# Virus-host co-evolution under a modified nuclear genetic code

**Among eukaryotes with modified nuclear genetic codes, viruses are unknown. However, here we provide evidence of an RNA virus that infects a fungal host (*Scheffersomyces segobiensis*) with a derived nuclear genetic code where CUG codes for serine. The genomic architecture and phylogeny are consistent with infection by a double-stranded RNA virus of the genus *Totivirus*. We provide evidence of past or present infection with totiviruses in five species of yeasts with modified genetic codes. All but one of the CUG codons in the viral genome have been eliminated, suggesting that avoidance of the modified codon was important to viral adaptation. Our mass spectroscopy analysis indicates that a congener of the host species has co-opted and expresses a capsid gene from totiviruses as a cellular protein. Viral avoidance of the host's modified codon and host co-option of a protein from totiviruses suggest that RNA viruses co-evolved with yeasts that underwent a major evolutionary transition from the standard genetic code.**

1    Derek J. Taylor[1], Matthew J. Ballinger, Shaun M. Bowman, Jeremy A. Bruenn

2    Department of Biological Sciences, The State University of New York at Buffalo, Buffalo,

3    NY 14260, USA.

4

5    [1]Correspondence and requests for materials should be addressed to D. J. T. (e-mail:

6    *djtaylor@buffalo.edu*; phone: 716-645-2880).

7

**Introduction**

Crick (1968) declared the universal genetic code to be nearly immutable because change would cause 'mistakes' in so many of the proteins of a cellular life form. However, Crick also implied that viruses are a possible exception to this evolutionary 'freezing' process because viruses have but a few protein coding targets. It is now well-established that modified genetic codes have evolved from the universal genetic code at least 34 times (do Céu Santos & Santos 2012). Yet, the consequences of these shifts for virus-host co-evolution remain poorly understood (Holmes 2009; Shackelton & Holmes 2008). Indeed, some authors have proposed that genetic code variants evolved as an antiviral defense (Holmes 2009; Shackelton & Holmes 2008; Taylor & Bruenn 2009). As viruses must use the protein translation machinery of the host, differences in genetic codes could preclude viral transfers among hosts. An evolutionary leap as great as a genetic code change could allow hosts to escape the co-evolutionary struggle with viruses. In agreement with the antiviral hypothesis, viruses are unknown from organisms with modified nuclear genetic codes. Viruses are known to infect the mitochondrial genomes of fungi with alternative mitochondrial genetic codes, but the tremendous divergence of these viruses from known viruses that use the ancestral genetic code obscures their origins (Shackelton & Holmes 2008). Modes of viral adaptation to hosts with non-standard genetic codes remain mysterious.

The CUG codon of the "CTG yeasts" (a diverse monophyletic group that contains human pathogens such as *Candida albicans* and wood-digesting species such as *Scheffersomyces stipitis*) has been reassigned from leucine to serine such that the serine tRNA possesses both derived serine and ancestral leucine sequence motifs (Butler et al. 2009; Santos et al. 2011). This shift in the genetic code, which replaces a hydrophobic residue with a polar residue, interferes with protein folding and can affect surface residue function (Feketová et al. 2010; Rocha et al. 2011). As a consequence, CTG yeasts lack the CUG codon

33    from (>90%) functionally relevant positions in proteins (Rocha et al. 2011) and rarely receive

34    genes via horizontal transfer compared to fungi with the standard code (Fitzpatrick 2011;

35    Richards et al. 2011). Viruses face the same functional barrier posed by this modified code.

36    However, the recent finding of "fossils" of totiviruses in the nuclear genomes of CTG yeast

37    (Taylor & Bruenn 2009) could indicate that exogenous RNA viruses have adapted to a

38    modified nuclear genetic code but remain undetected.

39    Totiviruses have a double stranded RNA genome (from 4.5-8kb in size in fungi) and are

40    characterized by an overlapping open reading frame between the capsid gene and the RdRp

41    gene with a programmed ribosomal frameshift. The family (Totiviridae) that contains the

42    totiviruses is ancient with a broad eukaryotic host range and extensive co-evolution in the

43    fungi (Liu et al. 2012). Totiviruses 'snatch' the hosts' mRNA caps (modified guanines at the

44    5' end) with a unique binding mechanism involving at least five proposed residues of the coat

45    protein (Fujimura & Esteban 2011). Frequently, totiviruses in fungi are associated with a

46    satellite killer dsRNA virus that codes for a toxin (Bostian et al. 1984). The fossils of

47    totiviruses are best characterized in the CTG yeast, *S. stipitis* (Frank & Wolfe 2009; Taylor &

48    Bruenn 2009). There are four tandem copies of a capsid-like protein gene in the genome of *S.*

49    *stipitis*, but it is unknown if these copies are translated into proteins. We explored the

50    existence of RNA viruses in the CTG clade of yeasts and their possible mode of co-evolution

51    by attempting to isolate modern and fossil viral genomes and their products from *S. stipitis*

52    and its relatives *Scheffersomyces segobiensis* and *Scheffersomyces coipomoensis*.

53    **Materials & Methods**

54    **Cell Cultures**. Freeze-dried culture stock was obtained from the USDA ARS Culture

55    Collection for *Scheffersomyces segobiensis* (Santa Maria & Garcia Aser) Kurtzman & M.

56    Suzuki (2010) NRRL Y-11571 (Type strain) and for *Scheffersomyces coipomoensis* (Ramirez

57    & Gonzalez) Urbina & Blackwell, 2012 comb. nov. NRRL Y-17651 (Type strain). Yeast

58 cultures were grown in 150 ml of YPD broth (yeast extract 1%, peptone 2%, and dextrose

59 2%) with an inoculum (single colony) of cells from streaked YPD agar plates.

60 **dsRNA assay.** Total nucleic acids (depleted of ribosomal RNA) were extracted from whole

61 cells (Bruenn & Keitz 1976). We purified dsRNA using CF-11 chromatography (Franklin

62 1966). The results (see Fig. S1) were consistent over the two-year period of more than 10

63 assays, indicating a stable infection.

64 **Viral particle isolation.** The tentative totivirus and satellite virus dsRNAs from *S.*

65 *segobiensis* were isolated from a CsCl gradient fraction with a density of 1.40 g/cc (Fig. S2),

66 which is expected based on the known *Saccharomyces cerevisiae virus L-A*. Note that even

67 though totiviruses and their satellite viruses are separately encapsidated, there is an overlap in

68 their densities, since particles can encapsidate up to four copies of satellite dsRNA (making a

69 total of about 4.8 kbp, essentially the same as the totivirus at 4.6 kbp).

70 **PCR, RTPCR, Sanger and Next Generation Sequencing.** We extracted total RNA from

71 yeast cells using the Masterpure yeast RNA purification kit (Epicentre) and an RNAse-free

72 DNAse treatment. RNA-seq was used to sequence the putative RNA virus, examine the tRNA

73 expression of the host, test for host expression of known fungal viral sequences, and isolate

74 host protein-coding sequences for bioinformatics analysis. Ribosomal RNA species were

75 removed using the Ribo-zero Magnetic Gold kit (Epicentre). We then prepared an RNA-seq

76 library using the ScriptSeq™ v2 RNA-Seq Library Preparation Kit (Epicentre). Libraries

77 were quantified using the Agilent 2100 Bioanalyzer RNA 6000 Pico Chip and submitted to

78 the University at Buffalo Next Generation sequencing facility for RNA sequencing. The

79 facility carried out RNA-seq using 50-cycle paired-end runs on two flow cells of an Illumina

80 HiSeq 2000. *De novo* RNA sequence assembly was carried out in CLC Assembly Cell 4.06

81 (http://www.clcbio.com) on an Apple Macintosh Mac Pro Xeon 64-bit workstation. The

82 putative viral contigs were reassembled for mapping purposes using the reference assembly

83    algorithm and the *de novo* contigs as references. A total of 409665 reads were mapped to the

84    totivirus with an average coverage of 4404.80 times.

85    The putative viral sequence was confirmed by Sanger sequencing using random and specific

86    primers for cDNA library construction. We used the Takara 5'-Full RACE Core Set to expand

87    sequence. Sanger and next generation sequences were compared using Geneious version 5.6.3

88    created by Biomatters (available from http://www.geneious.com/).

89    Because endogenous RNA viruses of fungi can be fragmented and differ greatly in their

90    nucleotide sequences from known viruses (Taylor et al. 2009), PCR probes alone are often an

91    ineffective tool for paleoviral discovery. We therefore carried out 454 Life Sciences

92    (http://www.454.com) sequencing with GS FLX Titanium series reagents of a DNA library

93    from *Scheffersomyces coipomoensis*. This form of sequencing also permitted multigene

94    bioinformatics analysis of the host protein coding genes. Strain identity of the assembly was

95    confirmed by BLAST analysis (Altschul et al. 1990) of the nuclear ribosomal RNA sequences

96    that are known for *S. coipomoensis*. We obtained 614185 reads with about 252 Mb of aligned

97    bases. The assembly carried out in Newbler (http://www.454.com), yielded 14.7 Mb of

98    aligned bases (488 contigs) with an average peak depth of 12X.

99    To establish that the virus was coded by exogenous RNA and not by the DNA of the host, we

100   compared RT-PCR and PCR products. For RNA templates, DNase-treated extracts were

101   exposed to RT-PCR using the Qiagen one step RT-PCR kit. For DNA templates, nucleic acid

102   extracts were exposed to PCR by excluding reverse transcriptase from the RT-PCR protocol.

103   We amplified a fragment of the single copy xylose reductase gene as a positive control for the

104   PCR of DNA. Primers used were: segoxylF CTGTTCTGAACAGATCTACCGTGC (xylose

105   reductase), segoxylR AAGTATGGGTGGTGTTCAACTTGC (xylose reductase), SvLgap3F

106   CGCAATACGACCAGGAGATTG (RdRp of virus from *S. segobiensis*), and segoSvLgap3R:

107   GTACACCAAGGTTAGTAGACAAG (RdRp of virus from *S. segobiensis*). cDNA synthesis

108    was performed at 48° C for 30 minutes, followed by 15 minutes at 94° C for reverse

109    transcriptase deactivation and Taq activation. DNA only reactions were added to the thermal

110    cycler 2 mins before the end of the previous 94° C step to activate the Taq polymerase. PCR

111    amplification was done for 35 cycles of 94° C for 30 s, 48° C for 30 s, and 72° C for 1 min. A

112    final extension at 72° C for 10 min was performed. New sequences from this study have the

113    following Genbank accession numbers: KC610514, KC616419-KC616429.

114    **Bioinformatics and protein mass spectroscopy.** We obtained amino acid sequences from

115    totivirids using the BLAST blastp algorithm with the the capsid gene and the RdRp sequences

116    of *Saccharomyces cerevisiae virus L-BC(La)* as queries and E_values <1e-05.  We searched

117    the non-redundant (nr) peptide sequence database (National Center for Biotechnology

118    Information, Bethesda, USA) and the Department of Energy Joint Genome Institute (J.G.I.)

119    genome browser for matches. Fossil or paleoviral copies of Totivirus-like genes were

120    identified by significant BLAST tblastx hits (E_values <1e-05) of relevant NCBI databases

121    using the sequences of *Saccharomyces cerevisiae virus L-BC(La)*. Duplicated capsid gene

122    copies adjacent to the complete integrated viral genomes in the assemblies of *S. stipitis* and *D.*

123    *hansenii* were assumed to be paralogs (Taylor & Bruenn 2009). These duplicated paleoviruses

124    and closely related (i.e. phylogenetic sister viruses) co-infecting viral strains were omitted for

125    the phylogenetic analyses. Sequences were aligned using MAFFT (Katoh et al. 2009) with

126    default settings. We carried out maximum likelihood analyses with PhyML 3.0 as

127    implemented in Seaview 4.3.5 (Anisimova & Gascuel 2006; Gouy et al. 2010).  Model

128    optimization in Prottest (Abascal et al. 2005) indicated that the LG + invariable sites

129    parameter (I)+ gamma parameter for among-site rate variation (G) was the best fit under a

130    Bayesian Information Criterion (BIC). For reliability estimates we used SH-like approximate

131    likelihood ratio tests (Anisimova & Gascuel 2006). Searches were comprised of five random

132    starts under the subtree pruning and regrafting (SPR) algorithm and midpoint rooted.

133    Aguileta et al. (2008) found that concatenation of the two most informative genes from a

134    genomic scale assessment recovered an expected reference fungal phylogeny with strong

135    support. We used the approach of Taylor et al. (2009) who concatenated five of the most

136    phylogenetically informative fungal genes (Aguileta et al. 2008) to estimate fungal relations.

137    Accession numbers for the genes used (Minichromosome Maintenance protein 7[MCM7],

138    Kontroller of Growth[KOG1], Elongator complex subunit[ELP3], NAD-specific glutamate

139    dehydrogenase[GDH2], and acetolactate synthase [ILV2]) in the fungal analysis are presented

140    in Table S1. Data were collected from GenBank and from our newly sequenced cultures of *S.*

141    *coipomoensis* and *S. segobiensis*. Sequences were aligned in MAFFT and exposed to

142    maximum likelihood analyses in RAxML 7.3.2 (Stamatakis 2006) and in PhyML 3.0. Models

143    were partitioned by gene in RAxML using the best-fit models as indicated by

144    PartitionfinderProtein (Lanfear et al. 2012).

145    Relative synonymous codon usage (Sharp et al. 1986) and third position base composition for

146    yeasts and viruses was calculated using the CAIcal server (Puigbò et al. 2008). For viruses the

147    entire open reading frame of the genome was used in the calculations. For yeasts we used the

148    representative genes from the phylogenetic analysis. Species used for the yeast and viral

149    codon usage analyses are listed in Table S2. Bivariate plots of RCSU and base composition

150    were graphed using the R statistical programming language (Ihaka & Gentleman 1996).

151    FSfinder (Moon et al. 2004) was used to locate putative slippery sites and pseudoknots in the

152    totiviral genome. The $tRNA_{CAG}^{Ser}$ for *S. segobiensis* was folded according to the model for

153    *Candida albicans* (Santos et al. 2011) using the VARNA secondary structure visualization

154    program (Darty et al. 2009).

155    Structural information for *Saccharomyces cerevisiae virus L-A* is from the crystal structure of

156    the ScV L-A capsid protein (Naitow et al. 2002). For ScV L-BC (La) and SsV L, structural

157   information was predicted by the I-TASSER webserver using ScV L-A cap as a template

158   (Roy et al. 2010; Roy et al. 2012).

159   To examine protein expression of paleoviral copies we isolated crude protein from *S. stipitis*

160   and *S. cerevisiae* by French press and further isolated proteins migrating between 73 kDa and

161   92 kDa from 10% SDS-PAGE. Protein mass spectroscopy was carried out at the Seattle

162   Biomedical Research Institute Proteomics Core Facility.

163

**Results and Discussion**

164

165   Because the fossil viruses in yeast have a similar architecture to dsRNA totiviruses

166   (Taylor & Bruenn 2009), we carried out a specific chromatographic assay for dsRNA. We

167   detected dsRNA products in *S. segobiensis* with approximately the same gel-estimated size to

168   the totivirus (4.5 kb) and satellite virus of *S. cerevisiae* (1.2 kb, Fig. S1). Viral particles

169   containing both sizes of dsRNAs were also isolated by CsCl equilibrium gradient

170   centrifugation (Fig. S2). No such products were detected in *S. stipitis* or in *S. coipomoensis*. A

171   tblastn using the protein sequences of the two known totiviruses from *S. cerevisiae* as queries

172   revealed significant matches to a contig from a database of our RNA sequence assemblies

173   (extracted from *S. segobiensis* cells) of similar length to the dsRNA particles on the gel.

174   Assemblies from Sanger sequencing of the RNA virus agreed with the assembly using

175   Illumina RNA sequencing, but the 5' UTR was complete only in the Illumina assembly. The

176   assembled virus had the genomic architecture of totiviruses with overlapping capsid and

177   RdRp open reading frames flanked by 5' and 3' UTRs (Fig. 1). We identified a putative

178   slippery site for ribosomal frameshifting (GGGTTTT) at position 1981 that was

179   independently identified using the frameshift prediction software fsfinder (Moon et al. 2004).

180   The five sites identified as functionally important for cap-snatching in totiviruses are

181     conserved in the virus from *S. segobiensis* (Fig. S3), suggesting a cap-snatching mechanism

182     similar to those of well-studied totiviruses. These sites show weak conservation in the fossil

183     copies, consistent with the loss of host mRNA decapping in host-coded elements. The

184     successful PCR amplification of a single copy nuclear gene fragment (xylose reductase gene)

185     from the host genome indicates that the DNA template was of sufficient quality to detect

186     endogenous viral genes using our methods (Fig. S4). However, the primers nested within the

187     viral genome failed to PCR amplify a DNA copy from the host (*S. segobiensis*) genome, but

188     RT-PCR did amplify an RNA copy of the viral gene. The results support the existence of an

189     exogenous RNA virus in *S. segobiensis* with the genomic architecture of a totivirus.

190        Further evidence of affinity to totiviruses comes from sequence analysis of the virus in

191     *S. segobiensis*. A BLAST blastp analysis of the RdRp-like ORF yielded a conserved domain

192     match (E=1.56e-58) to RdRP_4, a viral RNA-directed RNA-polymerase family that includes

193     "RdRPs from Luteovirus, Totivirus and Rotavirus". The best expect value (E= 2e-137)

194     obtained was the totivirus, *Saccharomyces cerevisiae virus L-BC(La)*, with an identity of 37%

195     of residues. The RdRp gene phylogeny (Fig. 2A) positioned the tentative species

196     *Scheffersomyces segobiensis virus L* within the totiviruses and most closely related to

197     *Saccharomyces cerevisiae virus L-BC(La)*. Support values are strong enough to rule out

198     random error as an explanation for the evolutionary position of *Scheffersomyces segobiensis*

199     *virus L* within the totiviruses. The less conserved capsid gene tree (Fig. 2B) showed a similar

200     association, but with fewer outgroup sequences and more paleoviruses. We detected fossil

201     totiviruses in the genomes of *S. coipomoensis,* and *Pichia membranifaciens* and a putative

202     totivirus in the assembly of *Nadsonia fulvescens*. We deem the sequences from *Nadsonia* to

203     be putative viral sequences because they are present in the RNA-based EST libraries, but not

204     the DNA based genome assembly (see Liu et al. (2012) for a discussion of this approach to

205     dsRNA virus discovery).

206    As with other totiviruses, we found evidence that *Scheffersomyces segobiensis virus L*

207    has a putative killer satellite virus. We carried out a BLAST search to identify candidate RNA

208    contigs for the dsRNA band observed earlier. A sequence was obtained of the correct size (1.2

209    kbp) that had a significant match to the killer satellite K2 virus of *S. cerevisiae*. An RTPCR

210    with specific primers confirmed that the contig was not coded in the DNA genome of the

211    yeast. The existence of a satellite virus bolsters the evidence for a totivirus in the CTG clade.

212    The evolutionary and genomic evidence indicates that *Scheffersomyces segobiensis*

213    *virus L* originated from exogenous totiviruses that use a standard genetic code. The RdRp

214    permits the deepest assessment of evolution, and it reveals a derived position of the

215    *Scheffersomyces segobiensis virus L* at the tip of the standard code totivirids. The known

216    paleoviruses in the CTG clade are distantly related to *Scheffersomyces segobiensis virus L* and

217    lack the ability to produce the fusion gene product typical of totiviruses (Taylor & Bruenn

218    2009). Nor is there any evidence of a functioning totivirus genome being coded in the DNA

219    from known genomes of yeasts. Moreover, our RTPCR results indicate that the DNA genome

220    of the host, *S. segobiensis*, lacks coding sequences related to the viral genome. Finally,

221    successful endogenization of a virus that imparts a selective disadvantage by obligately

222    decapping host mRNA seems unlikely. Most of the successful paleoviruses in the CTG clade

223    possess functionally differing residues at the decapping sites. Because the CTG virus we

224    discovered has the conserved residues for such a decapping mechanism, an "escaped" genome

225    hypothesis requires that the unique decapping mechanism was lost in the host genome and

226    then re-evolved in the escaped viral genome – also unlikely. We conclude that adaptation to

227    the genetic code shift of hosts happened in exogenous viruses.

228    Our viral evolutionary trees are consistent with the jumping of viruses between hosts

229    with different genetic codes. Host products from the CTG clade are phylogenetically

230    interspersed with standard code sequences in both major clades and the virus/paleovirus

231    phylogenies bear little resemblance to the host relationships. For example, a virus from the

232    truffle (*Tuber aestivum*) is most closely related to the fossil virus in the CTG yeast *S. stipitis*

233    rather than the virus of another genus of basidiomycete, *Xanthophyllomyces*. Horizontal

234    transfer of paleoviruses could be a source of some evolutionary noise, but this process appears

235    rare in the CTG clade. The timescale of the totivirus evolution is difficult to estimate because

236    of the dearth of reliable fossil calibrations for fungi (Berbee & Taylor 2010; Rolland & Dujon

237    2011). However, our observed close sequence and structural similarity for rapidly evolving

238    capsid proteins of RNA viruses likely postdates the ancient split of ascomycetes with

239    basidiomycetes (452 MYA to 1400 MYA) and the origin of the CTG clade (> 150

240    MYA)(Berbee & Taylor 2010; Massey et al. 2003; Pesole et al. 1995).

241         The evolutionary position, CUG codon usage, and tRNA expression evidence are

242    consistent with the host, *S. segobiensis*, having a modified CTG genetic code. Our

243    phylogenetic analysis (Fig. 3A) of nuclear protein coding genes revealed strong support for

244    the placement of the host yeast within the CTG clade and as a sister species to the CTG

245    species *S. stipitis*. The monophyly of the CTG clade is well established in evolutionary

246    genomics (Butler et al. 2009; Louis et al. 2012; Wohlbach et al. 2011). The close sister group

247    relationship of *S. segobiensis* and *S. stipitis* has also been independently supported by several

248    studies using the ribosomal rRNA gene family (Cadete et al. 2012; Kurtzman 2010; Urbina &

249    Blackwell 2012). CUG comprises a substantial percentage of codons for leucine residues in

250    standard code yeasts, while being almost absent in the CTG yeasts at the homologous leucine

251    position (Fig. 3B). The usage of CUG in CTG yeast (including *S. segobiensis*) is

252    underrepresented compared to the standard code yeast (Fig. 3C). A BLAST search of RNA

253    contigs for the characteristic modified chimeric serine tRNA of *S. stipitis* that recognizes

254    CUG revealed a significant match in *S. segobiensis*.  That is, the sequence has both serine and

255    leucine identity sites (Fig. 3D). We failed to detect a "standard code" leucine tRNA. The pre-

256    tRNA species (Fig. S5) had unique mutations in the flanking regions from *S. stipitis*,

257    consistent with the modest divergence of a sister species.

258        As the genetic code shift had a profound functional effect on the proteins of yeasts, we

259    expected the virus to adapt to the shift. We found that *Scheffersomyces segobiensis virus L*

260    had but a single codon of the modified CUG type. This evolutionary loss of CUG codons

261    resulted in the lowest CUG frequency known among related mycoviruses, where the CUG

262    codon is generally overrepresented.  The sole CUG in *Scheffersomyces segobiensis virus L*

263    occurs at a position in the capsid protein that appears to be structurally unimportant (Fig. 1).

264    *Scheffersomyces segobiensis virus L* appears to have adapted to the host shift in genetic code

265    by eliminating functionally relevant CUG codons.

266        In plots of relative synonymous codon usage (RSCU) versus third position base

267    composition (Fig. 4A; Table S2), *Scheffersomyces segobiensis virus L* grouped with CTG

268    yeasts rather than with other totivirids. The same pattern of *Scheffersomyces segobiensis virus*

269    *L* grouping with CTG yeasts to the exclusion of other viruses was found for relevant leucine

270    codons (Figs. 3B,C). In *S. cerevisiae*, CUN codons are decoded by either tRNA-UAG or by

271    tRNA-GAG. But in the CTG clade yeast, CUN codons are decoded differently. Here a

272    derived tRNA-CAG decodes the reassigned CUG codon, while the remaining CUN family

273    codons are decoded by a single tRNA-IAG. The inosine interacts only weakly with CUA

274    compared to CUC and CUU in *C. albicans* (Massey et al. 2003), and seems the most likely

275    driving force behind reduction in CUA usage in CTG yeasts. We note that *Scheffersomyces*

276    *segobiensis virus L* lacks a bias in base composition at the third position for the examined

277    codons, suggesting that the observed codon usage bias in *Scheffersomyces segobiensis virus L*

278    is more complicated than base compositional shifts alone.

279        Our analysis of the genome of *Scheffersomyces segobiensis virus L* is consistent with

280    adaptation to offset the functional effects of the genetic code shift in the host. But, our results

281 also indicate that the endogenization of viral genes by host yeasts of both genetic codes is

282 more common than previously thought. Presently, it is unknown if these endogenous non-

283 retroviral genes function as proteins or are merely transcriptional noise. Here, we show that at

284 least one of these genes derived from a totivirus is expressed as a protein. Isolation of proteins

285 migrating between 73 kDa and 92 kDa from *S. stipitis* yielded approximately 535 proteins, as

286 estimated from the equivalent size range of proteins in *S. cerevisiae*. The mass spectroscopy

287 analysis was able to unambiguously identify the capsid protein (cap) of *Saccharomyces*

288 *cerevisiae virus L-A* in the *S. cerevisiae* control and 153 proteins in *S. stipitis*. By this method,

289 we were able to detect one of the four NIRV capsid polypeptides from *S. stipitis*. Figures S6A

290 and S6B show the distribution of tryptic peptides detected by mass spectroscopy from *S.*

291 *stipitis* capsid4 and *Saccharomyces cerevisiae virus L-A* cap. All of the *S. stipitis* virus-like

292 peptides had probabilities of 0.999 or greater. Our inability to detect the remaining cap

293 proteins from *S. stipitis* may be due to the lack of sensitivity of the method used. We were

294 able to detect a single tryptic peptide from *Saccharomyces cerevisiae virus L-BC(La)*cap,

295 which is present at about one tenth of the concentration of *Saccharomyces cerevisiae virus L-*

296 *A* cap. Our results indicate that the co-option of a non-retroviral RNA viral protein has

297 occurred as with retroviral proteins (e.g. syncytin genes (Feschotte & Gilbert 2012)).

298 Although expression of similar capsid proteins has an antiviral effect (Valle & Wickner 1993;

299 Yao & Bruenn 1995), the current function of the co-opted viral proteins in yeast is unknown.

**Conclusions**

301 Our discoveries indicate that a major evolutionary transition involving a change in the genetic

302 code of the fungi failed to result in permanent host escape from viruses. We found evidence

303 of present or past viral infection in five lineages of yeasts with a modified genetic code. Thus,

304 viral infection is likely widespread in the CTG clade of fungi. The mode of viral adaptation

305 recalls the prediction of Crick (1968) that some viruses would be less susceptible to

306  evolutionary "freezing" because they present a reduced protein target. The genomes of

307  totiviruses are among the smallest known for RNA viruses (Holmes, 2009). Still, even with a

308  small viral genome we found evidence that exogenous viral adaptation was associated with

309  the elimination of modified codons from functional positions. Our results also highlight the

310  value of recent paleovirological approaches to understanding virus-host biology (Aswad &

311  Katzourakis 2012; Feschotte & Gilbert 2012; Holmes 2011; Koonin 2010; Patel et al. 2011).

312  Fossil copies in the genomes of yeasts informed us about prior infections and made possible

313  our discovery of a virus adapted to a modified nuclear genetic code. Moreover, we found that

314  at least one of these "fossil" genes is co-opted by the yeast host and expressed as a protein.

315

321
322

322 **References**

323 Abascal F, Zardoya R, and Posada D. 2005. ProtTest: selection of best-fit models of protein

324         evolution. *Bioinformatics* 21:2104-2105.

325 Aguileta G, Marthey S, Chiapello H, Lebrun MH, Rodolphe F, Fournier E, Gendrault-

326         Jacquemard A, and Giraud T. 2008. Assessing the Performance of Single-Copy Genes

327         for Recovering Robust Phylogenies. *Systematic Biology* 57:613-627.

328 Altschul SF, Gish W, Miller W, Myers, EW, and Lipman DJ. 1990. Basic local alignment

329         search tool. *Journal of Molecular Biology* 215:403-410.

330 Anisimova M, and Gascuel O. 2006. Approximate likelihood-ratio test for branches: A fast,

331         accurate, and powerful alternative. *Systematic Biology* 55:539-552.

332 Aswad A, and Katzourakis A. 2012. Paleovirology and virally derived immunity. *Trends in*

333         *Ecology & Evolution*. 27: 627-636

334 Berbee ML, and Taylor JW. 2010. Dating the molecular clock in fungi - how close are we?

335         *Fungal Biology Reviews*:1-16.

336 Bostian KA, Elliott Q, Bussey H, Burn V, Smith A, and Tipper DJ. 1984. Sequence of the

337         preprotoxin dsRNA gene of type 1 killer yeast: multiple processing events produce a

338         two-component toxin. *Cell* 36:741-751.

339 Bruenn J, and Keitz B. 1976. The 5' ends of yeast killer factor RNAs are pppGp. *Nucleic*

340         *Acids Research* 3:2427-2436.

341 Butler G, Rasmussen MD, Lin MF, Santos MAS, Sakthikumar S, Munro CA, Rheinbay E,

342         Grabherr M, Forche A, Reedy JL et al. . 2009. Evolution of pathogenicity and sexual

343         reproduction in eight Candida genomes. *Nature* 459:657-662.

344 Cadete RM, Melo MA, Dussán KJ, Rodrigues RCLB, Silva SS, Zilli JE, Vital MJS, Gomes

345         FCO, Lachance M-A, and Rosa CA. 2012. Diversity and Physiological

346       Characterization of D-Xylose-Fermenting Yeasts Isolated from the Brazilian

347       Amazonian Forest. *PLoS One* 7:e43135.

348    Crick FHC. 1968. The origin of the genetic code. *Journal of Molecular Biology* 38:367-379.

349    Darty K, Denise A, and Ponty Y. 2009. VARNA: Interactive drawing and editing of the RNA

350       secondary structure. *Bioinformatics* 25:1974.

351    do Céu Santos M, and Santos M. 2012. Structural and molecular features of non-standard

352       genetic codes. In: Cannarozzi G, and Schneider A, eds. *Codon Evolution: Mechanisms*

353       *and Models*. Oxford: Oxford University Press, 258-270.

354    Feketová Z, Mašek T, Vopálenský V, and Pospíšek M. 2010. Ambiguous decoding of the

355       CUG codon alters the functionality of the Candida albicans translation initiation factor

356       4E. *Fems Yeast Research* 10:558-569.

357    Feschotte C, and Gilbert C. 2012. Endogenous viruses: insights into viral evolution and

358       impact on host biology. *Nature Reviews Genetics* 13:283-296.

359    Fitzpatrick DA. 2011. Horizontal gene transfer in fungi. *FEMS Microbiology Letters* 329:1-8.

360    Frank AC, and Wolfe KH. 2009. Evolutionary capture of viral and plasmid DNA by yeast

361       nuclear chromosomes. *Eukaryotic Cell* 8:1521-1531.

362    Franklin RM. 1966. Purification and properties of the replicative intermediate of the RNA

363       bacteriophage R17. *Proceedings of the National Academy of Sciences, USA* 55:1504-

364       1511.

365    Fujimura T, and Esteban R. 2011. Cap-snatching mechanism in yeast LA double-stranded

366       RNA virus. *Proceedings of the National Academy of Sciences, USA* 108:17667-17671.

367    Gouy M, Guindon S, and Gascuel O. 2010. SeaView version 4: A multiplatform graphical

368       user interface for sequence alignment and phylogenetic tree building. *Molecular*

369       *Biology and Evolution* 27:221-224.

370     Holmes EC. 2009. *The evolution and emergence of RNA viruses*. New York: Oxford

371         University Press.

372     Holmes EC. 2011. The evolution of endogenous viral elements. *Cell Host & Microbe* 10:368-

373         377.

374     Ihaka R, and Gentleman R. 1996. R: A language for data analysis and graphics. *Journal of*

375         *computational and graphical statistics* 5:299-314.

376     Katoh K, Asimenos G, and Toh H. 2009. Multiple Alignment of DNA Sequences with

377         MAFFT. *Methods Mol Biol* 537:39-64.

378     Koonin EV. 2010. Taming of the shrewd: novel eukaryotic genes from RNA viruses. *BMC*

379         *Biol* 8:2.

380     Kurtzman CP. 2010. Phylogeny of the ascomycetous yeasts and the renaming of Pichia

381         anomala to Wickerhamomyces anomalus. *Antonie van Leeuwenhoek* 99:13-23.

382     Lanfear R, Calcott B, Ho SYW, and Guindon S. 2012. PartitionFinder: combined selection of

383         partitioning schemes and substitution models for phylogenetic analyses. *Molecular*

384         *Biology and Evolution* 29:1695-1701.

385     Liu H, Fu Y, Xie J, Cheng J, Ghabrial SA, Li G, Yi X, and Jiang D. 2012. Discovery of Novel

386         dsRNA Viral Sequences by In Silico Cloning and Implications for Viral Diversity,

387         Host Range and Evolution. *PLoS One* 7:e42147.

388     Louis VL, Despons L, Friedrich A, Martin T, Durrens P, Casarégola S, Neuvéglise C,

389         Fairhead C, Marck C, and Cruz JA. 2012. Pichia sorbitophila, an Interspecies Yeast

390         Hybrid, Reveals Early Steps of Genome Resolution After Polyploidization. *G3:*

391         *Genes| Genomes| Genetics* 2:299-311.

392

393      Massey SE, Moura G, Beltrão P, Almeida R, Garey JR, Tuite MF, and Santos MAS. 2003.

394          Comparative evolutionary genomics unveils the molecular mechanism of

395          reassignment of the CTG codon in Candida spp. *Genome Research* 13:544-557.

396      Moon S, Byun Y, Kim HJ, Jeong S, and Han K. 2004. Predicting genes expressed via− 1

397          and+ 1 frameshifts. *Nucleic Acids Research* 32:4884-4892.

398      Naitow H, Tang J, Canady M, Wickner RB, and Johnson JE. 2002. L-A virus at 3.4 A

399          resolution reveals particle architecture and mRNA decapping mechanism. *Nat Struct*

400          *Biol* 9:725-728.

401      Patel MR, Emerman M, and Malik HS. 2011. Paleovirology--ghosts and gifts of viruses past.

402          *Current Opinion in Virology* 1:304-309.

403      Pesole G, Lotti M, Alberghina L, and Saccone C. 1995. Evolutionary origin of nonuniversal

404          CUG (Ser) codon in some Candida species as inferred from a molecular phylogeny.

405          *Genetics* 141:903.

406      Puigbò P, Bravo IG, and Garcia-Vallve S. 2008. CAIcal: a combined set of tools to assess

407          codon usage adaptation. *Biology direct* 3:38.

408      Richards TA, Leonard G, Soanes DM, and Talbot NJ. 2011. Gene transfer into the fungi.

409          *Fungal Biology Reviews* 25:98-110.

410      Rocha R, Pereira PJB, Santos MAS, and Macedo-Ribeiro S. 2011. Unveiling the structural

411          basis for translational ambiguity tolerance in a human fungal pathogen. *Proceedings of*

412          *the National Academy of Sciences, USA* 108:14091-14096.

413      Rolland T, and Dujon B. 2011. Yeasty clocks: Dating genomic changes in yeasts.1-9.

414      Roy A, Kucukural A, and Zhang Y. 2010. I-TASSER: a unified platform for automated

415          protein structure and function prediction. *Nature protocols* 5:725-738.

416      Roy A, Yang J, and Zhang Y. 2012. COFACTOR: an accurate comparative algorithm for

417          structure-based protein function annotation. *Nucleic Acids Research* 40:W471-W477.
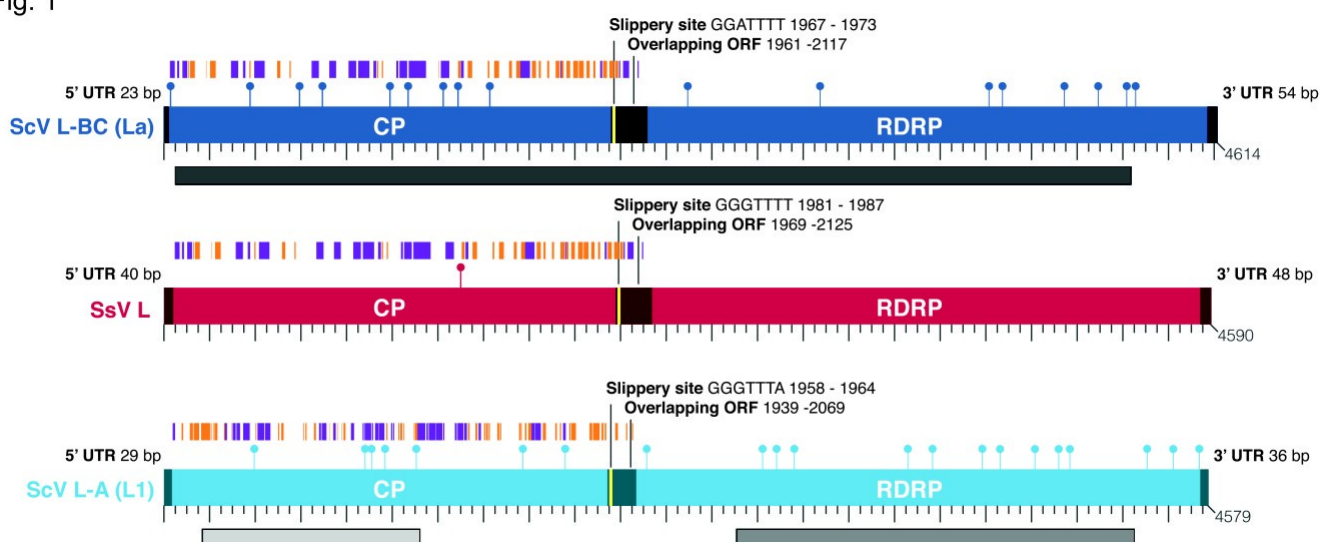
418    Santos MAS, Gomes AC, Santos MC, Carreto LC, and Moura GR. 2011. The genetic code of

419         the fungal CTG clade. *Comptes rendus - Biologies* 334:607-611.

420    Shackelton LA, and Holmes EC. 2008. The role of alternative genetic codes in viral evolution

421         and emergence. *Journal of Theoretical Biology* 254:128-134.

422    Sharp PM, Tuohy TMF, and Mosurski KR. 1986. Codon usage in yeast: cluster analysis

423         clearly differentiates highly and lowly expressed genes. *Nucleic Acids Research*

424         14:5125-5143.

425    Stamatakis A. 2006. RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses

426         with thousands of taxa and mixed models. *Bioinformatics* 22: 2688-2690.

427    Taylor DJ, and Bruenn J. 2009. The evolution of novel fungal genes from non-retroviral RNA

428         viruses. *BMC Biology* 7:88.

429    Urbina H, and Blackwell M. 2012. Multilocus Phylogenetic Study of the Scheffersomyces

430         Yeast Clade and Characterization of the N-Terminal Region of Xylose Reductase

431         Gene. *PLoS One* 7:e39128.

432    Valle RP, and Wickner RB. 1993. Elimination of L-A double-stranded RNA virus of

433         *Saccharomyces cerevisiae* by expression of *gag* and *gag-pol* from L-A cDNA clone.

434         *Journal of Virology* 67:2764-2771.

435    Wohlbach DJ, Kuo A, Sato TK, Potts KM, Salamov AA, Labutti KM, Sun H, Clum A,

436         Pangilinan JL, Lindquist EA et al. . 2011. Comparative genomics of xylose-

437         fermenting fungi for enhanced biofuel production. *Proceedings of the National*

438         *Academy of Sciences, USA*:1-6.

439    Yao W, and Bruenn JA. 1995. Interference with replication of two double-stranded RNA

440         viruses by production of N-terminal fragments of capsid polypeptides. *Virology*

441         214:215-221.

442

# Figure 1

Comparison of the genomic architecture of the newly discovered *Scheffersomyces segobiensis virus L* [ SsV L] (in red) that uses a modified nuclear genetic code with those of related totiviruses (*Saccharyomyces cerevisiae virus L-BC(La)* [ ScV L

The figures show overlapping reading frames for capsid (CP) and RNA dependent RNA polymerase genes (RdRp) that are typical of the double-stranded RNA totiviruses. Terminal UTRs (untranslated regions) and central overlapping reading frames are distinguished by dark colored shading. The positions of the ribosomal frameshift sites are indicated in yellow. Coding CUG codons are represented by lollipops. Capsid protein secondary structural domains are shown in purple (α-helices) and orange (β-sheets). BLASTp results for each of the previously known totiviruses to SsV L are shown as grey lines underlying the respective genomes, with darker shades indicating a lower expect value. The scale bar increments represent 50 nucleotides.

Fig. 1

# Figure 2

Evolutionary relationships of exogenous and endogenous totiviruses showing the derived and non-monophyletic positions of viruses and paleoviruses from CTG yeast.

Asterisks indicate paleoviral sequences. Numbers are support values (SH-like approximate likelihood tests estimated in PhyML 3.0). The phylograms are estimated using the maximum likelihood (ML) optimality criterion from alignments of predicted amino acid residues of (A) the RNA dependent RNA polymerase gene (RdRp) and (B) the capsid gene. Colored dashed boxes indicate the two major clades of totivirus-like sequences. Genbank Accession numbers are provided in parentheses.

Fig. 2

A



*Magnaporthe oryzae* virus 2 {164597962}
*Gremmeniella abietina* RNA virus L2 {50080151}
Aspergillus mycovirus 178 {161897699}
*Epichloe festucae* virus 1 {94536500}
*Coniothyrium minitans* mycovirus {78762704}
*Helicobasidium mompa* No17 dsRNA virus {33867952}
*Sphaeropsis sapinea* RNA virus 1 {9630962}
*Aspergillus foetidus* slow virus 1 {400131542}
*Beauveria bassiana* virus 1 {345108728}
*Magnaporthe oryzae* virus 1 {54193769}
*Helminthosporium victoriae* virus 190S {124484602}
*Tolypocladium cylindrosporum* virus 1 {315573170}
*Botryotinia fuckeliana* totivirus 1 {134141997}
*Saccharomyces cerevisiae* virus La {313104168}
*Scheffersomyces segobiensis* RNA virus 1
Black raspberry virus F {157939585}
*Saccharomyces cerevisiae* virus LA {NP_620495}
*Tuber aestivum* virus 1 {312233876}
*Scheffersomyces stipitis* {NC_009047}*
*Penicillium marneffei* {ABAR01000272}*
*Debaryomyces hansenii* {CR382134}*
*Xanthophyllomyces dendrorhous* virus L2 {383388877}
*Xanthophyllomyces dendrorhous* virus L1b {383388874}
*Rosellinia necatrix* quadrivirus 1 {374504767}
Amasya cherry disease associated mycovirus {99644784}

0.5

B



*Saccharomyces cerevisiae* virus La {NP_042580}
*Pichia membranifaciens* {v1.0 DOE Joint Genome Institute}*
*Scheffersomyces coipomoensis**
*Schizosaccharomyces pombe* {NC_003421}*
*Scheffersomyces segobiensis* RNA virus
Black raspberry virus F {ABU55398}
*Saccharomyces cerevisiae* virus LA {AAA50320}
*Nadsonia fulvescens* {ESTs DOE Joint Genome Institute}
*Tuber aestivum* virus 1 {ADQ54105}
*Candida parapsilosis**
*Scheffersomyces stipitis* {CP000501}*
*Penicillium marneffei* {ABAR01000142}*
*Debaryomyces hansenii* {CR382134}*
*Xanthophyllomyces dendrorhous* virus L2 {AFH09415}

0.5

# Figure 3

Evidence that the viral host, *Scheffersomyces segobiensis*, uses the modified genetic code of the "CTG" clade.

A) Midpoint-rooted maximum likelihood phylogram of CTG clade yeasts and related yeasts based on the protein sequences of the five most phylogenetically informative genes for fungi. Branches are labeled with support values from approximate likelihood ratio tests and nonparametric bootstrapping. Blue shading indicates standard code yeasts, while red shading indicates CTG code yeasts. Gray spheres indicate lineages with evidence of past or prior infections with totiviruses; B) Comparison of CUG codon positions in each taxon versus homologous amino acid residues in *Saccharyomyces cerevisiae*. Blue bars represent the percent of *S. cerevisiae* leucine residues that are coded by the CUG codon for each taxon, while red bars represent the percent of *S. cerevisiae* serine residues that are coded by CUG for each taxon. C) Boxplots showing relative synonymous codon usage (RSCU) for CUG codon usage by each taxon, the blue plot represents CUG RSCU values for non-CTG clade yeast, while the red plot represents CTG clade yeast including *S. segobiensis*. D) The secondary structure model of the CTG clade type of tRNA[SER] detected in *S. segobiensis*. The red shading indicates serine identifier sites, while blue shading indicates standard leucine identifier sites. The gray site is a typical guanine residue of the tRNA[SER] of CTG yeast that lowers the leucine amino-acylation efficiency.

A

B

C

D

*Candida glabrata*
*Saccharomyces mikatae*
*Saccharomyces paradoxus*
*Saccharomyces cerevisiae*
*Saccharomyces bayanus*
*Naumovozyma castellii*
*Kluyveromyces lactis*
*Ashbya gossypii*
*Komagataella pastoris*
*Ogataea angusta*
*Candida tenuis*
*Debaryomyces hansenii*
*Millerozyma farinosa*
*Candida tropicalis*
*Candida dubliniensis*
*Candida albicans*
*Lodderomyces elongisporus*
*Candida orthopsilosis*
*Candida parapsilosis*
*Spathaspora passalidarum*
*Scheffersomyces coipomoensis*
*Scheffersomyces stipitis*
*Scheffersomyces segobiensis*

CTG clade

1 / 100
1 / 100
1 / 100
1 / 100
1 / 100
0.79 / 79
1 / 100
1 / 99
1 / 100
1 / 100
1 / 99
1 / 100
1 / 100
1 / 100
1 / 100
1 / 100
1 / 100
1 /100
0.95 / 91
1 / 99
1 / 100

0.2

2.0
1.5
1.0
0.5
0.0

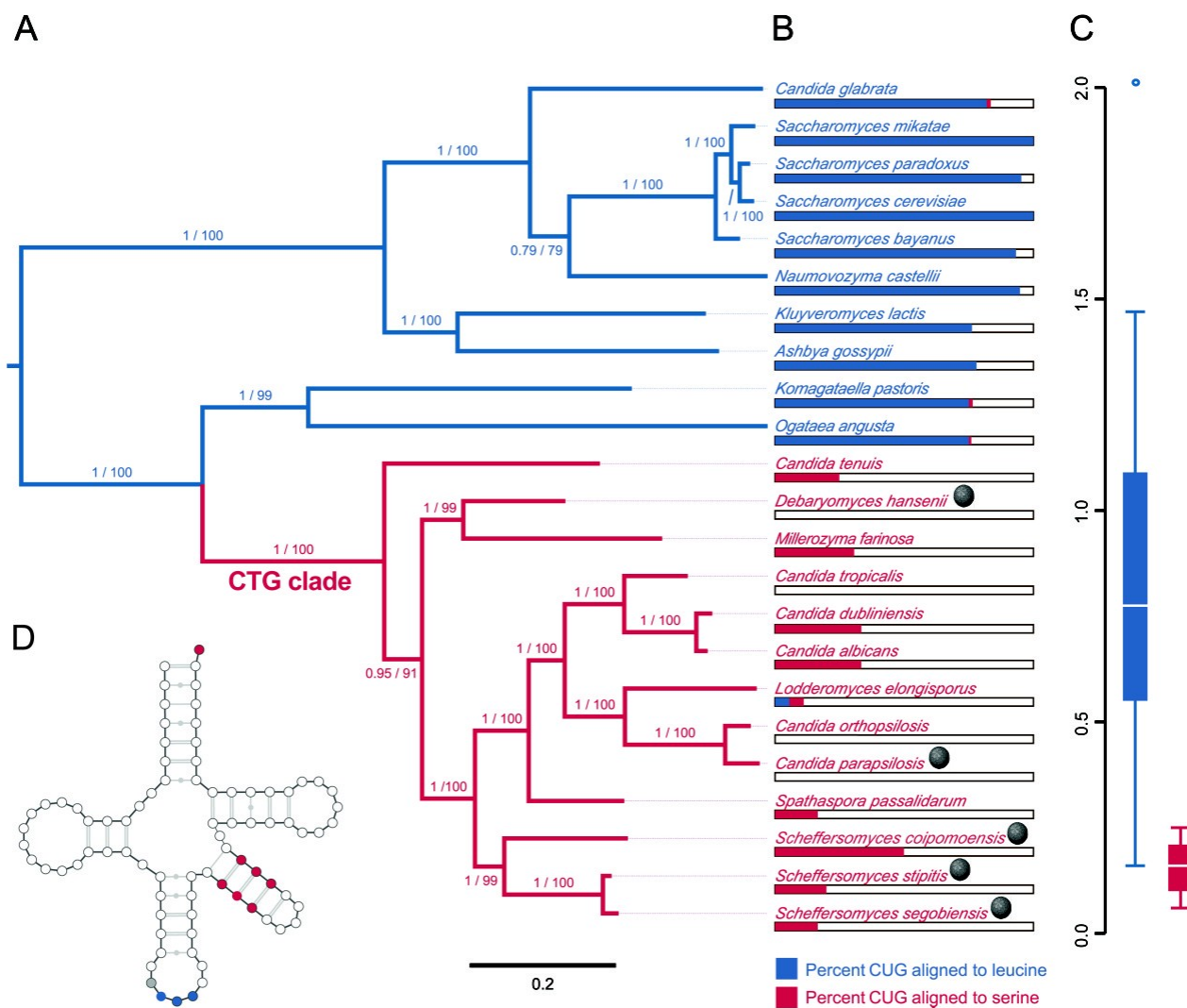■ Percent CUG aligned to leucine
■ Percent CUG aligned to serine

Fig. 3

# Figure 4

Bivariate plots of relative synonymous codon usage (RSCU) for serine and leucine versus third position base composition in yeast and their dsRNA viruses. CTG clade yeasts are shown as solid red points, and their viruses as hollow red points.

Standard code yeasts are shown as solid blue points and their viruses as hollow blue points. A) CUG is used by standard code yeasts and their viruses but avoided by CTG clade yeasts and *Scheffersomyces segobiensis virus L* . B) Leucine codon UUG is overused by CTG clade yeasts and *S. segobiensis virus L* relative to standard code yeasts and their viruses. C) Leucine codon CUA is used by standard code yeasts and their viruses, but avoided by CTG clade yeasts and *Scheffersomyces segobiensis virus L* . Comparison with %N3 suggests that these codon preferences are not attributable to nucleotide composition alone.