

PAL: A Perl Script for Rapidly Identifying the Active Site of Large Protein Families
Andres A. Larrea.^{1,3} and Pablo A. Larrea²

¹ University of Miami Miller School of Medicine Department of Biochemistry and Molecular Biology

² Unaffiliated

³ Corresponding author:
Department of Biochemistry and Molecular Biology
University of Miami Miller School of Medicine
1011 NW 15 St
Miami, FL 33136
ALarrea@med.miami.edu

Keywords: Multiple sequence alignment analysis; active site determination; protein domains

Software URL: <http://homepage.mac.com/amlarrea/FileSharing3.html>

Word Count: 899

1 In the post-genomic era, with an ever-increasing amount of sequence information, it is critical to
2 develop tools that assist in sifting through all the data to identify information amidst the noise.
3 The last ten years have seen an explosion in DNA sequence data. As of late March 2007, there
4 were 921 genomic sequence sets in the NCBI database (with 495 completely sequenced
5 genomes) (1). With the continuing rapid expansion in genome data, it is becoming increasingly
6 important to develop tools that assist in the rapid analysis of large numbers of sequences. The
7 ability to parse through large amounts of sequence data is vital as this allows us to identify
8 patterns that can be used to predict functions for unknown open reading frames (ORFs).
9 The two main steps in determining the functional sites in proteins are: (1) aligning a large set of
10 homologous sequences and (2) analyzing this large file to identify highly conserved residues. A
11 great deal of software exists for generating a quality alignment (*e.g.* Clustal (2-4), Muscle (5, 6),
12 T-Coffee (7-9), etc). Although, historically computation time has been a limiting factor in
13 aligning large files, with improvements in computer technology as well as in the algorithms,
14 some of these programs are able to generate quality multiple sequence alignment (MSA) files
15 very quickly. Generating an MSA file is, however, only the first step in analyzing proteins. The
16 real goal is to allow genetic variation and natural selection to point us in the direction of
17 important amino acid residues as evident in the degree of conservation and pattern of inheritance
18 in phylogenetic trees(10). Unfortunately there is no software able to quickly analyze large
19 multiple sequence alignment files. Most of the software that exists for post-alignment analysis is
20 Graphical User Interface (GUI) based (JalView (11), SeaView (12), ClustalX (4)). Although
21 GUIs have several advantages, looking at large amounts of data and identifying important
22 residues, is not one of them.

Once an MSA file has been created, it is important to be able to sift through it quickly to identify conserved residues that often mark functionally important protein features like enzyme active sites. To this end, we have written PAL, a small highly customizable Perl script. While the inputs and outputs are straight-forward (**Figure 1**) they can be looked at in a variety of ways to solve several problems faced by investigators as they analyze multiple sequence alignments:

1. Generating a high quality alignment is an iterative process. It is important to include a large number of sequences in a multiple sequence alignment, but sequences that are too divergent can introduce noise that obscures gene function, while identical sequences offer no predictive benefit. By analyzing the *nearest.txt* (**Figure 1G**) file, users can quickly identify which sequences should be removed before the next alignment cycle.

2. While enzyme specificity is evolutionarily flexible, novel chemistry is far more difficult to evolve and active site residues are predicted to be highly conserved (13). Once the number of sequences being used has been trimmed down and the alignment is finalized, the files *consensus.txt* (**Figure 1D**) and *nearconsensus.txt* (**Figure 1E**) can be used to predict putative enzyme active sites. These sites can then be verified experimentally using site-directed mutagenesis and enzymatic analysis approaches.

3. Although enzyme chemistry is difficult to evolve and reflected in highly conserved regions in a protein sequence, substrate specificity and protein-protein interaction surfaces are often sites of divergence between homologs. Examination of the file *letters.txt* (**Figure 1F**), output by PAL, can be used to identify highly divergent regions that may be involved in specificity.

4. With very large protein families, it can be helpful to divide the sequences into smaller sub-families for analysis. By using the *group.analysis* option in the inputs, users can choose to

divide the full family into sub-families with different amounts of homology. Each sub-family is then individually analyzed for conservation and divergence by PAL.

The only input requirement is an alignment in fasta format. Because of this, PAL can use any alignment software as well as manually edited alignments. What makes PAL so powerful is that it is able analyze very large alignment files in a short period of time (**Figure 2**).

The blessing of an ever-increasing amount of genomic data is also a curse that leads to large MSA files that are difficult to analyze using traditional graphical approaches. In this paper we present a Perl script that assists in the rapid identification of important residues in large protein families. With this software we have been able to identify a putative active site for a large exonuclease family containing over 300 members in less than 2 minutes. The sample outputs shown in **Figure 1** have formed the basis of ongoing work to understand the mechanism of the poorly understood ExoVII nuclease family of enzymes (Larrea A *et. al.*, unpublished data).

Software requirements and availability:

All software described here is open source and freely available upon request. Although all benchmarking has been done on a Mac, much of the programming was done on a Windows platform PC. Since the software is written in Perl it should run the same across all platforms.

Acknowledgements:

This work was supported by a pre-doctoral fellowship from the NIH (F31-GM70395 to AAL). Many thanks to Arun Malhotra and Richard Myers for their invaluable help during the beta testing of the software and during the preparation of the manuscript.

Competing Interest:

The authors have no competing interests.

1. **Cummings, L., L. Riley, L. Black, A. Souvorov, S. Resenchuk, I. Dondoshansky, T. Tatusova.**2002. Genomic BLAST: custom-defined virtual databases for complete and unfinished genomes. *FEMS Microbiol Lett* 2:133-8.
2. **Higgins, D.G., P.M. Sharp.**1988. CLUSTAL: a package for performing multiple sequence alignment on a microcomputer. *Gene* 1:237-44.
3. **Thompson, J.D., D.G. Higgins, T.J. Gibson.**1994. CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice. *Nucleic Acids Res* 22:4673-80.
4. **Thompson, J.D., T.J. Gibson, F. Plewniak, F. Jeanmougin, D.G. Higgins.**1997. The CLUSTAL_X windows interface: flexible strategies for multiple sequence alignment aided by quality analysis tools. *Nucleic Acids Res* 24:4876-82.
5. **Edgar, R.C..**2004. MUSCLE: a multiple sequence alignment method with reduced time and space complexity. *BMC Bioinformatics* 113.
6. **Edgar, R.C..**2004. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res* 5:1792-7.
7. **O'Sullivan, O., K. Suhre, C. Abergel, D.G. Higgins, C. Notredame.**2004. 3DCoffee: combining protein sequences and structures within multiple sequence alignments. *J Mol Biol* 2:385-95.
8. **Poirot, O., K. Suhre, C. Abergel, E. O'Toole, C. Notredame.**2004. 3DCoffee@igs: a web server for combining sequences and structures into a multiple sequence alignment. *Nucleic Acids Res Web Server issue*:W37-40.
9. **Poirot, O., E. O'Toole, C. Notredame.**2003. Tcoffee@igs: A web server for computing, evaluating and combining multiple sequence alignments. *Nucleic Acids Res* 13:3503-6.

10. **Livingstone, C.D., G.J. Barton.**1996. Identification of functional residues and secondary structure from protein multiple sequence alignment. *Methods Enzymol* 497-512.
11. **Clamp, M., J. Cuff, S.M. Searle, G.J. Barton.**2004. The Jalview Java alignment editor. *Bioinformatics* 3:426-7.
12. **Galtier, N., M. Gouy, C. Gautier.**1996. SEAVIEW and PHYLO_WIN: two graphic tools for sequence alignment and molecular phylogeny. *Comput Appl Biosci* 6:543-8.
13. **Petsko, G.A.**2000. Enzyme evolution. Design by necessity. *Nature* 6770:606-7.

Figure 1. PAL inputs and outputs

(A) Table of inputs and outputs for PAL. (B-G) Sample outputs from PAL analysis of an alignment with 345 sequences homologous to *E. coli* XseA (gi:16130434). (B) Portion of *names.txt*. Section shown only includes sequences numbered 181 to 191. This file references sequence numbers to sequence identifiers. (C) Portion of *lettercnt.txt*. This file shows amino acid or gap occurrence as a function of position along the alignment. For example, the top line indicates that 285 homologs out of 345 contain a gap at position 181. (D) Complete output from *consensus.txt*. This file lists residues that are 100% conserved among all 345 sequences. These residues are likely to form the active site and work is underway to verify this prediction experimentally. (E) Complete output from *nearconsensus.txt*. This file lists positions that only include amino acids defined as similar in the *near.match* file. (F) Portion of *letters.txt*. This is an alternate form of the data shown in *lettercnt.txt*. In this case each position along the alignment is listed with the residues found at that position in decreasing abundance. For position 181 in the alignment, a gap is seen most often, and Met is the next most abundant residue. The region shown (position 181-190 in the sequence) identifies a putative loop missing in most XseA homologs. (G) Portion of output from *nearest.txt*. Matrix of each sequence, with the pairwise value to itself, and to every other sequence. Pairwise values are obtained using the matrix chosen (in this case *blosum62.mtx*). Top ten scores are shown.

Figure 2. Benchmarking PAL

Benchmarking of PAL on G4 PPC and Intel C2D processors of alignments generated by ClustalW with either 145 or 345 sequences.

A

Filename	Type	Description
INPUTS		
XXX.fasta	alignment	aligned protein sequences in FASTA format
XXX.mtx	matrix	alignment matrix to be used. blosum62.mtx and identity.mtx are two supplied but any can be generated.
near.match	text	list of amino acids that are to be treated as equivalent
group.analysis	text	threshold of group definition to be used
OUTPUTS		
names.txt	list	list of the names of all the sequences read (B)
lettercnt.txt	matrix	matrix of position vs. amino acid. expanded form letters.txt (C)
consensus.txt	list	list of residues that are 100% conserved (D)
nearconsensus.txt	list	list of residues that "similar" as defined by near.match (E)
letters.txt	list	list of position vs. amino acid (F)
nearest.txt	matrix	matrix of top ten pairwise scores (G). In each case the sequence is first evaluated against itself and then against every other sequence. The top 10 values are listed first sequence number and then pairwise score.

B		C																				D										E																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																	
Num	Description																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																

