

1 Supplementary Materials for

2 **Mitochondrial introgression suggests extensive ancestral hybridization events**

3 **among *Saccharomyces* species**

4 **Authors:** David Peris, Armando Arias, Sandi Orlic, Carmela Belloch, Laura Pérez-
5 Través, Amparo Querol, Eladio Barrio

6
7
8 Correspondence to: david.perisnavarro@wisc.edu

9
10 **This PDF file includes:**

11 Supplementary Text

Table of Contents

12		
13		
14		
15	COX2 nucleotide statistics	<u>2</u>
16		
17	Maximum-Likelihood phylogenetic tree of COX2 3'end	
18	region is incongruent with the species tree due to a	
19	recombination hotspot	<u>2-3</u>
20		
21	ORF1 gene structure	<u>3-5</u>
22		
23	ORF1 detection of selection	<u>5-6</u>
24		
25	References	<u>6-8</u>
26		
27	Supplementary Tables and Figures	<u>8-33</u>
28		
29		
30		

31 **Supplementary Text**

32 **COX2 nucleotide statistics**

33 A total of five hundred and seventeen COX2 sequences of worldwide distributed wild
34 *Saccharomyces* strains and forty-nine sequences of domesticated natural hybrids were
35 sequenced or retrieved from public databases (Table S1, Figure 1). The COX2 alignment
36 is 585 bp, with 106 variable positions (18.11%) (Figure S2) where 43.4% of
37 polymorphisms are found in the last 89 nucleotides. 80 of the variable positions are
38 phylogenetically informative. 27 variable sites correspond to 0-fold degenerated, but only
39 10 of them are informative. 20 variable sites correspond to 2-fold degenerated positions
40 (19 informative and 1 singleton), being 5 of them non-synonymous substitutions. Finally,
41 42 variable sites are 4-fold degenerated and 38 of them are informative.

42 COX2 genetic diversity statistics are summarized in Table S4. The most diverse
43 sequences were found in three species, *S. cerevisiae*, *S. paradoxus* and *S. mikatae*. In
44 the case of *S. cerevisiae* nucleotide diversity values of 2.1% are so much higher than
45 values inferred for the nuclear genome, around 0.9% (Liti et al. 2009; Wang et al. 2012),
46 suggesting that recombination is increasing the values. For *S. eubayanus* and *S. uvarum*,
47 the COX2 nucleotide diversity values, 0.5 and 0.6% respectively, are lower than those
48 found for the nuclear genome, around 0.8% (Patagonia A-Patagonia B *S. eubayanus*)
49 and 1% (Holarctic-South America B *S. uvarum*) (Almeida et al. 2014; Peris et al. 2014) in
50 agreement with the low mutation rate of yeast mitochondrial genomes (Clark-Walker
51 1991). COX2 values for *S. kudriavzevii* strains are similar to those inferred for the nuclear
52 genome, 1.11% (Hittinger et al. 2010).

53

54 **Maximum-Likelihood phylogenetic tree of COX2 3'end region is incongruent with** 55 **the species tree due to a recombination hotspot**

56 COX2 alignment was split in two segments based on the most common recombination
57 point. We reconstructed the Maximum-Likelihood tree with PhyML v3.0 (Guindon and
58 Gascuel 2003) using the best fitted substitution model inferred by jModeltest (Posada
59 2008). Phylogenetic trees comparison with the *Saccharomyces* species tree was
60 performed in Tree Puzzle v5.2 (Schmidt et al. 2002) implemented with the
61 conservative Shimodaira-Hasegawa (Shimodaira and Hasegawa 1999) and Expected

62 Likelihood weights (ELW) (Strimmer and Rambaut 2002) tests. Both COX2 5'end
63 segment and the species tree agree each other; however, the COX2 3'end was best
64 explained by the inferred ML using that sequence region, and significantly rejected the
65 topology of the species tree (p -value $< 1 \cdot 10^{-5}$), suggesting incongruence between the
66 inferred ML COX2 3'end phylogenetic tree and the species tree.

67 A set of COX2 representative sequences of each haplotype was the input of
68 RDPv3.44 package (Martin et al. 2010). RDP includes six methods to detect
69 recombination: RDP (Martin and Rybicki 2000), Bootscanning (Salminen et al. 1995),
70 MaxChi (Smith 1992), Chimaera (Smith 1992), GeneConv (Padidam et al. 1999) and Sis-
71 scan (Gibbs et al. 2000). Default settings were set up with a statistical significance
72 threshold of p -value < 0.05 , with Bonferroni correction for multiple comparisons.

73 Codons (nucleotide sequences 526-528 and 535-537) in the COX2 3' end show
74 several convergent nucleotide substitutions among *S. cerevisiae* and *S. paradoxus*
75 strains. This finding might be the result of a “patchy-tachy” effect due to different
76 substitution rate in those positions (Sun et al. 2011). Homoplasies can drive to
77 recombinant false positives. One confounding effect can be observed in L1528 and
78 CECT11757 strains which is difficult to know if they were truly introgressed or they
79 suffered two different recombination events. A recent study shown transfers from
80 European *S. paradoxus* GC48 cluster to L1528. Also the presence of type II *ORF1*
81 segments supports the introgression scenario and a double recombination scenario.

82

83 ***ORF1* gene structure**

84 52 different haplotypes were found in our representative strains (Table S1). *ORF1*
85 start codon is nineteen nucleotides inside the 3' end of COX2 gene. The average GC-
86 composition of the *ORF1* is extremely low, 18%. Eleven strains have a GTG as a start
87 codon, which is uncommon in yeast mitochondria: haplotypes M2-M7 and M9-M12. The
88 translation of the *ORF1* gene predicts fourteen strains having a premature stop codon:
89 haplotypes M2, M3, M7-M11, M17-M20, M23 and M46. *ORF1* sequence length range
90 from 1363bp (ZA17 strain) to 1516bp (VRB strain), translated to 454 and 505,
91 respectively. Note that we sequenced a partial *ORF1* gene, 45 nucleotides before the last
92 base pair annotated for the S288c *ORF1* gene.

93 *ORF1* secondary structure and domains (LAGLIDADG and NUMOD) were annotated
94 in Jalview according to Dalgaard *et al.* (1997) and using the Conserved Domains tool
95 in NCBI (<http://www.ncbi.nlm.nih.gov/Structure/cdd/wrpsb.cgi>) (Marchler-Bauer *et al.*
96 2009). WebLogo profiles of LAGLIDADG and NUMOD1 domains were done in WebLogo
97 2.8.2 tool (<http://weblogo.berkeley.edu/>). Three different domains were found in all
98 *ORF1* sequences: two LADGLIDADG (P1 and P2) and one NUMOD1 (Figure 4, S10 and
99 S11). Two non-*Saccharomyces ORF1* were detected by PSI-BLAST (Altschul *et al.*
100 1997) showing some similarities with the *Saccharomyces ORF1* sequences. We
101 generated a new alignment in Jalview 4.0.b2 (Waterhouse *et al.* 2009) including the
102 *Williopsis saturnus* var. *suaveolans* (*ORF1* and *ORF3*), and one from *Kazachstania*
103 *servazii* (SasefMp08) sequences (Figure S11). The sequence from *K. servazii* was found
104 to be closely related to the *Saccharomyces ORF1* sequences (Figure S5). Although, *W.*
105 *saturnus* is phylogenetically fairly distant from *Saccharomyces* genus (James SA *et al.*
106 1998), they have been diverging since 235mya, we observed some conservation among
107 LAGLIDADG and NUMOD aminoacids when compared with *Saccharomyces* and
108 *Kazachstania ORF1* sequences (Figure S10 and S11).

109 Differences in size among *ORF1* gene sequences were due to the presence of GC
110 insertions (Figure S5) and AT repeats (Table S3). Seven different GC clusters insertion
111 points were found along the *ORF1* alignment. GC cluster 1 is observed in VRB; three
112 types of GC cluster 2 were found in CBS435, CECT11757 and 120MX (Figure S5 and
113 S12); GC cluster 3 is shown in CBS10644; GC cluster 4 is shown in VRB; GC cluster 5
114 is inserted in *ORF1* from CBS435, and was structurally similar to GC cluster 3 of
115 CBS10644; GC cluster 6 and 7 were the most frequent among *Saccharomyces* strains
116 (Figure S9). 40 *S. cerevisiae* and one hybrid *S. cerevisiae* x *S. kudriavzevii* (AMH) have
117 the GC cluster 6, and 36 *S. cerevisiae*, 2 hybrids *S. cerevisiae* x *S. kudriavzevii*, 2 Far
118 Eastern and 2 American *S. paradoxus* shown the GC cluster 7 in their *ORF1* sequences
119 (Figure S9). Some GC clusters sequences were found inverted (Figure S5). GC cluster 7
120 of the two American *S. paradoxus* (120MX and CBS5313) and three *S. cerevisiae*, two
121 from Japan and one from USA (CBS435, Y9 and YPS606), showed identical GC cluster
122 sequence (Figure S9).

123 *ORF1* GC clusters were classified according to Zamaroczy and Bernardi (1986),
124 except for GC cluster 2. GC cluster 1 and 4 are similar to a1 family, and GC cluster 3 and
125 5 were similar to a4 family. In the case of GC cluster 1, 2, 4, and 5 are on the opposite
126 strand. As Séraphin *et al.* (1987) and Weiller *et al.* (1989) described, the GC clusters of
127 the *ORF1* gene were flanked by TAG and AGGAG, or CTA and CTCCT when cluster was
128 inserted in the opposite strand (Figure S12). These conserved nucleotides were flanked
129 by A+T rich sequences. Flanking sequences TAG and AG (CTA and CT) are conserved
130 in most of the sequences with and without GC clusters. All GC clusters in *ORF1* belong
131 to group M1 (Weiller *et al.* 1989). Conserved flanking regions might help to the
132 transposition of GC clusters.

133 Differences in size were also due to the presence of AT repeats. Tandem repeat
134 sequences for *COX2* and *ORF1* genes were detected and described using `Tandem`
135 `Repeat Finder` software (Benson 1999). Twenty-one different A+T rich sequences that
136 repeated at least twice were found in the *ORF1* alignment (Table S3). The length of A+T
137 rich tandem repeats ranged from three nucleotides to twenty five nucleotides. The most
138 repeated sequence (AAT) was repeated ten times in haplotypes M1, M12, M22, M39,
139 M40, and M50. The A+T rich tandem repeats were located near to GC clusters (Figure
140 4). In *COX2* gene, we found three A+T rich sequences repeated twice (Table S3).

141

142 ***ORF1* detection of selection**

143 The single likelihood ancestor counting (SLAC), fixed effects likelihood (FEL), and
144 random effects likelihood (REL) methods (Pond and Frost 2005), implemented in the
145 `HYPHY` web based version, `Datamonkey` (Delport *et al.* 2010), were used to detect the
146 signatures of selection operating on *ORF1* protein gene. NJ phylogenetic trees were
147 reconstructed under the REV (GTR) substitution model. Two different approaches were
148 performed to describe selection signatures. In the first approach, a complete *ORF1*
149 alignment, without GC Insertions and indels, partitioned by recombinant sections detected
150 by GARD, was used. In a second study, the two LADGLIDADG and the NUMOD domains
151 were analyzed independently.

152 Aminoacid positions were described as being under positive or purifying selection
153 when significant values were generated by two of the three tested methods. Codon-

154 specific selection pressure along the sequences (i.e. site specific dN-dS) was measured
155 and *p*-values were estimated at each site. We analyzed 417 *ORF1* codons of the total
156 458. A 61 codons, corresponding to 15% of the total were found to be under purifying
157 selection, where 43 of the 61 codons were found in LAGLIDADGs and NUMOD1 domains
158 (Figure S10). A functional characterization might help to understand how important these
159 aminoacids are for the activity of this homing endonuclease.

160

161 **References**

- 162 Almeida P, Gonçalves C, Teixeira S *et al.* 2014. A Gondwanan imprint on global diversity
163 and domestication of wine and cider yeast *Saccharomyces uvarum*. *Nat Commun.* 5:
164 Altschul SF, Madden TL, Schaffer AA, Zhang J, Zhang Z, Miller W, Lipman DJ. 1997.
165 Gapped BLAST and PSI-BLAST: a new generation of protein database search programs.
166 *Nucl Acids Res.* 25:3389-3402.
- 167 Benson G. 1999. Tandem repeats finder: a program to analyze DNA sequences. *Nucl*
168 *Acids Res.* 27:573-580.
- 169 Clark-Walker GD. 1991. Contrasting mutation rates in mitochondrial and nuclear genes
170 of yeasts versus mammals. *Curr Genet.* 20:195-198.
- 171 Dalgaard JZ, Klar AJ, Moser MJ, Holley WR, Chatterjee A, Mian IS. 1997. Statistical
172 modeling and analysis of the LAGLIDADG family of site- specific endonucleases and
173 identification of an intein that encodes a site-specific endonuclease of the HNH family.
174 *Nucl Acids Res.* 25:4626-4638.
- 175 de Zamaroczy M, Bernardi G. 1986. The GC clusters of the mitochondrial genome of
176 yeast and their evolutionary origin. *Gene.* 41:1-22.
- 177 Delport W, Poon AFY, Frost SDW, Kosakovsky Pond SL. 2010. Datamonkey 2010: a
178 suite of phylogenetic analysis tools for evolutionary biology. *Bioinformatics.* 26:2455-
179 2457.
- 180 Gibbs MJ, Armstrong JS, Gibbs AJ. 2000. Sister-Scanning: a Monte Carlo procedure for
181 assessing signals in recombinant sequences. *Bioinformatics.* 16:573-582.
- 182 Guindon S, Gascuel O. 2003. A simple, fast, and accurate algorithm to estimate large
183 phylogenies by Maximum Likelihood. *Syst Biol.* 52:696-704.

184 Hittinger CT, Gonçalves P, Sampaio JP, Dover J, Johnston M, Rokas A. 2010.
185 Remarkably ancient balanced polymorphisms in a multi-locus gene network. *Nature*.
186 464:54-58.

187 James SA, Roberts IN, Collins MD. 1998. Phylogenetic heterogeneity of the genus
188 *Williopsis* as revealed by 18S rRNA gene sequences. *Int J Syst Bacteriol*. 48:591-596.

189 Leducq JB, Charron G, Samani P *et al*. 2014. Local climatic adaptation in a widespread
190 microorganism. *Proc R Soc Lond B Biol Sci*. 281:

191 Liti G, Carter DM, Moses AM *et al*. 2009. Population genomics of domestic and wild
192 yeasts. *Nature*. 458:337-341.

193 Livingstone C, Barton GJ. 1993. Protein sequence alignments: a strategy for the
194 hierarchical analysis of residue conservation. *Comput Appl Biosci*. 9:745-756.

195 Marchler-Bauer A, Anderson JB, Chitsaz F *et al*. 2009. CDD: specific functional
196 annotation with the Conserved Domain Database. *Nucleic Acids Res*. 37:D205-D210.

197 Martin D, Rybicki E. 2000. RDP: detection of recombination amongst aligned sequences.
198 *Bioinformatics*. 16:562-563.

199 Martin DP, Lemey P, Lott M, Moulton V, Posada D, Lefevre P. 2010. RDP3: a flexible
200 and fast computer program for analyzing recombination. *Bioinformatics*. 26:2462-2463.

201 Padidam M, Sawyer S, Fauquet CM. 1999. Possible emergence of new Geminiviruses
202 by frequent recombination. *Virology*. 265:218-225.

203 Peris D, Sylvester K, Libkind D, Gonçalves P, Sampaio JP, Alexander WG, Hittinger CT.
204 2014. Population structure and reticulate evolution of *Saccharomyces eubayanus* and its
205 lager-brewing hybrids. *Mol Ecol*. 23:2031-2045.

206 Pond SLK, Frost SDW. 2005. Datamonkey: rapid detection of selective pressure on
207 individual sites of codon alignments. *Bioinformatics*. 21:2531-2533.

208 Posada D. 2008. jModelTest: phylogenetic model averaging. *Mol Biol Evol*. 25:1253-
209 1256.

210 Salminen MO, Carr JK, Burke DS, McCutchan FE, Garcia-Martinez J. 1995. Identification
211 of breakpoints in intergenotypic recombinants of HIV type 1 by bootscanning. *AIDS Res*
212 *Hum Retroviruses*. 11:1423-1425.

213 Schmidt HA, Strimmer K, Vingron M, von Haeseler A. 2002. TREE-PUZZLE: maximum
214 likelihood phylogenetic analysis using quartets and parallel computing. *Bioinformatics*.
215 18:502-504.

216 Séraphin B, Simon M, Faye G. 1987. The mitochondrial reading frame RF3 is a functional
217 gene in *Saccharomyces uvarum*. *J Biol Chem*. 262:10146-10153.

218 Shimodaira H, Hasegawa M. 1999. Multiple comparisons of log-likelihoods with
219 applications to phylogenetic inference. *Mol Biol Evol*. 16:1114-1116.

220 Smith JM. 1992. Analyzing the mosaic structure of genes. *J Mol Evol*. 34:126-129.

221 Strimmer K, Rambaut A. 2002. Inferring confidence sets of possibly misspecified gene
222 trees. *Proceedings Biological sciences / The Royal Society*. 269:137-142.

223 Sun S, Evans BJ, Golding GB. 2011. "Patchy-Tachy" leads to false positives for
224 recombination. *Mol Biol Evol*. 28:2549-2559.

225 Wang QM, Liu WQ, Liti G, Wang SA, Bai FY. 2012. Surprisingly diverged populations of
226 *Saccharomyces cerevisiae* in natural environments remote from human activity. *Mol Ecol*.
227 21:5404-5417.

228 Waterhouse AM, Procter JB, Martin DMA, Clamp M, Barton GJ. 2009. Jalview Version 2-
229 -a multiple sequence alignment editor and analysis workbench. *Bioinformatics*. 25:1189-
230 1191.

231 Weiller G, Schueller CME, Schweyen RJ. 1989. Putative target sites for mobile G+C rich
232 clusters in yeast mitochondrial DNA: Single elements and tandem arrays. *Mol Gen Genet*.
233 218:272-283.

234

235 **Supplementary Tables and Figures**

236 **Table S1.** Strain information.

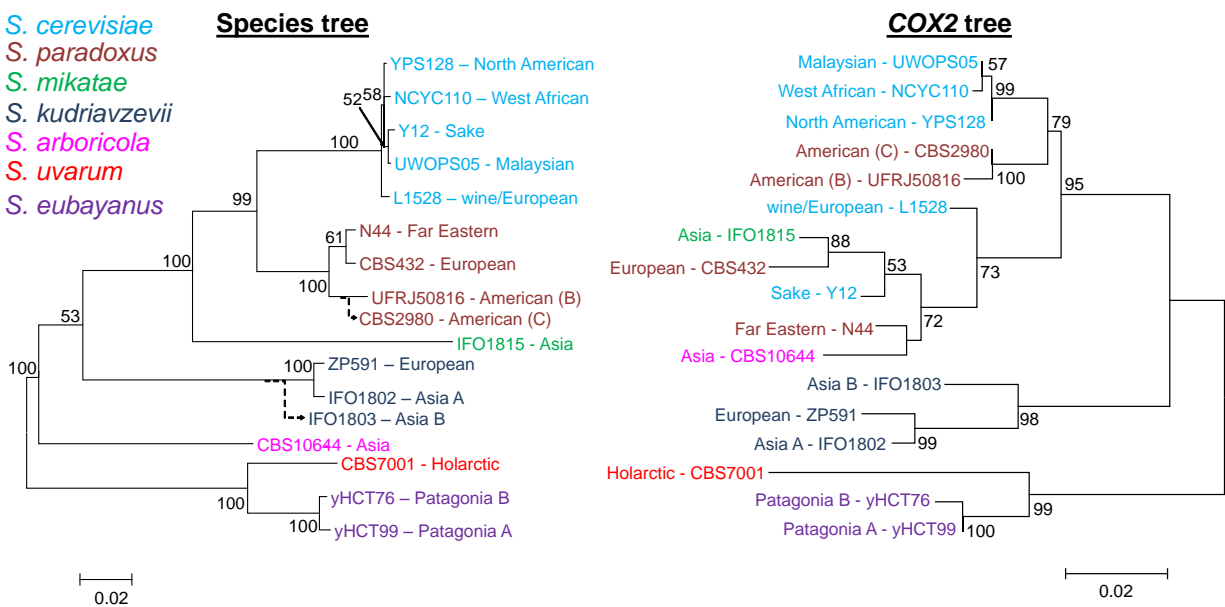
237 **Table S2.** Primer pairs used in this study.

238 **Table S3.** Distribution of AT tandem repeats.

239 **Table S4.** Summary statistics for all *Saccharomyces* and each species using COX2.

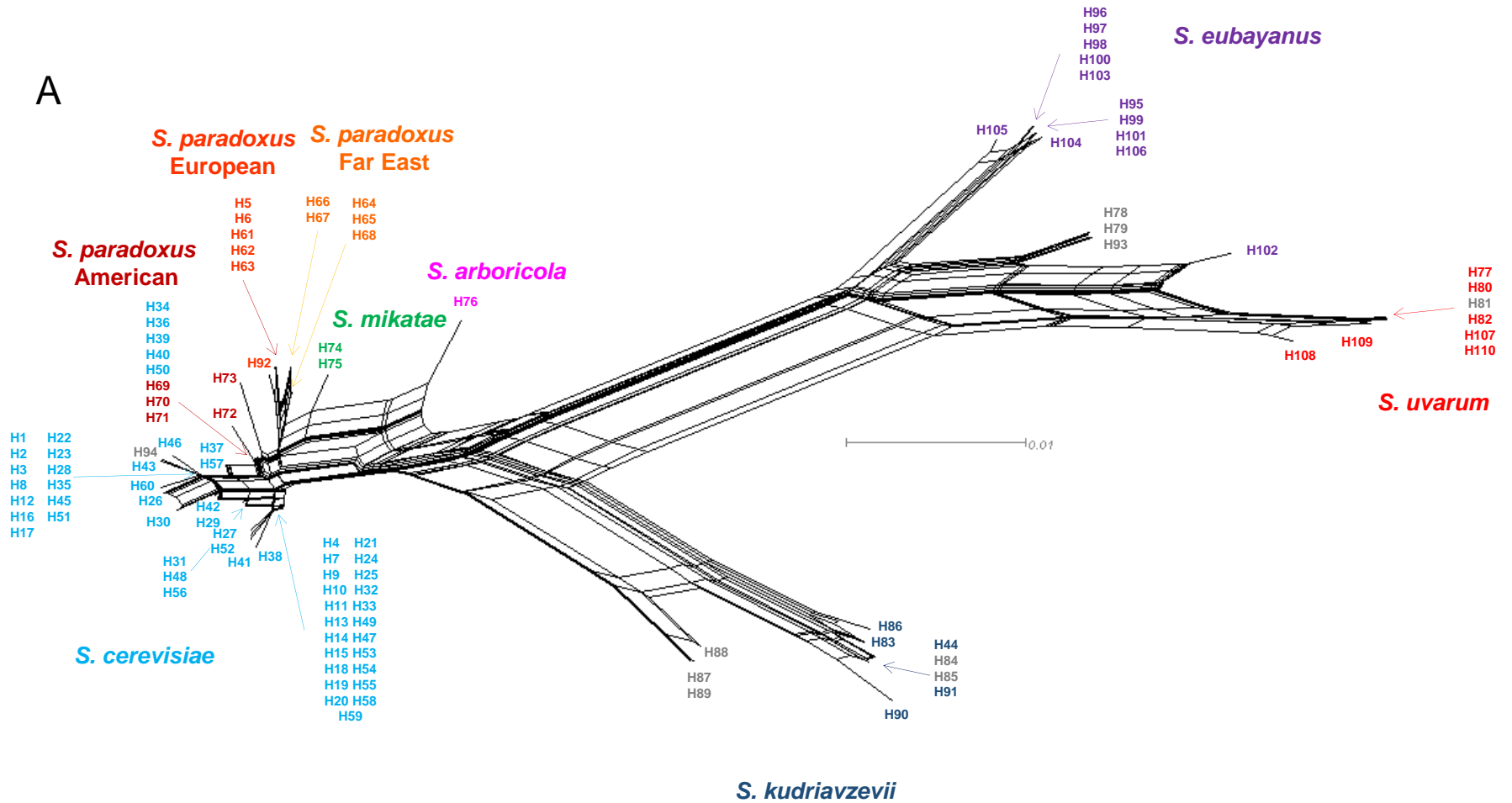
240

241 **Figure S1. Maximum-Likelihood phylogenetic trees.**



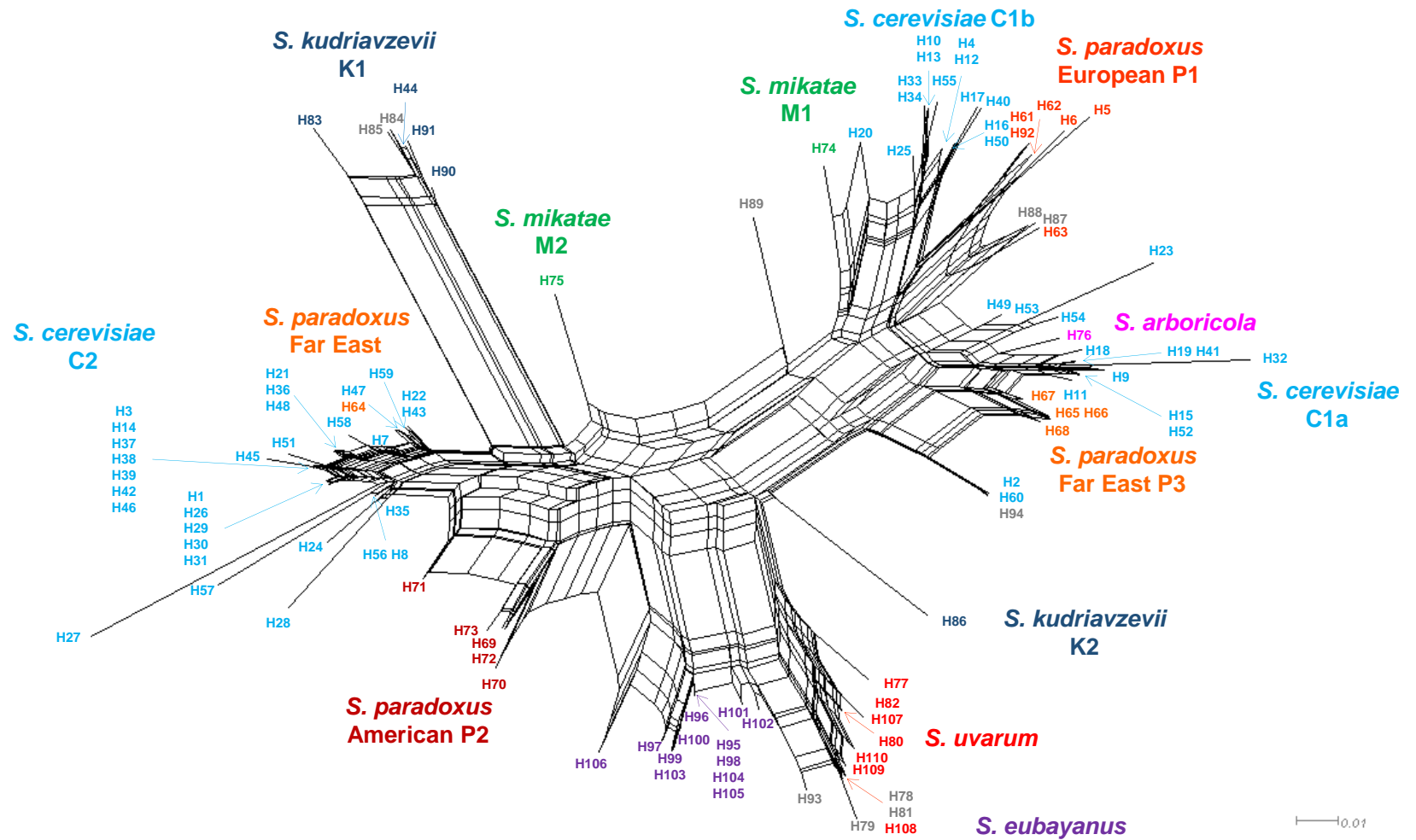
242
 243 ML phylogenetic tree using a concatenated alignment of *RIP1*, *MET2* and *FUN14* (left),
 244 and *COX2* (right) are represented. Strains were colored according to the species
 245 designation. CBS2980 and IFO1803 strain positions in the phylogenetic species tree were
 246 inferred according to Leducq *et al.* (2014) and Hittinger *et al.* (2010). Bootstrap values
 247 above 50 are given for each branch. Scales are given in nucleotide substitution per site.
 248

379 Figure S3. COX2 Neighbor-Net phylogenetic network for each COX2 segment.



380

B

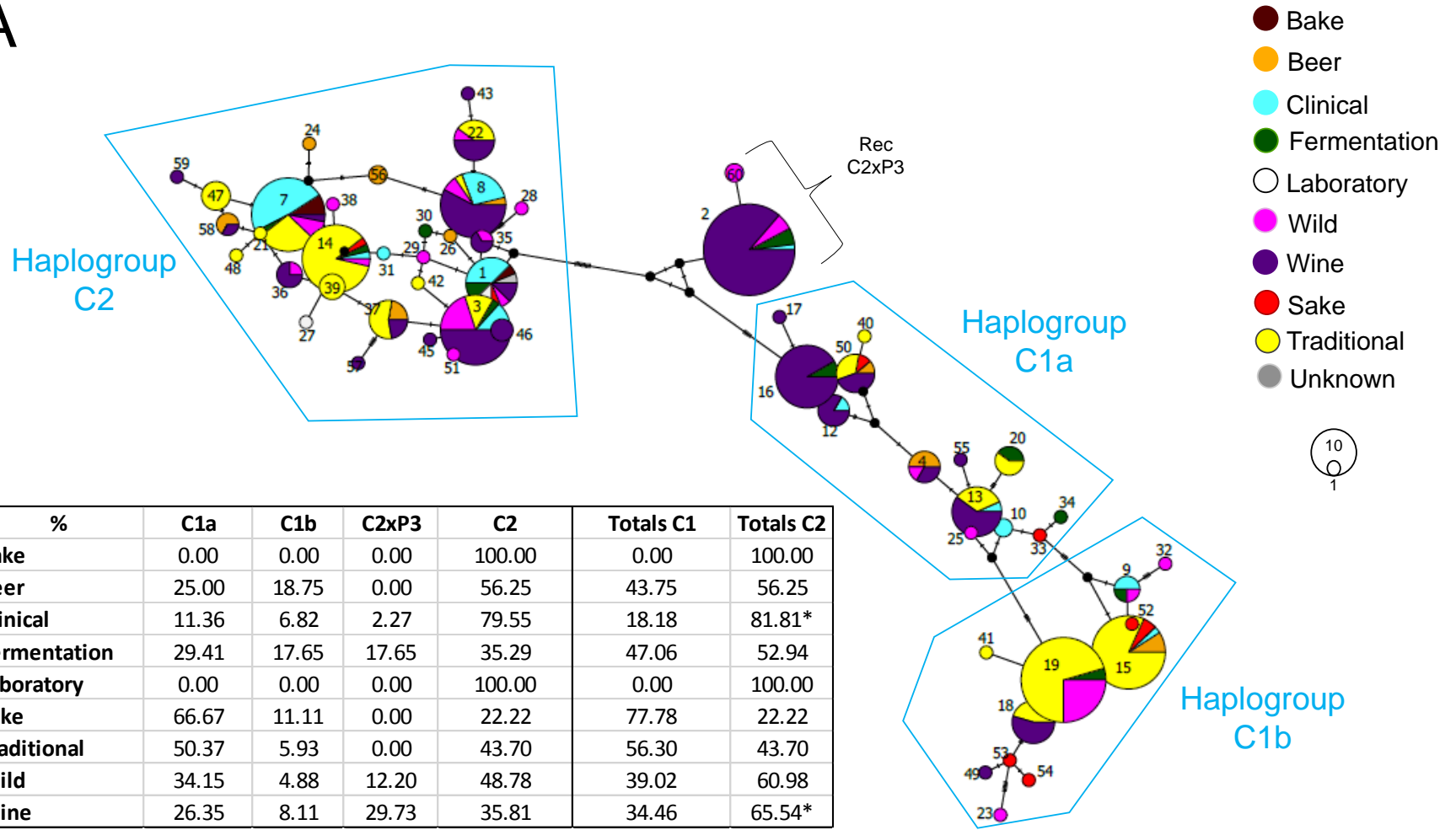


381

382 Neighbor-net phylogenetic networks for COX2 5' end and 3' end segments are represented in A) and B), respectively. Haplotypes
 383 were colored according to each species designation. Haplotypes exclusively from hybrids are colored in grey. Scale is given in
 384 nucleotide substitution per site.

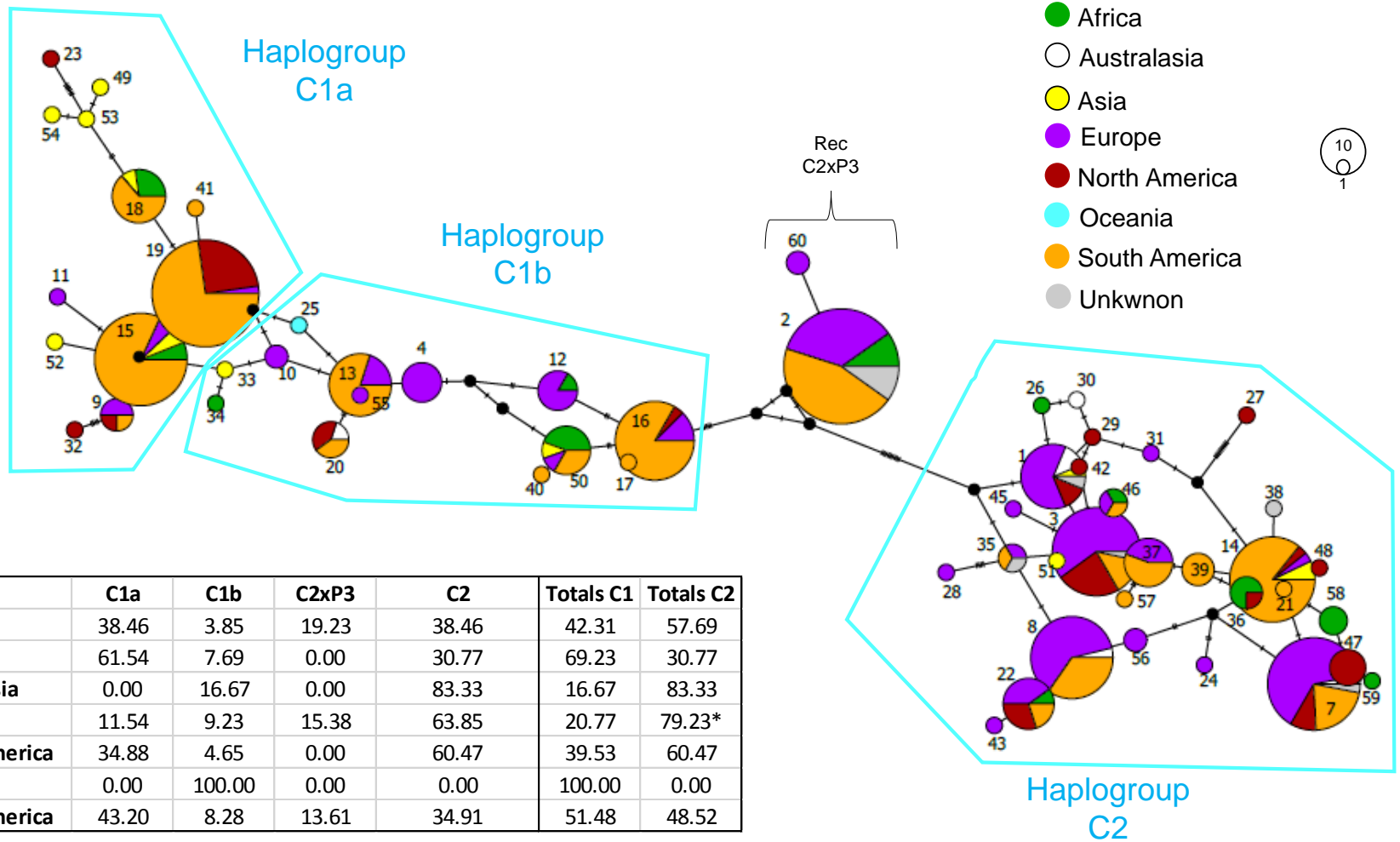
385 **Figure S4. *S. cerevisiae* COX2 MJ networks.**

A



386

B

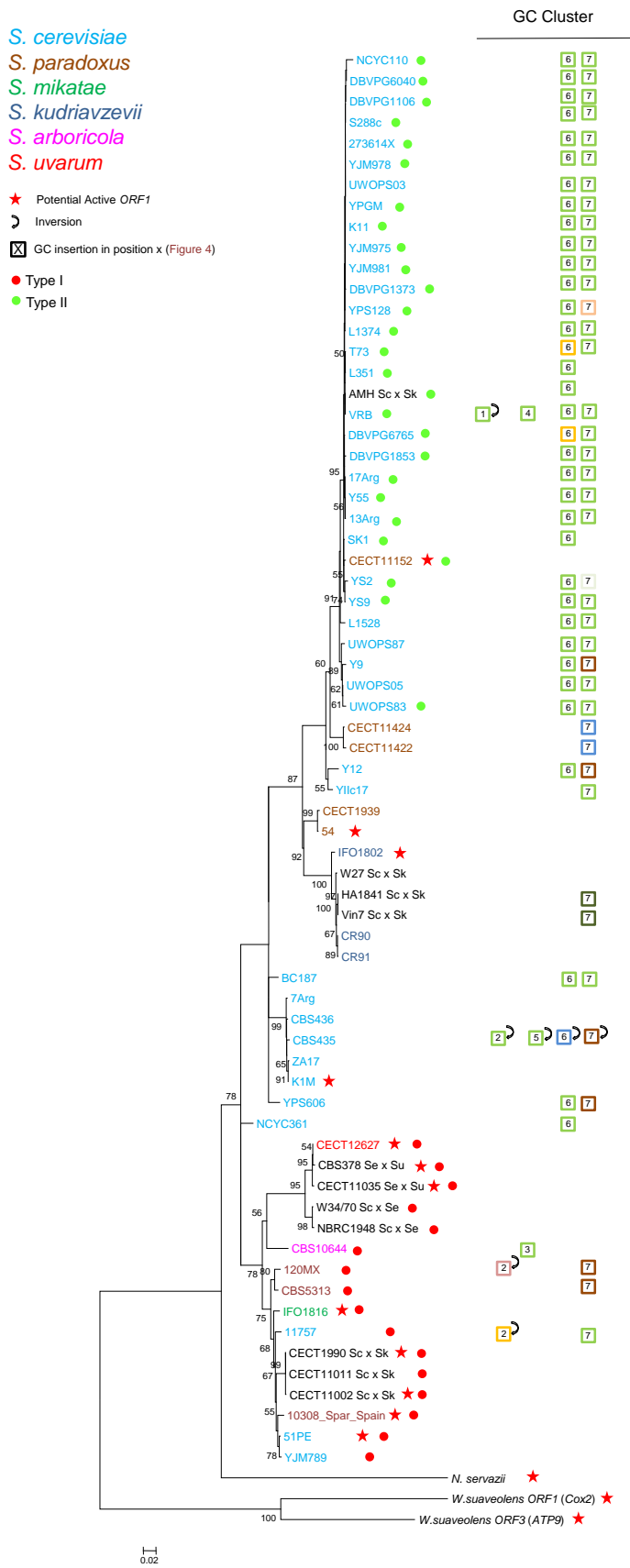


387

388 *S. cerevisiae* COX2 MJ networks were represented in A) and B) where haplotype pie charts were colored according to the
 389 isolation source and continent of isolation, respectively. Circles sizes represents the number of sequences in a haplotype and is
 390 scaled according to the legend. Number of mutations from one haplotype to another are indicated by lines in the edges connecting

391 the haplotypes. A table showing the distribution in each haplogroup or recombinant group by isolation source or continent is also
392 displayed in A) and B), respectively.

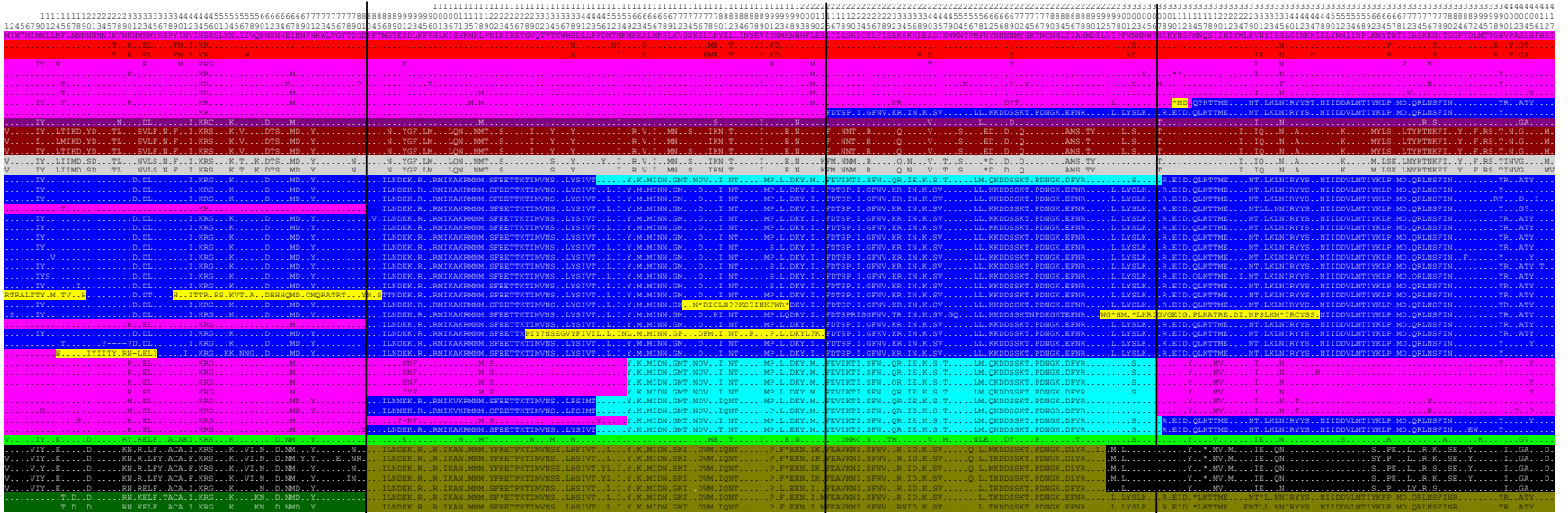
Figure S5. ORF1 Neighbor-Joining phylogenetic tree.



0.02

395 A NJ phylogenetic tree was reconstructed for the *ORF1* sequence alignment. Strains
396 were colored according to the species designation. Bootstrap values above 50 are given
397 for each branch. A red star represents a potential active *ORF1* based on the absence of
398 GC clusters and premature stop codons which disrupt the coding sequence. Type I and
399 Type II *ORF1* sequences, detected as non-recombinant, were colored in red and green,
400 respectively. GC clusters in the *ORF1* sequence were represented by squares and the
401 number represents the position in the *ORF1* alignment (see Figure 4). An arrow indicates
402 that a particular GC cluster was inverted to infer the GC cluster family. Square colors
403 represents GC cluster similarity according to the NJ trees from Figure S9 and S12.

Figure S6. *ORF1* aminoacid polymorphic sites.



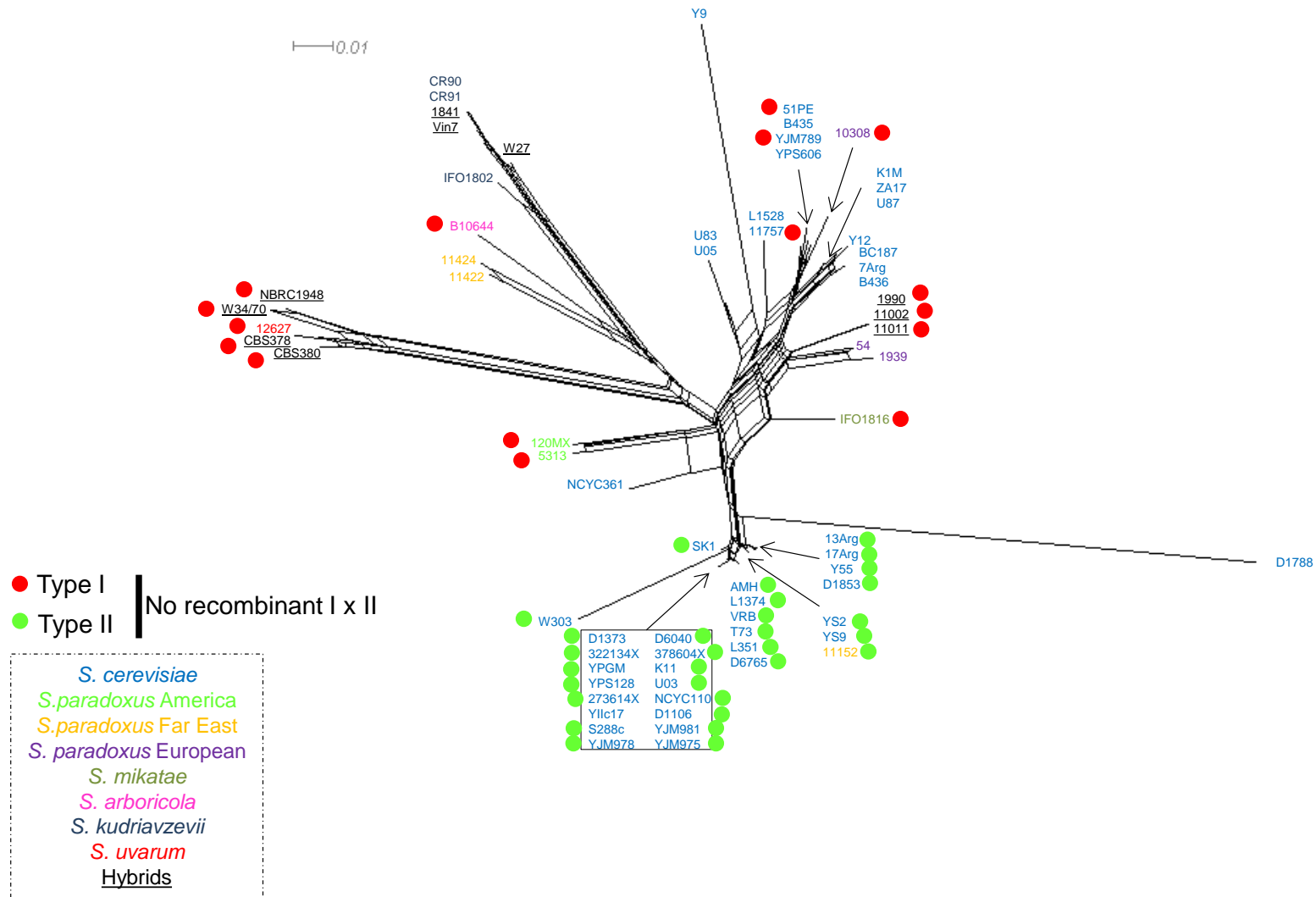
404

405 Variable *ORF1* aminoacid positions among haplotypes. Haplotype number are colored according to the species designation, as in Figure
 406 2. Alignment aminoacid positions for each polymorphic site are also shown. The COX2 Haplogroup designation is shown. Sequences
 407 were colored according to their similarity as inferred from Figure S7. Symbols * and ^ represent type I and type II *ORF1* sequences,
 408 respectively. Regions colored in yellow are sequences from unknown source. Lines indicate the sites corresponding to each alignment
 409 partition to reconstruct the *ORF1* phylogenetic networks by segments (Figure S7).

410 **Figure S7. ORF1 NN phylogenetic networks by segments.**

A

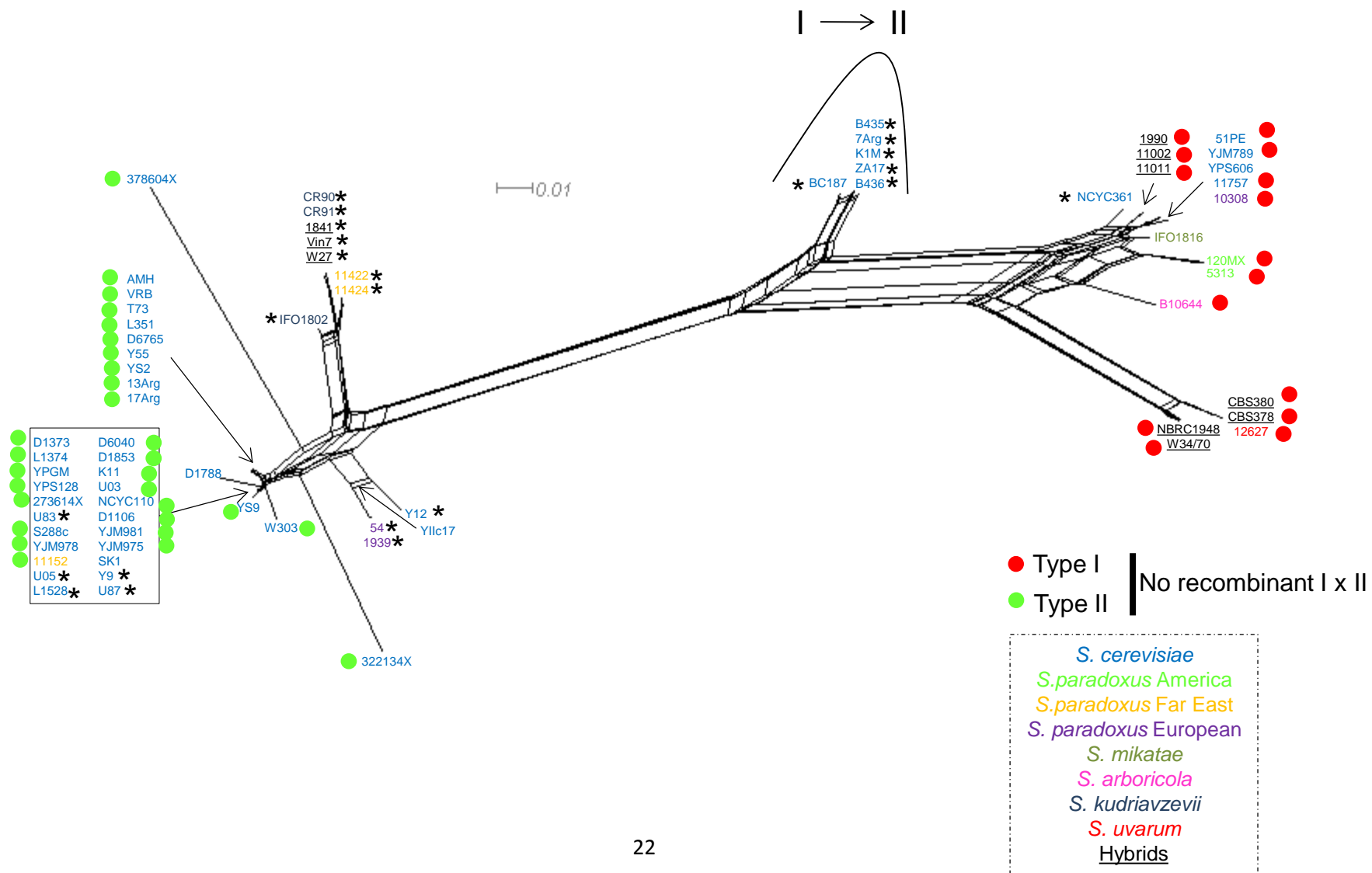
497-End-246*



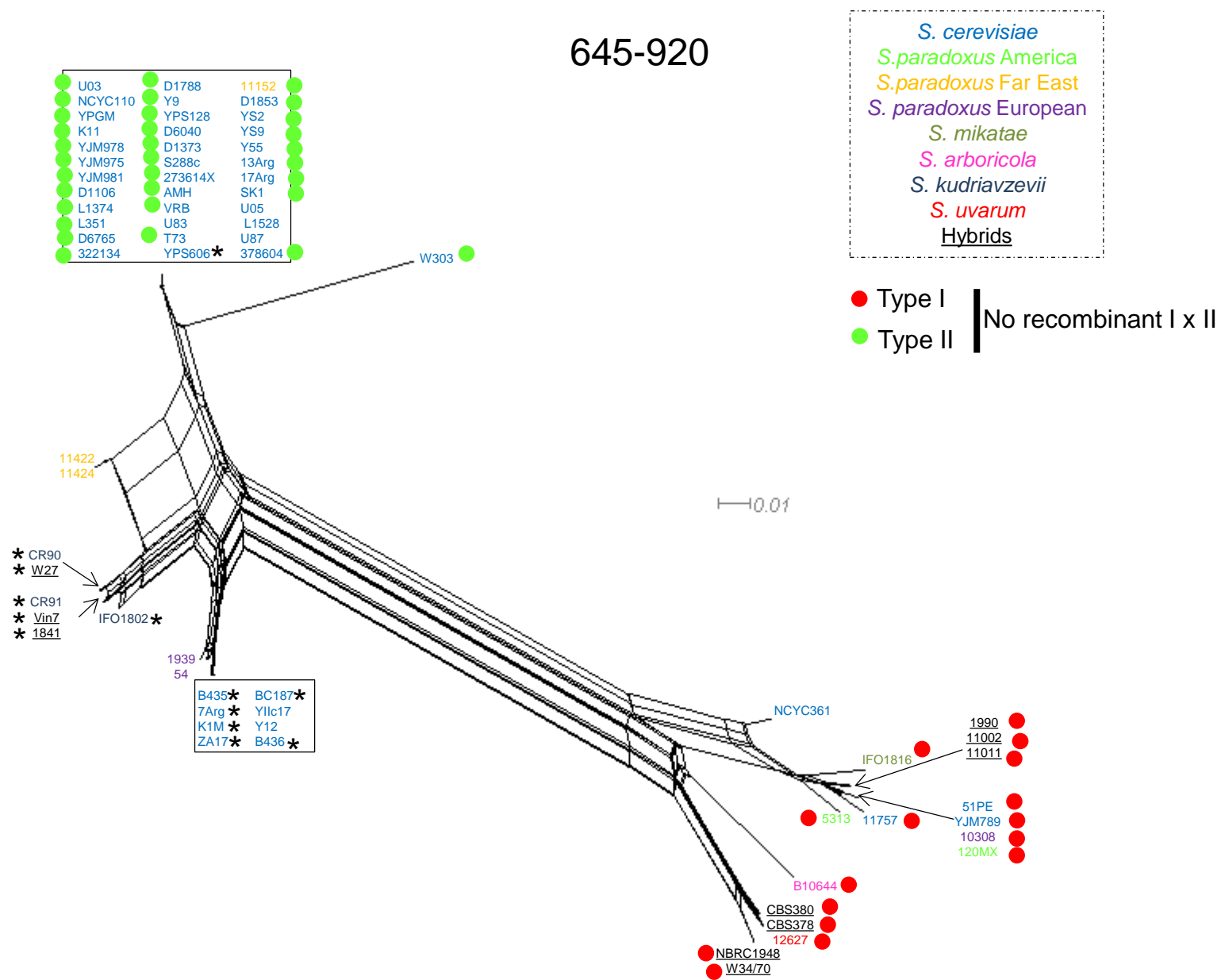
411

B

247-644

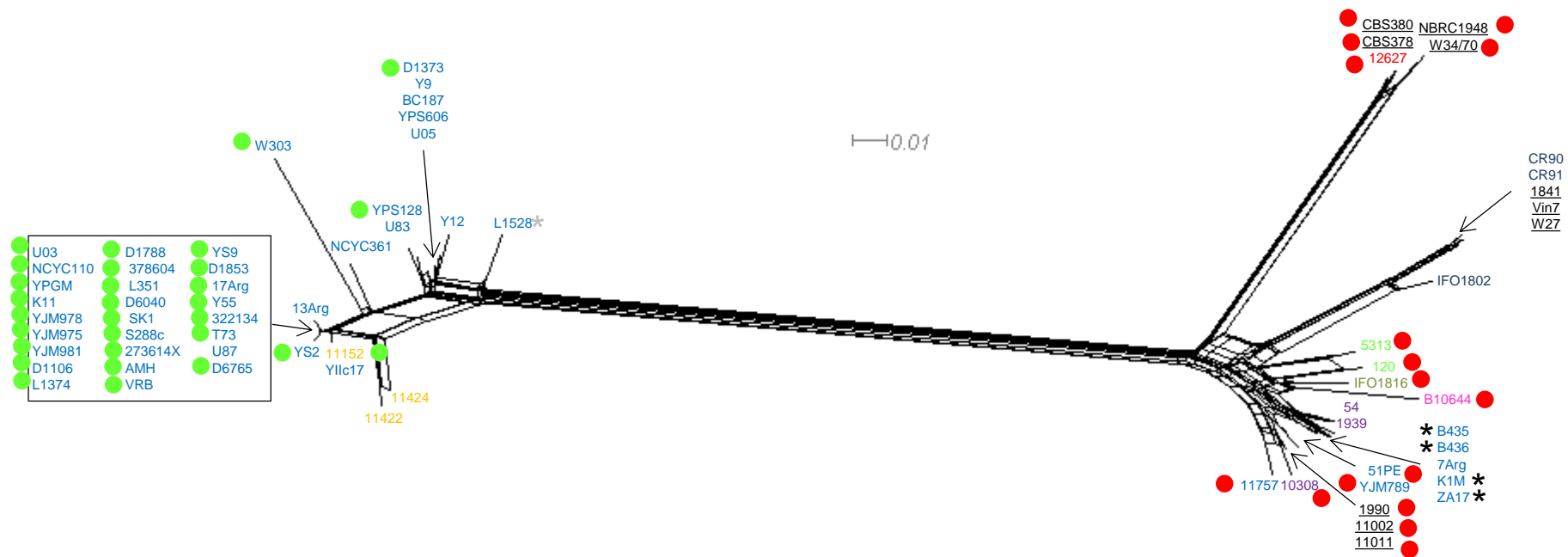


C



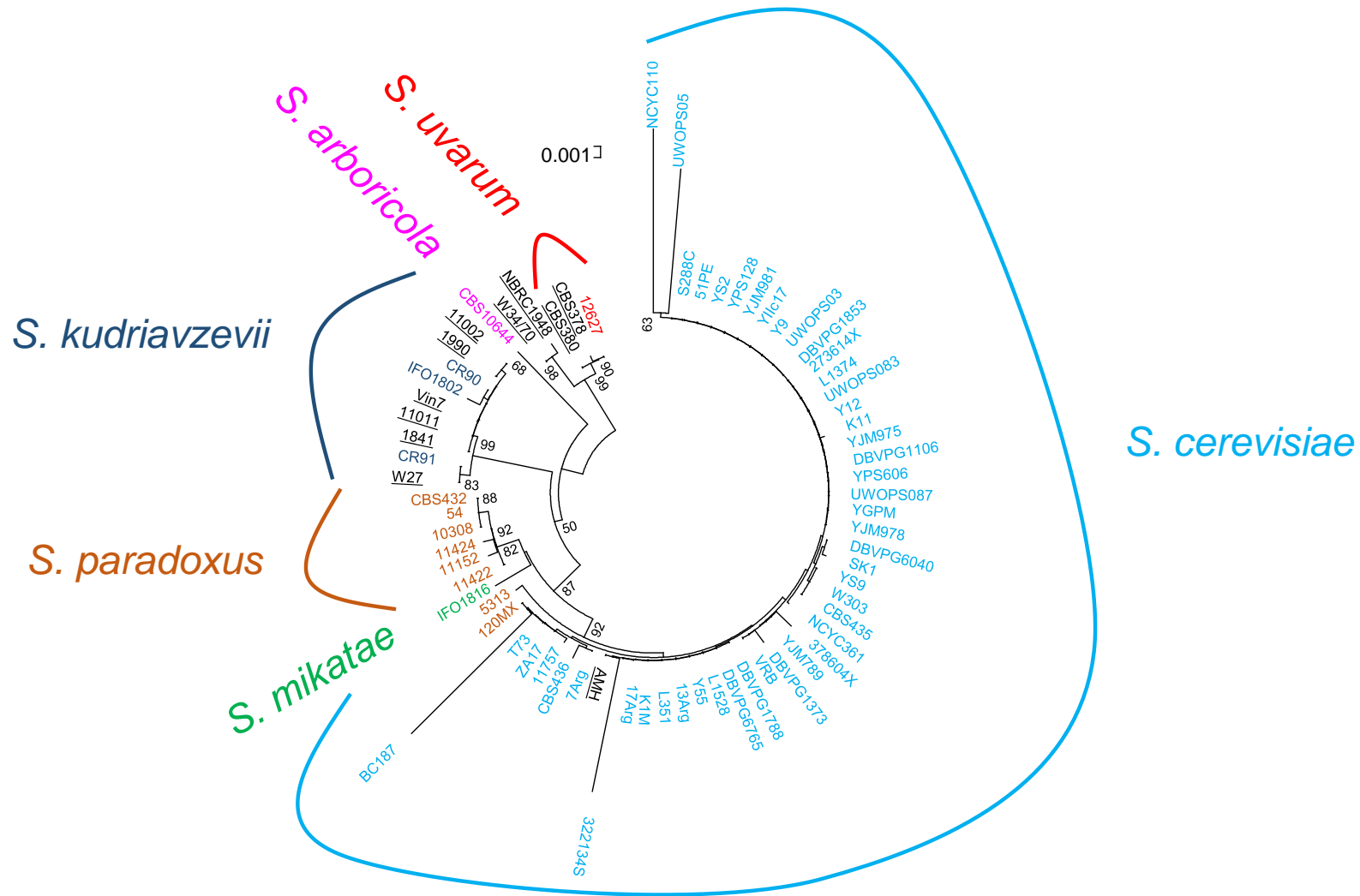
D

921-1251



415 NN phylogenetic networks. for each *ORF1* alignment partitions, inferred by GARD and RDP, are shown in A), B), C) and D). A) shows
416 the NN phylogenetic network corresponding to the *COX2* 497 nucleotide position until the position 246 of the *ORF1* alignment. Nucleotide
417 positions from *ORF1* alignment are shown in figures B-D. Scale bar are given in nucleotide substitution per site. Type I and Type II *ORF1*
418 sequences, detected as non-recombinant, are represented by red and green circles, respectively. Sequences are colored to each
419 species designation. In figure B) some *S. cerevisiae* sequences are grouped in a I -> II group, indicating that the region used for inferring
420 that NN phylonetwork contains a recombination point for those particular strains driving to the ambiguous position. Asterisks highlight
421 sequences which clustering changed from one *ORF1* type to another due to its recombinant character.

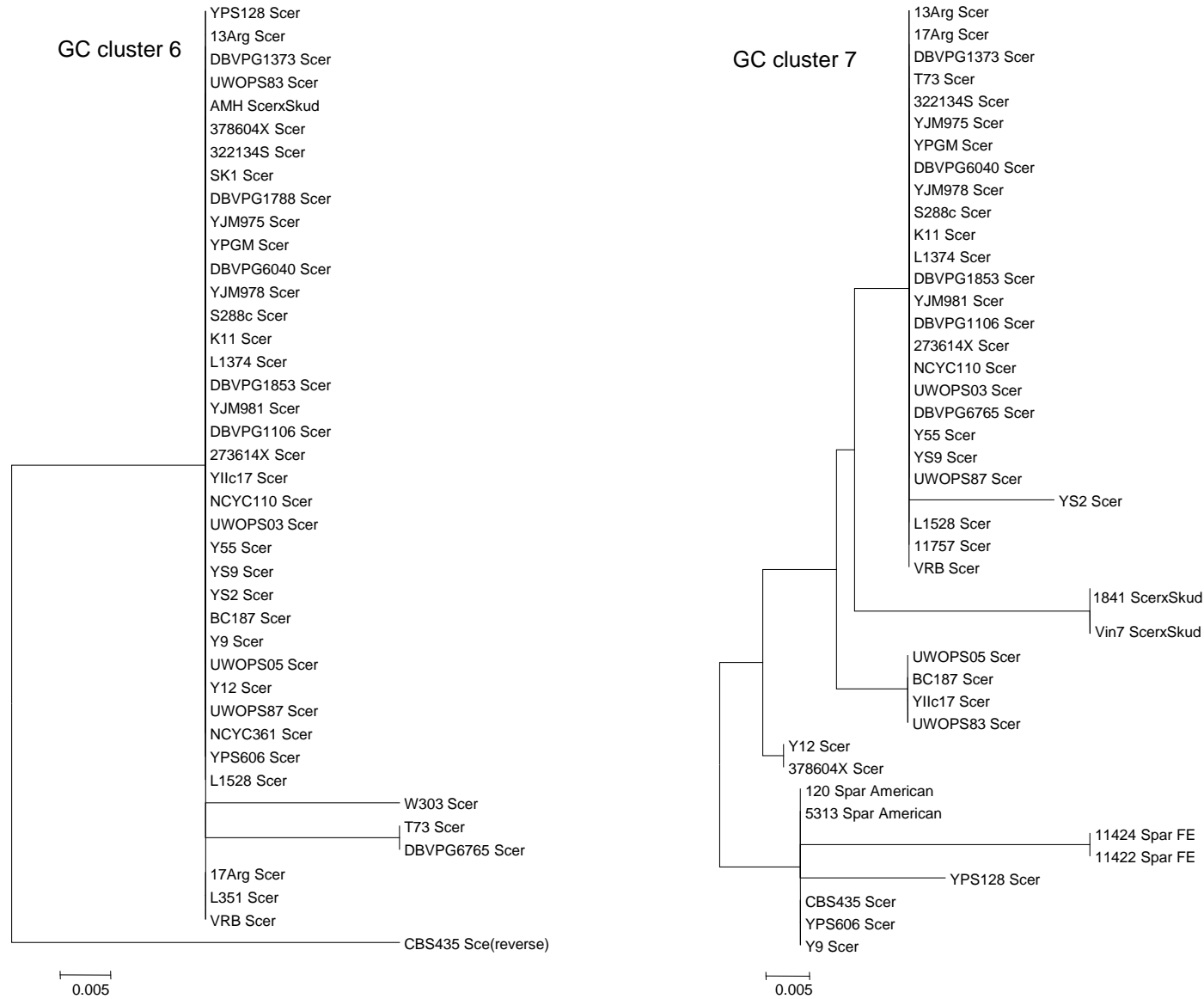
422 **Figure S8. COX3 NJ phylogenetic tree.**



423

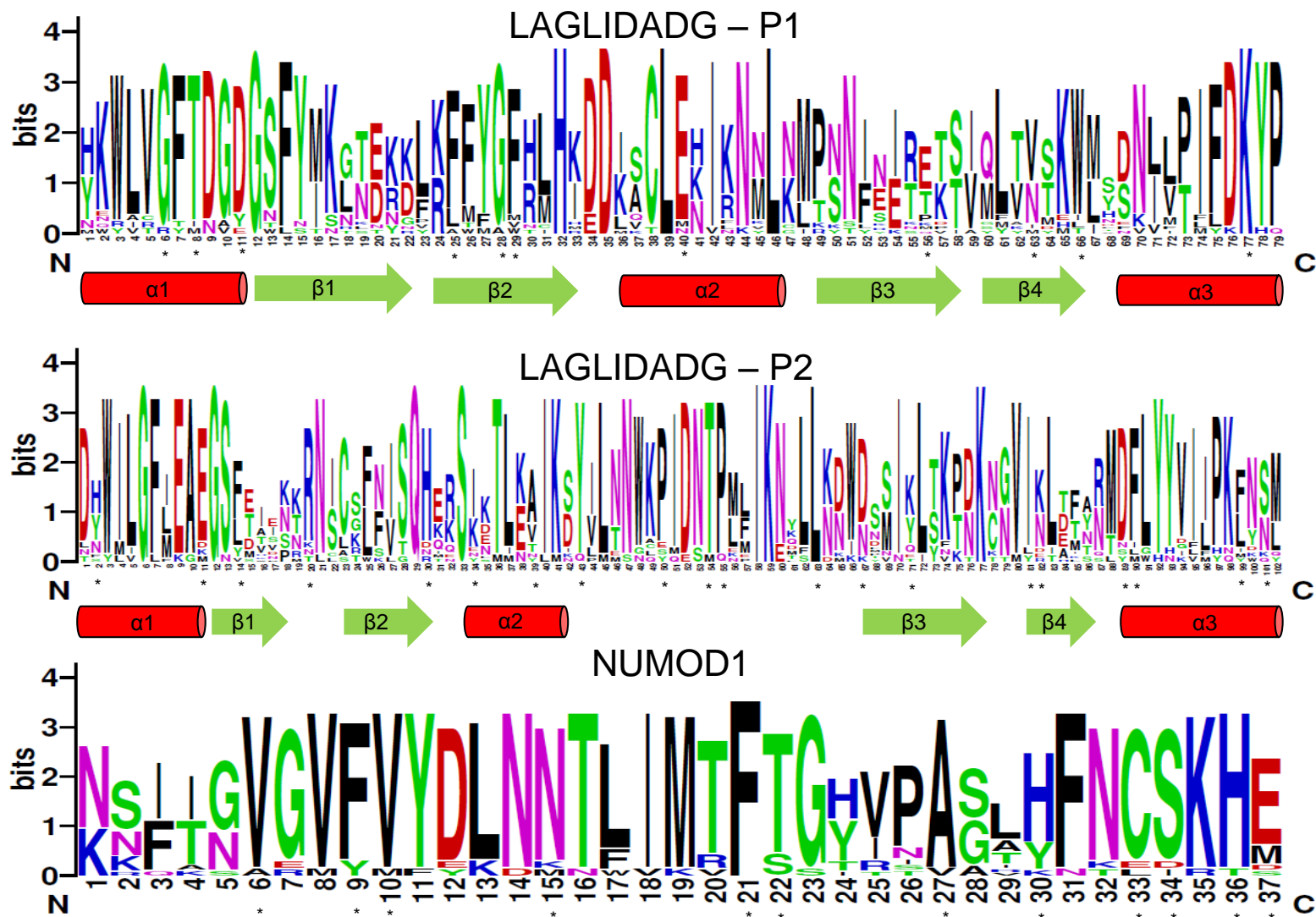
424 A COX3 Neighbor-Joining phylogenetic tree is shown. Strains are colored according to the species designation. Bootstrap values above
425 50 are given for each branch. Scale is given in nucleotide substitution per site.

426 **Figure S9. GC cluster 6 and 7 Neighbor-Joining phylogenetic trees.**



428 To classify the GC cluster according to sequence similarity we reconstructed the NJ phylogenetic tree of GC cluster 6 and 7. This
429 classification is shown by numbers in Figure S5. Scale bars represent number of substitutions per site.

430 Figure S10. Weblogo representation of LAGLIDADG and NUMOD domains.



431

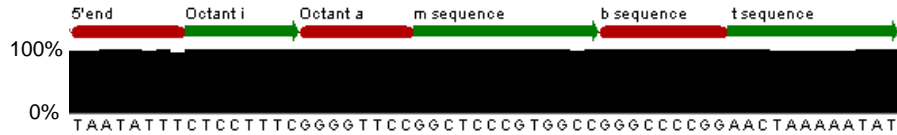
432 Weblogo (<http://weblogo.berkeley.edu/>) of the three homing endonuclease domains for *Saccharomyces ORF1* genes plus *SasefMp08*
433 from *Kazchastania servazii*, and *ORF1* and *ORF3* from *Williopsis saturnus* var. *suaveolens* are represented. Cylinders and arrows
434 represent α -helix and β -sheet, as previously described (Dalgaard et al. 1997). Asterisk symbol indicates the codons for those aminoacids
435 under purifying selection detected by DataMonkey.

436 **Figure S11. COX2 and ORF1 aminoacid alignment (Supplementary File S11).**

437 The complete *COX2* and *ORF1* aminoacid alignments are shown. Protein secondary
438 structure and domains are indicated for *ORF1*. Jalview also display the alignment quality
439 (based on BLOSUM 62), physicochemical conservation calculated according to
440 Livingstone and Barton (1993), and consensus sequence.

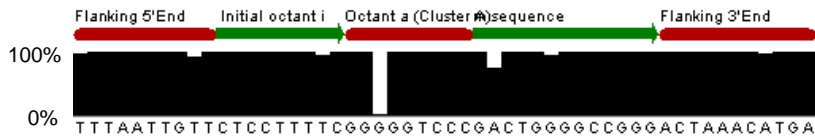
GC Cluster 6

ATTAATTATAAT	ATT-----	-----GAATTATAT
TATAATATT---	-----	-----GAATTATAT
--ATATAATTC-	-----	---AATATATAATT
TATAATATTTCT	-----	---AAAAATATTATT
TATAATATTCCT	-----	---AAAAATATTATT
-----CT	-----	AAATAAAATTAATT
TATAATAAT-AG	TTCCGGGGCCCGCCAC-GGAGCCGGAACCCGAAAGGAG	-AATTATAT----
-ATATAATT-CT	CCTTTCGGGGTTCCGGCTCC-GTGGCCGGGGCCCGGAACT	-ATTATTATA----*
-ATATAATT-CT	CCTTTCGGGGTTCCGGCTCC-GTGGCCGGGGCCCGGAACT	-ATTATTATAATT
TATAATATTCCT	CCTTTCGGGGTTCCGGCTCCCGTGGCCGGGGCCCGGAACT	-AAAAATATTATT
TATAATATTTCT	CCTTTCGGGGTTCCGGCTCCCGTGGCCGGGGCCCGGAACT	-AAAAATATTATT



GC Cluster 7

TAATTGTTCA	-----	AAACATATGATT
TAATTGTTCA	-----	AAACACATGATT
TAATTGCTCA	-----	AAACACATGATT
TAATTGTTCT	-----	AAACATGAAATC
TAATTGTTCT	-----	AAACATGAAATT
TAATTGTTCT	-----	AAACATGAGATT
TAATTGCTCT	-----	AAACATGAAATT
TAATTGCTCT	-----	AAACATGATATT
TAATTGCTCT	CCTTTTCGG-GGTCCCGACTGGGGCCGGGACT	AAACATGATATT
TAATTGTTCT	CCTTTTCGG-GGTCCCGACTGGGGCCGGGACT	AAACATGAAATT
TAATTGTTCT	CCTTTTCGG-GGTCCCGCTGGGGCCGGGACT	AAACATGAAATT
TAATTGTTCT	CCTTTTCGG-GGTCCCGACTGGGGCCGGGACT	AAACTGAAATT
TAATTGTTCT	CCTTTTCGG-GGTCCCGACTGGGGCCGGGACT	AAACATGAGATT
TAATTGTTCT	CCTTTTCGGGGTCCCGCTGGGGCCGGGACT	AAACATGAGATT
TAATTGTTCT	CCTTTTCGG-GGTCCCGCTGAGCCGGGACT	AAACATGAGATT
TAATTGTTCT	CCTTTTCGG-GGTCCCGCTGGGGCCGGGACT	AAACATGAGATT



443

444

445

446

447

The Flanking regions were the seven GC clusters were found are shown. In the case of GC cluster 6 and 7 a consensus sequence is shown and bars show the percentage of sequences showing a particular nucleotide substitution. An arrow indicates that the GC cluster sequence was inverted to annotate its structure. Question marks indicate an unknown GC cluster structure.