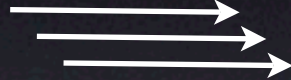# Data Integration using the Distributed Annotation System (DAS)

Andreas Prlić, Ewan Birney, Tony Cox, Thomas A. Down, Rob Finn, Stefan Gräf, David Jackson, Andreas Kähäri, Eugene Kulesha, Roger Pettett, James Smith, Jim Stalker, Tim J. P. Hubbard

- what is DAS

- what do we do with it

- DAS registration server

- latest developments

Integration of personal data into bioinf. resources

- Integration of annotations from external sources into local applications

- online access to most recent data versions
  - no need for local installations

# DAS, how it works

Dowell, Jokerst, Allen, Eddy, Stein
BMC Bioninformatics 2001

# a few principles...

- Clients are "intelligent" (few)
- Servers are simple and easy to set up (many)
- (most of) data is precalculated
- libraries for server and client
- multiple programming languages

# http://www.ensembl.org

**e! Ensembl**

Ensembl v39 - Jun 2006

## Use Ensembl to...

- Run a BLAST search
- Search Ensembl
- Data mining [BioMart]
- Export data
- Download data

## Docs and downloads

- Information
- What's New
- About Ensembl
- Ensembl data
- Software

## Other links

- Home
- Sitemap
- Vega
- *Pre!* Pre Ensembl
- *e!* View previous release of page in Archive!
- *e!* Stable Archive! link for this page
- *e!* Archive! sites
- Trace server

**Sanger** **EBI**

*Rattus norvegicus* 3.4 update

## What's New in Ensembl 39

- ▸ **Honeybee dropped from Ensembl** (*Apis mellifera*)
- ▸ **New Opossum assembly and genebuild** (*Monodelphis domestica*)
- ▸ **New Mouse assembly and genebuild** (*Mus musculus*)
- ▸ **New Ciona savignyi assembly and genebuild** (*Ciona savignyi*)
- ▸ **New Zebrafish assembly and genebuild** (*Danio rerio*)

More news...

## About Ensembl

Ensembl is a joint project between EMBL - EBI and the Sanger Institute to develop a software system which produces and maintains automatic annotation on selected eukaryotic genomes. Ensembl is primarily funded by the Wellcome Trust.

This site provides free access to all the data and software from the Ensembl project. Click on a species name to browse the data.
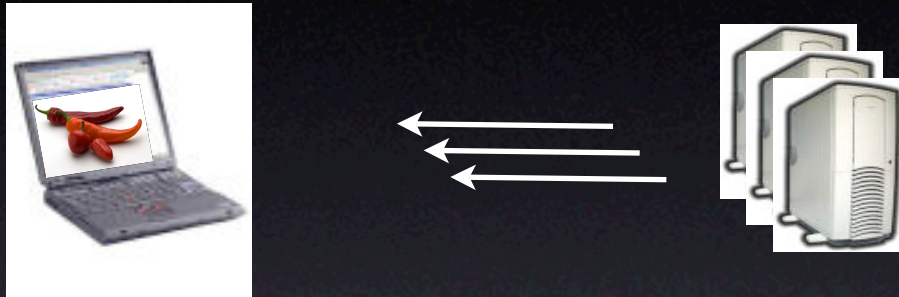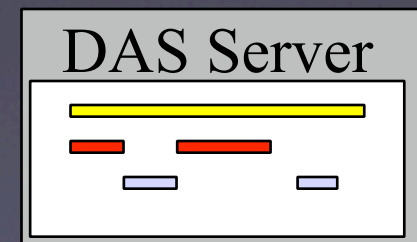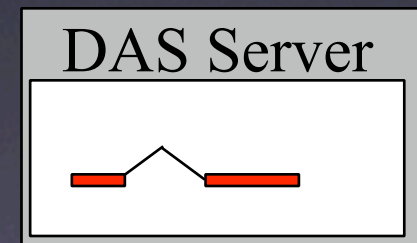
Access to all the data produced by the project, and to the software used to analyse and present it, is provided free and without constraints. Some data and software may be subject to third-party constraints.

For all enquiries, please contact the Ensembl HelpDesk (helpdesk@ensembl.org).

## Other sites using the Ensembl system

- ▸ EBI Genome Reviews database

## Mammalian genomes

*Homo sapiens* — NCBI 36 | Vega

*Pan troglodytes* — PanTro 1.0

*Macaca mulatta* — MMUL 0.1 | *pre!*

*Mus musculus* — **UPDATED!** NCBI m36 | Vega

*Rattus norvegicus* — RGSC 3.4 | *pre!*

*Oryctolagus cuniculus* — **Pre!** RABBIT

*Canis familiaris* — CanFam 1.0 | Vega | *pre!*

*Bos taurus* — Btau 2.0

*Dasypus novemcinctus* — **Pre!** ARMA

*Loxodonta africana* — **Pre!** BROAD E1

*Echinops telfairi* — **Pre!** TENREC

*Monodelphis domestica* — **UPDATED!** MonDom 4

## Other species

*Gallus gallus* — WASHUC 1

*Xenopus tropicalis* — JGI 4.1

*Danio rerio* — **UPDATED!** Zv 6 | Vega

*Takifugu rubripes* — FUGU 4.0

*Tetraodon nigroviridis* — TETRAODON 7

*Gasterosteus aculeatus* — **Pre!** BROAD S1

*Oryzias latipes* — **Pre!** **NEW!** MEDAKA 1

*Ciona intestinalis* — JGI2

*Ciona savignyi* — **UPDATED!** CSAV 2.0

*Drosophila melanogaster* — BDGP 4

*Anopheles gambiae* — AgamP3

*Aedes aegypti* — **Pre!** AEDES 1

*Caenorhabditis elegans* — WS 150

*Saccharomyces cerevisiae* — SGD 1

- > 20 vertebrates / model organism

- 5 mill. page impressions / week
- 100  mirrors/internal installations worldwide
- open source
- used for other species as well
- MySQL
- 5-10 G / species + 100 G multi species data

Add your own uses Registry

Now: Mostly Cloudy, 11° C    Thu: 11° C

Ensembl v37: Homo sapiens Pep...

- Transcript information
- Exon information
- Protein information
- Export protein data

**Chromosome 17**
**74,183,768 - 74,289,900**

- View of Chromosome 17
- Graphical view
- Graphical overview
- Export information about region
- Export sequence as FASTA
- Export EMBL file
- Export Gene info in region
- Export SNP info in region
- Export Vega info in region

**Use Ensembl to...**

- Run a BLAST search
- Search Ensembl
- Data mining [BioMart]
- Upload your own data
- Export data
- Download data

**Docs and downloads**

- Information
- What's New
- About Ensembl
- Ensembl data
- Software

# Linking protein structure to e! Peptide view

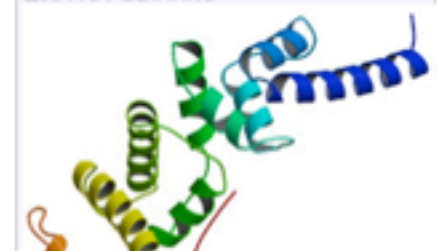| | |
|---|---|
| **Prediction Method** | Genes were annotated by the Ensembl automatic analysis pipeline using either a GeneWise/Exonerate model from a database protein or a set of aligned cDNAs followed by an ORF prediction. GeneWise/Exonerate models are further combined with available aligned cDNAs to annotate UTRs (For more information see V.Curwen et al. , Genome Res. 2004 14:942-50.) |
| **InterPro** | IPR001472 Bipartite nuclear localization signal - [View other genes with this domain]<br>IPR000904 SEC7-like - [View other genes with this domain]<br>IPR001849 Pleckstrin-like - [View other genes with this domain] |
| **Protein Family** | ENSF00000001251 : CYTOHESIN<br>This cluster contains 4 Ensembl gene member(s) |
| **Protein Features** | |
| **Protein Sequence** | |

Peptide
Coiled coils
Pfam          Sec7          PH
NLS_BP        PH
Prosite profiles
Sec7
DAS PDB_Spice
ENSP - PDB ma..
Scale (aa)      0      40      80      120      1
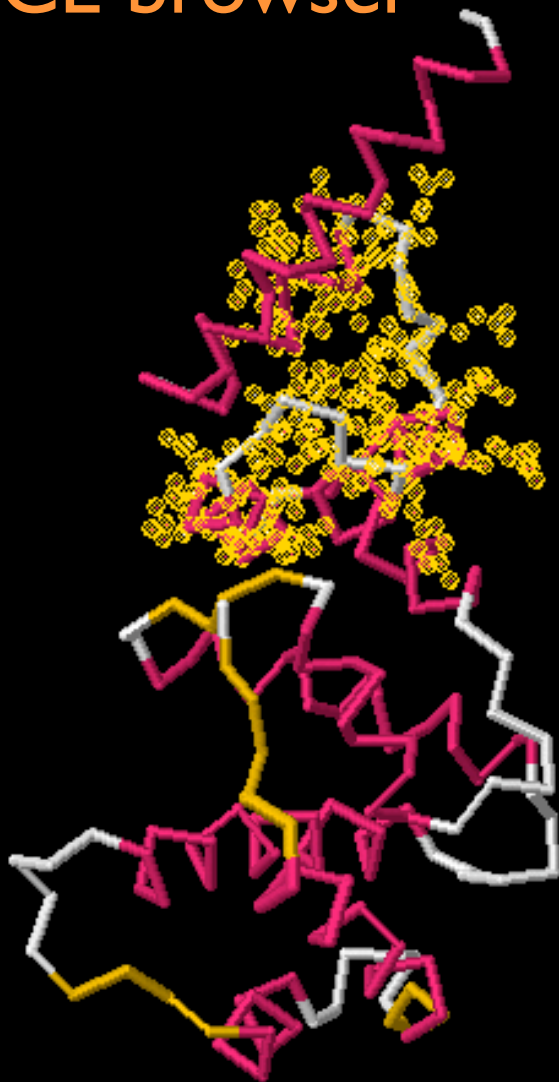
PDB: 1pbv, via UniProt Q99418 PDB: ARNO    ✕
ID: ENSP - PDB mapping to 1pbv.
TYPE: ENSP - PDB mapping
METHOD: Compara, MSD-Mapping of UniProt to PDB
LINK:
http://das.sanger.ac.uk/registry/showspice.jsp?pdb=1pbv.
NOTE: PDB: 1pbv. 85%id 49%coverage via UniProt
Q99418 PDB:ARNO

MEEDGSYVPSDLTWEIDQILENIRRRKQEILADIQRLKDEIAEVANEIIHLGSTEEPKHM
QRNKQVAMGRKKINPDPKKGIQFLIENDLLKNTCEDIAQFLYKGEGLNKTAIGDYLGEPD
ETNIQVLHAIVELHLEITTDLNLVQALRQILWSFRLPGIAQKIDRMMEAIAQRYCQCNNGVI
QSTVTCYVLSFAIIMLNTSLDHPNVKIKPTVERIANHGIHDGGDLPIILLRHLYESIK
HEPFKIPIDDGHDLTHTFFHPDREGWLLKLGGRVKTVKRDSYILTHNCLYYFIYTTDKIP
RGIIPLENLSIREVEDSKKPHCFELYIPDHKDQVIKACKTEADGRVVIGHHFTGSQLPDD
IKEEWIKCIKAAISRDPFYEMEAARKKKVSSTKRM

Exon alternating       Residue overlap
text colour            splice site

Find:    Find Next    Find Previous    Highlight all    Match case

http://www.efamily.org.uk/software/dasclients/spice

See exon structure mapped onto 3D

PDB                                          1A17.
                                                    159

dssp
SECSTRUC
BEND

cath
Cath Domain

s3dm
ALPHA-BET

UniProt                              ∧∨        P53041

uniprot

description
CHAIN
SECSTRUC
REPEAT
REGION
METAL
ACT_SITE
CONFLICT
Smart

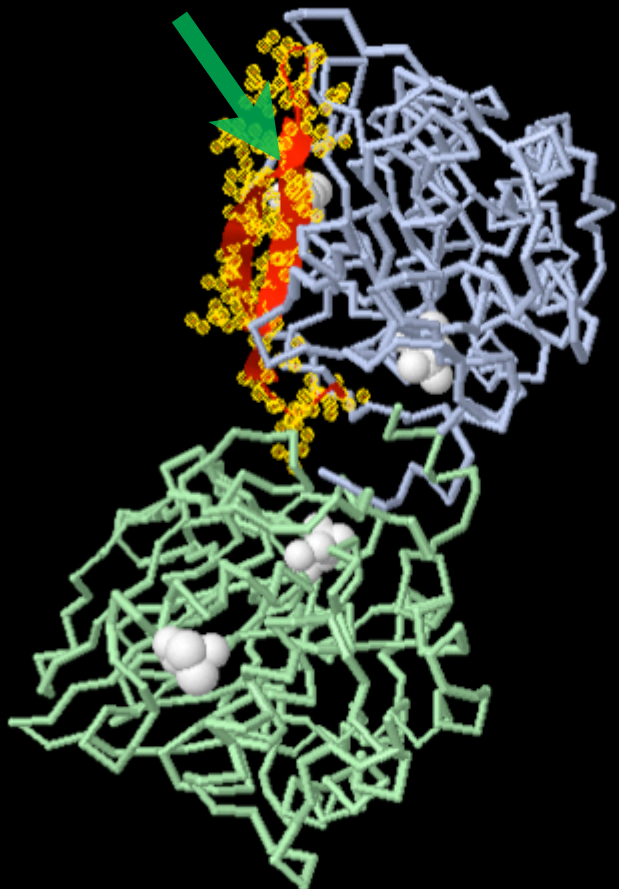ENSP                                 ∧∨    ENSP00000012443

ENSP - PDB

MisPred

hsa35pep

exon
exon
sSNP
nsSNP

Click

enter RASMOL like command...

Show SNPs

interact with Menu & RASMOL

RASMOL commands

Zoom

DAS commands

Structure

Features

Alignment

Sequence

auto install
latest version
send arguments

Java Web Start

DAS registry

SPICE

Meta information
about DAS servers

DAS registry

SPICE

# The DAS registration server

| | | | | | |
|---|---|---|---|---|---|
| DS_110 | | dssp | cmbi4.cmbi.ru.nl/das/dssp/ | features | PDBresnum,Protein Structure |
| DS_111 | | cath | cathwww.biochem ... 0/das/cath_pdb/ | features | PDBresnum,Protein Structure |
| DS_112 | | structure | das.sanger.ac.uk/das/structure/ | structure | PDBresnum,Protein Structure |
| DS_113 | | alig_pdb_sp | das.sanger.ac.uk/das/msdpdbsp/ | alignment | UniProt,Protein Sequence<br>PDBresnum,Protein Structure |
| DS_114 | | signalp | genome.cbs.dtu.dk:9000/das/signalp/ | types<br>features | UniProt,Protein Sequence |
| DS_115 | | netphos | genome.cbs.dtu.dk:9000/das/netphos/ | types<br>features | UniProt,Protein Sequence |
| DS_116 | | netoglyc | genome.cbs.dtu.dk:9000/das/netoglyc/ | types<br>features | UniProt,Protein Sequence |
| DS_117 | | tmhmm | genome.cbs.dtu.dk:9000 | | |

http://das.sanger.ac.uk/registry/

# DAS registration server

- allows to "publish" DAS servers & share with community

- communicates with clients

- regularly checks servers, sends notification
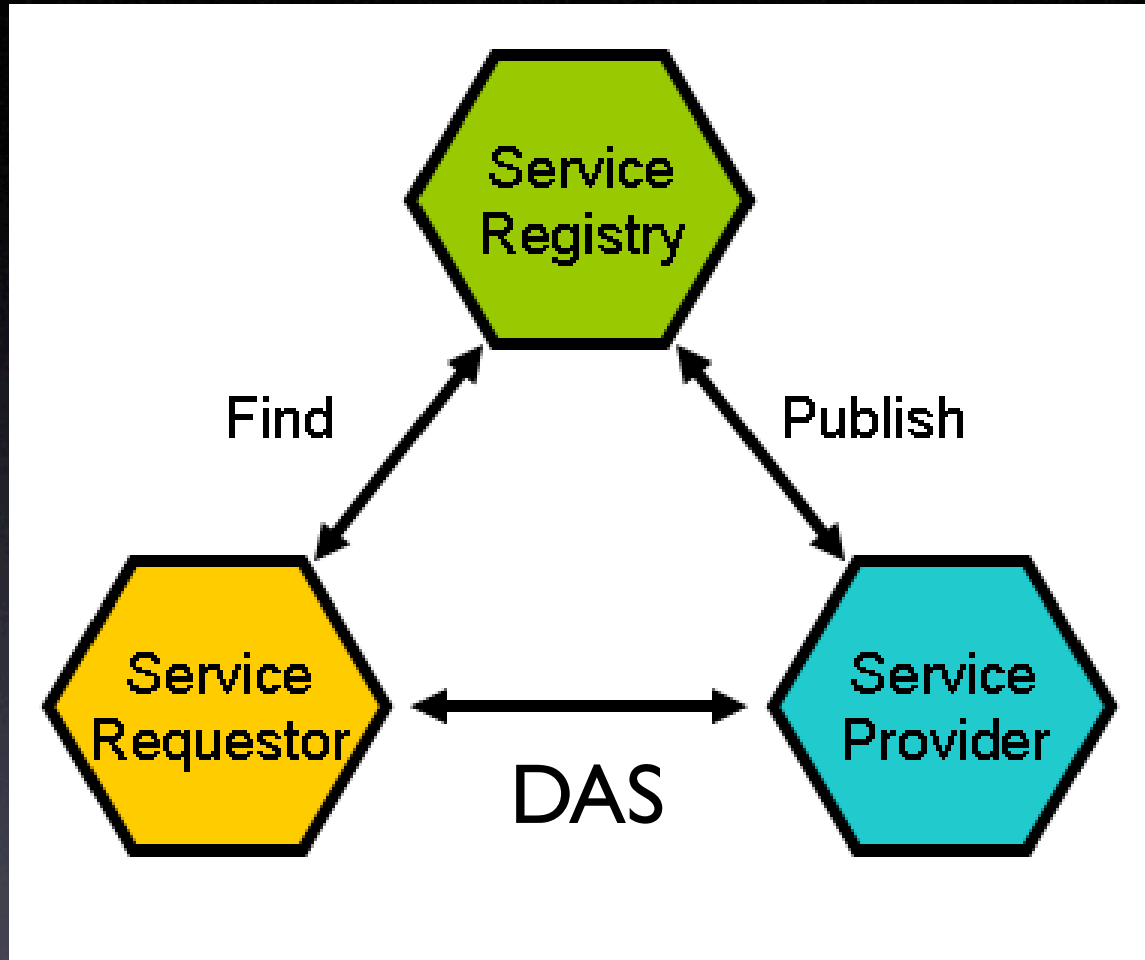
# What is the glue?

- "Coordinate Systems"
  - Authority
  - Type of data
  - Version
  - Organism (optional)

# Clients and Coordinate Systems

- Ensembl - most of the views can display DAS sources from multiple CS

- SPICE -  PDB, UniProt, Ensp

- Dasty - UniProt

DAS registration server



e.g. Ensembl, SPICE

a DAS source

the DAS - SOA

registered DAS servers over time

111 DAS sources
26 institutions
12 countries

eFamily

BIOSAPIENS NETWORK

+ others

# DAS - issues

- inconsistent implementations

- no consistent use annotation types

- error handling

- searches not possible - in DAS/1

- open sharing of data - low security

# http://sisyphus.mrc-cpe.cam.ac.uk

## SISYPHUS

Submit Query

### Structural alignments for proteins with non-trivial relationships

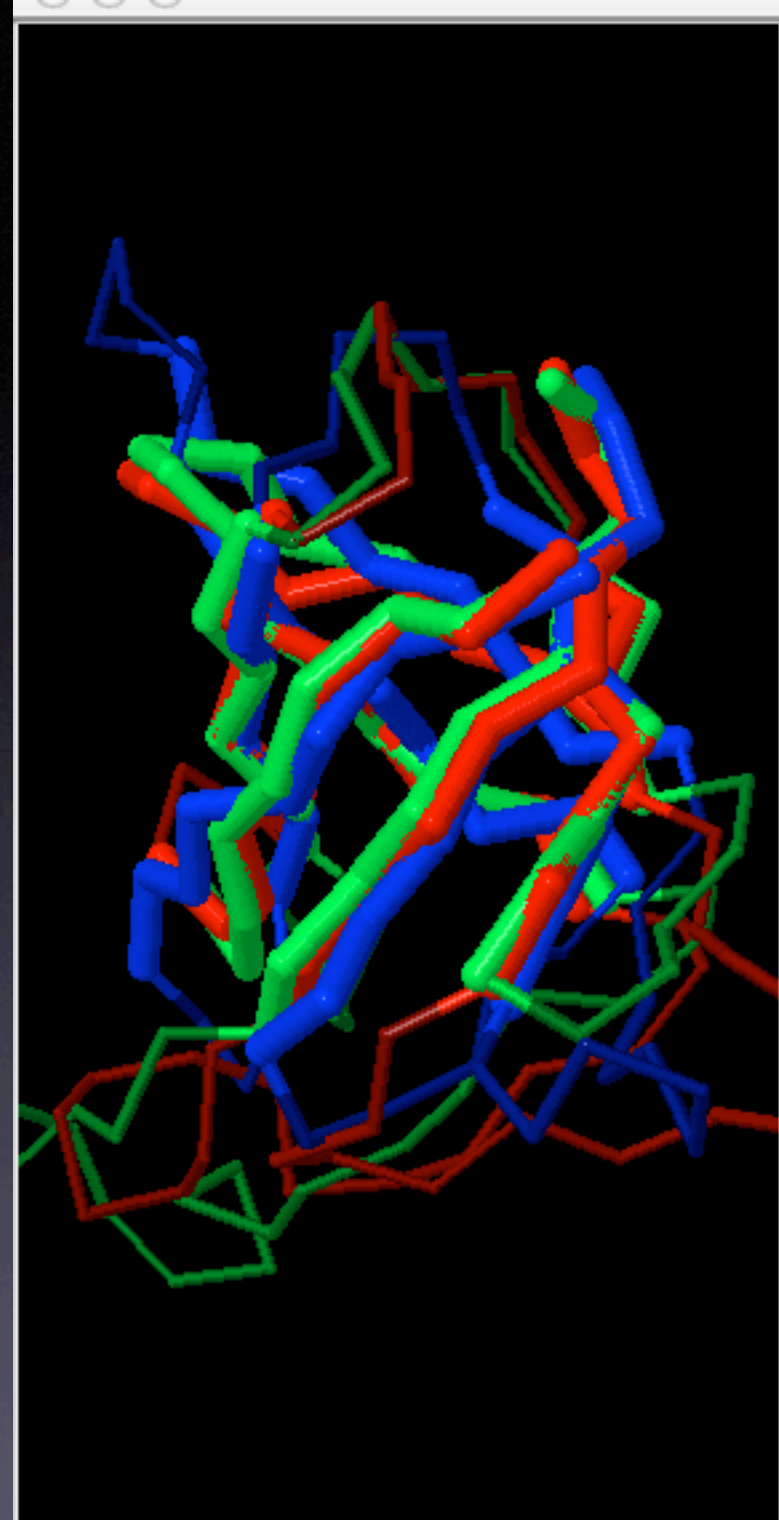*Sisyphus: in greek: crafty*

### Authors:

Antonina Andreeva, Andreas Prlic, Tim Hubbard, Alexey Murzin

```
SATVGIFVDagsraenvknNGTAHFLEHLAFkgtqnr-----------pqQGIELeienig---shLNAYTsr----eNTVYYAKSLq--edIPKAVDI
TTHIAIALEgvswsap--dYFVALATQAIVGnwdraigtgtn------spSPLAVaasqngslansYMSFStsyadsgLWGMYIVTDsnehnVRLIVNE
FSALGLYIDagsrfegrnlKGCTHILDRLAFkstehv----------egRAMAEtlellg---gnYQCTSsr----eNLMYQASVFn--qdVGKMLQL
LFHIQIGFEglpidhp--dIYALATLQTLLGgggsfsaggpgkgmysrlyTHVLNqy----yfvenCVAFNhsysdsgIFGISLSCIp--qaAPQAVEV
TCTVGVWIDagsryeseknNGAGYFVEHLAFkgtknr----------pgNALEKevesmg---ahLNAYStr----eHTAYYIKALs--kdLPKAVEL
LAHVAIAVEgpgwahp--dNVALQVANAIIGhydctygggah------lsSPLASiaatn-klcqsFQTFNicyadtgLLGAHFVCDhms--IDDMMFV
TCTVGVWIDagsryeseknNGAGYFLEHLAFkgtknr----------pgNALEKevesmg---ahLNAYSsr----eHTAYYIKALs--kdVPKAVEL
LAHVAIAVEgpgwahp--dLVALQVANAIIGhydrtygglh------ssSPLASiavtn-klcqsFQTFSicysetgLFGFYFVCDrms--IDDMMFV
TASVGVVFGsgaanenpynNGVSNLWKNIFL-----------------sKENSAvaakeg---laLSSNIsr----dFQSYIVSSLp--gsTDKSLDF
KAWISLAVEgepvnsp--nYFVAKLAAQIFGsynafepasrl------qgIKLLDniqey-qlcdnFNHFSlsykdsgLWGFSTATRnvt-mIDDLIHF
ISTLAVKVHggsryat--kDGVAHLLNRFNFqntntr----------saLKLVResellg---gtFKSTLdr----eYITLKATFLk--ddLPYYVNA
DSVAAIGIPvn---ka--sLAQYEVLANYLTsal---------------SELSGli---------SSAKLdkftdggLFTLFVRDQ-dsavVSSNIKK
LVHAAIVAEsaaigga--eANAFSVLQHVLG-----------------ANPHVkrgnp----fdVSAFNasysdsgLFGFYTISQaay--AGQVIKA
GSTIGVFIKagsryenssnLGTSHLLRLASSlttkga-----------ssFKITRgieavg---gkLSVEStr----eNMAYTVECLr--ddVEILMEF
ASRIGLFIKagsryensnnLGTSHLLRLASSlttkga-----------ssFKITRgieavg---gkLSVEStr----eNMAYTVECLr--ddVDILMEF
LVHAALVAEsaaigsa--eANAFSVLQHVLGagphvkrgsna------tSSLYQavakgvhgpfov----asysds--FYTISQa--GDVIKA
```

eFamily

| | |
|---|---|
| ☑ | 1efc.A_A301:310_A326 |
| ☐ | 1b23.P_P313:322_P338 |
| ☐ | 1exm.A_A313:322_A33 |
| ☐ | 1d2e.A_A349:358_A37 |
| ☐ | 1f60.A_A335:344_A35 |
| ☐ | 1skq.A_A323:332_A34 |
| ☐ | 1kk0.A_A322:331_A35 |
| ☑ | 1s0u.A_A349:358_A37 |
| ☐ | 1r5b.A_A557:566_A57 |
| ☐ | 1pj5.A_A744:753_A760 |
| ☐ | 1nrk.A_A246:255_A26 |
| ☐ | 1wos.A_A279:288_A29 |
| ☑ | 1vlo.A_A279:288_A294 |

**PDB**    1EFC.A

- dssp
  - SECSTRUC
  - BEND
  - 3HELIX
  - BRIDGE
- cath
- s3dm

**UniProt**    P02990

- uniprot
  - description
  - INIT_MET
  - CHAIN
  - MOD_RES
  - SECSTRUC
  - NP_BIND
  - MUTAGEN
  - VARIANT
  - TIGRFAMs
  - TIGRFAMs

**ENSP**    0

- superfam
- cbs_sort
- cbs_ptm
- cbs_func

enter RASMOL like command...

- Alignment DAS:

- rotation matrices, shift vectors

- range information (optional)

http://www.jalview.org  A. Waterhouse, J. Procter, G. Barton

# Acknowledgments

- T. Down, T. Hubbard

- Web Team, E. Kulesha, R. Pettett,  T. Cox

- eFamily Project

- S. Gräf,  A. Kahari, BioSapiens

- A. Murzin,  A. Andreeva

- R. Finn, H.Hotz,  A.Ahmed

- Jmol, Biojava, MSD, everybody who sets up DAS servers