

1 **Contribution of recombination and selection to**
2 **molecular evolution of *Citrus tristeza virus***

3
4 Susana Martín ^{1,2†}, Adrián Sambade ^{1‡}, Luis Rubio ¹, María C. Vives ¹, Patricia
5 Moya ¹, José Guerri ¹, Santiago F. Elena ², Pedro Moreno ^{1*}

6
7 Authors Susana Martín and Adrián Sambade contributed equally to this work.

8
9 ¹ *Centro de Protección Vegetal y Biotecnología, Instituto Valenciano de*
10 *Investigaciones Agrarias, Moncada, 46113 Valencia, Spain*

11 ² *Instituto de Biología Molecular y Celular de Plantas (CSIC-UPV), 46022*
12 *Valencia, Spain*

13
14 Running title: *Citrus tristeza virus molecular evolution*

15
16 * Corresponding author: Instituto Valenciano de Investigaciones Agrarias, Ctra.
17 Moncada-Náquera Km. 4.5, Moncada, 46113 Valencia, Spain. Phone: 34-96-
18 3424069, Fax: 34-96-3424001. *E-mail address:* pmoreno@ivia.es

19
20 † Present address: Instituto de Biología Molecular y Celular de Plantas (CSIC-
21 UPV), Campus UPV CPI 8E, Ingeniero Fausto Elio s/n, 46022 Valencia, Spain

22
23 ‡ Present address: Department of Cell and Developmental Biology, John Innes
24 Centre, Norwich NR4 7UH, UK

25
26 SUMMARY: 212 words

27 MAIN TEXT AND FIGURE LEGENDS: 5496 words

28 NUMBER OF TABLES: 2

29 NUMBER OF FIGURES: 3

30

31 **Summary**

32

33 The genetic variation of *Citrus tristeza virus* (CTV) was analyzed comparing the
34 predominant sequence variants in seven genomic regions (p33, p65, p61, p18,
35 p13, p20, and p23) of 18 pathogenically distinct isolates from seven different
36 countries. Analyses of the selective constraints acting on each codon suggest
37 that most regions were under purifying selection. Phylogenetic analysis show
38 diverse patterns of molecular evolution for different genomic regions. A first
39 clade composed by isolates genetically close to the reference mild isolates
40 T385 or T30 was inferred from all genomic regions. A second clade, mostly
41 comprising virulent isolates, was defined from regions p33, p65, p13, p20, and
42 p23. For regions p65, p61, p18, p13, and p23 a third clade that mostly included
43 South American isolates could not be related with any reference genotype.
44 Phylogenetic relationships among isolates did not reflect their geographical
45 origin, suggesting significant gene flow between geographically distant areas.
46 Incongruent phylogenetic trees for different genomic regions suggested
47 recombination events, an extreme that was supported by several
48 recombination-detecting methods. A phylogenetic network incorporating the
49 effect of recombination showed an explosive radiation pattern for the evolution
50 of some isolates and grouped isolates by virulence. Taken together, the above
51 results suggest that negative selection, gene flow, sequence recombination,
52 and virulence may be important factors driving CTV evolution.

53

54

55

56 *Keywords:* CTV; genetic variability; maximum likelihood; recombination;
57 selective constraints; virus evolution

58

59 Introduction

60

61 *Citrus tristeza virus* (CTV) is a closterovirus, family *Closteroviridae*, with
62 two capsid proteins of 25 and 27 kDa, coating ~97 and ~3% of the virion length,
63 respectively (Febres *et al.*, 1996; Satyanarayana *et al.*, 2004). The single-
64 stranded and positive-sense CTV genomic RNA (gRNA) is about 20 kb in size
65 and contains 12 open reading frames (ORFs) potentially encoding at least 19
66 proteins (Karasev *et al.*, 1995). ORFs 1a and 1b, encoding replication-related
67 proteins, are translated from the gRNA, whereas the 10 3'-proximal ORFs,
68 encoding proteins p33, p6, p65, p61, p27, p25, p18, p13, p20, and p23, are
69 expressed via 3' co-terminal subgenomic RNAs (Hilf *et al.*, 1995). Protein p6
70 may operate as a membrane anchor (Satyanarayana *et al.*, 2000); proteins p65
71 (a homologue of the HSP70 heat shock proteins), p61 and the two coat proteins
72 are involved in virion assembly (Satyanarayana *et al.*, 2000); p20 accumulates
73 in amorphous inclusion bodies (Gowda *et al.*, 2000), and p23, an RNA-binding
74 protein (López *et al.*, 2000), controls asymmetrical accumulation of plus and
75 minus strands during RNA replication (Satyanarayana *et al.*, 2002) and is
76 involved in symptom expression (Ghorbel *et al.*, 2001, Fagoaga *et al.*, 2005).
77 Proteins p23, p20 and p25 act as RNA silencing suppressors (Lu *et al.*, 2004).
78 The functions of p33, p13 and p18 remain unknown.

79 CTV is primarily dispersed by propagation of infected buds, and then it is
80 locally spread aphids. CTV-induced symptoms include *i*) decline of citrus
81 species propagated on sour orange (*Citrus aurantium* L.) rootstock, *ii*) yellowing
82 and growth cessation of sour orange, lemon (*C. limon* (L.) Burn. f.) or grapefruit
83 (*C. paradisi* Macf.) seedlings (seedling yellows), or *iii*) stunting, stem pitting, and
84 poor yield of different citrus varieties regardless the rootstock used (Moreno *et al.*,
85 2008). The molecular mechanisms involved in symptom expression are still
86 unknown.

87 As for other RNA viruses, genetic variation has been observed in CTV
88 isolates resulting from the error-prone nature of RNA polymerases and selection
89 pressures (Domingo & Holland, 1994), superinfection of field trees with
90 divergent CTV variants (Rubio *et al.*, 2001), genetic drift after transmission to
91 new hosts (Albiach-Martí *et al.*, 2000a; D'Urso *et al.*, 2000; Ayllón *et al.*, 2006),
92 or recombination (Rubio *et al.*, 2001; Vives *et al.*, 2005). Characterization of the

93 genetic structure of viral populations and factors contributing to their evolution
94 may help understanding important features like the outbreak of new epidemics
95 or virulence changes in current isolates (Fernández-Cuartero *et al.*, 1994;
96 Escriu *et al.*, 2000). These studies have practical implications in virus control,
97 since durability of host resistance largely depends on genetic variability of the
98 virus (García-Arenal & McDonald, 2003).

99 Previously we compared the predominant sequence variants of gene *p23*
100 from 18 CTV isolates of different origins and pathogenicity characteristics
101 (Sambade *et al.*, 2003). Phylogenetic analyses showed that sequence variants
102 predominant in mild isolates (causing mild to moderate symptoms in Mexican
103 lime (*C. aurantifolia* (Christ.) Swing.) and sometimes decline of sweet orange
104 (*C. sinensis* (L.) Osb.) propagated on sour orange rootstock) and those
105 predominant in virulent isolates (additionally inducing seedling yellows and stem
106 pitting in grapefruit or sweet orange) clustered separately. To gain further
107 insight into the mechanisms of CTV evolution, we analyzed the genetic variation
108 and phylogenetic relationships in seven gRNA regions of these isolates and
109 sought for recombination between divergent sequence variants. Our analyses
110 showed variable selection pressures along the gRNA and frequent
111 recombination events. Apparently, CTV variants cluster within at least two
112 evolutionarily divergent lineages.

113

114 **Materials and Methods**

115

116 *Virus isolates*

117

118 The CTV isolates used in this study were from Argentina (C-268-2, C-269-
119 6 and C-270-3), Brazil (Barão B, Cald-CB, Galego 50 and Val-CB), France (K),
120 Florida (T36 and T55), Israel (VT), Japan (T388), and Spain (T32, T300, T305,
121 T312, T346, and T385), and their pathogenicity characteristics have been
122 described (Sambade *et al.*, 2002). These isolates were classified into 5
123 biogroups according to symptoms induced in a panel of indicator hosts
124 (Garnsey *et al.*, 2005). In short, the K isolate is asymptomatic in all hosts
125 (biogroup 0); isolates T32, T55 and T385, induce symptoms in Mexican lime
126 (biogroup 1), and T300, T312, and T346 also cause decline of sweet orange

127 grafted on sour orange rootstock (biogroup 2); T36, Galego 50 and C-268-2
128 additionally induce seedling yellows (biogroup 3); while the remaining isolates,
129 in addition to the latter symptoms, cause stem pitting on grapefruit (Barão B, C-
130 269-6, C-270-3, and VT; biogroup 4) or on grapefruit and sweet orange (T305,
131 T388, Cald-CB, and Val-CB; biogroup 5).

132

133 *cDNA synthesis, cloning and sequencing*

134

135 cDNA of regions located in genes *p33*, *p65*, *p61*, *p18*, *p13*, *p20*, and *p23*
136 of the CTV gRNA was synthesized by reverse transcription (RT) and PCR
137 amplification using double-stranded RNA (dsRNA)-rich preparations (Moreno *et*
138 *al.*, 1990) as template and appropriate primers (Supplementary data 1). Primers
139 amplify 41%-44% of *p33*, *p65* and *p61* genes, 76% of *p18* and 87-99% of *p13*,
140 *p20* and *p23*. RT-PCR was performed in a 25 µL reaction mix containing: 20
141 mM Tris-HCl, pH 8.4, 50 mM KCl, 500 µg/mL bovine serum albumin, 3 mM
142 MgCl₂, 4 mM each of dATP, dCTP, dGTP, and dTTP, 1 µM of each primer, 20
143 U of SuperScript™ II reverse transcriptase, 1 U of RNaseOut and 1 U of Taq
144 DNA polymerase (Invitrogen, USA). The reaction proceeded in an air thermal
145 cycler (Idaho Technologies, USA) using 30 min at 46 °C for RT, 2 min at 94 °C,
146 40 cycles of 5 s at 94 °C, 5 s at 55 °C and 30 s at 72 °C, and a final step of 2
147 min at 72 °C. The resulting RT-PCR products were cloned in the pGEM-T
148 vector (Promega, USA) (Sambrook *et al.*, 1989).

149 The sequence variants predominant in each CTV isolate were selected by
150 single-strand conformation polymorphism (SSCP) analysis (Rubio *et al.*, 2001).
151 For this purpose, ten clones from each cDNA product were PCR-amplified as
152 described above, and the DNA synthesized was SSCP-analyzed in the same
153 gel as the RT-PCR product from which the clones were obtained (Sambade *et*
154 *al.*, 2003). Clones whose DNA strands co-migrated with the most intense DNA
155 bands of the starting RT-PCR product were sequenced. The nucleotide
156 sequence of the cDNA clones selected was determined in both directions with
157 an ABI PRISM 3100 DNA sequence analyzer (Applied Biosystems). The
158 sequence of the virulent CTV isolate NUagA from Japan (AB046398) was used
159 for comparisons.

160 The nucleotide sequences obtained have been deposited in GenBank
161 under accession numbers FM955890- FM956002.

162

163 *Sequence analyses*

164

165 Nucleotide sequences were translated to proteins using GENEDOC and
166 multiple protein alignments were performed with MUSCLE program (Karl &
167 Hugh, 1997; Edgar, 2004). Nucleotide alignments were then obtained by
168 concatenating codons with the amino acid alignment as guide. Sites containing
169 insertions were removed from all subsequent analyses.

170 Nucleotide substitution models for different CTV regions were inferred
171 using the model selection tool available at the DATAMONKEY server
172 (<http://www.datamonkey.org>) of the HYPHY package (Kosakovsky-Pond &
173 Frost, 2005a). Genetic distances and substitution parameters were calculated
174 by maximum likelihood with the TREE-PUZZLE 5.2 program (Schmidt *et al.*,
175 2002) assuming that sites had heterogeneous substitution rates described by a
176 gamma distribution with eight classes. The best amino acid substitution model
177 (lowest *AIC* value among competing models) was inferred with PROTTEST
178 (Abascal *et al.*, 2005), available at <http://darwin.uvigo.es/software/prottest.html>.
179 Detection of codons under selection was done using the fixed effects maximum
180 likelihood (FEL) method of the HYPHY package (Kosakovsky-Pond & Frost,
181 2005b). Recombination was detected with the GARD program available at the
182 DATAMONKEY server using the HKY85 substitution model and a beta-gamma
183 distribution with four classes for rate variation. Further confirmation of
184 recombination events and identification of parental sequences were performed
185 with the RDP3 package (Martin *et al.*, 2005a) that incorporates the
186 recombination-detecting algorithms GENECONV (Padidam *et al.*, 1999),
187 BOOTSCAN (Salminen *et al.*, 1995; Martin *et al.*, 2005b), MAXCHI (Smith,
188 1992; Posada & Crandall, 2001), CHIMAERA (Posada & Crandall, 2001),
189 SISCAN (Gibbs *et al.*, 2000), 3SEQ (Boni *et al.*, 2007), and RDP (Martin &
190 Rybicki, 2000), using their default parameter values. Average nucleotide
191 distances of CTV isolates were calculated using MEGA 4.0 software (Tamura *et al.*,
192 2007), after testing homogeneity of pattern substitution among lineages.
193 Evolutionary distances were estimated by the composite maximum likelihood

194 method assuming that substitution rates among sites fitted a gamma
195 distribution.

196

197 *Phylogenetic analysis*

198

199 Protein maximum likelihood trees were inferred with PROTTEST (Abascal
200 *et al.*, 2005); and significance for the nodes was estimated with 1000 bootstrap
201 replicates using PHYML program (Guindon & Gascuel, 2003), available at
202 <http://phylemon.bioinfo.cipf.es/> (Tárraga *et al.*, 2007). Nucleotide maximum
203 likelihood trees were constructed by the sequential addition method using
204 HYPHY (Kosakovsky-Pond *et al.*, 2005), with the HKY85 substitution model and
205 a gamma distribution with six classes for rate heterogeneity. Topology
206 comparisons were performed with Shimodaira & Hasegawa (1999) test
207 implemented in the DNAML program of the PHYLIP 3.67 package (Felsenstein,
208 2005). Phylogenetic trees were drawn using MEGA 4.0 (Tamura *et al.*, 2007).
209 The ratio of non-synonymous to synonymous substitution rates in different
210 branches of maximum likelihood trees was estimated using the codon-based
211 genetic algorithm implemented in the GA-BRANCH program (Kosakovsky-Pond
212 & Frost, 2005c) available at the DATAMONKEY server. Split-decomposition
213 analysis of concatenated CTV sequences was performed using the
214 SPLITSTREE program with default parameters (Huson, 1998).

215

216 **Results**

217

218 *Phylogenetic relationships between CTV isolates*

219

220 A predominant sequence variant was observed for each isolate and gRNA
221 region, except for isolates Galego 50, that showed two p13 variants, and C-268-
222 2 that had three variants in p33 and p61. Deduced amino acid sequences were
223 used to infer phylogenetic trees, including the NUagA sequence in these
224 regions for comparison (Fig. 1). At least two clades were observed in most
225 regions: one of them (clade I), comprising the sequence variants of mild isolates
226 T32, T55, T300, T312, and T385, and T346 in p20, was supported by bootstrap
227 values of 86.3% (p33), 67.5% (p65), 96.1% (p61), 65.5% (p18), 75.9% (p13),

228 99.5% (p20) and 94.4% (p23). A second clade (clade II) enclosing biogroup 5
229 plus a variable set of isolates of biogroup 4, was supported by bootstrap values
230 of 96.8% (p33), 87.2% (p23) 76.9% (p13) and 56.2% (p65), albeit in regions
231 p33 and p65 a biogroup 3 isolate (Galego 50) was in the same cluster .
232 Although support for clade II was less robust than for clade I, the virulent
233 isolates T388, T305 and NUagA (NUagA type group) were closely related to
234 each other and distantly related to clade I in all regions, forming a stable
235 nucleus within clade II. Thus, phylogenetic relationships reflected to some
236 degree pathogenicity characteristics of the isolates. In p18 all South American
237 isolates, except Galego 50 grouped together into a distinct clade (clade III).

238 While isolates defining clade I and NUagA type isolates clustered together
239 regardless the genomic region, other isolates showed incongruent phylogenetic
240 relationships for different regions (Fig. 1). The virulent isolates Cald-CB, Barão
241 B and Val-CB clustered together, and closely related to NUagA type isolates,
242 only in p33, p65 and p23. In p33, p65 and p61, isolate K was genetically closer
243 to clade I than to NUagA, but in other regions it was located between them. The
244 two major p13 variants found in Galego 50 were divergent, with the variant
245 Galego 50A grouping with isolates C-270-3 and Barão B and the variant Galego
246 50B being closer to NUaGA group. Similarly, for isolate C-268-2, one of the
247 three p33 variants was close to NUaGA and the other two were separated from
248 clade I and NUaGA, and two of the three p61 variants were close to clade I, and
249 the third was separated from both groups (Fig. 1).

250 The incongruent phylogenetic relationships observed for some isolates in
251 different genomic regions suggested that their gRNA might have originated from
252 recombination events between diverged sequences.

253

254 *Genetic variation and selective pressures in different genomic regions*

255

256 The average nucleotide distance for p33 (0.1641 ± 0.0196) (Table 1) was
257 significantly higher than those of regions p65, p18, p13 p20, and p23 (ranging
258 from 0.0741 ± 0.0240 to 0.1145 ± 0.0188), and average distance for p13
259 (0.0741 ± 0.0240) was lower than those of p33 and p61 (Model II ANOVA:
260 $F_{6,131}=48.4279$, $p < 0.0001$ and Tukey-Kramer *post hoc* test at 95% confidence
261 level). To evaluate the selective constraints operating in each region, codons

262 under selection were detected using the FEL method (Supplementary data 2).
263 Since d_N and d_S estimates are sensitive to the effect of recombination, we
264 preliminarily screened the different genomic regions with the GARD tool and
265 found significant recombination signals in *p61*, *p20* and *p23* genes (positions
266 283, 252 and 239, respectively). Alignments for these three genes were split in
267 the corresponding non-recombinant regions. The number of negatively selected
268 sites and the mean normalized d_N-d_S values for each genomic region, after
269 correcting the significance p values for multiple comparisons of the same null
270 hypothesis using a false discovery rate (FDR) of 5%, are in Table 1. A total of
271 150 codons were subjected to significant purifying selection, with ratios of
272 negatively selected codons ranging from 10% for *p20* to 15.87% for *p33*.
273 Although these ratios were not significantly different ($\chi^2=5.145$, 6 d.f., $p=0.525$),
274 the strength of negative selection estimated by the mean normalized d_N-d_S
275 value did differ among genomic regions (Kruskal-Wallis test: $H=53.4090$, 6 d.f.,
276 $p<0.0001$), this difference being entirely driven by the less negative d_N-d_S
277 estimated for *p33* (-2.1098 ± 0.2505) relative to the average value for the other
278 six regions (-7.6334 ± 0.4393) (Dunn's *post hoc* test $p<0.05$). Only codons 203
279 ($d_N-d_S=3.6789$) and 244 ($d_N-d_S=6.6412$) in *p61* showed a significant signature
280 of positive selection.

281

282 *Branch-specific analysis of the ratio of non-synonymous to synonymous* 283 *substitution rates in CTV phylogenies*

284

285 To get deeper insights into the selective pressures acting at the protein
286 level during the CTV evolution, branch-specific ratios of d_N to d_S rates (ω) were
287 estimated using a genetic algorithm that identified models with variable
288 numbers of ω categories per lineage that fitted better to data than the single-
289 ratio (all lineages evolving with equal ω) or the fully saturated (each lineage
290 evolves at different ω) models (Table 2). Although ω values and the proportion
291 of associated branches varied in different regions and periods of CTV
292 diversification, most classes had $\omega<1$ and most branches were assigned to
293 these classes (100% in *p61* and over 74% in the other regions). Figure 2b
294 illustrates a simplified comparison of branch-specific ω values in phylogenetic

295 trees grouping these values in four selection categories: *i*) strong negative
296 selection ($\omega < 0.1$); *ii*) moderate negative ($0.1 < \omega < 0.4$); *iii*) weak negative
297 ($0.4 < \omega < 1$); and *iv*) positive selection ($\omega > 1$). In most cases, internal branches
298 connecting the groups of mild and virulent isolates, particularly clade I and
299 NUagA group were associated to $\omega < 1$, indicating that divergence of these
300 genotypes occurred under negative selection pressure. In p65 purifying
301 selection was strong for all internal branches, whereas in p33, p18 and p13,
302 periods of strong and moderate selection alternated. While in p33 negative
303 selection was moderate for lineages leading to mild and virulent groups and
304 strong in the internal branch leading to the out-groups of C-270-3 and the
305 ancestor of T36 and C-268-2B, the opposite was true in p18, with strong
306 selection for lineages leading to the mild and virulent groups and moderate in
307 the lineage leading to C-270-3 and other South American isolates. In p61, p20
308 and p23 selection was moderate for lineages leading to virulent or mild groups
309 and to C-270-3, and in p23 it was weak or even positive after diversification of
310 the cluster formed by K and clade I. Positive selection occurred during limited
311 periods in all regions but p61, and except for p13, it was observed only in
312 terminal branches, with highest frequency being observed in p33. For some
313 isolates signature of positive or weak negative selection was detected in most
314 regions, i.e., T385 (p33, p61, p13, p20, and p23), T312 (p33, p65, p61, p18,
315 and p23) or VT (p33, p65, p61, and p18).

316

317 *Frequent recombination events in CTV genomes*

318

319 Sequences were first examined for recombination using the GARD tool
320 that identifies recombination breakpoints when the likelihood of phylogenetic
321 trees inferred for the partitioned alignments is significantly higher than that
322 obtained for the non-partitioned alignment. Due to computational limits and to
323 avoid arbitrary assembling of C-268-2 sequence variants, alignments of pair-
324 wise concatenated regions corresponding to adjacent genes were used as
325 input. A total of nine recombination breakpoints were detected: six of them
326 located in the boundaries of different regions, suggesting recombination events
327 somewhere between the analyzed fragments, and three located within p61, p20

328 and p23 regions (positions 283, 252 and 239 of the corresponding region).
329 Since GARD does not require different topologies, a model containing partitions
330 could outperform a non partitioned one if both share the same topology due to
331 best fit of branch lengths. To test for topological differences, maximum
332 likelihood trees for the non recombinant regions defined by GARD (HKY85
333 substitution model and sequential addition method for topology inference) were
334 compared using the Shimodaira & Hasegawa (1999) test that compares the
335 goodness of a set of competing phylogenetic trees to describe the evolution of a
336 given alignment. Maximum likelihood estimates of the substitution parameters
337 (transversion-transition rates ratio; shape parameter of the gamma distribution
338 of substitution rates per site, and relative substitution rates) were inferred using
339 TREE-PUZZLE. This analysis confirmed the differences in tree topology
340 inferred for different genomic regions, although in p61, p20 and p23 the trees
341 inferred from the entire region were not significantly worse than those inferred
342 from each partition defined by GARD, suggesting that partitions established by
343 GARD were mostly due to differences in branch length ($p < 0.05$).

344 To assess the frequency and extension of recombination during CTV
345 diversification, the seven genomic regions were concatenated and
346 recombination events were identified using the RDP3 program that implements
347 several recombination-detecting methods (Martin *et al.*, 2005a), using default
348 setting parameters for the subset of fast detection methods GENECONV,
349 BOOTSCAN, MAXCHI, CHIMAERA, SISCAN, 3SEQ and RDP. Isolate C-268-2
350 was excluded to avoid arbitrary assembling of its p33 and p61 variants. A total
351 of 14 recombination events were detected by at least one of the methods, but
352 only those predicted by at least four different methods were accepted (Fig. 2c)
353 and assignment of parental and daughter sequences was confirmed by
354 constructing maximum likelihood trees (Fig. 2b). For example, isolate C-269-6
355 grouped with NUagA group in regions p33, p65, p23 and p20, but with C-270-3
356 in p61 and p18 and as outgroup in p13 (Fig. 2b). A recombination involving a
357 NUagA ancestor as major parental, and the p61 and p18 regions from a C-270-
358 3 ancestor, was identified by the seven methods used (Fig. 2c). Recombination
359 events involving the same parentals were also detected for Cald-CB and Val-
360 CB in p18, Barão B in p18 and p13 and Galego 50A in p13. Isolates Barão B
361 and Val-CB, that grouped together and distant from other isolates in p20,

362 presented in this region an additional recombination between a NUagA-type
363 and an unknown ancestor not represented in the alignment. Galego 50 was
364 closely related to VT in most regions (p33, p65, p61, p18, and p20), but not in
365 p23 and p13. In this latter region the variant Galego 50B grouped with VT and
366 the variant Galego 50A with C-270-3, as a result of a recombination between
367 VT and C-270-3 ancestors (Fig. 2b and 2c). In p23 Galego 50 was close to
368 isolates C-270-3 and C-268-2, and RDP3 predicted a recombination between a
369 Cald-CB ancestor as major parental, and an unknown isolate providing the 3'
370 end of p20 and p23. The non assignment of C-270-3 as minor parental was due
371 to the use of UPGMA trees in RDP3 default analysis (not shown), but inspection
372 of maximum likelihood trees indicated that likely VT and C-270-3 are the major
373 and minor parentals, respectively (Fig. 2b and 2c).

374 Isolate T346 was close to clade I in regions p33, p61 and p23, to C-270-3
375 in p65, and in intermediate positions in p20, p18 and p13, an incongruence that
376 was compatible with a recombination event in the p65 region between T312 and
377 C-270-3 ancestors as major and minor parental, respectively. The phylogenetic
378 relationships of isolates T346 and T36 widely varied among regions: while in
379 p33, p65 and p23 both isolates were divergent (genetic distances from 0.1180
380 to 0.2447), in p61, p18 and p13 they were closely related (genetic distances
381 from 0.0090 to 0.0535); and two recombination events involving regions p61
382 and p13 were detected between ancestors of these isolates. In p13 and p18
383 T346 and K clustered in the same group. These results and the genetic
384 distances observed in these regions are compatible with a recombination
385 between K and T36 in p18 and p13 and later recombination between T346 and
386 T36 in p13, and with other possibly older recombination between T346 and T36
387 in p61.

388 Finally, isolate C-268-2 contained three major diverged variants in p33 and
389 p61 but it was monomorphic for the other regions. Furthermore, it clustered with
390 C-270-3 in p65, p18, p20 and p23 regions, but close to clade I in p13,
391 suggesting that this isolate may be the result of multiple recombination events.

392

393 *A phylogenetic network for CTV isolates*

394

395 Due to the recombinant nature of CTV genomes, bifurcating phylogenetic
396 trees do not properly reflect the actual evolutionary history of different isolates,
397 since one isolate may be directly linked to more than one ancestral sequence.
398 To provide a more accurate representation of those relationships a phylogenetic
399 network was constructed from the concatenated alignment of the seven regions
400 by the split-decomposition method implemented in SPLITSTREE (Huson,
401 1998). Again isolate C-268-2 was excluded for the reasons given above (Fig. 3).

402 The largest splits divided CTV isolates into two groups: one formed by
403 isolates of biogroups 0, 1 and 2 (with the exception of T346) and the other by
404 isolates of biogroups 3, 4 and 5 (with the exception of T36). Within the mild
405 group, isolates T55 from Florida, and T312, T300 and the ancestor of T32 and
406 T385 from Spain showed a radiation pattern. In the second group including
407 isolates from South America, Japan, Israel and Spain, Cald-CB and the
408 ancestors of Barão B and Val-CB, and NUagA group, diverged from a common
409 ancestor in a star-like manner, whereas the other isolates had a more complex
410 phylogeny. Isolate C-269-6 was connected to the node joining isolates of
411 biogroup 5 and to the ancestor of C-270-3, consistent with its recombinant
412 nature (between NUagA and C-270-3) revealed by RDP3 analysis. Galego 50A
413 was connected to Galego 50B and to the common ancestor of both variants and
414 VT, also in agreement with recombination analysis. C-270-3 was connected to
415 the T346 and C-269-6 ancestors giving further support to previous finding that
416 these isolates likely arose from a recombination between a mild (T346) or a
417 severe (C-269-6) major ancestor and C-270-3. Finally, the ancestors of isolates
418 K, T36 and T346, that showed variable phylogenetic relationships (Fig. 2b),
419 were interconnected in the network in a complex pattern (Fig. 3), supporting
420 recombination between them detected by RDP3 (Fig. 2c).

421

422 **DISCUSSION**

423

424 The genetic variation and evolutionary factors shaping CTV populations
425 were studied comparing the predominant sequence variants in seven genomic
426 regions of 18 isolates from different geographical origin and pathogenicity
427 characteristics. It was assumed that pathogenicity would be largely associated
428 to the major sequence variant since: 1) virions obtained from a cDNA clone of

429 the major component of isolate T36 induced the symptoms characteristics of
430 this isolate (Satyanarayana *et al.*, 1999, 2001), and *ii*) in citrus plants
431 successively co-inoculated with a mild and a virulent CTV isolate, symptom
432 onset was associated with predominance of the sequence variant characteristic
433 of the virulent isolate (Sambade *et al.*, 2002, 2007).

434 Analysis of selective pressures acting on different codons suggests that all
435 regions examined are mostly subjected to purifying selection with only two
436 codons in p61 being positively selected. Purifying selection measured as
437 normalized d_N-d_S showed similar proportion of selected sites and selection
438 intensity among regions, except for p33 that had less intense selection. Less
439 negative d_N-d_S value of p33 indicates that in this region more nonsynonymous
440 substitutions are allowed, coherently this region also had higher evolutionary
441 distances. Although net selection pressure was similar in the other regions,
442 branch-specific analysis showed variable strength of selection depending on the
443 genomic region and period of CTV diversification, i.e., this pressure was strong
444 in the diversification of clade I and NUagA type isolates in p65 but moderate or
445 weak in p23.

446 Data available on the functional domains of CTV proteins are still limited,
447 and the function of proteins p33, p13 and p18, that are dispensable for CTV
448 infection and movement (Tatineni *et al.*, 2008), is unknown. Albeit selective
449 pressures are less intense in p33, the fraction of selected sites is similar in all
450 regions, indicating selective constraints to amino acid changes and providing
451 candidate positions to test in functional studies. Negative selection was
452 expected in p65, p61, p20, and p23, considering their role in the CTV biology.
453 Genes *p65* and *p61* are part of a conserved five-gene block encoding proteins
454 involved in virion assembly and movement (Satyanarayana *et al.*, 2000; Dolja *et al.*,
455 2006). The p65 region analyzed here encodes 5 of the 8 motifs conserved
456 among HSP70 proteins (Pappu *et al.*, 1994). Proteins p20 and p23 act as
457 silencing suppressors (Lu *et al.*, 2004). Within the p20 region, amino acids I38,
458 Y113, R130, L137, S141, and L159 are strictly conserved among silencing
459 suppressors of closteroviruses (Reed *et al.*, 2003). Gene *p23*, that was
460 completely sequenced, contains the RNA-binding domain and putative zinc
461 finger required for asymmetrical accumulation of positive and negative RNA
462 strands, with conserved residues C68, C71, H75, and C85 being involved in this

463 activity (López *et al.*, 2000; Satyanarayana *et al.*, 2002). These amino acids
464 were encoded by invariable codons in all CTV isolates, with the exception of
465 L137 in p20 and C71 in p23, that were found under significant purifying
466 selection using cut-off values of $p=0.5$ and $p=0.1$, respectively. The observation
467 that protein p23 is a pathogenicity determinant in citrus (Ghorbel *et al.*, 2001;
468 Fagoaga *et al.*, 2005) is consistent with separation of the mild and virulent CTV
469 isolates in phylogenetic analysis.

470 Phylogenetic analysis showed that the mild isolates T32, T55, T300, T312,
471 and T385 form a clade (clade I), supported by bootstrap values >70% in six
472 regions, thus defining a CTV lineage which also includes isolate T30 from
473 Florida and probably others from Colombia and Taiwan (Ruiz-Ruiz *et al.*, 2006;
474 Albiach-Martí *et al.*, 2000b). A second clade (clade II) including a variable set of
475 virulent isolates was supported by bootstrap values >75% in three regions,
476 suggesting a higher recombination frequency for those isolates. Within clade II,
477 isolates T305, T388 and NUagA were closely related to each other and distantly
478 related to clade I in all regions, and together with other virulent isolates showed
479 a star-like evolution pattern from a common ancestor in the phylogenetic
480 network. These data suggest that virulent isolates could represent a second
481 CTV lineage distantly related to clade I that would also include isolate T318A
482 from Spain (Ruiz-Ruiz *et al.*, 2006). A third clade including Brazilian and
483 Argentinean isolates genetically related with isolate C-270-3 was observed in
484 several genomic regions. Sequence comparisons showed that several clones of
485 two Colombian isolates released on GeneBank were closely related to C-270-3,
486 while others were closely related to severe isolates in p23. However, the latter
487 were related to clade I or to severe isolates but not to C-270-3 in p33 (not
488 shown). Although no complete gRNA sequence from South American isolates is
489 currently available, these findings are compatible with the diversification of a
490 third CTV lineage and frequent recombinations in this area. All lineages
491 included closely related CTV variants from distant locations, a circumstance
492 that, together with the radiation pattern observed for diversification of some CTV
493 isolates, provide further support to previous suggestion that genetic flow has
494 likely occurred (Rubio *et al.*, 2001).

495 Most isolates showed incongruent phylogenetic relationships in different
496 regions, suggesting frequent recombination events, a possibility that was

497 supported by recombination-detecting methods and by a split-decomposition
498 phylogenetic network. Nine of the 19 isolates compared were recombinant,
499 particularly, most isolates from Brazil (between ancestors of NUagA or VT and
500 C-270-3) and Argentina (C-268-2 appears as a mosaic of mild, severe and C-
501 270-3-type sequence variants), in agreement with results obtained for these and
502 other Argentinean isolates (Iglesias *et al.*, 2008). Recombination involving *p18*
503 was so frequent that all South American isolates except Galego 50 formed in
504 this region a monophyletic group diverged from the other isolates. A similar
505 grouping was described for genes *p25* and *p27*, located upstream *p18*, for
506 Argentinean isolates (Iglesias *et al.*, 2008). The variable position of K, T346 and
507 T36 in phylogenetic trees, their lack of association with other isolates and the
508 network topology suggest that they might represent CTV lineages with a more
509 complex history involving recombination events among their ancestors and
510 possibly with genotypes unrelated with those analyzed here.

511 The phylogenetic network grouped CTV isolates in two major clusters
512 separated by long splits: one of them comprising isolates of biogroups 0, 1 and
513 2 (except for isolate T346) that induce mild to moderate symptoms in the most
514 sensitive hosts and cause symptomless infections in grapefruit and sweet
515 orange seedlings, and the other, including isolates of biogroups 3, 4 and 5
516 (except for isolate T36) that incite severe symptoms in sensitive hosts and are
517 also pathogenic on grapefruit or sweet orange. Sequence separation between
518 mild and virulent CTV isolates is consistent with a different host response after
519 infection, as indicated by specific changes induced in the citrus transcriptome
520 by both types of isolates (Gandía *et al.*, 2007). This separation suggests that
521 virulence might be an important evolutionary factor shaping CTV populations.

522 Homologous recombination seems a common process in some plant RNA
523 viruses, particularly potyviruses (Chare & Holmes, 2005) and bromoviruses
524 (Codoñer & Elena, 2008). It has been postulated that recombination might
525 prevent accumulation of deleterious mutations in small populations and/or allow
526 a faster adaptation to changing environments (Lai, 1992; Roossinck, 1997;
527 García-Arenal *et al.*, 2001). On the other hand, simulation studies showed that
528 recombination in RNA viruses is more likely to create combinations of
529 deleterious mutations than purge them from genomes, thus causing fitness
530 reductions (Holmes, 2003). Our results and previous analyses (Hilf *et al.*, 2009)

531 suggest that RNA recombination is a major factor in CTV variation and likely
532 plays a role in its evolution. Potential factors contributing to frequent
533 recombination in CTV include: *i*) dispersal of diverged CTV genotypes in the
534 same area by movement of infected buds, *ii*) the long life of citrus trees
535 providing many opportunities for repeated infections with diverged sequence
536 variants (Rubio *et al.*, 2001; Weng *et al.*, 2007; Gomes *et al.*, 2008), *iii*) the
537 presence in infected cells of multiple viral RNA species produced during CTV
538 replication that likely facilitates recombination events (Hilf *et al.*, 1995; Yang *et*
539 *al.*, 1997; Ayllón *et al.*, 1999; Gowda *et al.*, 2003), and *iv*) the large size of CTV
540 genome that may accumulate 2-3 mutations per genome and replication round.
541 Recombination might help maintaining functional genomes even if many
542 nonfunctional recombinants were produced (Allison *et al.*, 1990).

543 Summarizing, analyses of genetic variation in the 3' half of CTV genome
544 suggest that at least two different lineages might have evolved, and that
545 selective pressure against amino acid changes, gene flow, homologous
546 recombination and perhaps virulence, may be important factors in CTV
547 evolution and in shaping CTV populations.

548

549 **Acknowledgements**

550

551 We are indebted to S. M. Garnsey (University of Florida-C.R.E.C., Lake
552 Alfred) for kindly providing freeze-dried citrus tissue infected with isolates T36,
553 T55, K, and VT from the international collection of exotic citrus pathogens
554 maintained at the quarantine facilities of the USDA in the Beltsville Agricultural
555 Research Center (Maryland, USA), and to S. Gago-Zachert (Universidad de La
556 Plata, Argentina) and M. A. Machado (Centro APTA Citros "Sylvio Moreira",
557 Cordeirópolis, SP, Brazil) for providing dsRNA-rich preparations of the
558 Argentinean and Brazilian isolates, respectively. We are also thankful to M. E.
559 Martínez and M. Boil for technical assistance in the laboratory and to J. Piquer
560 for excellent care of plants. S. Martín and A. Sambade were recipient of
561 fellowships from the Spanish Ministerio de Ciencia e Innovación and Generalitat
562 Valenciana, respectively. This work was supported in part by grants AGL2004-
563 05099/AGR and AGL2007-61885/AGR (work at IVIA) and BFU2006-14819-
564 C02-01/BMC (work at IBMCP) from the Ministerio de Ciencia e Innovación.

565

566 **Reference List**

567

568 **Abascal , F., Zardoya R. & Posada, D. (2005).** ProtTest: selection of best-fit
569 models of protein evolution. *Bioinformatics* **21**, 2104-2105.

570 **Albiach-Martí, M. R., Guerri J., Hermoso de Mendoza, A., Laigret F.,**
571 **Ballester-Olmos, J. F. & Moreno P. (2000a).** Aphid transmission alters the
572 genomic and defective RNA populations of *Citrus tristeza virus* isolates.
573 *Phytopathology* **90**, 134-138.

574 **Albiach-Martí, M. R., Mawassi, M., Gowda, S., Satyanarayana, T., Hilf, S.**
575 **Shanker, M. E., Almira, E. C., Vives, M. C., López, C. & other authors**
576 **(2000b).** Sequences of *Citrus tristeza virus* separated in time and space are
577 essentially identical. *J Virol* **74**, 6856-6865.

578 **Allison, R., Thompson, C. & Ahlquist, P. (1990).** Regeneration of a functional
579 RNA virus genome by recombination between deletion mutants and
580 requirement for *Cowpea chlorotic mottle virus* 3a and coat genes for systemic
581 infection. *Proc Natl Acad Sci U S A* **87**, 1820-1824.

582 **Ayllón, M. A., López, C., Navas-Castillo, J., Mawassi, M., Dawson, W. O.,**
583 **Guerri, J., Flores, R. & Moreno, P. (1999).** New defective RNAs from *Citrus*
584 *tristeza virus*: evidence for a replicase-driven template switching mechanism in
585 their generation. *J Gen Virol* **80**, 817-821.

586 **Ayllón, M. A., Rubio, L., Sentandreu, V., Moya, A., Guerri, J. & Moreno, P.**
587 **(2006).** Variations in two gene sequences of *Citrus tristeza virus* after host
588 passage. *Virus Genes* **32**, 119-128.

589 **Boni, M. F., Posada, D. & Feldman, M. W. (2007).** An exact nonparametric
590 method for inferring mosaic structure in sequence triplets. *Genetics* **176**, 1035-
591 1047.

592 **Chare, E. R. & E. C. Holmes. (2005).** A phylogenetic survey of recombination
593 frequency in plant RNA viruses. *Arch Virol* **151**, 933-946.

594 **Codoñer, F. M. & Elena, S. F. (2008).** The promiscuous evolutionary history of
595 the *Bromoviridae* family. *J Gen Virol* **89**, 1739-1747.

596 **Dolja, V. V., Kreuze, J. F. & Valkonen, J.P.T. (2006).** Comparative and
597 functional genomics of closteroviruses. *Virus Res* **117**, 38–51.

598 **Domingo, E. & Holland, J. J. (1994).** Mutation rates and rapid evolution of
599 RNA viruses. In *The evolutionary biology of viruses*, pp. 161-184. Edited by S.
600 S. Morse. Raven Press. USA : New York.

601 **D'Urso, F., Ayllón, M. A., Rubio, L., Sambade, A., Hermoso de Mendoza, A.,**
602 **Guerri, J. & Moreno, P. (2000).** Contribution of uneven distribution of genomic
603 RNA variants of *Citrus tristeza virus* (CTV) within the plant to changes in the
604 viral population following aphid transmission. *Plant Pathol* **49**, 288-294.

605 **Edgar, R. C. (2004).** MUSCLE: a multiple sequence alignment method with
606 reduced time and space complexity. *Bioinformatics* **5**, 113.

607 **Escriu, F., Fraile, A. & García-Arenal, F. (2000).** Evolution of virulence in
608 natural populations of the satellite RNA of *Cucumber mosaic virus*.
609 *Phytopathology* **90**, 480-485.

610 **Fagoaga, C., López, C., Moreno, P., Navarro, L., Flores, R. & Peña, L.**
611 **(2005).** Viral-like symptoms induced by the ectopic expression of the p23 gene
612 of *Citrus tristeza virus* are citrus specific and do not correlate with the
613 pathogenicity of the virus strain. *Mol Plant-Microbe Interact* **18**, 435-445.

614 **Febres, V. J., Ashoulin, L., Mawassi, M., Frank, A., Bar-Joseph, M.,**
615 **Manjunath, K. L., Lee, R. F. & Niblett, C. L. (1996).** The p27 protein is present
616 at one end of *Citrus tristeza virus* particles. *Phytopathology* **86**, 1331-1335.

617 **Felsenstein, J., 2005.** PHYLIP (Phylogeny Inference Package) version 3.6.
618 Distributed by the author. Department of Genome Sciences, University of
619 Washington, Seattle.

620 **Fernández-Cuartero, B., Burgyan, J., Aranda, M. A., Salanki, K., Moriones,**
621 **E. & García-Arenal, F. (1994).** Increase in the relative fitness of a plant virus
622 RNA associated with its recombinant nature. *Virology* **203**, 373-377.

623 **Gandía, M., Conesa, A., Ancillo, G., Gadea, J., Forment, J., Pallás, V.,**
624 **Flores, R., Duran-Vila, N., Moreno, P. & Guerri, J. (2007).** Transcriptional
625 response of *Citrus aurantifolia* to infection by *Citrus tristeza virus*. *Virology* **367**,
626 298–306.

627 **García-Arenal, F., Fraile, A. & Malpica, J. M. (2001).** Variability and genetic
628 structure of plant virus populations. *Annu Rev Phytopathol* **39**, 157-186.

629 **García-Arenal, F. & McDonald, B. A. (2003).** An analysis of the durability of
630 resistance to plant viruses. *Phytopathology* **93**, 941-952.

631 **Garnsey, S. M., Civerolo, E. L., Gumpf, D. J., Paul, C., Hilf, M. E., Lee, R. F.,**
632 **Brlansky, R. H., Yokomi, R. K., & Hartung, J. S. (2005).** Biological
633 characterization of an international collection of *Citrus tristeza virus* (CTV)
634 isolates. In *Proceedings of the 16th Conference of the International*
635 *Organization of Citrus Virologists*, pp. 75-93, Edited by M. E. Hilf, N. Duran-Vila
636 & M. A. Rocha-Peña. Riverside, CA: IOCV. <http://www.ivia.es/iocv/>
637 **Ghorbel, R., López, C., Fagoaga, C., Moreno, P., Navarro, L., Flores, R. &**
638 **Peña, L. (2001).** Transgenic citrus plants expressing the *Citrus tristeza virus*
639 p23 protein exhibit viral-like symptoms. *Mol Plant Pathol* **2**, 27-36.
640 **Gomes, C. P., Nagata, T., de Jesus, W. C., Neto, C. R., Pappas, G. J. &**
641 **Martin, D.P. (2008).** Genetic variation and recombination of RdRp and HSP 70h
642 genes of *Citrus tristeza virus* isolates from orange trees showing symptoms of
643 citrus sudden death disease. *Virology* **16**, 5-9.
644 **Gibbs, M. J., Armstrong, J. S. & Gibbs, A. J. (2000).** Sister-Scanning: a
645 Monte Carlo procedure for assessing signals in recombinant sequences.
646 *Bioinformatics* **16**, 573-582.
647 **Gowda, S., Satyanarayana, T., Davis, C. L., Navas-Castillo, J., Albiach-**
648 **Martí, M. R., Mawassi, M., Valkov, N., Bar-Joseph, M., Moreno, P. &**
649 **Dawson, W. O. (2000).** The p20 gene product of *Citrus tristeza virus*
650 accumulates in the amorphous inclusion bodies. *Virology* **274**, 246-254.
651 **Gowda, S., Ayllón, M. A., Satyanarayana, T., Bar-Joseph, M. & Dawson W.**
652 **O. (2003).** Transcription strategy in a Closterovirus: a novel 5'-proximal
653 controller element of *Citrus tristeza virus* produces 5'- and 3'- terminal
654 subgenomic RNAs and differs from 3' open reading frame controller elements. *J*
655 *Virology* **77**, 340-352.
656 **Guindon, S. & Gascuel, O. (2003).** A simple, fast, and accurate algorithm to
657 estimate large phylogenies by maximum likelihood. *Syst Biol* **52**, 696-704.
658 **Hilf, M. E. (2009).** Two distinct evolutionary pathways for *Citrus tristeza virus*:
659 Recombination defines two gene modules and provides for increased genetic
660 diversity in a narrow host range plant virus. In *Proceedings of the 17th*
661 *Conference of the International Organization of Citrus Virologists* (in press).
662 **Hilf, M. E., Karasev, A. V., Pappu, H. R., Gumpf, D. J., Niblett, C. L.,**
663 **Garnsey, S. M. (1995).** Characterization of *Citrus tristeza virus* subgenomic
664 RNAs in infected tissue. *Virology* **208**, 576-582.

665 **Holmes, E. C. (2003).** Error thresholds and the constraints to RNA virus
666 evolution. *Trends Microbiol* **11**, 543-546.

667 **Huson, D. H. (1998).** Splits Tree: analyzing and visualizing evolutionary data.
668 *Bioinformatics* **14**, 68-73.

669 **Iglesias, N. G., Gago-Zachert, S. P., Robledo, G., Costa, N., Plata, M. I.,**
670 **Vera, O., Grau, O. & Semorile, L. C. (2008).** Population structure of *Citrus*
671 *tristeza virus* from field Argentinean isolates. *Virus Genes* **36**, 199-207.

672 **Karasev, A. V., Boyko, V. P., Gowda, S., Nikolaeva, O. V., Hilf, M. E.,**
673 **Koonin, E. V., Niblett, C. L., Cline, K., Gumpf, D. J. & other authors (1995).**
674 Complete sequence of the *Citrus tristeza virus* RNA genome. *Virology* **208**, 511-
675 520.

676 **Karl, B. N. & Hugh B. N. (1997).** GeneDoc: a tool for editing and annotating
677 multiple sequence alignments. www.nrbcs.org/gfx/genedoc/ebinet.htm.

678 **Kosakovsky-Pond & S. L., Frost, S. D. W. (2005a).** Datamonkey: rapid
679 detection of selective pressure on individual sites of codon alignments.
680 *Bioinformatics* **21**, 2531-2533.

681 **Kosakovsky-Pond, S. L. & Frost, S. D. W. (2005b).** Not so different after all: a
682 comparison of methods for detecting amino acid sites under selection. *Mol Biol*
683 *Evol* **22**, 1208-1222.

684 **Kosakovsky-Pond, S. L. & Frost, S. D. W. (2005c).** A genetic algorithm
685 approach to detecting lineage-specific variation in selection pressure. *Mol Biol*
686 *Evol* **22**, 478-485.

687 **Kosakovsky-Pond, S. L., Frost, S. D. W. & Muse S. V. (2005).** HyPhy:
688 hypothesis testing using phylogenies. *Bioinformatics* **21**, 676-679.

689 **Lai, M. M. (1992).** RNA recombination in animal and plant viruses. *Microbiol*
690 *Rev* **56**, 61-79.

691 **López, C., Navas-Castillo, J., Gowda, S., Moreno, P. & Flores, R. (2000).**
692 The 23-kDa protein coded by the 3'-terminal gene of *Citrus tristeza virus* is an
693 RNA-binding protein. *Virology* **269**, 462-470.

694 **Lu, R., Folimonov, A., Shintaku, M., Li, W. X., Falk, B. W., Dawson, W. O. &**
695 **Ding, S. W. (2004).** Three distinct suppressors of RNA silencing encoded by a
696 20-kb viral RNA genome. *Proc Natl Acad Sci U S A* **101**, 15742-15747.

697 **Martin, D. & Rybicki, E. (2000).** RDP: detection of recombination amongst
698 aligned sequences. *Bioinformatics* **16**, 562-563.

699 **Martin, D. P., Williamson, C. & Posada, D. (2005a).** RDP2: recombination
700 detection and analysis from sequence alignments. *Bioinformatics* **21**, 260-262.

701 **Martin, D. P., Posada, D., Crandall, K. A. & Williamson, C. (2005b).** A
702 modified bootscan algorithm for automated identification of recombinant
703 sequences and recombination breakpoints. *AIDS Res Hum Retroviruses* **21**, 98-
704 102.

705 **Moreno, P., Guerri, J. & Muñoz, N. (1990).** Identification of Spanish strains of
706 *Citrus tristeza virus* by analysis of double-stranded RNA. *Phytopathology* **80**,
707 477-482.

708 **Moreno, P., Ambrós, S., Albiach-Martí, M. R., Guerri, J. & Peña L. (2008).**
709 Plant diseases that changed the world - *Citrus tristeza virus*: a pathogen that
710 changed the course of the citrus industry. *Mol Plant Pathol* **9**, 251-268.

711 **Padidam, M., Sawyer, S. & Fauquet, C. M. (1999).** Possible emergence of
712 new geminiviruses by frequent recombination. *Virology* **265**, 218-225.

713 **Pappu, H. R., Karasev, A. V., Anderson, E. J., Pappu, S. S., Hilf, M. E.,
714 Febres, V. J., Eckloff, R. M., McCaffery, M., Boyko, V. & Gowda, S. (1994).**
715 Nucleotide sequence and organization of eight 3' open reading frames of the
716 *Citrus tristeza closterovirus* genome. *Virology* **199**, 35-46.

717 **Posada, D. & Crandall, K. A. (2001).** Evaluation of methods for detecting
718 recombination from DNA sequences: Computer simulations. *Proc Natl Acad Sci*
719 *U S A* **98**, 13757-13762.

720 **Reed, J. C., Kasschau, K. D., Prokhnevsky, A. I., Gopinath, K., Pogue, G.
721 P., Carrington, J. C. & Dolja, V. V. (2003).** Suppressor of RNA silencing
722 encoded by *Beet yellows virus*. *Virology* **306**, 203-209.

723 **Roossinck, M. J. (1997).** Mechanisms of plant virus evolution. *Annu. Rev.*
724 *Phytopathology* **35**, 191-209.

725 **Rubio, L., Ayllón, M. A., Kong, P., Fernández, A., Polek, M., Guerri, J.,
726 Moreno, P. & Falk, B. W. (2001).** Genetic variation of *Citrus tristeza virus*
727 isolates from California and Spain, evidence for mixed infections and
728 recombination. *J Virol* **75**, 8054-8062.

729 **Ruiz-Ruiz, S., Moreno, P., Guerri, J. & Ambrós, S. (2006).** The complete
730 nucleotide sequence of a severe stem pitting isolate of *Citrus tristeza virus* from
731 Spain, comparison with isolates from different origins. *Arch Virol* **151**, 387-398.

732 **Salminen M. O., Carr J. K., Burke D. S., McCutchan, F. E. (1995).**
733 Identification of breakpoints in intergenotypic recombinants of HIV type 1 by
734 Bootscanning. *AIDS Res Hum Retroviruses* **11**, 1423-1425.

735 **Sambade, A., Rubio, L., Garnsey, S. M., Costa, N., Müller, G. W., Peyrou,**
736 **M., Guerri, J. & Moreno, P. (2002).** Comparison of viral RNA populations of
737 pathogenically distinct isolates of *Citrus tristeza virus*, application to monitoring
738 cross-protection. *Plant Pathol* **51**, 257-265.

739 **Sambade, A., López, C., Rubio, L., Flores, R., Guerri, J. & Moreno, P.**
740 **(2003).** Polymorphism of a specific region in gene p23 of *Citrus tristeza virus*
741 allows discrimination between mild and severe isolates. *Arch Virol* **148**, 2325-
742 2340.

743 **Sambade, A., Ambrós, S., López, C., Ruiz-Ruiz, S., Hermoso de Mendoza,**
744 **A., Flores, R., Guerri J. & Moreno, P. (2007).** Preferential accumulation of
745 severe variants of *Citrus tristeza virus* in plants co-inoculated with mild and
746 severe variants. *Arch Virol* **152**, 1115-1126.

747 **Sambrook, J., Fritsch, E. F. & Maniatis, T. (1989).** Molecular Cloning: a
748 Laboratory Manual, 2nd ed. Cold Spring Harbor, NY: Cold Spring Harbor
749 Laboratory.

750 **Satyanarayana, T., Gowda, S., Boyko, V. P., Albiach-Martí, M. R., Mawassi,**
751 **M., Navas-Castillo, J., Karasev, A. V., Dolja, V., Hilf, M. E. & other authors**
752 **(1999).** An engineered closterovirus RNA replicon and analysis of heterologous
753 terminal sequences for replication. *Proc Natl Acad Sci U S A* **96**, 7433-7438.

754 **Satyanarayana, T., Gowda, S., Mawassi, M., Albiach-Martí, M. R., Ayllón, M.**
755 **A., Robertson, C., Garnsey, S. M. & Dawson, W. O. (2000).** Closterovirus
756 encoded HSP70 homolog and p61 in addition to both coat proteins function in
757 efficient virion assembly. *Virology* **278**, 253-265.

758 **Satyanarayana, T., Bar-Joseph, M., Mawassi, M., Albiach-Martí, M. R.,**
759 **Ayllón, M. A., Gowda, S., Hilf, M. E., Moreno, P., Garnsey, S. M. & Dawson,**
760 **W. O. (2001).** Amplification of *Citrus tristeza virus* from a cDNA clone and
761 infection of citrus trees. *Virology* **280**, 87-96.

762 **Satyanarayana, T., Gowda, S., Ayllón, M. A., Albiach-Martí, M. R.,**
763 **Rabindran, S. & Dawson, W. O. (2002).** The p23 protein of *Citrus tristeza virus*
764 controls asymmetrical RNA accumulation. *J Virol* **76**, 473-483.

765 **Satyanayanana, T., Gowda, S., Ayllón, M. A. & Dawson, W. O. (2004).**
766 Closterovirus bipolar virion: evidence for initiation of assembly by minor coat
767 protein and its restriction to the genomic RNA 5' region. *Proc Natl Acad Sci USA*
768 **101**, 799–804.

769 **Schmidt, H .A., Strimmer, K., Vingron, M. & von Haeseler A. (2002).** TREE-
770 PUZZLE: maximum likelihood phylogenetic analysis using quartets and parallel
771 computing. *Bioinformatics* **18**, 502-504.

772 **Shimodaira, H. & Hasegawa, M. (1999).** Multiple comparisons of log-
773 likelihoods with applications to phylogenetic inference. *Mol Biol Evol* **16**, 1114-
774 1116.

775 **Smith, J. M. (1992).** Analyzing the mosaic structure of genes. *J Mol Evol* **34**,
776 126-129.

777 **Tamura, K., Dudley, J., Nei, M. & Kumar, S. (2007).** MEGA4, Molecular
778 Evolutionary Genetics Analysis (MEGA) software version 4.0. *Mol Biol Evol* **24**,
779 1596-1599.

780 **Tárraga, J., Medina, I., Arbiza, L., Huerta-Cepas, J., Gabaldón, T., Dopazo,**
781 **J. & Dopazo, H. (2007).** Phylemon: a suite of web tools for molecular evolution,
782 phylogenetics and phylogenomics. *Nucl Acids Res* **35**, W38-42.

783 **Tatineni, S., Robertson, C. J., Garnsey, S. M., Bar-Joseph, M., Gowda, S. &**
784 **Dawson, W. O. (2008).** Three genes of *Citrus tristeza virus* are dispensable for
785 infection and movement throughout some varieties of citrus trees. *Virology* **376**,
786 297-307.

787 **Vives, M. C., Rubio, L., Sambade, A., Mirkov, T. E., Moreno P. & Guerri, J.**
788 **(2005).** Evidence of multiple recombination events between two RNA sequence
789 variants within a *Citrus tristeza virus* isolate. *Virology* **331**, 232-237.

790 **Weng, Z., Barthelson, R., Gowda, S., Hilf, M. E., Dawson, W.O., Galbraith,**
791 **D. W. & Xiong, Z. (2007).** Persistent infection and promiscuous recombination
792 of multiple genotypes of an RNA virus within a single host generate extensive
793 diversity. *PLoS ONE* **9**, e917.

794 **Worobey, M. & Holmes E. C. (1999).** Evolutionary aspects of recombination in
795 RNA viruses. *J Gen Virol* **80**, 2535-2543.

796 **Yang, G. A., Mawassi, M., Gofman, R., Gafny, R. & Bar-Joseph, M. (1997).**
797 Involvement of a subgenomic mRNA in the generation of a variable population
798 of defective *Citrus tristeza virus* molecules. *J Virol* **71**, 9800-9802.

799 **Figure legends**

800

801 **Figure 1.** Unrooted protein maximum likelihood phylogenetic trees of CTV
802 genomic regions p33, p65, p61, p18, p13, p20 and p23. Isolates with congruent
803 phylogenetic relationships are in bold italics. Circles beside the isolate names
804 indicate their pathogenicity characteristics (biogroups 1-2, 4 and 5), with
805 biogroups 0 (isolate K) and 3 left without circle. Geographical origin is indicated
806 by double daggers (Florida), underlined names (South America), rectangular
807 boxes (Spain), ovals (Japan) or no mark for isolates K (France) and VT (Israel).
808 Scale bars indicate number of changes per position for a unit branch length.
809 Bootstrap values for significance of nodes are indicated by asterisks (***) 90-
810 100%, ** 70-89%, * 50-69%).

811

812 **Figure 2.** Recombination analysis of concatenated CTV sequences. a) Layout
813 of the CTV genome with the regions analyzed indicated as black boxes. ORFs
814 are represented by empty boxes with indication of the encoded protein. Motifs
815 protease (Pro), methyl transferase (MT) and helicase (HEL) of the p349 protein
816 are also indicated. b) Unrooted nucleotide maximum likelihood phylogenetic
817 trees of seven CTV genomic regions (see Fig. 1). Isolates with congruent
818 phylogenetic relationships are in bold italics. In each region, isolates detected
819 as recombinant by RDP3 programs are highlighted with solid ovals and others
820 showing incongruent phylogenetic relationships with dotted ovals. Values for
821 branch-specific non-synonymous to synonymous class substitution rates (ω) are
822 indicated with colors. c) Recombination hypotheses generated by at least four
823 algorithms of the RDP3 program and further refined by inspection of maximum
824 likelihood trees. Concatenated alignments are outlined at the top; long colored
825 boxes represent CTV concatenated sequences (isolate code above the box)
826 and internal pale colored segments indicate recombinant regions; the major
827 parental for each recombinant sequence is indicated below the isolate code,
828 and the minor, by short boxes below the pale colored segments.

829

830 **Figure 3:** Phylogenetic network of CTV isolates constructed by split-
831 decomposition analysis. Circles beside the isolate names as in Figure 1.

832

833 **Table 1.** Nucleotide distances and frequency and strength of negative selection
 834 at protein level in different CTV genomic regions.

Genomic regions	Average nucleotide distance \pm SD	Number of negatively selected codons	$d_N-d_S \pm$ SEM [†]
p33	0.1641 \pm 0.0196**	20 (0.1587)	-2.1098 \pm 0.2505**
p65	0.1091 \pm 0.0160	30 (0.1224)	-6.7162 \pm 0.6392
p61	0.1250 \pm 0.0160	35 (0.1471)	-7.2201 \pm 1.1737
p18	0.0955 \pm 0.0170	12 (0.1034)	-8.6499 \pm 1.2031
p13	0.0741 \pm 0.0240*	16 (0.1403)	-6.3229 \pm 1.0417
p20	0.1145 \pm 0.0188	16 (0.1000)	-7.8506 \pm 0.9940
p23	0.1018 \pm 0.0143	21 (0.1005)	-9.0407 \pm 1.7731

835 SD: Standard deviation for nucleotide distance computed by the bootstrap
 836 method (500 replicates).

837 [†] d_N-d_S : Average of normalized values of the difference between non-
 838 synonymous and synonymous substitutions of selected codons estimated by
 839 the fixed effects maximum likelihood method (FEL) and their standard errors
 840 (SEM).

841 ** Statistically different to the rest of values ($p < 0.05$).

842 * Statistically different to p33 and p61 values ($p < 0.05$).

843

844 **Table 2.** Lineage-specific analysis of selective pressures in seven genomic
 845 regions of CTV.

	Genomic regions						
	p33	p65	p61	p18	p13	p20	p23
ω^*	1.029 (25%)	1.327 (14%)		21.785 (10%)	∞ (12%)	∞ (1%)	∞ (2%)
	0.253 (38%)	0.211 (22%)	0.641 (21%)	0.265 (64%)			0.461 (40%)
			0.122 (67%)		0.208 (71%)	0.155 (54%)	0.143 (57%)
	0.087 (36%)	0.043 (50%)		0.013 (27%)		0.030 (44%)	
		0.000 (14%)	0.000 (12%)		0.000 (17%)		

846 * ω . d_N/d_S class ratios for non-synonymous to synonymous substitution rates; the
 847 proportion of branches assigned to each class is shown in parenthesis.

Figure 1

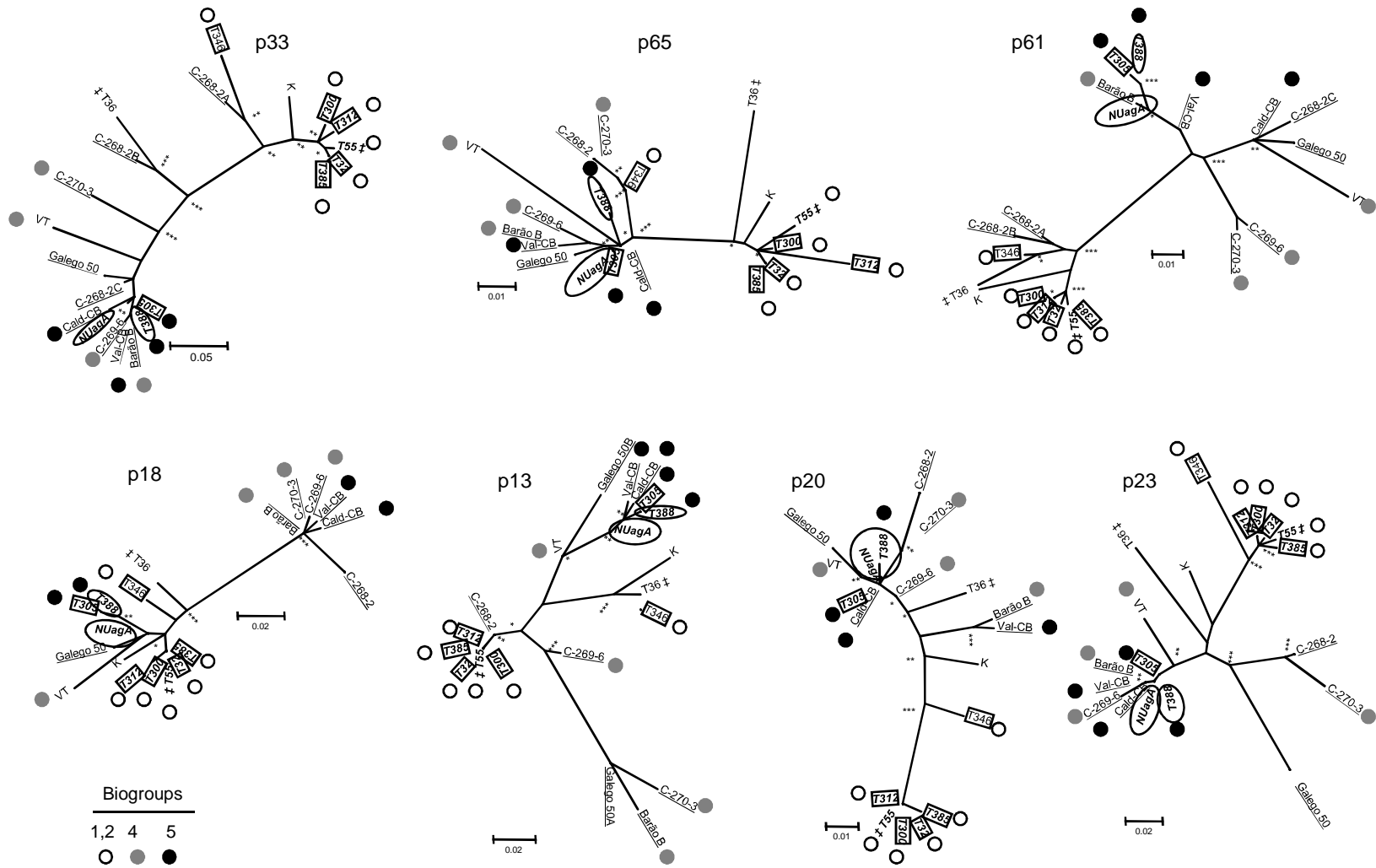


Figure 2

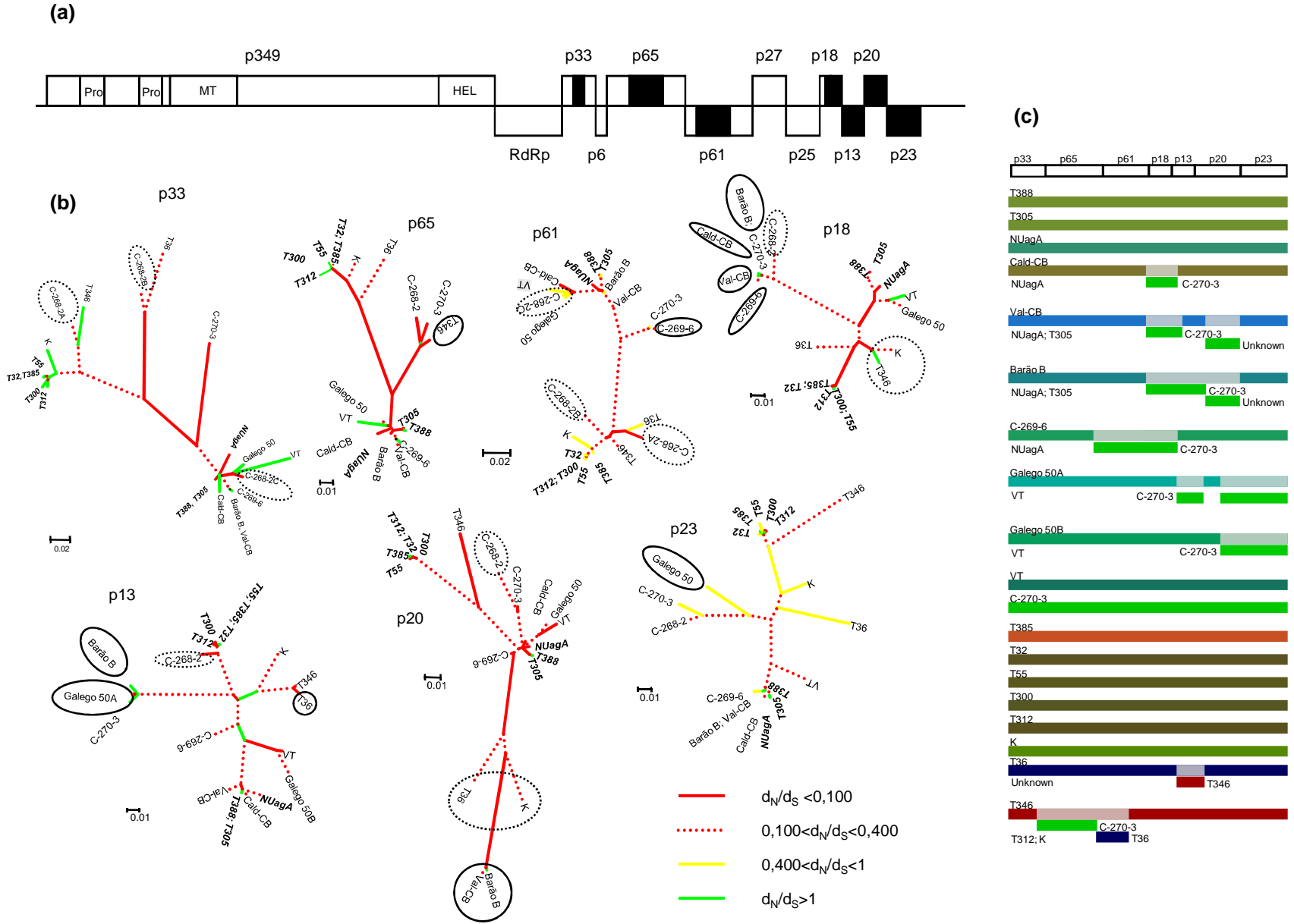


Figure 3

