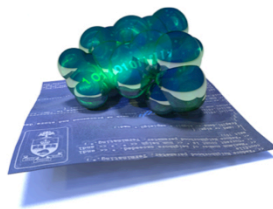A
BIOINFORMATICS
COURSE

# PROTEIN STRUCTURE DOMAINS

BORIS STEIPE

DEPARTMENT OF BIOCHEMISTRY – DEPARTMENT OF MOLECULAR GENETICS
UNIVERSITY OF TORONTO
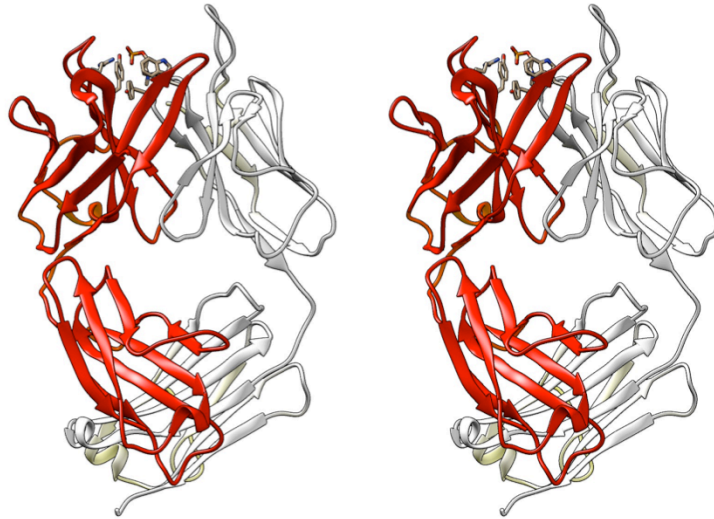
# Large proteins are usually composed of compact, semi-independent modules.

Reasons:
Folding efficiency
Modularity


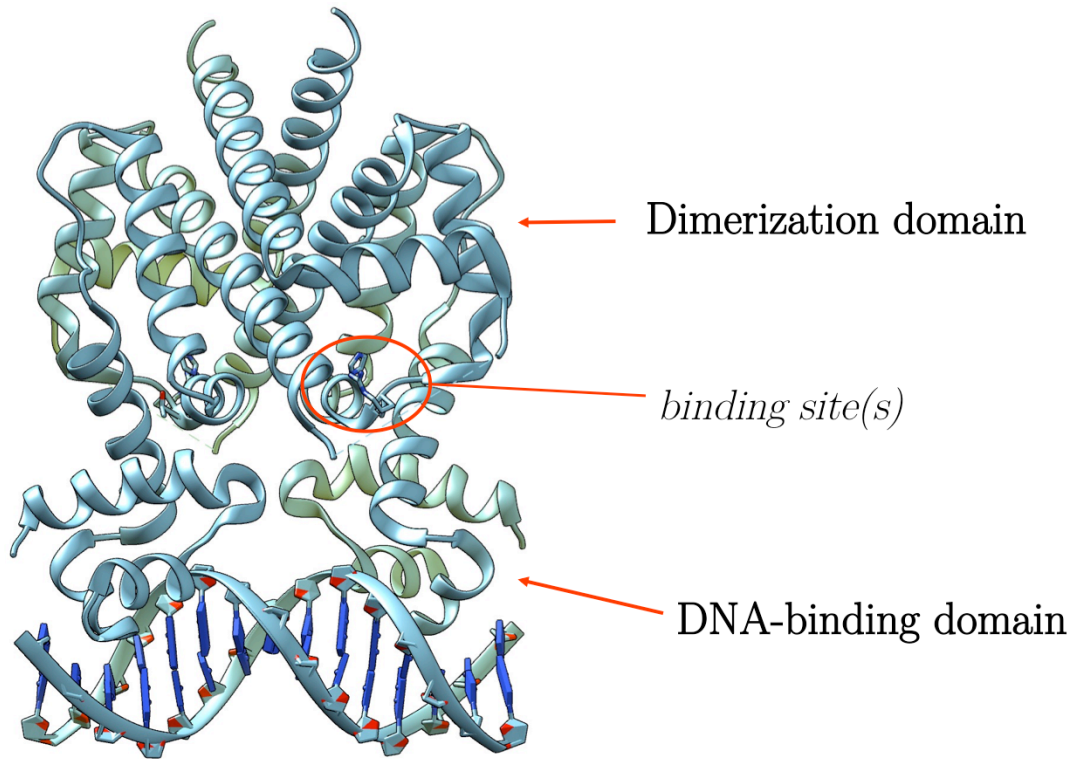
2MCP
(Immunoglobulin Fab fragment with bound hapten - phosphocholine)

Domains are ubiquitous in proteins and – although the idea of a domain is purely conceptual – they supply an intriguing link between sequence and structure in evolution.
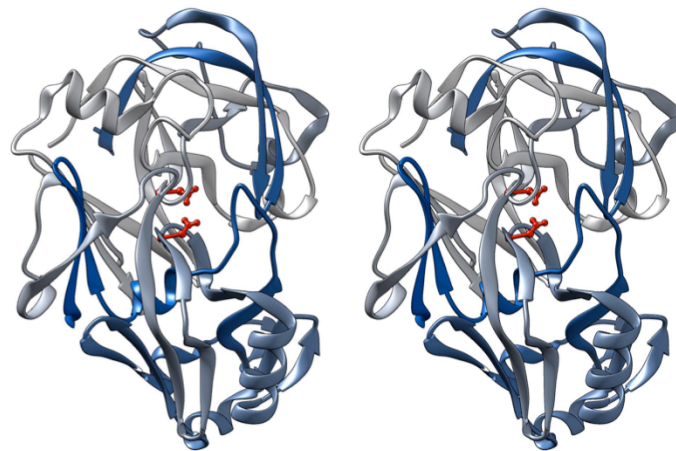
Dimerization domain

*binding site(s)*

DNA-binding domain

1QPI (Tet-repressor-Operator complex)

Many domains are empirically found to have distinct functions in proteins. The example above is the Tet-repressor. This observation is highly non-trivial and reflects on the nature of the process of evolution. To explain this fact requires to consider what selective *advantages* can be gained from compartmentalizing function in domains, rather than distributing it over the entirety of the structure, **or** understanding how this fact can be a consequence of the *process* of evolution, or of constraints that arise from evolutionary *mechanisms.*
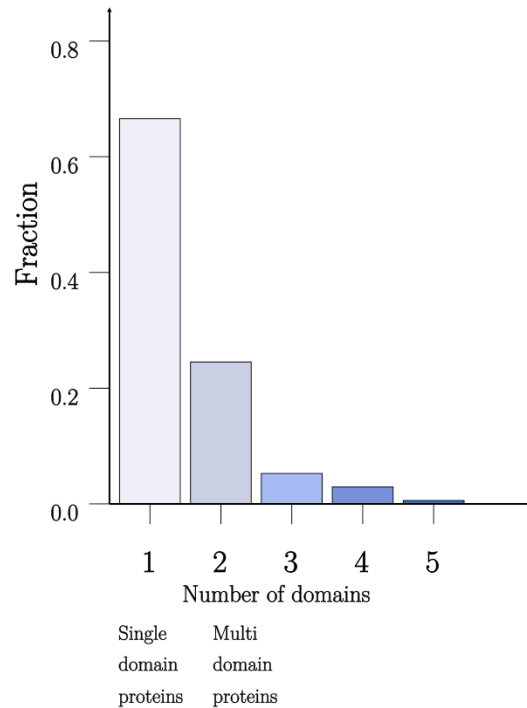
(Lysosomal Aspartic Protease, 1LYB)

But it is not always the case that distinct domains correspond to discrete functions: in this example, the active site of the aspartic protease cathepsin D is shared between two domains – in fact the relative motion of the domains appears to be important for the catalytic mechanism.

Number of domains in 787 representative proteins used as the basis for the CATH database of protein structure architectures:

(However CATH proteins represent a biased sample: in general, large multi-domain proteins, are often flexible and hard to crystallize.)
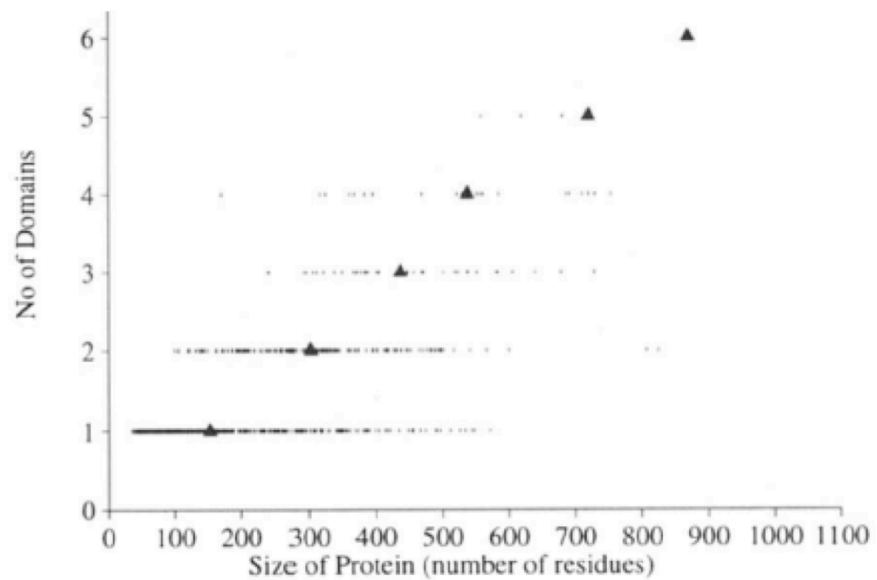


*after Jones et al. (1998) Protein Science **7**:233*

A quarter of all PDB structuires contain multi-domain proteins. The fraction of multi-domain proteins in living organism is higher, and it is higher (in general) in eukaryotes than in prokaryotes. Multi domain proteins are less likely to crystallize than single domain proteins: frequently domains are well-defined in their internal structure, but mobile relative to erach other and this mobility is detrimental to crystal growth.

Non-random relationship between domain number and chain length in the 787 representative proteins used as the basis for the CATH database

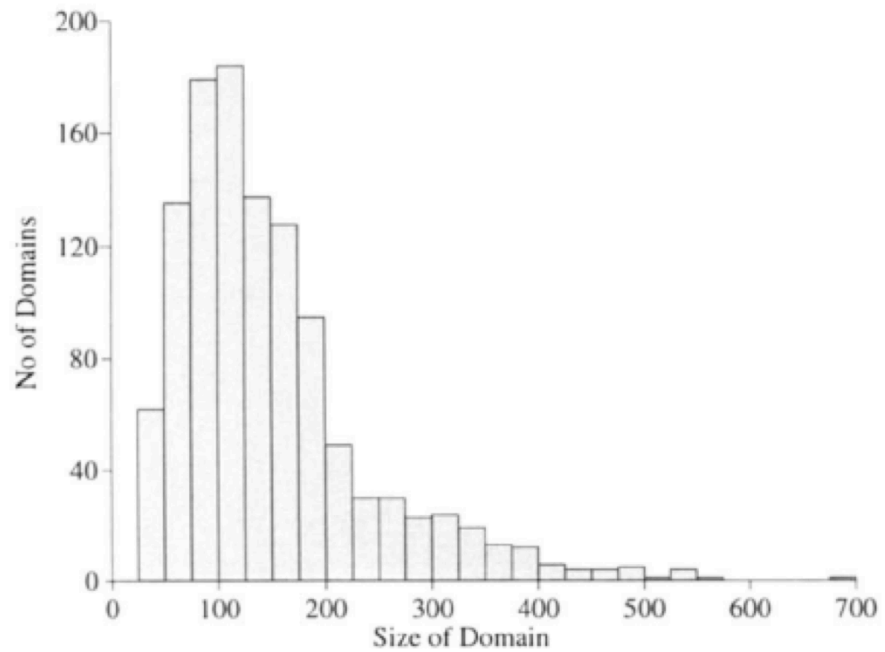Jones *et al.* (1998)
*Protein Science* **7**:233



As is to be expected, a larger fraction of longer protein chains contains more than one domain. Interestingly, the distribution correlates well with a characteristic size of properly folding domains, which is – from statistical mechanical considerations – on the order of 110 amino acids.

Domain size
in the 787
representative
proteins used
as the basis
for the CATH
database

Jones *et al.* (1998)
*Protein Science* **7**:233



Indeed, this length of 110 amino acids is approximately what we observe in nature. It results from a balance between the need to make multiple stabilizing interactions (enthalpy) and not loose too much entropy upon folding from a disordered unfolded state into the single conformation of the native state.

To ...

... identify regions of the polypeptide chain that fold independently, that are stable on their own
*(folding units; initiation sites for folding)*

... identify gene fusion or gene insertion events
from analysis of the 3D structure
*(understand evolutionary history)*

... understand protein mechanism as an additive/cooperative result of domain function
*(CDART, SMART - domain architecture)*

... allow for meaningful structural classification of proteins
*(SCOP, CATH classifications)*

Domains can be used for sequence analysis in many ways.

Possible definitions are based on independently inherited (sub)sequences (**sequence domain**), modular protein functions (**functional domain**), **folding unit** or atomic contacts (**structural domain**).

## Domain: A part of structure that can fold irrespective of the presence of other parts of structure

But: what is measured is commonly sequence, function, or structure - NOT FOLDING!

Different domain *definitions* suggest different *algorithms* to identify domains.

The separation of a structure into domains requires the arbitrary definition of thresholds in a continuum of possibilities.
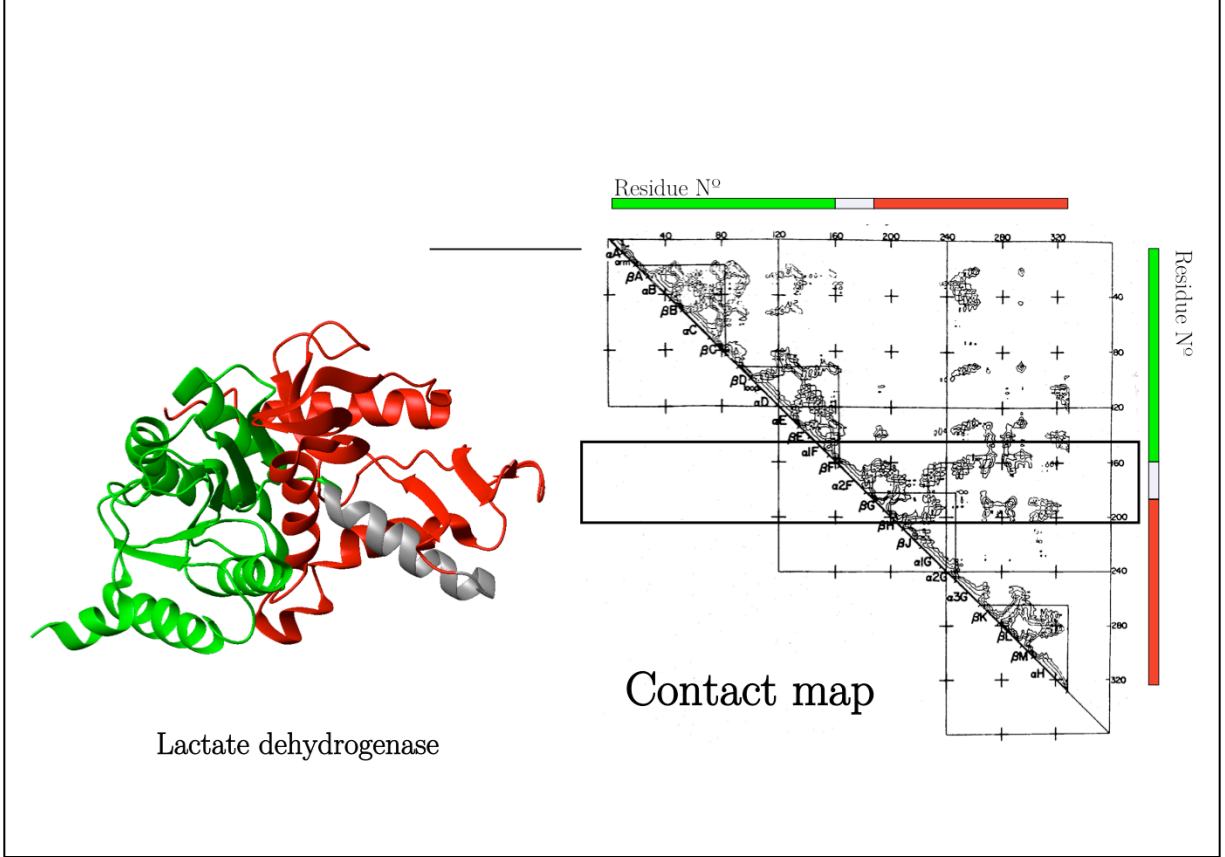
*Principle:*

Interactions between residues *within* domains are more extensive than *between* domains

"Interactions" can be quantified with many different metrics e.g. by counting inter-atomic contacts or computing buried surface area.
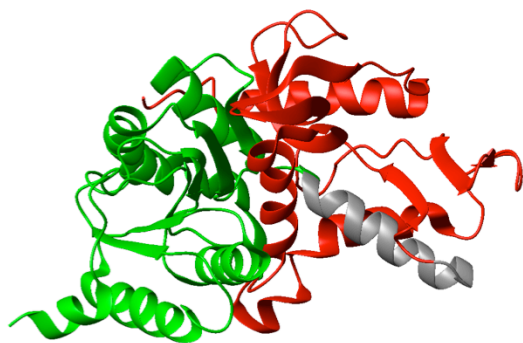
The definition of domains in structures is similar to the definition of clusters.
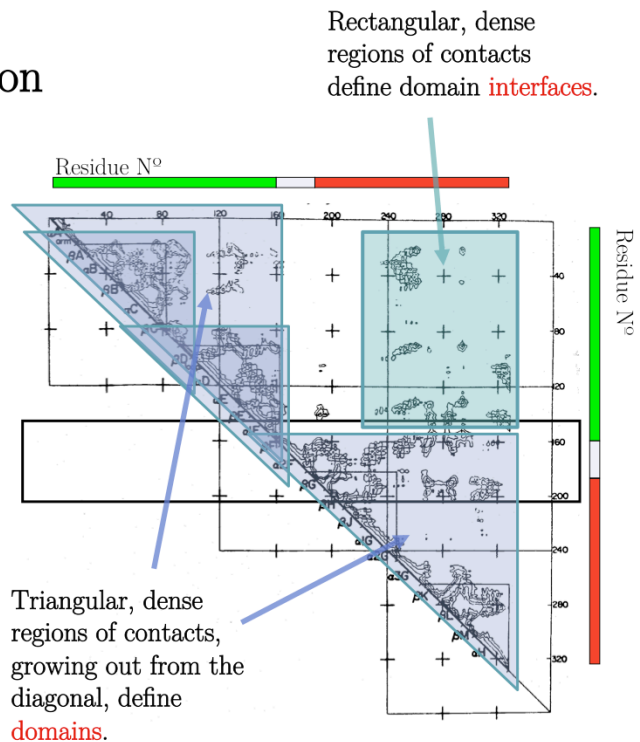
Lactate dehydrogenase

Contact map

Contact maps can illustrate this principle. A contact map plots the distance between residue pairs on a square grid.

Lactate dehydrogenase

Portions of structure that have a large number of intra-domain contacts show up as density growing out from the diagonal. Here, the red and green colours identify two distinct domains of LDH. Note that both of these domains clearly show additional subdomain structure.

Protein structure is composed hierarchically!

This is one of the problems of accurately defining domains: there may not be a natural level at which the hierarchy can be decomposed into structurally or functionally meaningful units.

However, domains much smaller than 80 amino acids or so are unlikely to fold independently.

CATH

CATH is a (largely) automated, authoritative, hierarchical classification of all PDB domains. It uses DETECTIVE, DOMAK and PUU to "chop" full-length protein structures into domains, with some manual curation applied to conflicting cases, then it uses CATHEDRAL and SSAP to find which ...

... **Class**

... **Architecture**

... **Topology**, and

... **Homology** family the domain can be classified into.

http://www.cathdb.info

Jones *et al.* (1998) Domain assignment for protein structures using a consensus approach: Characterization and analysis. *Protein Science* **7**:233-242

## DOMAK:

Maximal interactions within each unit -
minimal interaction between units (domains).

- Assign value to each parwise type of contact
- Arbitrarily split protein and calculate sum of values for both domains
- Large split values correspond to distinct domains

Siddiqui AS & Barton GJ (1995) Continuous and discontinuous domains: an algorithm for the automatic generation of reliable protein domain definitions. *Protein Sci.* **4(5)**:872-884.

## DETECTIVE:

Domains (as folding units) have recognizable hydrophobic cores.

- Define distinct hydrophobic cores by considering secondary structure, side-chain accessibility and side-chain side-chain contacts: residues are part of a core when they are in regular secondary structure and have buried side-chains that make non-polar contacts with each other
- Each core and its shell define a distinct domain

Swindells MB (1995) A procedure for detecting structural domains in proteins. *Protein Sci.* **4(1)**:103-112.

## PUU (Parser for Protein Unfolding Units):

Maximal interactions within each unit -
minimal interaction between units (domains).

- Build contact matrix for residues
- Solve eigenvalue problem related to "strength" of residue interfaces
  (this groups residues by interactions)
- Search for all reasonable bisections (guided by physical criteria that
  identify units of sufficient internal stability) of residue groups.
- Recursive bisections build folding "tree".

Holm L & Sander C (1994) Parser for protein folding units. *Proteins*
**19**(3):256-268

## CATHEDRAL (Sequential Structure Alignment Program)

Fast prescreening of domains against a database

- CATHEDRAL uses a fast graph-matching algorithm on graphs of secondary structure elements and their distances to prescreen a domain against a database
- Significantly matching domains are aligned and refined using SSAP

Pearl FMG *et al.* (2003) The CATH database: an extended protein family resource for structural and functional genomics
*Nucleic Acids Res.* **31(1)**:452-455.

SSAP  (Sequential Structure Alignment Program)
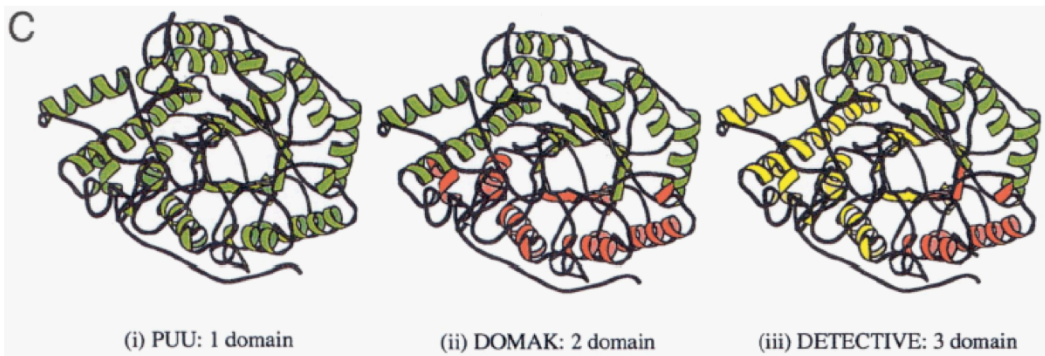
**Accurate structural alignment**

- SSAP is a general-purpose structural alignment program
- Residues are described by a vector of distance to other residues
- Alignments are sought that minimize the difference between distances

Orengo CA & Taylor WR (1996) SSAP: Sequential structure alignment

program for protein structure comparison

*Methods in Enzymology.* **266**:617-635.

TIM barrels - beta amylase 1BTC:



(i) PUU: 1 domain      (ii) DOMAK: 2 domain      (iii) DETECTIVE: 3 domain

Jones *et al.* (1998) Domain assignment for protein structures using a consensus approach:
Characterization and analysis. *Protein Science* **7**:233-242

Domain classification programs do not always agree. The three classifiers used for CATH all disagree on the 1BTC TIM-barrel structure.
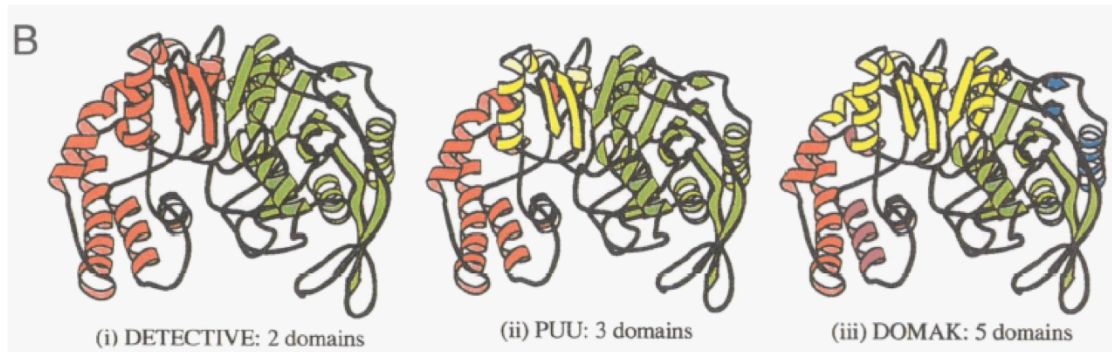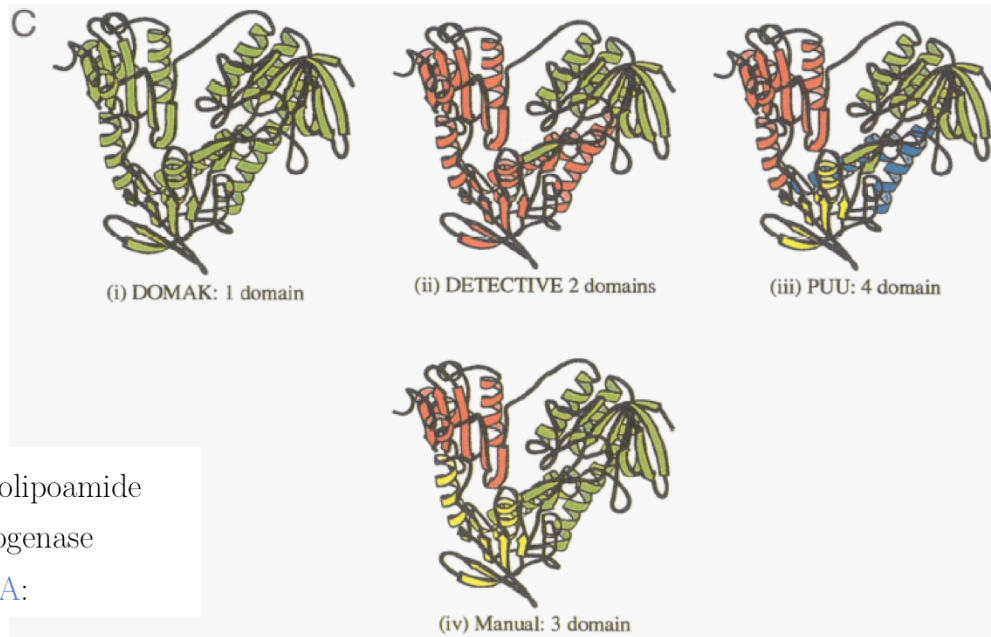
beta propeller - neuraminidase 1NSC_A:



(i) PUU: 1 domain     (ii) DETECTIVE: 2 domain     (iii) DOMAK: 3 domain

Jones *et al.* (1998) Domain assignment for protein structures using a consensus approach:
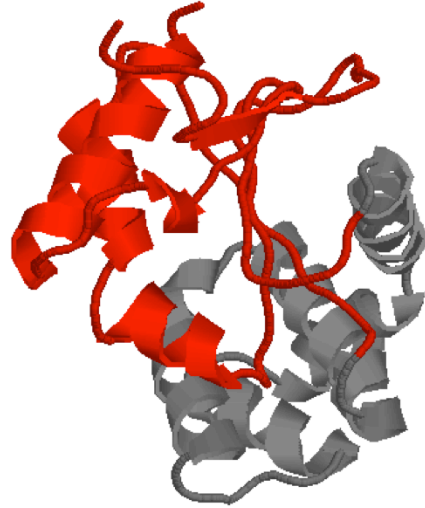Characterization and analysis. *Protein Science* **7**:233-242

Adenylosuccinate synthetase 1ADE_A:



(i) DETECTIVE: 2 domains    (ii) PUU: 3 domains    (iii) DOMAK: 5 domains

Jones *et al.* (1998) Domain assignment for protein structures using a consensus approach:
Characterization and analysis. *Protein Science* **7**:233-242

## DOMAIN DEFINITION EXAMPLES: CATH CONSENSUS



(i) DOMAK: 1 domain

(ii) DETECTIVE 2 domains

(iii) PUU: 4 domain

Dehydrolipoamide
dehydrogenase
1LPF_A:

(iv) Manual: 3 domain

Jones *et al.* (1998) Domain assignment for protein structures using a consensus approach:
Characterization and analysis. *Protein Science* **7**:233-242

CATH is built from domains that are defined from a consensus of its three core
classifiers – with manual intervention if necessary.

Expert intervention and curation is extremely helpful to maintain the quality of the
database – but anything that needs to be done manually **does not scale**.

COMPLICATIONS OF DOMAIN DEFINITION: DOMAIN INSERTIONS

Domain insertion

2TRX.PDB       1A2J.PDB

Thioredoxin       Protein disulfide isomerase

Sometimes the definition of domains is made additionally difficult by the complicated biology of natural proteins. Protein disulfide isomerase has a thioredoxin fold in principle – but there is a complete independent domain inserted into one of its loops.

11BG (Bull seminal ribonuclease)

Domain swapping leads to elements of the *same sequence* being integrated into the structure of *different domains*. This requires a strained connection region, otherwise the higher local concentration of the *intra*-domain contact would override the possibility of *inter*-domain interactions.

The result is a very tight and often essentially irreversible association that requires complete unfolding of the domains to undo.

The top and bottom parts of the image correspond to well structured, compact domains. However both of these domains incorporate a helix from the respectively other sequence of this homodimer.

cf.

Wodak, Malevanets & MacKinnon (2015) The Landscape of Intertwined Associations in Homooligomeric Proteins. Biophys J 109(6):1087-1100.

DOMAIN FAMILIES

2IMM - Immunoglobulin like

1TPH - beta/alpha barrel

1UBI - beta grasp

1CKA - SH3-like barrel

3CHY - Flavodoxin like

1FXD - Ferredoxin like

1SNC - OB-fold

1NFN - 4-helix bundle

Example:

The eight most frequent SCOP Superfolds

http://scop.mrc-lmb.cam.ac.uk/

(Last update 2009)

Domains are found in many families – but the boundaries between families may be somewhat fluid as it may be not obvious when two domains are actually different. Moreover, even within families, we may find proteins that - as far as we can tell - are not actually related to each other, but have arrived at their particular fold through **convergent** evolution. Nevertheless, assigning folds to families allows us to bring some order to the zoo of possibilities and this underlies approaches to organize and retrieve domains in databases.

The examples above are taken from the highly curated SCOP database, which unfortunately has not been updated since 2009.
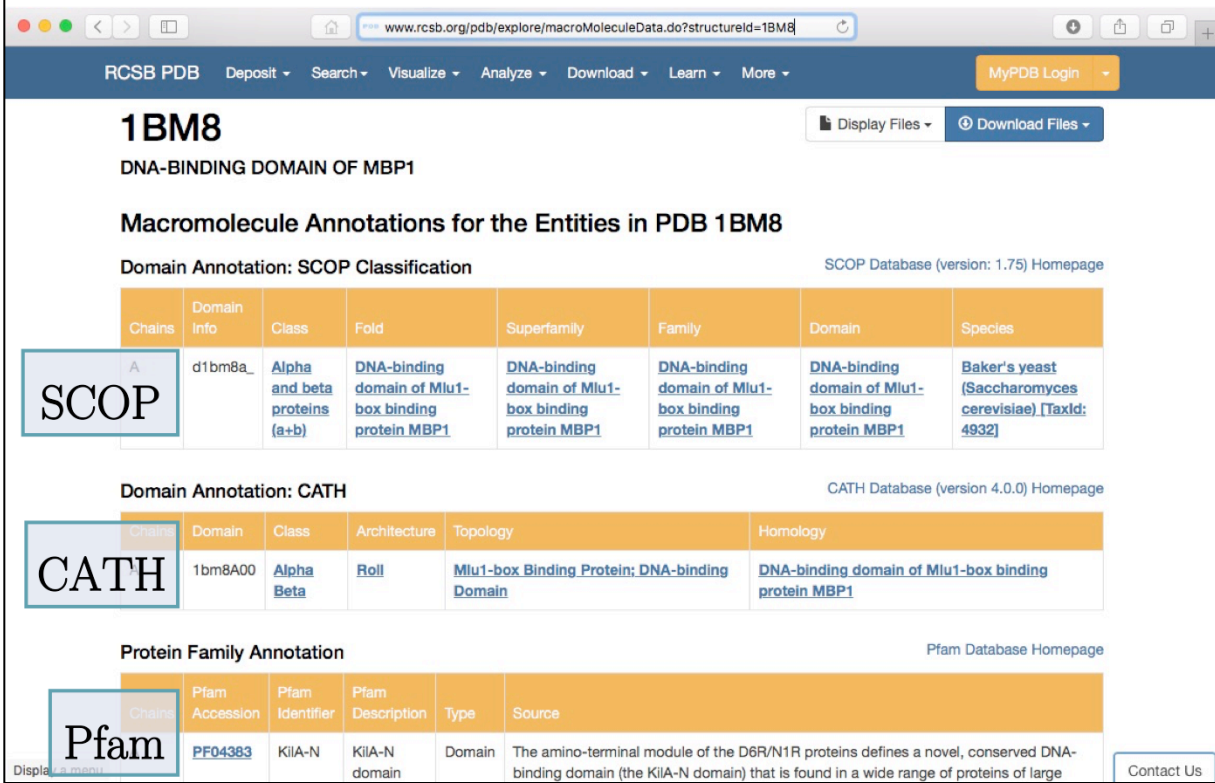
CDD

CDD is the integrated domain analysis and search tool at the NCBI. It imports structural domain definitions from SCOP and CATH, as well as sequence domains from **Pfam**, thus providing an integrated access to sequence- as well as structural domains. It's results are available via cross references in all typical sequence resources, as well as in BLAST searches *etc.*

http://www.ncbi.nlm.nih.gov/Structure/cdd/cdd.shtml

Marchler-Bauer *et al.* (2011) CDD: a Conserved Domain Database for the functional annotation of proteins. *Nucleic Acids Research* **39(DB issue)**:D225-D229

DOMAIN DATABASES

Unified access via PDB:

The PDB provides cross references into the domain databases for each of its entries. SCOP and CATH are databases for structural domains, Pfam defines sequence domains.

http://steipe.biochemistry.utoronto.ca/abc

BORIS . STEIPE@UTORONTO.CA

DEPARTMENT OF BIOCHEMISTRY & DEPARTMENT OF MOLECULAR GENETICS
UNIVERSITY OF TORONTO, CANADA