

FASTA Herder: a web application to trim protein sequence sets

Caroline Louis-Jeune¹, Miguel A. Andrade-Navarro², and Carol Perez-Iratxeta^{*1,3}

¹Ottawa Hospital Research Institute, 501 Smyth Road, Ottawa, ON K1H 8L6, Canada

²Max Delbrück Center for Molecular Medicine, Robert-Rössle-Str. 10, 13125 Berlin, Germany

³Department of Cellular and Molecular Medicine, Faculty of Medicine, University of Ottawa, Ottawa, ON, Canada

*Corresponding author's e-mail address: cpereziratxeta@gmail.com

Published online: 23 July 2015 (version 1); 15 September 2015 (version 2)

Cite as: Louis-Jeune et al., ScienceOpen Research 2015 (DOI: 10.14293/S2199-1006.1.SOR-LIFE.A67837.v2)

Reviewing status: Please note that this article is under continuous review. For the current reviewing status and the latest referee's comments please click [here](#) or scan the QR code at the end of this article.

Primary discipline: Life sciences

Secondary discipline: Quantitative & Systems biology, Bioinformatics & Computational biology

Keywords: FASTA, multiple sequence alignment, evolutionary conservation, BLAST, sequence similarity, protein primary structure

ABSTRACT

The ever increasing number of sequences in protein databases usually turns out large numbers of homologs in sequence similarity searches. While information from homology can be very useful for functional prediction based on amino acid conservation, many of these homologs usually have high levels of identity among themselves, which hinders multiple sequence alignment computation and, especially, visualization. More generally, high redundancy reduces the usability of a protein set in machine learning applications and biases statistical analyses. We developed an algorithm to identify redundant sequence homologs that can be culled producing a streamlined FASTA file. As a difference from other automatic approaches that only aggregate sequences with high identity, our method clusters near-full length homologs allowing for lower sequence identity thresholds. Our method was fully tested and implemented in a web application called FASTA Herder, publicly available at <http://fh.ogic.ca/>.

INTRODUCTION

Multiple sequence alignment (MSA) remains the most important analytic tool to assess evolutionary relations between proteins and to determine the conserved regions of the sequence that usually harbor structural and functional properties. However, the algorithms that produce MSAs are not perfect and further manual curation carried through by visual inspection is necessary [10]. Inspection is also required to detect or confirm some features of the aligned proteins, such as the presence of a structural domain. Notwithstanding the availability of excellent MSA visualization tools [7], the growth of the protein database is making difficult these tasks, as well as the calculation of the alignment itself (e.g. Uniref100 [9], which currently contains over 23 million sequences, has nearly doubled every two years since 2008). Removing highly redundant sequences in a set of homologs helps MSA and its

interpretation as well as reducing biases in the protein set when it is taken as a sample for statistical analysis or for machine learning applications. In a previous work, we conceived an algorithm to cluster the whole protein database by aggregating near-full length homologs [5]. Instead of clustering sequences based just on high sequence identity, our approach allows low but significant levels of identity according to BLAST [1] and is very restrictive in requiring this similarity to be full length without large unmatched portions between the compared sequences. This is to ensure that when two sequences are aggregated they have the same domain structure, which increases the likelihood of both having a similar function [6]. We have implemented this algorithm in a web tool to quickly organize and/or trim a set of proteins in FASTA format.

APPROACH

Given a set of proteins, we use a greedy algorithm to group similar sequences together. First, we pick the longest sequence as the query sequence and we compare it with BLAST to the subset of sequences that either have the same length as the query or are shorter by less than n amino acids (aa). We will retain BLAST matches covering the full length of the query or shorter by less than $n / 2$ aa from each end. The value of n depends on the length of the query sequence, q , and is by default equal to 32, if $q \geq 200$; 20, if $200 > q \geq 100$; 10, if $100 > q \geq 60$; 7, if $60 > q \geq 40$; 5, if $40 > q \geq 20$, and 2 if $q < 20$. Sequences are considered to match the query sequence if they are sufficiently similar to the query as detected by a BLAST single hit (either continuous or gapped) that covers all (or almost all) of the protein. An identity threshold of 24.8% was used for BLAST hits if $q \geq 80$, else a function of q equal to $290.15 \cdot q - 0.562$ for shorter hits, following Sander and Schneider [8]. Query and matched sequences are aggregated and removed from the set and the longest sequence among

the remaining ones is chosen as a new query. This procedure is repeated until no sequences remain.

BENCHMARKS

To validate the reduction in number of sequences in FASTA files by our method, we used the OrthoBench, a published benchmark set of manually curated protein families that was designed to assess orthologs' detection methods [11] and is available at <http://eggno.embl.de/orthobench/>. The OrthoBench consists of 7,981 proteins grouped into 100 protein families classified according to the rate of evolution, domain architecture, low complexity/repeats, lineage-specific loss/duplication, horizontal gene transfer, or alignment quality. We excluded family RefOG065 because among its 19 members 15 of them are also found in RefOG002, and we worked with 7,962 unique sequences in 99 groups (Supplementary Tables 1 and 2). We sought to evaluate how our algorithm avoids clustering together sequences from different families despite allowing very low sequence identity. For that, we noted the level of reduction of the number of sequences and the number of errors (when sequences belonging to different families are clustered together) as performance indicators. The overall method performance with the default parameters was good, reducing the number of sequences down to 14.41% of the original FASTA size with 60 misclassifications that affect all the same cluster. They are all members of group RefOG086, Membrane GTPase involved in stress response (87 sequences), clustering with a member of the RefOG035 family, ATP-binding cassette (318 sequences). For detailed compression values on individual families, see Supplementary Table 3. Our default parameters are conservative. We tested how different thresholds of tolerance of length differences between sequences allowed to cluster together would affect performance. Results indicate that moderately increasing tolerance increases the reduction in number of sequences without largely affecting the number of clusters affected and misclassifications (Table 1). Some proteins contain low-complexity regions (LCRs), regions with repetitive sequences and little diversity in their amino acid composition. Hypothesizing that the length of LCRs will be extremely variable within protein families and might cause misclassification of orthologs in different clusters, we tested how detecting and masking these from the original sequences would affect the results. For conservative parameters of LCR detection, our algorithm achieved a similar reduction in the number of sequences without increasing the number of errors (Table 1). As this is heavily dependent on particular protein sets, we decided to provide LCR detection and filtering as a user option. The proportion of misclassifications are only indicative because in the OrthoBench, proteins from different families can be highly divergent, and therefore the number of mistakes could be expected to be smaller than in real-life applications. Indeed the rate of compression for individual families depends obviously of the degree of evolutionary divergence of the sequences, especially when that would involve large changes

Table 1. FASTA Herder on the Orthobench.

Tolerance <i>n</i>	Compression (%)		Affected clusters		Misclassifications	
	FH	FH with LCR	FH	FH with LCR	FH	FH with LCR
0	14.41	13.32	1	1	60	1
4	13.61	12.48	1	1	62	1
16	11.65	10.53	1	1	63	3
20	11.07	9.97	1	1	65	3
48	8.76	7.92	3	3	29	4

Tolerance *n*: Additional number of amino acids permitted over the default values (*n* = 0 corresponds to the default parameters). Compression (%): number of clusters as percentage of the original FASTA file. FH: without using SEG, FH with LCR: using SEG to detect LCRs with parameters (*W* = 5 and *K1* = 1.0). Affected clusters: number of clusters with misclassified sequences. Misclassifications: total number of misclassified sequences.

Table 2. Compression averages of the Orthobench individual families by class.

Class	Families	Compression (%)	Similarity range (%)	
			Min	Max
Alignment quality (high)	3	35.88	34.2	99.8
Alignment quality (low)	8	37.46	27.6	99.8
Domain shuffling/evolution	4	56.78	34.9	99.3
Horizontal gene transfer	6	6.49	36.1	100
LCRs/repeats	9	31.84	29.3	99.7
Multigene families/ paralogy	6	16.38	22.6	99.9
Pangenomes	10	16.27	41.7	99.6
pN/pS (high)	8	14.10	42.8	100
pN/pS (low)	6	11.78	44.9	99.9
Random families	29	33.17	29.0	99.4
Speed of evolution (high)	8	59.11	28.6	88.6
Speed of evolution (low)	1	7.69	88.2	100

Class: class name. Families: number of families in class. Compression (%): average of compression for class families. Similarity range (%): Averages of min and max similarity values for class families.

in protein length or domain shuffling (Table 2). We compared FASTA Herder to two other published methods that use similar greedy clustering strategies, CD-HIT [4] and kClust [3]. Both methods are very fast and geared toward the compression of large databases. We calculated the compression rates and errors at the lowish identity thresholds that FASTA Herder uses for the same benchmark set. Without restrictions to the coverage, as the minimum of the percentage of the aligned length for the longest sequence, both tools produced very compressed results but usually more errors (Table 3). Limiting the admissible clusters to high levels of coverage produces comparable results for the three tools (Table 3). Another previously published web tool that possess a capability similar to this work is PISCES [12], a protein database that culls the Protein Data Bank to build the largest possible set of structures that comply with identity cut-offs for sequence and structure. More related to our tool, PISCES

Table 3. CD-HIT and kClust in the Orthobench.

	Identity (%)	Coverage (%)	Compression (%)	Affected clusters	Misclassifications
CD-HIT	25	20	2.84	4	9
		40	3.18	3	5
		60	4.11	2	4
	30	90	8.74	2	4
		20	3.52	3	5
		40	3.88	2	4
		60	4.57	1	3
		90	8.99	1	3
		20	3.17	4	204
kClust	25	40	3.74	2	4
		60	4.71	2	4
		90	11.22	1	22
	30	20	3.92	3	22
		40	4.42	2	4
		60	5.25	1	3
		90	11.25	1	23

Identity (%): sequence identity thresholds. Coverage (%): Alignment coverage of the leader sequence. Compression (%): percentage of the original size. Affected clusters: number of clusters with misclassified sequences. Misclassifications: total number of misclassified sequences.

accepts as well a FASTA file and culls out sequences by a single identity threshold that can be provided by the user. However PISCES compares all to all sequences, it is much slower and accepts smaller inputs than the other tools.

FASTA Herder takes a protein set in FASTA file format and quickly clusters it (e.g. the OrthoBench is clustered in under 2 minutes). Optionally, LCRs, as detected with SEG [13], can be ignored before clustering. The server also permits to adjust the stringency of the clustering based on the allowed difference in length between sequences to be clustered together, although it may be important to stick to conservative levels when culling sequences from a MSA. We use the BLAST implementation of the BLAST+ suite [2].

CONCLUSION

We have implemented an ease-of-use web application to quickly reduce the redundancy of a protein set by clustering homologs of comparable lengths. Although our work is somewhat related to the field of orthology prediction methods (discussed in Ref [11]), its purpose and scope are different. Orthology detection encompasses the use of whole genomes to identify proteins derived from a single ancestral sequence through speciation events. The aim of our tool is to reduce potentially high redundancy in a FASTA file that may contain proteins belonging to one or more families. This would speed MSA calculation, facilitate MSA inspection, and remove biases that may affect any statistical analysis or computational application involving that set of proteins. In conclusion, we have created a tool to remove redundant sequences from sets of sequences. Simplification of MSAs was our main goal, motivated by the increase in redundant sequences in the databases due to genomic sequencing projects. Our method is simple, works in a matter of seconds with large datasets that are unmanageable by other methods, and includes a number

of options (LCR filtering, length thresholds) that make it very flexible.

FUNDING

This work was supported (in part) by a Government of Ontario Ministry of Research and Innovation (ORF-RE05-084) grant.

SUPPLEMENTARY MATERIAL

Supplementary material is freely available [here](#).

REFERENCES

- [1] Altschul SF, Madden TL, Schaffer AA, Zhang J, Zhang Z, Miller W, Lipman DJ. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res.* 1997;25(17):3389–3402. doi:10.1093/nar/25.17.3389
- [2] Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL. BLAST+: architecture and applications. *BMC Bioinformatics.* 2009;10:421. doi:10.1186/1471-2105-10-421
- [3] Hauser M, Mayer CE, Söding J. kClust: fast and sensitive clustering of large protein sequence databases. *BMC Bioinformatics.* 2013;14:248. doi:10.1186/1471-2105-14-248
- [4] Huang Y, Niu B, Gao Y, Fu L, Li W. CD-HIT Suite: a web server for clustering and comparing biological sequences. *Bioinformatics.* 2010;26(5):680–2. doi:10.1093/bioinformatics/btq003
- [5] Perez-Iratxeta C, Palidwor G, Andrade-Navarro MA. Towards completion of the Earth's proteome. *EMBO Rep.* 2007;8(12):1135–1141. doi:10.1038/sj.embor.7401117
- [6] Ponting CP, Schultz J, Copley RR, Andrade MA, Bork P. Evolution of domain families. *Adv Protein Chem.* 2000;54:185–244.
- [7] Procter JB, Thompson J, Letunic I, Creevey C, Jossinet F, Barton GJ. Visualization of multiple alignments, phylogenies and gene family evolution. *Nat Methods.* 2010;7(Suppl 3):16–25. doi:10.1038/nmeth.1434
- [8] Sander C, Schneider R. Database of homology-derived protein structures and the structural meaning of sequence alignment. *Proteins.* 1991;9(1):56–68. doi:10.1002/prot.340090107

- [9] Suzek BE, Huang H, McGarvey P, Mazumder R, Wu CH. UniRef: comprehensive and non-redundant UniProt reference clusters. *Bioinformatics*. 2007;23(10):1282–1288. doi:10.1093/bioinformatics/btm098
- [10] Thompson JD, Linard B, Lecompte O, Poch O. A comprehensive benchmark study of multiple sequence alignment methods: current challenges and future perspectives. *PLoS ONE*. 2011;6(3):e18093. doi:10.1371/journal.pone.0018093.t001
- [11] Trachana K, Larsson TA, Powell S, Chen WH, Doerks T, Muller J, Bork P. Orthology prediction methods: a quality assessment using curated protein families. *Bioessays*. 2011;33(10):769–780. doi:10.1002/bies.201100062
- [12] Wang G, Dunbrack RL. PISCES: recent improvements to a PDB sequence culling server. *Nucleic Acids Res*. 2005;33(Web Server issue):W94–98.
- [13] Wootton JC, Federhen S. Analysis of compositionally biased regions in sequence databases. *Meth Enzymol*. 1996;266: 554–571.

COMPETING INTERESTS

The authors declare no competing interests.

PUBLISHING NOTES

© 2015 Louis-Jeune et al. This work has been published open access under Creative Commons Attribution License **CC BY 4.0**, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited. Conditions, terms of use and publishing policy can be found at www.scienceopen.com.

Please note that this article may not have been peer reviewed yet and is under continuous post-publication peer review. For the current reviewing status please click [here](#) or scan the QR code on the right.



scienceOPEN.com
research+publishing network