *Yarrowia* is within the typical range for hemiascomycetous yeasts,
42 but its genome size is 1.7 times larger than that of *S. cerevisiae,* which contains approximately the
43 same number of genes.

44 Comparisons of the potential mechanisms of initiation of DNA replication in the budding yeasts, notably
45 the predictions of DNA sequence-specific origins of DNA replication, have implicated major differences
46 between *Y. lipolytica* and *S. cerevisiae* (9). *S. cerevisiae* and a small clade of highly related budding
47 yeasts have origins of DNA replication that are composed of extensive DNA sequence-specific
48 elements that constitute a functional origin of DNA replication (10, 11)). The origin DNA elements in *S.*
49 *cerevisiae* are recognized by the Origin Recognition Complex (ORC) and Cdc6, and two essential
50 initiator subunits of ORC confer base-specific interactions via an inserted alpha-helix in the Orc4
51 subunit (Orc4-IH) and a loop in the Orc2 subunit (Orc2-loop) (9, 12). In contrast, *Yarrowia*, like nearly
52 all other fungi and all animals and plants, lack these DNA sequence-specific recognition domains in
53 Orc2 and Orc4, suggesting that it may have a more relaxed DNA sequence-specificity at its origins of
54 DNA replication. Moreover, *S. cerevisiae* has lost RNA interference (RNAi) mechanisms but has gained
55 Silent Information Regulator (SIR) proteins (Sir1, Sir3, and Sir4) that function in gene silencing and
56 suppression of recombination of repetitive DNA sequences such as ribosome DNA (rDNA) and
57 telomeric DNA (13–15). Interestingly, *Y. lipolytica* lacks both RNAi and SIR-dependent gene-silencing
58 proteins, including the RNAi silencing proteins such as Dicer and Argonaute (Ago) and SIR proteins,
59 except for Sir2 (9)). All eukaryotes harbor the Sir2 gene that encodes an NAD-dependent histone
60 deacetylase (16). Of relevance in *S. cerevisiae* and its related budding yeasts such as *Kluyveromyces*
61 *lactis,* the interaction between ORC and SIR proteins and a role for ORC silencing the mating type
62 genes (17–21)). It is, therefore, not known how *Y. lipolytica* silences gene expression or suppresses
63 recombination of repetitive rDNA and telomeric DNA.

64 One possible explanation for the occurrence DNA sequence-specific origins of DNA replication in some
65 budding yeast species, such as *S. cerevisiae*, is that these organisms have lost much of the intergenic
66 DNA and lack introns, therefore they possess a very gene-dense genome relative to their genome size.
67 The presence of DNA sequence-defined origins in gene-rich organisms, such as *S. cerevisiae,* could
68 provide an advantage in recruiting ORC to intergenic sites within these species, thereby avoiding
69 conflicts between DNA transcription and the initiation of DNA replication, which can result in genome
70 instability (9, 22)). As a result, organisms like *S. cerevisiae*, with high gene density and smaller
71 intergenic regions, have evolved an efficient mechanism to ensure that the replication complex can find
72 appropriate sites for initiation and avoid initiating DNA replication in a transcribed region. By gaining a
73 deeper understanding of how DNA replication is initiated in a variety of species, including human cells
74 and in diverse yeasts such as *Y. lipolytica*, further insights into how origins of DNA replication are
75 located in the genome and the replication strategies in eukaryotic cells will become apparent (23)). For
76 this reason, the complete genome sequences and assemblies of two *Y. lipolytica* strains of opposite
77 mating type were performed to assist in subsequent studies of whole genome DNA replication and
78 gene silencing mechanisms. Both long-read sequencing using PacBio and Oxford nanopore methods
79 and short-read Illumina-based methods of both genomic DNA and cDNA were used to generate and
80 assemble the genome of the two *Y. lipolytica* strains.

81 Previously, the main reference genome for *Y. lipolytica* was strain CLIB122 (also called E150), a
82 derivative from a mating between a French isolate W29 and an American isolate YB423-12 (CBS 6124-

83    2) (see Figure 1), that was obtained through short read Sanger sequencing (6, 24)). A number of other
84    strains have been shot-gun sequenced and their genomes compared, resulting in a recent summary of
85    these genome comparisons (25)). The estimated size of the six-chromosome genome was ~21 Mb
86    (24)). However, these assemblies contained gaps, and the telomeric ends as well as rDNA repeats
87    could not be integrated into the genome assembly due to their repetitive nature. The CLIB122/E150
88    strain and its derivatives, including the commonly used PO1f strain (Figure 1) (26), are the main strains
89    used for industrial purposes. A high-quality, near-contiguous genome assembly of a distantly related *Y.*
90    *lipolytica* strain DSM 3286 (a German strain) was obtained using a combination of long-read and short-
91    read genomic DNA sequencing (27). This allowed the characterization of the repetitive rDNA and
92    telomeric regions and the observation that rDNA clusters are located near the telomeric regions.
93    Additionally, the genetic and phenotypic diversity of 56 haploid strains of *Y. lipolytica* was investigated
94    by sequencing of a diverse set of *Y. lipolytica* strains collected from various geographical and biological
95    origins, and included revision of the version of the E150 *Y. lipolytica* strain genome sequence and
96    annotation (25).
97
98    In this study, we have used both long and short-read sequencing of genomic DNA and cDNA copies of
99    RNA transcripts to precisely compare the genomes of genetically related *MATA* and *MATB* strains that
100   are distant from the DSM 3286 strain isolated in Germany and other geographically diverse strains that
101   have been sequenced (25, 28)). Chromosomal rearrangements were observed comparing the multiple
102   isolates. The sequences allowed annotation of the genome and revealed the presence of many
103   repeated DNA sequences such as transposable elements, including LTR-retrotransposons, LINE
104   elements, and DNA transposons from various families, which are distributed variably among strains.
105   Moreover, the distribution of and number of rDNA repeats was analyzed.
106

107   **MATERIAL AND METHODS**

108   Strains: Two related *Yarrowia lipolytica* strains of opposite mating type were obtained from Richard A.
109   Rachubunski, University of Alberta, Canada, and single clones were isolated and used for both genome
110   and transcript sequencing. The strains were 22301-5 (*MATB*) and E122 (*MATA*, alternatively called
111   CLIB120) (5, 29, 30)  (Figure 1). E122 is *MATA*, *ura3-302*, *leu2-270*, *lys8-11* and is related to the *MATB*
112   strain E150 (CLIB122 and is the current DNA sequence reference strain (6, 24)). Strain 22301-5 is
113   *MATB, his-1, uras-302, leu2-270,* (Figure 1).
114

115   **DNA and RNA Preparation and Sequencing**
116

117   *RNA:* To isolate the high molecular weight RNA, the TRIzol Plus RNA Purification Kit from Thermo
118   Fisher Scientifics was used. Briefly, 5ml liquid of yeast culture was grown to an $OD^{600}$ of 2.0 and cells
119   were pelleted and lysed with TRIzol™ reagent according to the user manual. Following lysis, the RNA
120   present in the sample was bound to the PureLink RNA Mini Kit Spin Cartridge (12183018A, Invitrogen)
121   where it was washed to remove contaminants. Lastly, eluted RNA was stored in 50-microliter aliquots
122   at -80°C.  To generate a short-read RNA, a Direct-zol RNA purification kit from Zymo Research (Cat #
123   R2050) was used according to the user manual. Long-read sequencing of RNA transcripts was
124   performed as follows. The RNA was prepared using the ONT SQK-PCS109 kit according to the
125   manufacturer's instructions, and it was loaded onto a PromethION P24 system with a PROM-0002 flow

126  cell. Base-calling was performed using the live hac base-calling guppy version 3.2.10. Two cells of
127  *MATA* and *MATB* were run for each experiment.

129  *DNA:* For the generation of high molecular weight DNA and ultra-long nanopore sequencing reads,
130  cells in a 100 ml culture grown overnight at 30°C in Yeast extract, Peptone and Dextrose (YPD) were
131  harvested, washed in sterile distilled water and incubated for 2 hr at 37° C in 10 ml SEB buffer (0.9 M
132  sorbitol, 0.1 M EDTA, 0.8% β -mercaptoethanol) containing 5 mg Zymolyase 20T (Sunrise Science
133  products, CAT#N0766391). Protoplast formation was monitored by phase contrast microscopy. The
134  protoplasts were then harvested, resuspended in 3 ml TE Buffer (Tris-EDTA, pH 8.0), then 300 μl 10%
135  SDS was added, and the samples were incubated at 65°C for 30 minutes. 1 ml of 5 M potassium
136  acetate was added, and the samples were kept on ice for 1 h. The supernatant was recovered after
137  centrifugation, and DNA was precipitated by adding 0.1 volume of 3 M sodium acetate and 2.5 volumes
138  of ethanol at −20°C for at least 1 h. The DNA was recovered by centrifugation and resuspended in 3 ml
139  TE (31)). Then 100 μg/mL of proteinase K along with 50 μg/mL RNase A were added and the samples
140  were incubated at 37°C for 3 hrs. After centrifugation for 45 min at 12,000 × g and 4°C, the supernatant
141  was collected and transferred to a 2-ml Eppendorf tube. Samples were then extracted two more times
142  with phenol/chloroform/isoamyl alcohol and one final time with chloroform. To precipitate DNA, 2-2.5
143  volumes of 100 % freeze-cold ethanol was added to the aqueous phase along with 1/10 volume of 3 M
144  sodium acetate, mixing by inversion, and samples were incubated at −20°C for at least 1 h. The DNA
145  was recovered by centrifugation for 20 min at 12,000 × g and 4°C, and the pellet was subsequently
146  washed three times with 2 ml of 80% (vol/vol) ethanol. The pellet was then air dried and dissolved in
147  100 μl of Tris-EDTA.

149  DNA fragment length was assessed for molecular weight distributions of genomic DNA samples were
150  evaluated using a Femto Pulse pulse-field capillary electrophoresis system (Agilent). >5ug of DNA was
151  size selected via SRE XS (Circulomics). The full reaction was repaired and end prepped with NEBNext
152  FFPE DNA Repair Buffer and Ultra II End prep kit (NEB).  The reaction was cleaned up with 1X
153  Ampure beads and precipitated with ethanol. DNA was bound to ONT adapter from the SQK-LSK109
154  kit (ONT) via NEBnext Quick T4 ligation module (NEB). DNA was resuspended in SQB buffer (ONT)
155  and loading beads (ONT)  and sequenced on one PromethION 24 cell PROM0002 with a three-day run
156  time.

158  For short-read DNA sequencing, the YeaStar Genomic DNA Kit from Zymo Research (CAT# D2002)
159  was used according to the user manual. DNA sequencing libraries were prepared per the
160  manufacturer's instructions with a Kapa DNA hyperprep kit (Roche CAT #KK8504).  It was loaded on
161  an Illumina MiSeq with a PE150 v2 format.

163  **DNA Sequence Assembly**
164  *Long-Read processing:* The unprocessed long-reads were produced using Guppy v.5 base-caller from
165  Oxford Nanopore Technologies (https://github.com/nanoporetech).  To assemble the reads, the long-
166  read assembly pipeline Flye v. 2.9-b1774 (https://github.com/fenderglass/Flye) (32)  was used in nano-
167  hq mode, which is intended for high-quality reads (<5% error rate). The minimum overlap between
168  reads was set to 7KB. The pipeline was run with five iterations of polishing.

169    *Short-Read processing:* The paired-end reads were trimmed using Cutadapt v.3.7 (33). Cutadapt
170    removes adapter sequences from high-throughput sequencing reads (33)).  BWA v.0.7.17-r1188 (34)
171    was used to index the long-read assembly and align the trimmed short-reads to the assembly. The
172    alignments were sorted and indexed using Samtools v.1.14 (35).  Pilon v1.24
173    (https://github.com/broadinstitute/pilon) (36)  was used for polishing the long-read assembly with the
174    aligned short-reads. We obtained exactly one contig per chromosome and mitochondria for *MATB* and
175    an extra contig for *MATA*. To scaffold the extra contig, we used RaGOO (37)  and the assembled
176    *MATB* as the reference genome sequence.

177    **Transcript assembly**
178
179    For the transcriptome assembly, three different transcriptome assemblies were combined: one from the
180    short-read sequencing, a second one from the Nanopore long cDNA reads, and a third combining both.
181    The short-reads were first trimmed using Trimmomatic v.0.38 (38)).  The trimmed reads were aligned to
182    the genome using hisat2 (v.2.2.1).  A short-read only transcriptome was then assembled using Stringtie
183    v.1.3.6 (39)). The long cDNA reads were *de novo* assembled using Oxford Nanopore's Workflow
184    Transcriptomes (wf-transcriptomes) pipeline (v1.1.1). A third transcriptome that used long- and short-
185    reads was assembled using TASSEL (40)  (https://github.com/kainth-amoldeep/TASSEL). These three
186    transcriptomes were then combined using gff compare (v.0.12.2) (41)  and the output was used as
187    input to the annotation pipeline.

188    **Gene annotations**

189    The MAKER-P (v.3.0) (42) pipeline was used to annotate protein-coding genes in the two strains
190    22301-5 (*MATB*) and E122 (*MATA*). As evidence, we used all annotated proteins from *Yarrowia*
191    *lipolytica* (budding yeasts) downloaded from the NCBI protein database. These protein sequences were
192    clustered using CDHit-est (v4.6) (43) with parameters (-c 0.95 -n 10 -d 0 -M 3000 -t 1). For transcript
193    evidence, the combined transcriptome from gff compare was used and transcript assembly.The
194    assembled transcripts were checked and filtered for intron retention using Suppa (v2) (44)). For gene
195    prediction, we used Augustus (v3.3) (45, 46) trained on *Y. lipolytica* and FGENESH
196    (http://www.softberry.com) trained on *S. pombe*, respectively. Repeat masking was done using
197    repeatmasker (RepeatMasker Open-4.0) with the ensembl repeat annotation pipeline using parameters
198    (-nolow -gccalc  -species "Fungi"  -engine ncbi). The repeat masked genomes were used in annotation
199    with MAKER-P evidence. Additional improvements to structural annotations were done using PASA
200    (v2.3.3) (47) using the assembled transcriptome and fungal EST from NCBI using query
201    (EST[Keyword]) AND fungi [Organism]. Gene identifiers were assigned using existing nomenclature
202    schema established for *Yarrowia* for each strain (6, 24)). Functional domain identification was
203    completed with InterProScan (v5.38-76.0) (48)). TRaCE (49) was used to assign canonical transcripts
204    based on domain coverage, protein length, and similarity to transcripts assembled by StringTie
205    (v1.3.4a) (39)). Finally, the gene annotations were imported to ensembl core databases, verified, and
206    validated for translation using the ensembl API (50)).

207 Genome comparative analysis was done using 5 *Yarrowia lipolytica* strains including 15 closely related
208 species and outgroups providing the foundation for building protein-based gene trees based on the
209 EnsemblCompara pipeline (51)).

**RESULTS**

211 A hybrid assembly approach was employed by pairing Illumina short-read DNA sequences and high-
212 quality long-read DNA sequences with much-improved base-calling using NANOPORE Technologies'
213 Guppy 5 base caller with very high mean coverage, as indicated in Figure 2A and Table 1. Longer
214 contiguous and quasi-contiguous genome assemblies for *MATB* and *MATA* were obtained, compared
215 to previous assemblies of the related French isolates. The improvements over the reference assembly
216 for strain CLIB122-E150 included the incorporation of telomeric repeats on each chromosome, as well
217 as a decrease in the estimated number of missing essential gene markers according to BUSCO
218 assessment. We assessed the quality of the assemblies using BUSCO version 4.1.1 (52) , which
219 employs AUGUSTUS as the gene predictor in genome mode on the *Saccharomycetes* lineage set of
220 2137 essential genes**.** The results are shown in Figure 2B and indicate a completeness index of 96.8%.
221 *Yarrowia* appears to be missing genes that are present in other *Saccharomycetes,* or the genes are
222 diverged enough so that  BUSCO does not detect them.

223 A dot plot produced with chromeister (release 1.5.a.) (53)  compared the E122 *MATA* and 22301-5
224 *MATB* genomes, revealing very high similarity between these two related genomes, as expected
225 (Figure 3A). The similarity score indicated a divergence of just 0.003. In contrast, a similar comparison
226 between 22301-5 *MATB* and the German isolate DSM 3286 showed considerable genome
227 rearrangements and a reduced similarity score of 0.0031 (Figure 3B), as noted previously (27)). When
228 the genes that were expressed in E122 *MATA* but not 22301-5 *MATB* were analyzed, they were in and
229 surrounding the *MATA* locus and included *Sla2,* encoding an actin binding protein and *Apn2* encoding a
230 DNA-(apurinic or apyrimidinic site) lyase, both genes that flank the *MATA* locus (54, 55)**,** as well as the
231 *MATA1* and *MATA2* mating type genes (55)).

**Gene annotations and comparative analysis**

234 The structural gene annotation pipeline identified 7,728 and 7,769 genes in *Yarrowia lipolytica* strains
235 E122 (*MATA*) and 22301-5 (*MATB*), respectively (Table 2). This gene count surpasses that of
236 previously reported *Yarrowia* strains DSM 3286 (27) with 6,467 protein-coding genes and the reference
237 strain CLIB-122 (56) with 6,448 protein-coding genes. To further evaluate annotation quality, we utilized
238 the Annotation Edit Distance (AED) score generated by MAKER-P (42)). An AED score of 0 indicates
239 genes supported by evidence, while a score of 1 indicates a lack of evidence. Employing mRNA and
240 homology evidence, as described in the methods, to calculate AED scores yielded 6,297 and 6,212
241 genes with some evidence (AED score < 1) in strains E122 *MATA* and 22301-5 *MATB*, respectively. In
242 contrast, strains DSM 3286 and CLIB-122 (E150) exhibited 6,236 and 6,387 protein-coding genes with
243 AED scores <1 (Table 2).

245 When comparing annotation features between the strains E122 *MATA* and 22301-5 *MATB* using genes
246 filtered by AED <1, it was observed that 78% of the protein-coding genes contained a single exon. The

247   shortest introns were ~40 base pairs and the longest intron was 6782 base pairs (Figure 4C and D).
248   Moreover, we found more genes in E122 *MATA* and 22301-5 *MATB* with multiple introns compared to
249   the previously analyzed strains (Figure 4A). The median gene lengths in E122 *MATA* and 22301-5
250   *MATB* were higher than previously estimated, reflecting their increased intron counts compared to the
251   other two strains (Table 2). *Yarrowia lipolytica* genomes are known to be intron-rich, with previous
252   estimates of 15% of genes containing introns, which is 4 times that of *S. cerevisiae* (8)). Intron-
253   containing genes in E122 *MATA*, 22301-5 *MATB* and the DSM 3286 strain represented ~20% of the
254   protein-coding genes, highlighting the improvement in genome assembly using long-read technology
255   compared to CLIB-122. 80% of these intron-containing genes were mono-intronic, compared to 20%
256   that were multi-intronic (with up to five introns). The internal exons of the multi-intronic genes were
257   mostly short compared to 1st intron (Figure 4B and D).
258
259   A total of 9453 unfiltered orthologous genes were found between E122 *MATA*, 22301-5 *MATB*, the
260   DSM 3286 and the CLIB-122 strains (Figure 5A), of which 5,975 were core genes (that were found in
261   all four strain genome sequences (Figure 5B, left hand set). These core genes were closer to the
262   number of core genes (6,042) detected from 7 *Yarrowia lipolytica* strains and slightly lower than the
263   pan-genome genes (6,528) detected with 54 strains (25). Prior studies also suggest that *Yarrowia*
264   *lipolytica* exhibited lower genetic diversity since the core genes were barely different than the pan-
265   genome (25). The *MATA* and *MATB* strains had 1,204 unique gene ortholog groups present in both
266   strains but not present in DSM 3286 or CLIB122 (Figure 5B, second set from left). Subsequent Gene
267   Ontology (GO) analysis of this subset identified only 14 genes with associated GO terms and GO
268   enrichment analysis highlighted significant enrichment in molecular processes including 2 iron, 2 sulfur
269   cluster binding (GO:0051537), methylmalonate-semialdehyde dehydrogenase (acylating) activity
270   (GO:0004491), oxidoreductase activity (GO:0016491) and TBP-class protein binding (GO:0017025)
271   (Figure 5C) (57), as well as biological processes (GO:0006352), for DNA-templated transcription
272   initiation (Figure 5D).
273
274   **Analysis of Repeat Sequences**

275   Using RepeatMasker software, a comparison of the repeated DNA sequences in the E122 *MATA* and
276   22301-5 *MATB* strains, as well as the reference CLIB122 strain and the DSM 3286 strain, showed
277   many more repeats in the two genomes sequenced here (Figure 6A AND 6B). In particular, there is an
278   increase in the number of RNA repeats, LTRs, LINE and SINE repeats in the E122 *MATA* and 22301-5
279   *MATB* strains compared to the DSM 3286 strain. It is not clear whether this difference is due to
280   biological variation or to technical issues with genome sequencing and analysis, but it is likely the latter.
281   The distribution of repeat elements in the E122 *MATA* and 22301-5 *MATB* strains showed similar
282   profiles, with RNA repeats making up about 30%, simple repeats around 35-40%, and LTRs at 15-18%.
283   Both have smaller proportions of Type I Transposons (LINE and SINE) and Type II Transposons, each
284   accounting for about 3-7% of the total repeats. DSM 3286, on the other hand, has a higher percentage
285   of simple repeats (~55%) and a lower proportion of RNA repeats (~10%), with LTRs making up 20% of
286   its repeat content. CLIB122 is distinct with 30% of its repeats being LTRs and about 10% RNA repeats.
287   Across all strains, low complexity regions and unknown repeats remain minimal, each contributing
288   around 1-2% of the total repeats. In summary, while simple and RNA repeats dominate the repeat
289   landscape in these yeast strains, there is significant variability in the proportion of LTRs and other

7

290  transposon types, particularly between DSM 3286, CLIB122, whereas the repeats are more similar in
291  the E122 *MATA* and 22301-5 *MATB* strains. This variability in repeat element distribution probably
292  reflects a combination of technical differences as well as biological variation in genomic evolution
293  among the strains.

294  As previously observed (27) , the rDNA repeats consisting of the 18S and 28S genes were located at
295  the ends of chromosomes B (right end), C (both ends), E (right end) and F (both ends), and lie adjacent
296  to the telomeres (Figure 7A). The 5S rDNA genes are scattered throughout the genome on every
297  chromosome. For the E122 *MATA* and 22301-5 MATB strains, the calculated size of these repeats in
298  kilobase pairs (yellow bar) and the number of rDNA repeats (blue bar) for each region of the genome is
299  shown in Figure 7B.

**DISCUSSION**

301  Heterothallic yeast, like *Yarrowia lipolytica*, typically engage in outcrossing, where genetic material is
302  exchanged between individuals of different mating types. This may lead to greater genetic diversity and
303  adaptability to changing environments and contribute to divergence and speciation (58)).  In contrast,
304  homothallic yeast, such as *S. cerevisiae*, primarily engage in selfing since they can switch their mating
305  type through a gene conversion process initiated by the HO endonuclease (59) , where recombination
306  occurs within the same individual. This process may result in the fixation of beneficial alleles or the
307  accumulation of deleterious mutations, potentially leading to lower genetic diversity (60)). These
308  differences in reproductive strategies lead to distinct recombination pathways in the two types of yeast
309  (61)). However, a recent comparison of 56 shot-gun sequenced strains showed a very low level of
310  genetic diversity, indicating that *Y. lipolytica* may be a species that has recently emerged (25)).

311  *Y. lipolytica* exhibits a remarkably low rate of mating and spore viability between different lineages due
312  to chromosomal rearrangements, which may contribute to its poor fertility. Chromosomal
313  rearrangements in *Y. lipolytica* could have been caused by crossing-over events facilitated by the
314  different types of transposable elements present in the organism (27)). Yeast genomes contain mobile
315  genetic elements, such as transposons and retrotransposons, which can translocate within the
316  genome. These elements can be inserted into new locations within the genome or cause chromosomal
317  reorganization by combining with various regions. The prevalence of repeated sequences found in the
318  E122 *MATA* and 22301-5 *MATB* strains may also play a role in rearranging and evolving the genome of
319  this yeast species, consistent with the notion that transposable elements and other repetitive elements
320  can be significant contributors to genome evolution (25, 27)).

321  We found a higher number of introns in the E122 *MATA* intron (Figure 4). and 22301-5 *MATB* strains
322  compared to the DSM 3286 and CLIB122 strains, with most genes having a single In particular, we
323  observed a higher proportion of genes with more than one intron. This intron-rich genome may enable
324  the production of multiple protein isoforms from a single gene, offering *Yarrowia* the ability to rapidly
325  adapt to changing environments or industrial processes. This could prove especially valuable
326  for *Yarrowia*, as it frequently operates in diverse and challenging growth conditions.

327  The enhanced genome assembly and annotation of *Yarrowia lipolytica* strain E122 *MATA* and 22301-5
328  *MATB*, was made possible by employing a hybrid sequencing approach that combined the precision of

329 Illumina short reads with the depth of Oxford Nanopore long reads and advanced base calling with
330 Guppy 5. This high-quality assembly allowed us to capture telomeric regions, rDNA repeats and
331 improve the completeness of essential gene markers, achieving a BUSCO score of 96.8%. These
332 results establish a strong foundation for further functional and comparative studies on *Y. lipolytica* and
333 its applications.

## Gene Annotation and Genetic Diversity

335 Our comparative gene analysis revealed significant differences between the sequenced E122 *MATA*
336 and 22301-5 *MATB* strains and other previously studied *Yarrowia* strains, such as DSM 3286 and
337 CLIB-122. The increased gene count in E122 *MATA* and 22301-5 *MATB* and the presence of
338 alternative isoforms highlight a potentially broader genetic repertoire and greater regulatory complexity
339 in these strains. This added complexity may reflect adaptive mechanisms developed in response to
340 specific environmental and industrial conditions. For instance, the increased median gene length,
341 attributed to a higher intron count, suggests unique gene structures that could enhance regulatory
342 flexibility, supporting more intricate metabolic or stress-response pathways.

343 The core gene analysis indicates that E122 *MATA* and 22301-5 *MATB* share 5,975 core genes with
344 other *Yarrowia* strains, consistent with prior findings of limited genetic diversity within *Yarrowia lipolytica*
345 (25)). However, identifying 1,204 unique ortholog groups in *MATA* and *MATB* suggests subtle genomic
346 differences that could contribute to strain-specific phenotypes. Gene Ontology (GO) enrichment
347 analysis of these unique genes emphasizes metabolic processes and transcription initiation. These
348 functions are advantageous in environmental settings where efficient resource utilization and
349 adaptability to stressful conditions are beneficial.

## Repeat Elements and Genomic Evolution

351 The investigation into repetitive DNA elements highlights more repeats in E122 *MATA* and 22301-5
352 *MATB*, particularly in RNA, LTRs, LINE, and SINE elements, compared to DSM 3286 and CLIB122
353 (Figure 6B). The distinct repeat profiles observed in E122 *MATA* and 22301-5 *MATB*—with a balance
354 of RNA and simple repeats making up 30-40% of the genome—suggest unique genomic architectures
355 that may influence adaptation. The higher proportions of certain repeat types, particularly simple and
356 RNA repeats, could facilitate rapid genomic changes, enhancing adaptability in dynamic environments
357 like industrial fermentation.

358 These repeat variations among strains may indicate different genomic stability and plasticity strategies.
359 For example, the high LTR content in DSM 3286 may signify historical transposon activity, promoting
360 genomic rearrangements. In contrast, the more balanced and stable repeat landscape in E122 *MATA*
361 and 22301-5 *MATB* suggests a refined evolutionary adaptation that could confer resilience in industrial
362 contexts. The conserved nature of certain repeat types across strains, such as LINE elements,
363 suggests shared functional roles across *Yarrowia* lineages, whereas the unique repeat profiles of E122
364 *MATA* and 22301-5 *MATB* reflect strain-specific evolutionary pressures.

## rDNA Repeats and Size

366 Analysis of the distribution of rDNA repeats in the E122 *MATA* strain revealed distinct patterns in repeat
367 length and counts across multiple chromosome regions, adjacent to the telomeres as shown in Figure
368 7A. Regions chrC-R, chrE-R, and chrF-R exhibit the most extended total rDNA lengths (~11 KB) with
369 moderate counts, suggesting these regions contain larger rDNA repeats or a higher density of sizeable
370 elements. In contrast, chrC-L has a shorter total length (~8 KB) and a lower count, indicating fewer and
371 potentially smaller rDNA repeats. This heterogeneity could indicate region-specific roles or stability
372 requirements for rDNA within the MAT-A strain. We note that the location of the rDNA repeats in the
373 two strains analyzed herein is essentially the same as in DSM 3286, but the estimated number of rDNA
374 repeats differs. We suggest that this difference reflects both technical and biological variation.

375 The stability of the rDNA and telomeric repeats needs explanation. In *S. cerevisiae*, the SIR proteins
376 play an important role in the maintenance of the repeats by preventing recombination (14, 62). *Y.*
377 *lipolytica* lacks the SIR proteins, except for *SIR2*, which is present in all eukaryotes, and it also lacks
378 genes encoding RNAi components that in other eukaryotes suppress gene expression in
379 heterochromatin (16, 63)). This raises the interesting issue of how the rDNA and telomeric repeats
380 resist recombination and thus maintain stability.

381 **GC Content and Evolutionary Implications**. The overall GC content in *Y. lipolytica* is higher than
382 in *S. cerevisiae*, which aligns with the significant evolutionary divergence between these species,
383 estimated at around 300 million years. The consistent GC content across E122 *MATA*, 22301-5 *MATB*,
384 and DSM 3286 (48.9%) compared to the low GC content of the mitochondrial chromosome (22.59%)
385 suggests differences in selective pressures and genome organization between nuclear and
386 mitochondrial genomes. This higher GC content may have implications for DNA stability, transcription
387 efficiency, and DNA replication dynamics, offering insights into the evolutionary and functional
388 constraints on the *Yarrowia* genome. For example, the well-characterized origins of DNA replication in
389 *S. cerevisiae* are AT-rich. We are analyzing the genome replication and origins of DNA replication in
390 *Yarrowia* to determine if the genome has GC-rich origins of DNA replication that are more akin to the
391 GC-rich origins in human cells. The more complete genome sequences of the E122 *MATA* and 22301-
392 5 *MATB* strains should facilitate the analysis of genome replication patterns and mechanisms.

393 **Conclusions and Future Directions**

394 Our findings provide a comprehensive view of the genomic landscape and diversity within *Yarrowia*
395 *lipolytica* strains *MATA* and *MATB*, laying the groundwork for further research in functional genomics
396 and strain optimization. The variability in repeat elements, the distinct genomic organization, and the
397 elevated gene complexity observed in *MATA* and *MATB* highlight the evolutionary and functional
398 divergence within *Y. lipolytica*. The insights gained from future studies of basic molecular biology in *Y.*
399 *lipolytica* will contribute to our understanding of the molecular underpinnings that enable *Yarrowia*
400 *lipolytica* to thrive in highly varied environments, ultimately advancing strain development for
401 biotechnology.

402 **DATA AVAILABILITY**
403 The DNA sequence and annotation data are available at Dryad
404 https://datadryad.org/share/cvheElifiD1l5ooEar-cRwTXDDAmKF8pr3ayrUHA8xg The data are:

405 **_Files and variables_**

406 File: matB_final.fasta    Description**:** Strain 22301-5 Genome assembly

407 File: matA_final.fasta    Description**:** Strain E122 Genome assembly

408 File: matA_AEDcln.gff   Description:   Strain E122 Transcriptome assembly and annotation

409 File: matB_AEDcln.gff   Description:   Strain 22301-5 Transcriptome assembly and annotation

410

411

412

413

414 **AUTHOR CONTRIBUTIONS**

415 N.Z. did the experiments and the DNA sequencing using the Cold Spring Harbor Laboratory

416 Sequencing and Technology Core facility. O. E D., K.P., Z.L. and D.W. performed data analysis. B.S.

417 conceived the project and oversaw all aspects of the research. N.Z. K.P. D.W. and B.S. wrote the

418 paper.

419

420 **ACKNOWLEDGEMENTS**

424

425 **FUNDING**

432

433 **CONFLICT OF INTEREST**

434 The authors declare no conflict of interest.

435

436

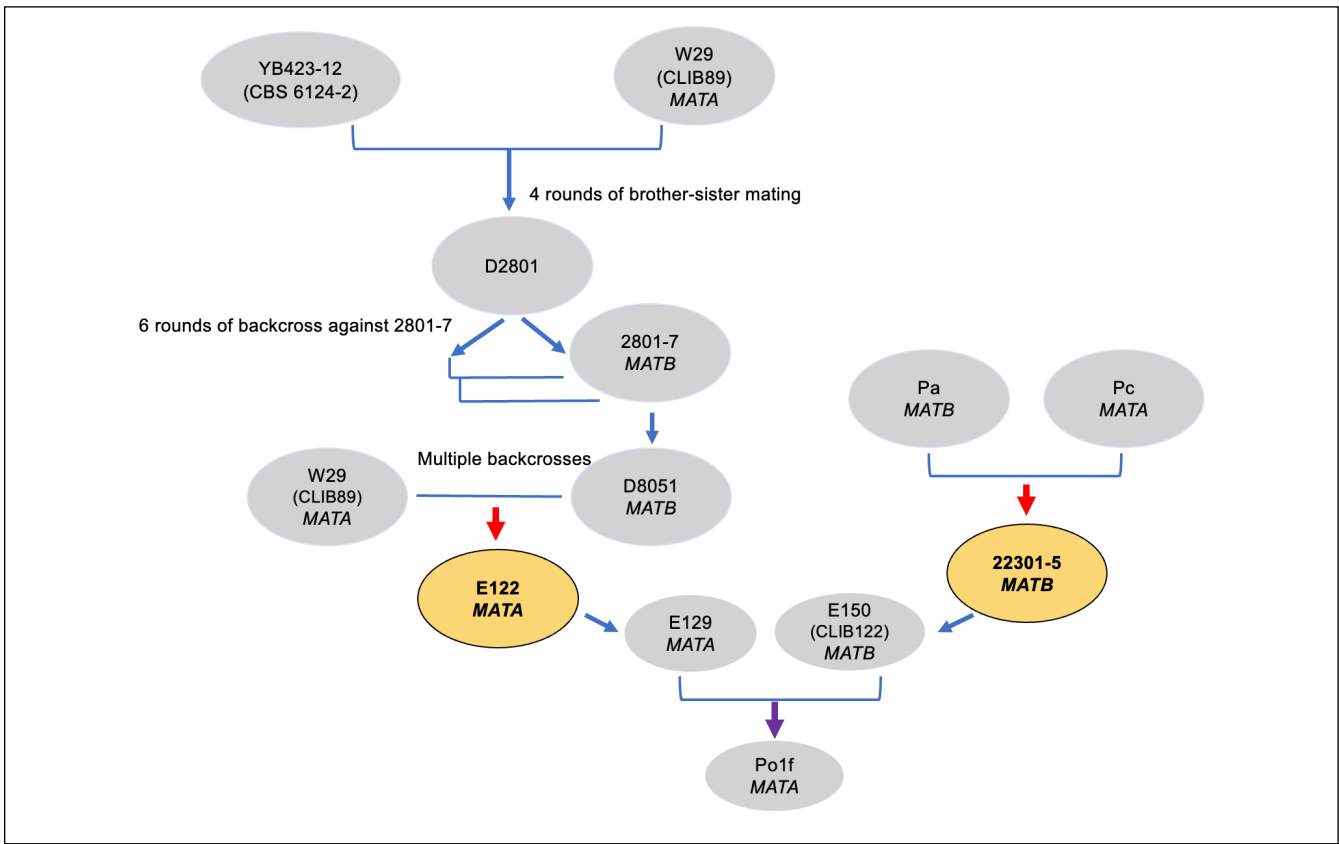437  **Figures and Tables**

438



439
440  **Figure 1**: **Breeding and backcrossing strategy for strain development in *Yarrowia lipolytica*.** The
441  flowchart illustrates the lineage and mating strategy employed to derive key strains of *Y. lipolytica*.
442  YB423-12 lys1.13 and W29 (CLIB89 *MATA*) underwent four rounds of brother-sister mating to generate
443  the intermediate strain D2801. D2801 was subjected to six rounds of backcrossing against 2801-7
444  *MATB*, resulting in the development of the strain D8051 *MATB*. Parallel strategies were employed
445  using Pa *MATB* and Pc *MATA*, which were crossed to form 22301-5 *MATB*. E122 *MATA* was derived
446  from multiple backcrosses and further developed into strains like E129 *MATA*, E150 (CLIB122 *MATB*),
447  and Po1f *MATA*, widely used for research and industrial applications. Color-coded ovals indicate key
448  final strains derived from these processes (e.g., E122 *MATA* and 22301-5 *MATB*). Blue arrows
449  represent mating and backcrossing steps; red arrows indicate the lack of a *Ura3* marker, and purple
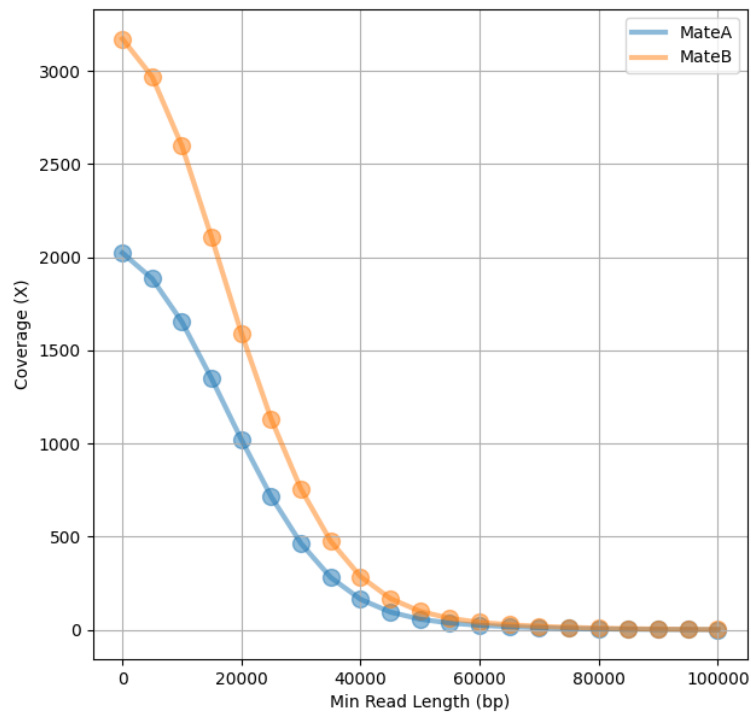450  arrows indicate the lack of *Xpr2* and *Axp* genes.
451

|  | E122 *MATA* | 22301-5 *MATB* |
|---|---|---|
| Total Read Length | 35055999261 | 53578162586 |
| Mean coverage | 1630 | 2364 |
| Reads N50/N90 | 19835 / 6596 | 19758 / 6756 |
| Total Length | 21019611 | 21008502 |
| Fragments N50 | 3712330 | 3688210 |
| Fragments | 8 | 7 |
| Largest fragment | 4317224 | 4320808 |
| Total length | 20313536 | 21008502 |

452
453
454 **Table 1**:**Comparative genome assembly statistics for *Yarrowia lipolytica* strains E122**
455 **(*MATA*) and 22301-5 (*MATB*).**
456 Total Read Length: Total bases sequenced for each strain; Mean Coverage: Average sequencing
457 coverage (sequencing depth); Reads N50/N90: Median (N50) and 90th percentile (N90) read
458 lengths; Total Length: Total assembled genome length in base pairs;
459 Fragments N50: Fragments: Number of assembled fragments; Largest Fragment: Largest assembled
460 fragment in base pairs;  Total Length (Alternate): Total length in base pairs across all fragments. These
461 statistics provide a detailed comparison of genome assembly quality, highlighting the structural integrity
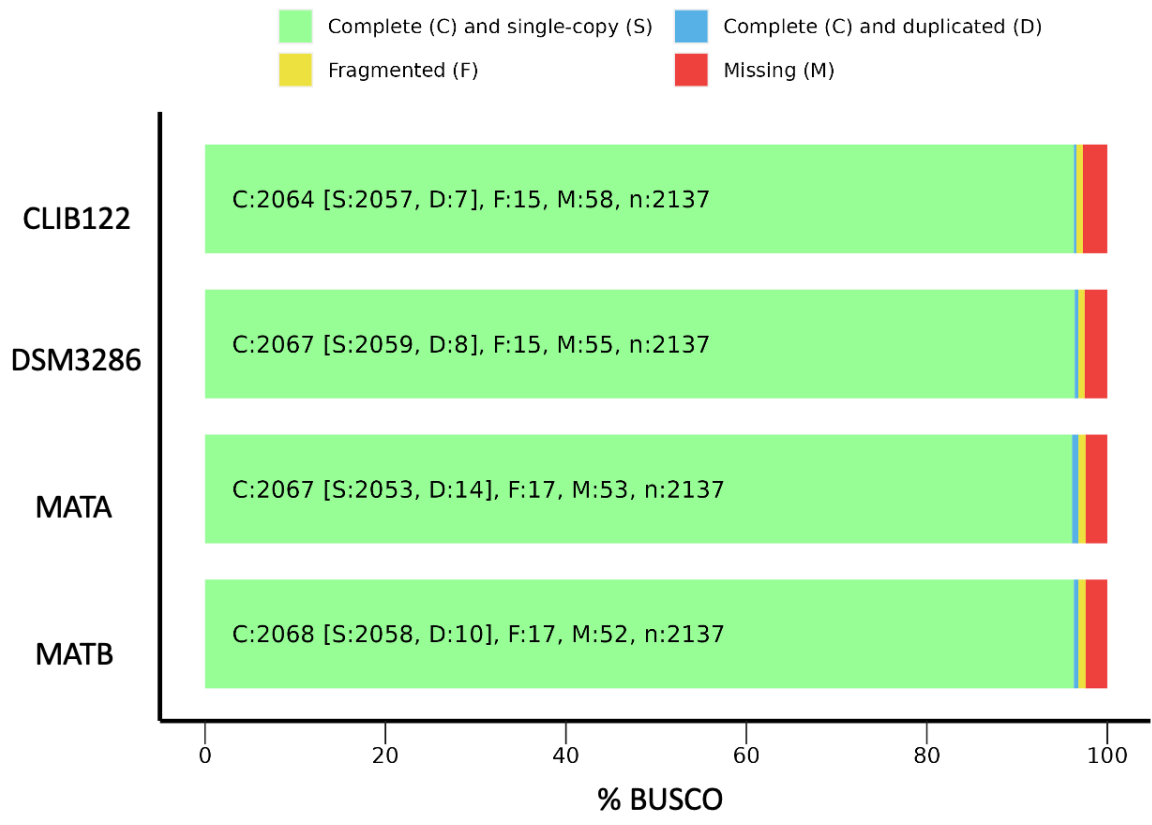462 of the two *Yarrowia lipolytica* strains.
463

464

**Figure 2. Comparison of genome assembly metrics between *Yarrowia lipolytica* strains E122 (*MATA*) and 22301-5 (*MATB*). A.** Coverage versus minimal read length distribution: Coverage (X) is plotted against minimal read lengths for both strains. The blue line represents E122 *MATA*, while the orange line represents 22301-5 *MATB*. The higher coverage for 22301-5 *MATB* indicates a more profound sequencing effort compared to E122 *MATA*. The gradual decline in coverage with increasing read length reflects the expected distribution of read sizes. **B.** BUSCO analysis for genome completeness across four *Yarrowia lipolytica* strains (CLIB122, DSM 3286, E122 *MATA*, and 22301-5 *MATB*). C: Number of complete BUSCO genes, divided into S: Single-copy genes and D: Duplicated genes. F: Number of fragmented BUSCO genes. M: Number of missing BUSCO genes. N: Total number of BUSCO groups analyzed (2137). Each bar represents the percentage distribution of these categories for a strain. Most BUSCO groups are complete and single-copy (light green), reflecting high genome assembly quality. Slight variations in duplicated (blue), fragmented (yellow), and missing (red) categories highlight subtle differences in genome assemblies among the strains.
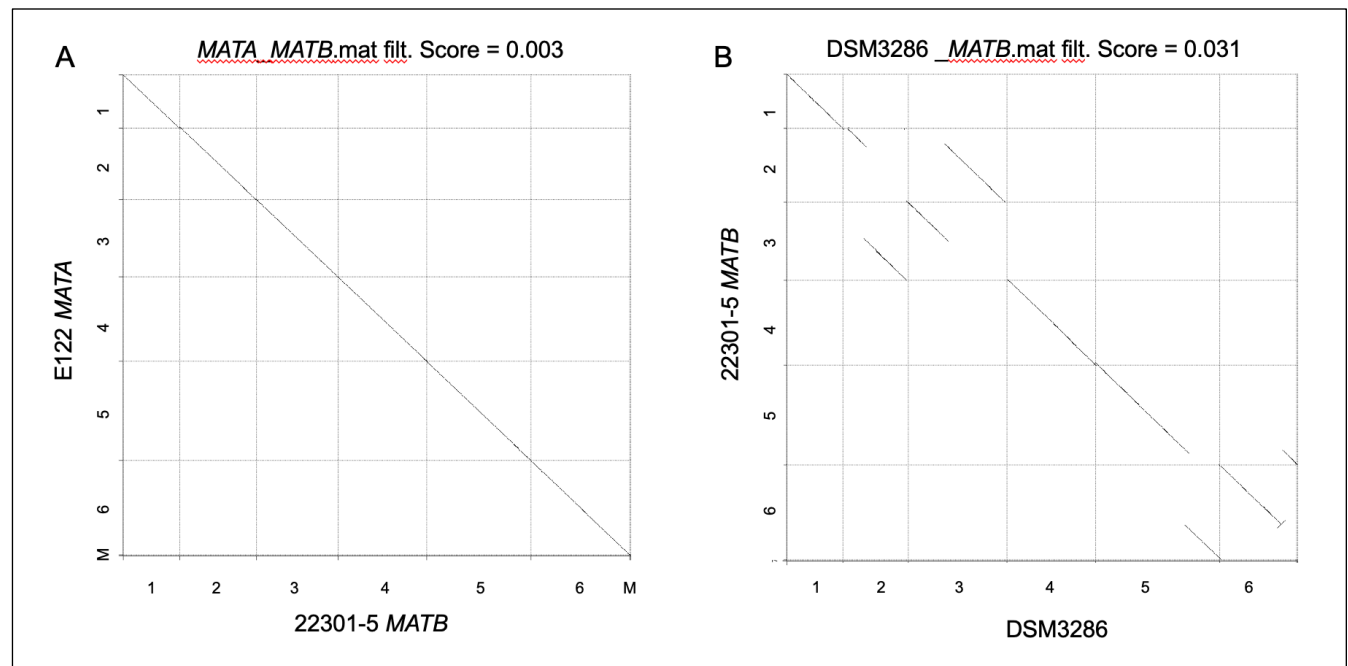
**Figure 3.** Assembly alignment comparison: **A.** Dot plot of the final assemblies for E122 *MATA* and 22301-5 *MATB* shows synteny between the two genomes. The diagonal alignment highlights the high degree of sequence conservation, with a calculated score of 0.003, indicating minor structural or sequence variations between the two genomes. This data highlights the quality of sequencing and genome assembly while comparing structural differences between these two strains of *Y. lipolytica*. **B.** Dot plot of the sequences for 22301-5 *MATB* and DSM 2386. The diagonal alignment highlights the genome rearrangements between the two strains, thereby reflecting a lower calculated score of 0.031.

|  | E122 MATA | 22301-5 MATB | DSM3286 | CLIB89-W29 | CLIB122 |
|---|---|---|---|---|---|
| Gene count | 6,633 | 6,649 | 6,467 | 7,934 | 6,389 |
| Gene length(median) | 1,413 | 1,398 | 1,257 | 1,077 | 1,245 |
| Exon count | 8,576 | 8,658 | 7,915 | 10,335 | 7,460 |
| Exon length(median) | 1,026 | 1,279 | 1,020 | 690 | 1,059 |
| Intron count | 1,948 | 2,013 | 1,448 | 2,401 | 1,071 |
| Intron length(median) | 224 | 220 | 196 | 61 | 206 |
| Peptide count | 6,677 | 6,702 | 6,467 | 7,934 | 6,389 |
| Peptide length(median) | 394 | 395 | 402 | 342 | 403 |
| Exons per transcript(Avg) | 1.3 | 1.3 | 1.2 | 1.3 | 1.2 |
| Single-exon gene count(%) | 5,024 (75.7) | 5,023 (75.5) | 5,263 (81.4) | 5,852 (73.6) | 5,442 (85.2) |

489

490 **Table 2: Comparative Gene and Transcriptomic Features of *Yarrowia lipolytica* Strains.**
491 This table presents the genomic and transcriptomic characteristics of four *Yarrowia lipolytica* strains
492 (E122 *MATA*, 22301-5 *MATB*, DSM3286, and CLIB122), highlighting differences in gene structure and
493 composition. The *Yarrowia lipolytica* strains exhibit notable variations in gene structure, with 22301-5
494 *MATB* having the highest gene count (7,769) and DSM3286 the lowest (6,439), while E122 *MATA* has
495 the longest median gene length (1,473 bp). The high proportion of single-exon genes (78.2–85.2%) and
496 low exons per transcript (~1.2–1.3) suggest a predominantly condensed genetic structure, with
497 CLIB122 having the most compact gene architecture and E122 MATA showing greater structural
498 complexity. This data provides insights into structural genomic variation and transcriptional complexity
499 across these *Y. lipolytica* strains.
500
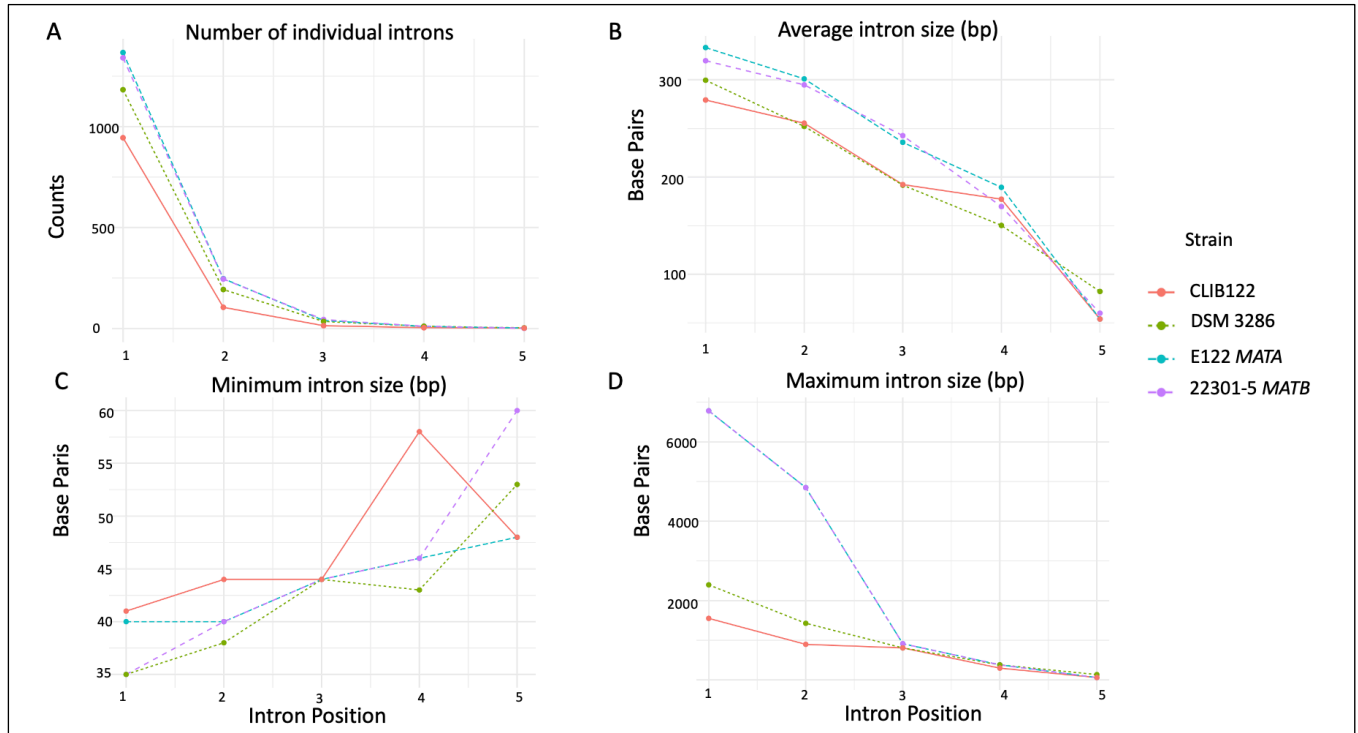
501
502
503
504
505 **Figure 4. Comparative analysis of intron properties across *Yarrowia lipolytica* strains**
506 **CLIB122, DSM 3286, E122 *MATA*, and 22301-5 *MATB*.**
507 **A.** A line plot of the number of individual introns observed at each gene position. The first intron position
508 is the most prevalent across all strains, with a sharp decline in frequency for subsequent positions. **B.**
509 Line plot showing the average size of introns for each strain across different intron positions. 22301-5
510 *MATB* exhibits the highest average intron size for more 3' intron positions compared to other strains. **C.**
511 Line plot of the minimum size of introns at each intron position within the genes, with E122 *MATA* and
512 22301-5 *MATB* showing an increase in size for more 3' intron positions compared to other strains. **D.**
513 Line plot showing the maximum intron size for each intron position within the genes. E122 *MATA* and
514 22301-5 *MATB* shows the highest number of introns for the initial positions, followed by a steep decline.
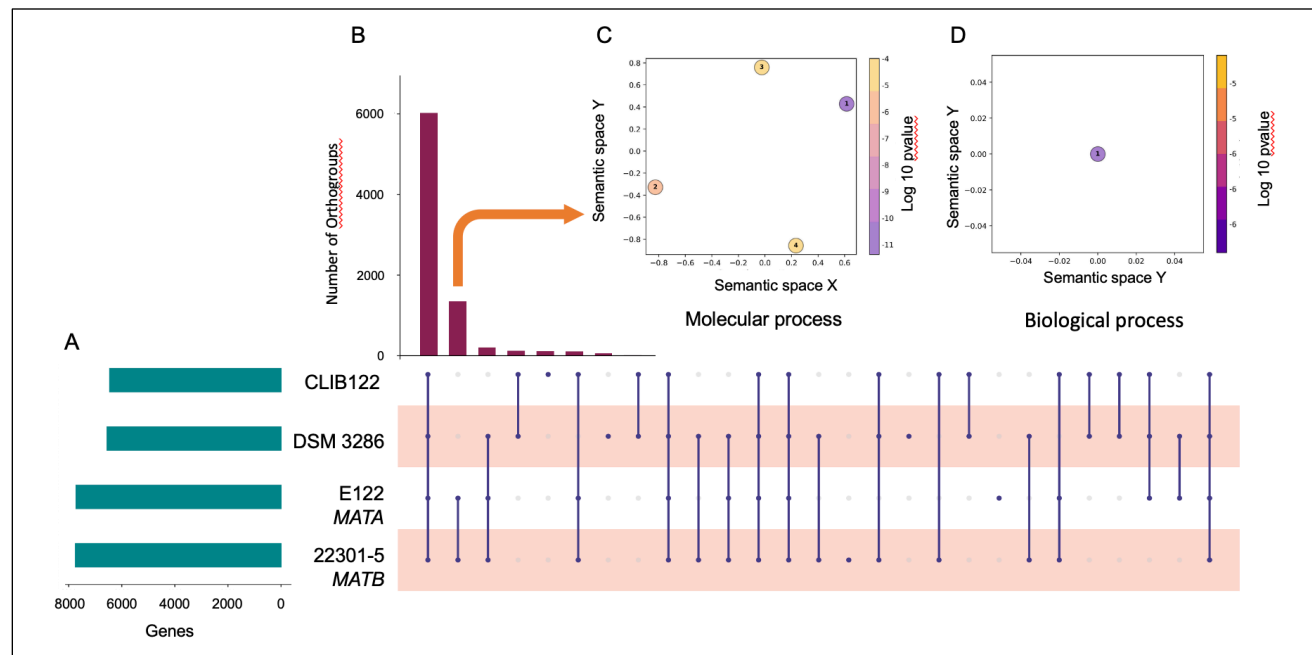515
516

**Figure 5. Comparative analysis of orthologous gene groups, functional enrichment, and gene distributions among four *Yarrowia* strains: CLIB122, DSM 3286, E122 *MATA*, and 22301-5 *MATB*. A. Left**, the green bar graph shows the number of predicted genes using an unfiltered analysis for each of the four strains**. Right,** a comparison of gene orthology across the strains, with vertical lines connecting shared orthologous genes shared between the indicated strains. Unconnected points indicate unique contributions. **B.** Bar Chart showing the number of genes in each of the orthologous groups shown in A, right panel, that are shared among the strains. Most orthologs are shared across all strains (left most bar in panel B). The second bar in panel B represents orthologs shared between E122 *MATA* and 22301-5 *MATB* but not found in DSM 3286 or CLIB122 sequences. C. Insert shows a scatter plot showing functional enrichment analysis in the molecular process category identifies significant GO terms such as "methylmalonate-semialdehyde dehydrogenase activity" (1), "TBP-class protein binding" (2), and "oxidoreductase activity" (3). Points are colored by the significance (log10 p-value) **D.** Insert shows a scatter plot of Biological Process): Enrichment in the biological process category highlighted "DNA-templated transcription initiation" as a key process. This multi-dimensional analysis underscores the conserved and divergent functional pathways and genetic architecture of these *Yarrowia* strains
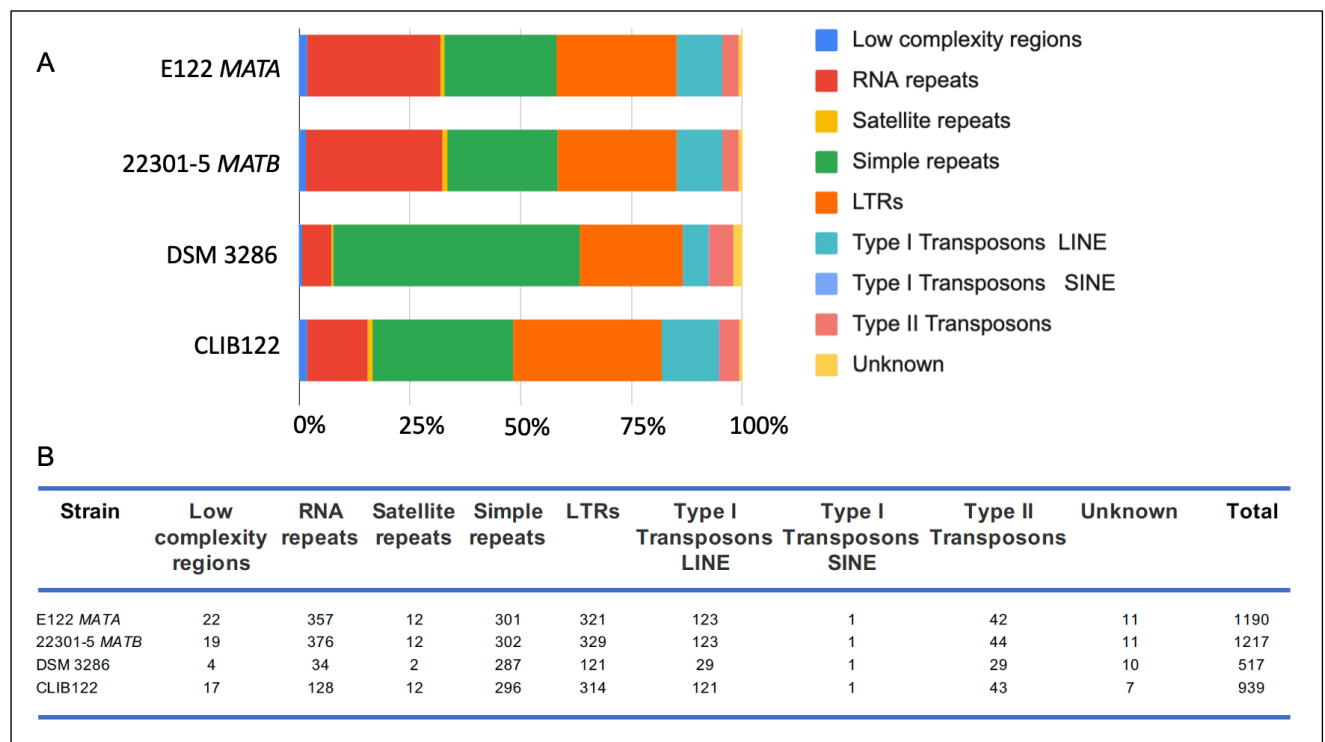
A.

| Strain | Low complexity regions | RNA repeats | Satellite repeats | Simple repeats | LTRs | Type I Transposons LINE | Type I Transposons SINE | Type II Transposons | Unknown | Total |
|---|---|---|---|---|---|---|---|---|---|---|
| E122 *MATA* | 22 | 357 | 12 | 301 | 321 | 123 | 1 | 42 | 11 | 1190 |
| 22301-5 *MATB* | 19 | 376 | 12 | 302 | 329 | 123 | 1 | 44 | 11 | 1217 |
| DSM 3286 | 4 | 34 | 2 | 287 | 121 | 29 | 1 | 29 | 10 | 517 |
| CLIB122 | 17 | 128 | 12 | 296 | 314 | 121 | 1 | 43 | 7 | 939 |



**Figure 6. Comparative repeat composition in the genomes of *Yarrowia lipolytica* strains *MATA*, *MATB*, DSM 3286, and CLIB122. A.** The stacked bar chart represents the proportional distribution of various repeat classes, including: Low complexity regions (blue); RNA repeats (red); Satellite repeats (yellow); Simple repeats (green); LTRs (orange); Type I Transposons (LINE) (teal); Type I Transposons (SINE) (light blue); Type II Transposons (pink); Unknown repeats (yellow). Each bar represents the total percentage of genomic content attributed to these repeat types in the corresponding strain, highlighting variations in transposable elements and repetitive sequences across strains. This analysis reveals genome-wide repeat diversity and relative abundance, including transposons and low-complexity regions. **B.** The counts of the various repeat elements described in A in the four strains.
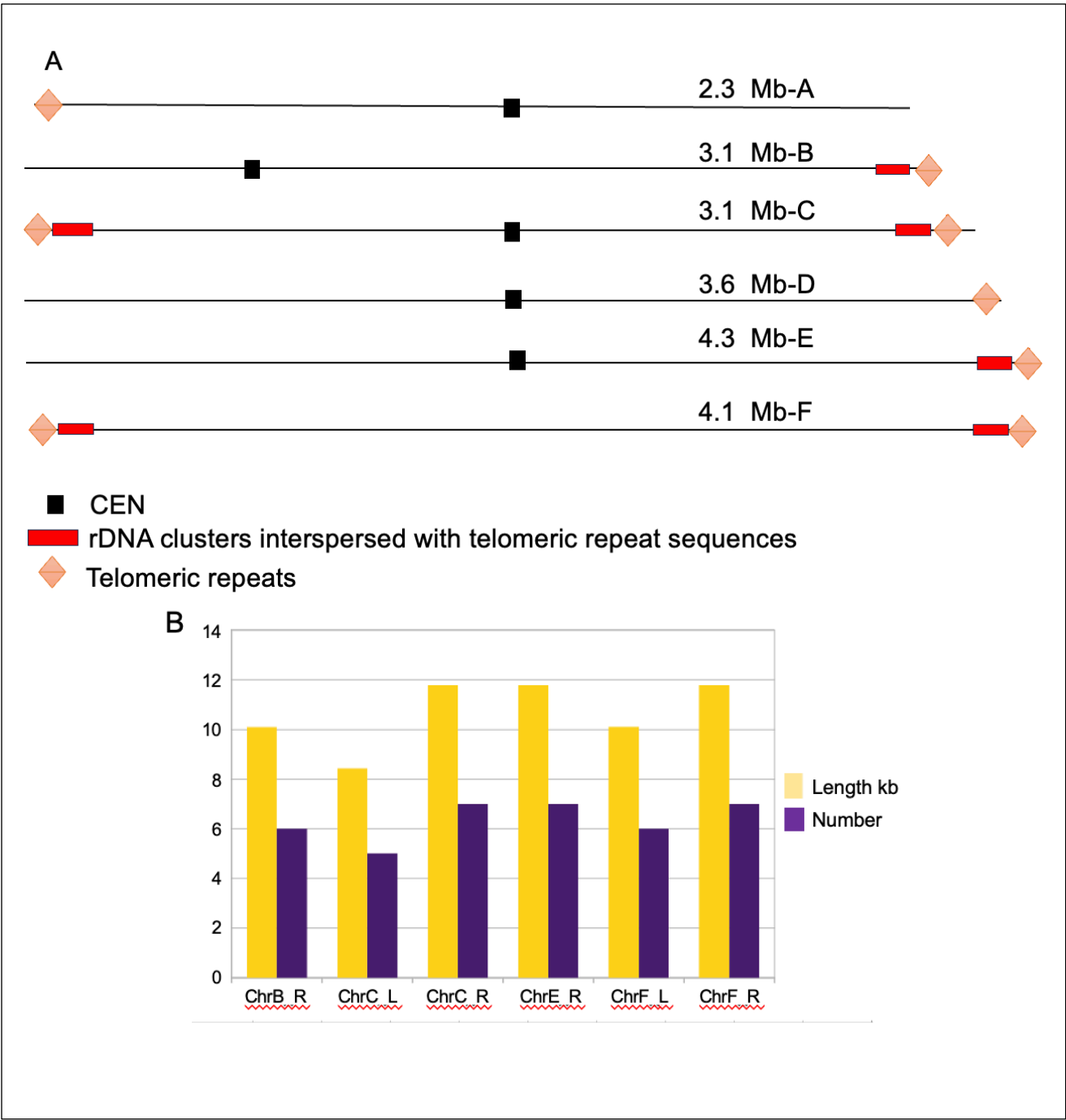
20

**Figure 7. Analysis of rDNA Repeats.**
**A.** Genome landscape of six chromosomes of *Y. lipolytica* mapped rDNA (red bar) and telomere (orange diamond) sequences. The known centromere sequences are shown with black bars. **B.** This bar chart compares the Length (in kilobases, yellow bars) and Count (purple bars) of annotated features for specific regions across chromosomes in the genome. The x-axis represents distinct chromosome regions, including ChrB-R (Right arm of Chromosome B), ChrC-L (Left arm of Chromosome C), ChrC-R (Right arm of Chromosome C), ChrE-R (Right arm of Chromosome E), ChrF-L (Left arm of Chromosome F), ChrF-R (Right arm of Chromosome F). The yellow bars indicate the cumulative length of the regions (in kilobases), while the purple bars indicate the total count of repeats identified within these regions.

**REFERENCES**

1. Gonçalves,F.A.G., Colen,G. and Takahashi,J.A. (2014) Yarrowia lipolytica and Its Multiple Applications in the Biotechnological Industry. *Sci. World J.*, **2014**, 476207.

2. Mamaev,D. and Zvyagilskaya,R. (2021) Yarrowia lipolytica : A multitalented yeast species of ecological significance. *Fems Yeast Res*, 10.1093/femsyr/foab008.

3. Bankar,A.V., Kumar,A.R. and Zinjarde,S.S. (2009) Environmental and industrial applications of Yarrowia lipolytica. *Appl. Microbiol. Biotechnol.*, **84**, 847.

4. Groenewald,M., Boekhout,T., Neuvéglise,C., Gaillardin,C., Dijck,P.W.M. van and Wyss,M. (2014) Yarrowia lipolytica: Safety assessment of an oleaginous yeast with a great industrial potential. *Crit. Rev. Microbiol.*, **40**, 187–206.

5. Barth,G. and Gaillardin,C. (1997) Physiology and genetics of the dimorphic fungus Yarrowia lipolytica. *Fems Microbiol Rev*, **19**, 219–237.

6. Magnan,C., Yu,J., Chang,I., Jahn,E., Kanomata,Y., Wu,J., Zeller,M., Oakes,M., Baldi,P. and Sandmeyer,S. (2016) Sequence Assembly of Yarrowia lipolytica Strain W29/CLIB89 Shows Transposable Element Diversity. *Plos One*, **11**, e0162363.

7. Fumasoni,M. and Murray,A.W. (2020) The evolutionary plasticity of chromosome metabolism allows adaptation to constitutive DNA replication stress. *Elife*, **9**, e51963.

8. Mekouar,M., Blanc-Lenfle,I., Ozanne,C., Silva,C.D., Cruaud,C., Wincker,P., Gaillardin,C. and Neuvéglise,C. (2010) Detection and analysis of alternative splicing in Yarrowia lipolytica reveal structural constraints facilitating nonsense-mediated decay of intron-retaining transcripts. *Genome Biol*, **11**, R65.

9. Hu,Y., Tareen,A., Sheu,Y.-J., Ireland,W.T., Speck,C., Li,H., Joshua-Tor,L., Kinney,J.B. and Stillman,B. (2020) Evolution of DNA replication origin specification and gene silencing mechanisms. *Nat Commun*, **11**, 5175.

10. Marahrens,Y. and Stillman,B. (1992) A yeast chromosomal origin of DNA replication defined by multiple functional elements. *Science (New York, NY)*, **255**, 817–823.

11. Bell,S.P. and Labib,K. (2016) Chromosome Duplication in Saccharomyces cerevisiae. *Genetics*, **203**, 1027–1067.

12. Lee,C.S.K., Cheung,M.F., Li,J., Zhao,Y., Lam,W.H., Ho,V., Rohs,R., Zhai,Y., Leung,D. and Tye,B.-K. (2021) Humanizing the yeast origin recognition complex. *Nat Commun*, **12**, 33.

13. Rusche,L.N. and Hickman,M.A. (2020) Sex in Fungi. 10.1128/9781555815837.ch11.

14. Rusche,L.N., Kirchmaier,A.L. and Rine,J. (2003) The establishment, inheritance, and function of silenced chromatin in Saccharomyces cerevisiae. *Annual review of biochemistry*, **72**, 481–516.

592    15. Hickman,M.A., Froyd,C.A. and Rusche,L.N. (2011) Reinventing heterochromatin in budding yeasts: Sir2 and
593        the origin recognition complex take center stage. *Eukaryotic cell*, **10**, 1183–1192.

594    16. Smith,J.S., Brachmann,C.B., Celic,I., Kenna,M.A., Muhammad,S., Starai,V.J., Avalos,J.L., Escalante-
595        Semerena,J.C., Grubmeyer,C., Wolberger,C., *et al.* (2000) A phylogenetically conserved NAD+-dependent
596        protein deacetylase activity in the Sir2 protein family. *Proc. Natl. Acad. Sci.*, **97**, 6658–6663.

597    17. Bell,S.P., Kobayashi,R. and Stillman,B. (1993) Yeast Origin Recognition Complex Functions in Transcription
598        Silencing and DNA Replication. *Science*, **262**, 1844–1849.

599    18. Hickman,M.A. and Rusche,L.N. (2010) Transcriptional silencing functions of the yeast protein Orc1/Sir3
600        subfunctionalized after gene duplication. *Proceedings of the National Academy of Sciences*, **107**, 19384–
601        19389.

602    19. Hou,Z., Bernstein,D.A., Fox,C.A. and Keck,J.L. (2005) Structural basis of the Sir1–origin recognition
603        complex interaction in transcriptional silencing. *Proc. Natl. Acad. Sci.*, **102**, 8489–8494.

604    20. Grunstein,M. and Gasser,S.M. (2013) Epigenetics in Saccharomyces cerevisiae. *Cold Spring Harbor
605        Perspectives in Biology*, **5**, a017491.

606    21. Maria,H. and Rusche,L.N. (2022) The DNA replication protein Orc1 from the yeast Torulaspora delbrueckii is
607        required for heterochromatin formation but not as a silencer-binding protein. *Genetics*, **222**, iyac110.

608    22. Hu,Y. and Stillman,B. (2023) Origins of DNA replication in eukaryotes. *Mol Cell*, **83**, 352–372.

609    23. Hyrien,O., Guilbaud,G. and Krude,T. (2025) The double life of mammalian DNA replication origins. *Genes
610        Dev.*, **39**, 304–324.

611    24. Devillers,H. and Neuvéglise,C. (2019) Genome Sequence of the Oleaginous Yeast Yarrowia lipolytica H222.
612        *Microbiol. Resour. Announc.*, **8**, 10.1128/mra.01547-18.

613    25. Bigey,F., Pasteur,E., Połomska,X., Thomas,S., Coq,A.-M.C.-L., Devillers,H. and Neuvéglise,C. (2023)
614        Insights into the Genomic and Phenotypic Landscape of the Oleaginous Yeast Yarrowia lipolytica. *J. Fungi*, **9**,
615        76.

616    26. Liu,L. and Alper,H.S. (2014) Draft Genome Sequence of the Oleaginous Yeast Yarrowia lipolytica PO1f, a
617        Commonly Used Metabolic Engineering Host. *Genome Announc.*, **2**, e00652-14.

618    27. Luttermann,T., Rückert,C., Wibberg,D., Busche,T., Schwarzhans,J.-P., Friehs,K. and Kalinowski,J. (2021)
619        Establishment of a near-contiguous genome sequence of the citric acid producing yeast Yarrowia lipolytica
620        DSM 3286 with resolution of rDNA clusters and telomeres. *Nar Genom Bioinform*, **3**, lqab085-.

621    28. Devillers,H., Brunel,F., Połomska,X., Sarilar,V., Lazar,Z., Robak,M. and Neuvéglise,C. (2016) Draft Genome
622        Sequence of Yarrowia lipolytica Strain A-101 Isolated from Polluted Soil in Poland. *Genome Announc.*, **4**,
623        e01094-16.

624    29. Madzak,C. (2021) Yarrowia lipolytica Strains and Their Biotechnological Applications: How Natural
625        Biodiversity and Metabolic Engineering Could Contribute to Cell Factories Improvement. *J Fungi*, **7**, 548.

626  30. Gaillardin,G.B. and C. and Wolf (2019) Yarrowia lipolytica.

627  31. Mansour,S., Bailly,J., Landaud,S., Monnet,C., Sarthou,A.S., Cocaign-Bousquet,M., Leroy,S., Irlinger,F. and
628      Bonnarme,P. (2009) Investigation of Associations of Yarrowia lipolytica, Staphylococcus xylosus, and
629      Lactococcus lactis in Culture as a First Step in Microbial Interaction Analysis. *Appl. Environ. Microbiol.*, **75**,
630      6422–6430.

631  32. Kolmogorov,M., Yuan,J., Lin,Y. and Pevzner,P.A. (2019) Assembly of long, error-prone reads using repeat
632      graphs. *Nat. Biotechnol.*, **37**, 540–546.

633  33. Martin and M. (2011) Cutadapt removes adapter sequences from high-throughput sequencing reads.
634      *EMBnet.journal [S.I.]*.

635  34. Li,H. and Durbin,R. (2009) Fast and accurate short read alignment with Burrows–Wheeler transform.
636      *Bioinformatics*, **25**, 1754–1760.

637  35. Danecek,P., Bonfield,J.K., Liddle,J., Marshall,J., Ohan,V., Pollard,M.O., Whitwham,A., Keane,T.,
638      McCarthy,S.A., Davies,R.M., *et al.* (2021) Twelve years of SAMtools and BCFtools. *GigaScience*, **10**,
639      giab008.

640  36. Walker,B.J., Abeel,T., Shea,T., Priest,M., Abouelliel,A., Sakthikumar,S., Cuomo,C.A., Zeng,Q., Wortman,J.,
641      Young,S.K., *et al.* (2014) Pilon: An Integrated Tool for Comprehensive Microbial Variant Detection and
642      Genome Assembly Improvement. *PLoS ONE*, **9**, e112963.

643  37. Alonge,M., Soyk,S., Ramakrishnan,S., Wang,X., Goodwin,S., Sedlazeck,F.J., Lippman,Z.B. and Schatz,M.C.
644      (2019) RaGOO: fast and accurate reference-guided scaffolding of draft genomes. *Genome Biol.*, **20**, 224.

645  38. Bolger,A.M., Lohse,M. and Usadel,B. (2014) Trimmomatic: a flexible trimmer for Illumina sequence data.
646      *Bioinformatics*, **30**, 2114–2120.

647  39. Pertea,M., Pertea,G.M., Antonescu,C.M., Chang,T.-C., Mendell,J.T. and Salzberg,S.L. (2015) StringTie
648      enables improved reconstruction of a transcriptome from RNA-seq reads. *Nat. Biotechnol.*, **33**, 290–295.

649  40. Kainth,A.S., Haddad,G.A., Hall,J.M. and Ruthenburg,A.J. (2023) Merging short and stranded long reads
650      improves transcript assembly. *PLOS Comput. Biol.*, **19**, e1011576.

651  41. Pertea,G. and Pertea,M. (2020) GFF Utilities: GffRead and GffCompare. *F1000Research*, **9**, ISCB Comm J-
652      304.

653  42. Campbell,M.S., Holt,C., Moore,B. and Yandell,M. (2014) Genome Annotation and Curation Using MAKER
654      and MAKER-P. *Curr. Protoc. Bioinform.*, **48**, 4.11.1-4.11.39.

655  43. Huang,Y., Niu,B., Gao,Y., Fu,L. and Li,W. (2010) CD-HIT Suite: a web server for clustering and comparing
656      biological sequences. *Bioinformatics*, **26**, 680–682.

657  44. Trincado,J.L., Entizne,J.C., Hysenaj,G., Singh,B., Skalic,M., Elliott,D.J. and Eyras,E. (2018) SUPPA2: fast,
658      accurate, and uncertainty-aware differential splicing analysis across multiple conditions. *Genome Biol.*, **19**, 40.

659   45. Solovyev,V., Kosarev,P., Seledsov,I. and Vorobyev,D. (2006) Automatic annotation of eukaryotic genes,
660       pseudogenes and promoters. *Genome Biol.*, **7**, S10.

661   46. Stanke,M., Diekhans,M., Baertsch,R. and Haussler,D. (2008) Using native and syntenically mapped cDNA
662       alignments to improve de novo gene finding. *Bioinformatics*, **24**, 637–644.

663   47. Haas,B.J., Delcher,A.L., Mount,S.M., Wortman,J.R., Jr,R.K.S., Hannick,L.I., Maiti,R., Ronning,C.M.,
664       Rusch,D.B., Town,C.D., *et al.* (2003) Improving the Arabidopsis genome annotation using maximal transcript
665       alignment assemblies. *Nucleic Acids Res.*, **31**, 5654–5666.

666   48. Jones,P., Binns,D., Chang,H.-Y., Fraser,M., Li,W., McAnulla,C., McWilliam,H., Maslen,J., Mitchell,A.,
667       Nuka,G., *et al.* (2014) InterProScan 5: genome-scale protein function classification. *Bioinformatics*, **30**, 1236–
668       1240.

669   49. Olson,A.J. and Ware,D. (2021) Ranked choice voting for representative transcripts with TRaCE.
670       *Bioinformatics*, **38**, 261–264.

671   50. Stabenau,A., McVicker,G., Melsopp,C., Proctor,G., Clamp,M. and Birney,E. (2004) The Ensembl Core
672       Software Libraries. *Genome Res.*, **14**, 929–933.

673   51. Vilella,A.J., Severin,J., Ureta-Vidal,A., Heng,L., Durbin,R. and Birney,E. (2009) EnsemblCompara
674       GeneTrees: Complete, duplication-aware phylogenetic trees in vertebrates. *Genome Res.*, **19**, 327–335.

675   52. Simão,F.A., Waterhouse,R.M., Ioannidis,P., Kriventseva,E.V. and Zdobnov,E.M. (2015) BUSCO: assessing
676       genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics*, **31**, 3210–3212.

677   53. Pérez-Wohlfeil,E., Diaz-del-Pino,S. and Trelles,O. (2019) Ultra-fast genome comparison for large-scale
678       genomic experiments. *Sci. Rep.*, **9**, 10274.

679   54. Rosas-Quijano,R., Gaillardin,C. and Ruiz-Herrera,J. (2008) Functional analysis of the MATB mating-type
680       idiomorph of the dimorphic fungus Yarrowia lipolytica. *Current Microbiol*, **57**, 115–120.

681   55. Kurischko,C., Schilhabel,M.B., Kunze,I. and Franzl,E. (1999) The MAT A locus of the dimorphic yeast
682       Yarrowia lipolytica consists of two divergently oriented genes. *Mol. Gen. Genet. MGG*, **262**, 180–188.

683   56. Dujon,B., Sherman,D., Fischer,G., Durrens,P., Casaregola,S., Lafontaine,I., Montigny,J. de, Marck,C.,
684       Neuvéglise,C., Talla,E., *et al.* (2004) Genome evolution in yeasts. *Nature*, **430**, 35–44.

685   57. McCarthy,C.G.P. and Fitzpatrick,D.A. (2019) Pangloss: A Tool for Pan-Genome Analysis of Microbial
686       Eukaryotes. *Genes*, **10**, 521.

687   58. Lee,S.C., Ni,M., Li,W., Shertz,C. and Heitman,J. (2010) The Evolution of Sex: a Perspective from the Fungal
688       Kingdom. *Microbiol. Mol. Biol. Rev.*, **74**, 298–340.

689   59. Haber,J.E. (2012) Mating-Type Genes and MAT Switching in Saccharomyces cerevisiae. *Genetics*, **191**, 33–
690       64.

691   60. Hanson,S.J. and Wolfe,K.H. (2017) An Evolutionary Perspective on Yeast Mating-Type Switching. *Genetics*,
692       **206**, 9–32.

693   61. Dujon,B.A. and Louis,E.J. (2017) Genome Diversity and Evolution in the Budding Yeasts
694       (Saccharomycotina). *Genetics*, **206**, 717–750.

695   62. Kueng,S., Oppikofer,M. and Gasser,S.M. (2013) SIR proteins and the assembly of silent chromatin in budding
696       yeast. *Annual review of genetics*, **47**, 275–306.

697   63. Blander,G. and Guarente,L. (2004) The SIR2 Family of Protein Deacetylases. *Biochemistry*, **73**, 417–435.

698
699