

Version dated: November 4, 2019

TESTING AND VISUALISING COMPOSITIONAL HOMOGENEITY

## Software for Detecting Heterogeneous Evolutionary Processes across Aligned Sequence Data

LARS S JERMIIN<sup>1,2,3</sup>, DAVID R LOVELL<sup>4</sup>, BERNHARD MISOF<sup>5</sup>, PETER G FOSTER<sup>6</sup>, JOHN  
ROBINSON<sup>7</sup>

<sup>1</sup>*Land & Water, CSIRO, Acton, ACT 2601, Australia;*

<sup>2</sup>*Research School of Biology, Australian National University, Acton, ACT 2601, Australia;*

<sup>3</sup>*School of Biology & Environmental Science, University College Dublin, Belfield, Dublin 4, Ireland;*

<sup>4</sup>*Earth Institute, University College Dublin, Belfield, Dublin 4, Ireland;*

<sup>5</sup>*Data61, CSIRO, Acton, ACT 2601, Australia;*

<sup>6</sup>*School of Electrical Engineering & Computer Science, Queensland University of Technology, Brisbane,  
QLD 4001, Australia;*

<sup>7</sup>*Zoologisches Forschungsmuseum Alexander Koenig, 53113 Bonn, Germany;*

<sup>8</sup>*Department of Life Sciences, Natural History Museum, London, SW7 5BD, UK;*

<sup>9</sup>*School of Mathematics & Statistics, University Sydney, Sydney, NSW 2006, Australia*

**Correspondence:** Lars S Jermiin, Research School of Biology, Australian National  
University, Canberra, ACT, 0200, Australia; E-mail: lars.jermiin@anu.edu.au.

**Abstract:** Most model-based molecular phylogenetic methods assume that the sequences  
diverged on a tree under homogeneous conditions. If evolution occurred under these  
conditions, then it is unlikely that the sequences would become compositionally  
heterogeneous. Conversely, if the sequences are compositionally heterogeneous, then it is  
unlikely that they have evolved under homogeneous conditions. We present methods to detect  
and analyse heterogeneous evolution in aligned sequence data and to examine—visually and  
numerically—its effect on phylogenetic estimates. The methods are implemented in three  
programs, allowing users to better examine under what conditions their phylogenetic data  
may have evolved.

**Keywords:** Evolution under stationary conditions; Matched-pairs test of symmetry; PP plot;  
Heat map; Historical signal; compositional signal; compositional distance; networks.

30 Most model-based molecular phylogenetic methods assume that the sequences of  
31 nucleotides or amino acids have evolved along the edges of a single bifurcating tree. Often, the  
32 methods also assume that the evolutionary processes operating at the variable sites of these  
33 data (i.e., the sites that are free to evolve) can be approximated by independent and  
34 identically-distributed (*iid*) Markovian processes. Furthermore, it is often assumed that the  
35 evolutionary processes were stationary, reversible and homogeneous (SRH) (for details, see  
36 Bryant et al. 2005; Jayaswal et al. 2005; Ababneh et al. 2006a,b; Jermiin et al. 2017), with the  
37 term homogeneity implying time-homogeneity (i.e., a constant rate of change between two  
38 points in time).

39 In practice, when DNA has evolved under these conditions, commonly-used  
40 phylogenetic methods are likely to identify the correct topology (Huelsenbeck and Hillis 1993;  
41 Hillis et al. 1994a,b). However, the same methods may not be capable of identifying the  
42 correct topology when DNA has evolved under more complex conditions (Huelsenbeck and  
43 Hillis 1993; Hillis et al. 1994a,b; Ho and Jermiin 2004; Jermiin et al. 2004). One reason for  
44 this failure is that the strength of the *historical signal* (i.e., the signal in DNA that is due to  
45 the order and time of divergence events) decays over time (Ho and Jermiin 2004) whereas the  
46 strength of the *non-historical signals* (Grundy and Naylor 1999) may increase over time (Fig.  
47 1). This may lead to situations, where the non-historical signals—individually or jointly—may  
48 become stronger than the historical signal (Ho and Jermiin 2004). Unless phylogenetic  
49 methods are able to distinguish historical signals from non-historical signals, the latter may be  
50 misinterpreted as being part of the historical signal. This is because the non-historical signals  
51 are also *phylogenetic signals*.

52 The non-historical signal is a mixed bag of signals that may arise over time due to  
53 temporal variations in site- and lineage-specific evolutionary processes. For example, when the  
54 homologous sites in a pair of sequences evolve under different conditions, evolutionary  
55 processes cannot be homogeneous, and compositional heterogeneity across the sequences may  
56 arise. When this happens, there is a *compositional signal* in the data (Fig. 2). On the other  
57 hand, when compositional heterogeneity is found across an alignment of homologous  
58 sequences, there is evidence of evolution under non-stationary conditions.

59 Several methods have been developed to detect compositional heterogeneity across  
60 homologous sequences (reviewed in Jermiin et al. (2004, 2009)), but doubt remains about  
61 what method is most appropriate (cf. Jermiin et al. (2004) and Duchêne et al. (2017)). To  
62 resolve this matter and to empower concerned users of phylogenetic methods, we present  
63 software to detect and visualise compositional heterogeneity across aligned sequence data. The  
64 software also facilitates assessment of the impact of compositional heterogeneity on inferred  
65 phylogenetic trees and networks.

## METHODOLOGY

67

## Background

68 Consider a nucleotide sequence that evolves over the edges of a rooted tree (Fig. 3), and  
69 assume that the 90 sites in this sequence evolve under *iid* conditions. At time  $t_0$ , the ancestral  
70 sequence, **Seq0**, evolves along an ancestral edge in the tree (Fig. 3a). At time  $t_1$ , the sequence  
71 meets a bifurcation in the tree, and it becomes two identical sequences, **Seq1** and **Seq2** (Fig.  
72 3b). At time  $t_2$ , the two sequences have evolved further under independent evolutionary  
73 processes (Fig. 3c), so they are unlikely to be the same. The sequences at  $t_0$ ,  $t_1$  and  $t_2$  are  
74 shown in Figure 3d.

75 Methodologically, the challenge now is to extract as much information as possible from  
76 the alignment of **Seq1'** and **Seq2'** (e.g., to infer the time elapsed since the bifurcation at  $t_1$ ).  
77 One way to extract information from such a data set is to consider the ratio of the number of  
78 sites where the sequences differ to the total number of sites compared. This yields a metric  
79 called the  $p$  distance (the  $p$  distance between **Seq1'** and **Seq2'** is 42/90). Another way to do  
80 this is to use a divergence matrix (**N**). For **Seq1** and **Seq2**, we get:

$$\mathbf{N}(t_1) = \begin{array}{c|cccc} & A & C & G & T \\ \hline A & 21 & 0 & 0 & 0 \\ C & 0 & 23 & 0 & 0 \\ G & 0 & 0 & 24 & 0 \\ T & 0 & 0 & 0 & 22 \end{array},$$

81 while for **Seq1'** and **Seq2'** we get:

$$\mathbf{N}(t_2) = \begin{array}{c|cccc} & A & C & G & T \\ \hline A & 14 & 7 & 10 & 2 \\ C & 3 & 13 & 2 & 1 \\ G & 1 & 1 & 11 & 4 \\ T & 7 & 3 & 1 & 10 \end{array}.$$

82 The only difference between **N** and the alignments in Figure 3d is that information about the  
83 order of sites in the alignment is lost in the divergence matrix. However, as these sites are  
84 assumed to have evolved independently, this loss of information is of no consequence for most  
85 commonly-used phylogenetic methods.

86 Given **N**, we can obtain the  $p$  distance or any other evolutionary distance, like the F81  
87 distance (Felsenstein 1981). Likewise, we can determine whether two sequences have diverged  
88 under homogeneous conditions. If the distributions of  $X_1$  and  $X_2$  are equal, then the  
89 sequences will have evolved under homogeneous conditions. Assuming evolution under  
90 homogeneous conditions, the divergence matrix should be approximately symmetrical (i.e., if  
91  $\mathbf{N} = \{n_{ij}\}$ , then  $E(n_{ij}) = E(n_{ji}) \forall i, j$ ;  $E$  denotes the expected value).

92

### The Matched-pairs Test of Symmetry

93 The matched-pairs test of symmetry is suitable for testing whether  $E(n_{ij}) = E(n_{ji})$ . It is  
94 computed using:

$$X_B^2 = \sum_{i < j} \frac{(n_{ij} - n_{ji})^2}{n_{ij} + n_{ji}}, \quad (1)$$

95 which, assuming homogeneous conditions, is asymptotically distributed as a  $\chi^2$  variate on  
96  $\nu = c \times (c - 1)/2$  degrees of freedom, where  $c$  denotes the number of unique letters in the  
97 sequences' alphabet (for DNA,  $c = 4$ ). Given  $X_B^2$  and  $\nu$ , it is easy to obtain the probability of  
98 getting a test statistic that equals or exceeds  $X_B^2$ , given  $\nu$  (i.e.,  $p = P(\chi_\nu^2 \geq X_B^2)$ ). In this  
99 regard, it is worth remembering that  $X_B^2 = X_S^2 + X_A^2$ , where  $X_S^2$  is the test statistic from the  
100 matched-pairs test of marginal symmetry (Stuart 1955) while  $X_A^2$  is the test statistic from the  
101 matched-pairs test of internal symmetry (Ababneh et al. 2006b). It is also worth pointing out  
102 that if for any of the comparisons  $n_{ij} + n_{ji} = 0$ , the entry is ignored and  $\nu$  is reduced by 1.

103 The matched-pairs test of symmetry was devised by Bowker (1948) and introduced to  
104 molecular phylogenetics by Tavaré (1986). Subsequent attempts to promote this test as the  
105 best approach to test homogeneity of the evolutionary processes (Lanave and Pesole 1993;  
106 Waddell and Steel 1997; Waddell et al. 1999; Ababneh et al. 2006b) were largely unsuccessful,  
107 with one opponent stating that the test “is hardly necessary because typical phylogenetic  
108 datasets are large and can reject the null hypothesis with ease” (Yang 2014). That is an odd  
109 statement, as it recommends ignoring a reason for systematic error. More recently, Duchêne et  
110 al. (2017) used a test described by Foster (2004), which tests the fit of the compositional  
111 component of the (stationary, in this case) model to the data. This test uses a contingency  
112 table made up of  $c$  marginal sums. However, unlike the standard  $r \times c$  contingency table test  
113 of homogeneity, the test statistic is not compared to the  $\chi^2$  distribution but to a simulated  
114 null distribution obtained on the basis of the tree and the (possibly non-stationary) model of  
115 evolution being tested. In other words, it is a test of model fit—it needs to be used after the  
116 tree and model of evolution have been specified. Thus, it is akin to the Goldman-Cox test of  
117 goodness-of-fit (Goldman 1993), which uses simulations to assess the significance of a statistic.

118 While the test used by Foster (2004) tests marginal compositions, it ignores the  
119 homology statements that alignments represent. The impact of doing so can be dramatic, as  
120 the following example reveals. The three divergence matrices, left to right, are the products of  
121 increasingly dissimilar evolutionary processes:

$$\mathbf{N}_1 = \begin{bmatrix} 40 & 10 & 20 & 30 \\ 10 & 40 & 30 & 20 \\ 20 & 30 & 40 & 10 \\ 30 & 20 & 10 & 40 \end{bmatrix} \quad \mathbf{N}_2 = \begin{bmatrix} 40 & 10 & 20 & 30 \\ 30 & 40 & 10 & 20 \\ 20 & 30 & 40 & 10 \\ 10 & 20 & 30 & 40 \end{bmatrix} \quad \mathbf{N}_3 = \begin{bmatrix} 40 & 0 & 0 & 60 \\ 60 & 40 & 0 & 0 \\ 0 & 60 & 40 & 0 \\ 0 & 0 & 60 & 40 \end{bmatrix}.$$

122 In the first case ( $\mathbf{N}_1$ ), there is no evidence that the evolutionary processes might have been  
123 different, while in the other cases ( $\mathbf{N}_2$  and  $\mathbf{N}_3$ ), the evidence of that is clearer. However, it is  
124 also clear that the three matrices have the same marginal distribution, so Foster's (2004) test  
125 cannot detect this type of lineage-specific heterogeneity in the evolutionary processes. Foster's  
126 (2004) test is similar to Stuart's (1955) matched-pairs test of marginal symmetry. If the aim is  
127 to test the fit between tree, model and data, then it would be appropriate to use Foster's  
128 (2004) test or the Goldman-Cox test of goodness of fit (Goldman 1993). On the other hand, if  
129 the aim is to test whether sequences are consistent with the assumption of evolution under  
130 stationary conditions, then Stuart's (1955) matched-pairs test of marginal symmetry is  
131 recommended (Ababneh et al. 2006a). Stuart's (1955) matched-pairs test of marginal  
132 symmetry, like Bowker's (1948) matched-pairs test of symmetry, assumes aligned data but not  
133 a tree or model, so it is useful for screening phylogenetic data *before* they are analysed. On  
134 the other hand, Foster's (2004) test is applicable *after* this analysis, can be used with  
135 non-stationary models, and is not restricted to sequence pairs.

136

### *The PP Plot*

137 If we wish to apply the matched-pairs test of symmetry to an alignment with more than two  
138 sequences, then the problem of multiple comparisons arises. For example, if a data set contains  
139 22 sequences, then there will be  $22 \times 21/2 = 231$   $p$ -values to interpret, one for each pair of  
140 sequences. However, the  $p$ -values are not independent, so they must be interpreted jointly.  
141 This can be done using a PP-plot, which displays observed  $p$ -values against expected  $p$ -values.  
142 If evolution occurred under homogeneous conditions, then the 231  $p$ -values will be distributed  
143 as a uniform random variable on (0,1). Given this expectation, we can evaluate whether the  
144 data set, as a whole, meets the assumption of evolution under homogeneous conditions.

145 To demonstrate the merits of the PP plot, we analysed an alignment of simulated  
146 nucleotides generated under time-reversible conditions on a 22-tipped tree (Fig. 4a). The PP  
147 plot in Figure 4b shows the result from data generated under the null hypothesis. As  
148 expected, the 231 dots are distributed along the diagonal, with  $\sim 5\%$  of them (12) below 0.05  
149 (i.e., the horizontal line in Fig. 4b). None of the observed  $p$ -values fell below the 5%  
150 family-wise error rate (i.e.,  $0.05/231 = 0.000216$ ). The PP plot shows the distribution to  
151 expect when the data have evolved under homogeneous conditions. This interpretation is  
152 consistent with those in Schweder and Spjøtvoll (1982) and Vera-Ruiz et al. (2014).

153

### *The Heat Map*

154 A PP-plot that deviates noticeably from that shown in Fig. 4b (e.g., the dots are not  
155 distributed along the diagonal; more than 5% of the observed  $p$ -values are below 0.05; the  
156 smallest observed  $p$ -value is below a 5% family-wise error rate), suggests that some of the  
157 sequences have evolved under heterogeneous conditions. However, the PP plot cannot identify

158 the ‘offending’ sequences, but a color-coded heat map with the observed  $p$ -values can. Figure  
159 4c shows the heat map corresponding to the data in Figure 4b. Each pixel is color-coded  
160 according to the  $p$ -value for the corresponding pair of sequences. Most of the pixels are white  
161 because the  $p$ -values are  $\geq 0.05$ . Some pixels are yellow, but none of them are darker; this is  
162 consistent with the condition under which the sequences were generated.

163 When a heat map differs noticeably from that in Figure 4c, it allows us to identify  
164 sequences that are unlikely to have evolved under the same conditions. For example, if all but  
165 one of the sequences evolved under homogeneous conditions, then that would result in a heat  
166 map where a row and/or column has darker pixels. The color of a pixel depends on the  
167 probability that the corresponding pair of sequences have evolved under homogeneous  
168 conditions. A dark row and/or column identifies an offending sequence, which then can be  
169 removed if it is insignificant to the phylogenetic question. Figure 6 of Jayaswal et al. (2014)  
170 shows such a heat map (in this case the offending sequences could not be removed).

171 When two or more sequences are regarded as offending, we might ask whether the data  
172 can be grouped into subsets of sequences that are consistent with evolution under  
173 homogeneous conditions. To do so, one simply needs to permute the rows and columns of the  
174 heat map, or reorder the sequences in the alignment before analysing the data again. Figures  
175 6 and 7 of Jermiin et al. (2017) show two heat maps for the same data, obtained before and  
176 after a permutation of the rows and columns of the heat map. In the first figure, several small  
177 sets of sequences appear to have evolved under homogeneous conditions. However, the second  
178 figure reveals that many of these subsets can be merged into larger subsets of sequences that  
179 appear to have evolved under different homogeneous conditions. In summary, the PP plot and  
180 heat map provide researchers an opportunity to survey their data far more thoroughly before  
181 model selection and phylogenetic analysis.

## 182 *Compositional Distances*

183 A compositional signal may arise when sequences diverge under non-homogeneous conditions.  
184 If such a signal emerges, its amplitude can be measured using distance metrics that quantify  
185 departure from symmetry of a divergence matrix. Compositional distances are appropriate for  
186 vectors of non-negative values that carry information in their relative (not absolute) amounts  
187 (Aitchison 1986; Egozcue and Pawlowsky-Glahn 2011), like those in a divergence matrix.  
188 Compositional distances may be used to infer trees and networks, revealing relationships  
189 based solely on compositional differences. These trees and networks may uncover a  
190 compositional signal’s potential impact on phylogenetic estimates.

191 Given  $\mathbf{N}$  (for nucleotides):

$$\mathbf{N} = \begin{bmatrix} n_{11} & n_{12} & n_{13} & n_{14} \\ n_{21} & n_{22} & n_{23} & n_{24} \\ n_{31} & n_{32} & n_{33} & n_{34} \\ n_{41} & n_{42} & n_{43} & n_{44} \end{bmatrix}$$

192 we can define two vectors that relate to the off-diagonal elements of the upper and lower  
193 triangles:

$$\begin{cases} \mathbf{Y} = \{y_k\} = (n_{12}, n_{13}, n_{14}, n_{23}, n_{24}, n_{34}) \\ \mathbf{Z} = \{z_k\} = (n_{21}, n_{31}, n_{41}, n_{32}, n_{42}, n_{43}) \end{cases}$$

194 Given  $\mathbf{Y}$  and  $\mathbf{Z}$  for a  $c$ -state alphabet (e.g.,  $c = 20$  for protein), it is possible to compute three  
195 compositional distances:

$$\delta_{EFS} = \sqrt{\sum_{i=1}^l (y_k - z_k)^2}, \quad (2)$$

$$\delta_{AFS} = \sqrt{\frac{1}{2 \times l} \sum_{a=1}^l \sum_{b=1}^l \left( \log \frac{y_a}{y_b} - \log \frac{z_a}{z_b} \right)^2}, \quad (3)$$

and

$$\delta_{CFS} = \sqrt{\frac{X_B^2}{\nu}}. \quad (4)$$

196 Here,  $\delta_{EFS}$ ,  $\delta_{AFS}$ , and  $\delta_{CFS}$  respectively denote the Euclidean distance, Aitchison's (1986)  
197 distance, and a distance metric closely related to Bowker's (1948) matched-pairs test of  
198 symmetry, and  $l$  is the number of elements in  $\mathbf{Y}$  and  $\mathbf{Z}$ . The Euclidean distance measures the  
199 distance between two points in Euclidean space, taking no account of sign or scale, so they are  
200 not appropriate for count data. One more appropriate metric is that of Aitchison (1986); for a  
201 comparison of these distance metrics, see Lovell et al. (2011). One undesirable property of  
202  $\delta_{AFS}$  is that it is zero when  $n_{ij}/n_{ji}$  is constant, and will be small if this is even approximately  
203 so. Because of this, Aitchison's (1986) distance is not suitable for data used to measure lack of  
204 symmetry in divergence matrices. Instead, we may use  $\delta_{CFS}$ , which has the advantage of  
205 being able to accommodate that comparisons between different pairs of sequences may be  
206 associated with different degrees of freedom ( $\nu$ ). Note that  $\delta_{CFS} \geq 0.0$ , and that  $\delta_{CFS}$  is not  
207 an evolutionary distance in the sense that the LogDet (Lockhart et al. 1994; Steel 1994) or  
208 paralinear (Lake 1994) distances are.

209

### *The Nature of Bias in Phylogenetic Estimates*

210 It is difficult to detect bias in phylogenetic estimates from real sequence data, but it is well  
211 known that bias may manifest itself in at least two ways:

- 212 1. The topology of the tree (or network) is affected, implying that the length of at least  
213 some of the edges (or weights of some of the splits; ‘weight’ is analogous with length, in  
214 the sense of Huson and Bryant (2006)) also will be affected, or
- 215 2. The topology is unaffected but the length of the edges in the tree (or the weights of the  
216 splits in the network) may be affected.

217 Both of these biases are cause for concern, even if only the topology is of interest, because the  
218 topology is a discrete entity, whose accuracy often is dependent on the accuracy of the  
219 estimates of the other parameters. The challenge is to get all the estimates as accurate as  
220 possible without increasing the variance or the bias of these estimates (Dziak et al. 2019). In  
221 other words, both over- and under-parameterisation of the data should be avoided.

### 222 *Visualising the Effect of Compositional Heterogeneity on Trees and Networks*

223 Given a distance matrix  $\mathbf{D}_{CFS}$  with estimates of  $\delta_{CFS}$ , we may infer a *compositional tree*,  $\mathcal{T}$ ,  
224 and a *compositional network*,  $\mathcal{N}$ . This can be done by using programs like FastME (Lefort et  
225 al. 2015) and SplitsTree4 (Huson and Bryant 2006). Such structures display the relationships  
226 among sequences based solely on compositional distances, so they should not be interpreted as  
227 if they were phylogenetic trees or phylogenetic networks. Sequences that are compositionally  
228 similar may not be close in an evolutionary sense, and sequences that are compositionally  
229 dissimilar may not be distantly-related in an evolutionary sense. The advantage of using  
230 data-display networks to reveal conflicting signals in phylogenetic data has already been  
231 demonstrated by Morrison (2010), so it will not be reiterated here.

232 Consider a data set that has been found to violate the phylogenetic assumption of  
233 evolution under homogeneous conditions. In such a case, one might wish to know whether the  
234 compositional signal has become so strong that it might bias a phylogenetic estimate, unless it  
235 is properly accounted for.

236 To demonstrate the benefit of using  $\mathcal{T}$  and  $\mathcal{N}$ , we analysed an alignment of five 16S  
237 rRNA sequences from bacteria, first analysed phylogenetically by Embley et al. (1993) and  
238 then by Galtier and Gouy (1995), Mooers and Holmes (2000), Foster (2004), and Jayaswal et  
239 al. (2005, 2007). For these data,  $\mathbf{D}_{EFS}$ ,  $\mathbf{D}_{AFS}$  and  $\mathbf{D}_{CFS}$  are

<i>Aquifex</i>	0.0000	0.0120	0.0436	0.0461	0.0119
<i>Thermotoga</i>	0.0120	0.0000	0.0431	0.0447	0.0098
<i>Bacillus</i>	0.0436	0.0431	0.0000	0.0043	0.0391
<i>Deinococcus</i>	0.0461	0.0447	0.0043	0.0000	0.0418
<i>Thermus</i>	0.0119	0.0098	0.0391	0.0418	0.0000



<i>Aquifex</i>	0.0000	1.1104	2.8378	2.2533	0.8754
<i>Thermotoga</i>	1.1104	0.0000	3.0770	2.8549	0.9365
<i>Bacillus</i>	2.8378	3.0770	0.0000	0.1914	2.4787
<i>Deinococcus</i>	2.2533	2.8549	0.1914	0.0000	3.0188
<i>Thermus</i>	0.8754	0.9365	2.4787	3.0188	0.0000

and

<i>Aquifex</i>	0.0000	1.2805	3.2136	3.0451	0.9379
<i>Thermotoga</i>	1.2805	0.0000	3.3176	3.2677	1.0064
<i>Bacillus</i>	3.2136	3.3176	0.0000	0.3382	2.9026
<i>Deinococcus</i>	3.0451	3.2677	0.3382	0.0000	3.1461
<i>Thermus</i>	0.9379	1.0064	2.9026	3.1461	0.0000

240 respectively (the values in  $\mathbf{D}_{EFS}$  and  $\mathbf{D}_{AFS}$  were obtained using Homo v1.3: Rouse et al.

241 2013). The three matrices differ, reflecting the differences between Equations 2, 3 and 4.

242 Interestingly, the elements of  $\mathbf{D}_{EFS}$ ,  $\mathbf{D}_{AFS}$  and  $\mathbf{D}_{CFS}$  appear to be highly correlated (i.e.,  
243 carrying quite similar information), but this is not always the case (e.g., if  $\mathbf{Y} \propto \mathbf{Z}$ ).

244 Figures 5a and 5b shows a BioNJ tree (Gascuel 1997) and a Neighbor-Net (Bryant and  
245 Moulton 2004), both inferred from  $\mathbf{D}_{CFS}$  using SplitsTree4 (Huson and Bryant 2006). The  
246 compositional tree ( $\mathcal{T}$ ) has a long internal edge (marked † in Fig. 5a) that separates  
247 *Deinococcus* and *Bacillus* from the other three species. The same appears to be the case for  
248 the compositional network ( $\mathcal{N}$ ) in Fig. 5b. Indeed,  $\mathcal{N}$  is very *treelike*, because the split  
249 marked † in Figure 5b is 18.6 times longer than the second-longest alternative (marked ‡). In  
250 other words,  $\mathcal{T}$  and  $\mathcal{N}$  corroborate what is already known about these five sequences:  
251 *Deinococcus* and *Bacillus* are compositionally distinct from the other three species (Galtier  
252 and Gouy 1995; Jayaswal et al. 2005). However, in many other studies, such knowledge is not  
253 available or heeded. This is where compositional trees or networks become useful; not only do  
254 the topologies of  $\mathcal{T}$  and  $\mathcal{N}$  identify the compositionally most similar sequences, they also  
255 reveal where the biggest differences are—and as compositional differences grow, so do the  
256 length of edges in  $\mathcal{T}$  and splits in  $\mathcal{N}$ . Importantly, compositional networks are able to reveal  
257 conflicting information in multiple sequence alignments that compositional trees cannot reveal  
258 (because the latter are constrained to be acyclic graphs: Penny et al. 1992). Therefore, during  
259 the exploratory phase of assessing compositional heterogeneity, using  $\mathcal{N}$  may be better than  
260 using  $\mathcal{T}$ .

### 261 *Congruence between Phylogenetic and Compositional Trees*

262 Interestingly, the split observed between *Deinococcus* and *Bacillus* and the other three species  
263 (Fig. 5) is also found in optimal phylogenetic trees inferred under different time-reversible

264 Markovian models of sequence evolution (Jayaswal et al. 2005, 2007). At least two  
265 explanations may be given for this congruence of splits:

- 266 1. The historical and compositional signals in the data are aligned, implying that the  
267 historical signal is augmented by the compositional signal. A consequence of this is that  
268 the inferred topology may be correct. However, estimates of edge lengths may still be  
269 biased; this could lead to bias in estimates of divergence dates.
- 270 2. The historical and compositional signals are not aligned, implying that the historical  
271 signal might be undermined by the compositional signal. This would entail that the  
272 phylogenetic methods, unless specifically designed to accommodate a compositional  
273 signal, might misinterpret the compositional signal, as if it were the historical signal, and  
274 return a phylogenetic estimate with biases in both topology and edge lengths.

275 In the first explanation, the compositional signal may be stronger than the historical signal  
276 but because the two signals are aligned, this has no adverse effect on the inferred topology; on  
277 the contrary, it may help us to identify the correct topology. In the second explanation, both  
278 the strength and the complexity of the compositional signal are likely to contribute to bias in  
279 phylogenetic estimates. Importantly, the identities and lengths of internal edges in the true  
280 tree are both factors contributing to the success or failure of phylogenetic inference (Jermin  
281 et al. 2004), but neither of these factors is known (except for in simulation-based studies).

282 The problem with these two explanations is that they apply equally well to many  
283 studies of compositionally heterogeneous phylogenetic data sets and that we do not know  
284 which one is right. It is not wise to argue that other phylogenetic estimates corroborate a  
285 current phylogenetic hypothesis, unless bias due to model misspecification has been ruled out  
286 for *all the data sets* being compared. In the present case, the matter was resolved by analysing  
287 the alignment using a model that was heterogeneous over the tree (Foster 2004) and by using  
288 the general Markov model of sequence evolution (Jayaswal et al. 2007). However, this is rarely  
289 done.

### 290 *Testing for Similarity between Phylogenetic and Compositional Trees*

291 Often phylogenetic data contain more than five sequences and it may be less clear (than e.g.,  
292 Fig. 5) whether a compositional signal contributed adversely to a phylogenetic estimate. In  
293 such cases, it may be useful to compare the phylogenetic tree ( $\mathcal{T}_r$  — inferred directly from the  
294 sequence alignment) and the compositional tree ( $\mathcal{T}_c$  — inferred from the corresponding matrix  
295 of compositional distances ( $\mathbf{D}_{CFS}$ )). In such instances, the distance between  $\mathcal{T}_r$  and  $\mathcal{T}_c$  must  
296 first be obtained. Reviewing the performance of tree-comparison metrics, Kuhner and Yamato  
297 (2015) found that Nye et al.'s (2006) metric, which is based on topology only, is superior for  
298 dissimilar trees. Their metric,  $\delta_{Align}$ , which measures how well two trees align to each other,  
299 was revealed to be better than four other tree-distance metrics, including the Robinson and

300 Foulds (1981) metric and the Path Difference metric (Williams and Clifford 1971; Penny et al.  
301 1982).

302 When comparing  $\mathcal{T}_r$  and  $\mathcal{T}_c$ , a critical question is whether they are more similar, or  
303 dissimilar, to one another than random trees are to each other. If the evolutionary process of  
304 sequence data is modelled accurately, there is no reason to presume that  $\mathcal{T}_r$  and  $\mathcal{T}_c$  will be  
305 more similar, or dissimilar, to one another, than two random trees are. Thus, we may  
306 formulate a testable null hypothesis.  $H_0$ :  $\mathcal{T}_r$  and  $\mathcal{T}_c$  are neither more similar, or dissimilar, to  
307 each other than random trees are.

308 To execute this test, we first calculate  $\delta_{Align}$  for  $\mathcal{T}_r$  and  $\mathcal{T}_c$ . Next, we generate, say,  
309 2000 random trees and partition them into 1000 pairs. For each pair, we calculate  $\delta_{Align}^*$ ,  
310 where the ‘star’ signals that this is an estimate obtained from random trees. Finally, the  
311 distribution of  $\delta_{Align}^*$  values is charted and the value of  $\delta_{Align}$  for  $\mathcal{T}_r$  and  $\mathcal{T}_c$  is matched to this  
312 distribution. If the value of  $\delta_{Align}$  falls well within the distribution of  $\delta_{Align}^*$ , then the  
313 topologies of  $\mathcal{T}_r$  and  $\mathcal{T}_c$  are random with respect to each other; otherwise, they are more  
314 similar (e.g., if  $\delta_{Align} < \delta_{Align}^*$  for all pairs) or dissimilar (e.g., if  $\delta_{Align} > \delta_{Align}^*$  for all pairs) to  
315 each other than random trees are.

316 The method is illustrated in the biological example (below).

### 317 *Software*

318 The methods described above are implemented in three programs.

319 *Homo*.—Homo v2.0 is a complete re-development of previous versions of Homo (Rouse et al.  
320 2013; <http://www.csiro.au/Homo>). Unlike the previous version, this one is written in C++  
321 and designed for command line execution. Homo v2.0 includes corrections of errors found in  
322 the previous version, so Homo v1.3 should no longer be used. For each sequence pair, Homo  
323 executes the matched-pairs test of symmetry and returns:

- 324 • The probability ( $p$ ) of getting the test statistic by chance (assuming evolution under  
325 homogeneous conditions),
- 326 • Euclidean distance ( $\delta_{EFS}$ ) from full compositional symmetry of  $\mathbf{N}$ ,
- 327 • Euclidean distance ( $\delta_{EMS}$ ) from marginal compositional symmetry of  $\mathbf{N}$ ,
- 328 • Our distance ( $\delta_{CFS}$ ) from full compositional symmetry of  $\mathbf{N}$ .

329 If any of the observed  $p$  values is below the 5% family-wise error rate, the program prints a  
330 warning to the user on the terminal. Homo is executed using the following commands:

```
331 homo <infile> <b|f> <1|...|31>
```

332 or

```
333 homo <infile> <b|f> <1|...|31> > README
```

334 where `infile` is a text file with an alignment of characters in the fasta format, `b|f` refers to  
335 whether a brief or full report of the results should be provided, and `1|...|31` refers the data  
336 type and how these data should be analysed. If `b` is used, Homo prints one line with key  
337 statistics to the user terminal; if `f` is used, it prints five files with the values of  $p$  and  $\delta$ . A  
338 summary of the results is also be printed to the terminal.

339 Homo is designed to analyse alignments of nucleotides, di-nucleotides, codons, 10- and  
340 14-state genotypes, and amino acids. If the `infile` contains sequences of:

- 341 • Single nucleotides (4-state alphabet), the sequences may be recoded into six 3-state  
342 alphabets or seven 2-state alphabets,
- 343 • Di-nucleotides (16-state alphabet; i.e.,  $AA, AC, \dots, TG, TT$ ), the sequences may be  
344 divided into alignments with 1st or 2nd position sequences,
- 345 • Codons (a 64-state alphabet; i.e.,  $AAA, AAC, \dots, TTG, TTT$ ), the sequences may be  
346 divided into three alignments with di-nucleotide sequences and three alignments with  
347 single-nucleotide sequences,
- 348 • Amino acids (a 20-state alphabet), the letters may be recoded to a 6-state alphabet.  
349 This type of recoding was recently used to study early evolution of animals (Feuda et al.  
350 2017). Other types of recoding amino acids have been used (Kosiol et al. 2004; Susko  
351 and Roger 2007) but are not considered.

352 The 10- and 14-state genotype data cater for diploid and triploid genomes. For example, if a  
353 locus in a diploid genome contains nucleotides  $A$  and  $G$ , then the genotype sequence will  
354 contain an  $R$  at that locus. There are 10 distinguishable genotypes for each locus in diploid  
355 genomes and 14 for every locus in triploid genomes. For further detail about the data types  
356 and how the data may be analysed, simply type:

357 `homo`

358 on the command line and follow the instructions.

359 The output files from Homo fall into two categories: `.csv` files and `.dis` files. The  
360 `_Summary.csv` file contains all the estimates obtained for each pair of sequences. It can be  
361 opened and viewed by using, for example, Microsoft Excel. The `_Pvalues.csv` file contains all  
362 the  $p$  values set out in a format that can be read by HomoHeatMapper (see below). The three  
363 `.dis` files contain the  $\delta_{CFS}$ ,  $\delta_{EFS}$  and  $\delta_{EMS}$  values, and can be analysed further using  
364 FastME (Lefort et al. 2015) and SplitsTree4 (Bryant and Moulton 2004).

365 *HomoHeatMapper*.—HomoHeatMapper v1.0 is designed to generate a color-coded heat map  
366 from the `_Pvalues.csv` file. The colors used range from white (corresponding to  $p \geq 0.05$ ) to  
367 black (corresponding to  $p < 5 \times 10^{-11}$ ). HomoHeatMapper is written in Perl and can be  
368 executed using the following command:

369 `HomoHeatMapper -i <infile> -<t|f>`

370 where `infile` must be the `_Pvalues.csv` file and where `t` and `f` stand for triangle and full,  
371 respectively. The output is an `.svg` file with a heat map in scalable vector graphics format.  
372 This file can be opened and processed using Adobe Illustrator.

373 *RandTree*.—*RandTree* v1.0 is designed to generate random bifurcating trees from a set of  
374 labels. Starting from a rooted or unrooted tree with two or three tips, respectively, the tree is  
375 allowed to grow by randomly selecting tips, which will become bifurcating nodes in the tree.  
376 The probability that a tip is chosen equals  $1/n$ , where  $n$  is the number of tips in the growing  
377 tree. Thus, the probability of selecting a given tip in a 16-leaf tree is 0.0625. Having obtained  
378 a random unlabelled tree, the labels are distributed randomly across the tips.

379 *RandTree* is a command-line tool written in C++. It is executed using:

380 `randtree <infile> <r|u> <trees>`

381 where `infile` is the text file with an unique taxon label on each line, `r|u` refers to whether  
382 the random trees should be rooted or unrooted, and `trees` refers to the number of random  
383 trees to generate. Trees generated by *RandTree* are printed in the Newick format to a text file,  
384 which can be used by other phylogenetic programs.

## 385 BENCHMARKING

386 Recently, Naser-Khdour et al. (2020) applied the matched-pairs tests of symmetry (Bowker  
387 1948), marginal symmetry (Stuart 1955), and internal symmetry (Ababneh et al. 2006b) to a  
388 panel of 35 published phylogenetic data sets with the aim to measure the prevalence and  
389 impact of model misspecification. Applying an implementation of these tests in IQ-TREE  
390 (Nguyen et al. 2015), their research revealed widespread evidence of evolution under non-SRH  
391 conditions, and that this appeared to impact the accuracy of phylogenetic estimates of these  
392 data inferred assuming evolution under SRH conditions. This observation complements that  
393 of a previous simulation-based study on the adverse impact of compositional heterogeneity on  
394 phylogenetic estimates (Jermin et al. 2004).

395 We benchmarked *Homo* by comparing the result from the matched-pairs test of  
396 symmetry to those from the matched-pairs tests of symmetry, marginal symmetry, and  
397 internal symmetry, as implemented in *TestSym* (Ababneh et al. 2006b) and in *IQ-TREE*  
398 (Nguyen et al. 2015). In addition, we compared the result to that Foster’s (2004) test of  
399 homogeneity, as implemented in *p4*. We considered the alignment of `Seq1'` and `Seq2'` (Fig.  
400 3d), and asked whether it is reasonable to assume that `Seq1'` and `Seq2'` diverged under  
401 homogeneous conditions (i.e.,  $X_1 = X_2$ ). The divergence matrix,  $\mathbf{N}(t_2)$ , with its marginal  
402 frequencies, is reproduced here:

14	7	10	2	33
3	13	2	1	19
1	1	11	4	17
7	3	1	10	21
25	24	24	17	90

403 Table 1 shows the  $p$  values from different implementations of the matched-pairs tests of  
404 symmetry, marginal symmetry, and internal symmetry. As expected, Homo returned a  $p$  value  
405 identical to those returned by TestSym and IQ-TREE.

406 Foster’s (2004) test is very similar to Stuart’s (1955) matched-pairs test of marginal  
407 symmetry, so results obtained from the former test should be compared to those obtained  
408 from the latter. Assuming that evolution occurred under the GTR model, Foster’s (2004) test  
409 returned a probability of 0.150 and, if it had occurred under the F81 model, 0.157. In  
410 summary, Foster’s (2004) test returned lower probabilities than that from Stuart’s (1955)  
411 matched-pairs test of marginal symmetry (Table 1), most likely because the two tests used  
412 different approaches to assess the same null hypothesis.

413 Next, we compared the times taken by Homo v2.0 and Homo v1.4 to complete an  
414 analysis of the same data. To do so, we analysed an amino-acid alignment from Butler et al.  
415 (2009). These data—18 sequences and 412,814 sites—were analysed on a MacBook Air  
416 (Processor name: Intel Core i5; Processor speed: 1.6 GHz). Homo 2.0 completed the survey in  
417 0.43 s while Homo 1.4 completed it in 143.26 s; that is a 341-fold speedup. When Homo v2.0  
418 was used in `b` mode, the essential output was returned in 0.317 s. In conclusion, Homo v2.0 is  
419 well-tuned for large phylogenomic data sets.

## 420 BIOLOGICAL EXAMPLE

421 To illustrate the insights that may be gained by using the software presented in this paper, we  
422 surveyed an alignment of amino acids from Butler et al. (2009). The data matrix is the one  
423 used in the previous section.

### 424 *The Survey*

425 The PP plot in Figure 6a reveals that this data set is unlikely to have evolved under  
426 homogeneous conditions, but a single dot at the righthand side of the plot suggests that at  
427 least one pair of sequences have evolved under similar conditions. The heat map in Figure 6b  
428 shows that these two sequences come from *Saccharomyces cerevisiae* and *S. paradoxus*. The  
429 summary statistics for the 153 (non-independent) comparisons show that the smallest  $p$ -value

430 was 0.0, and that 99.3% of the  $p$ -values are below the 5% family-wise error rate. In summary,  
431 we conclude that the alignment has a strong compositional signal and that only two of the 18  
432 sequences appear to have evolved under the same conditions. Compositional heterogeneity is  
433 clearly a pronounced feature of these data, so it would be wise to consider this feature  
434 carefully when analysing the data phylogenetically. We note that the large number of sites  
435 here can produce very small  $p$ -values corresponding to small deviations from homogeneity.

### 436 *The Impact*

437 An obvious question arising from this discovery is whether the compositional signal is  
438 phylogenetic (i.e., whether it, on its own, is able to produce what essentially looks like a  
439 phylogenetic tree). To address this question, we analysed the data using the network- and  
440 tree-based methods described above.

441 Figure 7 depicts the compositional network inferred from the  $\mathbf{D}_{CFS}$  matrix derived  
442 from the multiple sequence alignment of amino acids published by Butler et al. (2009). The  
443 network is highly complex and treelike, with several internal splits many times longer than the  
444 alternative splits. This feature implies that the phylogenetic tree reported by Butler et al.  
445 (2009) may be affected by a strong and complex compositional signal.

446 To determine whether this is the case, we compared the tree published by Butler et al.  
447 (2009) (Fig. 8a) to the compositional tree inferred from  $\mathbf{D}_{CFS}$  (Fig. 8b). The important thing  
448 to observe here is that five of the internal edges in the two trees are identical. There is no  
449 reason to expect the two trees to be more similar or dissimilar to each other than any pair of  
450 random trees, so there may be reason to question the accuracy of the phylogenetic tree  
451 inferred by Butler et al. (2009). To ascertain whether there is reason for such concern, we  
452 compared the two trees statistically.

453 In practice, we computed  $\delta_{align}$  for the two trees in Figure 8 as well as  $\delta_{Align}^*$  for 999  
454 pairs of randomly-generated 18-tipped trees. The latter estimates were needed to generate the  
455 null distribution. Figure 8c shows that the  $\delta_{Align}$  value for the two trees lies well below the  
456 distribution of  $\delta_{Align}^*$  values for the randomly-generated trees, implying that the trees are  
457 significantly more alike than random trees are (two-tailed test,  $p < 0.002$ ). Therefore, we may  
458 now conclude that the tree topology published by Butler et al. (2009) is affected by the  
459 presence of a compositional signal in the alignment of amino acids. In other words, the tree in  
460 Figure 1 of Butler et al. (2009) may not reflect the evolution of these 18 species.

### 461 AVAILABILITY

462 Homo v2.0 is available from <http://www.github.com/ljjermin/Homo.v2.0/>.

463 HomoHeatMapper is available from <http://www.github.com/ljjermin/HomoHeatMapper/>.

464 RandTree v1.0 is available from <http://www.github.com/ljjermin/RandTree.v1.0/>.

465

## ACKNOWLEDGEMENTS

466 We are grateful to Mary Kuhner for processing the data depicted in Figure 8c and to Kenneth  
467 Wolfe for constructive feedback.

468

## BIBLIOGRAPHY

- 469 Ababneh F., Jermiin L.S., Robinson J. 2006a. Generation of the exact distribution and  
470 simulation of matched nucleotide sequences on a phylogenetic tree. *J. Math. Model. Algor.*  
471 5:291–308.
- 472 Ababneh F., Jermiin L.S., Ma C., Robinson J. 2006b. Matched-pairs tests of homogeneity  
473 with applications to homologous nucleotide sequences. *Bioinformatics* 22:1225–1231.
- 474 Aitchison J. 1986. *The Statistical Analysis of Compositional Data*. Chapman and Hall,  
475 London.
- 476 Bowker, A.H. 1948. A test for symmetry in contingency tables. *J. Am. Stat. Assoc.*  
477 43:572–574.
- 478 Bryant, D., Moulton V. 2004. Neighbor-Net: An agglomerative method for the construction of  
479 phylogenetic networks. *Mol. Biol. Evol.* 21:255-265.
- 480 Bryant, D., Galtier N., Poursat M.-A. 2005 Likelihood calculation in molecular phylogenetics.  
481 In: Gascuel O., Editor, *Mathematics evolution and phylogeny*, Oxford University Press. Inc.,  
482 New York, p 33–62.
- 483 Butler G., Rasmussen M.D., Lin M.F., Santos M.A.S., Sakthikumar S., Munro C.A., Rheinbay  
484 E., Grabherr M., Forche A., Reedy J.L., Agrafioti I., Arnaud M.B., Bates S., Brown A.J.P.,  
485 Brunke S., Costanzo M.C., Fitzpatrick D.A., de Groot P.W.J., Harris D., Hoyer L.L., Hube  
486 B., Klis F.M., Kodira C., Lennard N., Logue M.E., Martin R., Neiman A.M., Nikolaou E.,  
487 Quail M.A., Quinn J., Santos M.C., Schmitzberger F.F., Sherlock G., Shah P., Silverstein  
488 K.A.T., Skrzypek M.S., Soll D., Staggs R., Stansfield I., Stumpf M.P.H., Sudbery P.E.,  
489 Srikantha T., Zeng Q.D., Berman J., Berriman M., Heitman J., Gow N.A.R., Lorenz M.C.,  
490 Birren B.W., Kellis M., Cuomo C.A. 2009. Evolution of pathogenicity and sexual  
491 reproduction in eight *Candida* genomes. *Nature* 459:657–662.
- 492 Duchêne D.A, Duchêne S, Ho S.Y.W. 2017. New statistical criteria detect phylogenetic bias  
493 caused by compositional heterogeneity. *Mol. Biol. Evol.* 34:1529–1534.
- 494 Duchêne D.A, Duchêne S, Ho S.Y.W. 2018. Differences in performance among test statistics  
495 for assessing phylogenomic model adequacy. *Genome Biol. Evol.* 10:1375–1388.
- 496 Dziak, J.J, Coffman D.L., Lanza S.T., Li R., Jermiin L.S. 2019. Sensitivity and specificity of  
497 information criteria. *Brief. Bioinformatics* <https://doi.org/10.1093/bib/bbz016>.



- 498 Egozcue J.J., Pawlowsky-Glahn V. 2011. Basic concepts and procedures. In:  
499 Pawlowsky-Glahn V., Buccianti A., editors. Compositional Data Analysis. Chichester, John  
500 Wiley and Sons, p. 12–28.
- 501 Emsley T.M., Thomas R.H., Williams R.A.D. 1993. Reduced thermophilic bias in the 16S  
502 rDNA sequence from *Thermus ruber* provides further support for a relationship between  
503 *Thermus* and *Deinococcus*. Syst. Appl. Microbiol. 16:25–29.
- 504 Felsenstein J. 1981. Evolutionary trees from DNA sequences: a maximum likelihood approach.  
505 J. Mol. Evol. 17:368–376.
- 506 Felsenstein J. 2005. PHYLIP (Phylogeny Inference Package). Distributed by the author  
507 (<http://evolution.gs.washington.edu/phylip.html>).
- 508 Feuda R., Dohrmann M, Pett W, Philippe H, Rota-Stabelli O, Lartillot N, Wörheide G.,  
509 Pisani D. 2017. Improved modeling of compositional heterogeneity supports sponges as sister  
510 to all other animals. Curr. Biol. 27:3864–3870.
- 511 Foster P.G. 2004. Modeling compositional heterogeneity. Syst. Biol. 53:485–495.
- 512 Galtier N., Gouy M. 1995. Inferring phylogenies from DNA sequences of unequal base  
513 compositions. Proc. Natl. Acad. Sci. U.S.A. 92:11317–11321.
- 514 Gascuel O. 1997. BIONJ: an improved version of the NJ algorithm based on a simple model of  
515 sequence data. Mol. Biol. Evol. 14, 685-695.
- 516 Goldman N. 1993. Statistical tests of models of DNA substitution. J. Mol. Evol. 36:182–198.
- 517 Grundy W.N., Naylor G.J.P. 1999. Phylogenetic inference from conserved sites alignments. J.  
518 Exp. Zool. 285:128–139.
- 519 Hillis D.M., Huelsenbeck J.P., Cunningham C.W. 1994a. Application and accuracy of  
520 molecular phylogenies. Science 264:671–677.
- 521 Hillis D.M., Huelsenbeck J.P., Swofford D.L. 1994b. Hobgoblin in phylogenetics. Science  
522 269:363–364.
- 523 Ho S.Y.W., Jermiin L.S. 2004. Tracing the decay of the historical signal in biological sequence  
524 data. Syst. Biol. 53:623–637.
- 525 Huelsenbeck JP, Hillis DM. 1993. Success of phylogenetic methods in the four-taxon case.  
526 Syst. Biol. 42:247–264.
- 527 Huson D.H., Bryant D. 2006. Application of phylogenetic networks in evolutionary studies.  
528 Mol. Biol. Evol. 23:254–267.
- 529 Jayaswal V., Jermiin L.S., Robinson J. 2005. Estimation of phylogeny using a general Markov  
530 model. Evol. Bioinformatics 1:62–80.
- 531 Jayaswal V., Robinson J., Jermiin L.S. 2007. Estimation of phylogeny and invariant sites  
532 under the general Markov model of nucleotide sequence evolution. Syst. Biol. 56:155–162.

- 533 Jayaswal V., Wong T.K.F., Robinson J, Poladian L. Jermiin L.S. 2014. Mixture models of  
534 nucleotide sequence evolution that account for heterogeneity in the substitution process across  
535 sites and across lineages. *Syst. Biol.* 63:726–742.
- 536 Jermiin L.S., Ho S.Y.W., Ababneh F., Robinson J., Larkum A.W.D. 2004. The biasing effect  
537 of compositional heterogeneity on phylogenetic estimates may be under-estimated. *Syst. Biol.*  
538 53:638–643.
- 539 Jermiin L.S., Ho J.W.K., Lau K.-W., Jayaswal V. 2009. Seqvis: A tool for detecting  
540 compositional heterogeneity among aligned nucleotide sequences. In: Posada D., editor.  
541 *Bioinformatics for DNA Sequence Analysis*. Humana Press, New York. p. 65–91.
- 542 Jermiin L.S., Jayaswal V., Ababneh F.M., Robinson J. 2017. Identifying optimal models of  
543 evolution. In: Keith J., editor. *Bioinformatics: data, sequence analysis, and evolution*. Vol. 1.  
544 (2nd Edition) Humana Press, New York. p. 379–420.
- 545 Kosiol C., Goldman N., Buttimore N.H. 2004. A new criterion and method for amino acid  
546 classification. *J. Theor. Biol.* 228:97–106.
- 547 Kuhner M.K., Yamato J. 2015. Practical performance of tree comparison metrics. *Syst. Biol.*  
548 64:205–214.
- 549 Lake J.A. 1994. Reconstructing evolutionary trees from DNA and protein sequences:  
550 paraligner distances. *Proc. Natl. Acad. Sci. USA*, 91:1455–1459.
- 551 Lanave C., Pesole G. 1993. Stationary MARKOV processes in the evolution of biological  
552 macromolecules. *Binary* 5, 191–95.
- 553 Lefort V., Desper R., Gascuel O. 2015. FastME – A comprehensive, accurate and fast  
554 distance-based phylogeny inference program. *Mol. Biol. Evol.* 32:2798-800.
- 555 Lockhart P.J., Steel M.A., Hendy M.D., Penny D. 1994. Recovering evolutionary trees under a  
556 more realistic model of sequence evolution. *Mol. Biol. Evol.*, 11:605–612.
- 557 Lovell D.R., Müller W., Taylor J., Zwart A., Helliwell C. 2011. Proportions, percentages,  
558 ppm: Do the molecular biosciences treat compositional data right? In: Pawlowsky-Glahn V.  
559 & Buccianti A., editors. *Compositional Data Analysis: Theory and Applications*. John Wiley  
560 & Sons Inc, pp. 191–207.
- 561 Mooers A.Ø., Holmes E.C. 2000. The evolution of base composition and phylogenetic  
562 inference. *Trends Ecol. Evol.* 15:365–369
- 563 Morrison D.A. 2010. Using data-display networks for exploratory data analysis in  
564 phylogenetic studies. *Mol. Biol. Evol.* 27:1044–1057.
- 565 Naser-Khdour S., Minh B.Q., Zhang W., Stone E.A., Lanfear R. 2020. The prevalence and  
566 impact of model violations in phylogenetic analysis. *Genome Biol. Evol.*  
567 <https://doi.org/10.1093/gbe/evz193>

- 568 Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ. 2015. IQ-TREE: A fast and effective  
569 stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol. Biol. Evol.*,  
570 32:268–274.
- 571 Nye, T.M.W., Liò P., Gilks W.R. 2006. A novel algorithm and web-based tool for comparing  
572 two alternative phylogenetic trees. *Bioinformatics* 22:117–119.
- 573 Penny D., Foulds L.R., Henny M.D. 1982. Testing the theory of evolution by comparing  
574 phylogenetic trees constructed from five different protein sequences. *Nature* 297:197–200.
- 575 Penny D., Henny M.D., Steel M.A. 1992. Progress with methods for constructing evolutionary  
576 trees. *Trends in Ecology and Evolution* 7:73–79
- 577 Rambaut A., Grassly N.C. 1997. Seq-Gen: an application for the Monte Carlo simulation of  
578 DNA sequence evolution along phylogenetic trees. *CABIOS* 13:235–238.
- 579 Robinson D.F., Foulds L.R. 1981 Comparison of phylogenetic trees. *Math. Biosci.* 53:131–147.
- 580 Rouse G.W., Jermin L.S., Wilson N.G., Eeckhaut I., Lanterbecq D., Oji T., Young C.M.,  
581 Browning T., Cisternas P., Helgen L.E., Stuckey M., Messing C.G. 2013. Fixed, free, and  
582 fixed: the fickle phylogeny of extant Crinoidea (Echinodermata) and their Permian-Triassic  
583 origin. *Mol. Phylogenet. Evol.* 66:161–181.
- 584 Sand A., Holt M.K., Johansen J., Brodal G.S., Mailund T., Pedersen C.N.S. 2014. tqDist: a  
585 library for computing the quartet and triplet distances between binary or general trees.  
586 *Bioinformatics* 30:2079–2080.
- 587 Schweder T., Spjøtvoll E. 1982. Plots of P-values to evaluate many tests simultaneously.  
588 *Biometrika* 69:493–502.
- 589 Steel M.A. 1994. Recovering a tree from the leaf colourations it generates under a Markov  
590 model. *Appl. Math. Lett.*, 7:19–23.
- 591 Stuart A. 1955. A test for homogeneity of the marginal distributions in a two-way  
592 classification. *Biometrika*, 42:412–416.
- 593 Susko E., Roger A.J. 2007. On reduced amino acid alphabets for phylogenetic inference. *Mol.*  
594 *Biol. Evol.* 24:2139–2150.
- 595 Tavaré S. 1986. Some probabilistic and statistical problems on the analysis of DNA sequences.  
596 *Lect. Math. Life Sci.*, 17:57–86.
- 597 Vera-Ruiz V.A., Lau K.W., Robinson J., Jermin L.S. 2014. Statistical tests to identify  
598 appropriate types of nucleotide sequence recoding in molecular phylogenetics. *BMC Bioinf.* 15  
599 (Suppl. 2):S8.
- 600 Waddell P.J., Steel M.A. 1997. General time reversible distances with unequal rates across  
601 sites: mixing  $\Gamma$  and inverse Gaussian distributions with invariant sites. *Mol. Phylogenet.*  
602 *Evol.* 8:398–414.

- 603 Waddell, P.J., Cao Y., Hauf J., Hasegawa M. 1999. Using novel phylogenetic methods to  
604 evaluate mammalian mtDNA, including amino acid-invariant sites-LogDet plus site stripping,  
605 to detect internal conflicts in the data, with special reference to the positions of hedgehog,  
606 armadillo, and elephant. *Syst Biol* 48:31–53.
- 607 Williams W.T., Clifford H.T. 1971. On the comparison of two classifications of the same set of  
608 elements. *Taxon* 20:519–522.
- 609 Yang Z. 2014. *Molecular Evolution: A Statistical Approach*. Oxford University Press, Oxford.

Table 1: Probabilities of obtaining the site-pattern distribution in  $\mathbf{N}(t_2)$  by chance, assuming symmetry, marginal symmetry, and internal symmetry of the evolutionary processes. The probabilities were obtained using Homo v2.0, TestSym (Ababneh et al. 2006b) and IQ-TREE (Naser-Khdour et al. 2020).

Matched-pairs test of	Homo v2.0	TestSym v1.0	IQ-TREE v1.7
Symmetry	0.0213	0.0213	0.0213
Marginal symmetry	–	0.1836	0.1836
Internal symmetry	–	0.0183	0.0183

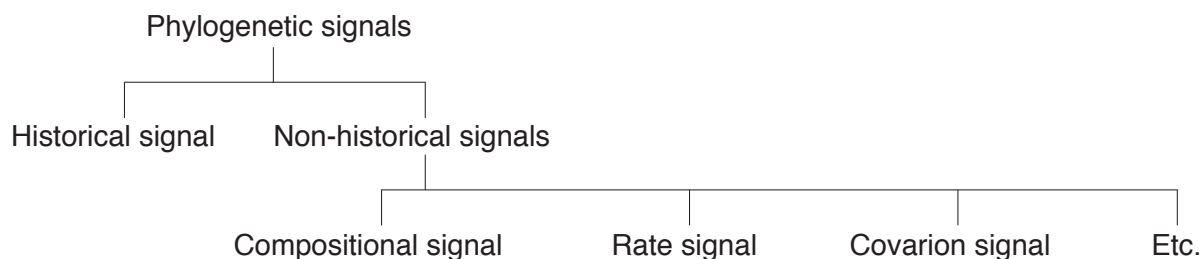
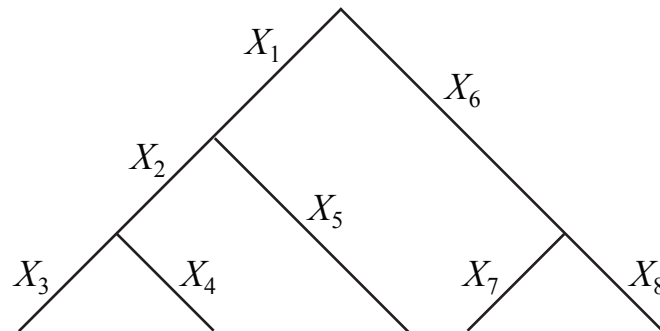


Figure 1: The phylogenetic signals (i.e., signals in phylogenetic data that, on their own, can generate a phylogeny), partitioned into some of its constituent components. Phylogenetic studies often aim to extract a historical signal from phylogenetic data. However, the accuracy of these studies depends not only on how decayed the historical signal is (Ho and Jermiin 2004) but also on whether non-historical signals have arisen over the course of time. The non-historical signals include the compositional signal (caused by non-homogeneous site patterns in the data), the rate signal (caused by independently evolving sites evolving at different rates), the covarion signal (caused by sites not evolving independently). Non-historical signals may bias phylogenetic estimates unless properly accounted for.



---

**Condition 1**      $X_i = X_j$  for all  $i \neq j$

**Implication**

Compositional heterogeneity unlikely to arise

---

**Condition 2**      $X_i \neq X_j$  for any  $i \neq j$

**Implication**

Compositional heterogeneity may arise

Figure 2: The phylogenetic challenge, illustrated using a nucleotide sequence evolving over a rooted 5-tipped tree with eight Markovian processes (i.e.,  $X_1, \dots, X_8$ ) distributed over the edges. Each site in the sequence evolving over this tree is governed by these eight edge-specific Markov processes. If  $X_i = X_j$  for all  $i \neq j$ , compositional heterogeneity across the descendant sequences is unlikely to arise. Otherwise, it may arise.

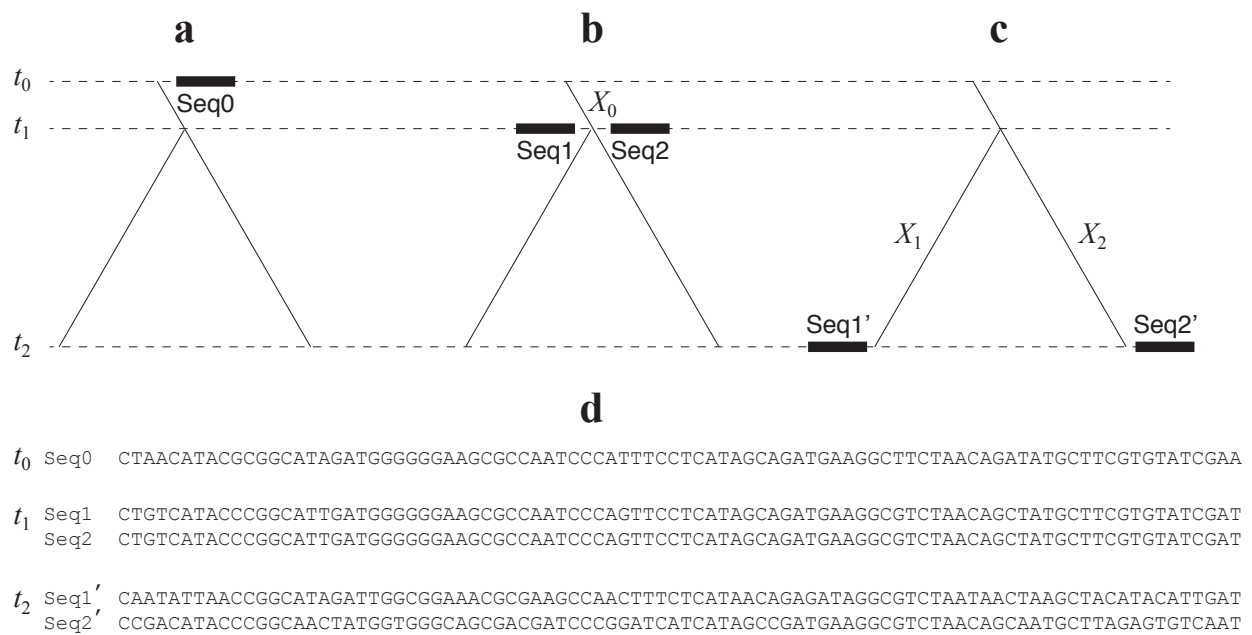


Figure 3: Rooted phylogenetic tree with the ancestral sequence evolving along the root edge (a) and, later on, at the start (b) and the end (c) of the bifurcation. The evolutionary processes operating over the three edges are marked  $X_0$ ,  $X_1$  and  $X_2$ . The corresponding sequences from the three points in time (i.e.,  $t_0$ ,  $t_1$  and  $t_2$ ) are shown in panel d.

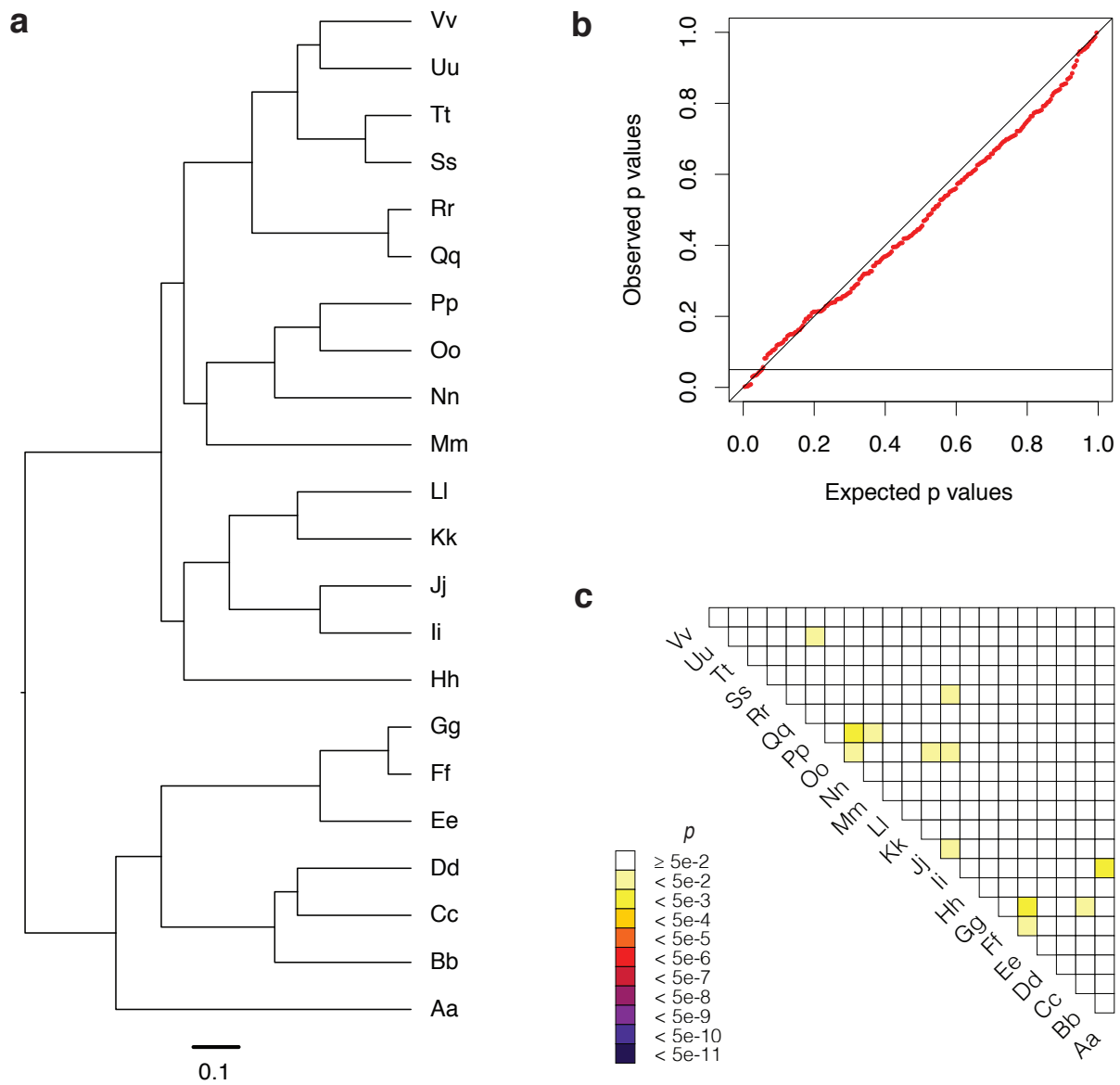


Figure 4: Using Seq-Gen (Rambaut and Grassley 1997), nucleotide sequences containing 100,000 sites were generated by simulation on a 22-tipped tree (a) under the GTR (Tavaré 1986) model of sequence evolution with the following parameters:  $S = [0.8, 0.4, 0.2, 0.1, 0.05, 0.025]$ ,  $\pi = [0.4, 0.3, 0.2, 0.1]$ ,  $pI = 0.15$ , and a continuous  $\Gamma$  distribution with  $\alpha = 2.7$ . The resulting sequences were then analysed using the matched-pairs test of symmetry. The resulting 231  $p$ -values were finally presented in a PP plot (b) and in a heat map (c).



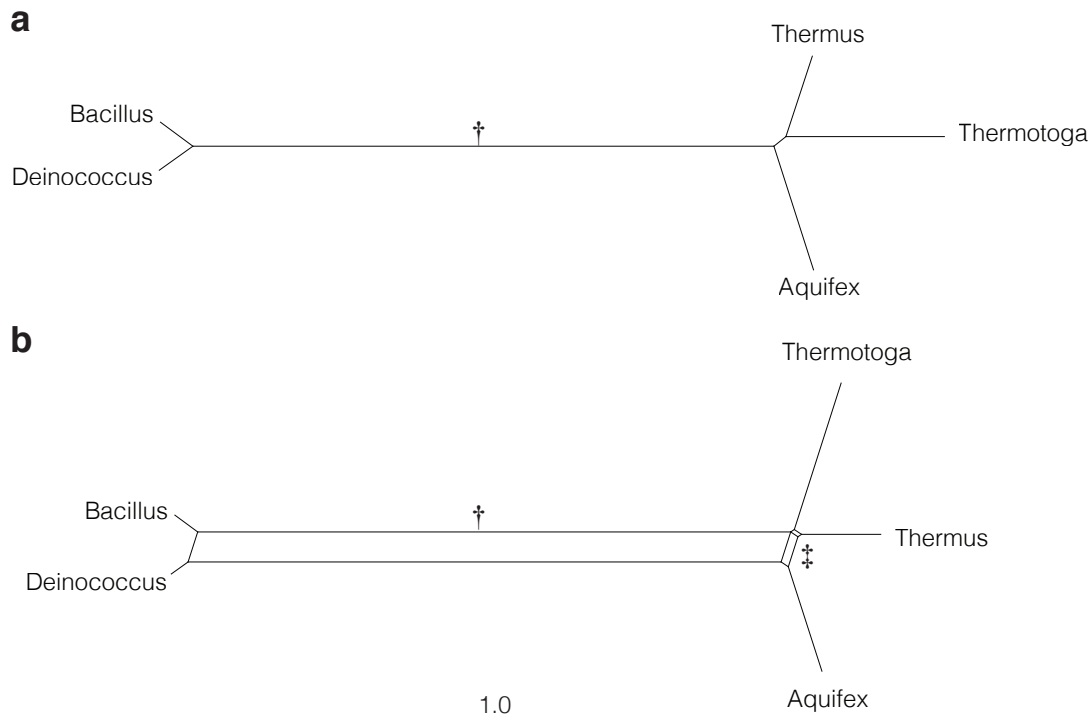


Figure 5: A compositional tree (a) and a compositional network (b), inferred from a matrix of compositional distances ( $D_{CFS}$ ) obtained from an alignment of bacterial 16S rRNA sequences. The tree and the network are drawn to scale. The characters † and ‡ point to splits that are referred to in the text.

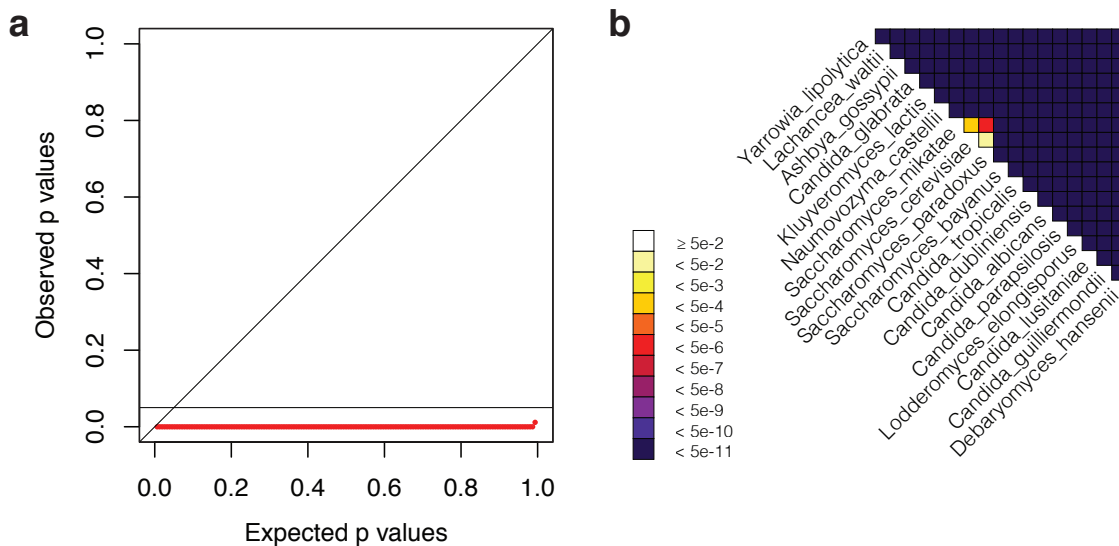


Figure 6: Visual output from our study of the alignment of amino acids from Butler et al. (2009). (a) PP plot showing that the data set, as a whole, is unlikely to have evolved under homogeneous conditions. (b) Heat map identifying the least offending sequences (*Saccharomyces cerevisiae* and *S. paradoxus*). In Butler et al. (2009), *Lachancea waltii* was called *Kluyveromyces waltii* and *Naumovozyma castellii* was called *S. castellii*.

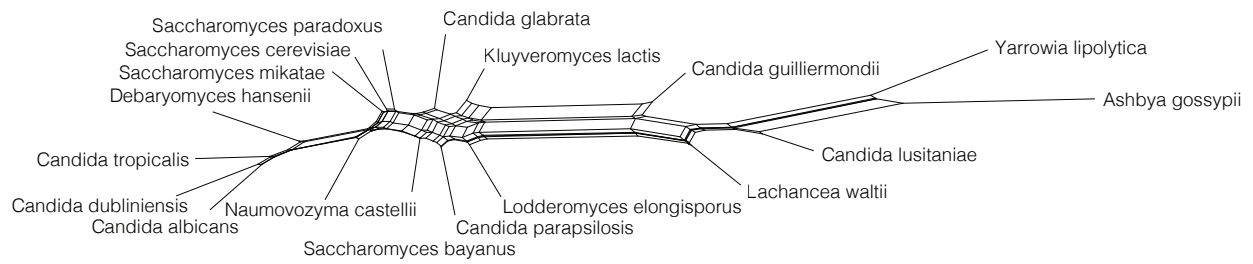


Figure 7: A compositional network inferred by SplitsTree4 (Huson and Bryant 2006) from a matrix of compositional distances ( $\mathbf{D}_{CFS}$ ) obtained from an alignment of amino acids by Butler et al. (2009). In Butler et al. (2009), *Lachancea waltii* was called *Kluyveromyces waltii* and *Naumovozyma castellii* was called *S. castellii*.

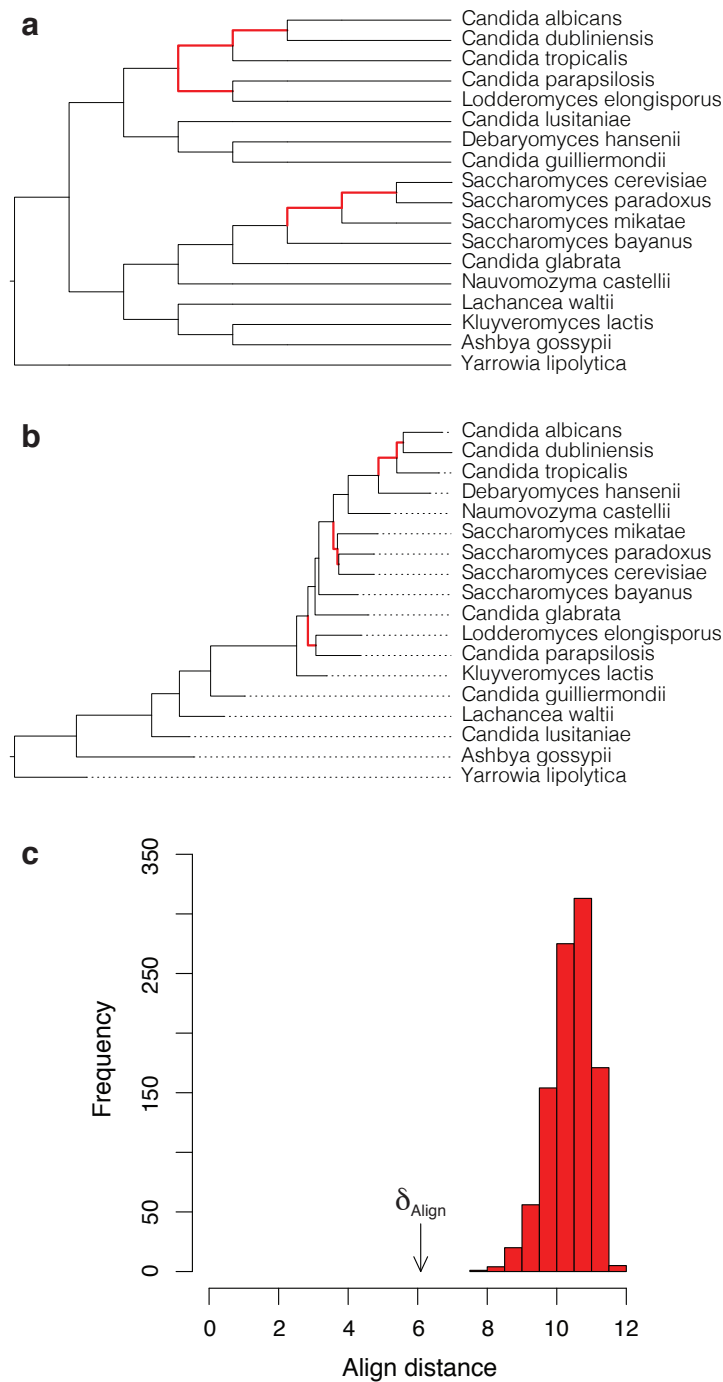


Figure 8: Figure with (a) the tree topology inferred by Butler et al. (2009), (b) the tree topology inferred from  $\mathbf{D}_{CFS}$  using the least-squares distance method implemented in PHYLIP (Felsenstein 2005), and (c) a histogram with the align distance between the trees in panels a and b (arrow) and between 999 randomly-generated pairs of 18-tipped trees (red bars). A similar result was obtained using the quartet distance (Sand et al. 2014). Identical splits in the two trees are highlighted using bold red edges. In Butler et al. (2009), *Lachancea waltii* was called *Kluveromyces waltii* and *Naumovozya castellii* was called *S. castellii*.