

# 1 **A minimum reporting standard for multiple sequence alignments**

2

3 Thomas KF Wong<sup>1,2</sup>, Subha Kalyaanamoorthy<sup>1,3</sup>, Karen Meusemann<sup>4,5,6</sup>, David K Yeates<sup>4</sup>, Bernhard  
4 Misof<sup>5</sup>, Lars S Jermiin<sup>1,2,7,8,\*</sup>

5

6 <sup>1</sup>Land & Water, CSIRO, Canberra, Australian Capital Territory, Australia.

7 <sup>2</sup>Research School of Biology, Australian National University, Canberra, Australian Capital Territory,  
8 Australia.

9 <sup>3</sup>Department of Chemistry, University of Waterloo, Waterloo, Ontario, Canada.

10 <sup>4</sup>Australian National Insect Collection, CSIRO National Research Collections Australia, Canberra,  
11 Australian Capital Territory, Australia.

12 <sup>5</sup>Zoologisches Forschungsmuseum Alexander Koenig, Bonn, Germany.

13 <sup>6</sup>Evolutionssystematik & Ökologie, Institut für Biologie I, Albert-Ludwigs-Universität Freiburg,  
14 Freiburg (Brsg.), Germany.

15 <sup>7</sup>School of Biology and Environmental Science, University College Dublin, Belfield, Ireland.

16 <sup>8</sup>Earth Institute, University College Dublin, Belfield, Ireland.

17

18 \*Correspondence should be addressed to L.S.J. ([lars.jermiin@anu.edu.au](mailto:lars.jermiin@anu.edu.au))

19

20 **Multiple sequence alignments (MSAs) play a pivotal role in studies of molecular sequence data,**  
21 **but nobody has developed a minimum reporting standard (MRS) to quantify the completeness**  
22 **of MSAs. We present an MRS that relies on four simple completeness metrics. The metrics are**  
23 **implemented in AliStat, a program developed to support the MRS. A survey of published MSAs**  
24 **illustrates the benefits and unprecedented transparency offered by the MRS.**

25

26 MSAs are widely used during annotation and comparison of molecular sequence data, allowing us  
27 to identify medically important substitutions<sup>1</sup>, infer the evolution of species<sup>2</sup>, detect lineage- and  
28 site-specific changes in the evolutionary processes<sup>3</sup> and use ancestral sequence reconstruction to  
29 engineer new enzymes<sup>4</sup>. There is a wide range of computational tools to obtain MSAs, and two of  
30 these (i.e., Clustal W<sup>5</sup> and Clustal X<sup>6</sup>) are now among the 100 most cited papers in science<sup>7</sup>.

31

32 MSAs often have gaps inserted between the nucleotides or amino acids of some of the sequences.  
33 These gaps are inserted to maximize the homology of residues from different sequences. A correct  
34 MSA is necessary for accurate genome annotation, phylogenetic inference and ancestral sequence  
35 reconstruction. However, deciding where to put the alignment gaps may be more art than science.  
36 This is because homology is defined as similarity due to historical relationships by descent<sup>8</sup>. Most  
37 of these relationships belong to the unobservable distant past, so it is impossible to measure the  
38 accuracy of most MSAs inferred from real sequence data.

39

40 Without this ability, reporting the completeness of MSAs may be the best that can be achieved. So  
41 far, the only metric sometimes used is the *percent missing data* for a sequence<sup>9</sup> or an alignment<sup>10</sup>,  
42 but neither is sufficiently transparent and insightful. To ameliorate this, we developed a minimum  
43 reporting standard (MRS) for MSAs.

44

45 The MRS uses four metrics to quantify the completeness of different attributes of MSAs. Given an  
46 MSA with  $m$  sequences and  $n$  sites, we may compute four metrics:  $\mathcal{C}_a = x_a / (m \times n)$ ,  $\mathcal{C}_r = x_r / n$ ,  
47  $\mathcal{C}_c = x_c / m$  and  $\mathcal{C}_{ij} = x_{ij} / n$ , where  $x_a$  is the number of completely specified characters<sup>11</sup> in the  
48 MSA,  $x_r$  is the number of completely specified characters in the  $r$ th sequence of the MSA,  $x_c$  is

49 the number of completely specified characters in the  $c$ th column of the MSA and  $x_{ij}$  is the number  
50 of homologous sites with completely specified characters in both sequences ( $i$  and  $j$ ). In summary,  
51  $C_a$ ,  $C_r$ ,  $C_c$  and  $C_{ij}$  measure the completeness of the alignment, the  $r$ th sequence, the  $c$ th site, and  
52 the  $i$ th and  $j$ th sequences, respectively.

53

54 The first of these metrics ( $C_a$ ) is related to the *percent missing data* used previously, but it is also,  
55 as shown in **Figure 1a**, the least useful completeness metric considered here: Alignments A and B  
56 differ greatly, but they have the same  $C_a$  value (i.e., 0.7). The  $C_r$ ,  $C_c$  and  $C_{ij}$  metrics, on the other  
57 hand, are able to detect these differences. For example, the  $C_r$  values range from 0.3 to 1.0 for  
58 Alignment A and from 0.4 to 1.0 for Alignment B, raising greater concern, from a sequence-centric  
59 perspective, about Alignment A than about Alignment B. If we were to omit any sequence from  
60 Alignment A, then it would be sensible to omit the one with the smallest  $C_r$  value. The  $C_c$  values  
61 range from 0.2 to 1.0 for Alignment A and from 0.5 to 0.8 for Alignment B. Again, there is greater  
62 concern about Alignment A than about Alignment B (due to the lower  $C_c$  scores and the greater  
63 range of values). The  $C_{ij}$  values range from 0.3 to 1.0 for Alignment A and from 0.0 to 0.9 for  
64 Alignment B. There is cause for great concern if  $C_{ij} = 0.0$  is detected as it means that sequences  $i$   
65 and  $j$  have no shared homologous sites with completely specified characters in both sequences.  
66 Evolutionary distances between such sequences cannot be estimated unless the MSA contains at  
67 least one other sequence that overlaps both  $i$  and  $j$ . When such a case occurs, the evolutionary  
68 distance between sequences  $i$  and  $j$  is called *inferred by proxy*. Currently, the prevalence of this  
69 problem is unknown.

70

71 **Figures 1b** and **1c** reveal the distributions of  $C_r$  and  $C_c$  for Alignments A and B, offering additional  
72 insight into the alignments' completeness. Conveniently, the  $C_c$  scores may be used to selectively  
73 omit the least complete sites. This *masking of sites* in MSAs is popular in phylogenetics and many  
74 methods<sup>12-20</sup> are now available. Additional information can be obtained by analyzing heat maps  
75 generated from the  $C_{ij}$  values. **Figure 1d** shows the heat maps obtained from Alignments A and B.  
76 The most obvious things to note are that in Alignment A *Tagliatelle* stands out as being the least  
77 complete sequence whereas *Capellini* and *Spaghetti* share no homologous sites with completely

78 specified nucleotides in both sequences in Alignment B. Although this was easy to detect in **Figure**  
79 **1**, it will be more difficult to do if  $n$  and/or  $m$  were larger, as is typically the case in phylogenomic  
80 data.

81

82 The benefits offered by the new completeness metrics are clear, but embedding figures like those  
83 in **Figure 1** in publications may be impractical. Alternatively, the essential details may be reported  
84 in a table (**Table 1**) or in one line (e.g., Alignment B:  $m = 10$ ,  $n = 100$ ,  $C_a = 0.7$ ,  $C_r = [0.4, 1.0]$ ,  
85  $C_c = [0.4, 0.8]$ , and  $C_{ij} = [0.0, 0.9]$ ). The closer to 1.0 the four  $C$  scores are, the more complete  
86 an alignment is. If, on the other hand, the values are closer to 0.0 than to 1.0, users may consider  
87 masking some of the sequences and/or sites before starting a phylogenetic analysis of the data.

88

89 Given their potential to inform researchers across a wide range of scientific disciplines, we argue  
90 that  $m$ ,  $n$ ,  $C_a$ ,  $C_r$ ,  $C_c$  and  $C_{ij}$  should be combined into what we henceforth call an MRS for MSAs,  
91 and that publications that report all of these values be labelled *compliant with the MRS for MSAs*.  
92 To our knowledge, this has never been done beforehand, leading to widespread ignorance about  
93 the MSAs that are relied upon in ground-breaking bio-medical research.

94

95 The MRS may be used to identify dubious MSAs in phylogenomic projects. These alignments occur  
96 in large phylogenomic projects because of a lack of awareness concerning a multitude of problems  
97 in the assembly, orthology assignment, and alignment procedures. Methods to identify odd MSAs  
98 are available<sup>2</sup>, but they are not yet fully reliable. The MRS enables a better outlier check for MSA.

99

100 Typically, MSAs comprise more sequences and sites than those in **Figure 1a**, so to facilitate using  
101 the MRS, we implemented AliStat, a fast, flexible, and user-friendly program for surveying MSAs.  
102 AliStat computes the  $C_a$ ,  $C_r$ ,  $C_c$ , and  $C_{ij}$  values from MSAs of nucleotides, di-nucleotides, codons  
103 and amino acids., AliStat lists the results on the command-line or in files that can be accessed by  
104 other programs.

105

106 The benefit of new MRS for MSAs is underlined in two surveys of large MSAs (**Table 2**). In the first  
107 case, surveying an MSA of the enzyme carboxyl/cholinesterase<sup>21</sup> revealed that some of the  $C_r$  and  
108  $C_c$  scores are closer to 0.0 than 1.0, and that at least two sequences have no homologous sites in  
109 common with completely specified characters in both sequences. Further inspection of the output  
110 files revealed large proportions of low  $C_r$  and  $C_c$  scores (**Supplementary material**), so it might be  
111 wise to mask some of the sequences and/or sites before phylogenetic analysis of these data.

112  
113 In the second case, surveying a massive concatenation of MSAs of nuclear genes<sup>22</sup> revealed a more  
114 complete alignment but also low  $C_r$ ,  $C_c$  and  $C_{ij}$  values. The presence of these values indicates that  
115 additional masking of this MSA might have been wise (see **Supplementary material**). For example,  
116 omitting the two least-complete sequences (i.e., the genera *Leucoptera* and *Pseudopostega*) could  
117 have been considered.

118  
119 The MRS for MSAs is a sound solution to a large and so-far-neglected problem: how do we report,  
120 as transparently and informatively as possible, the completeness of the MSAs used in bio-medical  
121 research? Better transparency about the completeness of MSAs is clearly needed, because MSAs  
122 represent a foundational cornerstone in many bio-medical research projects and, as revealed by  
123 the example in **Figure 1**, MSAs may look different but have the same percentage of missing data.  
124 So far, information on the completeness of MSAs used in bio-medical research has been largely  
125 absent, leaving readers unable to critically evaluate the merits of scientific discoveries made on  
126 the basis of MSAs. It is critical to recognize, and acknowledge, that many MSAs are the result of  
127 scientific procedures. Therefore, it is necessary to present the results of these procedures more  
128 transparently and comprehensively. Many scientific papers now include links to the MSAs used,  
129 but the MSAs are often so large that it is impossible to form a comprehensive picture about the  
130 completeness of these MSAs.

131  
132 Our MRS enables a radical change in scientific behavior, allowing authors to report their results  
133 more transparently, and readers the ability to critically assess discoveries made from analyses of  
134 sequence data stored in MSAs.

135

136 **METHODS**

137 Further details about AliStat are available in the online version of this paper.

138

139 **DATA**

140 Data, and code used to analyze the data, are available from <http://github.com/thomaskf/AliStat>.

141

142 **ACKNOWLEDGEMENTS**

143 We thank staff at the Australian National University and University College Dublin for their helpful  
144 feedback on the color scheme used in the heat map. Many of the respondents were color-blind.

145

146 **AUTHOR CONTRIBUTIONS**

147 L.S.J. conceived the project and wrote the first version of AliStat (in C) to conduct a pilot study of  
148 the merits of using the metrics. B.M. implemented the first Perl script to produce the heat maps,  
149 and T.K.F.W. implemented the final version of AliStat (in C++). S.K., K.M., D.K. Y. and L.S.J. tested  
150 the C++ version of AliStat, and provided constructive feedback on the software. L.S.J. drafted the  
151 paper with input from the other authors.

152

153 **COMPETING FINANCIAL INTERESTS**

154 The authors declare not competing financial interests.

155

156 **FIGURE LEGENDS**

157 **Figure 1.** Example, based on two multiple sequences alignments (**a**), illustrating the corresponding  
158 distributions of completeness scores for rows (**b**), columns (**c**), and pairs of sequences (**d**).

159

160 **Table 1.** Example of the MRS for the alignments in **Figure 1a**

Feature	Alignment A	Alignment B
Sequences	10	10
Sites	100	100
Alphabet	Nucleotides	Nucleotides
$\mathcal{C}_a$	0.7	0.7
$\mathcal{C}_r$ [min – max]	0.3 – 1.0	0.4 – 1.0
$\mathcal{C}_c$ [min – max]	0.1 – 1.0	0.4 – 0.8
$\mathcal{C}_{ij}$ [min – max]	0.3 – 1.0	0.0 – 0.9

161

162 **Table 2.** Example of the MRS for two published MSAs

Feature	Carboxyl/Colineesterase <sup>21</sup>	Lepidoptera <sup>22</sup>
Sequences	364	203
Sites	2,645	749,791
Alphabet	Amino acids	Amino acids
$\mathcal{C}_a$	0.2262	0.6422
$\mathcal{C}_r$ [min – max]	0.0106 – 0.5550	0.0609 – 0.9738
$\mathcal{C}_c$ [min – max]	0.0027 – 0.9972	0.0000 – 0.9655
$\mathcal{C}_{ij}$ [min – max]	0.0000 – 0.5550	0.0084 – 0.9672

163

164 **ONLINE METHODS**

165 The MRS for MSAs is first-of-a-kind, and AliStat, which enables compliance with the MRS for MSAs,  
166 is, to our knowledge, the first program to compute the four completeness scores presented above.  
167 AliStat is written in C++ and is available, under an CSIRO Open Source Software License Agreement  
168 (variation of the BSD / MIT License), from <http://github.com/thomaskf/AliStat/>).

169  
170 AliStat reads a text file with nucleotide, di-nucleotide, codon or amino-acid sequences, which are  
171 aligned and saved in the FASTA format. In other words, AliStat considers alphabets with four, 16,  
172 20, and 64 states. If the sequences comprise single nucleotides, the characters may be lumped to  
173 form six 3-state alphabets (i.e., CRT, AGY, ACK, GMT, AST and CGW) and seven 2-state alphabets  
174 (i.e., RY, KM, SW, AB, CD, GH and TV). If the 3- and 2-state alphabets are used, the letters R, Y, K,  
175 M, S, W, B, D, H and V are considered completely specified characters (unlike normal practice<sup>11</sup>).

176  
177 AliStat can be run in two modes: Brief mode or Full mode. Execution in brief mode is done using  
178 the following command:

179  
180 `alistat <infile> <data type> -b`

181  
182 and results in the following output format are printed to the terminal:

183  
184 File name, #seqs, #sites,  $C_a$ ,  $\max C_r$ ,  $\min C_r$ ,  $\max C_c$ ,  $\min C_c$ ,  $\max C_{ij}$ ,  $\min C_{ij}$

185  
186 The brief-mode execution was included to allow users to quickly obtain the essential values from a  
187 great number of alignments (e.g., when comparing genomes phylogenetically).

188  
189 The full-mode execution (default option) allows other options to be used and is intended when a  
190 more detailed examination of an MSA is required. For example, the `-t` option is used to indicate  
191 what types of  $C$  scores should be printed in output files, the `-m` option is used to set a threshold  
192 for masking sites, and the `-i` option is used to indicate that a heat map is needed. Other options



193 and how all of the options may be used are described in the AliStat manual. The same information  
194 can be obtained by typing

195

```
196     alistat -h
```

197

198 in the command-line.

199

200 The output files appear in the .txt, .csv, .R, .dis, .svg, and .fst formats, which can be processed by  
201 other software packages. The .txt file summarizes the results. The .csv files present the  $\mathcal{C}_c$  scores  
202 and may be examined using R. For example, if a user wishes to generate a histogram of the  $\mathcal{C}_c$   
203 scores, the Table\_2.csv file may be analyzed using the Histogram\_Cr.R file. In some cases, users  
204 may want to infer a tree or network based on the  $\mathcal{C}_{ij}$  score (or the  $\mathcal{J}_{ij}$  score, where  $\mathcal{J}_{ij} = 1.0 -$   
205  $\mathcal{C}_{ij}$ ). In such cases, .dis files may be analyzed by, for example, SplitsTree<sup>23</sup>. The heat map, which  
206 may be triangular or square, is stored in the .svg file and may be opened using Adobe Illustrator™.  
207 If the `-m` option is used, the original MSA is split into two, with all sites having a  $\mathcal{C}_c$  score larger  
208 than a user-specified threshold saved in a file called Mask.fst and the other sites saved in a file  
209 called Disc.fst. The two .fst files may be analyzed separately by other means (e.g., phylogenetic  
210 programs).

211

212

213 1 Higgs, D. R. & Wood, W. G. *Proc. Natl. Acad. Sci. USA* **105**, 11595-11596 (2008).

214 2 Misof, B. *et al. Science* **346**, 763-767 (2014).

215 3 Jayaswal, V., Wong, T. K. F., Robinson, J., Poladian, L. & Jermiin, L. S. *Syst. Biol.* **63**, 726-742  
216 (2014).

217 4 Wilding, M. *et al. Green Chem.* **19**, 5375-5380 (2017).

218 5 Thompson, J. D., Higgins, D. G. & Gibson, T. J. *Nucl. Acid. Res.* **22**, 4673-4680 (1994).

219 6 Thompson, J. D., Gibson, T. J., Plewniak, F., Jeanmougin, F. & Higgins, D. G. *Nucl. Acid. Res.*  
220 **25**, 4876-4882 (1997).

221 7 Van Noorden, R., Maher, B. & Nuzzo, R. *Nature* **514**, 550-553 (2014).

222 8 Morrison, D. A. *Syst. Bot.* **40**, 14-26 (2015).

223 9 Wiens, J. J. *Syst. Biol.* **52**, 528-538 (2003).

224 10 Driskell, A. C. *et al. Science* **306**, 1172-1174 (2004).

225 11 Cornish-Bowden, A. *Nucl. Acid. Res.* **13**, 3021-3030 (1985).

226 12 Castresana, J. *Mol. Biol. Evol.* **17**, 540-552 (2000).

227 13 Talavera, G. & Castresana, J. *Syst. Biol.* **56**, 564-577 (2007).

228 14 Dress, A. W. M. *et al. Algor. Mol. Biol.* **3**, 7 (2008).

229 15 Hartmann, S. & Vision, T. J. *BMC Evol. Biol.* **8**, 95 (2008).

230 16 Capella-Gutierrez, S., Silla-Martinez, J. M. & Gabaldon, T. *Bioinformatics* **25**, 1972-1973  
231 (2009).

232 17 Misof, B. & Misof, K. *Syst. Biol.* **58**, 21-34 (2009).

233 18 Kück, P. *et al. Front. Zool.* **7**, 10 (2010).

234 19 Criscuolo, A. & Gribaldo, S. *BMC Evol. Biol.* **10**, 210 (2010).

235 20 Wu, M. T., Chatterji, S. & Eisen, J. A. *PLoS One* **7**, e30288 (2012).

236 21 Pearce, S. L. *et al. BMC Biol.* **15**, 63 (2017).

237 22 Kawahara, A. Y. *et al. Proc. Natl. Acad. Sci. USA* **116**, 22657-22663 (2019).

238 23 Huson, D. H. & Bryant, D. *Mol. Biol. Evol.* **23**, 254-267 (2006).

**a****Alignment A**

```

Linguine      TCGCAGGATCGTATAGGAGGTGCTTTACGGCCAATATAAAGAGAGCCGTAAGAGTGATTCTGCCAAGCCCAACCTAACGCCGTTTTAAGAGCTGGGCTTC
Tagliatelle  -----ACCCGATCGTGCTGGTAGCCATATGCTCC
Fettuccine   -----GAGATAAGGAAACATGATTCCTTCCAACCTCAACGGATCGCTGCTTTTAGCGATATGCTCC
Capellini    -----AAGAGCAGGTGCCGCGATTCTGGTAATGTCATCAGTACGTCGCTTTTTTCAGAGGCGCTTC
Spaghetti    -----TATAGCAGATGCCGCGATTTCAGTTAGGGTCACTAGAACGTCGCTTTTAGAGATGCGCTTC
Vermicelli   -----TAGAGCAGGAGCCAAGTTTCTGCTAACATCATCGGACCATCGCTTTTTGGGATGCGCCTC
Lasagne      -----CCTATAAACCCAGACCACCAGGGTGACTTTTTAATCTTATCGGAACGTTACTGTAGAGCTGGGCACC
Farfallini   -----TAGCCCGCGCTTGGCAGAATAGTGCATGAGGCGTGATTCTGTTTGTAGTTTGAACAAACTTTTTGTTACTTAGCCGGAGTCC
Tortellini   -----TGACCCTCCGCCGATAAAACCAGTGGAGAGGCATGATTCTGGTCTACTTACCGAACGTAGCTATTAGAGTTAAGATCC
Ravioli      CCGTGGAGCGGCCTTGAAACTGACCACCCTCGACATATACAGGGCAAGAGACGTGATTCTGTTGCGCTCTTCTGACGTTGCTATTACCGTTGGGCACC

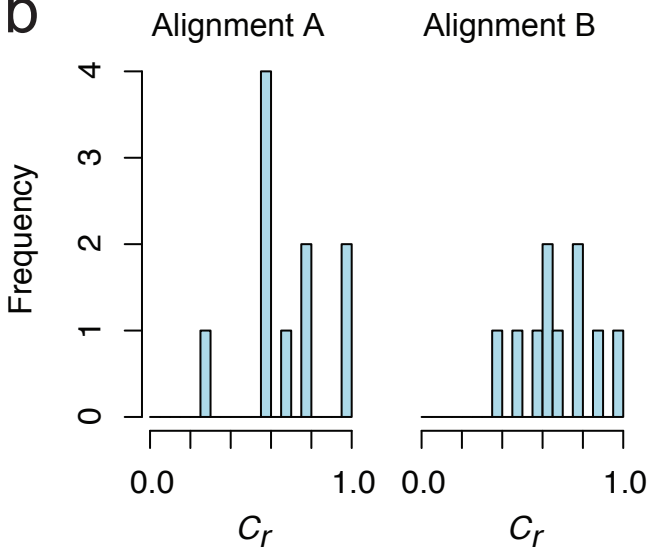
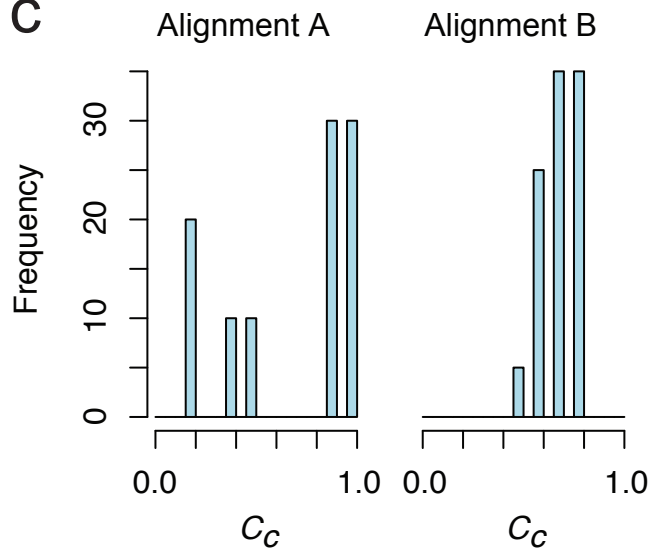
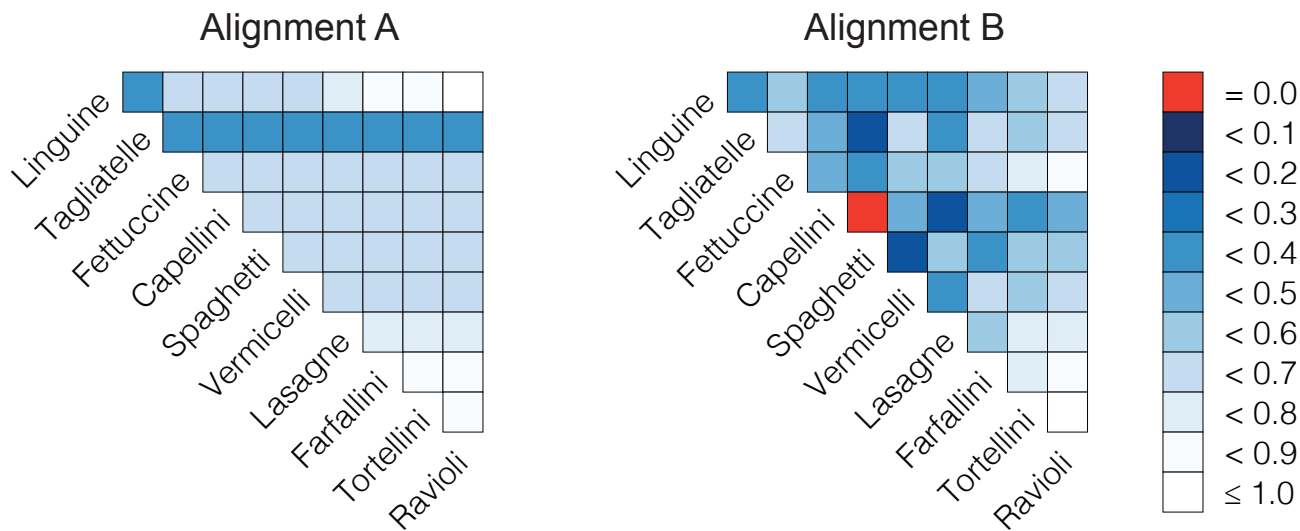
```

**Alignment B**

```

Linguine      CCAATGAACAAAACCGACCCAGGCCGACACGGTGA-----AACCTTTATCCACGGATCCCAGCATCAGA
Tagliatelle  GCAGG-----TAGTACTGGGCACACAACGAGCTTTCACGCGCATGTTTCATCCACGCGTTCGCAACATCGGA
Fettuccine   GCAGGCGGCTTTTTAGCCAGCCATA-----CTGGGCACACCACGAGCGTCATGCTCATGTTCAACCACGCGTTCGCGGCATCGGA
Capellini    -----CAGACGTGCTGAGATTGGTCCATGCGTTCGCCAGAATCGCA
Spaghetti    CCAGCGCGTCTCACAAACCAACCGACTCGACCGTGCCAAGATAAAAACCGGT-----
Vermicelli   -----TAAAACCGGTCACATGACGGCAGTCGTCGCTGGGATTTGTCCATGCGTTCGCCAGAATCTCA
Lasagne      CCAGGTCGGCCACCTAGCGACCCGGCACGTGGTGACTAGCTCATACTAGGCACACGACGATATTCGAGTT-----
Farfallini   -----ACGATACGCCTTGGATAAAGACTTTTCTGACAGACAACCTGAAGTCGTCAGGTCATTTCATCCCTTGCCTCCCATCAACCCA
Tortellini   CTAAGCGACGTAAGGGGCCCCCGATACCCCGGCGAGTACCAATACCCGGGATCGCGACGAGCGACGCTGTGTGATTTGTCGATTCGTTAG-----
Ravioli      CCACGTCACCTTACGGCCGGCCGAGACGACCACGAAATGCCAGTACTGGATACTATAAGCGTGTTTGTGTGTCAGTCAACCATTTCGATTGCATCATCGCA

```

**b****c****d**

1 **Supplementary Material:**

2

3 **A minimum reporting standard for multiple sequence alignments**

4

5 Thomas KF Wong<sup>1,2,#</sup>, Subha Kalyanamoorthy<sup>1,3,#</sup>, Karen Meusemann<sup>4,5,6</sup>, David K Yeates<sup>4</sup>,

6 Bernhard Misof<sup>5</sup>, Lars S Jermiin<sup>1,2,7,8,\*</sup>

7

8 <sup>1</sup>Land & Water, CSIRO, Canberra, Australian Capital Territory, Australia.

9 <sup>2</sup>Research School of Biology, Australian National University, Canberra, Australian Capital

10 Territory, Australia.

11 <sup>3</sup> Department of Chemistry, University of Waterloo, Waterloo, Ontario, Canada.

12 <sup>4</sup>Australian National Insect Collection, CSIRO National Research Collections

13 Australia, Canberra, Australian Capital Territory, Australia.

14 <sup>5</sup>Zoologisches Forschungsmuseum Alexander Koenig, Bonn, Germany.

15 <sup>6</sup>Evolutionsbiologie & Ökologie, Institut für Biologie I, Albert-Ludwigs-Universität Freiburg,

16 Freiburg (BrsG.), Germany.

17 <sup>7</sup>School of Biology and Environmental Science, University College Dublin, Belfield, Ireland.

18 <sup>8</sup>Earth Institute, University College Dublin, Belfield, Ireland.

19

20

21 #Joint first authors (these authors contributed equally to this work)

22

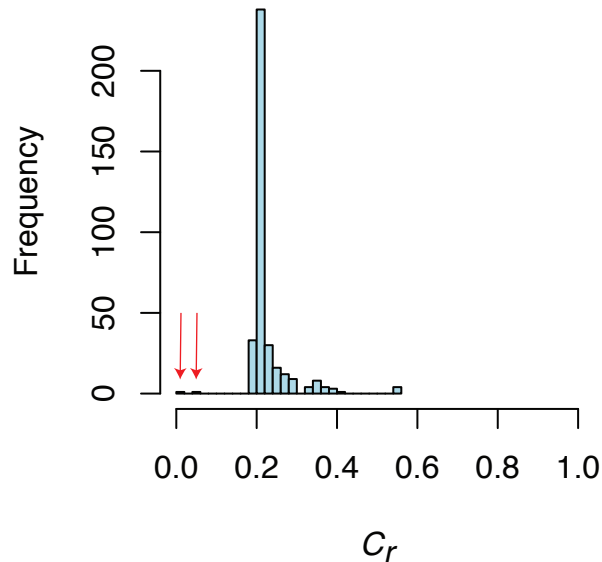
23 \*Correspondence should be addressed to L.S.J. ([lars.jermiin@anu.edu.au](mailto:lars.jermiin@anu.edu.au))

24

25 **Analysis of an alignment of carboxyl/cholinesterases (CCEs) from a paper by Pearce et al.<sup>1</sup>**

26 The alignment of amino acids used to annotate the CCE genes from *Helicoverpa armigera*, *H.*  
27 *zea*, *Manduca sexta* and *Bombyx mori* was surveyed using AliStat v1.11. The alignment  
28 comprised 364 sequences and 2645 sites. **Figure S1** presents the distribution of  $C_r$  values,  
29 with several sequences having values close to 0.0. Because the objective of the study by  
30 Pearce et al.<sup>1</sup> was to annotate the genes, it was not possible to remove any sequences from  
31 the data.

32



33

34

35

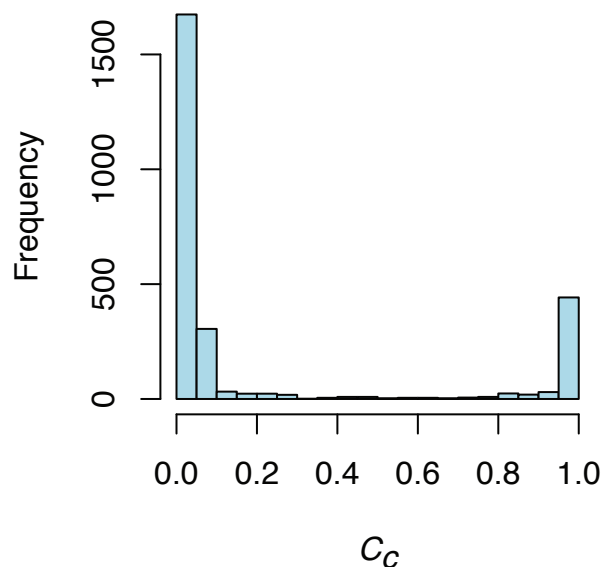
36

37

**Figure S1.** Histogram showing the distribution of  $C_r$  scores from the CCE genes. The arrows point to the lowest  $C_r$  scores.

38 **Figure S2** reveals the distribution of  $C_c$  scores, with a high proportion of sites with low  $C_c$   
39 values. Based on this distribution, sites with  $C_c \leq 0.5$  were masked in the study by Pearce et  
40 al.<sup>1</sup>.

41



42

43 **Figure S2.** Histogram showing the distribution  
44 of  $\mathcal{C}_c$  scores from the CCE genes.

45  
46 **Figure S3** shows the heat map of  $\mathcal{C}_{ij}$  scores, revealing a high proportion of sequence pairs  
47 with low  $\mathcal{C}_{ij}$  values. Because the aim of the study by Pearce et al.<sup>1</sup> was to annotate the genes,  
48 it was not possible to remove any sequences from the data.



50 **Figure S3.** Heat map showing the distribution of  $\mathcal{C}_{ij}$  scores from the CCE genes. The benefit of  
51 this heat map is realized by enlarging the image, at which point it become obvious that two  
52 sequences, labelled HzeaCCE016h and MsexCCE001s, have no homologous sites with  
53 unambiguous characters in both sequences.  
54

55  
56 Pearce et al.<sup>1</sup> used a threshold of  $\mathcal{C}_c = 0.5$  to mask the alignment of amino acids. **Table 1**  
57 reveals the effect of doing so. In this case, the table complies with the minimum reporting  
58 standard (MRS) for multiple sequence alignments (MSAs).  
59

60 **Table 1.** Example highlighting the effect of using  $\mathcal{C}_c =$   
61 0.5 to mask the sites in the alignment of CCEs.

Feature	Before masking	After masking
Sequences	364	364
Sites	2,645	546
Alphabet	Amino acids	Amino acids
$\mathcal{C}_a$	0.2262	0.9562
$\mathcal{C}_r$ [min – max]	0.0106 – 0.5550	0.0458 – 0.9890
$\mathcal{C}_c$ [min – max]	0.0027 – 0.9972	0.5055 – 0.9972
$\mathcal{C}_{ij}$ [min – max]	0.0000 – 0.5550	0.0000 – 0.9890

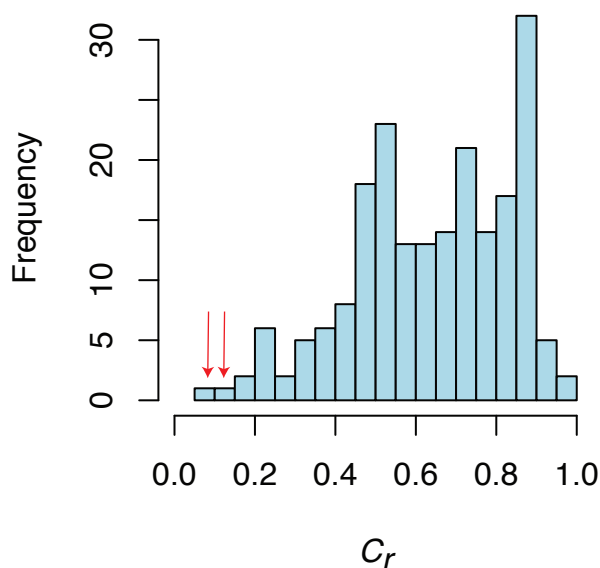
62

63 **Analysis of lepidopteran nuclear data from a study by Kawahara et al.<sup>2</sup>**

64 The alignment of amino acids labelled SA4\_aminoacid\_supermatrix\_resorted\_renamed.fas  
65 was surveyed using AliStat v1.11. The alignment comprised 203 sequences and 749,791 sites.

66 **Figure S4** reveals the distribution of  $\mathcal{C}_r$  values, with most sequences having values close to  
67 1.0. In this case, only a few sequences had very low  $\mathcal{C}_r$  scores.

68



69

70

71

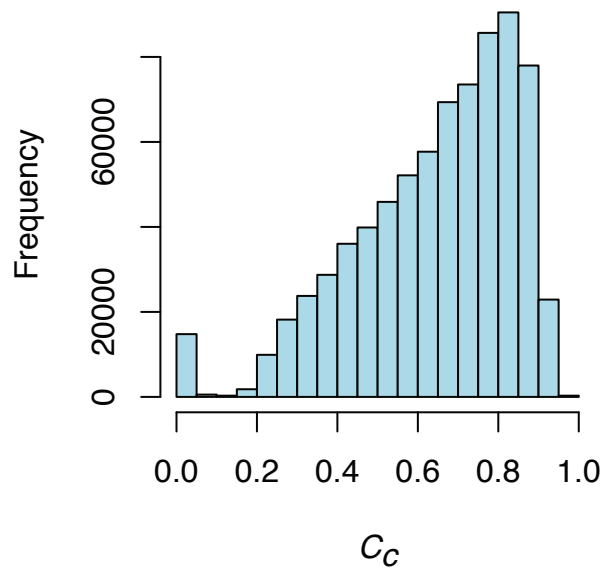
72

73

74

75 **Figure S5** reveals the distribution of  $\mathcal{C}_c$  scores, with a high proportion of sites with high  $\mathcal{C}_c$   
76 values. Based on this distribution, omitting sites with  $\mathcal{C}_c \leq 0.2$  might have been sufficient.

77



78  
79  
80  
81

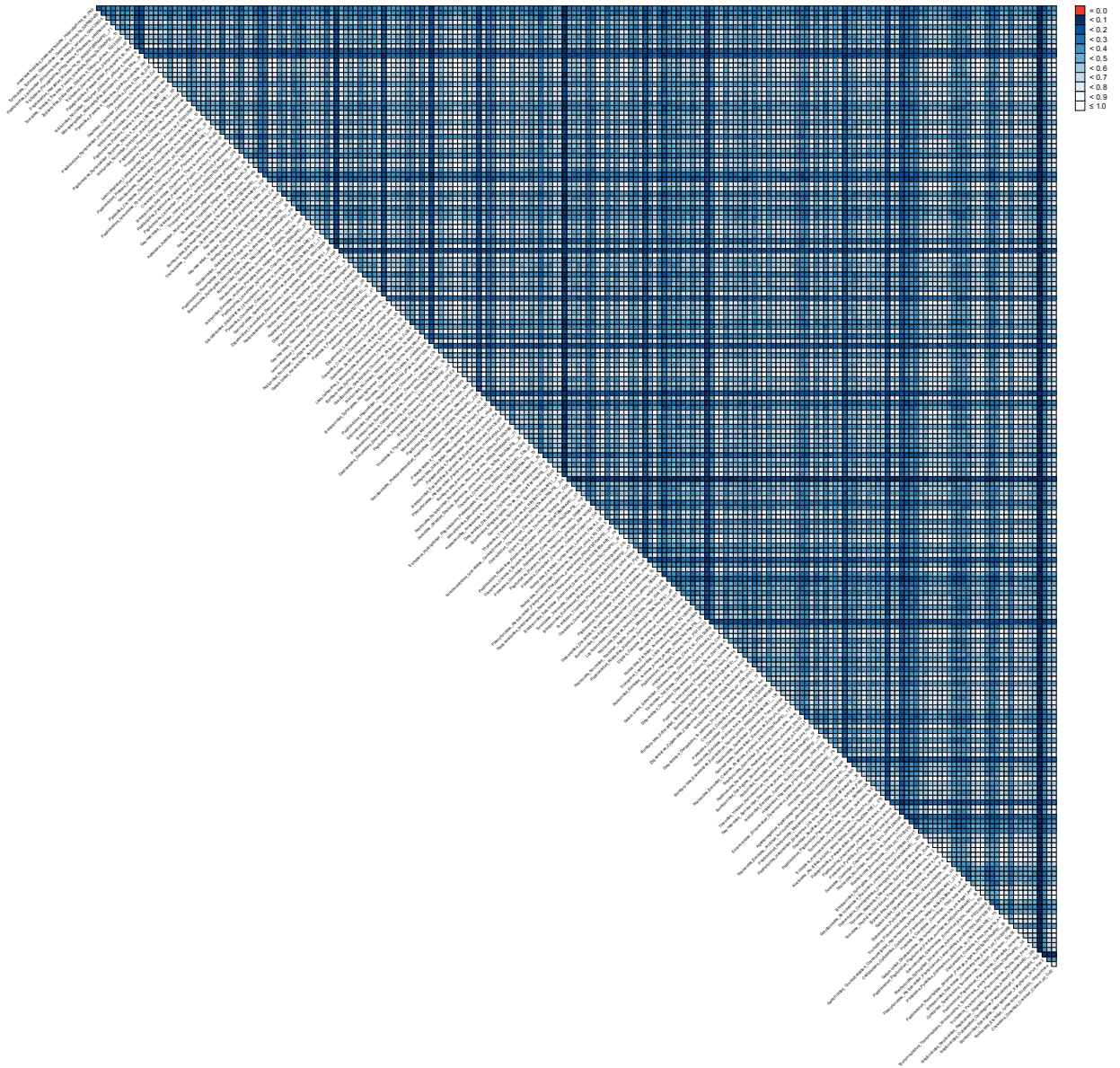
**Figure S5.** Histogram showing the distribution of  $C_c$  scores from the super-alignment (Kawahara et al.<sup>2</sup>).

82

83 **Figure S6** shows the heat map of  $C_{ij}$  scores, revealing a low proportion of sequence pairs with  
84 low  $C_{ij}$  values. In this case, the sequences found to be least complete are from the genera  
85 *Leucoptera* and *Pseudopostega*.

86





87  
88 **Figure S6.** Heat map showing the distribution of  $C_{ij}$  scores from the concatenate gene  
89 alignments. Again, the benefit of this heat map is realized by enlarging the image. In this case,  
90 the two most incomplete sequences are identified as being from the genera *Leucoptera* and  
91 *Pseudopostega*.  
92

### 93 **References**

- 94 1 Pearce, S. L. *et al. BMC Biol.* **15**, 63 (2017).  
95 2 Kawahara, A. Y. *et al. Proc. Natl. Acad. Sci. USA* **116**, 22657-22663 (2019).

96  
97