

# Pfam: multiple sequence alignments and HMM-profiles of protein domains

Erik L. L. Sonnhammer\*, Sean R. Eddy<sup>1</sup>, Ewan Birney<sup>2</sup>, Alex Bateman<sup>2</sup> and Richard Durbin<sup>2</sup>

Computational Biology Branch, National Center for Biotechnology Information, National Library of Medicine, Building 38A, Room 8N805, National Institutes of Health, Bethesda, MD 20894, USA, <sup>1</sup>Department of Genetics, Washington University School of Medicine, St Louis, MO 63110, USA and <sup>2</sup>Sanger Centre, Hinxton Hall, Cambridge CB10 1SA, UK

Received September 8, 1997; Revised and Accepted October 8, 1997

## ABSTRACT

**Pfam contains multiple alignments and hidden Markov model based profiles (HMM-profiles) of complete protein domains. The definition of domain boundaries, family members and alignment is done semi-automatically based on expert knowledge, sequence similarity, other protein family databases and the ability of HMM-profiles to correctly identify and align the members. Release 2.0 of Pfam contains 527 manually verified families which are available for browsing and on-line searching via the World Wide Web in the UK at <http://www.sanger.ac.uk/Pfam/> and in the US at <http://genome.wustl.edu/Pfam/> Pfam 2.0 matches one or more domains in 50% of Swissprot-34 sequences, and 25% of a large sample of predicted proteins from the *Caenorhabditis elegans* genome.**

## INTRODUCTION

A relatively small number of structural and functional domains are used in a large number of different proteins. Particularly for protein analysis and annotation in large-scale sequencing projects, there is a growing need for easily interpretable and sensitive detection of common protein domains. A protein containing one or more common domains can produce a morass of hundreds or thousands of BLAST hits when searching single sequence databases (e.g. GenBank, Swissprot, PIR). Although searches can be augmented by tools that condense and summarise results (1), satisfactory annotation of such proteins often becomes a time-consuming and error-prone process. Instead, a search of an organised database of protein domain families can produce more concise results which simplify annotation, domain parsing and functional prediction for a query sequence (2–5). Protein family databases are typically based on multiple sequence alignments of known family members. Conserved features can be recognised in the alignment and given higher weight in searches, which for distant similarities can often render the comparison more sensitive than pairwise alignment approaches.

We present here Pfam (6) release 2.0. Pfam was developed in order to use HMM-profile analysis to complement BLAST analysis in the *Caenorhabditis elegans* genome project. The main distinction between Pfam and most other protein family databases is that for all of Pfam, both the family definition and the search method span entire domains, including not only conserved motifs but also

less-conserved regions, insertions and deletions. HMM-profile methods allow variable conservation and insertions/deletions to be dealt with in a fairly robust way (7,8). Modelling of complete domains should facilitate more biologically meaningful sequence annotation, and, in some cases, more sensitive detection.

## DESCRIPTION OF THE DATA

For each protein domain family in Pfam, there are three important files. The *seed alignment* is a manually verified multiple alignment of a representative set of sequences (Fig. 1). An *HMM-profile* is built from the seed alignment for database searching and alignment purposes. A *full alignment* is generated automatically from the seed HMM-profile by searching Swissprot for all detectable members and aligning them to the HMM-profile. The distinction between seed and full alignments facilitates updating the database; the seed alignments are stable resources, whereas full alignments and HMM-profiles can be generated automatically for any new Swissprot (or other sequence database) release.

Each family has a name, a permanent accession number and a record of the methods used to identify the family members and create the alignments. There is also either a brief description of the usual function and structure of the domain, or (more often) links to other on-line documentation resources such as Prosite and Prints.

Both the seed and the full alignments are subjected to a small array of 'quality control' procedures, to verify that the alignments are sensible, that the HMM-detected sequences in the full alignment include all presumed members of the family in Swissprot and no other sequences, and that the family does not overlap with other Pfam families. The process of generating the Pfam family is iterated, if necessary, until all quality requirements are met.

Most Pfam families are based on, and cross-referenced to, corresponding Prosite or Prints entries. In many cases, however, the definition of which sequences belong to a family differs between the databases. This is a pragmatic consequence of the different search methods used. Prosite and Prints detection relies primarily on short conserved patterns corresponding to superfamily motifs. A Prosite pattern or Prints fingerprint may recognise a highly conserved motif shared amongst an otherwise highly diverged superfamily that Pfam splits into several families; conversely, Pfam may recognise a superfamily that Prosite and Prints classify into several distinct families with distinct motif signatures. For some protein domain families, there may be no motif sufficiently conserved to make a

\*To whom correspondence should be addressed. Tel: +1 301 435 5930; Fax: +1 301 480 9241; Email: [sonnhammer@ncbi.nlm.nih.gov](mailto:sonnhammer@ncbi.nlm.nih.gov)

```

ID SH2
AC PF00017
DE Src homology domain 2
AU Sonnhammer ELL
AL Clustalw
AM hmma -qr
SE Swissprot_feature_table
GA Bic_raw 25 hmms 20
DR PROSITE; PDOC50001;
DR SCOP; lsha; sf;
SQ 58
ABLL_CAEEL/179-254 WYHGKISRSDSEALICS..CITGSFLVRESETSIG...QYITSVRHDC.....RVFHYRINVDNTE..KMPITQEVKFRFLGELVHHH
BLK_MOUSE/117-198 WFFRTISRKDAERQLLAPMKNAGSFLIRESESNKG...AFSLSVKDIIT..TQGEVVKHYKIRSLDNG..GYISPRITFPPLQALVQHY
BTK_HUMAN/281-362 WYKXHMTRSQAEQLLQOE.GREGCVIVRDS.SKAGC...KYTVSVFAKSTGDPQGVIRHVVCSTFQS..QYVLAERKLFSTIPELINVH
CSW_MOUSE/111-186 WFHGNLSGKEAEKILIERGK.NGSFLVRESQSKPG...DFVLSVRETD.....KVTHVMIRNQDK...KYDVGGGESFCGLSELIDHY
CSW_DROME/66-81 WFHPTISGIEAEKLLQEQGF.DGSFLAIRLSSNPG...AFVLSVRRGN.....EVTHIKIQNGND...FFLDYGGKFAFLPELVQY
CTK_HUMAN/122-196 WFHGKISGAEVQQLQPP..EDGLFLVRESARHFG...DYVLCVDFGR.....DVTHYRVLHRDC...HLTIDEAVFFCNLMDMVVEHY
DRK_DROME/60-134 WYGRITRADAEKLLSN..KHGAFILIRISSESPG...DFSLSVKCPD.....GVQHFVLRDAQS..K.FFLWVVFVFNLSLMLVSHY
FER_HUMAN/460-531 WYHGAIPIREAEQELKK...QDGLVLRRESHGKPG...EYVLSVYSDG.....QRRHFIIQYVDN...MYRFFG.TGFSNIPQILIDH
FFS_DROME/438-510 WFHGLVLRREEVRLNN...DGLFLVRETIRNEES..QIVLSVCWNG.....H.KHFIVQTTGEG..NFRFFG.PPFASTIQELIMHO
FFS_FUSV/511-581 WYHGAIPIREAEQELLY...SGDGLVRESQKQ...EYVLSVLDG.....QPRHFIIQADN...LYRLED.DGLPTLLELIDH
FRK_HUMAN/116-193 WFFGAIIRSDAEKQLLYSENKTSFLIRESESKG...EFSLSVLDGA.....VVHFRIRKLDG..GFFLTRRIRFSLMEVSHY
GTPA_HUMAN/181-256 WFHGKIDRTIAEERLRQAGK.SGSLVLRSEDRPG...DFVLSVLSQMN.....VVHFRIRIAMCG..DYIYGG.BRFSSLSDLIGY
GTPA_HUMAN/351-426 WFHGKISRQAEYNLLMTVG.QVCSFLVRESNTPG...DYSLYFRNEN.....IQRFKICPTFNN...QPMGGRYNSIGDLIDHY
NCK_MOUSE/282-356 WYGVKTRHQAEMLNER.GHEGDFLRSESSPN...DFSVSLKAGQ.....KKRHFVQLKET...VYICGKRFSTIMEGLVDHY
P85A_HUMAN/624-698 WNVGSSNRNKAENLLRG..KRDGFLVRES.SKQG...CYACSVVVDG.....EVKHCVINKTATG..YGFAPENLYSLKELVSHY
P85B_BOVIN/618-692 WYVGNIRTOAEMLSG..KRDGFLVRES.SQRC...CYACSVVVDG.....DTKHCVIRKLDG...FLVTRRSTFQSLLELVDHY
PIP4_RAT/668-741 WYHSLTRAQAEHMLMVRPR.DGAFVLRKE.NEPN...SYAISFRAEG...KIKHCVRQEQG...TVMLGNSEFSDVLDVSHY
SEM5_CAEEL/60-136 WYHGKIDRTIAEERLRQAGK.SGSLVLRSEDRPG...DFVLSVLSQMN.....VVHFRIRIAMCG..DYIYGG.BRFSSLSDLIGY
SHC_HUMAN/378-449 WFHGKISRREAEALLQLN...GDFVRESSTTPG...QYVLTGLQSG.....QPKHLLVDPFG...VVRKTRHFRFESVSHLISYH
SRC1_DROME/162-244 WFFENVLRKEADKLLLAENPRGTFVLRSEHNP...GYSLSVKDWED.GGCVHVKHYIKPLDNG..GYIATQTFPSLQALVMAY
SRC2_DROME/214-292 WYVGVSRQRAESLLKQG.DKEGCVVVRKS.STKG...LYTSLHRTKVP...QSHVHYHVKQIARNC..EYVLSERKCCETIPDLINHY
SRK1_SROLA/122-199 WFLGKIRVRAEKMLNQSFNQVGSFLIRDSSTTPG...DFSLSVKDDQ.....RVRHYRVRKLENG..FLVTRRSTFQSLLELVDHY
SRK4_SROLA/122-199 WFFGQVRAEAEKRLMPPFNLLGSLVLRDSSTTPG...DFSLSVRDIID.....RVRHYRVRKLENG..TYFVTRRSTFQSLLELVDHY
STK_HYDAT/126-203 WYFGDVKRAEAEKRLMVRGLPSGTFVLRKAEATVG...NFSLSVRSDG.....SVKHYRVRKLDG...GYFITTRAPFNSLVEVQHY
SYK_HUMAN/15-92 WFFGNITREAEEDYLQVGGMSDGLVLRQESRNYLG...GFALSVAHGR.....KAHYTIERELNG..TYAATAGGKTHSPADLCHYH
SYK_FIG/163-238 WFHGKISRDESEQVILIGSKINGKFLIRAR..DNG...SYALGLLHEG...KVLHYRIDKDKTG..KLSIPGKNNFDLWQLVHY
WYCMNTRSKAEQLLRETE.DKEGGVWVRDS.SDPC...LYTVSLYTRFGGESSGFRHYHKEATSPKXVYLAEKHAFGSIPELIDHY
WYCMNTRSKAEQLLRETE.SKEGCVIVRDS.RHLC...SYTVSVFMGARRSTEAALKHYOIKKNDG...QWVVAERHAFQSIPELIDHY
WYAGPMEBAGAEGLITN..RSDGCTVLRQVKTDTA...EFALSIVKVV...EVKHIKIMTSGC...LYRTEKKAERGLLELVDHY
YES_XIPHE/159-241 WYFCKLSRQDTERLLLLGNERCTFLIRESETKGC...AYSLSLRDMDE.TKDNCKHYKIRKLDNG..GYITTRRQFMSLQMLVSHY
YKFL_CAEEL/20-101 WYHSLTRAEEDVFLQDN...NGDVVLSLDPKPEPRSYLLSVMNPKLDENSSVHKVIVNSVEN...KYFVNNNSFNTIQMLVSHY
ZAT70_HUMAN/163-239 WYHSLTRAEAEERKLYSGAQTGKFLLRPRK.EQG...TYALSILYCK...TVYHYLISQDKAG..KYCIPCEKCFDNLQMLVSHY
ZAT70_MOUSE/10-87 FFGSISRRAEAEHLLKLAGMDGLFLLRQCLRSGL...GYVLSLVRDV.....RHFHPTIERQLNG..TYAATAGGKARCCPAELCFQY
    
```

Figure 1. Example of a typical Pfam entry, the SH2 family. Shown is the flat file record including a reduced version of the seed alignment.

discriminative pattern or fingerprint. (Prosite is increasingly incorporating profiles for these families; these Prosite profiles are very similar to Pfam models.) Only the largest (>15 members) Prosite families were systematically used to construct Pfam entries. For smaller families, constructing an HMM-profile is of less value since the sensitivity is unlikely to improve relative to single-sequence searching, and because a small sample is often non-representative. Of the 71 Pfam families with no corresponding Prosite or Prints entry, 55 were 'discovered' as large clusters in Pfam-B (see below). 24 Pfam families contain links to other World Wide Web (WWW) protein family documentation resources, some of which were gleaned from the ProWeb server (9).

Pfam 2.0 contains 527 families, comprising 39 113 sequence segments and 6.8 million residues in the full alignments. All sequences were taken from Swissprot 34 (10). The alignments are on average 275 residues wide, including gaps. There are on average ~75 members per family in full alignments, and ~22 in seed alignments.

### Pfam-B

For comprehensiveness, all Swissprot sequences not in Pfam are clustered automatically by the program Domainer (2), which also constructs multiple alignments automatically and is the basis for the ProDom protein family database. The quality of these alignments tends to be low, but domain-based automated clustering is a convenient method of identifying large obvious families that need to be targeted for Pfam model construction. Although we do not stably maintain, annotate or produce HMM-profiles of these clusters, we make them available as Pfam-B. Pfam-B 2.0 contains 13 289 clusters, 62 611 sub-sequences, and 8.2 million residues. On average, alignments are 146 residues wide (including gaps) and contain five members.

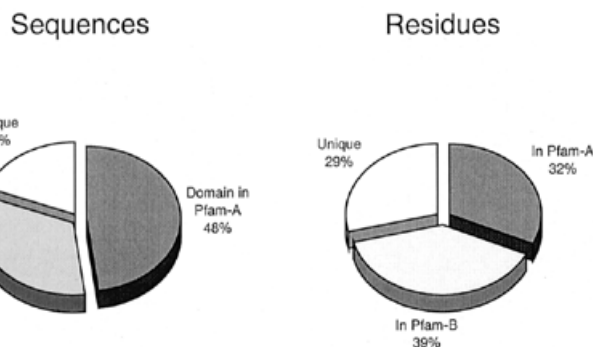


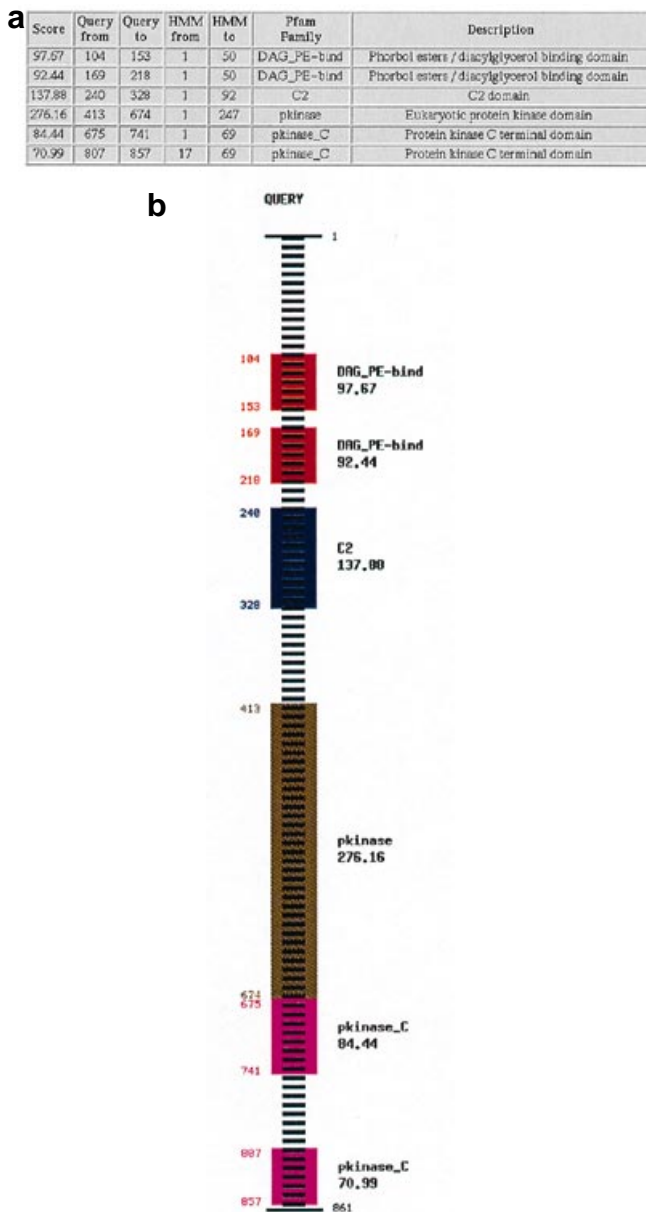
Figure 2. Pfam 2.0 contains domains from nearly half of all Swissprot 34 proteins. The automatic clusters in Pfam-B 2.0 contain domains from 33% of the Swissprot proteins that do not contain Pfam domains. When counting residue-by-residue, roughly a third of Swissprot is covered by Pfam and Pfam-B each. Pfam-B does not include proteins known to be fragments or segments shorter than 30 residues; the figures for unique sequences are therefore overestimated.

### Sequence database coverage

As shown in Figure 2, 48% of the sequences and 32% of the residues in Swissprot 34 are included in annotated Pfam alignments. If unannotated Pfam-B clusters are also taken into account, 81% of sequences and 71% of residues in Swissprot 34 are included in Pfam. In searches of a large and presumably unbiased set of predicted protein sequences from the *C.elegans* genome, 25% of sequences and 13% of residues show significant hits to Pfam HMM-profiles. The numbers are slightly lower for prokaryotic genomes.

### SEARCHING Pfam

The US and UK Pfam WWW servers provide users the ability to search query protein sequences against one, all, or a few Pfam



**Figure 3.** Tabular output (a) and schematic output (b) from a Pfam search with the *C.elegans* protein E01H11.1 as query. Both pictures were taken from the Washington University WWW server.

HMMs. Results are returned in tabular format, and both GIF- and Java-based graphical representations are available optionally. An example of the results from such a search is shown in Figure 3. Here, the *C.elegans* Kin-11 gene product (E01H11.1) is shown to possess a duplicated phorbol esters/diacylglycerol binding domain (DAG/PE-bind), a C2 domain, a protein kinase catalytic domain (pkkinase) and a duplicated domain frequently associated C-terminally to protein kinase domains (pkkinase\_C).

Users can also use Pfam HMM-profiles to search protein sequences locally using the freely available HMMER software package at <http://genome.wustl.edu/eddy/hmmer.html#hmmer>. For comparing genomic and EST data to Pfam HMM-profiles, the programs GeneWise and ESTWise (11) are available at <http://www.sanger.ac.uk/Software/Wise2/>

## WORLD WIDE WEB SERVERS, FTP ACCESS AND FORMAT

The Pfam home pages are <http://www.sanger.ac.uk/Pfam/> at the Sanger Centre in the UK and <http://genome.wustl.edu/Pfam/> at Washington University in the USA. The two servers are separately maintained and differ slightly in their services and capabilities, but are based on the same underlying Pfam database. Both servers support HMM searching, browsing of the family alignments and documentation and lookup of the domain organisation of proteins in Swissprot.

The entire database, including accessory data files such as Pfam schematics for Swissprot proteins, is also available as flat file format ASCII files by anonymous FTP at <ftp.sanger.ac.uk> and <genome.wustl.edu> in `/pub/databases/Pfam/`

The format of the Pfam alignment flat files is based on the EMBL/Swissprot two-character field labels. The following Pfam-specific labels are used: AL, alignment method of seed members; AM, alignment method of full alignment; AU, author responsible for the alignments; GA, gathering method/search program and cutoffs used to build full alignment; SE, source suggesting the seed members belong to the same family; SQ, number of sequences (and last line before the alignment starts). The alignment is in a simple format (Fig. 1) which consists of one line per subsequence containing the Swissprot sequence ID, start and end of the segment, and the aligned subsequence itself (no length limit). In the Pfam flat file, the corresponding Swissprot accession number is added to the right of each alignment line. Users of the Pfam database or WWW servers should cite this article as the appropriate reference.

## ACKNOWLEDGEMENTS

We thank Robert Finn for preparing most of the new families for Pfam 2.0, and Jose Aguilar for writing and maintaining the Washington University Pfam server. Pfam development in SRE's group is supported by grant R01-HG01363 from the NIH National Human Genome Research Institute. Pfam development at the Sanger Centre is supported by the Wellcome Trust.

## REFERENCES

- Sonnhammer,E.L.L. and Durbin,R. (1994) *Comput. Appl. Biosci.*, **10**, 301-307.
- Sonnhammer,E.L.L. and Kahn,D. (1994) *Protein Sci.*, **3**, 482-492.
- Attwood,T.K., Beck,M.E., Bleasby,A.J., Degtyarenko,K., Michie,A.D. and Parry-Smith,D.J. (1997) *Nucleic Acids Res.*, **25**, 212-217 [see also this issue (1998) *Nucleic Acids Res.* **26**, 304-308].
- Bairoch,A., Bucher,P. and Hofmann,K. (1997) *Nucleic Acids Res.*, **25**, 217-221.
- Henikoff,J.G., Pietrokovski,S. and Henikoff,S. (1997) *Nucleic Acids Res.*, **25**, 222-226 [see also this issue (1998) *Nucleic Acids Res.* **26**, 309-312].
- Sonnhammer,E.L.L., Eddy,S.R. and Durbin,R. (1997) *Proteins*, **28**, 405-420.
- Krogh,A., Brown,M., Mian,I.S., Sjoelander,K. and Haussler,D. (1994) *J. Mol. Biol.*, **235**, 1501-1531.
- Eddy,S.R. (1996) *Curr. Opin. Struct. Biol.*, **6**, 361-365.
- Henikoff,S., Endow,S.A. and Greene,E.A. (1996) *Trends Biochem. Sci.*, **21**, 444-445.
- Bairoch,A. and Apweiler,R. (1997) *Nucleic Acids Res.*, **25**, 31-36 [see also this issue (1998) *Nucleic Acids Res.* **26**, 38-42].
- Birney,E. and Durbin,R. (1997) In *ISMB-97; Proceedings Fifth International Conference on Intelligent Systems for Molecular Biology*. AAAI Press, Menlo Park, pp. 56-64.