1 **CIAlign - 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

## Abstract

### Background

Throughout biology, multiple sequence alignments (MSAs) form the basis of much investigation into biological features and relationships. These alignments are at the heart of many bioinformatics analyses. However, sequences in MSAs are often incomplete or very divergent, which leads to poorly aligned regions or large gaps in alignments. This slows down computation and can impact conclusions without being biologically relevant. Therefore, cleaning the alignment by removing these regions can substantially improve analyses. Manual editing of MSAs is very widespread but is time-consuming and difficult to reproduce.

### Results

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

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

2

32  **Conclusion**

33  CIAlign effectively removes problematic regions and sequences from MSAs and provides

34  novel visualisation options. This tool can be used to refine alignments for further analysis and

35  processing. The tool is aimed at anyone who wishes to automatically clean up parts of an

36  MSA and those requiring a new, accessible way of visualising large MSAs.

## Introduction

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

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

4

60 software can also take into account certain discrepancies in the alignment, for example

61 RaXML [9] can account for missing data in some columns and check for duplicate sequence

62 names and gap-only columns; similarly GUI based toolkits for molecular biology such as

63 MEGA [10] sometimes have options to delete or ignore columns containing gaps. However,

64 several common issues affect the speed, complexity and reliability of specific downstream

65 analyses but are not addressed by existing tools.

66 Clean and Interpret Alignments (CIAlign) is primarily intended to address four issues which

67 are commonly encountered when working with MSAs. Researchers in many fields regularly

68 edit MSAs by hand to address these issues, however as well as being extremely time

69 consuming, ensuring reproducibility with this approach is almost impossible and it cannot be

70 incorporated into an automated analysis pipeline.

71 The first issue we intend to address is that it is common for an MSA to contain more gaps

72 towards either end than in the body of the alignment. This problem occurs at both the

73 sequencing and alignment stage. For example, the ends of *de novo* assembled transcripts

74 tend to have lower read coverage [11] and therefore have a higher probability of mis-

75 assembly and therefore mis-alignment. MSAs created using these sequences therefore also

76 have regions of lower reliability towards either end. Similarly, both Sanger sequences and

77 sequences generated with Oxford Nanopore's long read sequencing technology, which are

78 often used directly in MSAs, tend to have lower quality scores at the either the beginning or

79 the end [12, 13, 14]. Automated removal of these regions from MSAs would therefore

80 increase the reliability of downstream analyses. Also, while generating an MSA, terminal gaps

81 complicate analysis, and the weighting of terminal gaps relative to internal gap opening and

5

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

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

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

104  Finally, very short partially overlapping sequences cannot always be reliably aligned using

105  standard global alignment algorithms. It is very common to remove these sequences,

106  manually or otherwise, prior to further analysis.

107  There are also several common issues in alignment visualisation. Large alignments can be

108  difficult to visualise and a small and concise but accurate visualisation can be useful when

109  presenting results, so this has been incorporated into the software. With many alignment

110  trimming tools it can be difficult to track exactly which changes the software has made, so a

111  visual output showing these changes could be helpful.

112  Finally, transparency is often an issue with bioinformatics software, with poor reporting of

113  exactly how a file has been processed [22–24]. CIAlign has been developed to process

114  alignments in a transparent manner, to allow the user to clearly and reproducibly report their

115  methodology.

116  CIAlign is freely available at github.com/KatyBrown/CIAlign.

117

## Materials and Methods

119  CIAlign is a command line tool implemented in Python 3. It can be installed either via pip3 or

120  from GitHub and is independent of the operating system. It has been designed to enable the

121  user to remove specific issues from an MSA, to visualise the MSA (including a markup file

122  showing which regions and sequences have been removed), and to interpret the MSA in

123  several ways. CIAlign works on nucleotide or amino acids alignments and will detect which of

124  these is provided. A log file is generated to show exactly which sequences and positions have

7

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

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

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

8

147 with the default parameters, unless otherwise specified. Additionally, the user can provide

148 parameters via a configuration file instead of via the command line.

149 CIAlign has been designed to maximise usability, reproducibility and reliability. The code is

150 written to be as readable as possible and all functions are fully documented. All functions are

151 covered by unit tests. CIAlign is freely available, open source and fully version controlled.

152

153 **Cleaning Alignments**

154 CIAlign consists of several functions to clean an MSA by removing commonly encountered

155 alignment issues. All of these functions are optional and can be fine-tuned using user

156 parameters. All parameters have default values. The available functions are presented here in

157 the order they are executed by the program. The order can have a direct impact on the

158 results, the functions removing positions that lead to the greatest disruptions in the MSA

159 should be run first as they potentially make removing more positions unnecessary and

160 therefore keep processing to a minimum. For example, divergent sequences often contain

161 many insertions compared to the consensus, so removing these sequences first reduces the

162 number of insertions which need to be removed. Sequences can be made shorter during

163 processing with CIAlign and therefore too short sequences are removed last.

164 Fig 1 shows a graphical representation of an example toy alignment before (Fig 1A) and after

165 (Fig 1B-1F) using each function individually. The remove gap only function is run by default

166 after every cleaning step, unless otherwise specified by the user.

167

168 *Remove Divergent*

9

169    For each column in the alignment, this function finds the most common nucleotide or amino

170    acid and generates a temporary consensus sequence. Each sequence is then compared

171    individually to this consensus sequence. Sequences which match the consensus at a

172    proportion of positions less than a user-defined threshold (default 0.65) are excluded from the

173    alignment (Fig 1B). It is recommended to run the `make_similarity_matrix` function to

174    calculate pairwise similarity before removing divergent sequences, in order to adjust the

175    parameter value for more or less divergent alignments.

176

177    *Remove Insertions*

178    In order to define a region as an insertion, an alignment gap must be present in the majority of

179    sequences and flanked by a minimum number of non-gap positions on either side, which can

180    be defined by the user (default 5). The minimum and maximum size of insertion to be

181    removed can also be defined by the user (default 3 and 200 respectively) (Fig 1C).

182

183    *Crop Ends*

184    Crop ends redefines where each sequence starts and ends, based on the ratio of the

185    numbers of gap and non-gap positions observed up to a given position in the sequence. It

186    then replaces all non-gap positions before and after the redefined start and end, respectively,

187    with gaps. This will be described for redefining the sequence start, however crop ends is also

188    applied to the reverse of the sequence to redefine the sequence end.

189    The number of gap positions separating every two consecutive non-gap positions is

190    compared to a threshold and if that difference is higher than the threshold, the start of the

191    sequence will be reset to that position. This threshold is defined as a proportion of the total

10

192    sequence length, excluding gaps, and can be defined by the user (default: 0.05) (Fig 1D, Fig

193    2).

194    The user can set a parameter that defines the maximum proportion of the sequence for which

195    to consider the change in gap positions (default: 0.1) and therefore the innermost position at

196    which the start or end of the sequence may be redefined. It is recommended to set this

197    parameter no higher than 0.1, since even if there are a large number of gap positions beyond

198    this point, this is unlikely to be the result of incomplete sequences (Fig 2).

199

200    *Remove short sequences*

201    Remove short sequences removes sequences which have less than a specified number of

202    non-gap positions,which can be set by the user (default: 50) (Fig 1E).

203

204    *Remove gap only columns*

205    Remove gap only removes columns that contain only gaps. These could be introduced by

206    manual editing of the MSA before using CIAlign or by running the functions above (Fig 1F).

207    The main purpose of the function is to clean the gap only columns that are likely to be

208    introduced after running any of the cleaning functions.

209

210    **Visualisation**

211    There are several ways of visualising the alignment, which both allow the user to interpret the

212    alignment and clearly show which positions and sequences CIAlign has removed. CIAlign can

213    also be used simply to visualise an alignment, without running any of the cleaning functions.

214    All visualisations can be output as publication ready image files.

11

215

*Mini Alignments*

217 CIAlign provides functionality to generate mini alignments, in which an MSA is visualised

218 using coloured rectangles on a single *x* and *y* axis, with each rectangle representing a single

219 nucleotide or amino acid (e.g. Fig 1, Figs 3-5). Even for large alignments, this function

220 provides a visualisation that can be easily viewed and interpreted. Many properties of the

221 resulting file (dimensions, DPI, file type) are parameterised. In order to minimise the memory

222 and time required to generate the mini alignments, the matplotlib imshow function [25] for

223 displaying images is used. Briefly, each position in each sequence in the alignment forms a

224 single pixel in an image object and a custom dictionary is used to assign colours. The image

225 object is then stretched to fit the axes.

226

*Sequence Logos*

228 CIAlign can generate traditional sequence logos [26] or sequence logos using rectangles

229 instead of letters to show the information and base / amino acid content at each position,

230 which can increase readability in less conserved regions.

231

**Interpretation**

233 Some additional functions are provided to further interpret the alignment, for example plotting

234 the number of sequences with non-gap residues at each position (the coverage), calculating a

235 pairwise similarity matrix, and generating a consensus sequence with various options.

236 Given the toy example shown in Fig 1A, running all possible cleaning functions will lead to the

237 markup plot shown in Fig 3A and the result shown in Fig 3B. In the markup plot each removed

12

238  part is highlighted in a different colour corresponding to the function with which it was

239  removed.

240

241  **Example Alignments**

242  Four example alignments are provided within the software directory to demonstrate the

243  functionality of CIAlign. Examples 1 and 2 use simulated sequences, examples 3 and 4 use

244  real biological sequences and are designed to resemble the type of complex alignment many

245  researchers encounter.

246  Example 1 is a very short alignment of six sequences which was generated manually by

247  creating arbitrary sequences of nucleotides that would show every cleaning function while

248  being as short as possible. This alignment contains an insertion, gaps at the ends of

249  sequences, a very short sequence and some highly divergent sequences.

250  Example 2 is a larger alignment based on randomly generated amino acid sequences using

251  RandSeq (a tool from ExPASy [27]) with an average amino acid composition, which were

252  aligned with MAFFT v7.407, under the default settings [4]. The sequences were adjusted

253  manually to reflect an alignment that would fully demonstrate the functionalities of CIAlign. It

254  consists of many sequences that align well, however there are again a few problems: one

255  sequence has a large insertion, one is very short, one is extremely divergent, and some have

256  multiple gaps at the start and at the end. For Example 3, putative mitochondrial gene

257  cytochrome C oxidase I (COI) sequences were identified by applying TBLASTN v2.9.0 [28] to

258  the human COI sequence (GenBank accession NC_012920.1, positions 5,904–7,445,

259  translated to amino acids), querying against 1,565 transcriptomic datasets from the NCBI

260  transcriptome shotgun assembly (TSA) database [29] under the default settings. 2,855

13

261 putative COI transcripts were reverse complemented where required, and those

262 corresponding to the COI gene of the primary host of the TSA dataset were identified using

263 the BOLD online specimen identification engine [30] (accessed 07/10/2019) querying against

264 the species level barcode records. The resulting 232 sequences were then aligned with

265 MAFFT v7.407, under the default settings [4].

266 For Example 4, 91 sequences were selected from Example 3 to be representative of as many

267 taxonomic families as possible and to exclude families with unclear phylogeny in the

268 literature. These sequences were aligned with MAFFT v7.407 under the default settings and

269 the alignment was refined with 1000 iterations. Robinson-Foulds distances of the resulting

270 trees were calculated using ete3 compare [31].

271 Materials and methods for benchmarking and for large-scale examples with biological data

272 are provided as Supplementary Materials and Methods.

273

## Results and Discussion

274

275 Here an example is presented and the visualisation functions are used to illustrate the

276 functionality of CIAlign. Results will differ when using different parameters and thresholds.

277 CIAlign was applied to the Example 2 alignment with the following options:

278 
```
python3 CIAlign.py --infile INFILE --outfile_stem OUTFILE_STEM --all
```

279 Using these settings on the alignment in Fig 4A results in the markup shown in Fig 4B and the

280 output shown in Fig 4C. The markup shows which function has removed each sequence or

281 position. The benefits of CIAlign are clear in this simulation – the single poorly aligned

282 sequence, the large insertion, very short sequences, and gap-only columns have been

14

283  removed, and the unreliably aligned end segments of the sequences have been cropped. The

284  resulting alignment is significantly shorter, which will speed up and simplify any further

285  analysis. The clear graphical representation makes it easy to see what has been removed, so

286  in the case of over-trimming the user can intervene and adjust functions and parameters.

287  In order to demonstrate the use of CIAlign on real biological sequences, an alignment was

288  generated based on the COI gene commonly used in phylogenetic analysis and DNA

289  barcoding [30]. As CIAlign addresses some common problems encountered when generating

290  an MSA based on *de novo* assembled transcripts, which tend to have a higher error rates at

291  transcript ends, gaps due to difficult to assemble regions and divergent sequences due to

292  chimeric connections between unrelated regions [11, 32], COI-like transcripts were identified

293  by searching the NCBI transcriptome shotgun assembly database. Aligning these transcripts

294  demonstrated several common problems – multiple insertions, poor alignment at the starts

295  and ends of sequences, and a few divergent sequences resulting in excessive gaps (Fig 5A).

296  This alignment was cleaned using the default CIAlign settings except the threshold for

297  removing divergent sequences was reset to 50%, as some of the sequences are from

298  evolutionarily distant species. Under these settings, CIAlign resolved several of the problems

299  with the alignment: the insertions and highly divergent sequences were removed and the

300  poorly aligned regions at the starts and ends of sequences were cropped (Fig 5B). One

301  sequence and 6,029 positions were removed from the alignment and a total of 2,446 positions

302  were cropped from the ends of 112 sequences. The processed alignment is 26.59% of the

303  size of the input alignment. However, a minimal amount of actual sequence data (as opposed

304  to gaps) was removed, with 85.70% of bases remaining.

15

305  A subset of this sequence set was selected to demonstrate the functionality of CIAlign in

306  streamlining phylogenetic analysis. 91 COI-like transcripts from different taxonomic families of

307  metazoa were selected from Example 3, incorporated into an MSA and cleaned using CIAlign

308  with the same settings as above (Supp. Fig 1). 1,437 positions were removed from the

309  alignment and a total of 289 positions were cropped from the ends of 17 sequences. The

310  processed alignment is 70.67% of the size of the input alignment and 96.52% of bases

311  remain. Phylogenetic trees were generated for the input alignment and for the alignment

312  processed with CIAlign, using PhyML [33] under the GTR model plus the default settings. For

313  the input alignment, PhyML used 138 MB of memory and took 532 seconds (on one Intel

314  Core i7-7560U core with 4 GB of RAM, running at 2.40 GHz). For the cleaned alignment, on

315  the same machine, PhyML used 109 MB of memory and took 243 seconds. The tree

316  generated with the input alignment (Supp. Fig 1D) had a Robinson-Foulds [34] difference

317  from a "correct" tree (generated manually based on the literature, Supp. Fig 1D) of 100.00

318  (normalised Robinson-Foulds 0.57, Quartet divergence [35] 0.159). The tree generated with

319  the cleaned alignment (Supp. Fig 1E) had a Robinson-Foulds difference from the correct tree

320  of 90.00 (normalised Robinson-Foulds 0.52, Quartet divergence 0.073) Therefore the tree

321  based on the CIAlign cleaned alignment was generated more quickly, used less memory, and

322  was more similar to the expected tree.

323

**Benchmarking with Simulated Data**

325  We performed a series of benchmarking analyses on simulated data, in order to test and

326  demonstrate the utility of the CIAlign cleaning functions, confirm the validity of our default

16

327  parameter settings and ensure that running these functions does not have unexpected

328  negative effects on downstream analyses.

329  First, CIAlign was tested using two tools (EvolvAGene [36], and INDELible [37]) which

330  generate sets of unaligned sequences alongside "true" alignments and phylogenies expected

331  to accurately represent the relationship between the sequences. We used these tools to

332  determine if cleaning a user generated alignment with CIAlign affects its distance from the

333  true alignment. Test alignments of the simulated data were created using four common

334  alignment algorithms. These alignments were then cleaned with CIAlign with relaxed,

335  moderate or stringent parameter settings (Supplementary Table 1). With relaxed CIAlign

336  settings, a median of 0.19% of correct pairs of aligned residues (POARs) [38] were removed,

337  for moderate settings 0.75% were removed and for stringent settings 3.76% (Fig 6A, Table 1).

338  For comparison, the median total proportion of residues removed was 1.72% for relaxed,

339  2.40% for moderate and 3.86% for stringent (Fig 6A, Table 1). The median proportions of gap

340  positions removed were much higher: 53-56% for all sets of parameters (Fig 6A, Table 1).

341  This shows that with relaxed and moderate settings, running CIAlign has a very minimal

342  impact on correctly aligned residues in the alignment, while a considerable amount of gaps

343  and noise are removed. The more stringent settings should be used cautiously, however even

344  with high stringency a large majority of correctly aligned residues remain and the majority of

345  gaps are removed.

346  Phylogenetic trees were generated for each of these alignments to determine if cleaning with

347  CIAlign impacts the distance between the true phylogenetic tree and a phylogenetic tree

348  based on a test alignment (Fig 6B, Table 1. The mean normalised Robinson-Foulds distance

349  [34] and Quartet divergence [35] between the test trees and true trees were virtually

17

350  unchanged by running CIAlign and none of the changes were statistically significant (p>0.05,

351  Mann Whitney U test) (Fig 6B, Table 1).

352

353  We also compared the input sequence for our simulations to consensus sequences based on

354  alignments with and without CIAlign cleaning. For all three stringency levels, CIAlign

355  increased the percentage nucleotide identity between the consensus sequence and the input

356  sequence by 4% to 5% (Fig 6C, Table 1). All of these changes are statistically significant

357  (p<0.001, Mann-Whitney U test).

358  The long-read sequencing simulation tool BadRead [39] was used to demonstrate the use of

359  CIAlign to remove common sources of error in long read sequencing data. Sequences were

360  generated to represent low, moderate and high quality Oxford Nanopore reads based on an

361  input genome, then aligned and cleaned with CIAlign with moderate settings (Supplementary

362  Table 1). Using CIAlign increased the identity between the alignment consensus and the input

363  sequence significantly for all read quality levels - by 6.57% for high quality reads, 9.51% for

364  moderate quality reads and 12.25% for poor quality reads (Fig 6D) (all p<0.001, Mann-

365  Whitney U test). For the high quality reads, the reads cleaned with CIAlign generated

366  consensus sequences almost identical to the input sequence, with a mean of 99.24% identity

367  (Fig 6D).

368  Full output tables for all three sets of simulations are available online at

369  github.com/KatyBrown/CIAlign/benchmarking/tables and the simulated data and alignments

370  at github.com/KatyBrown/benchmarking_data_CIAlign.

371

372  **Examples of Using CIAlign with Biological Data**

18

373 We also used CIAlign to clean real biological data from several online databases, in order to

374 test and demonstrate its usefulness in automated processing of different types of sequencing

375 data.

376 *Cleaning Pfam Alignments*

377 The Pfam database provides manually curated seed alignments for over 17,000 protein

378 families, plus much larger automatically generated full alignments containing sequences

379 identified by database searching [40]. CIAlign cleaning functions were applied to seed and full

380 alignments for 500 Pfam domains and consensus sequences were generated for both

381 alignments, before and after cleaning. Randomly selected sequences from the full alignment

382 were then compared to each consensus. For the full alignments, the mean identity between

383 the consensus sequence and the alignment sequences increased by 10.71% ($p < 0.001$,

384 Mann-Whitney U test) after cleaning with CIAlign (Fig 7A). For the seed alignments identity

385 also increased significantly, by 4.89% ($p < 0.001$, Mann-Whitney U test) (Fig 7A). After

386 running CIAlign, the full alignment consensus approaches the level of similarity to the

387 alignment sequences which is seen for seed alignment consensus, despite the full alignment

388 having undergone no manual curation (Fig 7A). Even for the curated seed alignments,

389 cleaning with CIAlign further increases the similarity between the consensus and the aligned

390 sequences.

391

392 *Removing Insertions and Deletions from Human Genes*

393 To demonstrate the ability of CIAlign to remove non-majority indels, we used data for 50

394 indels across over 150 individuals from the 1000 genomes project [41], which has annotated

395 insertions and deletions for individual human genomes. In all cases, CIAlign removed all

19

396 insertions present in a majority of samples and ignored all insertions present in a minority of

397 samples (Fig 7B).

398

399 *Removing Outliers*

400 CIAlign can also be used to remove clear outliers from an alignment, for example prior to

401 phylogenetic analysis. To illustrate this, we ran the CIAlign cleaning functions on data from the

402 mammalian 10K trees project [42]. Three single-gene trees were identified with clear outliers,

403 the 12S ribosomal gene from primates and the *APOB* and *RAG1* genes from Carnivora. The

404 issues with these trees are shown in Fig 7C and Supp Fig 2. CIAlign successfully removed

405 the outlying group, without removing any other sequences, in all three of these cases.

406

407 **Comparison with Other Software**

408 While the functionality of CIAlign has some overlaps with other software, for example Jalview

409 [43], Gblocks [7] and trimAl [8], the presented software can be seen as complementary to

410 these, with some different features and applications. Jalview is designed for manual curation

411 of alignments, but it is unsuitable for a simple overview of large alignments and does not

412 provide the option of editing automatically, which is useful in large batch applications and

413 ensures reproducibility. Gblocks is designed to choose blocks from an alignment that would

414 be suitable for phylogenetic analysis, which is too restrictive for many other purposes. Some

415 functionalities of trimAl overlap with those of CIAlign; however, trimAl is designed to

416 algorithmically define and remove any poorly aligned regions whereas CIAlign is designed to

417 remove specific MSA issues, as defined by the user, for different downstream applications.

418 For highly divergent alignments, trimAl can be too sensitive and remove useful regions.

20

419 CIAlign also provides additional visualisation options. Therefore, CIAlign should be seen as a

420 tool that aims to fill in the gaps that exist in currently available software.

421

## **Parameters**

423 Having as many parameters as possible to allow as much user control as possible gives

424 greater flexibility. However, this also means that these parameters should be adjusted, which

425 requires a good understanding of the cleaning functions and the MSA in question. CIAlign

426 offers default parameters selected to be often applicable based on our benchmarking

427 simulations and testing with different types of data. However, parameter choice highly

428 depends on MSA divergence and the downstream application. To choose appropriate values

429 it is recommended to first run CIAlign with all default parameters and then adjust these

430 parameters based on the results. Since the mini alignments show what has been removed by

431 which functions it is straightforward to identify the effect of each function and any changes to

432 the parameters which may be required.

433 New features are in progress to be added in the future, such as collapsing very similar

434 sequences, removing divergent columns, and making the colour scheme for the bases or

435 amino acids customisable.

436

## **Conclusion**

438 CIAlign is a highly customisable tool which can be used to clean multiple sequence

439 alignments and address several common alignment problems. Due to its multiple user options

440 it can be used for many applications. CIAlign provides clear visual output showing which

21

441 positions have been removed and for what reason, allowing the user to adjust the parameters

442 accordingly. A number of additional visualisation and interpretation options are provided.

443

## Availability

445 **Current release, v1.0.10**: doi.org/10.5281/zenodo.4650727

446 (corresponds to github.com/KatyBrown/CIAlign/releases/tag/v1.0.10)

447 **GitHub**: github.com/KatyBrown/CIAlign

448 **pip3**: pypi.org/project/cialign

449

## References

451 1. Boswell, RD. Sequence alignment by word processor. Trends Biochem Sci.
452 1987;12:279–80.

453 2. Higgins DG, Sharp PM. CLUSTAL: a package for performing multiple sequence
454 alignment on a microcomputer. Gene. 1988;73:237–44.

455 3. Edgar RC. MUSCLE: a multiple sequence alignment method with reduced time and
456 space complexity. BMC Bioinformatics. 2004;5:113.

457 4. Katoh K, Misawa K, Kuma K, Miyata T. MAFFT: a novel method for rapid multiple
458 sequence alignment based on fast Fourier transform. Nucleic Acids Res. 2002;30:3059–66.

459 5. Notredame C, Higgins DG, Heringa J. T-coffee: a novel method for fast and accurate
460 multiple sequence alignment11Edited by J. Thornton. J Mol Biol. 2000;302:205–17.

22

461  6. Needleman SB, Wunsch CD. A general method applicable to the search for similarities in the amino acid sequence of two proteins. J Mol Biol. 1970;48:443–53.

463  7. Talavera G, Castresana J. Improvement of phylogenies after removing divergent and ambiguously aligned blocks from protein sequence alignments. Syst Biol. 2007;56:564–77.

465  8. Capella-Gutiérrez S, Silla-Martínez JM, Gabaldón T. trimAl: a tool for automated alignment trimming in large-scale phylogenetic analyses. Bioinformatics. 2009;25:1972–3.

467  9. Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. Bioinforma Oxf Engl. 2014;30:1312–3.

469  10. Kumar S, Stecher G, Li M, Knyaz C, Tamura K. MEGA X: Molecular Evolutionary Genetics Analysis across Computing Platforms. Mol Biol Evol. 2018;35:1547–9.

471  11. Bushmanova E, Antipov D, Lapidus A, Prjibelski AD. rnaSPAdes: a de novo transcriptome assembler and its application to RNA-Seq data. GigaScience. 2019;8. doi:10.1093/gigascience/giz100.

474  12. Richterich P. Estimation of Errors in "Raw" DNA Sequences: A Validation Study. Genome Res. 1998;8:251–9.

476  13. Tyler AD, Mataseje L, Urfano CJ, Schmidt L, Antonation KS, Mulvey MR, et al. Evaluation of Oxford Nanopore's MinION Sequencing Device for Microbial Whole Genome Sequencing Applications. Sci Rep. 2018;8:1–12.

479  14. Magi A, Giusti B, Tattini L. Characterization of MinION nanopore data for resequencing analyses. Brief Bioinform. 2017;18:940–53.

481  15. Fitch WM, Smith TF. Optimal sequence alignments. PNAS. 1983;80:1382–6.

23

482    16. Tettelin H, Masignani V, Cieslewicz MJ, Donati C, Medini D, Ward NL, et al. Genome

483    analysis of multiple pathogenic isolates of Streptococcus agalactiae: implications for the

484    microbial "pan-genome." Proc Natl Acad Sci U S A. 2005;102:13950–5.

485    17. Hu B, Xie G, Lo C-C, Starkenburg SR, Chain PSG. Pathogen comparative genomics in

486    the next-generation sequencing era: genome alignments, pangenomics and metagenomics.

487    Brief Funct Genomics. 2011;10:322–33.

488    18. Sayyari E, Whitfield JB, Mirarab S. Fragmentary Gene Sequences Negatively Impact

489    Gene Tree and Species Tree Reconstruction. Mol Biol Evol. 2017;34:3279–91.

490    19. Schulz F, Roux S, Paez-Espino D, Jungbluth S, Walsh DA, Denef VJ, et al. Giant virus

491    diversity and host interactions through global metagenomics. Nature. 2020;578:432–6.

492    20. Käfer S, Paraskevopoulou S, Zirkel F, Wieseke N, Donath A, Petersen M, et al. Re-

493    assessing the diversity of negative strand RNA viruses in insects. PLoS Pathog. 2019;15.

494    doi:10.1371/journal.ppat.1008224.

495    21. Bäckström D, Yutin N, Jørgensen SL, Dharamshi J, Homa F, Zaremba-Niedwiedzka K,

496    et al. Virus Genomes from Deep Sea Sediments Expand the Ocean Megavirome and Support

497    Independent Origins of Viral Gigantism. mBio. 2019;10. doi:10.1128/mBio.02497-18.

498    22. Petyuk VA, Gatto L, Payne SH. Reproducibility and Transparency by Design. Mol Cell

499    Proteomics. 2019. doi:10.1074/mcp.IP119.001567.

500    23. Brito JJ, Li J, Moore JH, Greene CS, Nogoy NA, Garmire LX, et al. Recommendations

501    to      enhance      rigor      and      reproducibility      in      biomedical      research.      2020.

502    https://arxiv.org/abs/2001.05127v2.

24

503    24. Langille MGI, Ravel J, Fricke WF. "Available upon request": not good enough for

504    microbiome data! Microbiome. 2018;6:8.

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

506    26. Schneider TD, Stephens RM. Sequence logos: a new way to display consensus

507    sequences. Nucleic Acids Res. 1990;18:6097–100.

508    27. Gasteiger E, Gattiker A, Hoogland C, Ivanyi I, Appel RD, Bairoch A. ExPASy: The

509    proteomics server for in-depth protein knowledge and analysis. Nucleic Acids Res.

510    2003;31:3784–8.

511    28. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+:

512    architecture and applications. BMC Bioinformatics. 2009;10:421.

513    29. Transcriptome Shotgun Assembly Sequence Database. National Center for

514    Biotechnology Information, Bethesda, Maryland, USA. 2012.

515    https://www.ncbi.nlm.nih.gov/genbank/tsa/. Accessed 08 Oct 2019.

516    30. bold: The Barcode of Life Data System (http://www.barcodinglife.org).

517    https://www.ncbi.nlm.nih.gov/pmc/articles/PMC1890991/. Accessed 6 Apr 2020.

518    31. Huerta-Cepas J, Serra F, Bork P. ETE 3: Reconstruction, Analysis, and Visualization of

519    Phylogenomic Data. Mol Biol Evol. 2016;33:1635–8.

520    32. Liao X, Li M, Zou Y, Wu F-X, Yi-Pan, Wang J. Current challenges and solutions of de

521    novo assembly. Quant Biol. 2019;7:90–109.

522    33. Guindon S, Gascuel O. A simple, fast, and accurate algorithm to estimate large

523    phylogenies by maximum likelihood. Syst Biol. 2003;52:696–704.

25

34. Robinson DF, Foulds LR. Comparison of phylogenetic trees. Math Biosci. 1981;53:131–47.

35. Smith MR. Bayesian and parsimony approaches reconstruct informative trees from simulated morphological datasets. Biol Lett. 2019;15:20180632.

36. Hall BG. Simulating DNA coding sequence evolution with EvolveAGene 3. Mol Biol Evol. 2008;25:688–95.

37. Fletcher W, Yang Z. INDELible: a flexible simulator of biological sequence evolution. Mol Biol Evol. 2009;26:1879–88.

38. Thompson JD, Plewniak F, Poch O. A comprehensive comparison of multiple sequence alignment programs. Nucleic Acids Res. 1999;27:2682–90.

39. Wick RR. Badread: simulation of error-prone long reads. J Open Source Softw. 2019;4:1316.

40. Finn RD, Bateman A, Clements J, Coggill P, Eberhardt RY, Eddy SR, et al. Pfam: the protein families database. Nucleic Acids Res. 2014;42 Database issue:D222–30.

41. Auton A, Abecasis GR, Altshuler DM, Durbin RM, Abecasis GR, Bentley DR, et al. A global reference for human genetic variation. Nature. 2015;526:68–74.

42. Arnold C, Matthews LJ, Nunn CL. The 10kTrees website: A new online resource for primate phylogeny. Evol Anthropol Issues News Rev. 2010;19:114–8.

43. Waterhouse AM, Procter JB, Martin DMA, Clamp M, Barton GJ. Jalview Version 2--a multiple sequence alignment editor and analysis workbench. Bioinforma Oxf Engl. 2009;25:1189–91.

26

545

## Tables

547 **Table 1**

| CIAlign Stringency | Correct Pairs Removed | Nucleotides Removed | Gaps Removed | Positions Removed | Normalised Robinson-Foulds Distance | Quartet Divergence | Consensus Percentage Identity |
|---|---|---|---|---|---|---|---|
| None | - | - | - | - | 0.2475 | 0.1622 | 67.15 |
| Relaxed | 0.1900 | 1.725 | 53.35 | 7.430 | 0.2403 | 0.1626 | 71.47 |
| Moderate | 0.7500 | 2.400 | 53.42 | 8.130 | 0.2561 | 0.1665 | 71.50 |
| Stringent | 3.760 | 3.855 | 55.72 | 10.00 | 0.2518 | 0.1705 | 71.48 |

548

549 Table showing the effect of cleaning with CIAlign on simulated alignments. Correct pairs

550 removed, nucleotides removed, gaps removed and positions removed are median percentage

551 of the total in the test alignment which was removed by CIAlign. Normalised Robinson-Foulds

552 distance and Quartet divergence are the mean proportion similarity between the benchmark

553 tree and the tree generated based on the test alignment. Consensus percentage identity is

554 the mean alignment similarity between the consensus and the EvolvAGene input sequence.

555

## Figure Legends

557 **Fig 1**

558 **Mini alignments showing the main functionalities of CIAlign based on Example 1.**

559 **a** Input alignment before application of CIAlign, generated using the command "`CIAlign`

560 `--infile example1.fasta --plot_input`"*.* **b** Output alignment showing the

561 functionality of the remove divergent function, generated using the command "`CIAlign`

562 `--infile example1.fasta --remove_divergent --plot_output`". **c** Output

27

563    alignment showing the functionality of the remove insertions function, generated using the

564    command "`CIAlign --infile example1.fasta --remove_insertions`

565    `--plot_output`". **d** Output alignment showing the functionality of the crop ends function,

566    generated using the command "`CIAlign --infile example1.fasta --crop_ends`

567    `--plot_output`". **e** Output alignment showing the functionality of the remove short

568    sequences function, generated using the command "`CIAlign --infile`

569    `example1.fasta --remove_short --plot_output`". **f** Output alignment showing the

570    functionality of the remove gap only function, generated using the command "`CIAlign`

571    `--infile example1.fasta --plot_output`". Subplots were generated using the

572    `drawMiniAlignment` function of CIAlign.

573

574    **Fig 2**

575    **Crop ends diagram**

576    This manually created example illustrates how crop_ends works internally. The length of the

577    sequence shown is 111 including gaps and 80 excluding gaps (1). With a threshold of 10% for

578    the proportion of non-gap positions to consider for change in end positions, 8 positions at the

579    start and at the end, respectively, are being considered (illustrated by red crossbars). For

580    each of these, the number of preceding gaps is calculated (2). Then the change in gap

581    numbers (3) for every two consecutive non-gap positions is compared to the gap number

582    change threshold, which is 5%, i.e. 4 gaps, as a default value. Looking at the change in gap

583    numbers, the last change at each end equal to or bigger than the threshold is coloured in red.

28

584   This leads to redefining the start and the end of this example sequence to be where the

585   nucleotides are coloured in green.

586

587   **Fig 3**

588   **Mini alignments and legends showing further functionalities of CIAlign based on**

589   **Example 1**.

590   **a** Alignment showing the functionality of the plot markup function, generated using the

591   command "`CIAlign --infile example1.fasta --all`". The areas that have been

592   removed are marked up in different colours, each corresponding to a certain function of

593   CIAlign. **b** Output alignment after application of all functions of CIAlign combined, generated

594   using the command "`CIAlign --infile example1.fasta --all`". Subplots were

595   generated using the `drawMiniAlignment` function.

596

597   **Fig 4**

598   **Mini alignments showing the main functionalities of CIAlign based on Example 2.**

599   **a** Input alignment before application of CIAlign, generated using the command "`CIAlign`

600   `--infile example2.fasta --plot_input`". **b** Alignment markup showing areas that

601   were removed by CIAlign, generated using the command "*CIAlign --infile example2.fasta*

602   *--all*". **c** Output alignment after application of CIAlign, generated using the command

603   "`CIAlign --infile example2.fasta --all`". Subplots were generated using the

604   `drawMiniAlignment` function.

605

29

606 **Fig 5**

607 **Mini alignments showing the main functionalities of CIAlign based on Example 3**.

608 **a** Input alignment before application of CIAlign, generated using the command "`CIAlign`

609 `--infile example3.fasta --plot_input`". **b** Output alignment after application of

610 CIAlign, generated using the command "`CIAlign --infile example3.fasta --all`

611 `--remove_divergent_minperc 0.5`". Subplots were generated using the

612 `drawMiniAlignment` function.

613

614 **Fig 6**

615 **Metrics from benchmarking CIAlign with simulated data.**

616 **a** Box plots showing the impact of running CIAlign cleaning functions with relaxed (green, left

617 box), moderate (blue, middle box) and stringent (red, right box) parameter values on

618 alignments of sequences simulated using either EvolvAGene [36] or INDELible [37] (plots are

619 combined for the two tools). From left to right, the y-axis represents proportion of correctly

620 aligned pairs of residues [38] removed (identified by comparison with a benchmark

621 alignment), proportion of total nucleotides (i.e. non-gap positions) removed, proportion of

622 gaps removed, proportion of positions (gap or non-gap) removed. **b** Histograms showing the

623 distribution of normalised Robinson-Foulds distances [34] and Quartet divergence [35]

624 between benchmark trees and test trees without running CIAlign cleaning functions (orange)

625 and after running CIAlign with the three sets of parameter values, for trees based on

626 simulated sequences generated with EvolvAGene [36] (left two columns) and INDELible [37]

627 (right two columns). **c** Density plot showing the distribution of the percentage identity between

30

628  the input sequence to EvolvAGene [36] and a consensus sequence based on an alignment of

629  the simulated sequences generated by this tool, without running CIAlign (orange) and after

630  running CIAlign cleaning functions with the three sets of parameter values . **d** Density plots

631  showing the distribution of the percentage identity between the input sequence to BadRead

632  [39] and a consensus sequences generated with (blue) and without (orange) running CIAlign

633  cleaning functions for alignments of good (top), medium (middle) and poor (bottom) quality

634  simulated reads.

635

636  **Fig 7**

637  **Metrics from using CIAlign with biological data.**

638  **a** Left, density plots showing the distribution of percentage identity (top) and normalised

639  Needleman-Wunsch score (bottom) between samples of sequences from the Pfam [40] full

640  alignments and consensus sequences generated based on Pfam seed alignments without

641  (light blue) and with (light red) CIAlign cleaning and Pfam full alignments without (dark blue)

642  and with (dark red) CIAlign cleaning. Right, box plots showing the alignment total size (top)

643  and number of gaps (bottom) for these four alignments. **b** Left, bar chart showing the size of

644  insertions from the 1000 genomes data [41] used to test the ability of CIAlign to remove

645  insertions and deletions. Right, bar chart showing the proportion of sequences in which these

646  insertions were present in data from 162 individuals and whether they were (pink) or were not

647  (blue) removed by the CIAlign remove insertions function. **C** Left, phylogenetic tree based on

648  an alignment of sequences from the 10k trees project [42] for the 12s ribosomal gene in

649  primates. Colours represent known monophyletic groups of primates. Nodes have been

650  collapsed where multiple sequences from the same group formed a monophyletic clade.

31

651 Sequences annotated with circles were removed by CIAlign. Top-right, tree based on the

652 same alignment after cleaning with CIAlign, which removed the outlying group. Bottom-right,

653 mini alignments showing the effect of running CIAlign on this alignment.

654

655 **Supplementary Information**

656 **Supplementary Figure 1**

657 **Mini alignments and phylogenetic trees showing the application of CIAlign to**

658 **phylogenetic data, based on Example 4, a subset of Example 3. a** Input alignment before

659 application of CIAlign, generated using the command "*CIAlign --infile*

660 *example4.fasta --plot_input*". **b** Output alignment after application of CIAlign,

661 generated using the command "*CIAlign --infile example4.fasta --all*

662 *--remove_divergent_minperc 0.5*". Subplots were generated using the

663 "drawMiniAlignment function. **c** Phylogenetic tree generated manually using the literature to

664 show the current best estimate for the phylogenetic relationships between these 91 families of

665 metazoa. Relationships are based on the literature listed in the Supp. References. **d** PhyML

666 phylogenetic tree generated under the GTR model plus default settings on the input alignment

667 before application of CIAlign. **e** PhyML phylogenetic tree generated under the GTR model

668 plus default settings on the cleaned alignment after application of CIAlign. In (c-e) branch

669 colours correspond to the labelled phyla, coloured squares indicate class and bold text

670 indicates order. Common names are shown where available.

671

672 **Supplementary Figure 2**

32

673 **a** Left, phylogenetic tree based on an alignment of sequences from the 10k trees project [42]

674 for the *APOB* gene in Carnivora. Colours represent known monophyletic families of

675 Carnivora. Nodes have been collapsed where multiple sequences from the same family

676 formed a monophyletic clade. Sequences annotated with circles were removed by CIAlign.

677 Top-right, tree based on the same alignment after cleaning with CIAlign, which removed the

678 outlying group. Bottom-right, mini alignments showing the effect of running CIAlign on this

679 alignment. **b** As for **a**, but for the *RAG1* gene in Carnivora.

680

681 **Supplementary Table 1**

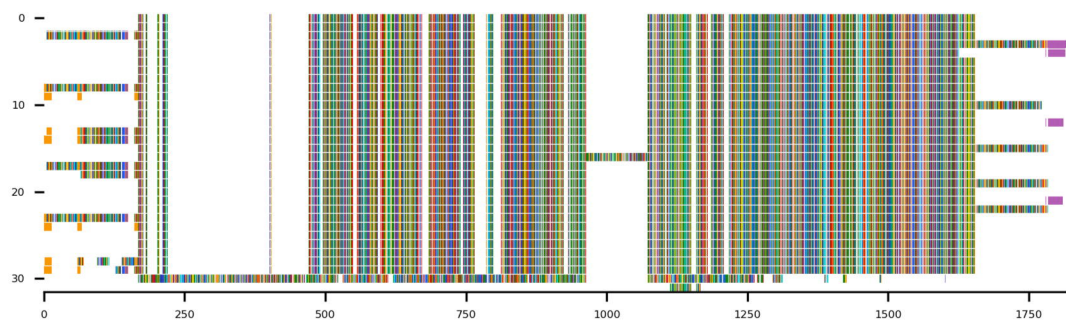682 Relaxed, moderate and stringent parameter settings used for benchmarking.

683

Figure 1



**Nucleotides**

A    C    T    G

A — Input

B — Remove Divergent

C — Remove Insertions

D — Crop Ends

E — Remove Short

F — Remove Gap Only

Figure 2

Figure 3

Figure 4



Amino Acids
Asp, Glu | Lys, Arg | Phe, Tyr | Gly | Ala | Trp | His | Pro
Met, Cys | Ser, Thr | Gln, Asn | Leu, Ile, Val

Mark-up
Cropped Ends
Too divergent
Insertions
Too short
Gap only

Figure 5



**Nucleotides**
A    C    T    G

Figure 6



**A** Correct Pairs Removed · Nt Removed · Gaps Removed · Positions Removed

**B** EvolvAGene · INDELible (None, Relaxed, Moderate, Stringent)

**C** Identity Between Original Sequence and Alignment Consensus

**D** % Identity: Good Quality Reads · % Identity: Medium Quality Reads · % Identity: Poor Quality Reads

Figure 7

**A**

Original

Cleaned with ClAlign

Original

Cleaned with ClAlign

**B**

Original

Cleaned with ClAlign

Original

Cleaned with ClAlign