

1    **CIAAlign - A highly customisable command line tool to clean, interpret and**  
2                    **visualise multiple sequence alignments.**

3                    Charlotte Tumescheit, Andrew E. Firth, Katherine Brown\*

4  
5    Department of Pathology, University of Cambridge, Cambridge, UK

6    \*corresponding author, kab84@cam.ac.uk

7  
8    **Keywords:** Multiple sequence alignment, alignment quality, Python tool, comparative  
9    genomics, transcriptomics, phylogenetics

## 10 **Abstract**

## 11 **Background**

12 Throughout biology, multiple sequence alignments (MSAs) form the basis of much  
13 investigation into biological features and relationships. These alignments are at the heart of  
14 many bioinformatics analyses. However, sequences in MSAs are often incomplete or very  
15 divergent, which leads to poorly aligned regions or large gaps in alignments. This slows down  
16 computation and can impact conclusions without being biologically relevant. Therefore,  
17 cleaning the alignment by removing these regions can substantially improve analyses.

18

## 19 **Results**

20 We present a comprehensive, user-friendly MSA trimming tool with multiple visualisation  
21 options. Our highly customisable command line tool aims to give intervention power to the  
22 user by offering various options, and outputs graphical representations of the alignment  
23 before and after processing to give the user a clear overview of what has been removed.

24 The main functionalities of the tool include removing regions of low coverage due to  
25 insertions, removing gaps, cropping poorly aligned sequence ends and removing sequences  
26 that are too divergent or too short. The thresholds for each function can be specified by the  
27 user and parameters can be adjusted to each individual MSA. CIAAlign is complementary to  
28 existing alignment trimming tools, with an emphasis on solving specific and common  
29 alignment problems and on providing transparency to the user.

30

## 31 **Conclusion**

32 CIAAlign effectively removes poorly aligned regions and sequences from MSAs and provides  
33 novel visualisation options. This tool can be used to improve the alignment quality for further  
34 analysis and processing. The tool is aimed at anyone who wishes to automatically clean up  
35 parts of an MSA and those requiring a new, accessible way for visualising large MSAs.

## 36 **Introduction**

37 Throughout biology, multiple sequence alignments (MSAs) of DNA, RNA or amino acid  
38 sequences are often the basis of investigation into biological features and relationships.  
39 Applications of MSAs include, but are not limited to transcriptome analysis, in which  
40 transcripts may need to be aligned to genes; RNA structure prediction, in which an MSA  
41 improves results significantly compared to predictions based on single sequences; and  
42 phylogenetics, where trees are usually created based on MSAs. There are many more  
43 applications of MSA at a gene, transcript and genome level involved in a huge variety of  
44 traditional and new approaches to genetics and genomics, many of which could benefit from  
45 the tool presented here.

46 An MSA typically represents three or more DNA, RNA or amino acid sequences, which  
47 represent partial or complete gene, transcript, protein or genome sequences. These  
48 sequences are aligned by inserting gaps between residues to bring more similar residues  
49 (either based on simple sequence similarity or an evolutionary model) into the same column,  
50 allowing insertions, deletions and differences in sequence length to be taken into account [1,  
51 2]. The first widely used automated method for generating MSAs was CLUSTAL [2] and more  
52 recent versions of this tool are still in use today, along with tools such as MUSCLE [3],  
53 MAFFT [4], T-Coffee [5] and many more. The majority of tools are based upon various  
54 heuristics used to optimise progressive sequence alignment using a dynamic programming  
55 based algorithm such as the Needleman-Wunsch algorithm [6]. It has been shown previously  
56 that removing divergent regions from an MSA improves the resulting phylogenetic tree [7].  
57 Various tools are available to remove or improve poorly aligned columns, including trimAl [8],  
58 Gblocks [7] and various refinement methods incorporated into alignment software [3, 4].

59 Some tree building software can also take into account certain discrepancies in the alignment,  
60 for example RaXML [9] can account for missing data in some columns and check for  
61 duplicate sequence names and gap-only columns; similarly GUI based toolkits for molecular  
62 biology such as MEGA [10] sometimes have options to delete or ignore columns containing  
63 gaps. However, several common issues affect the speed, complexity and reliability of specific  
64 downstream analyses but are not addressed by existing tools.

65 Clean and Interpret Alignments (CIA<sub>Align</sub>) is primarily intended to address four issues which  
66 are commonly encountered when working with MSAs. Although these issues are often trivial  
67 to solve via editing alignments by eye, manual editing tends to have poor reproducibility and  
68 becomes a bottleneck with batch analyses involving large numbers of MSAs. The need to  
69 resolve these issues depends on the downstream application of the MSA.

70 The first issue we intend to address is that it is common for an MSA to contain more gaps  
71 towards either end than in the body of the alignment. This problem occurs at both the  
72 sequencing and alignment stage. For example, the ends of *de novo* assembled transcripts  
73 tend to have lower read coverage [11] and therefore have a higher probability of mis-  
74 assembly and therefore mis-alignment. MSAs created using these sequences therefore also  
75 have regions of lower reliability towards either end. Similarly, both Sanger sequences and  
76 sequences generated with Oxford Nanopore's long read sequencing technology, which are  
77 often used directly in MSAs, tend to have lower quality scores at the either the beginning or  
78 the end [12, 13, 14]. Automated removal of these regions from MSAs would therefore  
79 increase the reliability of downstream analyses. Also, while generating an MSA, terminal  
80 gaps complicate analysis, and the weighting of terminal gaps relative to internal gap opening

81 and gap extension penalties can make a large difference to the resulting alignment [15]. This  
82 again leads to regions of ambiguity and therefore gaps towards the ends of the alignment.

83 Secondly, insertions or other stretches of sequence can be present in a minority of sequences  
84 in an MSA, leading to large gaps in the remaining sequences. For example, alignments of  
85 sections of bacterial genomes often result in long gaps representing genes which are absent  
86 in the majority of species. These gaps can be observed, for example, in multiple genome  
87 alignments shown in Tettelin et al. 2005 [16] for *Streptococcus agalactiae* and Hu et al. 2011  
88 [17] for *Burkholderia*, amongst others, which show many genes which are present in only a  
89 few genomes. While these regions are of interest in themselves and certainly should not be  
90 excluded from all further analysis, they are not relevant for every downstream analysis. For  
91 example, a consensus sequence for these bacteria would exclude these regions and their  
92 presence would increase the time required for phylogenetic analysis without necessarily  
93 adding any additional information. Large gaps in some sequences may also result from  
94 missing data, rather than true biological differences and, if this is known to be the case, it is  
95 often appropriate to remove these regions before performing phylogenetic analysis [18].

96 Thirdly, one or a few highly divergent sequences can heavily disrupt the alignment and  
97 therefore complicate downstream analysis. It is very common for an MSA to include one or a  
98 few outlier sequences which do not align well with the majority of the alignment. One example  
99 of this is metagenomic analyses identifying novel sequences in large numbers of datasets. It  
100 is common to manually remove phylogenetic outliers which are unlikely to truly represent  
101 members of a group of interest (see for example [19–21]) but this is not feasible when  
102 processing large numbers of alignments.

103 Finally, very short partially overlapping sequences cannot always be reliably aligned using  
104 standard global alignment algorithms. It is very common to remove these sequences,  
105 manually or otherwise, prior to further analysis.

106 There are also several common issues in alignment visualisation. Large alignments can be  
107 difficult to visualise and a small and concise but accurate visualisation can be useful when  
108 presenting results, so this has been incorporated into the software. With many alignment  
109 trimming tools it can be difficult to track exactly which changes the software has made, so a  
110 visual output showing these changes is generated.

111 Finally, transparency is often an issue with bioinformatics software, with poor reporting of  
112 exactly how a file has been processed [22–24]. CIAAlign has been developed to process  
113 alignments in a transparent manner, to allow the user to clearly and reproducibly report their  
114 methodology.

115 CIAAlign is freely available at [github.com/KatyBrown/CIAAlign](https://github.com/KatyBrown/CIAAlign).

116

## 117 **Materials and Methods**

118 CIAAlign is a command line tool implemented in Python 3. It can be installed either via pip3 or  
119 from GitHub and is independent of the operating system. It has been designed to enable the  
120 user to remove poorly aligned regions and sequences from an MSA, to visualise the MSA  
121 (including a markup file showing which regions and sequences have been removed), and to  
122 interpret the MSA in several ways. CIAAlign works on nucleotide or amino acids alignments  
123 and will detect which of these is provided. A log file is generated to show exactly which

124 sequences and positions have been removed from the alignment and why they were  
125 removed. Users can then adjust the software parameters according to their needs.

126 CIAAlign takes as its input any pre-computed MSA in FASTA format containing at least three  
127 sequences. Most MSAs created with standard alignment software will be of an appropriate  
128 scale, for example single or multi-gene alignments and whole genome alignments for many  
129 microbial species. Measurements on the runtime were conducted for MSAs created by  
130 randomly drawing equally probable nucleotides and adding gap regions such that each MSA  
131 has a certain proportion of gaps. When running CIAAlign with all core functions (cleaning  
132 functions and creating mini alignments for input, output and the markup) and for fixed gap  
133 proportions, the runtime scales quadratically with the size of the MSA, i.e. with  $n$  as the  
134 number of sequences and  $m$  the length of the MSA, the worst case time complexity is  
135  $O((nm)^2)$ . Further runtime measurements were taken for running CIAAlign with the core  
136 functions on an MSA of constant size with different numbers of gaps. The runtime decreases  
137 linearly with increasing numbers of gaps. It should be noted that, besides the size of the MSA  
138 and its gap content, the runtime is impacted by which combination of functions is applied. For  
139 very long MSAs the size of the final image becomes a limiting factor when creating a  
140 sequence logo, as the matplotlib library [25] has restrictions on the number of pixels in one  
141 object. We have provided detailed instructions about this limit in the “Guidelines for using  
142 CIAAlign” on the CIAAlign Github.

143 The path to the alignment file is the only mandatory parameter. Every function is run only if  
144 specified in the parameters and many function-specific parameters allow options to be fine-  
145 tuned. Using the parameter option `--all` will turn on all the available functions and run them



146 with the default parameters, unless otherwise specified. Additionally, the user can provide  
147 parameters via a configuration file instead of via the command line.

148

## 149 **Cleaning Alignments**

150 CIAAlign consists of several functions to clean an MSA by removing commonly encountered  
151 alignment issues. All of these functions are optional and can be fine-tuned using user  
152 parameters. All parameters have default values. The available functions are presented here in  
153 the order they are executed by the program. The order can have a direct impact on the  
154 results, the functions removing positions that lead to the greatest disruptions in the MSA  
155 should be run first as they potentially make removing more positions unnecessary and  
156 therefore keep processing to a minimum. For example, divergent sequences often contain  
157 many insertions compared to the consensus, so removing these sequences first reduces the  
158 number of insertions which need to be removed. Sequences can be made shorter during  
159 processing with CIAAlign and therefore too short sequences are removed last.

160 Fig 1 shows a graphical representation of an example toy alignment before (Fig 1A) and after  
161 (Fig 1B-1F) using each function individually. The remove gap only function is run by default  
162 after every cleaning step, unless otherwise specified by the user.

163

### 164 *Remove Divergent*

165 For each column in the alignment, this function finds the most common nucleotide or amino  
166 acid and generates a temporary consensus sequence. Each sequence is then compared  
167 individually to this consensus sequence. Sequences which match the consensus at a  
168 proportion of positions less than a user-defined threshold (default 0.75) are excluded from the

169 alignment (Fig 1B). It is recommended to run the `make_similarity_matrix` function to  
170 calculate pairwise similarity before removing divergent sequences, in order to adjust the  
171 parameter value for more or less divergent alignments.

172

### 173 *Remove Insertions*

174 In order to define a region as an insertion, an alignment gap must be present in the majority of  
175 sequences, flanked by a minimum number of non-gap positions on either side, which can be  
176 defined by the user (default 5). The minimum and maximum size of insertion to be removed  
177 can also be defined by the user (default 3 and 300 respectively) (Fig 1C).

178

### 179 *Crop Ends*

180 Crop ends redefines where each sequence starts and ends, based on the ratio of the  
181 numbers of gap and non-gap positions observed up to a given position in the sequence. It  
182 then replaces all non-gap positions before and after the redefined start and end, respectively,  
183 with gaps. This will be described for redefining the sequence start, however crop ends is also  
184 applied to the reverse of the sequence to redefine the sequence end.

185 The number of gap positions separating every two consecutive non-gap positions is  
186 compared to a threshold and if that difference is higher than the threshold, the start of the  
187 sequence will be reset to that position. This threshold is defined as a proportion of the total  
188 sequence length, excluding gaps, and can be defined by the user (default: 0.05) (Fig 1D, Fig  
189 2).

190 The user can set a parameter that defines the maximum proportion of the sequence for which  
191 to consider the change in gap positions (default: 0.1) and therefore the innermost position at

192 which the start or end of the sequence may be redefined. It is recommended to set this  
193 parameter no higher than 0.1, since even if there are a large number of gap positions beyond  
194 this point, this is unlikely to be the result of incomplete sequences (Fig 2).

195

#### 196 *Remove short sequences*

197 Remove short sequences removes sequences which have less than a specified number of  
198 non-gap positions, which can be set by the user (default: 50) (Fig 1E).

199

#### 200 *Remove gap only columns*

201 Remove gap only removes columns that contain only gaps. These could be introduced by  
202 manual editing of the MSA before using CAlign or by running the functions above (Fig 1F).  
203 The main purpose of the function is to clean the gap only columns that are likely to be  
204 introduced after running any of the cleaning functions.

205

### 206 **Visualisation**

207 There are several ways of visualising the alignment, which both allow the user to interpret the  
208 alignment and clearly show which positions and sequences CAlign has removed. CAlign can  
209 also be used simply to visualise an alignment, without running any of the cleaning functions.  
210 All visualisations can be output as publication ready image files.

211

#### 212 *Mini Alignments*

213 CAlign provides functionality to generate mini alignments, in which an MSA is visualised  
214 using coloured rectangles on a single x and y axis, with each rectangle representing a single

215 nucleotide or amino acid (e.g. Fig 1, Figs 3-5). Even for large alignments, this function  
216 provides a visualisation that can be easily viewed and interpreted. Many properties of the  
217 resulting file (dimensions, DPI, file type) are parameterised. In order to minimise the memory  
218 and time required to generate the mini alignments, the matplotlib imshow function [25] for  
219 displaying images is used. Briefly, each position in each sequence in the alignment forms a  
220 single pixel in an image object and a custom dictionary is used to assign colours. The image  
221 object is then stretched to fit the axes.

222

### 223 *Sequence Logos*

224 CAlign can generate traditional sequence logos [26] or sequence logos using rectangles  
225 instead of letters to show the information and base / amino acid content at each position,  
226 which can increase readability in less conserved regions.

227

### 228 **Interpretation**

229 Some additional functions are provided to further interpret the alignment, for example plotting  
230 the number of sequences with non-gap residues at each position (the coverage), calculating a  
231 pairwise similarity matrix, and generating a consensus sequence with various options.

232 Given the toy example shown in Fig 1A, running all possible cleaning functions will lead to the  
233 markup plot shown in Fig 3A and the result shown in Fig 3B. In the markup plot each removed  
234 part is highlighted in a different colour corresponding to the function with which it was  
235 removed.

236

### 237 **Example Alignments**

238 Four example alignments are provided within the software directory to demonstrate the  
239 functionality of CIAAlign.

240 Example 1 is a very short alignment of six sequences which was generated manually by  
241 creating arbitrary sequences of nucleotides that would show every cleaning function while  
242 being as short as possible. This alignment contains an insertion, gaps at the ends of  
243 sequences, a very short sequence and some highly divergent sequences.

244 Example 2 is a larger alignment based on randomly generated amino acid sequences using  
245 RandSeq (a tool from ExPASy [27]) with an average amino acid composition, which were  
246 aligned with MAFFT v7.407, under the default settings [4]. The sequences were adjusted  
247 manually to reflect an alignment that would fully demonstrate the functionalities of CIAAlign. It  
248 consists of many sequences that align well, however there are again a few problems: one  
249 sequence has a large insertion, one is very short, one is extremely divergent, and some have  
250 multiple gaps at the start and at the end.

251 For Example 3, putative mitochondrial gene cytochrome C oxidase I (COI) sequences were  
252 identified by applying TBLASTN v2.9.0 [28] to the human COI sequence (GenBank accession  
253 NC\_012920.1, positions 5,904–7,445, translated to amino acids), querying against 1,565  
254 transcriptomic datasets from the NCBI transcriptome shotgun assembly (TSA) database [29]  
255 under the default settings. 2,855 putative COI transcripts were reverse complemented where  
256 required, and those corresponding to the COI gene of the primary host of the TSA dataset  
257 were identified using the BOLD online specimen identification engine [30] (accessed  
258 07/10/2019) querying against the species level barcode records. The resulting 232 sequences  
259 were then aligned with MAFFT v7.407, under the default settings [4].

260 For Example 4, 91 sequences were selected from Example 3 to be representative of as many  
261 taxonomic families as possible and to exclude families with unclear phylogeny in the  
262 literature. These sequences were aligned with MAFFT v7.407 under the default settings and  
263 the alignment was refined with 1000 iterations. Robinson-Foulds distances of the resulting  
264 trees were calculated using ete3 compare [31].

265

## 266 **Results and Discussion**

267 Here an example is presented and the visualisation functions are used to illustrate the  
268 functionality of CIAAlign. Results will differ when using different parameters and thresholds.

269 CIAAlign was applied to the Example 2 alignment with the following options:

```
270 python3 CIAAlign.py --infile INFILE --outfile_stem OUTFILE_STEM --all
```

271 Using these settings on the alignment in Fig 4A results in the markup shown in Fig 4B and the  
272 output shown in Fig 4C. The markup shows which function has removed each sequence or  
273 position. The benefits of CIAAlign are clear in this simulation – the single poorly aligned  
274 sequence, the large insertion, very short sequences, and gap-only columns have been  
275 removed, and the unreliably aligned end segments of the sequences have been cropped. The  
276 resulting alignment is significantly shorter, which will speed up and simplify any further  
277 analysis. The clear graphical representation makes it easy to see what has been removed, so  
278 in the case of over-trimming the user can intervene and adjust functions and parameters.

279 In order to test CIAAlign on real biological sequences, an alignment was generated based on  
280 the COI gene commonly used in phylogenetic analysis and DNA barcoding [30]. As CIAAlign  
281 addresses some common problems encountered when generating an MSA based on *de novo*

282 assembled transcripts, which tend to have a higher error rates at transcript ends, gaps due to  
283 difficult to assemble regions and divergent sequences due to chimeric connections between  
284 unrelated regions [11, 32], COI-like transcripts were identified by searching the NCBI  
285 transcriptome shotgun assembly database. Aligning these transcripts demonstrated several  
286 common problems – multiple insertions, poor alignment at the starts and ends of sequences,  
287 and a few divergent sequences resulting in excessive gaps (Fig 5A). This alignment was  
288 parsed using the default CIAAlign settings except the threshold for removing divergent  
289 sequences was reset to 50%, as some of the sequences were from evolutionarily distant  
290 species. Under these settings, CIAAlign resolved several of the problems with the alignment:  
291 the insertions and highly divergent sequences were removed and the poorly aligned regions  
292 at the starts and ends of sequences were cropped (Fig 5B). One sequence and 6,029  
293 positions were removed from the alignment and a total of 2,446 positions were cropped from  
294 the ends of 112 sequences. The processed alignment is 26.59% of the size of the input  
295 alignment. However, a minimal amount of actual sequence data (as opposed to gaps) was  
296 removed, with 85.70% of bases remaining.

297 A subset of this sequence set was selected to demonstrate the functionality of CIAAlign in  
298 streamlining phylogenetic analysis. 91 COI-like transcripts from different taxonomic families of  
299 metazoa were selected from Example 3, incorporated into an MSA and cleaned using CIAAlign  
300 with the same settings as above (S1 Fig). 1,437 positions were removed from the alignment  
301 and a total of 289 positions were cropped from the ends of 17 sequences. The processed  
302 alignment is 70.67% of the size of the input alignment and 96.52% of bases remain.  
303 Phylogenetic trees were generated for the input alignment and for the alignment processed  
304 with CIAAlign, using PhyML [33] under the GTR model plus the default settings. For the input

alignment, PhyML used 138 MB of memory and took 532 seconds (on one Intel Core i7-7560U core with 4 GB of RAM, running at 2.40 GHz). For the cleaned alignment, on the same machine, PhyML used 109 MB of memory and took 243 seconds. The tree generated with the input alignment (S1D Fig) had a Robinson-Foulds difference from a “correct” tree (generated manually based on the literature, S1C Fig) of 100.00. The tree generated with the cleaned alignment (S1E Fig) had a Robinson-Foulds difference from the correct tree of 90.00. Therefore the tree based on the CIAAlign cleaned alignment was generated more quickly, used less memory, and was more similar to the expected tree.

While the functionality of CIAAlign has some overlaps with other software, for example Jalview [34], Gblocks [7] and trimAl [8], the presented software can be seen as complementary to these, with some different features and applications. Jalview is designed for manual curation of alignments, but it is unsuitable for a simple overview of large alignments and does not provide the option of editing automatically, which is useful in large batch applications and ensures reproducibility. Gblocks is designed to choose blocks from an alignment that would be suitable for phylogenetic analysis, which is too restrictive for many other purposes. Some functionalities of trimAl overlap with those of CIAAlign; however, trimAl is designed to algorithmically define and remove any poorly aligned regions whereas CIAAlign is designed to remove specific MSA issues, as defined by the user, for different downstream applications. For highly divergent alignments, trimAl can be too sensitive and remove useful regions. CIAAlign also provides additional visualisation options. Therefore, CIAAlign should be seen as a tool that aims to fill in the gaps that exist in currently available software.

Having as many parameters as possible to allow as much user control as possible gives greater flexibility. However, this also means that these parameters should be adjusted, which



328 requires a good understanding of the cleaning functions and the MSA in question. CIAAlign  
329 offers default parameters selected to be often applicable based on many example MSAs used  
330 during the development of the tool. However, parameter choice highly depends on MSA  
331 divergence and the downstream application. To choose appropriate values it is recommended  
332 to first run CIAAlign with all default parameters and then adjust them based on the results.  
333 Since the mini alignments show what has been removed by which functions it is  
334 straightforward to identify the effect of each function and any changes to the parameters  
335 which may be required.

336 New features are in progress to be added in the future, such as collapsing very similar  
337 sequences, removing divergent columns, and making the colour scheme for the bases or  
338 amino acids customisable.

339

## 340 **Conclusion**

341 CIAAlign is a highly customisable tool which can be used to clean multiple sequence  
342 alignments and address several common alignment problems. Due to its multiple user options  
343 it can be used for many applications. CIAAlign provides clear visual output showing which  
344 positions have been removed and for what reason, allowing the user to adjust the parameters  
345 accordingly. A number of additional visualisation and interpretation options are provided.

346

## 347 **References**

348 1. Boswell, RD. Sequence alignment by word processor. Trends Biochem Sci.  
349 1987;12:279–80.

- 350 2. Higgins DG, Sharp PM. CLUSTAL: a package for performing multiple sequence  
351 alignment on a microcomputer. *Gene*. 1988;73:237–44.
- 352 3. Edgar RC. MUSCLE: a multiple sequence alignment method with reduced time and  
353 space complexity. *BMC Bioinformatics*. 2004;5:113.
- 354 4. Katoh K, Misawa K, Kuma K, Miyata T. MAFFT: a novel method for rapid multiple  
355 sequence alignment based on fast Fourier transform. *Nucleic Acids Res*. 2002;30:3059–66.
- 356 5. Notredame C, Higgins DG, Heringa J. T-coffee: a novel method for fast and accurate  
357 multiple sequence alignment<sup>11</sup> Edited by J. Thornton. *J Mol Biol*. 2000;302:205–17.
- 358 6. Needleman SB, Wunsch CD. A general method applicable to the search for similarities in  
359 the amino acid sequence of two proteins. *J Mol Biol*. 1970;48:443–53.
- 360 7. Talavera G, Castresana J. Improvement of phylogenies after removing divergent and  
361 ambiguously aligned blocks from protein sequence alignments. *Syst Biol*. 2007;56:564–77.
- 362 8. Capella-Gutiérrez S, Silla-Martínez JM, Gabaldón T. trimAl: a tool for automated  
363 alignment trimming in large-scale phylogenetic analyses. *Bioinformatics*. 2009;25:1972–3.
- 364 9. Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of  
365 large phylogenies. *Bioinforma Oxf Engl*. 2014;30:1312–3.
- 366 10. Kumar S, Stecher G, Li M, Knyaz C, Tamura K. MEGA X: Molecular Evolutionary  
367 Genetics Analysis across Computing Platforms. *Mol Biol Evol*. 2018;35:1547–9.
- 368 11. Bushmanova E, Antipov D, Lapidus A, Prjibelski AD. rnaSPAdes: a de novo  
369 transcriptome assembler and its application to RNA-Seq data. *GigaScience*. 2019;8.  
370 doi:10.1093/gigascience/giz100.

- 371 12. Richterich P. Estimation of Errors in “Raw” DNA Sequences: A Validation Study.  
372 *Genome Res.* 1998;8:251–9.
- 373 13. Tyler AD, Mataseje L, Urfano CJ, Schmidt L, Antonation KS, Mulvey MR, et al.  
374 Evaluation of Oxford Nanopore's MinION Sequencing Device for Microbial Whole Genome  
375 Sequencing Applications. *Sci Rep.* 2018;8:1–12.
- 376 14. Magi A, Giusti B, Tattini L. Characterization of MinION nanopore data for resequencing  
377 analyses. *Brief Bioinform.* 2017;18:940–53.
- 378 15. Fitch WM, Smith TF. Optimal sequence alignments. *PNAS.* 1983;80:1382–6.
- 379 16. Tettelin H, Maignani V, Cieslewicz MJ, Donati C, Medini D, Ward NL, et al. Genome  
380 analysis of multiple pathogenic isolates of *Streptococcus agalactiae*: implications for the  
381 microbial “pan-genome.” *Proc Natl Acad Sci U S A.* 2005;102:13950–5.
- 382 17. Hu B, Xie G, Lo C-C, Starkenburg SR, Chain PSG. Pathogen comparative genomics in  
383 the next-generation sequencing era: genome alignments, pangenomics and metagenomics.  
384 *Brief Funct Genomics.* 2011;10:322–33.
- 385 18. Sayyari E, Whitfield JB, Mirarab S. Fragmentary Gene Sequences Negatively Impact  
386 Gene Tree and Species Tree Reconstruction. *Mol Biol Evol.* 2017;34:3279–91.
- 387 19. Schulz F, Roux S, Paez-Espino D, Jungbluth S, Walsh DA, Denef VJ, et al. Giant virus  
388 diversity and host interactions through global metagenomics. *Nature.* 2020;578:432–6.
- 389 20. Käfer S, Paraskevopoulou S, Zirkel F, Wieseke N, Donath A, Petersen M, et al. Re-  
390 assessing the diversity of negative strand RNA viruses in insects. *PLoS Pathog.* 2019;15.  
391 doi:10.1371/journal.ppat.1008224.

392 21. Bäckström D, Yutin N, Jørgensen SL, Dharamshi J, Homa F, Zaremba-Niedwiedzka K,  
393 et al. Virus Genomes from Deep Sea Sediments Expand the Ocean Megavirome and Support  
394 Independent Origins of Viral Gigantism. *mBio*. 2019;10. doi:10.1128/mBio.02497-18.

395 22. Petyuk VA, Gatto L, Payne SH. Reproducibility and Transparency by Design. *Mol Cell*  
396 *Proteomics*. 2019. doi:10.1074/mcp.IP119.001567.

397 23. Brito JJ, Li J, Moore JH, Greene CS, Nogoy NA, Garmire LX, et al. Recommendations  
398 to enhance rigor and reproducibility in biomedical research. 2020.  
399 <https://arxiv.org/abs/2001.05127v2>.

400 24. Langille MGI, Ravel J, Fricke WF. “Available upon request”: not good enough for  
401 microbiome data! *Microbiome*. 2018;6:8.

402 25. Hunter JD. Matplotlib: A 2D Graphics Environment. *Comput Sci Eng*. 2007;9:90–5.

403 26. Schneider TD, Stephens RM. Sequence logos: a new way to display consensus  
404 sequences. *Nucleic Acids Res*. 1990;18:6097–100.

405 27. Gasteiger E, Gattiker A, Hoogland C, Ivanyi I, Appel RD, Bairoch A. ExPASy: The  
406 proteomics server for in-depth protein knowledge and analysis. *Nucleic Acids Res*.  
407 2003;31:3784–8.

408 28. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+:  
409 architecture and applications. *BMC Bioinformatics*. 2009;10:421.

410 29. Transcriptome Shotgun Assembly Sequence Database. National Center for  
411 Biotechnology Information, Bethesda, Maryland, USA. 2012.  
412 <https://www.ncbi.nlm.nih.gov/genbank/tsa/>. Accessed 08 Oct 2019.

- 413 30. bold: The Barcode of Life Data System (<http://www.barcodinglife.org>).  
414 <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC1890991/>. Accessed 6 Apr 2020.
- 415 31. Huerta-Cepas J, Serra F, Bork P. ETE 3: Reconstruction, Analysis, and Visualization of  
416 Phylogenomic Data. *Mol Biol Evol.* 2016;33:1635–8.
- 417 32. Liao X, Li M, Zou Y, Wu F-X, Yi-Pan, Wang J. Current challenges and solutions of de  
418 novo assembly. *Quant Biol.* 2019;7:90–109.
- 419 33. Guindon S, Gascuel O. A simple, fast, and accurate algorithm to estimate large  
420 phylogenies by maximum likelihood. *Syst Biol.* 2003;52:696–704.
- 421 34. Waterhouse AM, Procter JB, Martin DMA, Clamp M, Barton GJ. Jalview Version 2--a  
422 multiple sequence alignment editor and analysis workbench. *Bioinforma Oxf Engl.*  
423 2009;25:1189–91.

424

## 425 **Figure Legends**

### 426 **Fig 1:**

427 Mini alignments showing the main functionalities of CIAAlign based on Example 1. **a** Input  
428 alignment before application of CIAAlign, generated using the command “CIAAlign --infile  
429 example1.fasta --plot\_input”. **b** Output alignment showing the functionality of the  
430 remove divergent function, generated using the command “CIAAlign --infile  
431 example1.fasta --remove\_divergent --plot\_output”. **c** Output alignment  
432 showing the functionality of the remove insertions function, generated using the command  
433 “CIAAlign --infile example1.fasta --remove\_insertions --plot\_output”. **d**  
434 Output alignment showing the functionality of the crop ends function, generated using the

435 command "CIAAlign --infile example1.fasta --crop\_ends --plot\_output". **e**  
436 Output alignment showing the functionality of the remove short sequences function,  
437 generated using the command "CIAAlign --infile example1.fasta --  
438 remove\_short --plot\_output". **f** Output alignment showing the functionality of the  
439 remove gap only function, generated using the command "CIAAlign --infile  
440 example1.fasta --plot\_output". Subplots were generated using the  
441 drawMiniAlignment function of CIAAlign.

442

#### 443 **Fig 2:**

444 This manually created example illustrates how crop\_ends works internally. The length of the  
445 sequence shown is 111 including gaps and 80 excluding gaps (1). With a threshold of 10% for  
446 the proportion of non-gap positions to consider for change in end positions, 8 positions at the  
447 start and at the end, respectively, are being considered (illustrated by red crossbars). For  
448 each of these, the number of preceding gaps is calculated (2). Then the change in gap  
449 numbers (3) for every two consecutive non-gap positions is compared to the gap number  
450 change threshold, which is 5%, i.e. 4 gaps, as a default value. Looking at the change in gap  
451 numbers, the last change at each end equal to or bigger than the threshold is coloured in red.  
452 This leads to redefining the start and the end of this example sequence to be where the  
453 nucleotides are coloured in green.

454

#### 455 **Fig 3:**

456 Mini alignments and legends showing further functionalities of CIAAlign based on Example 1. **a**  
457 Alignment showing the functionality of the plot markup function, generated using the

458 command “CIAAlign --infile example1.fasta --all”. The areas that have been  
459 removed are marked up in different colours, each corresponding to a certain function of  
460 CIAAlign. **b** Output alignment after application of all functions of CIAAlign combined, generated  
461 using the command “CIAAlign --infile example1.fasta --all”. Subplots were  
462 generated using the drawMiniAlignment function.

463

464 **Fig 4:**

465 Mini alignments showing the main functionalities of CIAAlign based on Example 2. **a** Input  
466 alignment before application of CIAAlign, generated using the command “CIAAlign --infile  
467 example2.fasta --plot\_input”. **b** Alignment markup showing areas that were removed  
468 by CIAAlign, generated using the command “CIAAlign --infile example2.fasta --all”. **c** Output  
469 alignment after application of CIAAlign, generated using the command “CIAAlign --infile  
470 example2.fasta --all”. Subplots were generated using the drawMiniAlignment  
471 function.

472

473 **Fig 5:**

474 Mini alignments showing the main functionalities of CIAAlign based on Example 3. **a** Input  
475 alignment before application of CIAAlign, generated using the command “CIAAlign --infile  
476 example3.fasta --plot\_input”. **b** Output alignment after application of CIAAlign,  
477 generated using the command “CIAAlign --infile example3.fasta --all --  
478 remove\_divergent\_minperc 0.5”. Subplots were generated using the  
479 drawMiniAlignment function.

480

481 **S1 figure**

482 Mini alignments and phylogenetic trees showing the application of CIAAlign to phylogenetic  
483 data, based on Example 4, a subset of Example 3. **a** Input alignment before application of  
484 CIAAlign, generated using the command `"CIAAlign --infile example4.fasta --`  
485 `plot_input"`. **b** Output alignment after application of CIAAlign, generated using the  
486 command `"CIAAlign --infile example4.fasta --all --`  
487 `remove_divergent_minperc 0.5"`. Subplots were generated using the  
488 `"drawMiniAlignment` function. **c** Phylogenetic tree generated manually using the literature to  
489 show the current best estimate for the phylogenetic relationships between these 91 families of  
490 metazoa. Relationships are based on the literature listed in the S1 References. **d** PhyML  
491 phylogenetic tree generated under the GTR model plus default settings on the input alignment  
492 before application of CIAAlign. **e** PhyML phylogenetic tree generated under the GTR model  
493 plus default settings on the cleaned alignment after application of CIAAlign. In (c-e) branch  
494 colours correspond to the labelled phyla, coloured squares indicate class and bold text  
495 indicates order. Common names are shown where available.



Figure 1

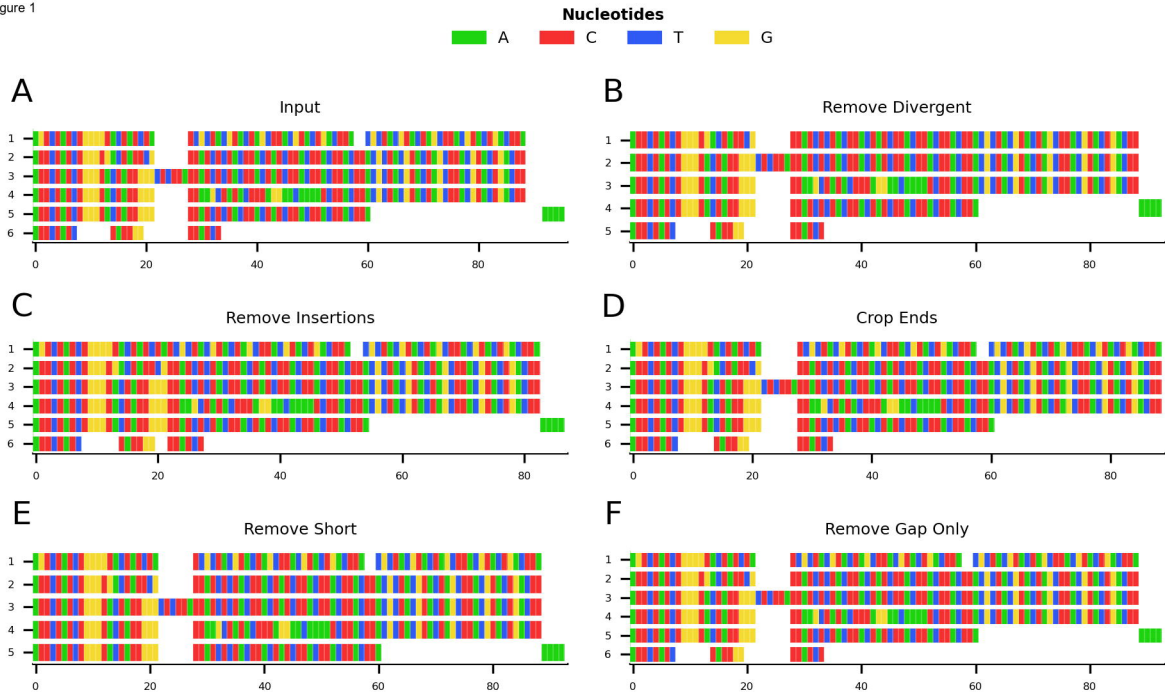


Figure 2

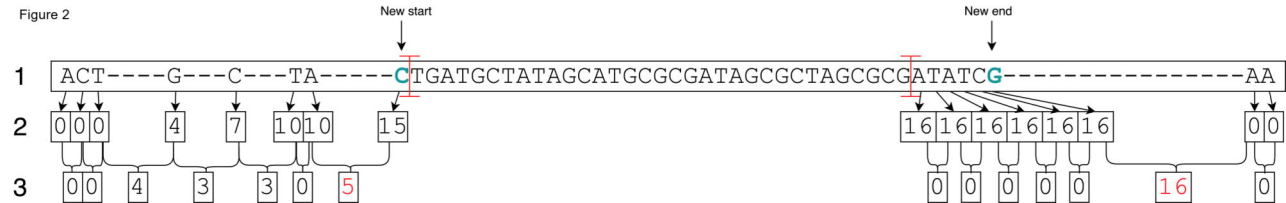


Figure 3

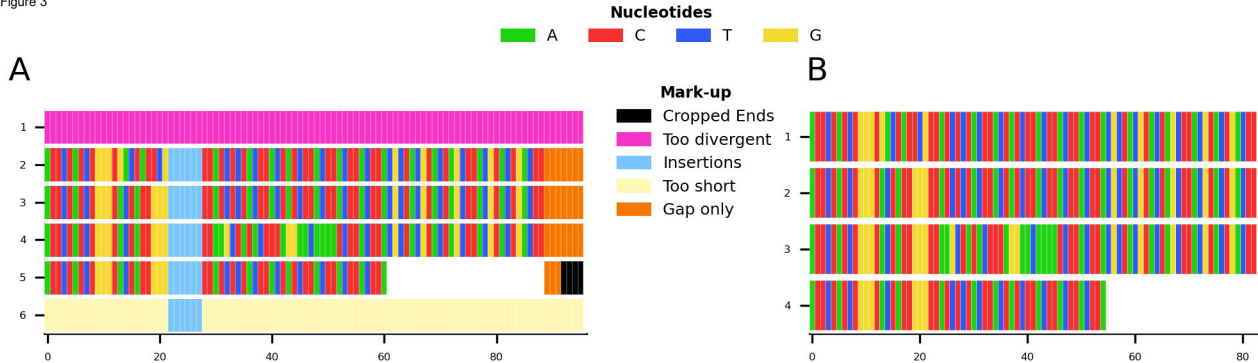
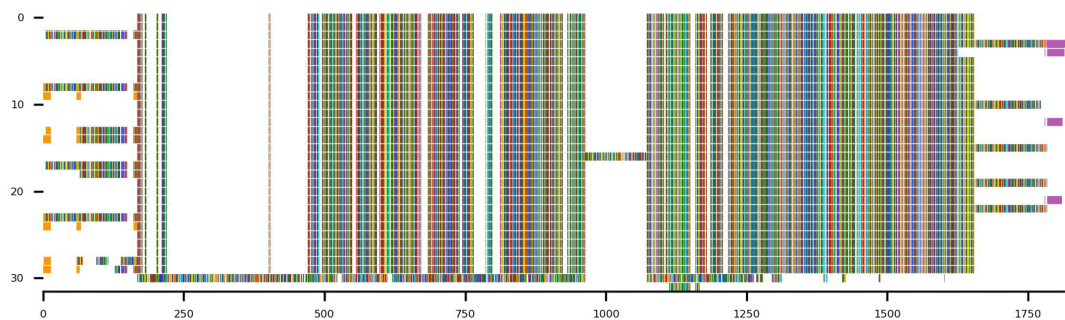


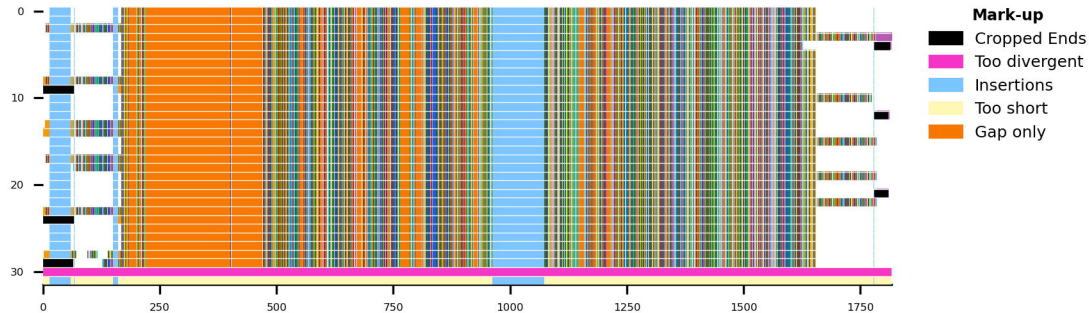
Figure 4



A



B



C

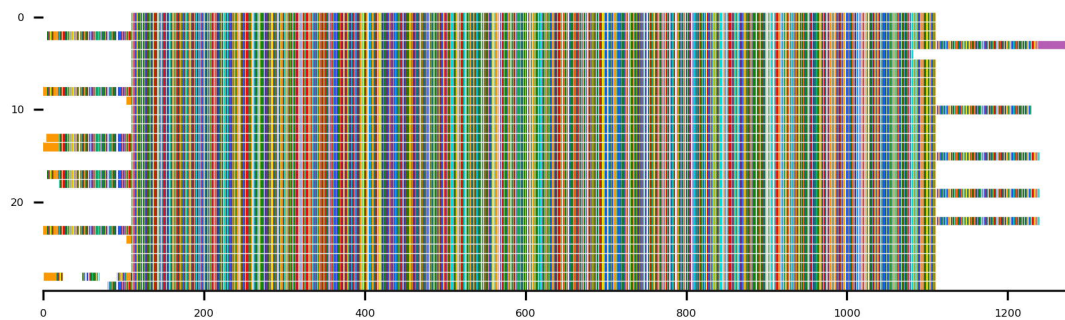


Figure 5

