

Table of Contents

13		
14		
15		
16	COX2 nucleotide statistics	<u>2</u>
17		
18	Maximum-Likelihood phylogenetic tree of COX2 3'end	
19	region is incongruent with the species tree due to a	
20	recombination hotspot	<u>2-3</u>
21		
22	ORF1 gene structure	<u>3-5</u>
23		
24	ORF1 detection of selection	<u>5-6</u>
25		
26	References	<u>6-8</u>
27		
28	Supplementary Tables and Figures	<u>8-35</u>
29		
30		
31		

32 **Supplementary Text**

33 **COX2 nucleotide statistics**

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

43 COX2 genetic diversity statistics are summarized in Table S3. The most diverse
44 sequences were found in three species, *S. cerevisiae*, *S. paradoxus* and *S. mikatae*. In
45 the case of *S. cerevisiae* nucleotide diversity values of 2.1% are so much higher than
46 values inferred for the nuclear genome, around 0.9% (Liti *et al.* 2009), suggesting that
47 recombination is increasing the values. For *S. eubayanus* and *S. uvarum*, the COX2
48 nucleotide diversity values, 0.5 and 0.6% respectively, are lower than those found for the
49 nuclear genome, around 0.8% (Patagonia A-Patagonia B *S. eubayanus*) and 1%
50 (Holarctic-South America B *S. uvarum*) (Almeida *et al.* 2014; Peris *et al.* 2014) in
51 agreement with the low mutation rate of yeast mitochondrial genomes (Clark-Walker
52 1991). COX2 values for *S. kudriavzevii* strains are similar to those inferred for the nuclear
53 genome, 1.11% (Hittinger *et al.* 2010).

54

55 **Maximum-Likelihood phylogenetic tree of COX2 3'end region is incongruent with** 56 **the species tree due to a recombination hotspot**

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

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

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

74 Codons (nucleotide sequences 526-528 and 535-537) in the COX2 3' end show
75 several convergent nucleotide substitutions among *S. cerevisiae* and *S. paradoxus*
76 strains. This finding might be the result of a “patchy-tachy” effect due to different
77 substitution rate in those positions (Sun *et al.* 2011). Homoplasies can drive to
78 recombinant false positives. One confounding effect can be observed in L1528 and
79 CECT11757 strains which is difficult to know if they were truly introgressed or they
80 suffered two different recombination events with two different *S. paradoxus* lineages, Far
81 Eastern and European. A recent study shown transfers from European *S. paradoxus*
82 GC42 cluster to *S. cerevisiae* L1528 (Wu and Hao 2015). Also the presence of type II
83 ORF1 segments supports the introgression scenario and a double recombination
84 scenario.

85

86 **ORF1 gene structure**

87 52 different haplotypes were found in our representative strains (Table S1). ORF1
88 start codon is nineteen nucleotides inside the 3' end of COX2 gene. The average GC-
89 composition of the ORF1 is extremely low, 18%. Eleven strains have a GTG as a start
90 codon, which is uncommon in yeast mitochondria: haplotypes M2-M7 and M9-M12. The
91 translation of the ORF1 gene predicts fourteen strains having a premature stop codon:
92 haplotypes M2, M3, M7-M11, M17-M20, M23 and M46. ORF1 sequence length range
93 from 1363bp (*S. cerevisiae* ZA17 strain) to 1516bp (*S. cerevisiae* VRB strain), translated

94 to 454 and 505, respectively. Note that we sequenced a partial *ORF1* gene, and the last
95 45 nucleotides were not sequenced in this study.

96 *ORF1* secondary structure and domains (LAGLIDADG and NUMOD) were annotated
97 in Jalview according to Dalgaard *et al.* (Dalgaard *et al.* 1997) and using the Conserved
98 Domains tool in NCBI (<http://www.ncbi.nlm.nih.gov/Structure/cdd/wrpsb.cgi>) (Marchler-
99 Bauer *et al.* 2009). WebLogo profiles of LAGLIDADG and NUMOD1 domains were done
100 in WebLogo 2.8.2 tool (<http://weblogo.berkeley.edu/>). Three different domains were
101 found in all *ORF1* sequences: two LADGLIDADG (P1 and P2) and one NUMOD1 (Figure
102 6, S9 and S10). Two non-*Saccharomyces ORF1* were detected by PSI-BLAST (Altschul
103 *et al.* 1997) showing some similarities with the *Saccharomyces ORF1* sequences. We
104 generated a new alignment in Jalview 4.0.b2 (Waterhouse *et al.* 2009) including the
105 *Cyberlindnera (Williopsis) saturnus* var. *suaveolans* (*ORF1* and *ORF3*), and one from
106 *Kazachstania servazii* (SasefMp08) sequences (Figure S10). The sequence from *K.*
107 *servazii* was found to be closely related to the *Saccharomyces ORF1* sequences (Figure
108 S3). Although, *C. saturnus* is phylogenetically fairly distant from *Saccharomyces* genus
109 (James SA *et al.* 1998), they have been diverging since 235mya, we observed some
110 conservation among LAGLIDADG and NUMOD aminoacids when compared with
111 *Saccharomyces* and non-*Saccharomyces ORF1* sequences (Figure S9 and S10).

112 Differences in size among *ORF1* gene sequences were due to the presence of GC
113 insertions (Figure S3) and AT repeats (Table S4). Seven different GC clusters insertion
114 points were found along the *ORF1* alignment. GC cluster 1 is observed in *S. cerevisiae*
115 VRB; three types of GC cluster 2 were found in *S. cerevisiae* CBS435, CECT11757, and
116 *S. paradoxus* 120MX (Figure S3 and S11); GC cluster 3 is shown in *S. arboricola*
117 CBS10644; GC cluster 4 is shown in *S. cerevisiae* VRB; GC cluster 5 is inserted in *ORF1*
118 of *S. cerevisiae* CBS435, and was structurally similar to GC cluster 3 of *S. arboricola*
119 CBS10644; GC cluster 6 and 7 were the most frequent among *Saccharomyces* strains
120 (Figure S3 and S12). 40 *S. cerevisiae* and one hybrid *S. cerevisiae* x *S. kudriavzevii*
121 (AMH) have the GC cluster 6, and 36 *S. cerevisiae*, 2 hybrids *S. cerevisiae* x *S.*
122 *kudriavzevii*, 2 Far Eastern and 2 American *S. paradoxus* shown the GC cluster 7 in their
123 *ORF1* sequences (Figure S3 and S11). Some GC clusters sequences were found
124 inverted (Figure S3 and S11). GC cluster 7 of the two American *S. paradoxus* (120MX

125 and CBS5313) and three *S. cerevisiae*, two from Japan and one from USA (CBS435, Y9
126 and YPS606), contained identical GC cluster sequence (Figure S12). The confirmation of
127 the American *S. paradoxus* mitochondrial background will shed light about the origin of
128 this GC cluster, which simplest scenario is a transfer from the *S. cerevisiae* mitochondrial
129 genome to *S. paradoxus ORF1*.

130 *ORF1* GC clusters were classified according to Zamaroczy and Bernardi (de
131 Zamaroczy and Bernardi 1986), except for GC cluster 2. GC cluster 1 and 4 are similar
132 to a1 family, and GC cluster 3 and 5 were similar to a4 family. In the case of GC cluster
133 1, 2, 4, and 5 are on the opposite strand. As Séraphin *et al.* (1987) and Weiller *et al.*
134 (1989) described, the GC clusters of the *ORF1* gene were flanked by TAG and AGGAG,
135 or CTA and CTCCT when the cluster was inserted in the opposite strand (Figure S11).
136 These conserved nucleotides were flanked by A+T rich sequences. Flanking sequences
137 TAG and AG (CTA and CT) are conserved in most of the sequences with and without GC
138 clusters. All GC clusters in *ORF1* belong to group M1 (Weiller *et al.* 1989). Conserved
139 flanking regions might help to the transposition of GC clusters.

140 Differences in size were also due to the presence of AT repeats. Tandem repeat
141 sequences for *COX2* and *ORF1* genes were detected and described using `Tandem`
142 `Repeat Finder` software (Benson 1999). Twenty-one different A+T rich sequences that
143 repeated at least twice were found in the *ORF1* alignment (Table S4). The length of A+T
144 rich tandem repeats ranged from three nucleotides to twenty-five nucleotides. The most
145 repeated sequence (AAT) was repeated ten times in haplotypes M1, M12, M22, M39,
146 M40, and M50. The A+T rich tandem repeats were located near to GC clusters (Figure
147 6). In the *COX2* gene, we found three A+T rich sequences repeated twice (Table S4).

148

149 ***ORF1* detection of selection**

150 The single likelihood ancestor counting (SLAC), fixed effects likelihood (FEL), and
151 random effects likelihood (REL) methods (Pond and Frost 2005), implemented in the
152 `HYPHY` web based version, `Datamonkey` (Delport *et al.* 2010), were used to detect the
153 signatures of selection operating on the *ORF1* gene. NJ phylogenetic trees were
154 reconstructed under the REV (GTR) substitution model. Two different approaches were
155 performed to describe selection signatures. In the first approach, a complete *ORF1*

156 alignment, without GC Insertions and indels, partitioned by recombinant sections detected
157 by GARD, was used. In a second study, the two LADGLIDADG and the NUMOD domains
158 were analyzed independently.

159 Aminoacid positions were described as being under positive or purifying selection
160 when significant values were generated by two of the three tested methods. Codon-
161 specific selection pressure along the sequences (i.e. site specific dN-dS) was measured
162 and *p*-values were estimated at each site. We analyzed 417 *ORF1* codons of the total
163 458. 61 codons, corresponding to 15% of the total, were found to be under purifying
164 selection, where 43 of the 61 codons were found in LAGLIDADGs and NUMOD1 domains
165 (Figure S9). A functional characterization might help to understand how important these
166 aminoacids are for the activity of this homing endonuclease.

167

168 **References**

169 Almeida P, Gonçalves C, Teixeira S *et al.* 2014. A Gondwanan imprint on global diversity
170 and domestication of wine and cider yeast *Saccharomyces uvarum*. *Nat Commun.* 5.

171 Altschul SF, Madden TL, Schaffer AA, Zhang J, Zhang Z, Miller W, Lipman DJ. 1997.
172 Gapped BLAST and PSI-BLAST: a new generation of protein database search programs.
173 *Nucl Acids Res.* 25:3389-3402.

174 Benson G. 1999. Tandem repeats finder: a program to analyze DNA sequences. *Nucl*
175 *Acids Res.* 27:573-580.

176 Clark-Walker GD. 1991. Contrasting mutation rates in mitochondrial and nuclear genes
177 of yeasts versus mammals. *Curr Genet.* 20:195-198.

178 Dalgaard JZ, Klar AJ, Moser MJ, Holley WR, Chatterjee A, Mian IS. 1997. Statistical
179 modeling and analysis of the LAGLIDADG family of site- specific endonucleases and
180 identification of an intein that encodes a site-specific endonuclease of the HNH family.
181 *Nucl Acids Res.* 25:4626-4638.

182 de Zamaroczy M, Bernardi G. 1986. The GC clusters of the mitochondrial genome of
183 yeast and their evolutionary origin. *Gene.* 41:1-22.

184 Delpont W, Poon AFY, Frost SDW, Kosakovsky Pond SL. 2010. Datamonkey 2010: a
185 suite of phylogenetic analysis tools for evolutionary biology. *Bioinformatics.* 26:2455-
186 2457.

187 Gibbs MJ, Armstrong JS, Gibbs AJ. 2000. Sister-Scanning: a Monte Carlo procedure for
188 assessing signals in recombinant sequences. *Bioinformatics*. 16:573-582.

189 Guindon S, Gascuel O. 2003. A simple, fast, and accurate algorithm to estimate large
190 phylogenies by Maximum Likelihood. *Syst Biol*. 52:696-704.

191 Hittinger CT, Gonçalves P, Sampaio JP, Dover J, Johnston M, Rokas A. 2010.
192 Remarkably ancient balanced polymorphisms in a multi-locus gene network. *Nature*.
193 464:54-58.

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

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

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

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

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

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

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

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

210 Peris D, Sylvester K, Libkind D, Gonçalves P, Sampaio JP, Alexander WG, Hittinger CT.
211 2014. Population structure and reticulate evolution of *Saccharomyces eubayanus* and its
212 lager-brewing hybrids. *Mol Ecol*. 23:2031-2045.

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

215 Posada D. 2008. jModelTest: phylogenetic model averaging. *Mol Biol Evol*. 25:1253-
216 1256.

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

220 Schmidt HA, Strimmer K, Vingron M, von Haeseler A. 2002. TREE-PUZZLE: maximum
221 likelihood phylogenetic analysis using quartets and parallel computing. *Bioinformatics*.
222 18:502-504.

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

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

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

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

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

232 Wang QM, Liu WQ, Liti G, Wang SA, Bai FY. 2012. Surprisingly diverged populations of
233 *Saccharomyces cerevisiae* in natural environments remote from human activity. *Mol Ecol*.
234 21:5404-5417.

235 Waterhouse AM, Procter JB, Martin DMA, Clamp M, Barton GJ. 2009. Jalview Version 2-
236 -a multiple sequence alignment editor and analysis workbench. *Bioinformatics*. 25:1189-
237 1191.

238 Weiller G, Schueller CME, Schweyen RJ. 1989. Putative target sites for mobile G+C rich
239 clusters in yeast mitochondrial DNA: Single elements and tandem arrays. *Mol Gen Genet*.
240 218:272-283.

241 Wu B, Hao W. 2015. A dynamic mobile DNA family in the yeast mitochondrial genome.
242 *G3: Genes|Genomes|Genetics*. 5:1273-1282.

243

244 **Supplementary Tables and Figures**

245 **Table S1.** Strain information.

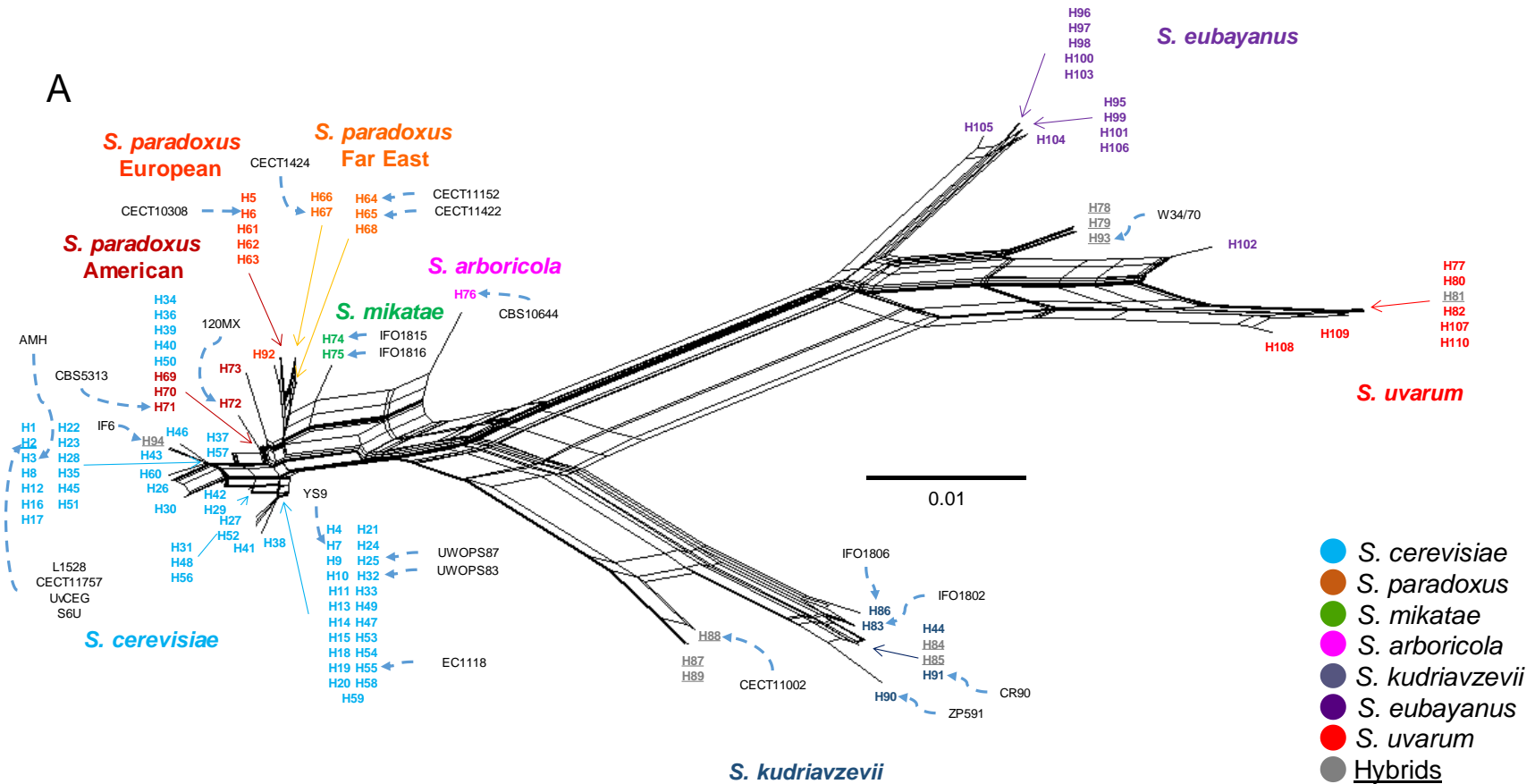
246 **Table S2.** Primer pairs used in this study.

247 **Table S3.** Summary statistics for all *Saccharomyces* and each species using COX2.

248 **Table S4.** Distribution of AT tandem repeats.

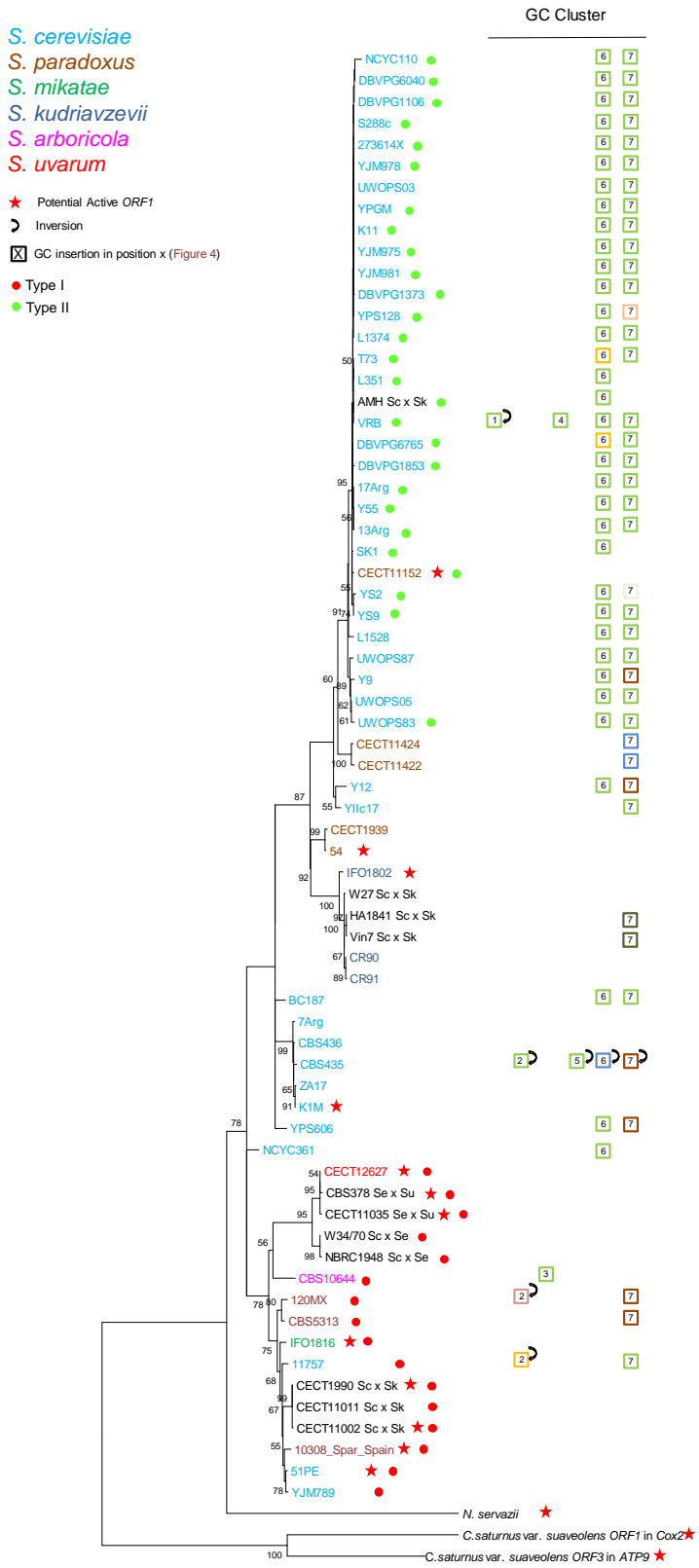
249

263 **Figure S2. COX2 Neighbor-Net phylogenetic network for each COX2 segment.**



264

270 **Figure S3. ORF1 Neighbor-Joining phylogenetic tree.**



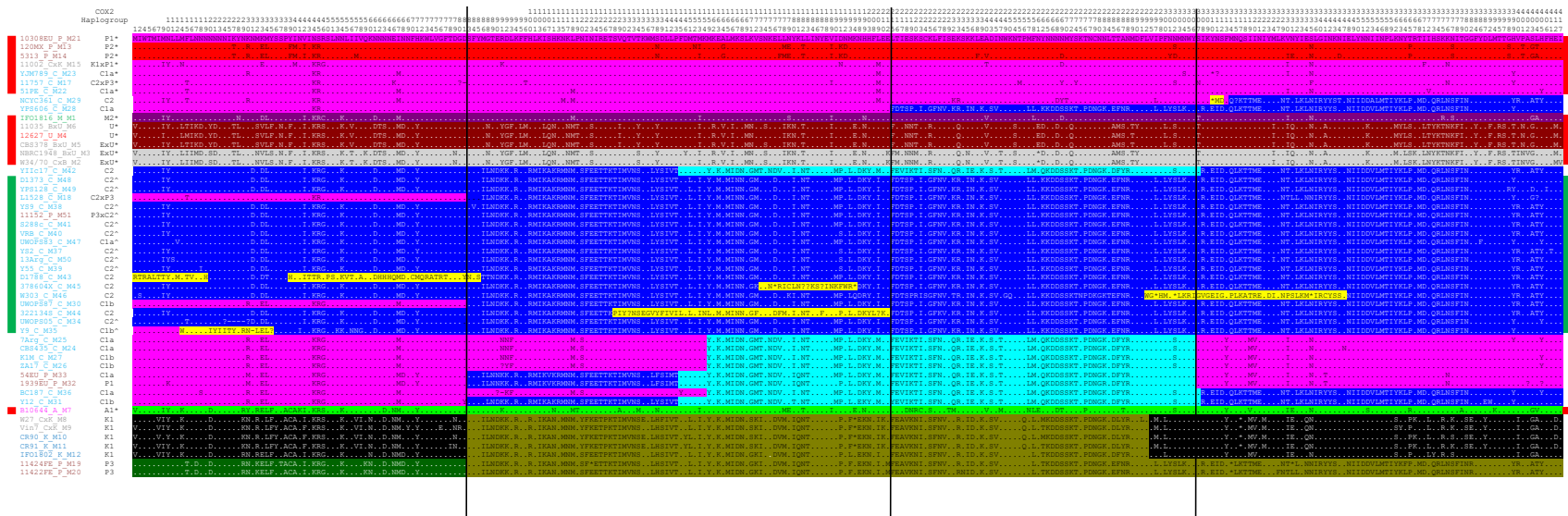
271

0.02

272 A NJ phylogenetic tree was reconstructed for the *ORF1* sequence alignment. Strains
273 were colored according to the species designation. Bootstrap values above 50 are given
274 for each branch. A red star represents a potential active *ORF1* based on the absence of
275 GC clusters, which could disrupt the coding sequence, and premature stop codons. Type
276 I and Type II *ORF1* sequences, detected as non-recombinant (represented by dots), were
277 colored in red and green, respectively. GC clusters in the *ORF1* sequence were
278 represented by squares and the number represents the position in the *ORF1* alignment
279 (see Figure 6). An arrow indicates that a particular GC cluster was inverted to infer the
280 GC cluster family. Square colors represent GC cluster similarity according to Figure S11,
281 and the NJ trees from Figure S12.

282

283 **Figure S4. *ORF1* aminoacid polymorphic sites.**



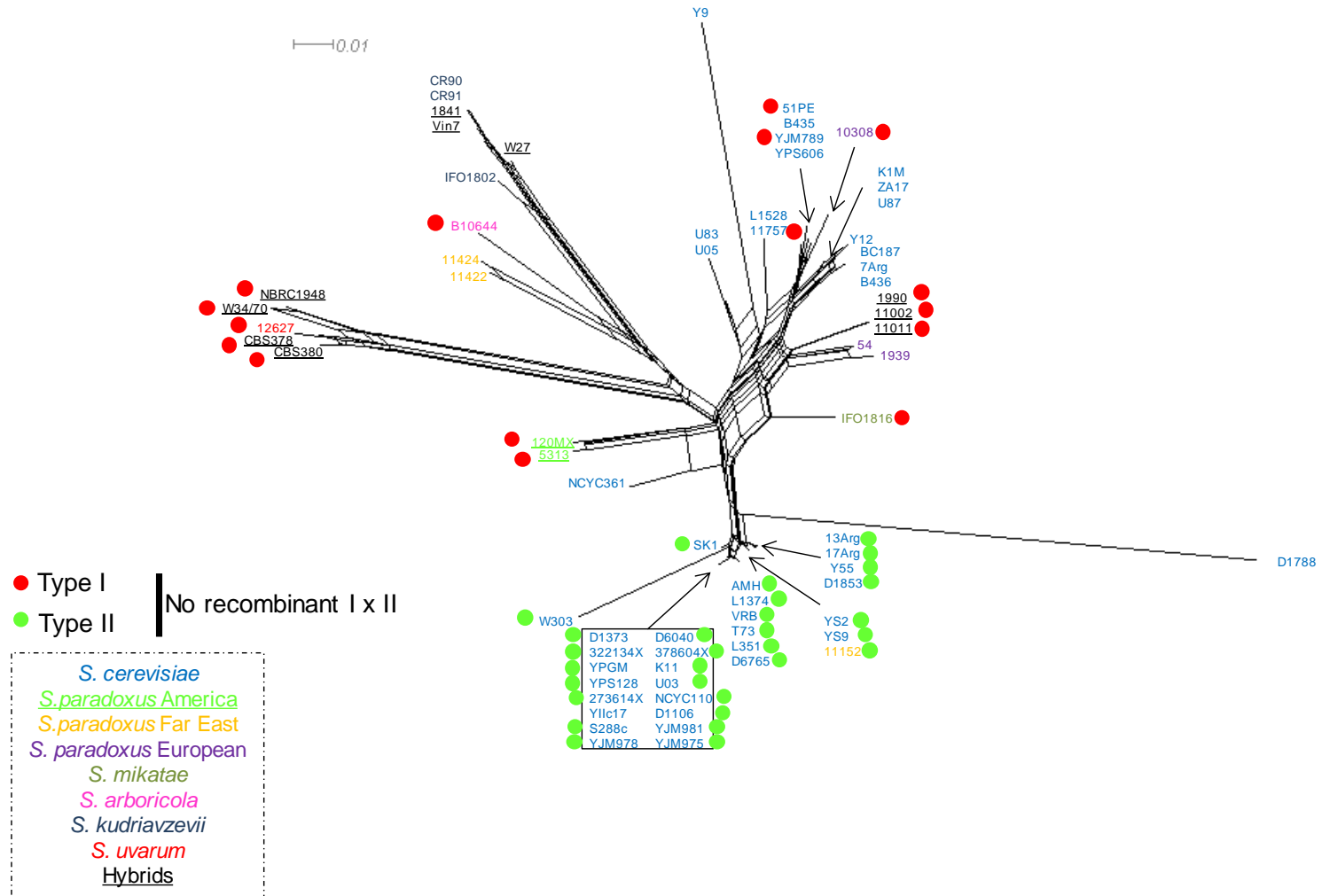
■ Type I
■ Type II

284
285 Variable *ORF1* aminoacid positions among haplotypes. Haplotype number are colored according to the species designation,
286 as in Figure 2. Alignment aminoacid positions for each polymorphic site are also shown. The *COX2* Haplogroup designation
287 is shown. Sequences were colored according to their similarity as inferred from Figure S5. Symbols * and ^ represent type
288 I and type II *ORF1* sequences, respectively. Regions colored in yellow are sequences from unknown source. Lines indicate
289 the sites corresponding to each alignment partition to reconstruct the *ORF1* phylogenetic networks by segments (Figure
290 S5).

291 **Figure S5. ORF1 NN phylogenetic networks by segments.**

A

COX2 3'end-246*

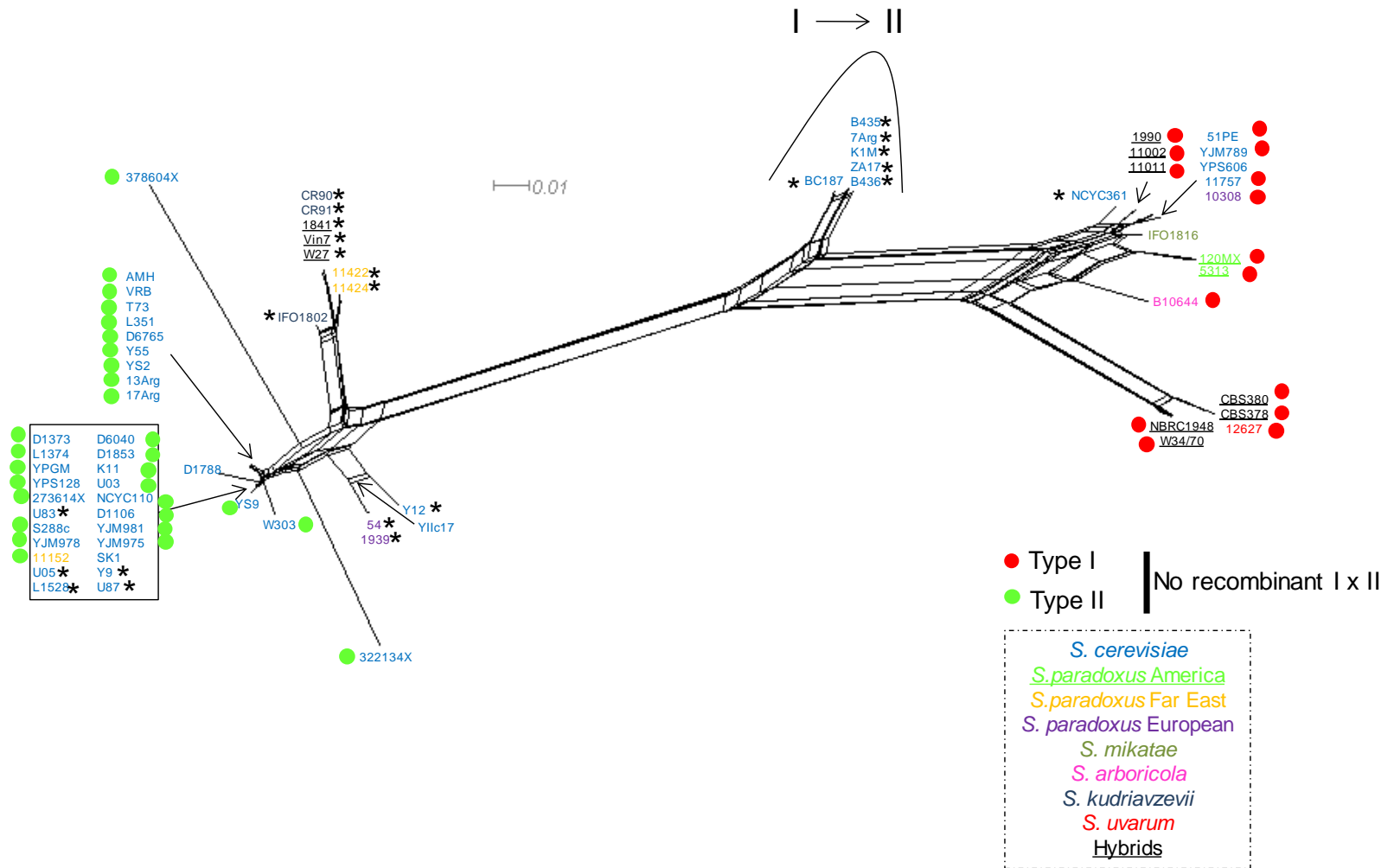


292

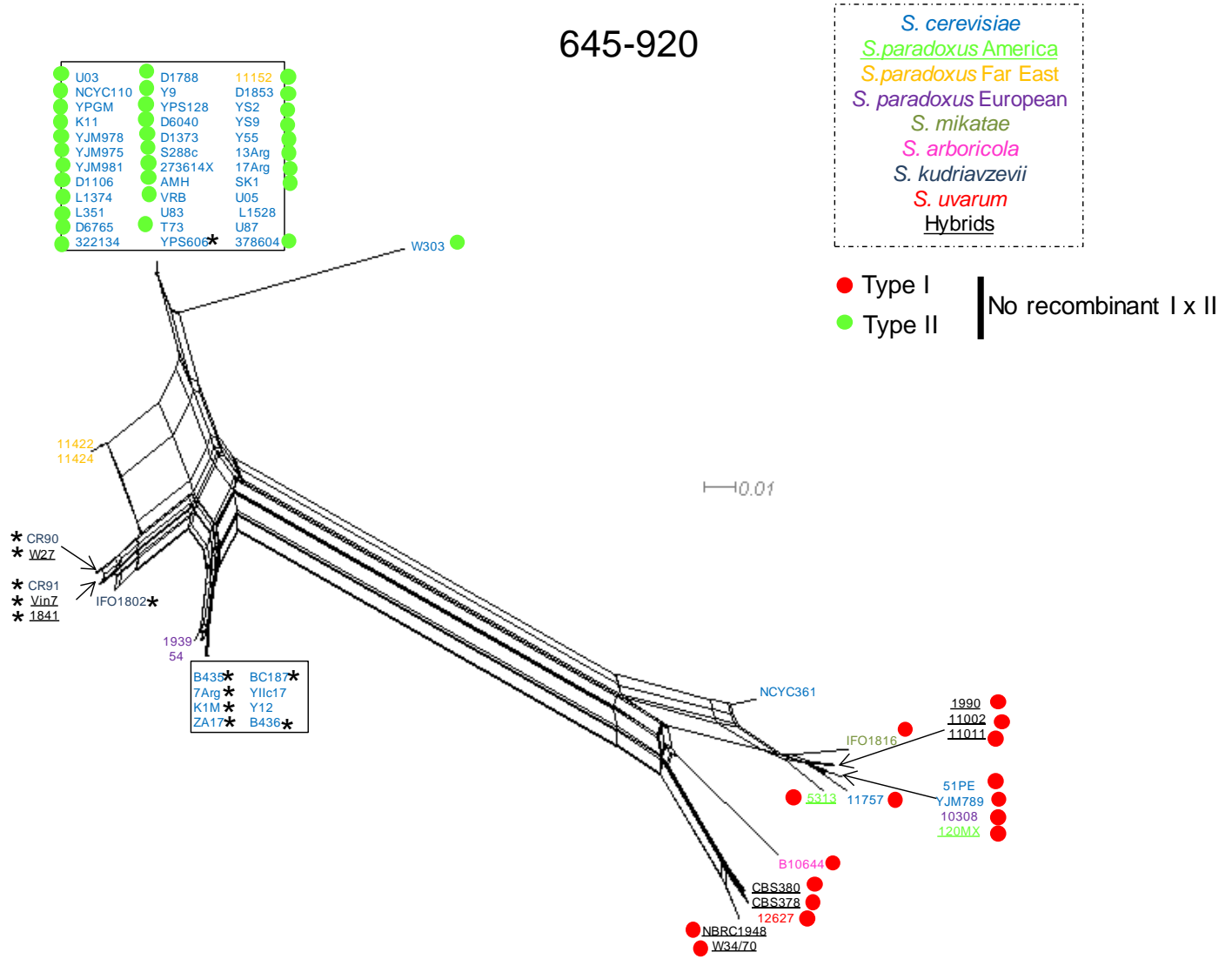
293

B

247-644



C



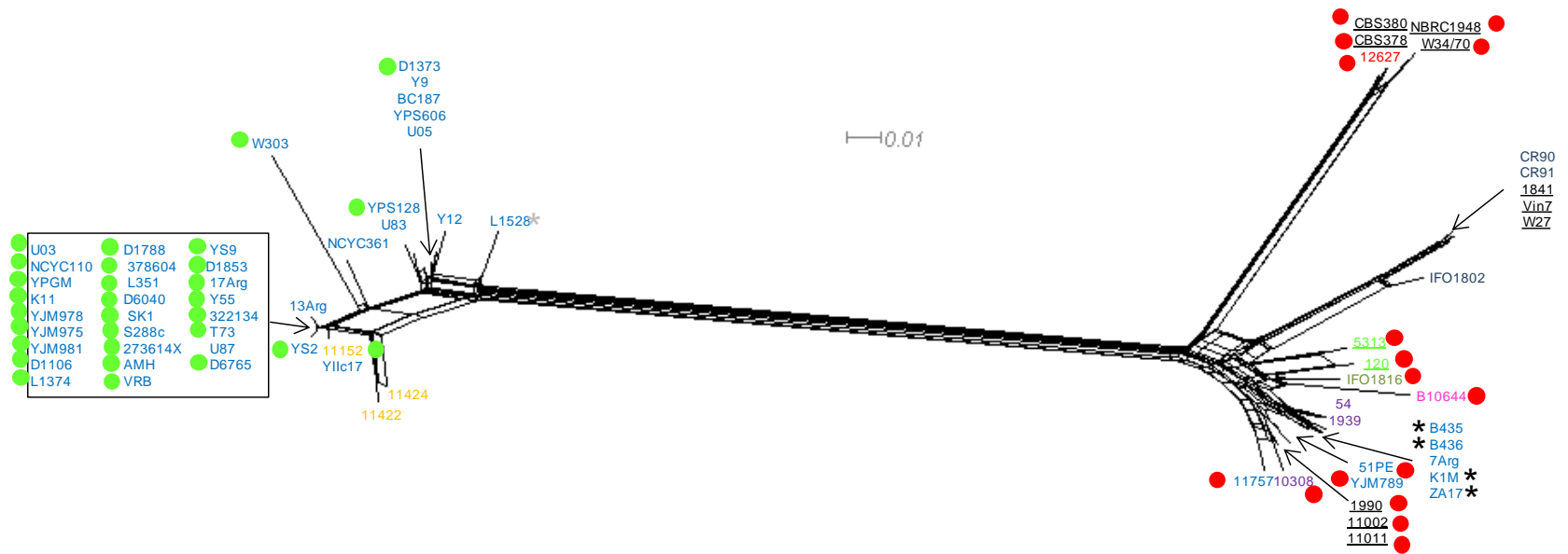
295

296

297

D

921-1251



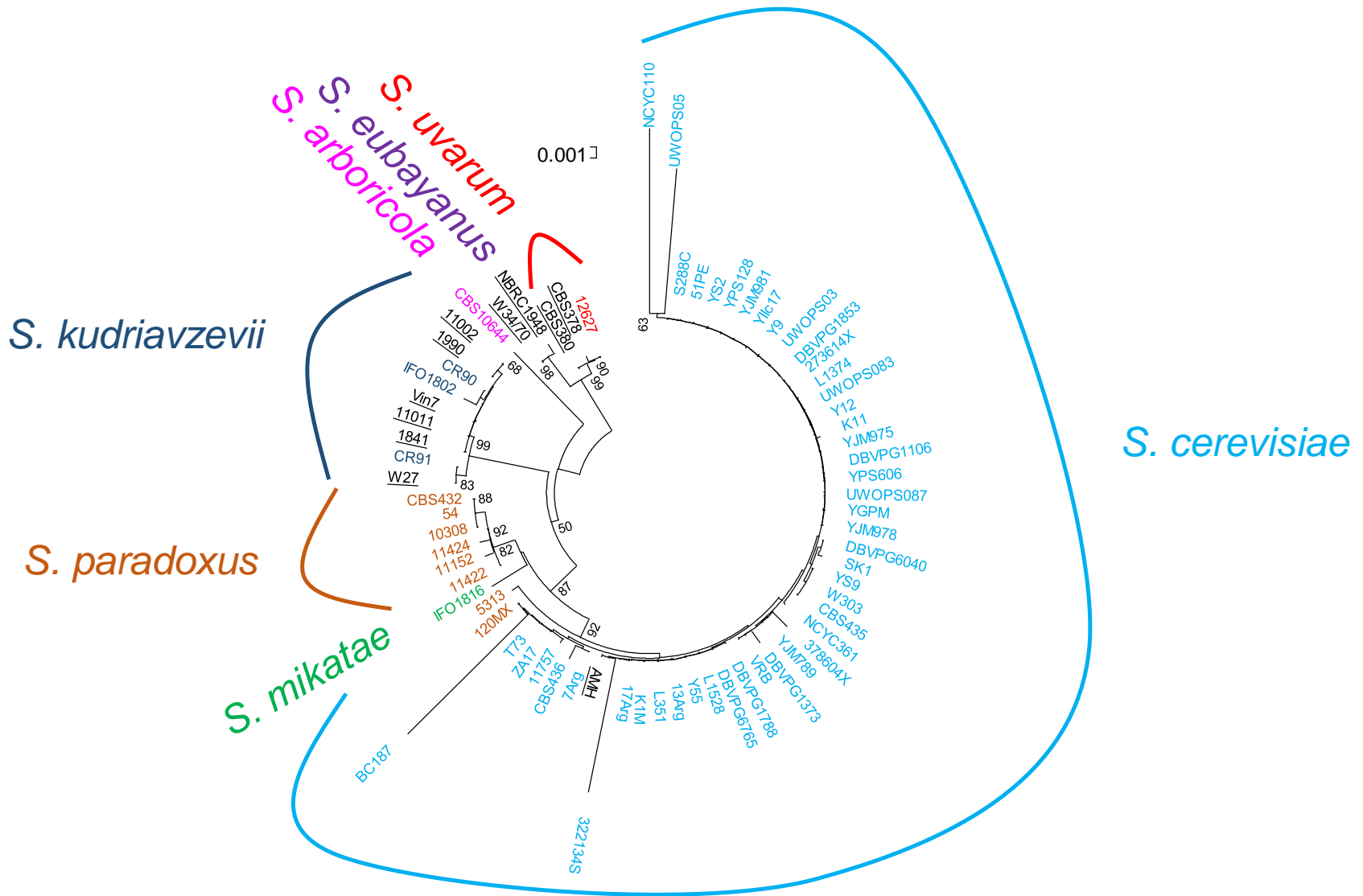
● Type I
 ● Type II

No recombinant I x II

- S. cerevisiae*
- S. paradoxus* America
- S. paradoxus* Far East
- S. paradoxus* European
- S. mikatae*
- S. arboricola*
- S. kudriavzevii*
- S. uvarum*
- Hybrids

299 NN phylogenetic networks. for each *ORF1* alignment partitions, inferred by *GARD* and *RDP*, are shown in A), B), C) and D).
300 A) shows the NN phylogenetic network corresponding to the region from *COX2* 3' end (Figure S1) to the position 246 of the
301 *ORF1* alignment. Nucleotide positions from *ORF1* alignment are shown in figures B-D. Scale bar are given in nucleotide
302 substitution per site. Type I and Type II *ORF1* sequences, detected as non-recombinant, are represented by red and green
303 circles, respectively. Sequences are colored to each species designation. In figure B) some *S. cerevisiae* sequences are
304 grouped in a I -> II group, indicating that the region used for inferring that NN phylonetwork contains a recombination point
305 for those particular strains driving to the ambiguous position. Asterisks highlight sequences which clustering changed from
306 one *ORF1* type to another due to its recombinant character.

307 **Figure S6. COX3 NJ phylogenetic tree.**

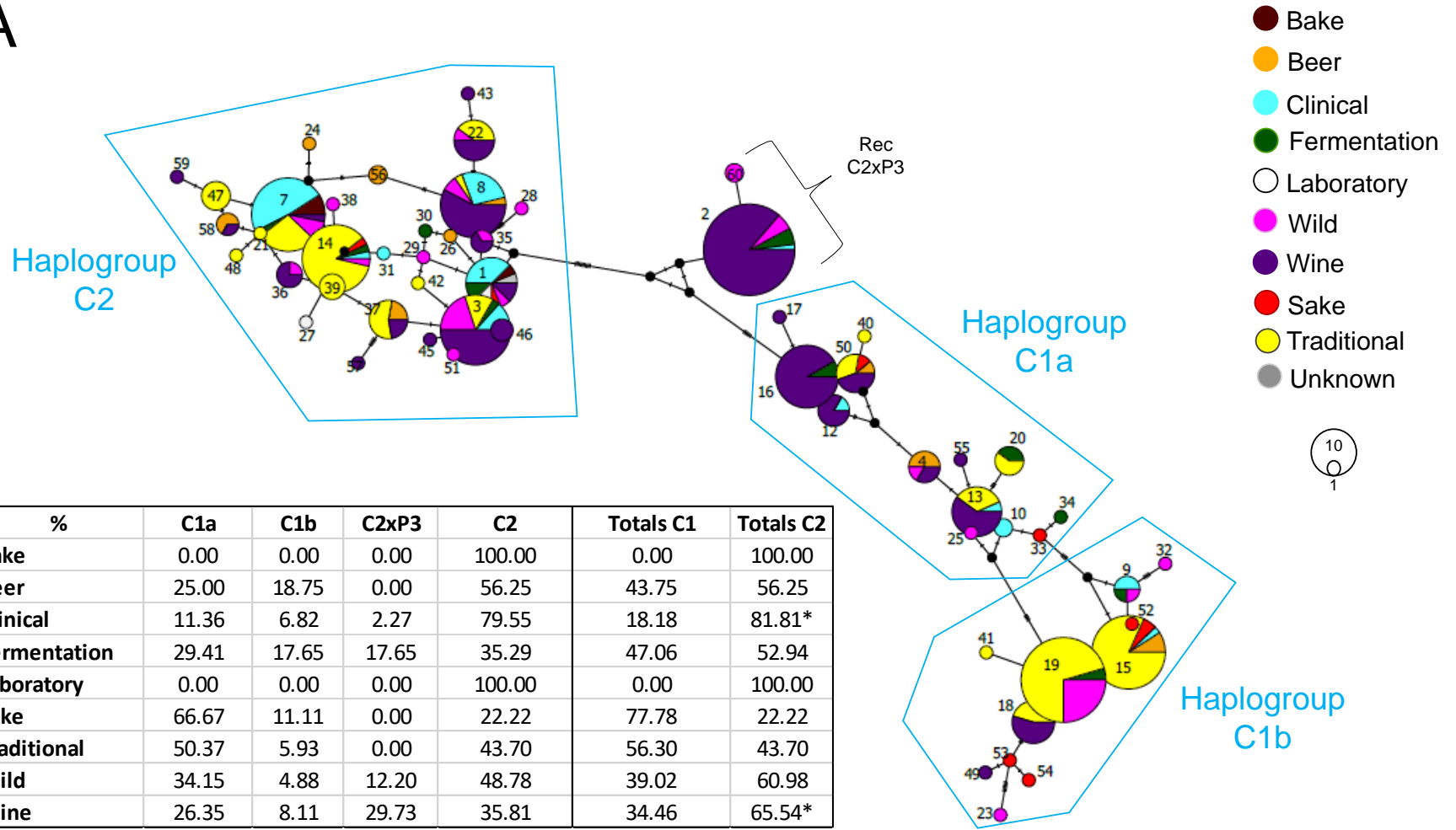


308

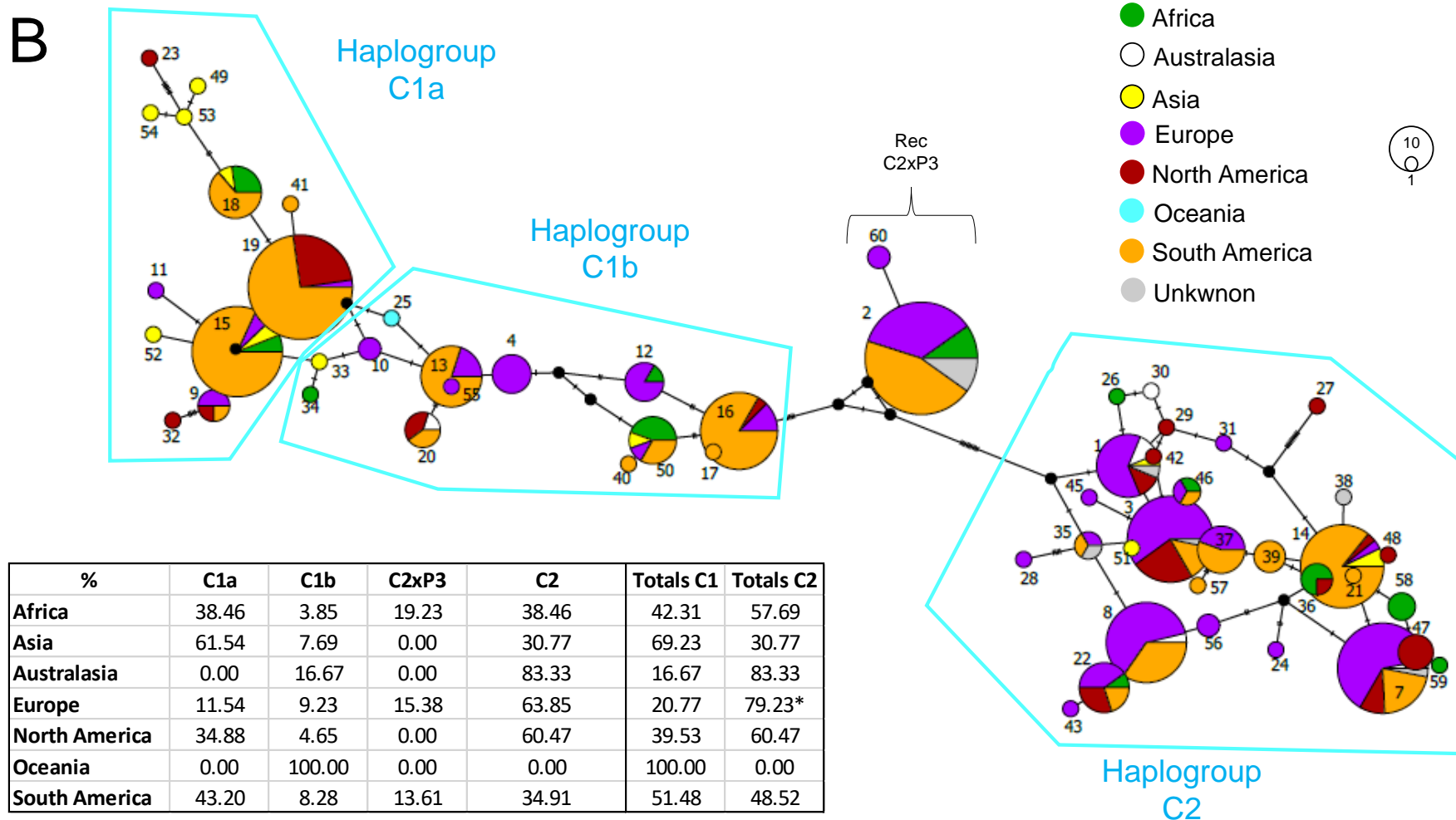
309 A *COX3* Neighbor-Joining phylogenetic tree is shown. Strains are colored according to the species designation. Bootstrap
310 values above 50 are given for each branch. Scale is given in nucleotide substitution per site.

311 **Figure S7. *S. cerevisiae* COX2 MJ networks.**

A



312



313

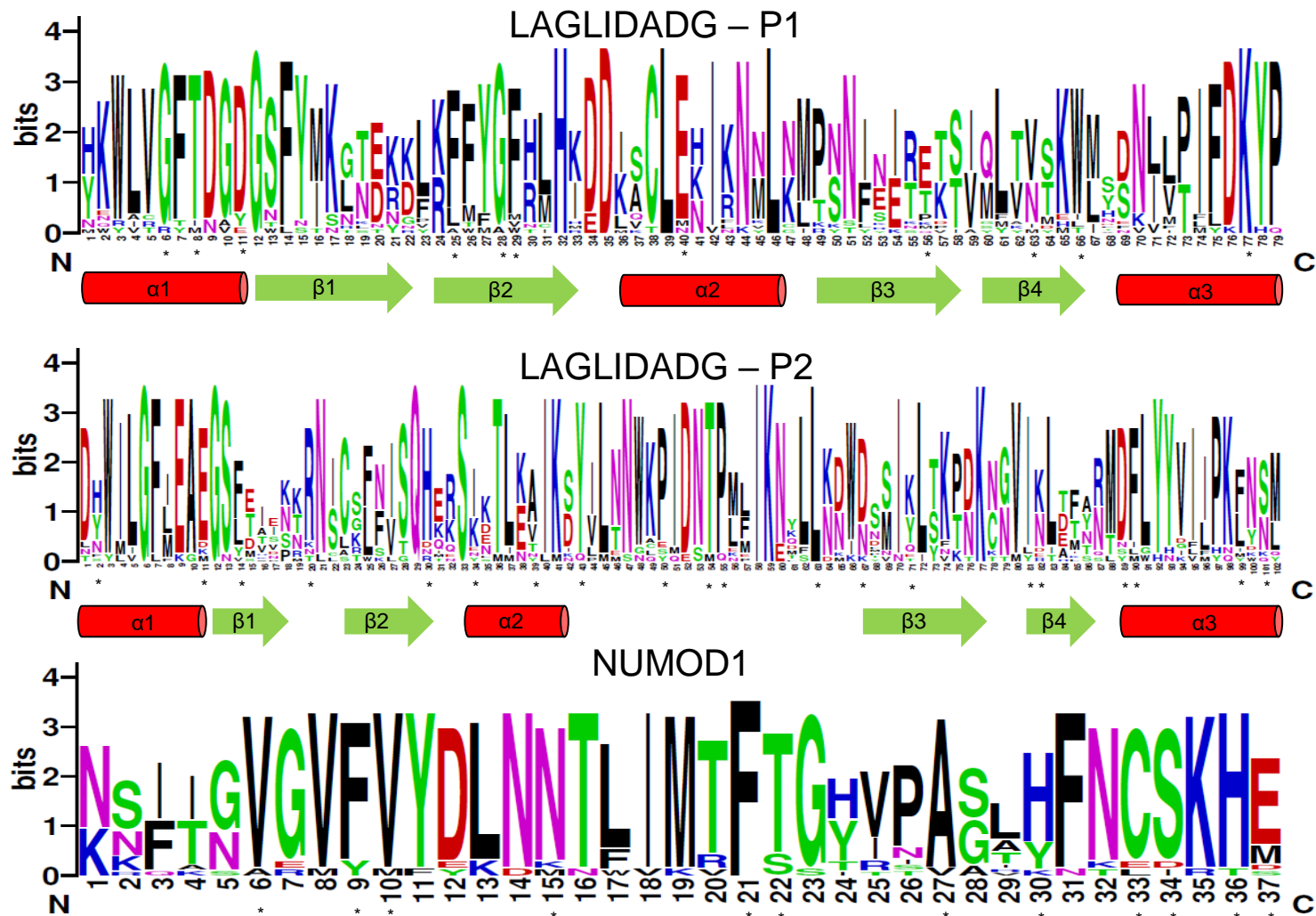
314 *S. cerevisiae* COX2 MJ networks were represented in A) and B) where haplotype pie charts were colored according to the
 315 isolation source and continent of isolation, respectively. Circle sizes represent the number of sequences in a haplotype and
 316 is scaled according to the legend. Number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

317 connecting the haplotypes. A table showing the strain percentage distribution in each haplogroup or recombinant group by
318 isolation source or continent is also displayed in A) and B), respectively.

322 A starting model of potential introgressions in *Saccharomyces* is described. Black, light blue, orange, green, dark blue, pink,
323 red and purple yeasts represent the yeast ancestors, *S. cerevisiae*, *S. paradoxus*, *S. mikatae*, *S. kudriavzevii*, *S. arboricola*,
324 *S. uvarum* and *S. eubayanus* yeasts, respectively. Small circles inside yeasts represent the mtDNA, which is colored
325 according to the species designation. A *S. paradoxus* yeast with a blue small circle indicates a potential inheritance of
326 mtDNA from *S. cerevisiae*. Red, green or red/green boxes represent type I, type II or recombinant *ORF1* sequences,
327 respectively. A question tag indicates dubious scenario, unknown *ORF1* or unknown mtDNA sequence due to the absence
328 of *COX3* sequence to support the mtDNA inheritance.

329

330 Figure S9. Weblogo representation of LAGLIDADG and NUMOD domains.



331
 332 Weblogo (<http://weblogo.berkeley.edu/>) of the three homing endonuclease domains for *Saccharomyces ORF1* genes plus
 333 *SasefMp08* from *Kazchastania servazii*, and *ORF1* and *ORF3* from *Cyberlindnera (Williopsis) saturnus* var. *suaveolens* are

334 represented. Cylinders and arrows represent α -helix and β -sheet, as previously described (Dalgaard *et al.* 1997). Asterisk
335 symbol indicates the codons for those aminoacids under purifying selection detected by DataMonkey.

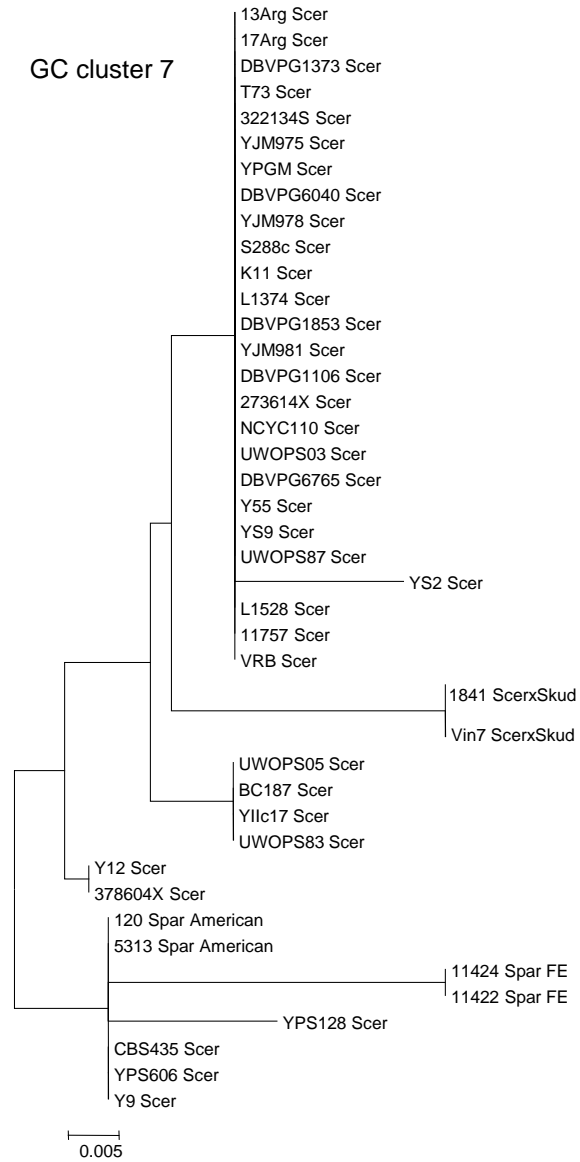
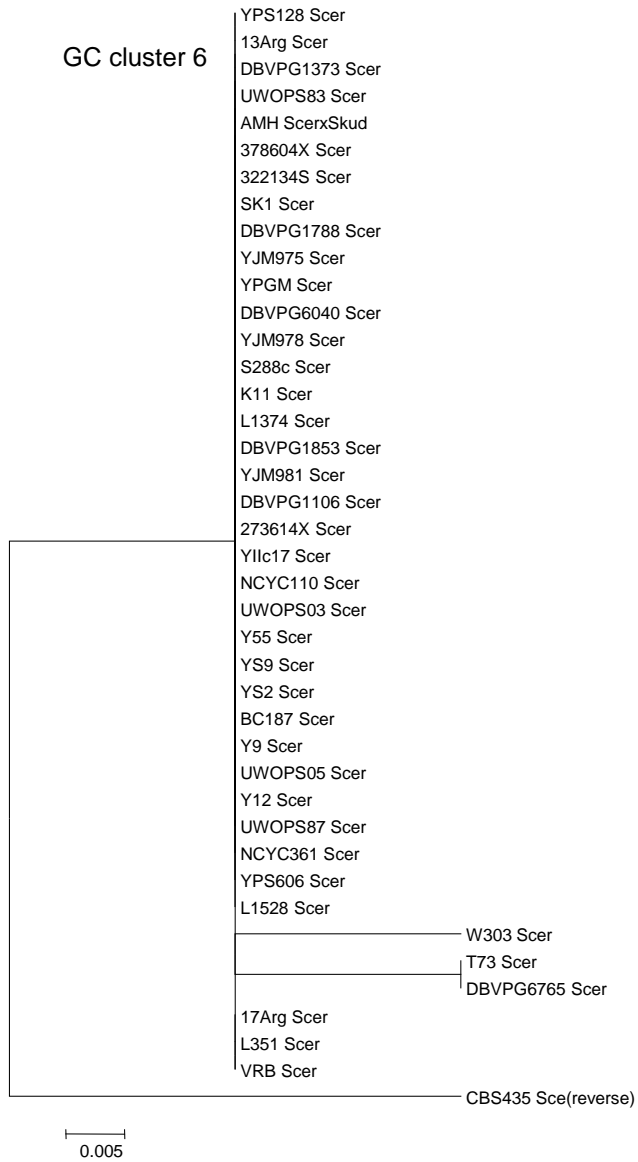
336

337 **Figure S10. COX2 and ORF1 aminoacid alignment (Supplementary File Figure S10).**

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

341

349 **Figure S12. GC cluster 6 and 7 Neighbor-Joining phylogenetic trees.**



351 To classify the GC cluster according to sequence similarity we reconstructed the NJ phylogenetic tree of GC cluster 6 and
352 7. This classification is shown by numbers in Figure S3. Scale bars represent number of substitutions per site.