# I - Summary

- <u>Introduction to protein domains</u>
- Domain databases

http://www.sanger.ac.uk/Software/Pfam/

# Protein Domains

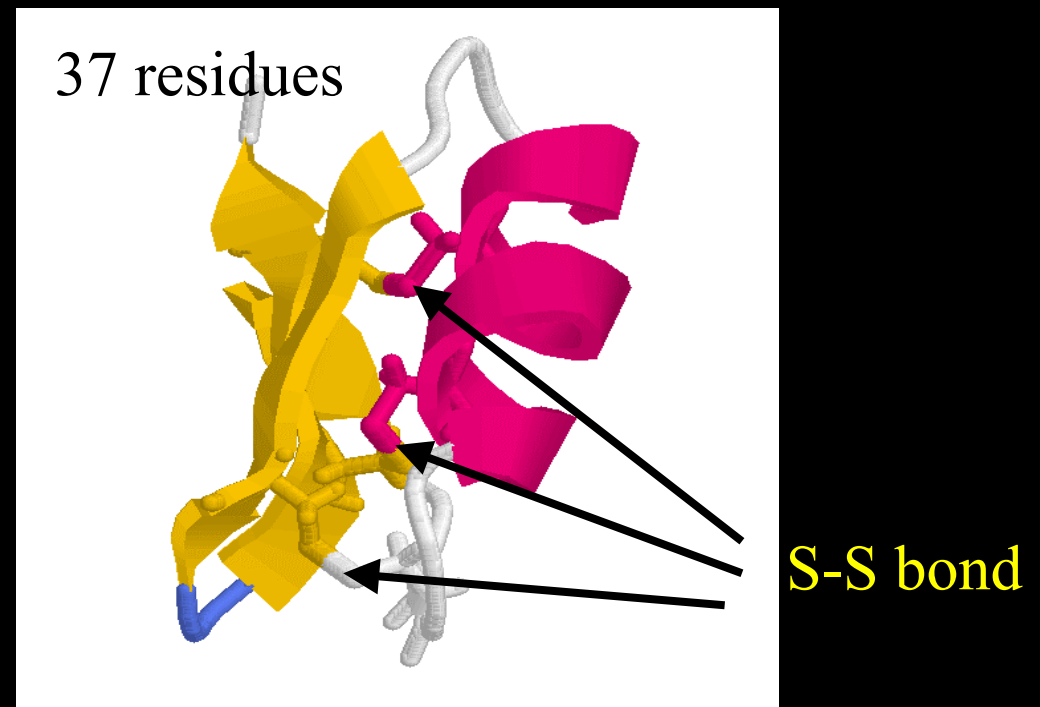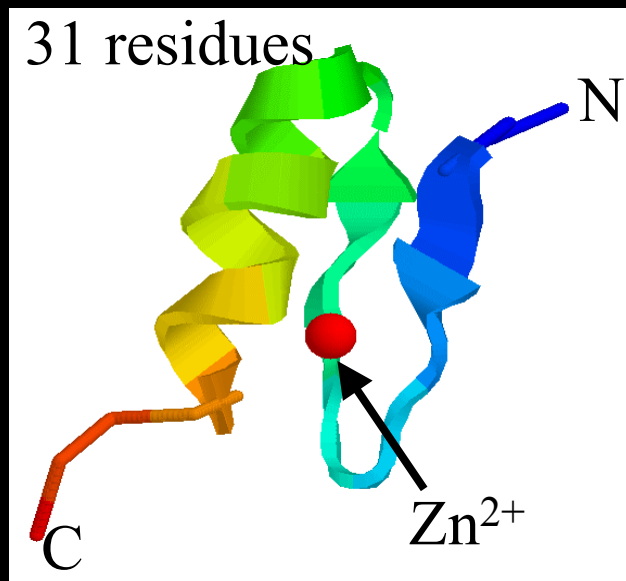- From a structural perspective protein domains are discrete units.

# What is a Domain?



- Defined by structure
- Domain boundaries can be inferred from careful sequence analysis
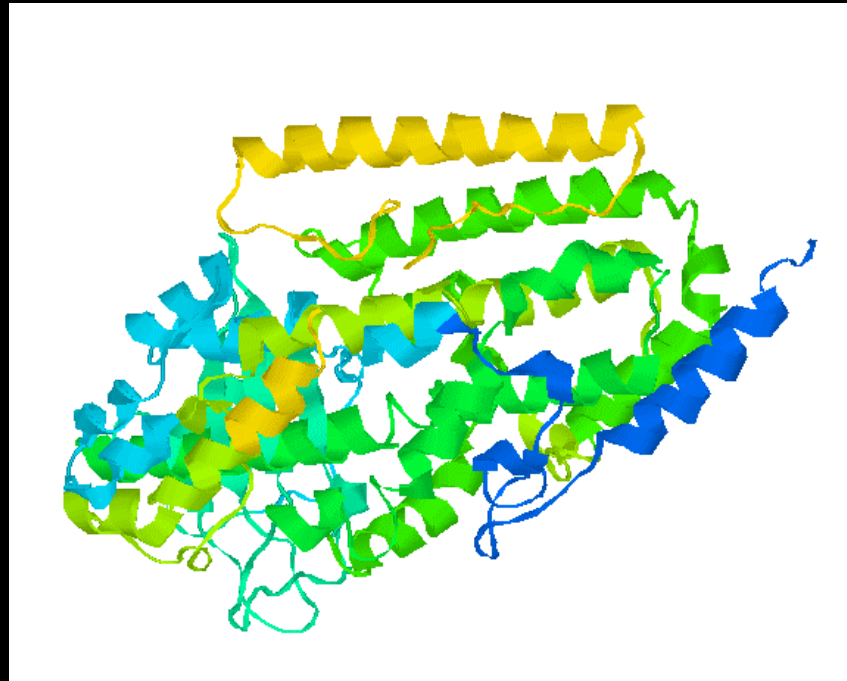- Domains are the currency of protein function

# Domains - size

- Domains can be 25 to 500 residues long

- Most are less than 200 residues.

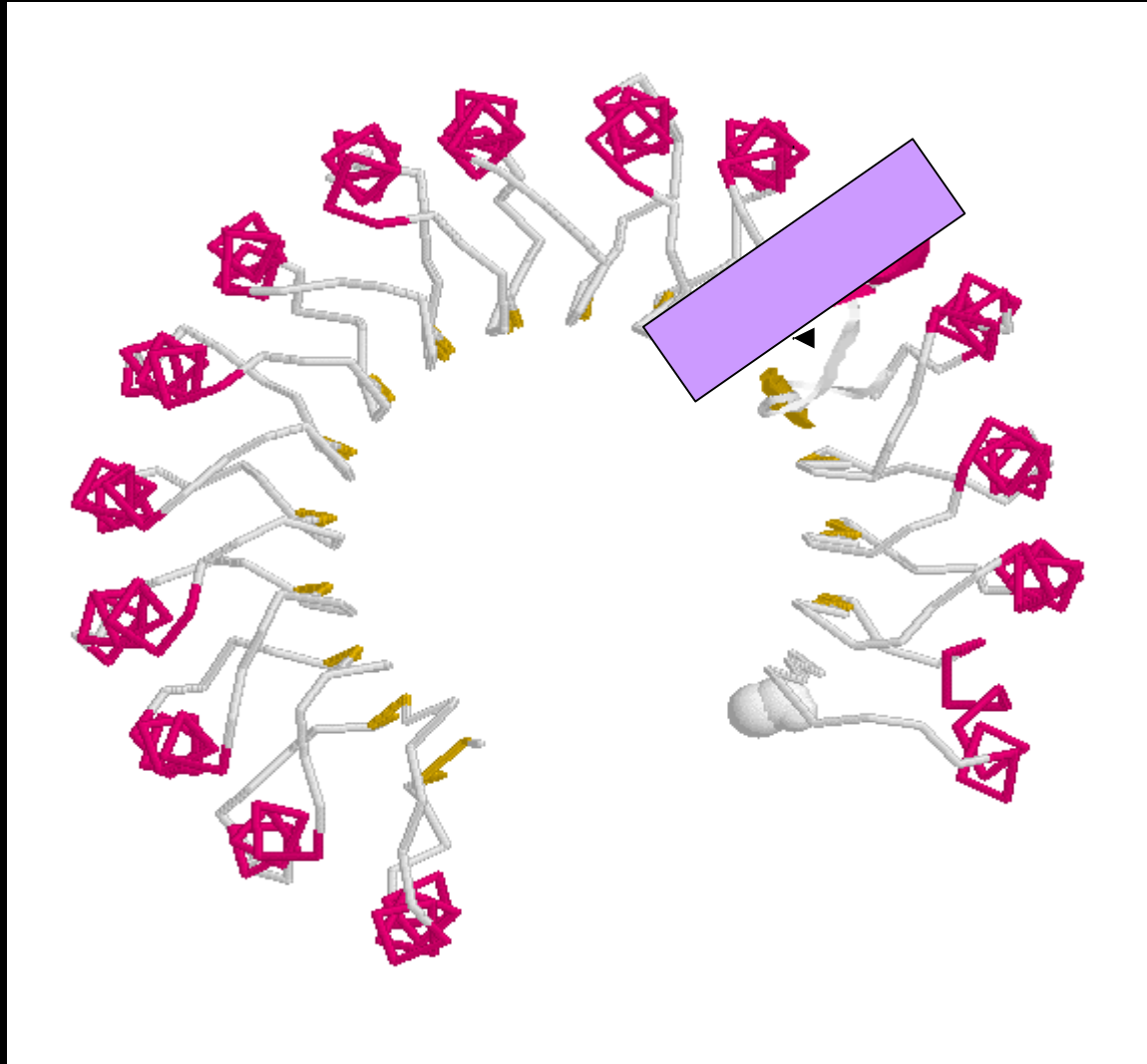- Domains can be smaller that 50 residues but these are stabilized by disulphide bonds or chelated metals.



31 residues

N

C

Zn$^{2+}$

37 residues

S-S bond

http://www.sanger.ac.uk/Software/Pfam/

# Example domains

- The lipoxygenase domain is a giant at 500 residues long.



http://www.sanger.ac.uk/Software/Pfam/

# Leucine Rich repeats



- A single repeat is not stable
- Multiple repeats are stable
- Each repeat is represented separately
- Unlimited number

http://www.sanger.ac.uk/Software/Pfam/

# WD40 repeats



- 7 repeats
- beta sheet per repeat
- Limited number (6-8)

# Structural domains

- Domains are most easily defined in known structures

- Several automatic programs available

- They don't always/often agree!

http://www.sanger.ac.uk/Software/Pfam/

# Defining domains from sequence

- Has been done successfully hundreds of times

- Cannot always be done

- Usually requires the domain to be mobile

http://www.sanger.ac.uk/Software/Pfam/

# Domains and structure determination

- Hard to get structure of complete protein

- Expressing smaller segments is easier

# Domain Hunting: CBS domains

- **Discovering new domains can reveal new biology**

CBS domains form energy-sensing modules whose binding of adenosine ligands is disrupted by disease mutations

See the related Commentary beginning on page 182.

John W. Scott,[1] Simon A. Hawley,[1] Kevin A. Green,[1] Miliea Anis,[1] Greg Stewart,[1] Gillian A. Scullion,[1] David G. Norman,[2] and D. Grahame Hardie[1]

[1]Division of Molecular Physiology, and
[2]Division of Biological Chemistry and Molecular Microbiology, Faculty of Life Sciences, Wellcome Trust Biocentre, University of Dundee, Dundee, Scotland, United Kingdom

J. Clin. Inv. 113:274-284.

http://www.sanger.ac.uk/Software/Pfam/

# Domain Hunting: RNAi



Cerrutti, Mian & Bateman. Trends Biochem Sci. 25:481-482  (2000)

http://www.sanger.ac.uk/Software/Pfam/

# I - Summary

- Introduction to protein domains
- <u>Domain databases</u>

http://www.sanger.ac.uk/Software/Pfam/

# Domain databases

- Many of the common domains have already been defined in domain databases.

- Advantages:
  - Pre-annotated domains
  - Easy interpretation of domain structure
  - Sensitivity can be higher

- The most used databases are:
  - Pfam                           - Prints
  - Prosite Profiles               - Blocks
  - SMART                          - ProDom

  http://www.sanger.ac.uk/Software/Pfam/

- Good coverage
- No specific bias
- Good graphical views
- Structural data in alignments
- No heirarchy

http://www.sanger.ac.uk/Software/Pfam/

- Domain collection by Ponting and Bork.
- Specialises in
  - Signaling domains
  - Extracellular domains
  - Nuclear domains
- Excellent quality families.
- Really nice graphics
- Coiled-coil, TM, low-complexity

http://www.sanger.ac.uk/Software/Pfam/

- Profiles
  - Sensitive
  - Low coverage (Good for signalling)
- Patterns
  - e.g.   N-{P}-[ST]-{P}
  - less sensitive
  - many false positives

http://www.sanger.ac.uk/Software/Pfam/

# Interpro

- Interpro is a database that presents Prosite, Prints, Prodom and Pfam domain.

- Annotation is a strong point



http://www.sanger.ac.uk/Software/Pfam/

# Conclusions

- Domains are the common currency of protein function

- Understanding the domain structure helps to understand the biology

- Domain databases are key labour saving tools

http://www.sanger.ac.uk/Software/Pfam/

# II - Summary

- <u>Introduction to Pfam</u>
- Protein Interactions
- Pfam Clans

http://www.sanger.ac.uk/Software/Pfam/

# Pfam: 8,000 families for the molecular biologist



Alex Bateman, Richard Durbin, Sean Eddy, Ajay Khanna, Rob Finn, Sam Griffiths-Jones, Jaina Mistry, John Tate, Volker Hollich and Erik Sonnhammer.

http://www.sanger.ac.uk/Software/Pfam/

# Annotating genomes



http://www.sanger.ac.uk/Software/Pfam/

# Family Pages



http://www.sanger.ac.uk/Software/Pfam/

# Family Pages



http://www.sanger.ac.uk/Software/Pfam/

# Pfam contains Alignments



http://www.sanger.ac.uk/Software/Pfam/

# The data deluge



http://www.sanger.ac.uk/Software/Pfam/

# Pfam contains:

SEED alignment
representative members

→ Profile-HMM
HMMer-2.0

↓ Search database

FULL alignment

Manually curated

Automatically made

http://www.sanger.ac.uk/Software/Pfam/

# Profiles, HMMs and PSSMs

- Complicated names - Simple idea

```
RU1A_HUMAN rrm1    SSATNAL
RU1A_HUMAN rrm2    VQAGAAR
SFR1_HUMAN rrm1    RDAEDAV
SXLF_DROME rrm1    MDSQRAI
                   |||||||
PABP_DROME rrm3    EAAEAAV
```

http://www.sanger.ac.uk/Software/Pfam/

# Pfam 18.0

- Pfam-A
  - 7,973 Curated families with annotation.

- Pfam-B
  - 100,000 families derived from Prodom.

A

http://www.sanger.ac.uk/Software/Pfam/

# Coverage
## Pfam Sequence Coverage



- Retire sometime between Sept 2012 and May 2033!

  http://www.sanger.ac.uk/Software/Pfam/

# Family Pages



http://www.sanger.ac.uk/Software/Pfam/

http://www.sanger.ac.uk/Software/Pfam/

- # Taxonomy information
  - ## Does your favourite thermophile have a member?



http://www.sanger.ac.uk/Software/Pfam/

http://www.sanger.ac.uk/Software/Pfam/

# II - Summary

- Introduction to Pfam
- Protein Interactions
- Pfam Clans

http://www.sanger.ac.uk/Software/Pfam/

translation factors

NMD factors

5'-to-3' mRNA decay factors

3'-to-5' mRNA decay factors

EJC factors

http://www.sanger.ac.uk/Software/Pfam/

From Izaurralde group home page

# Protein Interactions



SH2

SH3

# iPfam

# Protein Interactions



http://www.sanger.ac.uk/Software/Pfam/

# The Pfam Supercomputer



A typical computer console.

http://www.sanger.ac.uk/Software/Pfam/

# Complex Complexes



ATP synthase

Cytochrome bc1

# II - Summary

- Introduction to Pfam
- Protein Interactions
- <u>Pfam Clans</u>

http://www.sanger.ac.uk/Software/Pfam/

# What is Pfam?

- Database of Protein families

- Pfam defined by alignments & HMMs



- Flat classification

- Many families are related

http://www.sanger.ac.uk/Software/Pfam/

# Protein Space



Rossmann

Serine
Protease

Globins

Apartic
Protease

Immuno-
globulins

http://www.zbh.uni-hamburg.de/wurst/protspace

http://www.sanger.ac.uk/Software/Pfam/

# Protein Space



How Can We Identify
Related families in Proteins
Space ?

http://www.sanger.ac.uk/Software/Pfam/

# Pfam Clans

- Group together families
  - Structure databases
  - PRC
  - Overlaps
  - Literature

http://www.sanger.ac.uk/Software/Pfam/

# Structural Databases

Profile HMM Comparison

http://www.sanger.ac.uk/Software/Pfam/

"I think you should be more explicit here in step two."

http://www.sanger.ac.uk/Software/Pfam/

http://www.sanger.ac.uk/Software/Pfam/

# Browse clans



http://www.sanger.ac.uk/Software/Pfam/

# Clan relationships



http://www.sanger.ac.uk/Software/Pfam/

# HMM comparisons



http://www.sanger.ac.uk/Software/Pfam/

# Clan alignments



http://www.sanger.ac.uk/Software/Pfam/

# So How Are We Doing ?

- 172 Clans
  - Contains 1181 Pfams (15%)
  - Largest Clan is NADP_Rossmann
    - 53 families
    - Covers 4 SCOP sf
    - Added over 5000 domain hits to Pfam
  - Largest family without a structure is MFS
    - 19 families
    - Also added over 5000 domain hits to Pfam
  - 66% have a structure representative.
  - many families w/o structure can now be related to a structure

http://www.sanger.ac.uk/Software/Pfam/

# Conclusions

- Majority of proteins have Pfam domains

- Pfam helps to understand protein interactions

- Pfam clans give heirarchy

http://www.sanger.ac.uk/Software/Pfam/