1  **TITLE**

2

3

4  The genome of the CTG(Ser1) yeast S*cheffersomyces stipitis* is plastic

5

6  Samuel Vega Estevez[1], Andrew Armitage[2], Helen J. Bates[3], Richard J. Harrison[3] and

7  Alessia Buscaino[1#]

8

9  [1] University of Kent, School of Biosciences, Kent Fungal Group, Canterbury Kent,

10  CT2 7NJ, UK.

11  [2] Natural Resources Institute, University of Greenwich, Chatham Maritime, Kent,

12  ME4 4TB, UK.

13  [3] NIAB Cambridge Crop Research, 93 Lawrence Weaver Road , Cambridge

14   CB3 0LE

15

16  [#] Address correspondence to Alessia Buscaino

17  Email: A.Buscaino@kent.ac.uk

18

19

20

21

22

23

*Vega Estevez et al*

24 **ABSTRACT**

25 Microorganisms need to adapt to environmental changes, and genome plasticity can

26 lead to rapid adaptation to hostile environments by increasing genetic diversity.

27 Here, we investigate genome plasticity in the CTG(Ser1) yeast *Scheffersomyces*

28 *stipitis,* an organism with an enormous potential for second-generation biofuel

29 production. We demonstrate that *S. stipitis* has an intrinsically plastic genome and

30 that different *S. stipitis* isolates have genomes with distinct chromosome

31 organisation. Real-time evolution experiments show that *S. stipitis* genome plasticity

32 is common and rapid as extensive genomic changes with fitness benefits are

33 detected following *in vitro* evolution experiments. Hybrid MinION Nanopore and

34 Illumina genome sequencing identifies retrotransposons as major drivers of genome

35 diversity. Indeed, the number and position of retrotransposons is different in different

36 *S. stipitis* isolates, and retrotransposon-rich regions of the genome are sites of

37 chromosome rearrangements. Our findings provide important insights into the

38 adaptation strategies of the CTG (Ser1) yeast clade and have critical implications in

39 the development of second-generation biofuels. These data highlight that genome

40 plasticity is an essential factor to be considered for the development of sustainable

41 *S. stipitis* platforms for second-generation biofuels production.

42

43

44

45

46

47

*Vega Estevez et al*

**INTRODUCTION**

48

49    Eukaryotic genomes are often described as stable structures with well-preserved

50    chromosome organisation, and genome instability is viewed as harmful. However, an

51    increasing body of evidence demonstrates that eukaryotic microorganisms have a

52    plastic genome and genome instability is instrumental for rapid and reversible

53    adaptation to hostile environments (1–4). This is because genomic instability can

54    increase genetic diversity, allowing the selection of genotype(s) better adapted to a

55    new environment (5, 6). Repetitive DNA elements are major contributors to genome

56    plasticity as repeats can undergo inter and intra-locus recombination, resulting in

57    gene conversion, gross chromosomal rearrangements and segmental aneuploidies

58    (7). Transposable Elements (TE), a specific class of repetitive elements, alter

59    genome organisation by recombination-dependent mechanisms and by jumping to

60    new sites in the genome (8). TEs belong to two major classes: DNA transposons

61    (Class II) and retrotransposons (Class I). DNA transposons utilise a "cut and paste"

62    mechanism in which the parental element excises from its original location before

63    integrating elsewhere (9). In contrast, retrotransposons replicate through reverse

64    transcription of their RNA and integrate the resulting cDNA into another locus.

65    Retrotransposons can be further classified into Long Terminal Repeats (LTR)

66    retrotransposon and non-LTR retrotransposons (10). LTR retrotransposons are

67    characterised by two LTR sequences flanking an internal coding region containing

68    the genes encoding for the structural protein GAG and enzyme POL required for

69    reverse transcription and integration (11). While POL enzymes are conserved across

70    organisms, GAG proteins are poorly conserved (12). LINE elements are one of the

71    most abundant non-LTR retrotransposons and they are typically composed of a 5'

72    non-coding region, two ORFs (ORF1 and ORF2) and a 3' non-coding region that is

3

73  marked by a poly(A) tail (13). ORF1 proteins have a diverse amino acid sequence,

74  but they often contain a DNA-binding motif (14). ORF2 encodes endonuclease and

75  reverse transcriptase activity that are critical for transposition (15).

76

77  The CTG (Ser1) clade of fungi ,in which the CTG codon is translated as serine rather

78  than leucine, is an important group of ascomycetous yeasts featuring yeasts that

79  hold great promises in biotechnology, such as *Scheffersomyces stipitis*, and

80  dangerous human fungal pathogens, such as *Candida albicans* (16).

81  The CTG(Ser1) clade comprises several species with different lifestyles and

82  genomic organisations, including haploid and diploid species that colonise diverse

83  environments by reproducing sexually or para-sexually (16–19). One common

84  feature of CTG(Ser1) species is their ability to adapt remarkably well to extreme

85  environments (20). For example, CTG(Ser1) yeasts can grow on various carbon

86  sources and are highly tolerant to environmental changes such as changes in

87  osmolarity (16, 19, 20). It is well established that genome plasticity is a critical

88  adaptive mechanism in the human fungal pathogens *Candida albicans*, the most

89  studied CTG (Ser1)-clade member (4). In *C. albicans*, stress increases genome

90  instability by affecting the rate and type of genomic rearrangements (21). Different

91  classes of DNA repeats drive this genetic variation, including TEs, long repeats and

92  Major Repeat Sequences (MRS) (22–24). It is still unknown whether genome

93  plasticity is a general feature of the CTG(Ser1) clade and whether DNA repeats are

94  drivers for genome diversity across this yeast group.

95  This study investigates genome plasticity in *S. stipitis*, a CTG (Ser1)-clade yeast with

96  great potential for the eco-friendly and ethical production of second-generation

97  biofuels (25–27). Second-generation biofuels are generated by fermentation of

98    lignocellulose biomass, produced in large amounts (>1.3 billion tons produced

99    annually) as waste following agricultural and forestry processing operation (27).

100   Lignocellulose is a heteropolymer composed of fermentable hexose sugars, such a

101   glucose, and pentose sugars, such as xylose (28). The yeast *Saccharomyces*

102   *cerevisiae*, usually the organism of choice for industrial production of ethanol, is not

103   suitable for the production of second-generation ethanol because it cannot ferment

104   pentose sugars as it lacks specific transporters and enzymatic network important for

105   their metabolism (28). *S. stipitis* holds excellent potential for biofuel derived from

106   green waste because it is one of the few yeast species that can ferment both hexose

107   and pentose sugars (25–27). *S. stipitis* is a non-pathogenic haploid yeast that is

108   found in the gut of wood-ingesting beetles, in hardwood forests or areas high in

109   agricultural waste (29). Contrary to *C. albicans*, *S. stipitis* has a canonical sexual

110   cycle whereby mating of haploid cells generate diploid cells that undergo meiosis

111   and produce haploid spores (30). Although several *S. stipitis* natural isolates are

112   used for the optimisation of second-generation biofuels production, the genome of

113   only one strain (Y-11545) has been sequenced and assembled to the chromosomal

114   level (31). The Y-11545 genome has a size of 15.4 Million base pair (Mbp) organised

115   in 8 chromosomes and containing ~6000 protein-coding genes (31–33). *S. stipitis*

116   chromosomes are marked by regional centromeres composed of full-lenghts LTR

117   retrotrasposons (Tps5a, Tps5b and Tps5c) and non-coding, non-autonomous LARD

118   (large retrotransposon derivative) elements (31, 33).

119   To investigate the plasticity of the *S. stipitis* genome, we have taken several

120   complementary approaches. Firstly, we systematically identified *S. stipitis* DNA

121   repeats and investigated the genotypic diversity of 27 different *S. stipitis* natural

122   isolates collected from different environments. Secondly, we combined MinION

5

123    Nanopore with Illumina genome sequencing to generate a high-quality chromosome-

124    level sequence assembly of a second *S. stipitis* natural isolate (Y-7124) and

125    compared its genome structure to the reference Y-11545 genome. Lastly, we

126    performed *in vitro* evolution experiments and analysed *S. stipitis* genome

127    organisation changes following laboratory passaging under stress or unstressed

128    growth conditions. Thanks to this combined approach, we discovered that the *S.*

129    *stipitis* genome is plastic. Genome plasticity is not regulated by stress, however large

130    chromosome rearrangements are linked to adaptation to hostile environments. We

131    demonstrate that different *S. stipitis* natural isolates have distinct chromosomal

132    organisations and that transposable elements drive this extensive intra-species

133    genetic variation. Our findings have important implications for second-generation

134    biofuel production as genome plasticity is a paramount factor to be considered for

135    the successful development of superior biofuel-producer *S. stipitis* strains.

136

137

138                            **MATERIAL AND METHODS**

139    **Yeast strains and Growth Conditions**

140    Strains were obtained from the Agricultural Research Service (ARS) Collection, run

141    by the Northern Regional Research Laboratory (NRRL) (Peoria, Illinois, USA), or the

142    National Collection of Yeast Cultures (NCYC) (Norwich, United Kingdom) (**Table S1)**

143    and confirmed by sequencing (primers AB798 and AB799 of the 26S rDNA (D1/D2

144    domain) (34) (**Table S2**). Routine culturing was performed at 30 ºC with 200 rpm

145    agitation on Yeast Extract-Peptone-D-Glucose (YPD) media. Phenotypic and *in vitro*

146    evolution analyses were conducted on Synthetic Complete (SC) media containing

147    glucose (SC-G), xylose (SC-X), or a mixture of 60% glucose and 40% xylose (SC-

*Vega Estevez et al*

148    G+X). SC-G was used as a reference media as glucose is the preferred carbon

149    source for both the model system *S. cerevisiae* and *S. stipitis*, SC-X was used

150    because of *S. stipitis* unique ability to utilise xylose as a carbon source and SC G+X

151    was used because this sugar combination resemble the ratio found in lignocellulose

152    (28). Uridine (0.08 g/L in YPD and SC) and adenine hemisulfate (0.05 g/L in YPD)

153    were added as growth supplements. Solid media were prepared by adding 2%

154    agar.


**Contour-clamped homogeneous electric field (CHEF) electrophoresis**

156    Intact yeast chromosomal DNA was prepared as previously described (35).

157    Briefly, cells were grown overnight and spheroplast were prepared in an agarose

158    plug by treating cells (~ $OD_{600}$=7) with 0.6 mg/ml Zymolyase 100T (Amsbio #120493-

159    1) in 1% Low Melt agarose (Biorad® # 1613112). Chromosomes were separated in a

160    1% Megabase agarose gel (Bio-Rad) in 0.5X TBE using a CHEF DRII apparatus.

161    Run conditions as follows: 60-120s switch at 6 V/cm for 12 hours followed by a 120-

162    300s switch at 4.5 V/cm for 12 hours, 14 °C. Chromosomes were visualised by

163    staining the gel 0.5x TBE with ethidium bromide (0.5 µg/ml) for 30 minute, followed

164    by destaining in water for 30 minutes. Images were capture using a Syngene GBox

165    Chemi XX6 gel imaging system.


**Southern Blotting**

167    DNA from CHEF gels were transferred overnight to a Zeta-Probe GT

168    Membrane (Biorad®, #162-0196) in 20x SSC and crosslinked using UV (150 mJ).

169    Probing and detection of the DNA were conducted as previously described (36).

170    Briefly, probes were generated by PCR incorporation of DIG-11-dUTP into target

7

171    sequences following manufacturer's instructions (Roche). Chromosome 5-

172    chromosome 7 translocation was detected using primers AB1028 and AB1029

173    amplifying a 180 bp region of chromosome 5 (Chr5 nt: 448,855-449,034) in Y-11545

174    and in chromosome 7 of Y-7124 (Chr7 nt: 494,698-494,877) (**Table S2**). The

175    membrane was hybridised overnight at 42 ºC with DIG Easy Hyb (Roche®,

176    11603558001). The DNA was detected with anti-digoxigenin-Alkaline Phosphatase

177    antibody (Roche®, #11093274910) and CDP Star ready to use (Roche®,

178    #12041677001) according to manufacturer instructions.

179    **Phenotypic characterisation**

180        Growth analyses were performed using a plate reader (SpectrostarNano,

181    BMG labtech) in 96 well plate format at 30 °C for 48 hours in SC-G, SC-X or SC-

182    G+X. The growth rate ($\mu$, hours$^{-1}$) was calculated using: $\mu = (\ln(X_2)-\ln(X_1))/(t_2-t_1)$,

183    where: (i) $X_1$ is the biomass concentration ($OD_{600}$) at time point one ($t_1$, lag time) (ii)

184    $X_2$ is the biomass concentration ($OD_{600}$) at time point two ($t_2$, end of exponential

185    growth phase). The maximum OD (OD units) was determined with the MAX() from

186    Excel (Microsoft®). The lag time (minutes) was determined visually as the time in

187    which the exponential growth starts. Experiments were performed in 3 technical and

188    3 biological replicates. Heatmaps showing the average of 3 biological replicates were

189    generated by R using the library *pheatmap*. ANOVA test was performed to study

190    differences on growth rate, maximum OD and lag time between the strains. The

191    equality of variances presumption was tested using Levene's test, whereas the

192    normality of the data was tested by Shapiro-Wilk. When both assumptions were

193    satisfied, a Tukey's honest significant test was used to determine significant

194    differences between the natural isolates and the reference Y-11545 strain. When

*Vega Estevez et al*

195   assumption of equal variance were violated, one-way test was used to indicate

196   significance. In the case of equal variances, but a non-normal distribution of data, the

197   Kruskal-Wallis rank sum test was used to indicate statistical differences and

198   significance was determined by pairwise testing. A p-value lower than 0.05 was

199   considered significant for all these statistical tests.


200   **Adaptive Laboratory Evolution**

201        A single colony of the *S. stipitis* strain NRRL Y-7124 was grown overnight in 5

202   ml of YPD at 30 ºC, plated in YPD at a cell density of 100 and grown 48 hours at 30

203   ºC. 36 single colonies were streaked in two SC-G+X plates and grown at 30 ºC and

204   37 ºC, respectively and streaked daily for a total of 56 passages (8 weeks). The

205   karyotype variability of the colonies was assessed by CHEF electrophoresis.

206   Phenotypic differences were assessed by spotting assays. Strains with

207   rearrangements were grown overnight is SC-G+X and were diluted to an $OD_{600}=1$.

208   From this, five 1/10 serial dilutions were prepared and the cells were plated in SC-

209   G+X using a replica plater (Sigma Aldrich, R2383-1EA) and grown for 48 hours at

210   both 30 ºC and 37 ºC. Strains with no karyotypic modifications after evolution were

211   also used as control.


212   **Identification of DNA repeats**

213        Long sequences (>100 nucleotides) present more than once in the Y-11545

214   and Y-7124 genomes were identified by aligning each genome to itself using

215   BLASTN.  Repetitive elements (E < 1e-04) were manually verified using

216   IGV/SNAPGene, and clustered repeats were combined. This repeats dataset was

217   manually examined to further classify it as (a) related to transposable elements (b)

9

218    telomeric repeats, (c) centromeres (d) belonging to protein coding gene families and

219    (e) MRS repeats. Transposons were classified using established guidelines (10).

220    Briefly, LTR-transposons were identified by detecting two Long-terminal Repeat

221    sequences (size 260-430 nt) flanking an internal coding region. These potential LTR-

222    transposons were further annotated for the presence of the following marks: LTR

223    flanked by a TG and CA di-nucleotides, presence of a Primer Binding Site (PBS) with

224    homology to *S. stipitis* tRNAs (GtRNAdb (http://gtrnadb.ucsc.edu/index.html),

225    presence of a coding region with homology to *pol* gene and containing an Integrase

226    (INT), Reverse Transcriptase (RT) and RNAse H (RH) domain. Non-LTR LINE

227    transposons were identified by detection of coding regions homologous to LINE

228    retrotransposons ORF1 (containing a Zn-finger), ORF2 (containing an Endonuclease

229    and a Reverse Transcriptase domain) and terminal Poly-A sequence.

230    Retrotransposons were classified into different families based on sequence similarity

231    with a 90% cut-off. Terminal telomeric tandem repeats were identified using Tandem

232    Repeats Finder (37) with default parameters. Regional centromeres were identified

233    based on them being the only regions of the genome with a large retrotransposon

234    Tps5 cluster (~ 20-40 kb) as previously described (33). Gene families were identified

235    by extracting coding-regions from our repeats datasets and performing Clustal

236    Omega sequence alignment and PFAMs domain identification using

237    SMART(http://smart.embl.de) (38). The identified gene families were compared to

238    published information (39). The presence of MRS repeats was explored using

239    BLASTN and by searching for clusters of non-coding tandem repeats, a hallmark of

240    *C. albicans* MRS, with no-homology to retrotransposons and not-located at

241    chromosome ends. Sequence alignments were visualised with Jalview v2.11.1.0

242    (40). Phylogenetic trees were generated with phyloT : a phylogenetic tree generator

243    (biobyte.de) using default parameters and visualised with Itol (https://itol.embl.de/).

**Genome sequencing**

244

245        The genome of *S. stipitis* isolate Y-7124 was sequenced by Illumina short-

246    read and MinION long-read technologies.  To this end, DNA was extracted from an

247    overnight culture using the QIAGEN genomic tip 100/G kit (Qiagen®, #10243)

248    according to manufacturing protocol. For long-read sequencing, MinION (Oxford

249    Nanopore, Oxford UK) was performed on a DNA library prepared from size selected

250    gDNA. DNA fragments greater than 30 Kb were selected using a Blue Pippin (Sage

251    Science) and concentrated using Ampure beads. From this, a DNA library was

252    prepared using a Ligation Sequencing Kit 1D (SQK-LSK108) and run on the Oxford

253    Nanopore MinION flowcell FLOMIN 106D.  The same gDNA extract was also used

254    for the preparation of Illumina libraries. In this case, the DNA was sheared using the

255    Covaris M220 with microTUBE-50 (Covaris 520166) and size selected using the

256    Blue Pippin (Sage Science). The library was constructed using a PCR-free kit with

257    NEBNext End Repair (E6050S), NEBNext dA-tailing (E6053S) and Blunt T/A ligase

258    (M0367S) New England Biolabs modules. Sequencing was performed on a MiSeq

259    Benchtop Analyzer (Illumina) using a 2x300bp PE (MS-102-3003) flow cell.

**Genome assembly**

260

261    Base-calling and demultiplexing were conducted with Albacore v2.3.3 (available at

262    https://community.nanoporetech.com). Adapters and low-quality data were trimmed

263    using the eautils package fastq-mcf 1.04.636

264    (https://expressionanalysis.github.io/ea-utils/). On nanopore sequence data, adapter

265    trimming was performed with Porechop v.0.1.0 (https://github.com/rrwick/Porechop).

266    Genome assembly was completed using long reads, with read correction performed

267    with Canu v1.8 (41) followed by  assembly in SMARTdenovo github commit id

268    61cf13d (42)). The draft assembly was corrected using the corrected nanopore reads

269    through five rounds of Racon github commit 24e30a9 (43), and then by raw fast5

270    files using 10 rounds of Nanopolish v0.9.0 (44). Illumina sequencing reads were then

271    used to polish the resulting assembly through 10 rounds of Pilon v1.17 (45).

272    Following genome assembly, BUSCO v3 was run to assess evolutionary conserved

273    gene content (46), using the Saccharomycetales_odb9 gene database. The

274    Saccharomycetales database contains 1711 genes, which are therefore expected to

275    be present in *S. stipitis*. Of these, 1683 (98.36%) were identified in the Y-7124

276    assembly demonstrating a good level of completeness (>95%) (**Table S8**). Assembly

277    size and contiguity statistics were assessed using QUAST v4.5 (47). This initial

278    assembly of the nuclear genome contained 10 contigs. A chromosome level

279    assembly was produced by identification of overlapping regions between the contigs:

280    a 244 Kbp overlapping region between contig 7 and 2 led to the final assembly of

281    Chromosome 1, a 83 Kpb overlapping region between contig 9 and 10 led to the final

282    assembly of Chromosome 8.

283    **Genome annotation**

284       Genome annotation was performed using FUNGAP v1.0.1 (48) with fastq

285    reads from NCBI SRA accession SRR8420582 used as RNA-Seq training data and

286    protein sequences taken from NCBI assembly accession GCA_000209165.1 for *S.*

287    *stipitis* NRRL Y-11545 (CBS6054) used for example proteins. Protein fasta files were

288    extracted from predicted gene models using the yeast mitochondrial code (code 3)

289 and the alternative yeast nuclear code (code 12). Functional annotation of gene

290 models was performed through BLASTp searches vs all proteins from the NCBI

291 reference fungal genomes (downloaded 11[th] April 2020), retrieving the top-scoring

292 blast hit with an E-value $< 1 \times 10^{-30}$. These annotations were supplemented with

293 domain annotations from Interproscan v5.42-78.0 (49). The annotated genome was

294 submitted to NCBI, with submission files prepared using GAG v2.0.1

295 (http://genomeannotation.github.io/GAG.), Annie github commit 4bb3980

296 (http://genomeannotation.github.io/annie) and table2asn_GFF v1.23.377 (available

297 from https://ftp.ncbi.nih.gov/toolbox/ncbi_tools/converters/by_program/tbl2asn/).

## Comparative genomics

299 Whole-genome alignment between Y-7124 and Y-11545 was performed using

300 the nucmer tool from the MUMmer package v4.0 (50) with results visualised using

301 Circos v0.6 (51). Orthology analysis was performed between predicted proteins from

302 these isolates using OrthoFinder v2.3.11 (52), with results visualised using the

303 package VennDiagram in R (53).

304 Sequence variants were identified in Y-7124 through comparison to the Y-

305 11545 assembly. Short read sequence data for Y-7124 were aligned to the reference

306 genome using BWA v 0.7.15-r1140 (54), before filtering using using picardtools

307 v2.5.0 to remove optical duplicates (http://broadinstitute.github.io/picard/). SNP and

308 insertion/deletion (InDel) calling was performed using GATK4 (55). Low confidence

309 variants were then filtered using VCFtools v0.1.15 (56) using minimum mapping

310 quality of 40, phred quality of 30, read depth of 10 and genotype quality of 30. Effect

311 of variants on NRRL Y-11545 gene models was determined using SnpEff v4.2 (57).

312 **RESULTS**

313 **Classification of *S. stipitis* DNA repeats**

314 DNA repeats are drivers of genome variation. Understanding the repertoire of

315 repetitive elements associated with a genome is critical to gain insights into the

316 genome diversity of a specific organism. Comparative genomic analyses have

317 identified different repetitive elements in some CTG(Ser1) clade members, yet a

318 comprehensive survey of *S. stipitis* repetitive elements is lacking (18, 58). Therefore,

319 we sought to classify the major classes of repetitive elements associated with the Y-

320 11545 sequenced genome by aligning the genomic sequence of each strain to itself

321 and identifying long sequences (>100 nucleotides) present more than once in the

322 genome. The genomic position of these repeats was manually verified, and clustered

323 repeats were combined and categorised depending on their genomic position,

324 structure and sequence similarity. Our analyses identified known *S. stipitis* repeat-

325 rich loci such as centromeric transposon-clusters, the NUPAV sequence, an

326 integrated L-A ds-RNA virus, and several gene families (32, 33, 59). As observed in

327 other members of the CTG (Ser1)-clade, we did not detect any MRS repeats, a class

328 of repetitive elements found only in *C. albicans* and the closely related *C. dubliensis*

329 and *C. tropicalis* species (58, 60).  Here we focus on intra-chromosomal or inter-

330 chromosomal repeats that have not been described to date: non-centromeric TEs,

331 subtelomeric regions and telomeric repeats (**Fig 1**).

332 We identified six novel retrotransposon families scattered along chromosome arms:

333 3 LTR retrotransposons (*Ava*, *Bea* and *Caia*) and 3 LINE retrotransposons (*Ace*, *Bri*

334 and *Can*) (**Fig 1A**, **Table S3**). *Ava*, *Bea* and *Caia* have a similar structure where two

335 identical LTR sequences flank an internal domain. The internal domain contains two

336 ORFs: one encoding for a putative POL and one encoding for an *S. stipitis*-specific

337 protein that we named LTR-Associated Protein (Lap1 in *Ava*, Lap2 in *Bea* and Lap3

338 in *Caia*). Homology search failed to identify any GAG gene associated with the *Ava*,

339 *Bea* and *Caia* retrotransposons. As Gag proteins are poorly conserved among

340 different organisms, we hypothesise that the Lap proteins are Gag proteins.

341 *Ace*, *Bri* and *Can* are LINE elements composed of the Non-Coding regions NC-1 and

342 NC-2  surrounding an internal coding region encoding for a Pol enzyme and an *S.*

343 *stipitis*-specific LINE Associated protein (Linea1 in *Ace*, Linea2 in *Bri* and Linea3 in

344 *Can*). Linea1 and Linea2, but not Linea 3, have a Zinc-Finger DNA binding motif

345 (**Table S3**). Comparison across the CTG (Ser1)-clade revealed that *S. stipits* TE

346 repertoire is typical of this clade. Indeed, retrotransposons are common in this yeast

347 group: the genome of all species analysed contains LTR elements, whereas LINE

348 elements are present in 6/8 species (**Fig 1B, Table S4**). Similarly to other CTG-

349 (Ser1) clade yeast, we did not detect any DNA transposons integrated into the *S.*

350 *stipitis* genome (**Fig 1B, Table S4**).

351 Our repeat analysis demonstrates that the terminal sequences of *S. stipitis*

352 chromosomes are repeat-rich and composed of two elements with different degree

353 of repetitiveness: telomere proximal-repeats and subtelomeric regions.  The

354 telomeric repeats are non-canonical and composed of 24-nucleotide units repeated

355 in tandem. Each unit contains a TG motif reminiscent of typical telomeric repeats

356 (**Fig 1C**). *S. stipitis* subtelomeric regions (the ~30KB region adjacent to telomeric

357 repeats) are enriched in retrotransposon-derived elements. Indeed, DNA sequences

358 with homology to *Bea* LTR-retrotransposons and *Ace* LINE-elements are found in

359 5/16 subtelomeric regions (**Fig 1D**, **Table S3**). No full length-retrotransposons are

360 detected at these genomic locations. Subtelomeric regions contain several gene

361 families members, including gene encoding for ATP-dependent DNA helicases

15

362  (found in 7/16 subtelomeres), fungal-specific transcription factors (8/16

363  subtelomeres), MFS transporters (8/16 subtelomeres) and Agglutinine-like proteins

364  (11/16 subtelomeres) (**Fig 1D, Table S5**) (39). In summary, our analysis

365  demonstrates that the *S. stipitis* genome contains several classes of repetitive

366  elements that could be major

367  contributors of genome plasticity.

368

369  ***S. stipitis* natural isolates have distinct genomic organisations**

370  Having identified *S. stipitis* DNA repeats, our next step was to examine *S. stipitis*

371  phenotypic and genotypic diversity across a geographically diverse set of strains

372  (n=27) that were collected in different habitats (**Table S1** source NRRL and NCYC

373  collection), and that includes the sequenced Y-11545 strain (31). rDNA fingerprinting

374  confirm that all isolates belong to the *S. stipitis* species (D1/D2 domain of the S26S

375  rDNA similarity >99 %) (**Table S6**). Phenotypic analyses established that the natural

376  isolates vary in their ability to utilise and grow on different carbon sources. Indeed,

377  when compared to the reference Y-11545 strain, different natural isolates cultured in

378  Synthetic Complete media containing the hexose sugar Glucose (SC-G), the

379  pentose sugar Xylose (SC-X) or a mixture of both sugars as found in lignocellulose

380  (SC-G+X) display distinct growth rate, maximum culture density and lag phase (**Fig**

381  **2A** and **Table S7**). To determine whether the natural isolates have distinct genomic

382  organisations, we analysed their karyotype by chromosomes Contour-clamped

383  Homogenous Electric Field (CHEF) gel electrophoresis, a technique allowing

384  chromosome separation according to size. The CHEF electrophoresis analysis

385  reveals clear differences in chromosome patterns demonstrating that *S. stipitis*

386  natural isolates have a genome organised in different-sized chromosomes (**Figure**

16

*Vega Estevez et al*

387 **2B**). We concluded that intra-species phenotypic and genotypic variation is a

388 common feature of *S. stipitis*.

389

390 **Hybrid genomic sequencing identifies transposable elements as drivers of *S.***

391 ***stipitis* genome plasticity**

392 To date, only one *S. stipitis* isolate (Y-11545) has been sequenced and assembled

393 at chromosome level (31). To gain insights into *S. stipitis* genetic diversity, we

394 generated a chromosome-level sequence assembly of a second *S. stipitis* natural

395 isolate (Y-7124) by combining MinION Nanopore with Illumina genome sequencing.

396 This isolate was chosen because *(i)* karyotypic analysis reveals that its genomic

397 organisation is distinct from the genomic organisation of the reference strain Y-

398 11545, and *(ii)* Y-7124 is widely used both for industrial applications and for basic

399 research (61).

400 The Y-7124 genome was sequenced to 186.88x coverage resulting in a 15.69 Mb

401 assembly arranged in 10 contigs (**Table S8**). High accuracy reads from Illumina-

402 sequencing enabled the correction of errors that are associated with the MinION

403 technology. A final chromosome-level assembly was produced by manually

404 identifying overlapping regions between contigs.  Comparing the Y-7124 and Y-

405 11545 nucleotide sequences reveals that the two natural isolates overall share a

406 similar coding DNA sequence. The total number of SNPs between the two natural

407 isolates is 50,495 SNPs, equating to one variant every 306 bases. The majority of

408 these SNPs are synonymous changes (16,294 =74.25%), while ~25% (5,622) of

409 SNPs are missense and only (0.13% (28) are nonsense (**Table S9**). Despite this

410 high DNA sequence similarity, the Y-7124 genome is organised in eight

411 chromosomes with different sizes and organisations from that of Y-11545 (**Fig 3A**).

17

412 Comparison of the Y-7124 and Y-11545 genomes establishes that retrotransposons

413 are significant drivers of *S. stipitis* genome diversity as one of the most prominent

414 differences between the two genomes is the abundance and localisation of these

415 retrotransposons (**Fig 3B**). Indeed, the number of LTR and LINE non-centromeric

416 retrotransposons and transposons-derived repeats is greater in the Y-11545

417 reference genome compared to the Y-7124 genome: retrotransposons, solo LTR and

418 truncated LINE elements account for approximately 2% of the reference Y-11545

419 genome and only for ~1% of the Y-7124 genome (**Fig 3C**). We classified

420 retrotransposons loci present in both isolates (ancestral loci), those present in the

421 reference Y-11545 genome but absent in Y-7124 (deletion loci) and those not

422 present in the reference genome but present in a given strain (insertion loci). Out of

423 69 transposons loci, only ten ancestral loci (~15%) were detected in the two isolates.

424 These sites are likely to be inactive transposons or transposons that rarely

425 transpose. In addition, we detected 42 deletion loci (60 %) and 17 (24%) insertion

426 loci (**Fig 3D**). The presence of deletion and insertion loci suggests that *S. stipitis* LTR

427 transposons and LINE elements are active and competent of transposition. Although

428 active transposons can insert into genes to cause functional consequences (62), we

429 did not detect any TE-driven alteration in coding regions.

430

431 **Transposable Elements are sites of chromosome rearrangements**

432 Comparison of the Y-11545 and Y-7124 genome reveals that transposon-rich

433 regions are sites of complex chromosome rearrangements. Indeed, a transposon-

434 rich region is the breakpoint of a reciprocal translocation between chromosome 5

435 and chromosome 7. This translocation causes the size change of chromosome 5 $^{Y-}$

436 $^{7124}$ and chromosome 7 $^{Y-7124}$ detected by CHEF karyotyping (**Fig 4A**). Southern

437  analyses with a probe specific for chromosome 5 [Y-11545] confirms this finding (**Fig**

438  **S1**). The evolutionary history of Y-11545 and Y-7124 is unknown, and therefore it is

439  difficult to predict the molecular events underlying these genomic changes. However,

440  sequence analysis of the rearrangement breakpoint reveals that this structural

441  variation occurs in a genomic region that *(i)* contains homologous sequences

442  between chromosome 5 and 7 and *(ii)* is transposon-rich and contains two inverted

443  repeats on chromosome 7 (**Fig 4B**). A second significant difference between the

444  genome organisation of Y-11545 and Y-7124 is found at subtelomeric regions: these

445  regions differ in the number and organisation of subtelomeric gene families and in

446  the number of transposon-associated repeats (**Fig 4C**). Lastly, we detected a distinct

447  centromeres organisation where the numbers of *Tps5* retrotransposons, LTRs and

448  LARD regions differ between the two isolates (**Fig 4D**). The presence of transposons

449  and transposon-derived repeats associated with all these genomic locations strongly

450  suggest that retrotransposons have mediated the chromosomal rearrangement by

451  recombination-mediated mechanisms. Therefore, changes in transposons

452  organisation are responsible for the bulk of genomic changes identified in two

453  different natural isolates.

454

455  ***S. stipitis* real-time evolution leads to extensive genomic changes**

456  Our results demonstrate that intraspecies genetic diversity is common in *S. stipitis.*

457  However, as the evolutionary history of the analysed natural isolates are unknown, it

458  is difficult to predict whether the observed genomic changes are due to the selection

459  of rare genomic rearrangements events. To determine the time scale of *S. stipitis*

460  genome evolution, we investigated the genome organisation of 72 single colonies

461  passaged daily for 8 weeks (56 passages, ~672 divisions) in SC-G+X, as its sugar

462  composition resembles what found in lignocellulose (29) (**Fig 5A**). Strains were

463  grown at 30 °C, a temperature that does not lead to any growth defect, and 37 °C, a

464  stressful temperature that strongly inhibits *S. stipits* growth (**Fig 5B**). CHEF gel

465  electrophoresis was conducted to identify possible changes in the chromosome

466  organisation of the evolved strains. This analysis identifies genome rearrangements

467  in 19/36 strains evolved at 30 °C and 12/36 strains evolved at 37 °C (Blue and

468  Magenta- **Fig 5C**). Thus, changes in chromosome organisation were detected in the

469  presence (37 °C) or absence (30 °C)  of stress. To test whether chromosome

470  rearrangements are associated with a fitness benefit, we tested the ability of the

471  parental and 37 °C-evolved strains to grow in SC-G+X media at permissive (30 C)

472  and restrictive (37 C) temperature (**Fig 5D**). This analysis demonstrates that 37 °C-

473  evolved strains with no chromosomal rearrangement grow poorly at 37 °C (**Fig 5D**).

474  In contrast, 5/12 37 °C-evolved strains with chromosome rearrangements grow

475  better than the parental strain at this restrictive temperature (**Fig 5D**). This result

476  suggests that changes in chromosome organisation have an adaptive value. Thus,

477  genome plasticity is a defining feature of the *S. stipitis* genome, and its genome can

478  rapidly change in mitotic cells propagated *in vitro*. Our results strongly suggest that

479  the extensive genomic changes can lead to adaptation to hostile environments.

480

481  **DISCUSSION**

482  Here we demonstrate that the yeast *S. stipitis* has a plastic genome and that

483  genome plasticity is linked to adaptation to hostile environments. We show that non-

484  centromeric retrotransposons are significant drivers of *S. stipitis* genome diversity.

485  These findings have important implications for developing economically viable

486  second-generation biofuels and better understanding the CTG (Ser1)-clade biology.

*Vega Estevez et al*

**Retrotransposon are drivers of *S. stipitis* genome diversity**

487   Our repetitive sequence analysis demonstrates that *S. stipitis* has a DNA repeats

488   content typical of the CTG (Ser1)-clade including TEs, non-canonical terminal

489   telomeric repeats and subtlomeric regions. As observed in other members of the

490   CTG (Ser1)-clade (60), we did not detect any DNA-transposons or MRS repeats.

491   One of our major findings is that that non-centromeric retrotransposons are

492   significant drivers of *S. stipitis* genome diversity. Our data support the hypothesis

493   that *S. stipitis* TEs generate genome diversity via two distinct mechanisms:

494   transposition into new genomic locations and recombination-mediated chromosome

495   rearrangements. Indeed, we demonstrated that the number and genomic position of

496   non-centromeric retrotransposons vary between the Y-11545 and Y-7124 *S. stipitis*

497   isolates. Significantly, we did not detect transposon insertions into coding regions.

498   However, transposons might alter *S. stipitis* gene expression by inserting into gene

499   regulatory regions (62). We propose that *S. stipitis* transposons are active and

500   generate genome diversity by jumping into different genomic locations. Our data also

501   indicate that TEs can generate further genome diversity  though either homologous

502   recombination of nearly identical TE copies or by faulty repair of double-strand

503   breaks generated during transposable elements excision (62). Indeed, we find that

504   the translocation breakpoint between chromosome 5 and chromosome 7 is enriched

505   in retrotransposons. Furthermore, TE-rich subtelomeric regions and centromeres

506   have a distinct organisation in the two analysed isolates suggesting that the

507   transposons drive this genetic diversity. We hypothesise that transposons elements

508   cause the genetic variability observed during laboratory passaging. In the future, it

509   will be important to apply the hybrid genome sequencing approaches presented in

510   this study to dissect the nature of these rearrangements.

**Genome plasticity and production of second-generation biofuels**

513    One of our key findings is that the *S. stipitis* genome is intrinsically plastic and that

514    chromosome rearrangements are frequent events under stress or unstressed

515    conditions. Second-generation biofuels, generated by fermentation of agriculture and

516    forestry waste, have an enormous potential to meet future energy demands and

517    significantly reduce petroleum consumption. To meet the requirements for industrial

518    applications, second-generation biofuels need to be generated by microorganisms

519    that can efficiently utilise and ferment all the sugars found in lignobiomass (63).

520    Consequently, *S. stipitis* is one of the most promising yeast for producing second-

521    generation bioethanol as it can efficiently ferment both hexose and pentose sugars

522    (25, 26, 29). However, robust economically viable *S. stipitis* platforms still require

523    significant development as this organism struggles to survive under the harsh

524    environments generated during second-generation biofuel production. For example,

525    *S. stipitis* growth and fermentation is inhibited by the chemical pre-treatment required

526    to extract glucose and xylose from lignobiomass (61). Growth is also inhibited at high

527    ethanol concentration, and *S. stipitis* ferments xylose less efficiently than glucose.

528    Evolutionary engineering approaches under selective conditions (i.e. presence of

529    inhibitory compounds, high concentration of xylose or ethanol) have been applied to

530    isolate better performing *S. stipitis* strains (61).

531    Our data predict that the genetic make-up and associated improved phenotypes of

532    superior biofuel-producer strains are unstable and that the genetic drivers of

533    improved phenotypes might be lost over time. This hypothesis could explain why

534    short-read Illumina genome sequencing has failed to identify point mutations or

535    indels that could explain the superior performance of *S. stipitis* strains (64). It is also

536    possible that *S. stipitis* superior strains carry stable complex chromosomal

537   rearrangements with a breakpoint at DNA repeats. Such rearrangements could not

538   have been identified by Illumina sequencing as short sequenced fragments will not

539   resolve changes associated with long repetitive elements. Thus, economically viable

540   use of *S. stipitis* for second-generation biofuels production will require an in-depth

541   analysis of the genomic structures of superior strains.

542

543   **Genome plasticity in the CTG (Ser1)-clade**

544   The CTG-Ser1 clade is an incredibly diverse yeast group that includes many

545   important human pathogens and non-pathogenic species (17). Our data support the

546   hypothesis that genome plasticity is a general feature of the CTG (Ser1) yeast clade

547   as it has been observed in *C. albicans* and *S. stipitis* ((21, 22, 65) and this study),

548   two organisms with a very different lifestyle. Indeed, while *C. albicans* is a diploid

549   opportunistic human fungal pathogen that lives almost exclusively in the human host,

550   *S. stipitis* is a haploid non-pathogenic yeast found in the gut of wood-ingesting

551   beetles hardwood forests or areas high in agricultural waste (29, 66). Furthermore,

552   while *C. albicans* lacks a canonical sexual cycle and its associated meiosis, *S.*

553   *stipitis* has a canonical sexual cycle whereby mating of haploid cells generate diploid

554   cells that undergo meiosis and produce haploid spores (30).

555   Our results highlight that stress regulates genome plasticity differently in *C. albicans*

556   and *S. stipitis*. It has been demonstrated that stress exacerbates *C. albicans* genome

557   instability (21, 67). In contrast, we found that *S. stipitis* genome instability is not

558   regulated by stress as we detected a similar rate of chromosomal rearrangements

559   when cells are continuously passaged in unstress (30 ºC) or stress (37 ºC)

560   conditions. Importantly, we also demonstrated that the large genomic changes are

561    associated with fitness benefits suggesting that genome plasticity is instrumental for

562    adaptation to hostile environments.

563    In summary, our study demonstrates for the first time that *S. stipitis* genome is

564    plastic. Understanding the cause and effect of this extensive genome plasticity is of

565    paramount importance to understand the biology of the CTG(Ser1)-clade of fungi.

566

## DATA AVAILABILITY

568    This Whole Genome Shotgun project has been deposited at DDBJ/ENA/GenBank

569    under the accession JADGGA000000000. The version described in this paper is

570    version JADGGA010000000. Illumina and nanopore sequence data associated with

571    this work have been deposited on the Sequence Read Archive (SRA) under

572    BioProject PRJNA609885.

573

578

## CONFLICT OF INTERESTS DISCLOSURE

580    None declared.

581

585

## REFERENCES

586

587    1.    Burns KH. 2017. Transposable elements in cancer. Nat Rev Cancer 17:415–

588          424.

589    2.    Colnaghi R, Carpenter G, Volker M, O'Driscoll M. 2011. The consequences of

590          structural genomic alterations in humans: Genomic Disorders, genomic

591          instability and cancer. Semin Cell Dev Biol 22:875–885.

592    3.    Galhardo RS, Hastings PJ, Rosenberg SM. 2007. Mutation as a Stress

593          Response and the Regulation of Evolvability. Crit Rev Biochem Mol Biol

594          42:399–435.

595    4.    Buscaino. 2019. Chromatin-Mediated Regulation of Genome Plasticity in

596          Human Fungal Pathogens. Genes (Basel) 10:855.

597    5.    Hirakawa MP, Martinez D a, Sakthikumar S, Anderson MZ, Berlin A, Gujja S,

598          Zeng Q, Zisson E, Wang JM, Greenberg JM, Berman J, Bennett RJ, Cuomo C

599          a, Aviv R. 2015. Genetic and phenotypic intra-species variation in Candida

600          albicans 1–13.

601    6.    Peter J, De Chiara M, Friedrich A, Yue J-X, Pflieger D, Bergström A, Sigwalt A,

602          Barre B, Freel K, Llored A, Cruaud C, Labadie K, Aury J-M, Istace B,

603          Lebrigand K, Barbry P, Engelen S, Lemainque A, Wincker P, Liti G,

604          Schacherer J. 2018. Genome evolution across 1,011 Saccharomyces

605          cerevisiae isolates Species-wide genetic and phenotypic diversity. Nature.

606    7.    Aguilera A, García-Muse T. 2013. Causes of Genome Instability. Annu Rev

607          Genet 47:1–32.

608    8.    Fedoroff N. 2000. Transposons and genome evolution in plants. Proc Natl

609          Acad Sci U S A. National Academy of Sciences.

610    9.    Feschotte C, Pritham EJ. 2007. DNA Transposons and the Evolution of

611     Eukaryotic Genomes. Annu Rev Genet 41:331–368.

612  10.  Wicker T, Sabot F, Hua-Van A, Bennetzen JL, Capy P, Chalhoub B, Flavell A,

613     Leroy P, Morgante M, Panaud O, Paux E, SanMiguel P, Schulman AH. 2007.

614     A unified classification system for eukaryotic transposable elements. Nat Rev

615     Genet. Nature Publishing Group.

616  11.  Havecker ER, Gao X, Voytas DF. 2004. The diversity of LTR retrotransposons.

617     Genome Biol 5.

618  12.  Llorens C, Muñoz-Pomer A, Bernad L, Botella H, Moya A. 2009. Network

619     dynamics of eukaryotic LTR retroelements beyond phylogenetic trees. Biol

620     Direct 4:1–31.

621  13.  Finnegan DJ. 1997. Transposable elements: How non-LTR retrotransposons

622     do it. Curr Biol 7:245–248.

623  14.  Januszyk K, Li PWL, Villareal V, Branciforte D, Wu H, Xie Y, Feigon J, Loo JA,

624     Martin SL, Clubb RT. 2007. Identification and solution structure of a highly

625     conserved C-terminal domain within ORF1p required for retrotransposition of

626     long interspersed nuclear element-1. J Biol Chem 282:24893–24904.

627  15.  Luan DD, Korman MH, Jakubczak JL, Eickbush TH. 1993. Reverse

628     transcription of R2Bm RNA is primed by a nick at the chromosomal target site:

629     A mechanism for non-LTR retrotransposition. Cell 72:595–605.

630  16.  Papon N, Courdavault V, Clastre M. 2014. Biotechnological potential of the

631     fungal CTG clade species in the synthetic biology era. Trends Biotechnol

632     32:167–168.

633  17.  Gabaldón T, Naranjo-Ortíz MA, Marcet-Houben M. 2016. Evolutionary

634     genomics of yeast pathogens in the Saccharomycotina. FEMS Yeast Res 16.

635  18.  Butler G, Rasmussen MD, Lin MF, Santos MAS, Sakthikumar S, Munro CA,

26

636    Rheinbay E, Grabherr M, Forche A, Reedy JL, Agrafioti I, Arnaud MB, Bates S,

637    Brown AJP, Brunke S, Costanzo MC, Fitzpatrick DA, De Groot PWJ, Harris D,

638    Hoyer LL, Hube B, Klis FM, Kodira C, Lennard N, Logue ME, Martin R,

639    Neiman AM, Nikolaou E, Quail MA, Quinn J, Santos MC, Schmitzberger FF,

640    Sherlock G, Shah P, Silverstein KAT, Skrzypek MS, Soll D, Staggs R,

641    Stansfield I, Stumpf MPH, Sudbery PE, Srikantha T, Zeng Q, Berman J,

642    Berriman M, Heitman J, Gow NAR, Lorenz MC, Birren BW, Kellis M, Cuomo

643    CA. 2009. Evolution of pathogenicity and sexual reproduction in eight Candida

644    genomes. Nature 459:657–662.

645    19.    Wohlbach DJ, Kuo A, Sato TK, Potts KM, Salamov AA, LaButti KM, Sun H,

646    Clum A, Pangilinan JL, Lindquist EA, Lucas S, Lapidus A, Jin M, Gunawan C,

647    Balan V, Dale BE, Jeffries TW, Zinkel R, Barry KW, Grigoriev I V., Gasch AP.

648    2011. Comparative genomics of xylose-fermenting fungi for enhanced biofuel

649    production. Proc Natl Acad Sci 108:13212–13217.

650    20.    Krassowski T, Coughlan AY, Shen X-X, Zhou X, Kominek J, Opulente DA,

651    Riley R, Grigoriev I V., Maheshwari N, Shields DC, Kurtzman CP, Hittinger CT,

652    Rokas A, Wolfe KH. 2018. Evolutionary instability of CUG-Leu in the genetic

653    code of budding yeasts. Nat Commun 9:1887.

654    21.    Forche A, Abbey D, Pisithkul T, Weinzierl MA, Ringstrom T, Bruck D, Petersen

655    K, Berman J. 2011. Stress alters rates and types of loss of heterozygosity in

656    Candida albicans. MBio 2.

657    22.    Todd RT, Selmecki A. 2020. Expandable and reversible copy number

658    amplification drives rapid adaptation to antifungal drugs. Elife 9:1–33.

659    23.    Todd RT, Wikoff TD, Forche A, Selmecki A. 2019. Genome plasticity in

660    Candida albicans is driven by long repeat sequences. Elife 8.

661   24.   Hirakawa MP, Martinez DA, Sakthikumar S, Anderson MZ, Berlin A, Gujja S,

662         Zeng Q, Zisson E, Wang JM, Greenberg JM, Berman J, Bennett RJ, Cuomo

663         CA. 2015. Genetic and phenotypic intra-species variation in Candida albicans.

664         Genome Res 25:413–25.

665   25.   du Preez JC, van Driessel B, Prior BA. 1989. D-xylose fermentation by

666         Candida shehatae and pichia stipitis at low dissolved oxygen levels in fed-

667         batch cultures. Biotechnol Lett 11:131–136.

668   26.   du Preez JC, van Driessel B, Prior BA. 1989. Ethanol tolerance of Pichia

669         stipitis and Candida shehatae strains in fed-batch cultures at controlled low

670         dissolved oxygen levels. Appl Microbiol Biotechnol 30:53–58.

671   27.   Sims REH, Mabee W, Saddler JN, Taylor M. 2010. An overview of second

672         generation biofuel technologies. Bioresour Technol 101:1570–1580.

673   28.   Robak K, Balcerek M. 2018. Review of second generation bioethanol

674         production from residual biomass. Food Technol Biotechnol 56:174–187.

675   29.   Suh SO, Marshall CJ, McHugh J V., Blackwell M. 2003. Wood ingestion by

676         passalid beetles in the presence of xylose-fermenting gut yeasts. Mol Ecol

677         12:3137–3145.

678   30.   Melake T, Passoth V, Klinner U. 1996. Characterization of the genetic system

679         of the xylose-fermenting yeast Pichia stipitis. Curr Microbiol 33:237–242.

680   31.   Jeffries TW, Grigoriev I V., Grimwood J, Laplaza JM, Aerts A, Salamov A,

681         Schmutz J, Lindquist E, Dehal P, Shapiro H, Jin YS, Passoth V, Richardson

682         PM. 2007. Genome sequence of the lignocellulose-bioconverting and xylose-

683         fermenting yeast Pichia stipitis. Nat Biotechnol 25:319–326.

684   32.   Lynch DB, Logue ME, Butler G, Wolfe KH. 2010. Chromosomal G + C content

685         evolution in yeasts: Systematic interspecies differences, and GC-poor troughs

28

686        at centromeres. Genome Biol Evol 2:572–583.

687   33.   Coughlan AY, Wolfe KH. 2019. The reported point centromeres of

688        Scheffersomyces stipitis are retrotransposon long terminal repeats. Yeast

689        36:275–283.

690   34.   Villa-Carvajal M, Querol A, Belloch C. 2006. Identification of species in the

691        genus Pichia by restriction of the internal transcribed spacers (ITS1 and ITS2)

692        and the 5.8S ribosomal DNA gene. Antonie van Leeuwenhoek, Int J Gen Mol

693        Microbiol 90:171–181.

694   35.   Schwartz DC, Cantor CR. 1984. Separation of yeast chromosome-sized DNAs

695        by pulsed field gradient gel electrophoresis. Cell 37:67–75.

696   36.   Ketel C, Wang HSW, McClellan M, Bouchonville K, Selmecki A, Lahav T,

697        Gerami-Nejad M, Berman J. 2009. Neocentromeres form efficiently at multiple

698        possible loci in Candida albicans. PLoS Genet 5:e1000400.

699   37.   Benson G. 1999. Tandem repeats finder: A program to analyze DNA

700        sequences. Nucleic Acids Res 27:573–580.

701   38.   Letunic I, Bork P. 2018. 20 years of the SMART protein domain annotation

702        resource. Nucleic Acids Res 46:D493–D496.

703   39.   Jeffries TW, Grigoriev I V, Grimwood J, Laplaza JM, Aerts A, Salamov A,

704        Schmutz J, Lindquist E, Dehal P, Shapiro H, Jin Y-S, Passoth V, Richardson

705        PM. 2007. Genome sequence of the lignocellulose-bioconverting and xylose-

706        fermenting yeast Pichia stipitis. Nat Biotechnol 25:319–326.

707   40.   Waterhouse AM, Procter JB, Martin DMA, Clamp M, Barton GJ. 2009. Jalview

708        Version 2--a multiple sequence alignment editor and analysis workbench.

709        Bioinformatics 25:1189–91.

710   41.   Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. 2017.

711    Canu: Scalable and accurate long-read assembly via adaptive κ-mer weighting

712    and repeat separation. Genome Res 27:722–736.

713    42.    GitHub - ruanjue/smartdenovo: Ultra-fast de novo assembler using long noisy

714    reads.

715    43.    Vaser R, Sović I, Nagarajan N, Šikić M. 2017. Fast and accurate de novo

716    genome assembly from long uncorrected reads. Genome Res 27:737–746.

717    44.    Loman NJ, Quick J, Simpson JT. 2015. A complete bacterial genome

718    assembled de novo using only nanopore sequencing data. Nat Methods

719    12:733–735.

720    45.    Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, Cuomo

721    CA, Zeng Q, Wortman J, Young SK, Earl AM. 2014. Pilon: An Integrated Tool

722    for Comprehensive Microbial Variant Detection and Genome Assembly

723    Improvement. PLoS One 9:e112963.

724    46.    Simão FA, Waterhouse RM, Ioannidis P, Kriventseva E V., Zdobnov EM. 2015.

725    BUSCO: Assessing genome assembly and annotation completeness with

726    single-copy orthologs. Bioinformatics 31:3210–3212.

727    47.    Gurevich A, Saveliev V, Vyahhi N, Tesler G. 2013. QUAST: Quality

728    assessment tool for genome assemblies. Bioinformatics 29:1072–1075.

729    48.    Min B, Grigoriev I V, Choi I-G. 2017. FunGAP: Fungal Genome Annotation

730    Pipeline using evidence-based gene model evaluation. Bioinformatics

731    33:2936–2937.

732    49.    Jones P, Binns D, Chang HY, Fraser M, Li W, McAnulla C, McWilliam H,

733    Maslen J, Mitchell A, Nuka G, Pesseat S, Quinn AF, Sangrador-Vegas A,

734    Scheremetjew M, Yong SY, Lopez R, Hunter S. 2014. InterProScan 5:

735    Genome-scale protein function classification. Bioinformatics 30:1236–1240.

736  50.  Marçais G, Delcher AL, Phillippy AM, Coston R, Salzberg SL, Zimin A. 2018.

737       MUMmer4: A fast and versatile genome alignment system. PLOS Comput Biol

738       14:e1005944.

739  51.  Krzywinski M, Schein J, Birol I, Connors J, Gascoyne R, Horsman D, Jones

740       SJ, Marra MA. 2009. Circos: An information aesthetic for comparative

741       genomics. Genome Res 19:1639–1645.

742  52.  Emms DM, Kelly S. 2015. OrthoFinder: solving fundamental biases in whole

743       genome comparisons dramatically improves orthogroup inference accuracy.

744       Genome Biol 16:157.

745  53.  Chen H, Boutros PC. 2011. VennDiagram: A package for the generation of

746       highly-customizable Venn and Euler diagrams in R. BMC Bioinformatics 12:35.

747  54.  Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-

748       Wheeler transform. Bioinformatics 25:1754–1760.

749  55.  McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A,

750       Garimella K, Altshuler D, Gabriel S, Daly M, DePristo MA. 2010. The genome

751       analysis toolkit: A MapReduce framework for analyzing next-generation DNA

752       sequencing data. Genome Res 20:1297–1303.

753  56.  Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA,

754       Handsaker RE, Lunter G, Marth GT, Sherry ST, McVean G, Durbin R. 2011.

755       The variant call format and VCFtools. Bioinformatics 27:2156–2158.

756  57.  Cingolani P, Platts A, Wang LL, Coon M, Nguyen T, Wang L, Land SJ, Lu X,

757       Ruden DM. 2012. A program for annotating and predicting the effects of single

758       nucleotide polymorphisms, SnpEff: SNPs in the genome of Drosophila

759       melanogaster strain w1118; iso-2; iso-3. Fly (Austin) 6:80–92.

760  58.  Jackson AP, Gamble JA, Yeomans T, Moran GP, Saunders D, Harris D, Aslett

761    M, Barrell JF, Butler G, Citiulo F, Coleman DC, de Groot PWJ, Goodwin TJ,

762    Quail MA, McQuillan J, Munro CA, Pain A, Poulter RT, Rajandream M-A,

763    Renauld H, Spiering MJ, Tivey A, Gow NAR, Barrell B, Sullivan DJ, Berriman

764    M. 2009. Comparative genomics of the fungal pathogens Candida dubliniensis

765    and Candida albicans. Genome Res 19:2231–44.

766 59. Frank AC, Wolfe KH. 2009. Evolutionary capture of viral and plasmid DNA by

767    yeast nuclear Chromosomes. Eukaryot Cell 8:1521–1531.

768 60. Butler G, Rasmussen MD, Lin MF, Santos M a S, Sakthikumar S, Munro C a,

769    Rheinbay E, Grabherr M, Forche A, Reedy JL, Agrafioti I, Arnaud MB, Bates S,

770    Brown AJP, Brunke S, Costanzo MC, Fitzpatrick D a, de Groot PWJ, Harris D,

771    Hoyer LL, Hube B, Klis FM, Kodira C, Lennard N, Logue ME, Martin R,

772    Neiman AM, Nikolaou E, Quail M a, Quinn J, Santos MC, Schmitzberger FF,

773    Sherlock G, Shah P, Silverstein K a T, Skrzypek MS, Soll D, Staggs R,

774    Stansfield I, Stumpf MPH, Sudbery PE, Srikantha T, Zeng Q, Berman J,

775    Berriman M, Heitman J, Gow N a R, Lorenz MC, Birren BW, Kellis M, Cuomo

776    C a. 2009. Evolution of pathogenicity and sexual reproduction in eight Candida

777    genomes. Nature 459:657–62.

778 61. Slininger PJ, Shea-Andersh M a, Thompson SR, Dien BS, Kurtzman CP,

779    Balan V, da Costa Sousa L, Uppugundla N, Dale BE, Cotta M a. 2015.

780    Evolved strains of Scheffersomyces stipitis achieving high ethanol productivity

781    on acid- and base-pretreated biomass hydrolyzate at high solids loading.

782    Biotechnol Biofuels 8:1–27.

783 62. Bourque G, Burns KH, Gehring M, Gorbunova V, Seluanov A, Hammell M,

784    Imbeault M, Izsvák Z, Levin HL, Macfarlan TS, Mager DL, Feschotte C. 2018.

785    Ten things you should know about transposable elements 06 Biological

32

*Vega Estevez et al*

786        Sciences 0604 Genetics. Genome Biol 19:1–12.

787  63.  Jansen MLA, Bracher JM, Papapetridis I, Verhoeven MD, de Bruijn H, de Waal

788        PP, van Maris AJA, Klaassen P, Pronk JT. 2017. Saccharomyces cerevisiae

789        strains for second-generation ethanol production: from academic exploration to

790        industrial implementation. FEMS Yeast Res. Oxford Academic.

791  64.  Smith D, Quinlan A. 2008. Rapid whole-genome mutational profiling using

792        next-generation sequencing technologies. Genome Res 1638–1642.

793  65.  Selmecki A, Forche A, Berman J. 2010. Genomic plasticity of the human

794        fungal pathogen Candida albicans. Eukaryot Cell 9:991–1008.

795  66.  Fisher MC, Hawkins NJ, Sanglard D, Gurr SJ. 2018. Worldwide emergence of

796        resistance to antifungal drugs challenges human health and food security.

797        Science (80- ) 360:739–742.

798  67.  Freire-Benéitez V, Gourlay S, Berman J, Buscaino A. 2016. Sir2 regulates

799        stability of repetitive domains differentially in the human fungal pathogen

800        Candida albicans. Nucleic Acids Res https://doi.org/10.1093/nar/gkw594.

801  68.  Jackson AP, Gamble J a, Yeomans T, Moran GP, Saunders D, Harris D, Aslett

802        M, Barrell JF, Butler G, Citiulo F, Coleman DC, de Groot PWJ, Goodwin TJ,

803        Quail M a, McQuillan J, Munro C a, Pain A, Poulter RT, Rajandream M-A,

804        Renauld H, Spiering MJ, Tivey A, Gow N a R, Barrell B, Sullivan DJ, Berriman

805        M. 2009. Comparative genomics of the fungal pathogens Candida dubliniensis

806        and Candida albicans. Genome Res 19:2231–44.

807

808  **FIGURE LEGENDS**

809  **Figure 1**. Classification of non-centromeric *S. stipitis* repeats

810  **A)** Schematics of non-centromeric retrotransposons identified in this study.

33

811 For each transposon, subclass, superfamily and family is indicated. The organisation

812 of coding and non-coding sequences of each transposon is displayed. **B)** Cladogram

813 showing CTG (Ser1)-clade species with known transposable elements (this study

814 and (60, 68). The presence (V) or absence (X) of a TE is indicated. **C)** Sequence

815 alignment of telomeric terminal repeats in members of the CTG (Ser1)-clade (*C.*

816 *lusitaniae*, *S. stipitis*, *L. elongisporus*, *C. albicans* and *C. tropicalis*) (this study and

817 (60, 68). Consensus sequence to the *S. cerevisiae* telomeric repeats is indicated

818 (Magenta box). **D)** Schematics of gene family members associated with *S. stipitis*

819 subtelomeres (30 Kb from chromosome end).

820

821 **Figure 2**. Phenotypic and Genotypic Diversity in *S. stipitis*

822 **A)** Heatmaps comparing growth rate (Left), maximum OD (Middle) and lag time

823 (Right) for each *S. stipitis* natural isolate in comparison to the reference Y-11545

824 strain (blue). Analyses were performed in Glucose (G), Xylose (X) and

825 Glucose/Xylose (G+X) media. The heatmap data are the average of 3 biological

826 replicates. **B)** Karyotyping of *S. stipits* natural isolates by CHEF electrophoresis. Y-

827 11545 strain is highlighted in blue and the size of its eight chromosome is indicated.

828

829 **Figure 3**. Differences in TE distribution and organisation

830 **A)** The genomic organisation of Y-11545 and Y-7124 is distinct .*Left:* Schematics of

831 Y-11545 chromosome organisation. Chromosome (Chr) number and size (Mbp) is

832 indicated. *Middle:* Karyotyping of *S. stipits* Y-11545 and Y-7124 strains by CHEF

833 electrophoresis. *Right:* Schematics of Y-7124 chromosome organisation.

834 Chromosome (Chr) number and size (Mbp) is indicated **B)** Schematics of non-

835 centromeric transposon family distribution in Y-11545 (*left*) and Y-7124 (*right*) **C)**

836     Copy Number of full-length transposons (Left) and transposon-associated repeats

837     associated with the Y-11545 (dark grey) and Y-7124 (light grey) genome. **D)**

838     Percentage (%) of Ancestral, Deletion and Insertion sites associated with the Y-

839     11545 and Y-7124 genomes

840     **Figure 4**. Chromosome rearrangements between *S. stipitis* natural isolates

841     **A)** Circos plot displaying macrosynteny between Y-11545 (Left) and Y-7124 (Right).

842     Chromosome (Chr) number and size is indicated. Recriprocal translocation between

843     the two genomes is highlighted in purple and pink. **B)** Schematics of repetitive

844     sequences associated with the translocation junction in the Y-11545 (Left) and Y-

845     7124 (Right) genomes. **C)** Subtelomeric gene families and TEs distribution in the Y-

846     11545 and Y-7124 genomes **D)** Schematics of centromere organisation in the Y-

847     11545 (Left) and Y-7124 (Right)

848     **Figure 5** The *S. stipitis* genome is plastic following real-time evolution

849     **A)** Schematics of laboratory evolution strategy **B)** S. stipitis growth curve in SC G+X

850     liquid media at permissive (30 ºC) and restrictive (37 ºC) temperature **C)** Karyotype

851     organisation of *S. stipitis* colonies following 8 weeks of laboratory evolution at 30 ºC

852     and 37 ºC. Variation in the structure following laboratory evolution at 30 ºC (blue)

853     and 37 ºC (magenta) is indicated **D)** Serial dilution assay showing growth of parental

854     (P) and 37 ºC-evolved strains without (NO Rearrangements) and with

855     (Rearrangements) at 30 ºC and 37 ºC. The CHEF analysis strain number is indicated

856     (Magenta). * indicates colonies with a fitness advantage compared to the parental

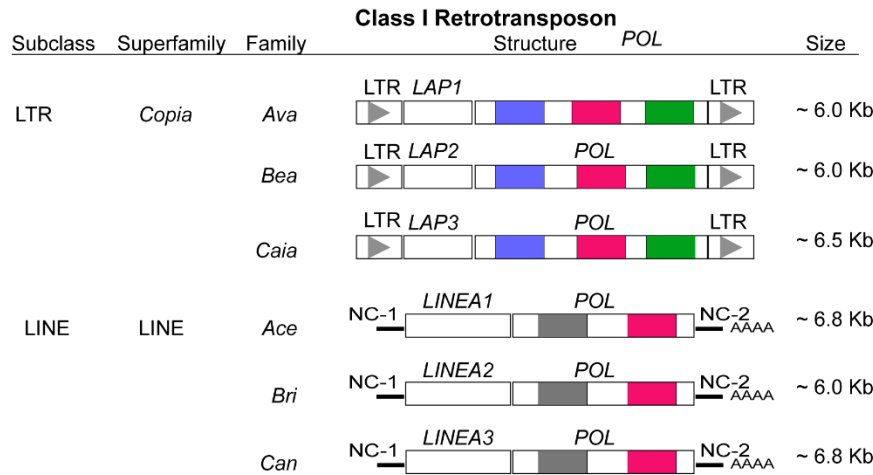857     strain.

858     **Figure S1.** Southern Blot analysis confirm the Chr5 / Chr7 translocation

859     Left: Schematics of chromosome 5 (Chr 5) and chromosome 7 (Chr 7) in Y-11545

860     (Left) and Y-7124 (Right). Reciprocal translocation is highlighted in purple and pink.

*Vega Estevez et al*

861    Southern Probe is indicated. Right: Southern Blot of Y-11545 and Y-7124

862    chromosomes separated by CHEF gel electrophoresis. Full chromosome profiling

863    (EtBr) and Southern Blot results (Southern) are indicated
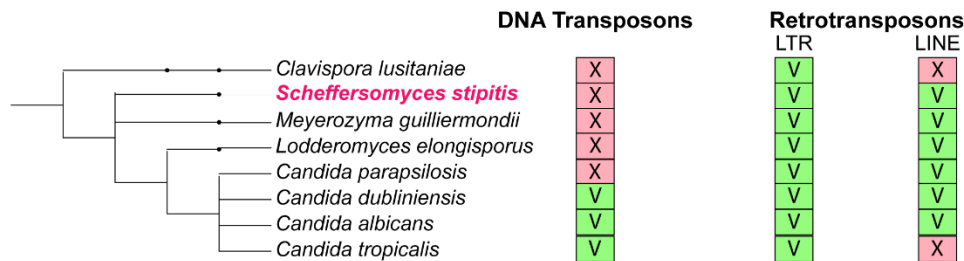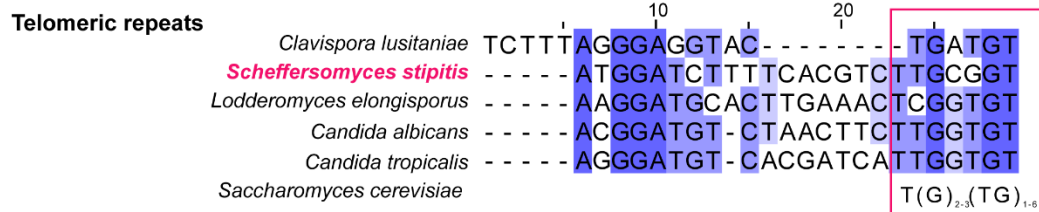
864

**Fig 1**

**A**



**B**



**C**



**D**



865

**Fig 2**

**A**



**B**



Chromosome Organisation : CHEF

**Fig 3**

**Fig 4**                                                                    *Vega-Estevez et al*

*Vega Estevez et al*

**Fig 5**

**A**

*Vega-Estevez et al*



**30 °C**  |---- 56 Passages ----|

36 Single Colonies

SC G+X

**37 °C**  |---- 56 Passages ----|

36 Single Colonies

**B**



**C**  Chromosome Organisation: CHEF Gel Electrophoresis

**30 °C**



**37 °C**



**D**



NO Rearrangements | Chromosome Rearrangements

30 °C

37 °C

869

41