

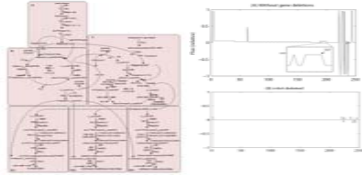


Genome scale Comparison of ligand binding sites  
in protein structures: Algorithms and applications  
in drug discovery

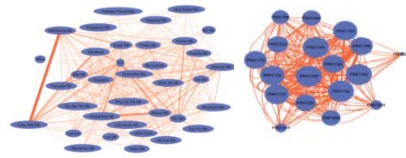
***Nagasuma Chandra***  
***Indian Institute of Science***  
***Bangalore, India***

# Nagasuma Chandra, Research Overview

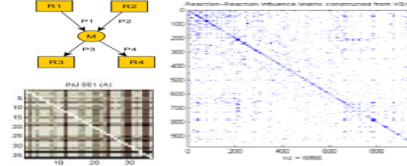
## Pathway Modelling



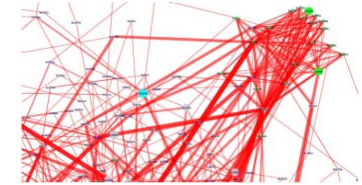
## Reactome Modelling



## Protein-Protein influences



## Interactome Modelling



**SYSTEMS  
BIOLOGY**

## Algorithms for Structure Analysis



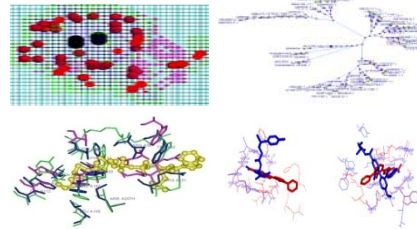
**STRUCTURAL  
BIOINFORMATICS**

## Drug Discovery

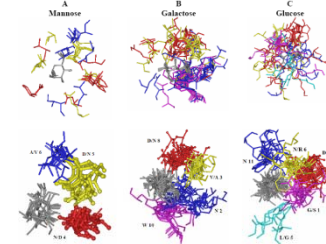
### Drug Target Identification



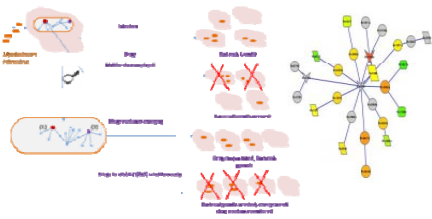
### Druggability Assessment



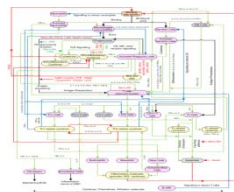
### Lead Identification



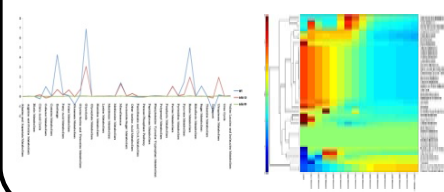
### Drug Resistance



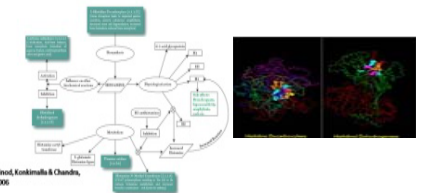
### Host-Pathogen Modelling



### Modelling drug effects

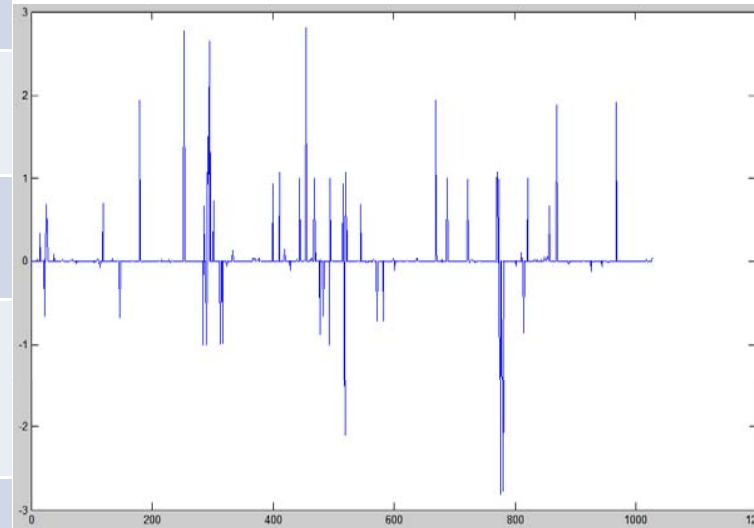


### Drug failure Analysis



# Modeling Metabolism in *M. tuberculosis*

|                           |      |
|---------------------------|------|
| Genes                     | 661  |
| Proteins                  | 543  |
| Reactions (Intra Systems) | 939  |
| Reactions (Exchange)      | 88   |
| Gene Association Reaction | 77%  |
| Metabolites               | 828  |
| Average Confidence Level  | 2.31 |



*Mtb Reactome Flux profile*

Growing bacteria *In silico* under different media

*In silico* gene deletions

Essential nutrients required for growth

Response to other Perturbations-  
Nutrient Uptake

## Insights obtained

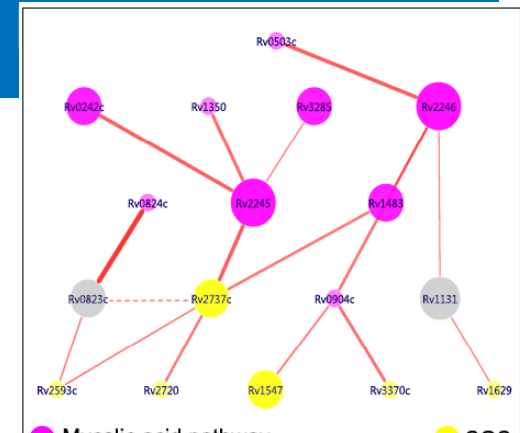
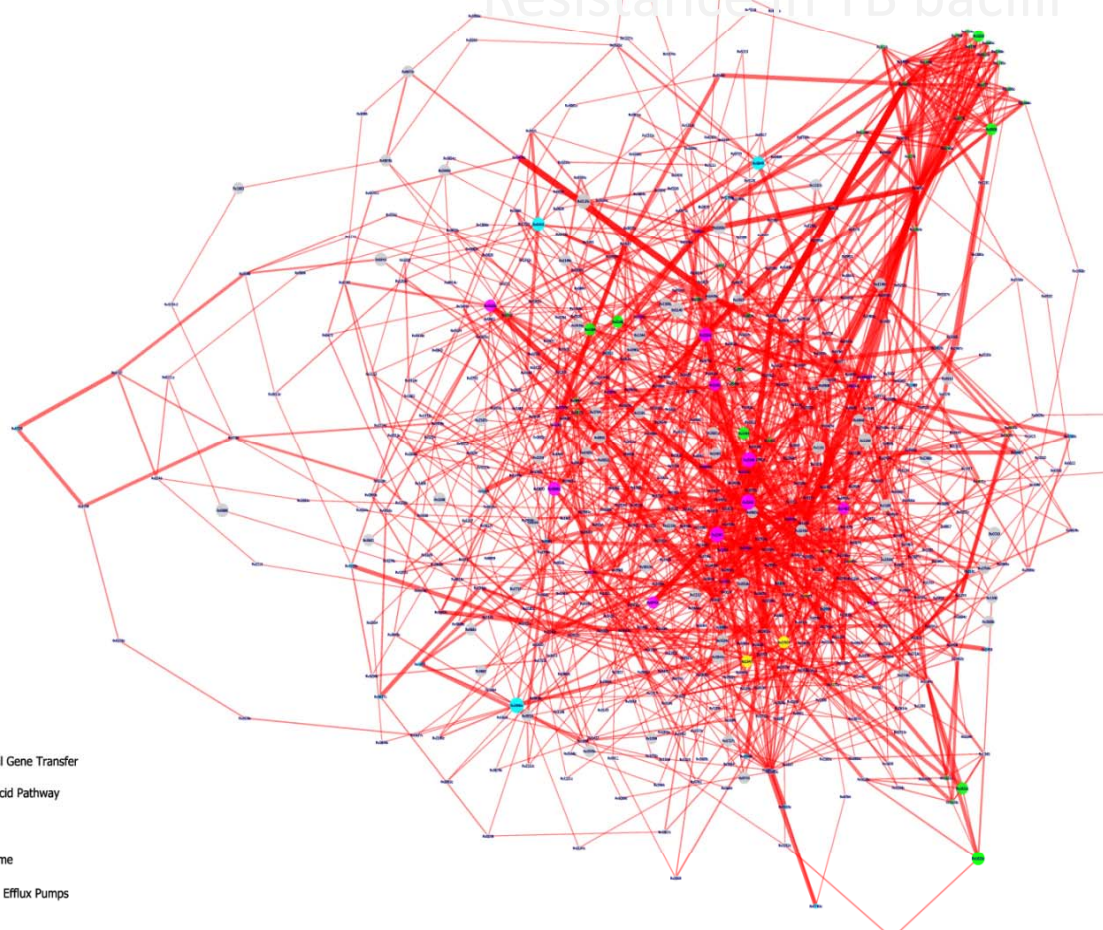
- Gene Essentiality: 220 essential genes
- Nutrient essentiality:
- Hard coupled reaction sets: groups of reactions that are forced to operate in unison due to mass conservation and connectivity constraints)
- Fatty Acid Metabolism & Lysine Metabolism

Raman *et al.*, 2008

# Drug Resistance Pathways

Abstraction of the flow of information that

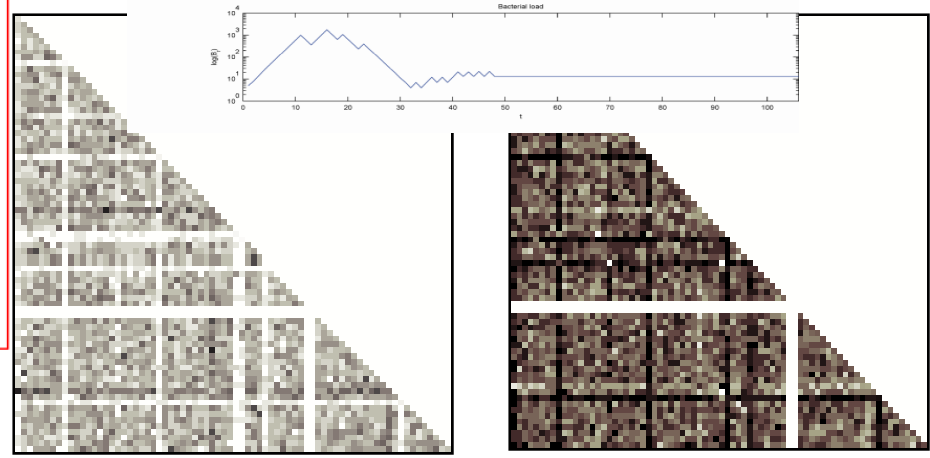
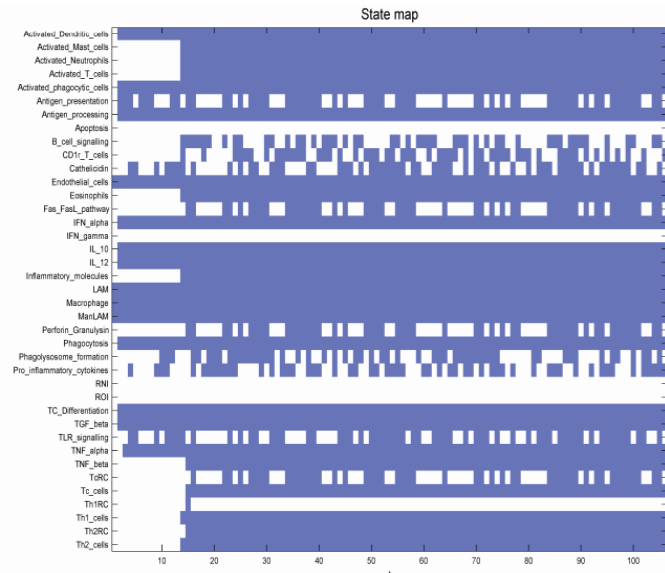
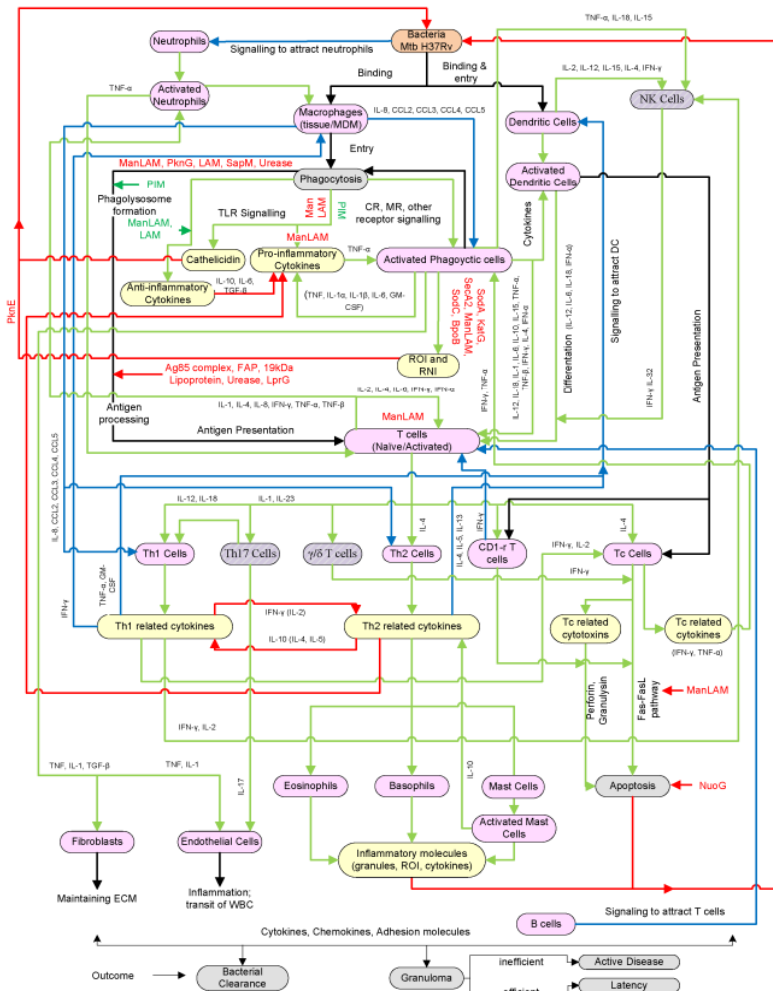
Resistance in TB bacilli



- Network of shortest paths from MAP to Resistance Genes
- 616 nodes and 1,683 edges
- Paths scored based on edge frequency, up-regulation of source and target nodes

*Raman and Chandra, 2008, BMC Microbiology*

# Host-Pathogen Interaction Network



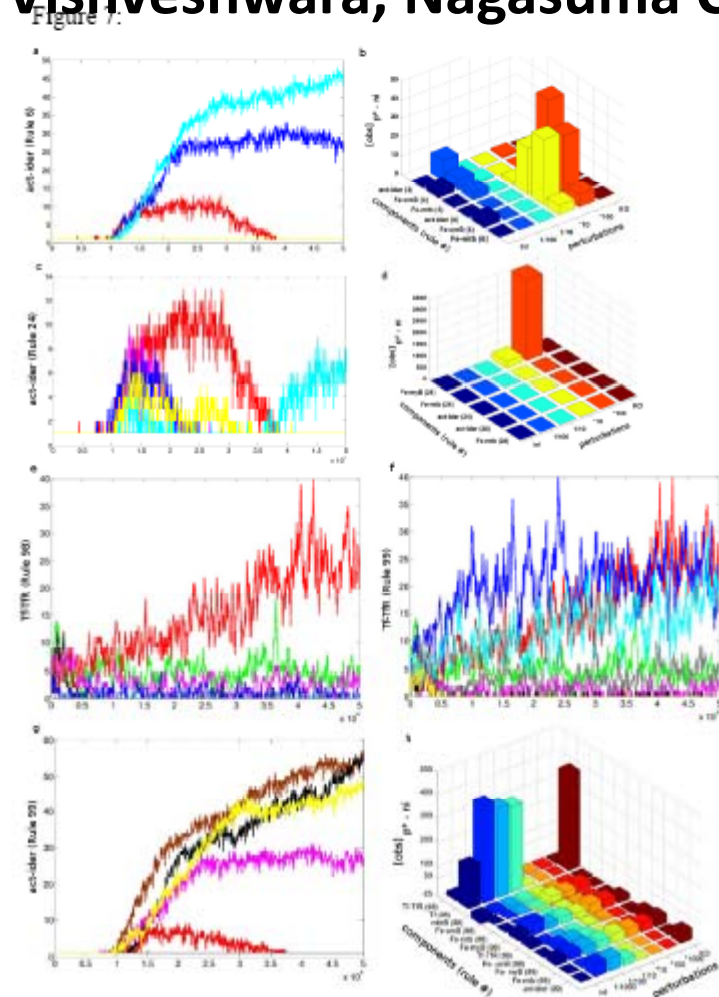
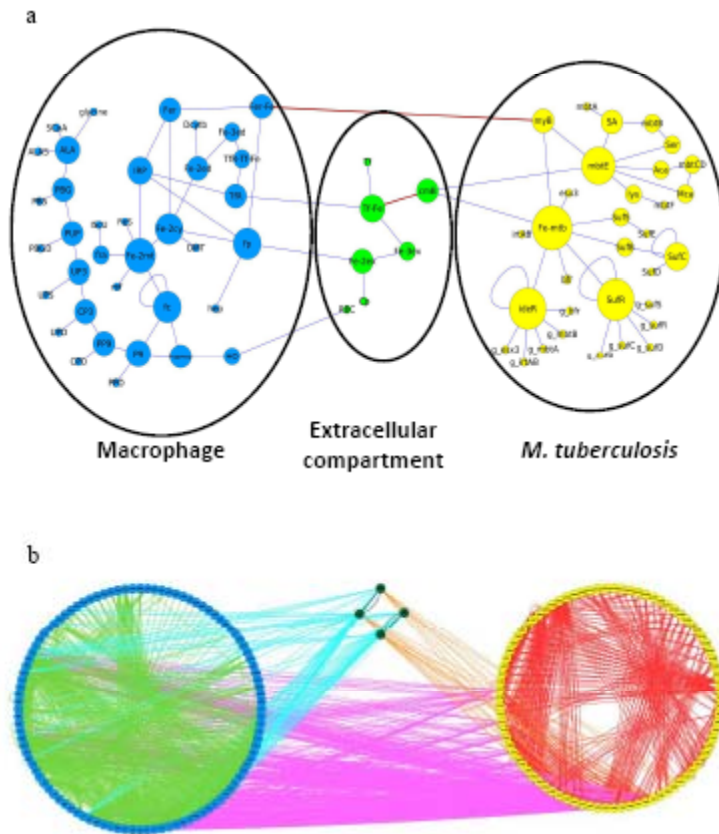
A Boolean model of HPIs developed, Simulations to capture a variety of scenarios  
 5 Raman, Bhat & Chandra, Mol. Biosyst, 2010



# Modelling iron homeostatis in M.tuberculosis

Soma Ghosh, KVS Prasad, Sarswathi Vishveshwara, Nagasuma Chandra

Figure 4:



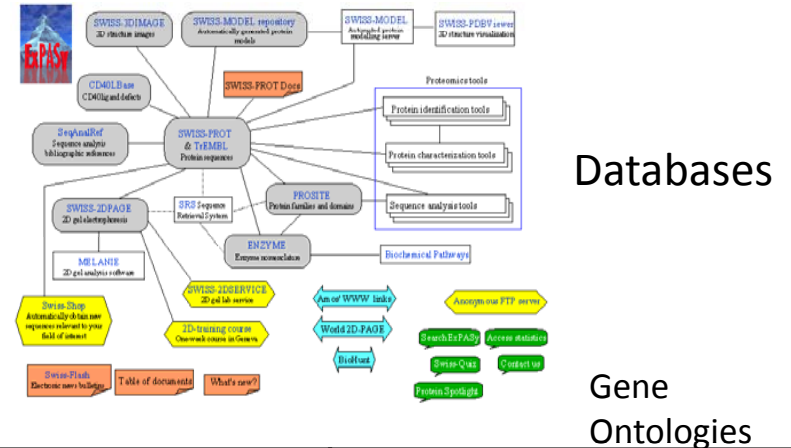
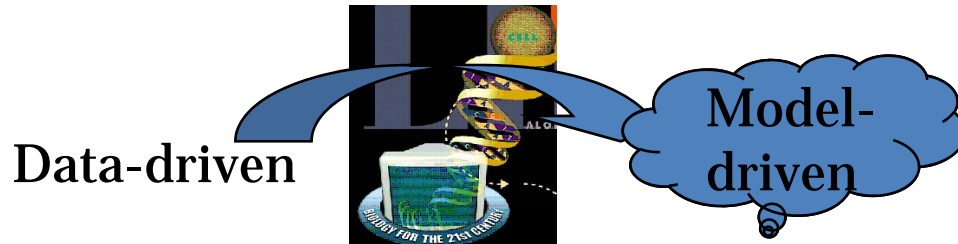
Ghosh et al., Mol Biosyst, 2011, In Press

# Integration of Systems Perspective with Structural level detail

Raman, Kalidas, Chandra, 2008

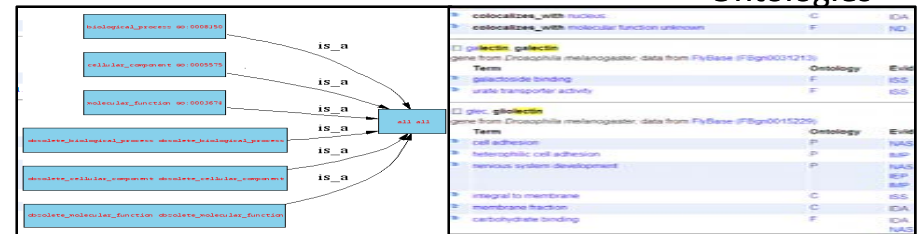
Model Driven Drug Discovery: Principles and Practices In: Biological Database Modeling, Artech House [ISBN 1596932589](#)

# Comparison of protein molecules



## How to compare?

- Annotation- keyword
- Function identification
  - Molecular level
  - Cellular level
- Sequence - Alignment
- Structure - Fold comparison

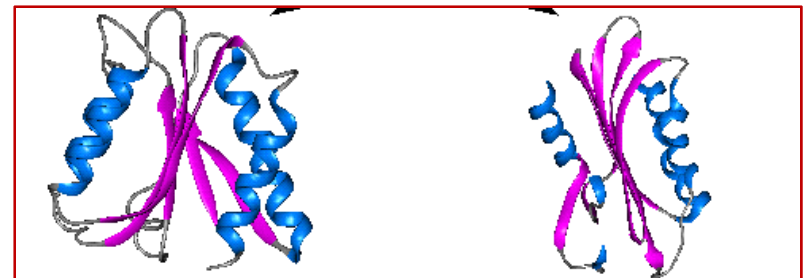


[15 DROME/5-142](#)  
[14 SHEEP/5-136](#)  
[1 HUMAN/5-137](#)  
[HUMAN/4-136](#)  
[HUMAN/4-136-SS](#)

```

SVNFEVNS.....PDCIDNITVYHFKVVVIFT.....QTIIEINY..KRNGE.WSSLH.QEK.
QVDFGIG.....TGQEGNIPFRFWYCDGM.....VVMNT..FTDGS.WQ..K.EEK.
EVNFEYTG.....MDEDSDIAFQFRLHFGH.....PAIMNS..RVFGI.WR..Y.EEK.
QVDFHTL.....MKEESDIVHFQVCFGR.....RVVMNS..REYGA.WK..Q.QVE.
EEEESS.....SSTTSCEEEEEEETT.....EEEEEE..EETT.EC..C.CEE.
    
```

Sequence Alignment



Structure Alignment

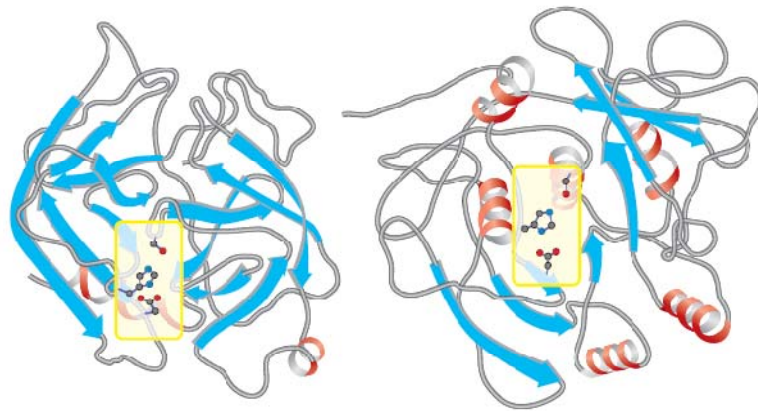


# Problems with these approaches

- **Dissimilar Sequences- Similar Function** → benzoylformate decarboxylase (BFD) and pyruvate decarboxylase (PDC)
- **Similar sequences- Different Function** → Steroid delta isomerase; nuclear transporter 2; scytalone dehydratase

**Similar Structures- Different Function** → triose phosphate isomerase and FMN-linked oxidoreductases

**Dissimilar Structures – Similar Function** → ATP binding proteins from different SCOp families; C-type lectin and bulb lectin



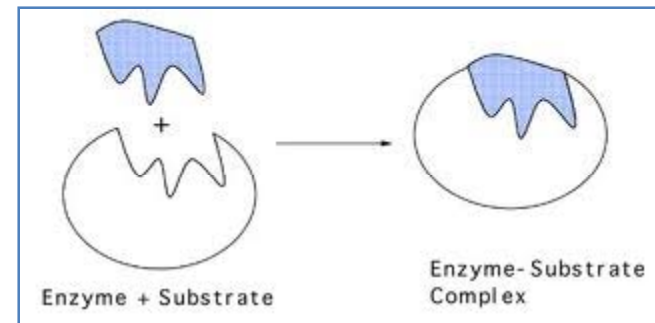
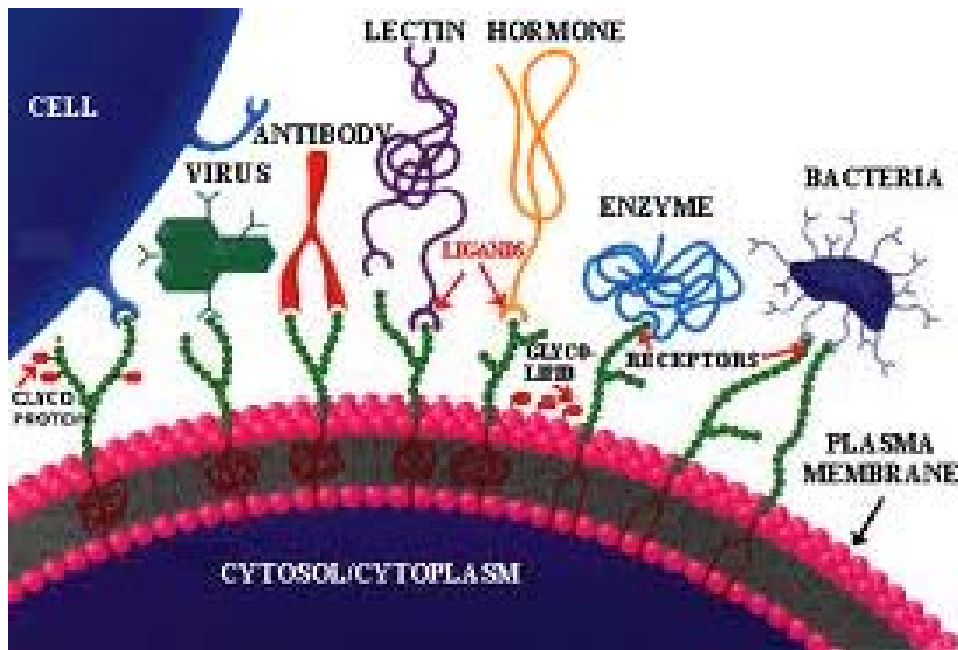
Chymotrypsin & Subtilisin

*What really matters for a protein molecule is its function and not what means it uses to achieve it !*

It's the **meaning** that counts....

*Whether two proteins can recognize the  
same molecules*

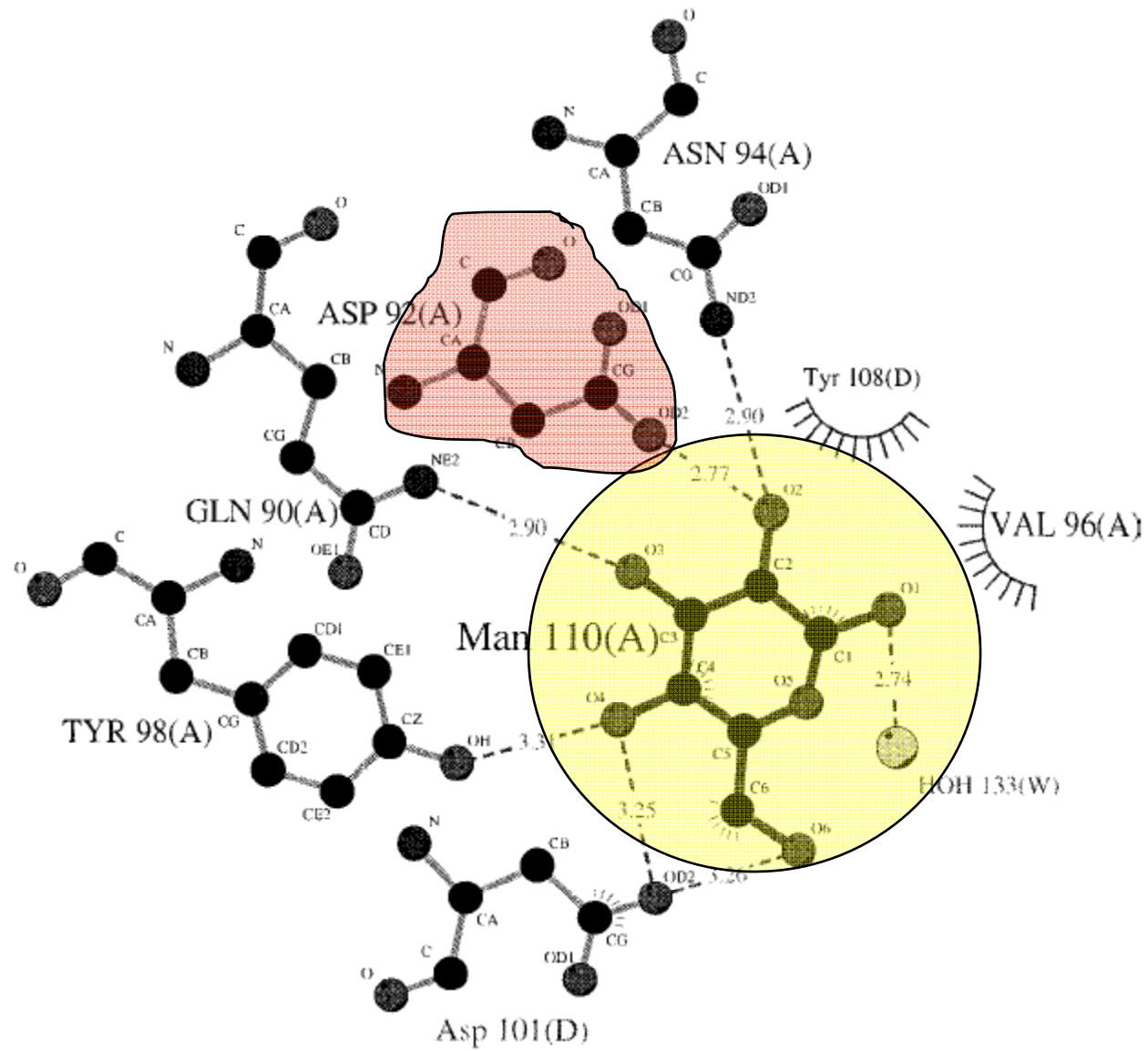
# Molecular recognition



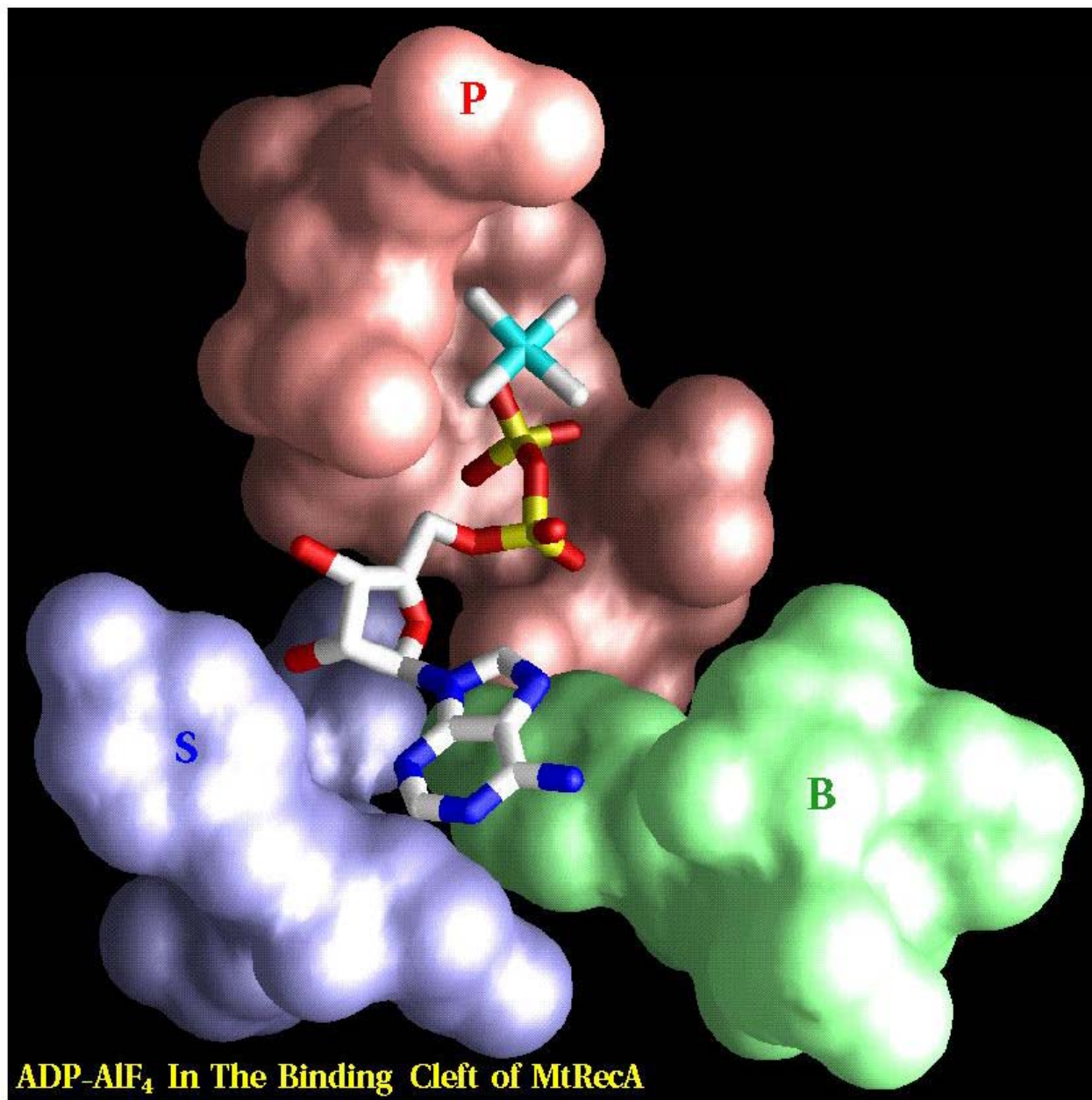
|      |   |     |                   |           |           |           |              |            |           |            |
|------|---|-----|-------------------|-----------|-----------|-----------|--------------|------------|-----------|------------|
| GARL | 1 | 39  | -YAGQSLDVEPYHLI   | <b>MQ</b> | <b>ED</b> | <b>CD</b> | <b>NLVLY</b> | DHS - TA   | VWTSNT    | DIPG - - - |
| GARL | 2 | 71  | -SNTDIPGKKGCKAV   | <b>LQ</b> | <b>SD</b> | <b>GN</b> | <b>FVY</b>   | DAEGRS     | LWASHS    | VRG - - -  |
| GARL | 3 | 103 | ASHSVRGN - GNYVLV | <b>LQ</b> | <b>ED</b> | <b>GN</b> | <b>VVIY</b>  | G - - - SD | IWSTNT    | YKGSAR -   |
| SNOW | 1 | 34  | -STGEFLNYGSFVFI   | <b>MQ</b> | <b>ED</b> | <b>CD</b> | <b>NLVLY</b> | DVD - KP   | IWATNT    | GGLS - - - |
| SNOW | 2 | 65  | ATNTGGLS - RSCFLS | <b>MQ</b> | <b>TD</b> | <b>DG</b> | <b>NLVVY</b> | NPSNKP     | IWASNT    | GGQ - - -  |
| SNOW | 3 | 97  | ASNTGGQN - GNYVCI | <b>LQ</b> | <b>KD</b> | <b>RN</b> | <b>VVIY</b>  | G - - - TD | RWATGT    | HTGLVG -   |
| HIPP | 1 | 1   | - - - - -         | <b>MQ</b> | <b>QD</b> | <b>CD</b> | <b>NLVLY</b> | DVD - KP   | - - - - - | - - - - -  |
| HIPP | 2 | 18  | ASNTGGLA - RGCHLS | <b>MQ</b> | <b>SD</b> | <b>GN</b> | <b>NLVVY</b> | SPSNSP     | IWASNT    | QGE - - -  |
| HIPP | 3 | 50  | ASNTQGEN - GNYVCI | <b>VQ</b> | <b>KD</b> | <b>RN</b> | <b>VVIY</b>  | G - - - TA | RWATGT    | YTGTVG -   |
| NARC | 1 | 35  | -STGQFLSYGSYVFI   | <b>MQ</b> | <b>ED</b> | <b>CD</b> | <b>NLVLY</b> | DVD - KP   | IWATNT    | GGLS - - - |
| NARC | 2 | 66  | ATNTGGLS - SDCHLS | <b>MQ</b> | <b>TD</b> | <b>DG</b> | <b>NLVVY</b> | SPQNKA     | IWASNT    | DGE - - -  |
| NARC | 3 | 98  | ASNTDGEN - GHFVCI | <b>LQ</b> | <b>KD</b> | <b>RN</b> | <b>VVIY</b>  | G - - - TD | RWATGT    | YTGA VG -  |
| ONIO | 1 | 26  | -YAGQSLVVEQYTFI   | <b>MQ</b> | <b>DD</b> | <b>CD</b> | <b>NLVLY</b> | EYS - TP   | IWASNT    | GVTG - - - |
| ONIO | 2 | 58  | -SNTGVTGKNGCRAV   | <b>MQ</b> | <b>AD</b> | <b>GN</b> | <b>FVY</b>   | DVKGRA     | VWASNS    | RRG - - -  |
| ONIO | 3 | 90  | ASNSRRGN - GNYILV | <b>LQ</b> | <b>KD</b> | <b>RN</b> | <b>VVIY</b>  | G - - - SD | IWSTGT    | YRKKVG -   |
| RAMS | 1 | 71  | -SNTGVSGRNGCRAV   | <b>MQ</b> | <b>AD</b> | <b>GN</b> | <b>FVY</b>   | DSNSRA     | VWASQS    | RRG - - -  |
| RAMS | 2 | 103 | ASQSRRGN - GNYILA | <b>LQ</b> | <b>ED</b> | <b>DR</b> | <b>NVVIY</b> | G - - - TD | IWSTGT    | YRRGVG -   |
| URSI | 1 | 38  | -YAGQSLVEEGPYKLI  | <b>MQ</b> | <b>ED</b> | <b>CD</b> | <b>NLVLY</b> | EHS - TP   | VWATNT    | GVTG - - - |
| URSI | 2 | 70  | -TNTGVTGRNGCKAV   | <b>MQ</b> | <b>AD</b> | <b>GN</b> | <b>FVY</b>   | DSNGRA     | VWASNS    | IKG - - -  |
| URSI | 3 | 102 | ASNSIKGN - GNYILV | <b>LQ</b> | <b>ND</b> | <b>DR</b> | <b>NVVIY</b> | G - - - SD | IWSTGT    | YIRGVG -   |
| LEEK | 1 | 41  | -YAGQSLDVEQYKFI   | <b>MQ</b> | <b>DD</b> | <b>CD</b> | <b>NLVLY</b> | EYS - TP   | IWASNT    | GVTG - - - |
| LEEK | 2 | 73  | -SNTGVTGKNGCRAV   | <b>MQ</b> | <b>KD</b> | <b>GN</b> | <b>FVY</b>   | DVNGRP     | VWATNS    | VRG - - -  |
| LEEK | 3 | 105 | ATNSVRGN - GNYILV | <b>LQ</b> | <b>QD</b> | <b>DR</b> | <b>NVVIY</b> | G - - - SD | IWSTGT    | YRRSAG -   |
| SHAL | 1 | 37  | -YAGQSLVEEQYTFI   | <b>MQ</b> | <b>DD</b> | <b>CD</b> | <b>NLVLY</b> | EYS - TP   | IWASNT    | GITG - - - |
| SHAL | 2 | 69  | -SNTGITGKNGCRAV   | <b>MQ</b> | <b>PD</b> | <b>DG</b> | <b>NFVY</b>  | NVKGRA     | VWASNS    | RRG - - -  |
| SHAL | 3 | 101 | ASNSRRGN - GNYILV | <b>LQ</b> | <b>KD</b> | <b>RN</b> | <b>VVIY</b>  | G - - - SD | IWSTGT    | YRKKVG -   |
| TULI | 1 | 54  | -YGGQSLTWESYTFI   | <b>MQ</b> | <b>TD</b> | <b>CD</b> | <b>NLVLY</b> | EGN - GP   | IWASGS    | NDLG - - - |
| TULI | 2 | 85  | ASGSNDLG - SGCYVT | <b>MQ</b> | <b>KD</b> | <b>DG</b> | <b>NLVIY</b> | SKSGNS     | VWASQT    | HQA - - -  |
| TULI | 3 | 117 | ASQTHQAE - GNYVLV | <b>LQ</b> | <b>KD</b> | <b>DR</b> | <b>NVVIY</b> | G - - - PS | LWATNT    | DQFSLT -   |
| CLIV | 1 | 38  | -SPGESLSHGRYVFI   | <b>MQ</b> | <b>ED</b> | <b>CD</b> | <b>NLVLY</b> | DVD - RP   | IWATNT    | GGIS - - - |
| CLIV | 2 | 69  | ATNTGGIS - HGCHLS | <b>MQ</b> | <b>AD</b> | <b>GN</b> | <b>NLVVY</b> | SQRNNP     | IWASNT    | GGE - - -  |
| CLIV | 3 | 102 | -SNTGGENDANYVLI   | <b>LQ</b> | <b>KD</b> | <b>DR</b> | <b>NVVIY</b> | G - - - PA | RWATGT    | YTGIVG -   |
| EPIP | 1 | 39  | -GTGGLALGGYIFKI   | <b>IQ</b> | <b>AD</b> | <b>CD</b> | <b>NLVLY</b> | DNN - RA   | VWASGT    | NGRA - - - |
| EPIP | 2 | 70  | ASGTNGRA - TDCILS | <b>MQ</b> | <b>RD</b> | <b>DG</b> | <b>NLVIY</b> | S - GSRV   | IWASNT    | NRQS - - - |
| EPIP | 3 | 101 | ASNTNRQS - GNYILI | <b>LQ</b> | <b>RD</b> | <b>DR</b> | <b>NVVIY</b> | DNSNNA     | IWATGT    | NVG - - -  |
| LIST | 1 | 38  | -NTGQSLTDGGNAFI   | <b>MQ</b> | <b>ED</b> | <b>CD</b> | <b>NLVLY</b> | ESS - RP   | TWASGT    | YHRG - - - |
| LIST | 2 | 69  | ASGTYHRG - SGCYLA | <b>MQ</b> | <b>ND</b> | <b>DG</b> | <b>NLVVY</b> | DNRNRA     | IWASQT    | DRQ - - -  |
| LIST | 3 | 102 | -SQTDRQNVGTVILI   | <b>LQ</b> | <b>KD</b> | <b>DH</b> | <b>NVVIY</b> | T - - - NP | IWATGT    | NRFGSP -   |
| CYMB | 1 | 40  | -NPGQSLTSGNVDLA   | <b>MQ</b> | <b>YD</b> | <b>CD</b> | <b>NLVLY</b> | DNG - KS   | IWSSGT    | YSGS - - - |
| CYMB | 2 | 71  | SSGTYGS - SGCYVT  | <b>LQ</b> | <b>TD</b> | <b>DG</b> | <b>NLVIY</b> | DNTNNP     | LWASNT    | SGE - - -  |
| CYMB | 3 | 103 | ASNTSGEN - GNYILI | <b>LQ</b> | <b>KD</b> | <b>DR</b> | <b>NLVIY</b> | S - - - HP | IWATGT    | NHAGSV -   |
| CUCK | 1 | 37  | -NTDGRLTNGELTLI   | <b>MQ</b> | <b>VD</b> | <b>CD</b> | <b>NLVLY</b> | HG - - -   | GWQSN     | TANNGRDC   |
| CUCK | 2 | 157 | -YGDGTL SARNHKL   | <b>MQ</b> | <b>GD</b> | <b>CD</b> | <b>NMVLY</b> | GGK - Y -  | GWQSN     | THGNGK - - |
| POLY | 1 | 40  | -PSGHS LNTGSYRLI  | <b>MQ</b> | <b>AD</b> | <b>CD</b> | <b>NLVVY</b> | DSG - KP   | VWATNT    | GGLA - - - |
| POLY | 2 | 103 | QTNTNEKE - DHYVLV | <b>LQ</b> | <b>QD</b> | <b>DR</b> | <b>NVVIY</b> | G - - - PA | VWATGS    | GPAVGL -   |
| ALOE | 1 | 11  | -HENQYISYGPYEFI   | <b>MQ</b> | <b>HD</b> | <b>CD</b> | <b>NLVLY</b> | ESG - NP   | TWASNT    | GGLA - - - |

Multiple alignment of bulb lectins indicating the mannose binding sequence motif

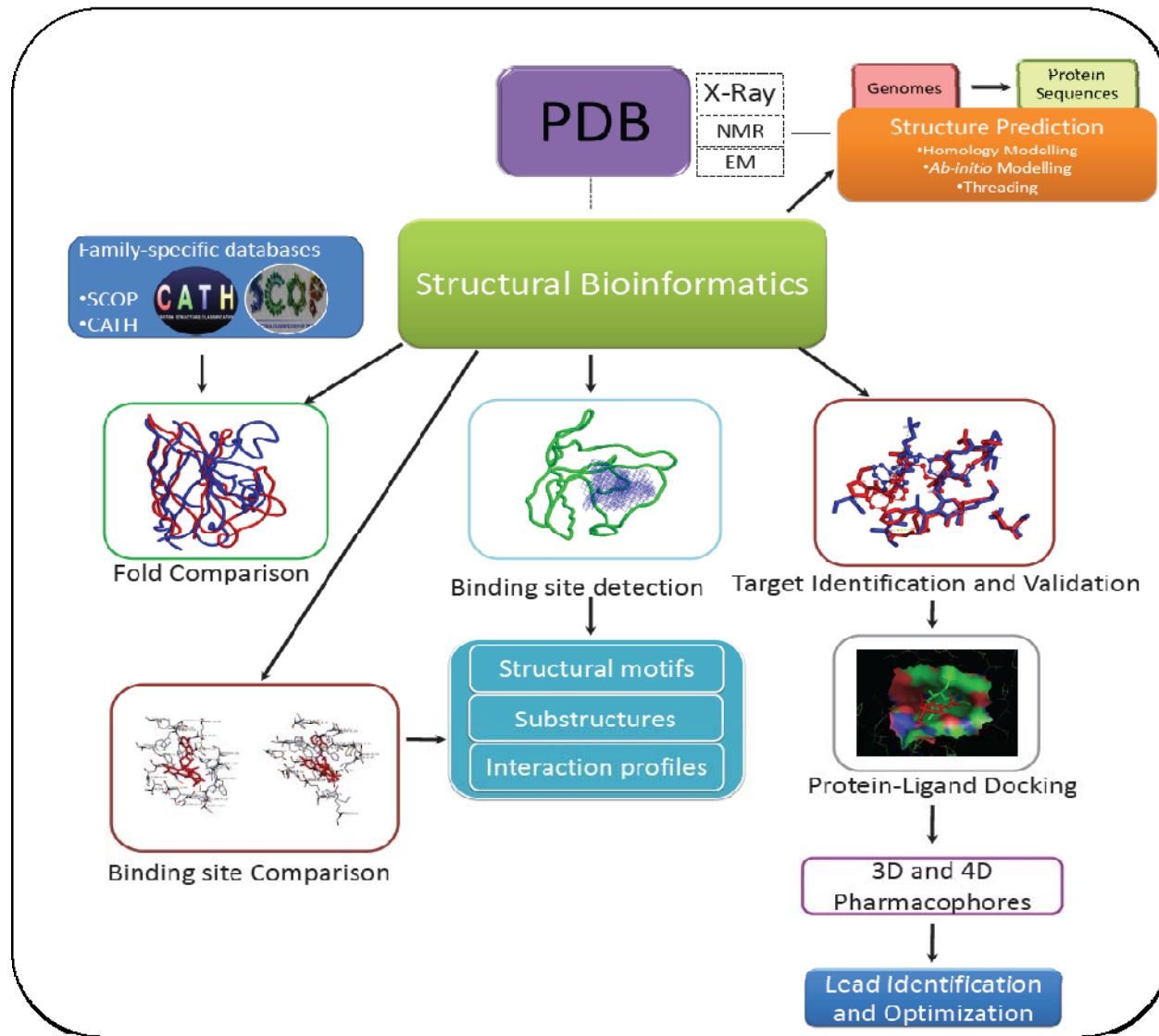




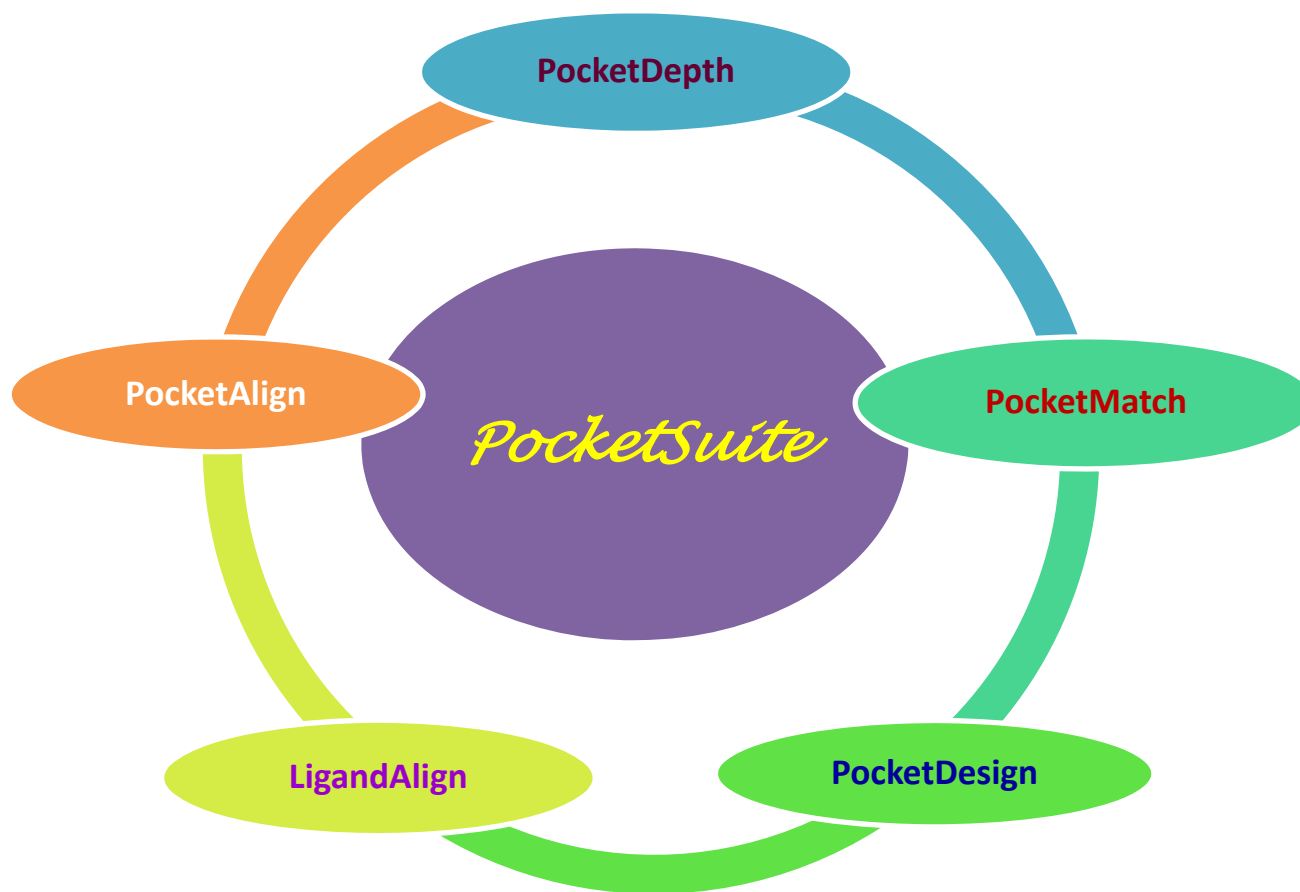
A structural motif specific for mannose recognition



# Structural Bioinformatics



## Structural Bioinformatics: *Development of 5 novel algorithms integrated into PocketSuite*



# *Binding site Prediction Methods:*

Homology-based methods

- Alignment with known sites
- Conservation

Sequence-based methods

- Motifs

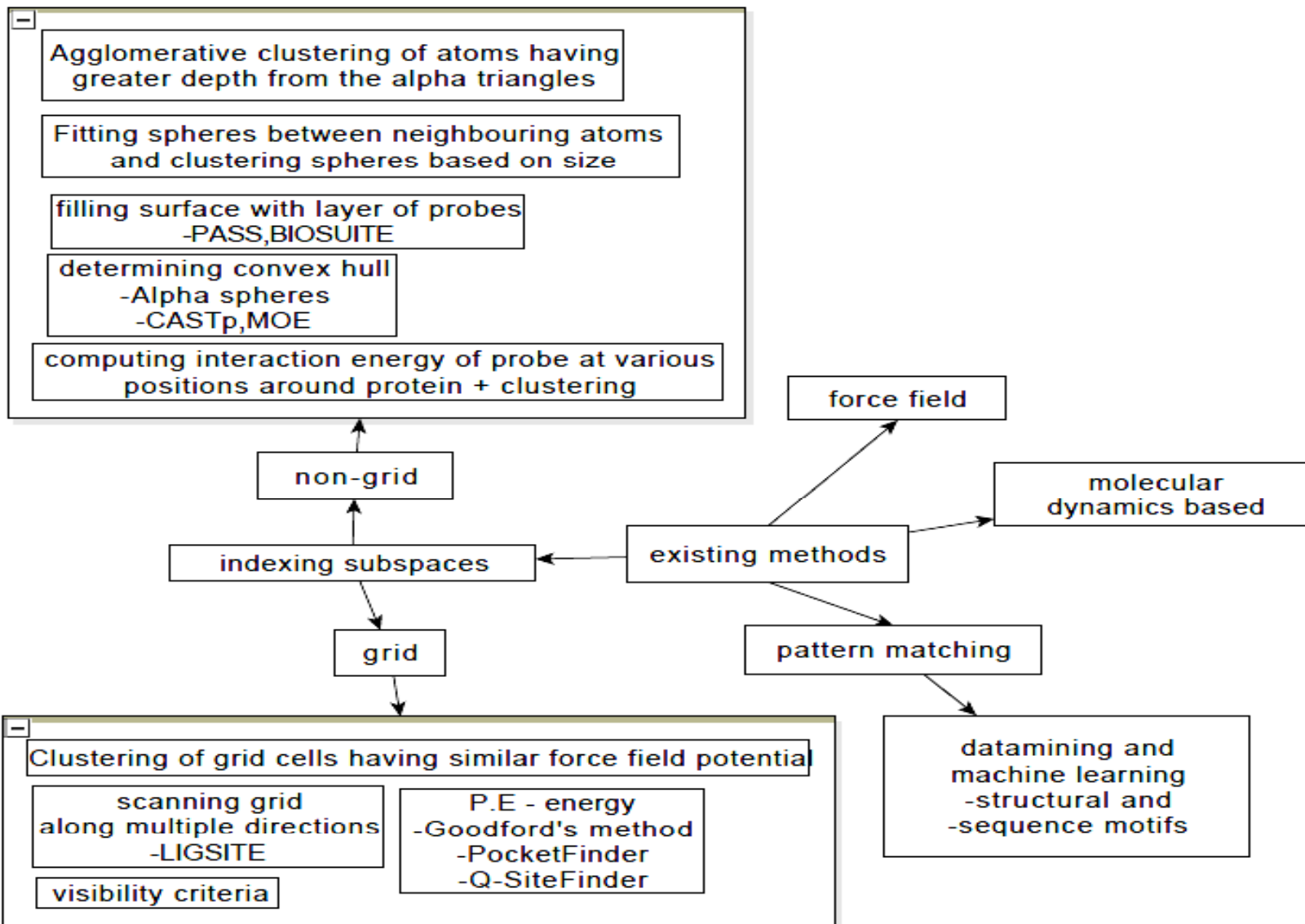
Structure-based methods

- Geometric Chemical
- Hybrid methods

Machine learning



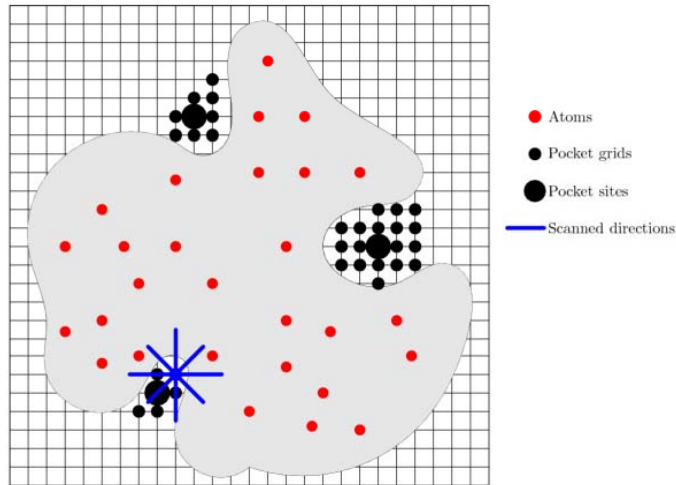
# Binding site prediction Methods



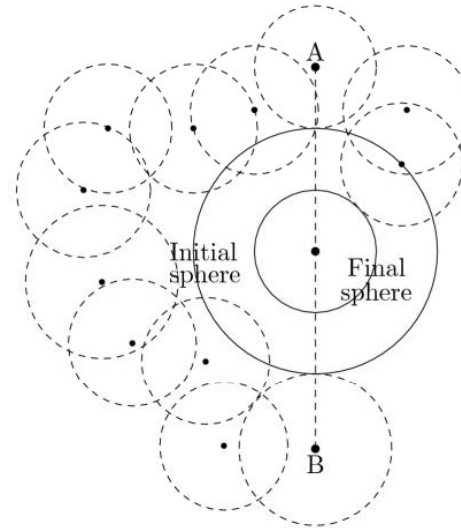
| Method                 | Algorithm  | Resource                             |
|------------------------|--|--------------------------------------|
| POCKET                 | Scan the grid along 3 dimensions   | (Levitt and Banaszak, 1992)          |
| Ligsite                | Scanning grid along 3 axes and 4 diagonals   | (Hendlich <i>et al.</i> , 1997)      |
| LigSite <sup>CSC</sup> | Similar to Ligsite but with residue conservation information for each set of residues to occur in site | (Huang and Schroeder, 2006)          |
| LigandFit              | Eraser to swipe the grid cells to demarcate cells belonging to a grove                                 | (Venkatachalam <i>et al.</i> , 2003) |
| PASS                   | Filling up surface of protein by multiple layers of probes and retaining probes with high burial count | (Brady and Stouten, 2000)            |
| CASTp                  | Fill the interatomic regions by spheres and cluster moderately sized spheres                           | (Liang <i>et al.</i> , 1998)         |
| VOIDOO                 | Similar to VOIDOO  | (Kleywegt and Jones, 1994)           |
| SURFNET                | Determine depressions on the surface of protein by placing spheres between pairs of atoms.             | (Glaser <i>et al.</i> , 2006)        |
| APROPOS                | Find clusters of atoms with depth from surface of protein  | (Peters <i>et al.</i> , 1996)        |
| Goodford's method      | Clustering of grid cells with higher energy values   | (Goodford, 1985)                     |
| Q-SiteFinder           |  | (Laurie and Jackson, 2005)           |
| PocketFinder           |  | (Jianghong <i>et al.</i> , 2005)     |

# Geometry based

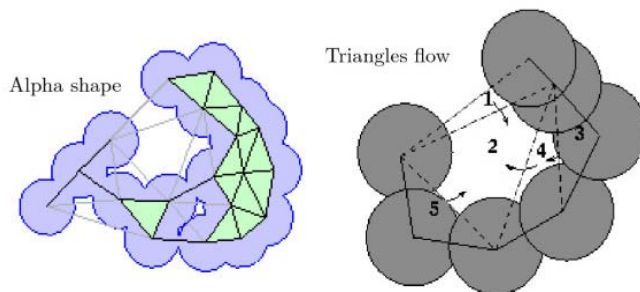
a. POCKET, LIGSITE, LIGSITE<sup>csc</sup>



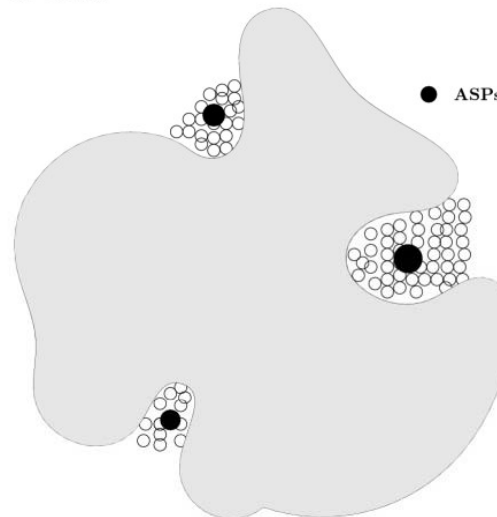
b. SURFNET



c. CAST



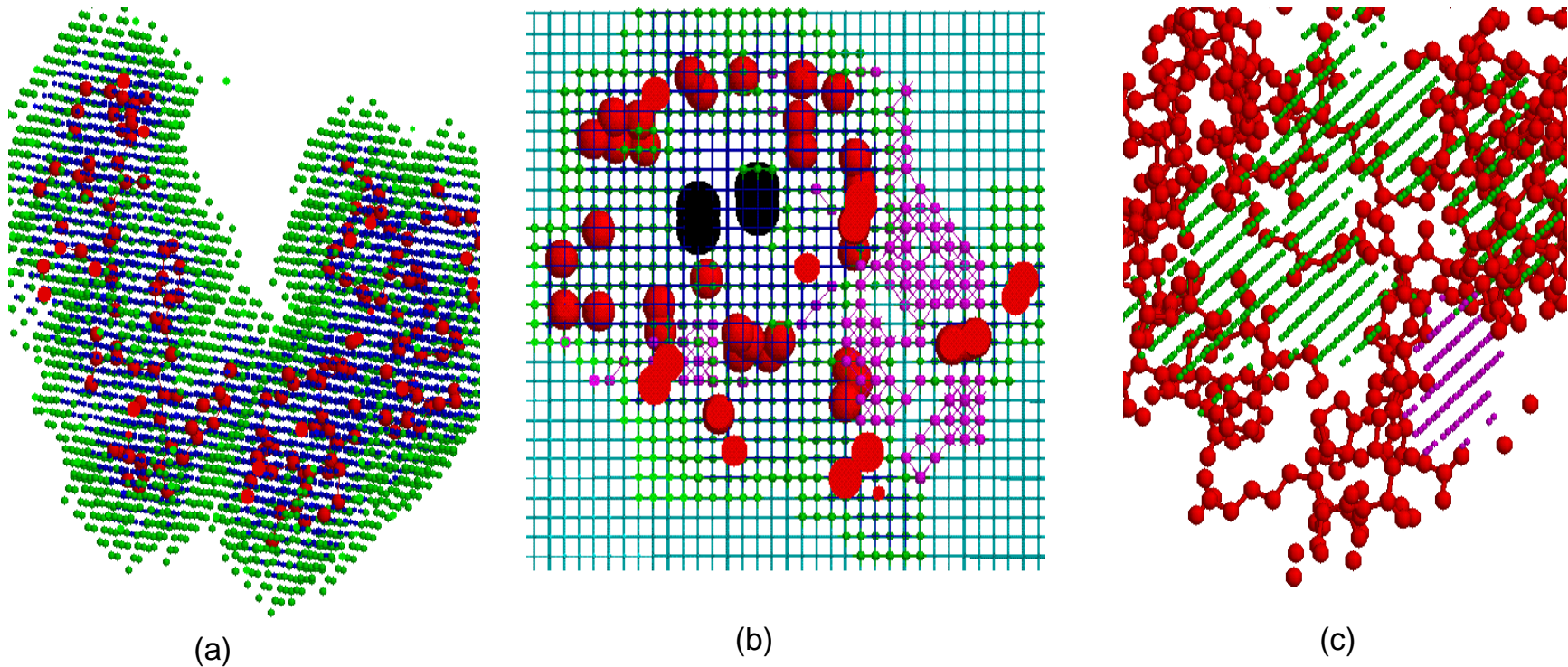
d. PASS



# PocketDepth

## Grid based binding site prediction method

Figure 1



# Grid Bar Generation

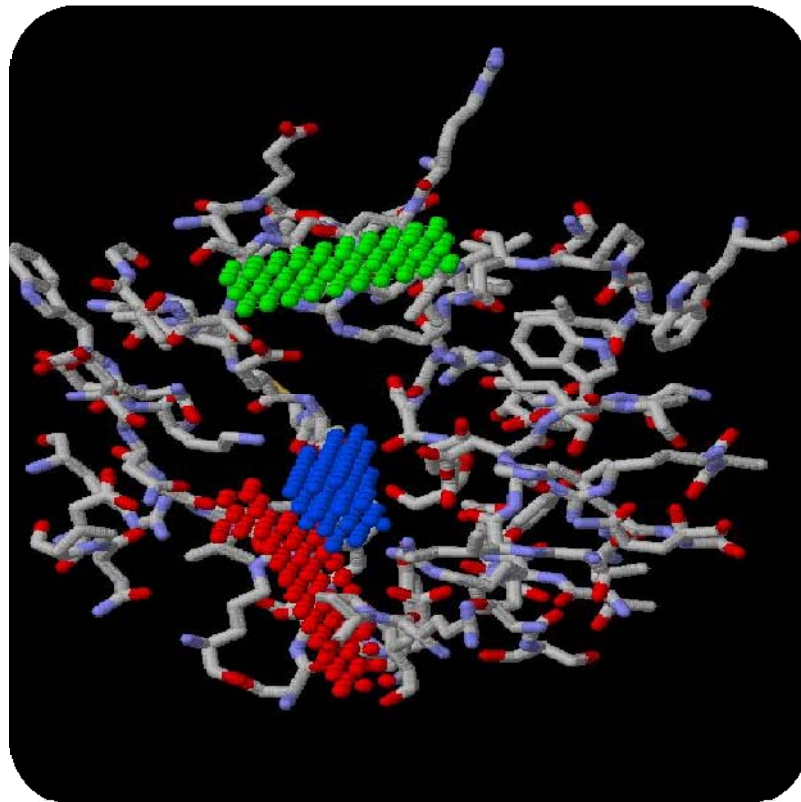
Fill inter-atomic (surface atoms) regions with grid bars.

A grid bar  $\{x, y\} \in A : GB(x, y) \subset G$  between pair of atoms  $x, y$

A grid bar is valid only if does not intersect an atom

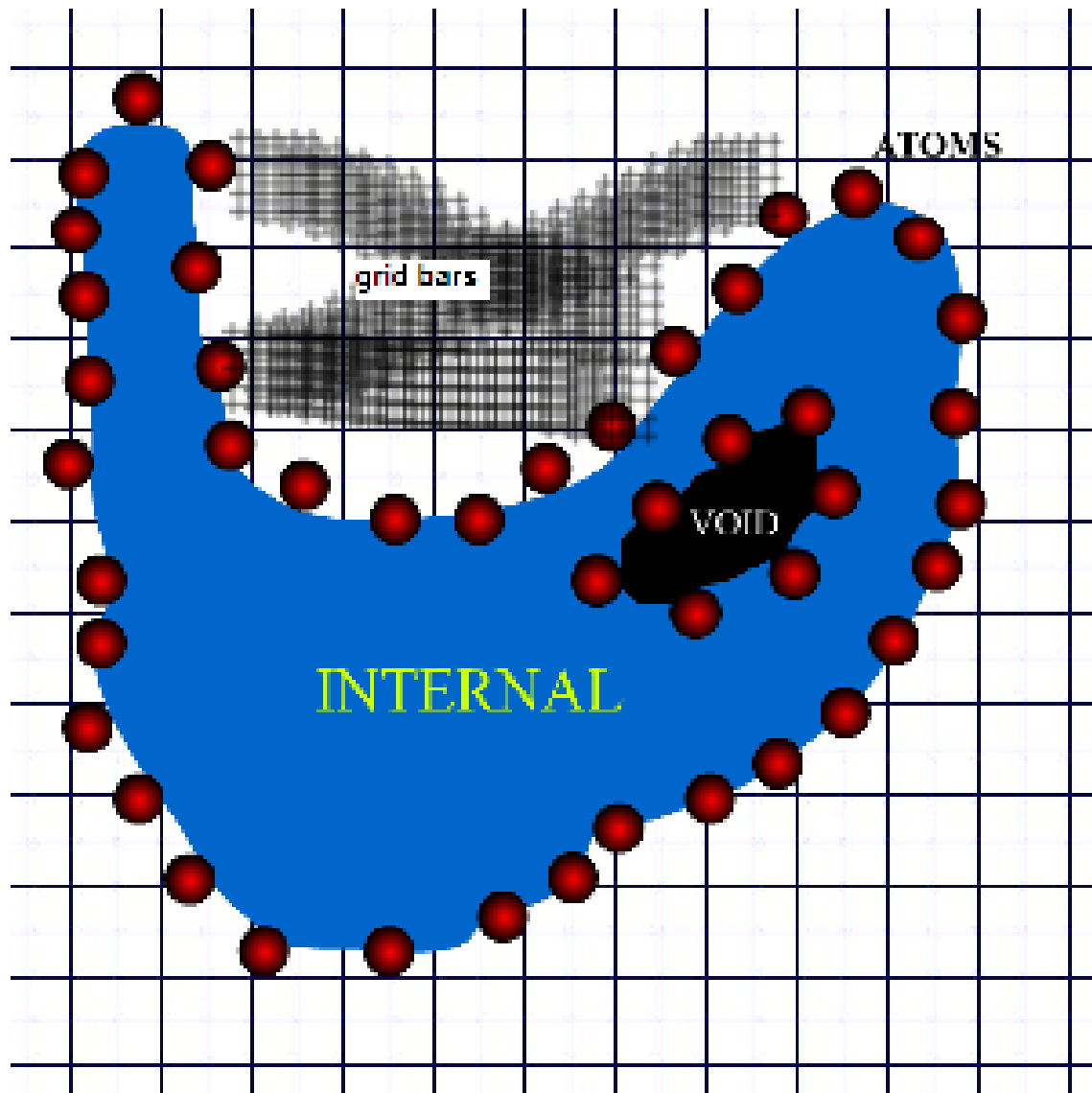
Obtain set of all valid grid bars

$\{(\forall a, b \in S)(\nexists c \in (A - S) : cell(c) \in GB(a, b))\}$





# PocketDepth



## Rendering of DepthFactor as temperature

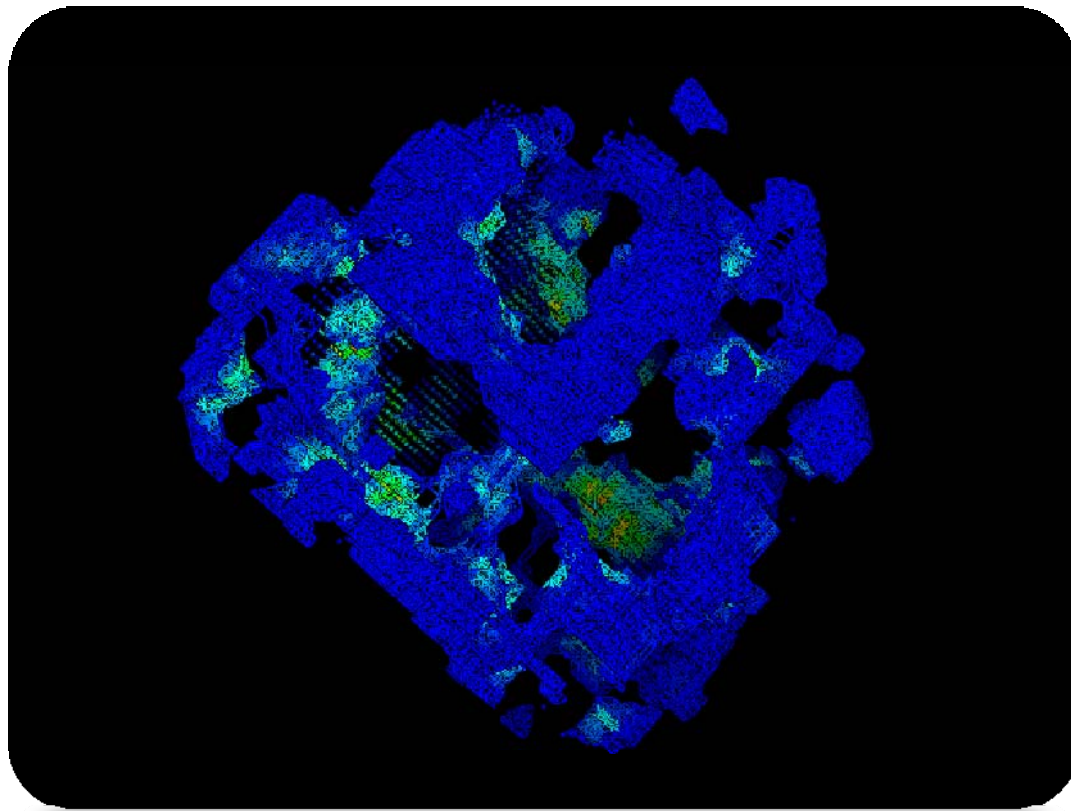
Update traversal counter, called Depth Factor, of each grid cell in a valid  $GB(a, b)$

$(\forall c \in GB(a, b)) c.depth \leftarrow c.depth + 1$

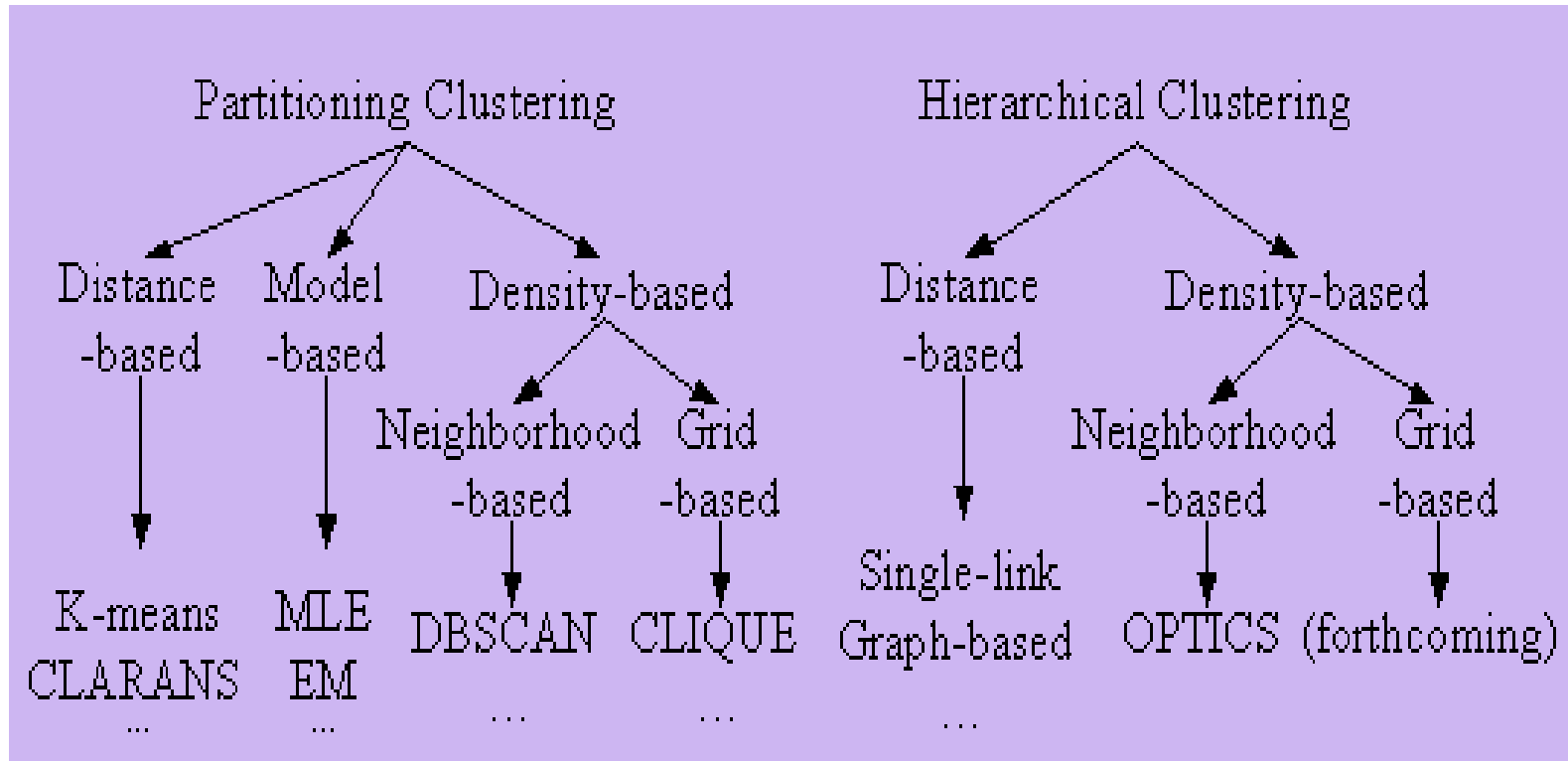
Cluster grid cells based on Depth Factor and spatial proximity (DBSCAN)

Partition the whole of the set of grid cells  $G$  into non-overlapping clusters

$S^C = C_1, \dots, C_n : C_i \cap C_j = \emptyset$  where  $S^C$  denotes a set of clusters  $1 \dots n$



# CLUSTERING METHODS



## DBSCAN clustering

FUNCTION DBSCAN(point  $p$ ,  $c$ ,  $N$ )

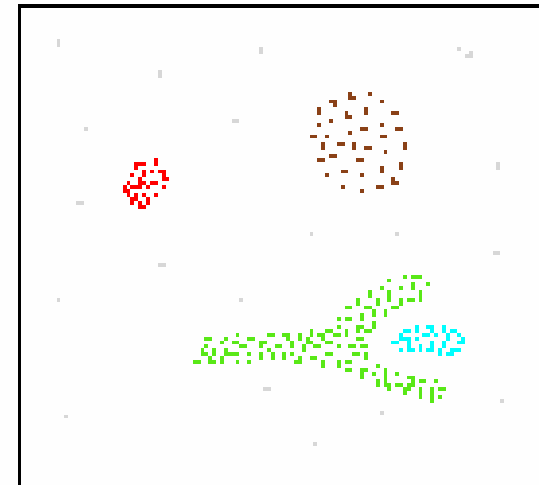
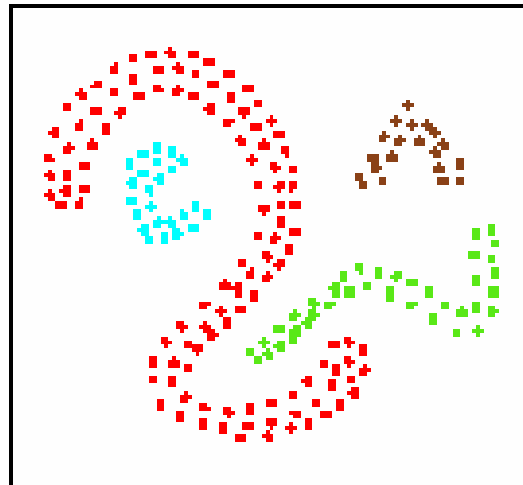
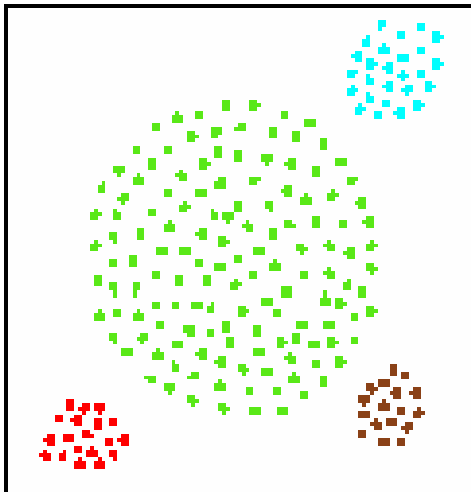
$p$  is a point and  $c$  is cluster number

**if**  $|S = \{q : d(q, p) \leq d_{threshold} \wedge q_c = \epsilon\}| \geq N$  **then**

$p_c \leftarrow c$

call DBSCAN( $q$ ) : ( $q \in S$ )

**end if**



# Clustering based on Depth Factor

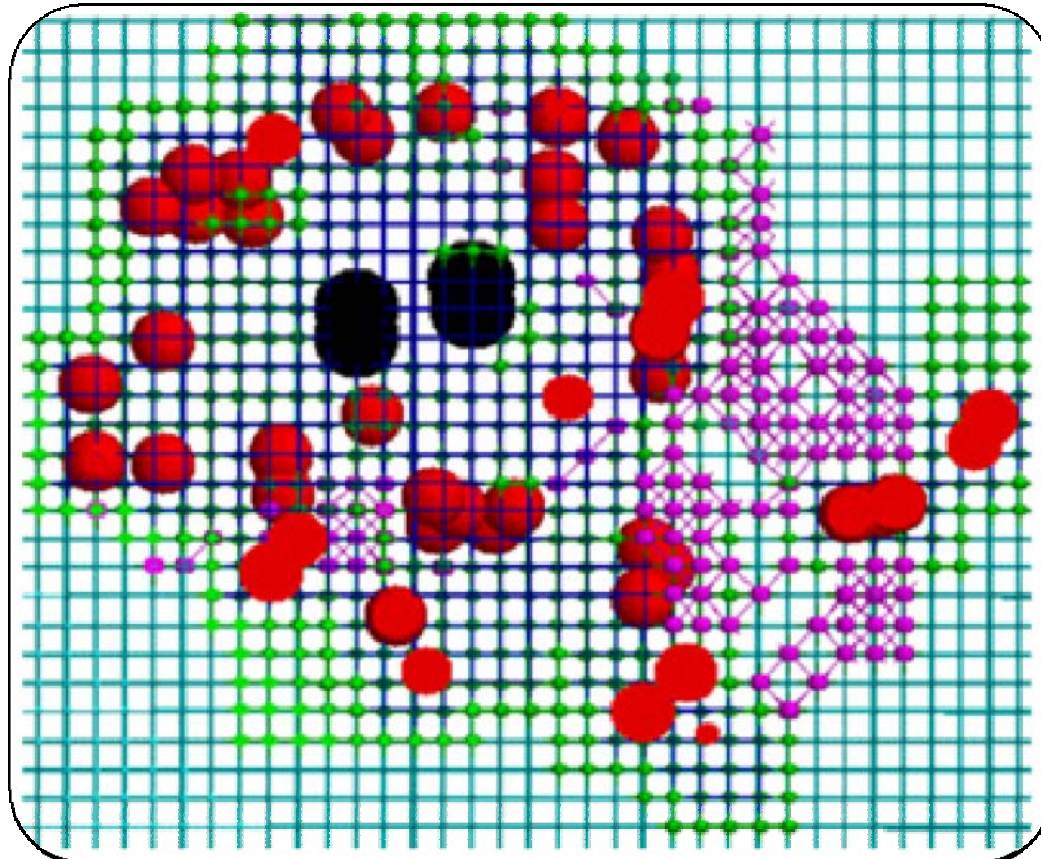
Each grid cell in a cluster  $C_i (\forall i \in [1..n])$  satisfies the depth and density requirements

$$(\forall (c \in C_i) : |\{(\forall c' \in C_i) \text{distance}(c, c') \leq \rho\}| \geq N \wedge DF(c) \geq \Delta$$

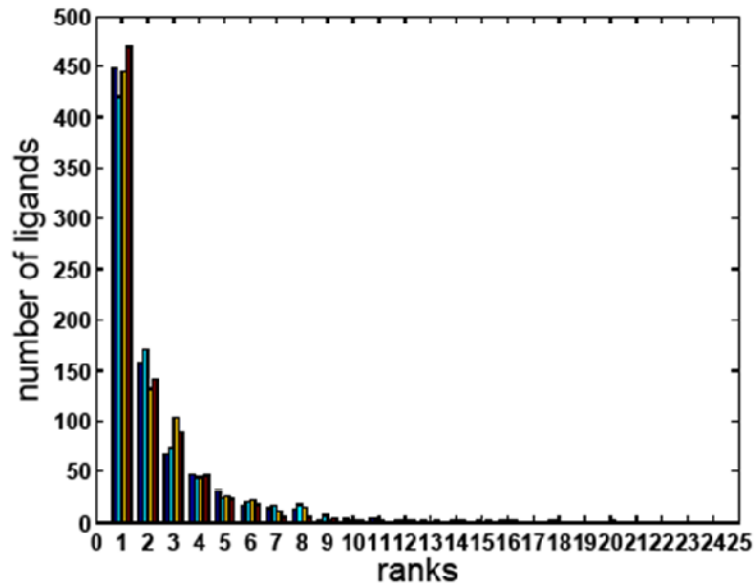
where  $\wedge$  denotes logical AND

where  $\rho, N$  are radius and number of points within radius (DBSCAN parameters);

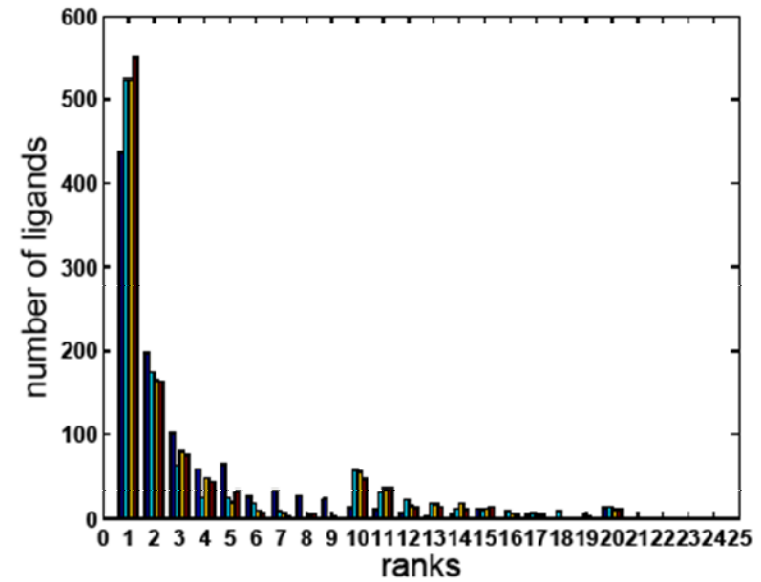
and  $DF(c)$  is the Depth Factor and  $\Delta$  is the imposed threshold



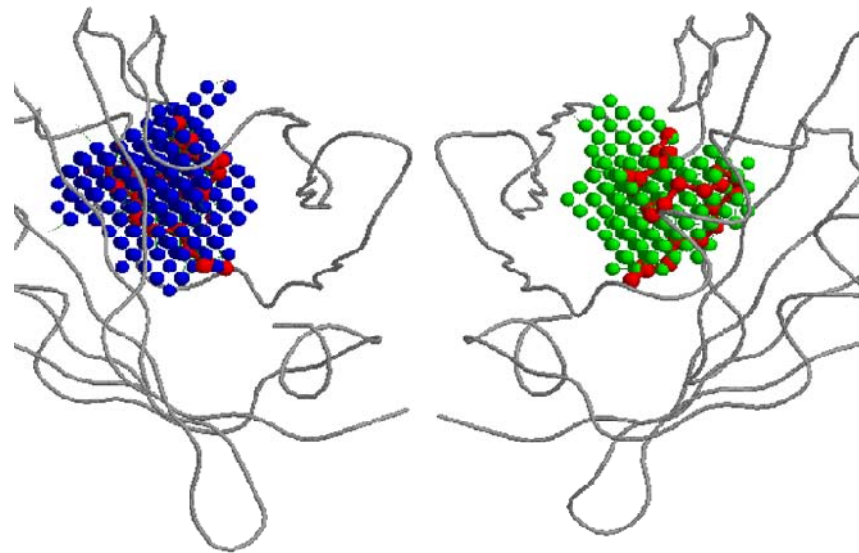
In 82% and 94% of proteins (*PDBBind* 1091) top 5 and 10 ranked clusters overlapped with crystal ligand



**(4A1)**



**(4A2)**

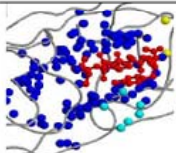
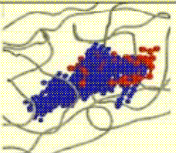
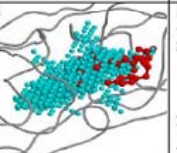
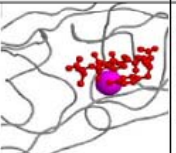
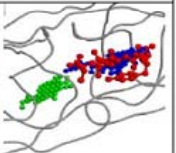

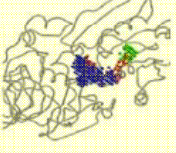
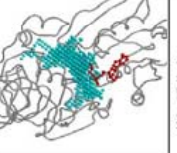

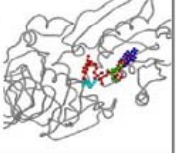
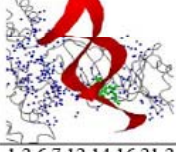
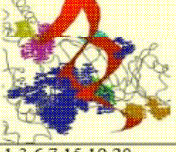
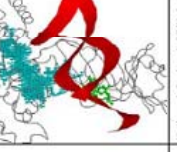
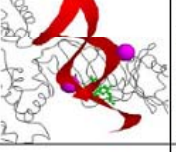
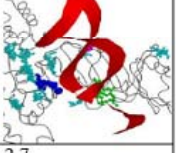
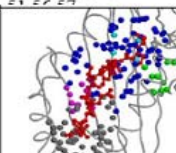
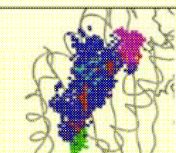
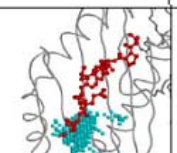
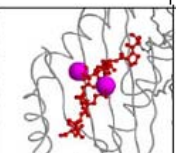
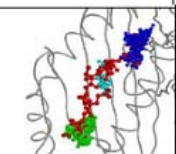
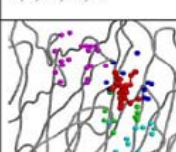
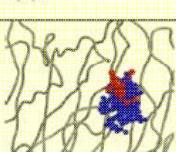
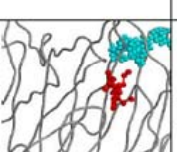
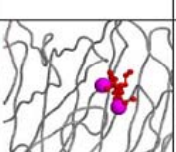
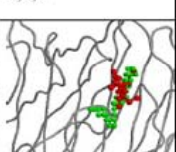
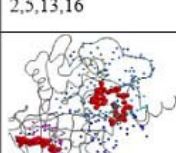
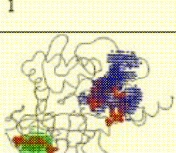
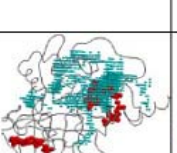
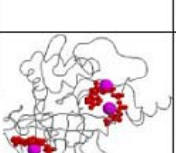
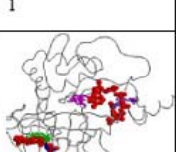


**(4B)**

In dimers  
ranks 1 and 2  
corresponded  
to ligand locations



# Binding Site Prediction Algorithms: PocketDepth Performance

| PDB   | CASTp   | PocketDepth   | LigandFit  | LigSite <sup>CSC</sup>  | QSiteFinder   |
|-------|---|---|--|---|---|
| 1sp5  |    |    |    |    |    |
| ranks | 1,5,12  | 1   | One  | One of top 10   | 1,2   |
| 1a72  |    |    |    |    |    |
| ranks | 2,11,13,17  | 1,3   |  |   | 2,7,10  |
| 1ais  |    |    |    |    |    |
| ranks | 1,2,6,7,12,14,16,21,25,28,34,39,41,51,52,57   | 1,3,6,7,15,19,20  |  |   | 2,7   |
| 1alm  |   |   |   |   |   |
| ranks | 2,3,13,14,28  | 2,3,5   |  |   | 1,2,6   |
| 2pel  |  |  |  |  |  |
| ranks | 2,5,13,16   | 1   |  |   | 1   |
| 2g88  |  |  |  |  |  |
| ranks | 1,2,3,4,7,30,48,50  | 1,2,5   |  |   | 2,5,7   |

Kalidas & Chandra, 2008  
JSB

# Binding Site Comparison

## Need..

Binding site comparison can

- Predict Important residues in a protein binding site
- Predict the function of a hypothetical protein.
- Predict the similarity between proteins.

## *Context of Binding Comparison...*

- Sequence or structural similarity && Not same molecular function
- Same function && No fold similarity

Nicola D.G & Richard M.J, (2006) J. Mol. Biol.

- Necessity of binding site comparison methods
  - Understanding protein function
  - Understanding side effects of drugs

## Challenges in site comparison

- **Point set superposition**
- Binding site → Set of points (atoms/residues)
- Determining **point-point correspondences**
- **Topology Undefined; Size of ‘match’ small & unknown** (‘Indels’ possible)
- Involves costly least squares evaluation of **rotation & translation** matrices
- **Many** possible correspondences

**Geometric Hashing; Maximal Common Sub-graph Search; Depth First Traversal** (incrementally determine correspondences)

*Comparison of a pair of binding sites  
involves three aspects:*

- (a) representation of each site as sorted lists of distances between chosen points,
- (b) alignment of two sets of distance lists and
- (c) choosing a scoring scheme for arriving at a final score



# Description of the site

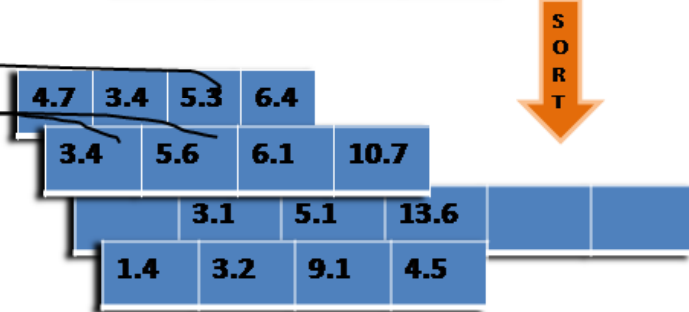
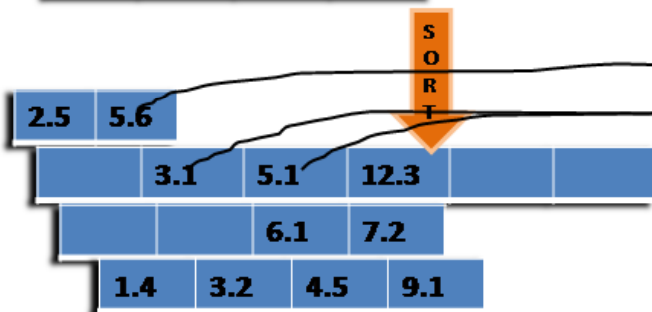
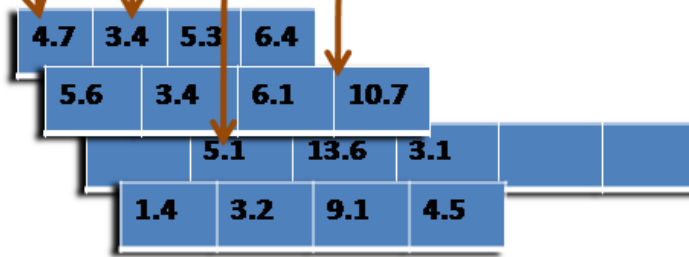
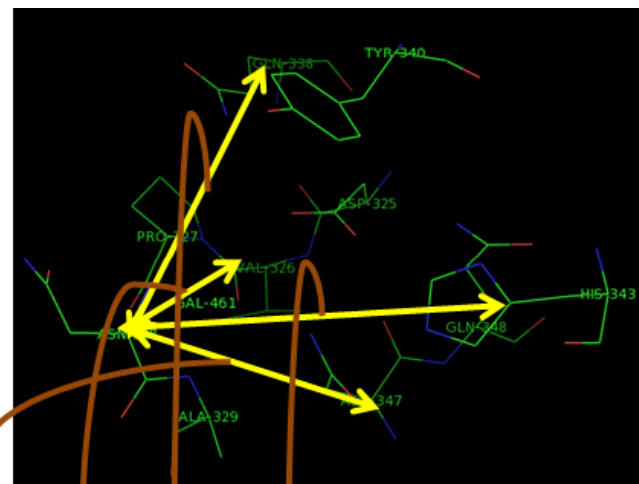
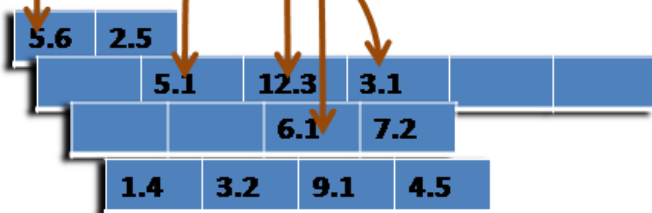
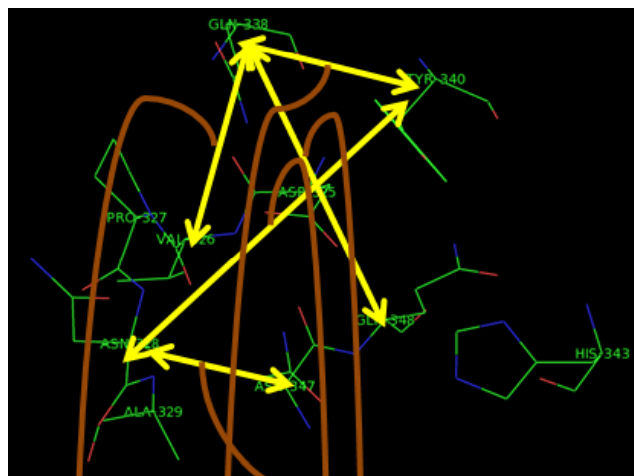
- Global features
  - Volume
  - Surface area
  - Number of polar/non-polar atoms/residues
- Shape Descriptors; Frame-invariant representations
  - Image moments
  - Spherical harmonics
  - All pair sorted distance sequences (*PocketMatch*)

(Morris et al., 2005, Bioinformatics ; Gold and Jackson, 2006, NAR; CavBase - Kuhn et al., 2007, CHEMMEDCHEM; PINTS – Stark et al., 2003, NAR; SPASM & RIGOR – Gerard et al., 1999, JMB; Binkowski et al., 2003, JMB; Morris et al., 2005, Bioinformatics; Nagano et al., 2002, JMB; Kunin et al., 2001, JMB; Campbell et al., 2003, An et al., 2005)

## Tools for Binding site comparison :

- **PocketMatch-** A new algorithm to compare binding sites in protein structures
- **CavBase**
- **SitesBase-** a database for structure-based protein–ligand binding site comparisons
- **CPASS** - Comparison of Protein Active-Site Structures
- **PINTS-** Patterns in Non-homologous Tertiary Structures
- **Spasm/RIGOR**
- **SMAP-WS Pairwise Comparison /SMAP-WS Database Search**
- **SiteSorter™** -N-by-N Binding Site Similarity Assessment
- **SLiC** -Site-Ligand Contact Analysis and Binding Mode Similarity Assessment
- **MAPPIS**(Multiple Alignment of Protein-Protein InterfaceS (PPIs))- Recognizes spatially conserved chemical interactions shared by a set of PPIs
- **MULTIBIND**(Multiple Alignment of Protein **B**inding Sites)- Recognizes Spatial Chemical Binding Patterns Common to a Set of Protein Structure

# PocketMatch Algorithm

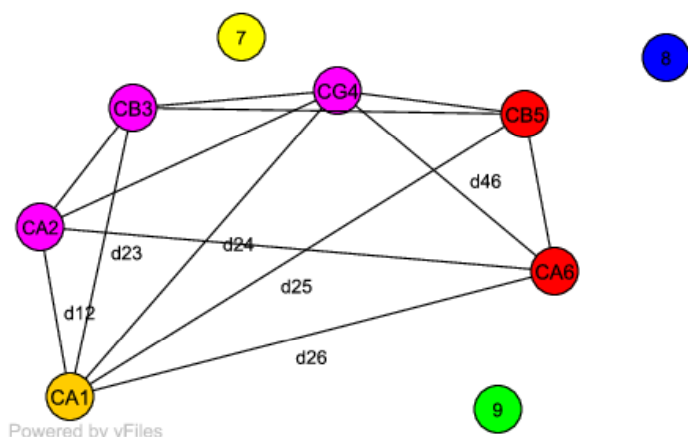


Number of matching distance elements

$$PMscore = \frac{\sum_{i=1}^{90} Count_i}{\text{maximum}(|S_1|, |S_2|)}$$

# PocketMatch Algorithm

<http://proline.physics.iisc.ernet.in/pocketmatch/>



Powered by yFiles

3 types of points (CA, CB, CNTR)

5 types of residue groups (AVILGP; KRH; DE; YFW; CSTQN)

$(3 * (3-1) / 2 + 3) * (5 * (5-1) / 2 + 5) \rightarrow 90$  lists

120 lists are possible and yielded similar results.

Representation of the binding site :

$$\begin{array}{l}
 \text{NGP} \\
 \text{NTP} \\
 \left\{ \begin{array}{l} ND_1 \\ d_1, \quad d_2, \quad \dots \quad d_j, \quad \dots \end{array} \right\} \\
 \left\{ \begin{array}{l} ND_2 \\ d_1, \quad d_2, \quad \dots \quad d_j, \quad \dots \end{array} \right\} \\
 \vdots \\
 \left\{ \begin{array}{l} ND_{90} \\ d_1, \quad d_2, \quad \dots \quad d_j, \quad \dots \end{array} \right\}
 \end{array}$$

Where, NGP : Number of pairs of group-types, NTP : Number of pairs to point-types,  $ND_i$  : Number of distances in the  $i^{\text{th}}$  bin,  $d_j$  : distance between  $j^{\text{th}}$  pair of points.

---

*Sub-routine 1* Alignment of a pair of sorted distance sequences

---

```

i=0; j=0; counter=0;
while (i ≤ m) ∧ (j ≤ n) do
  if |S1[i] - S2[j]| ≤ τ then
    i ← i + 1; j ← j + 1
    counter ← counter + 1
  else
    if S1[i] < S2[j] then
      i ← i + 1;
    else
      j ← j + 1;
    end if
  end if
end while

```

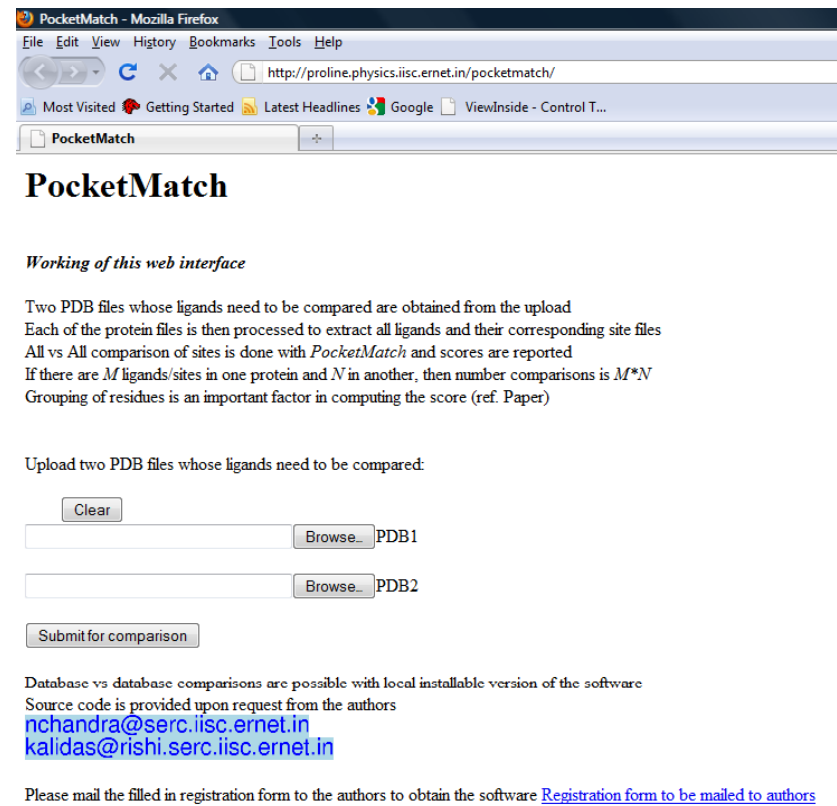
---

PocketMatch: A new algorithm to compare binding sites in protein structures

Kalidas Yeturu and Nagasuma Chandra, 2008, BMC Bioinformatics

# *PocketMatch* implementation

- Sites extracted → around (4Å) ligand/predicted pocket
- Complete residues → representative points → Sorted Distance lists
- MPI version (C language)
- Run on **IBM Bluegene** utilizing 1024 processors



The screenshot shows the PocketMatch web interface in a Mozilla Firefox browser window. The address bar displays the URL <http://proline.physics.iisc.ernet.in/pocketmatch/>. The page title is "PocketMatch".

## PocketMatch

*Working of this web interface*

Two PDB files whose ligands need to be compared are obtained from the upload. Each of the protein files is then processed to extract all ligands and their corresponding site files. All vs All comparison of sites is done with *PocketMatch* and scores are reported. If there are  $M$  ligands/sites in one protein and  $N$  in another, then number comparisons is  $M*N$ . Grouping of residues is an important factor in computing the score (ref. Paper).

Upload two PDB files whose ligands need to be compared:

PDB1

PDB2

Database vs database comparisons are possible with local installable version of the software. Source code is provided upon request from the authors.  
[nchandra@serc.iisc.ernet.in](mailto:nchandra@serc.iisc.ernet.in)  
[kalidas@rishi.serc.iisc.ernet.in](mailto:kalidas@rishi.serc.iisc.ernet.in)

Please mail the filled in registration form to the authors to obtain the software [Registration form to be mailed to authors](#)

# Perturbation studies

## Validation with respect to random perturbation of positions of site-points

Random perturbations of site points for (a) ligand(PP8) with 54 PMScores for perturbed sites with respect to its original site for different extents of perturbations(RMSD) are shown at different values of (1.0-green,0.5-red,0.25-cyan,0.125-blue,0.01-yellow)

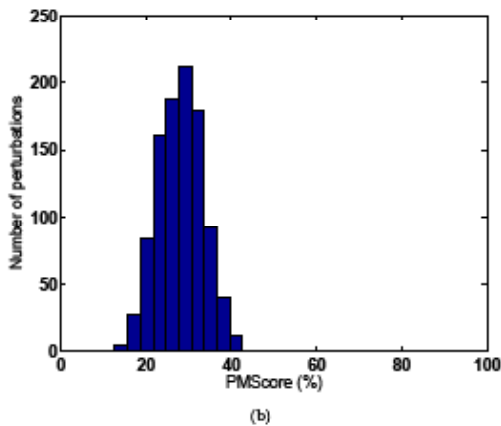
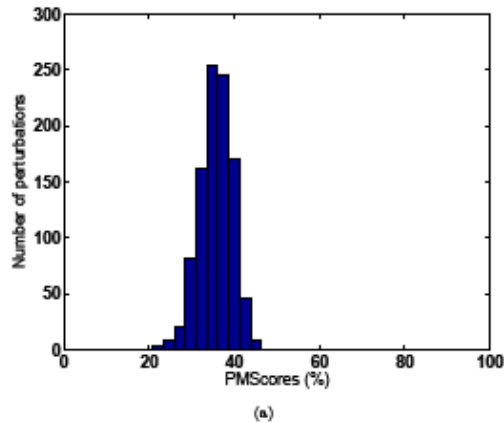
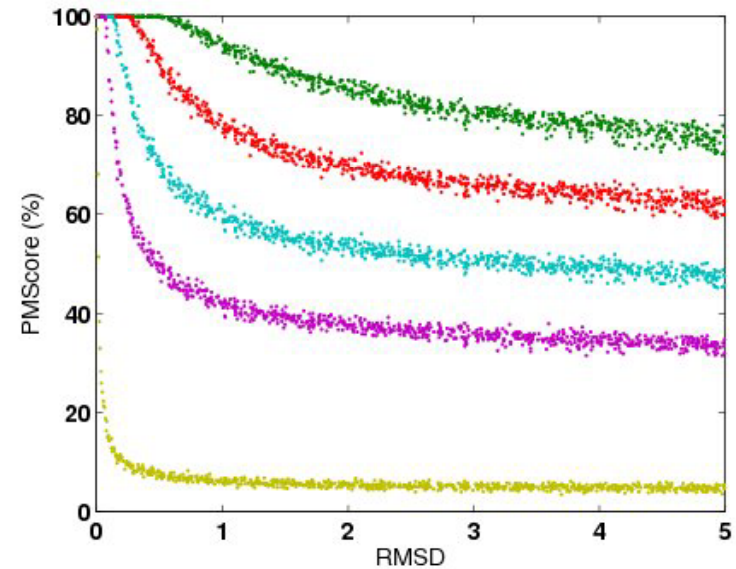


Figure 6:

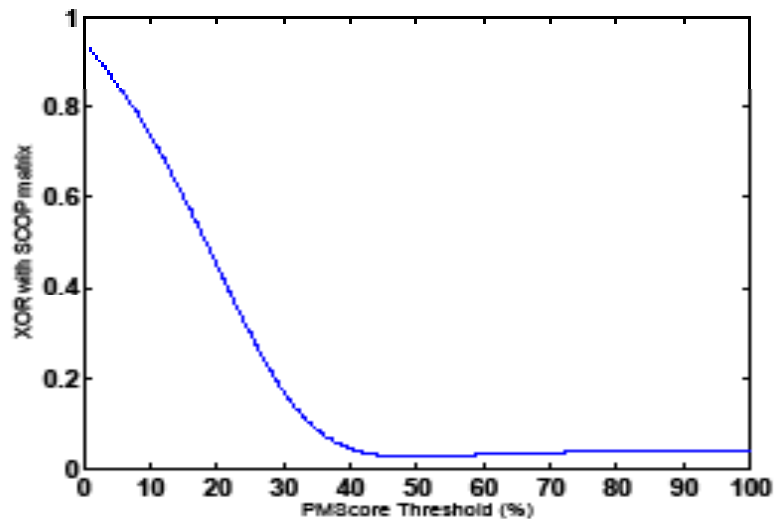
### Type perturbation

Superposition of sites with (a) High PMScores (80.9% for 1H8H-ATP and 1W0K-ADP) and (b) low PMScores (25.8% for 1H8H-ATP and 1H8H-ADP)

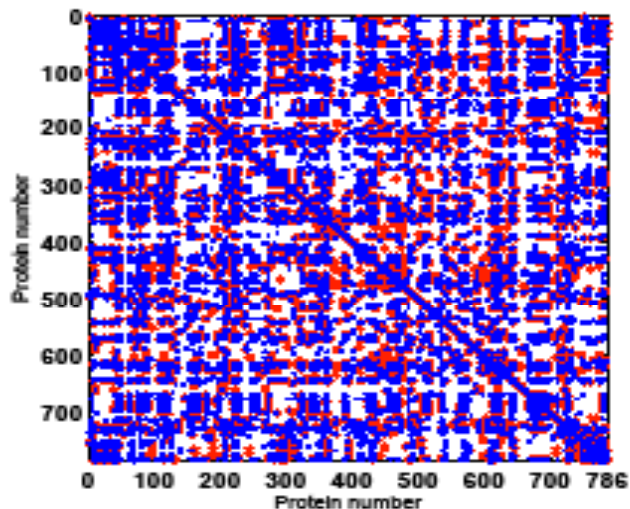




# Validation

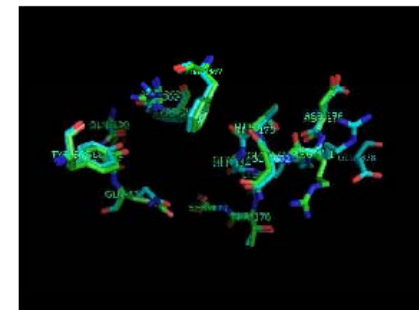


(a)

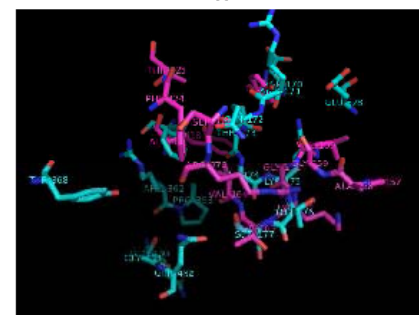


(b)

**SCOP VS PM**

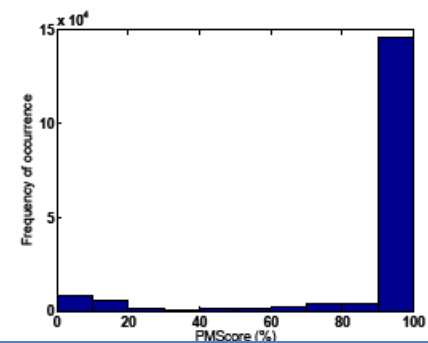


(a)



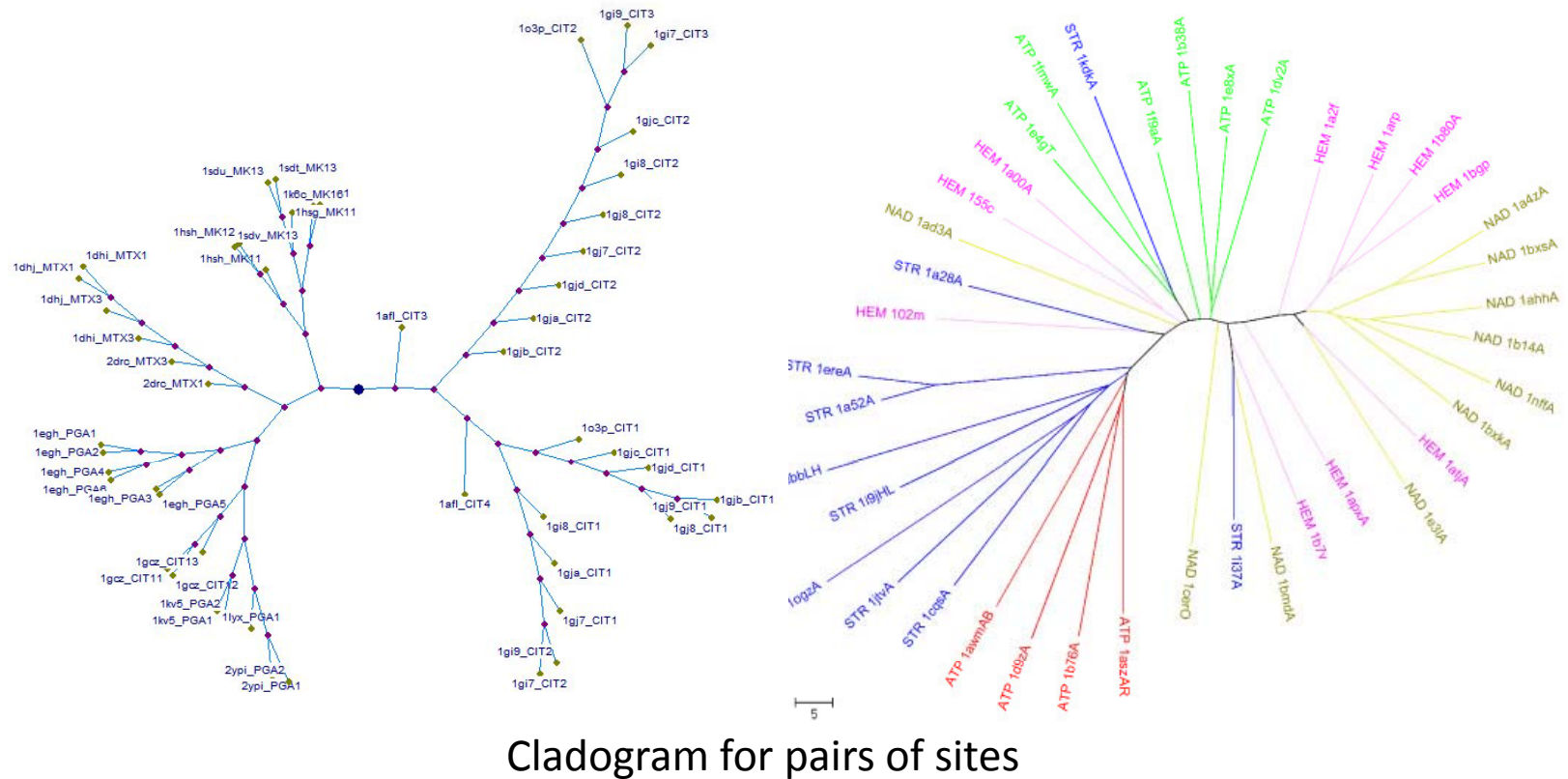
(b)

**ATP-ADP Similar and dissimilar sites**



**Known similarities in tetramers**

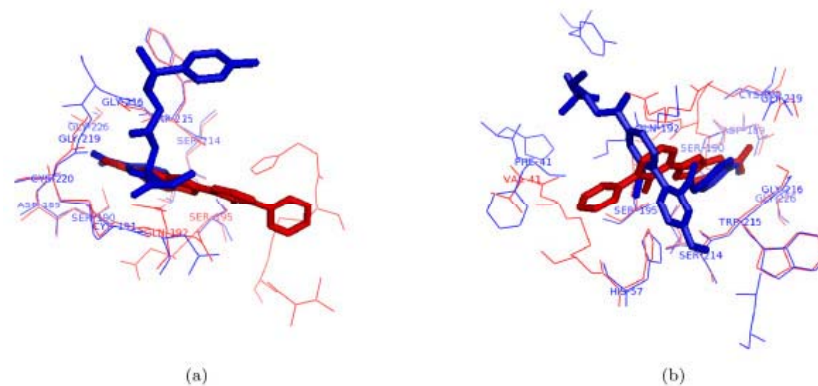
# Cladogram based validation



Cladogram for pairs of sites

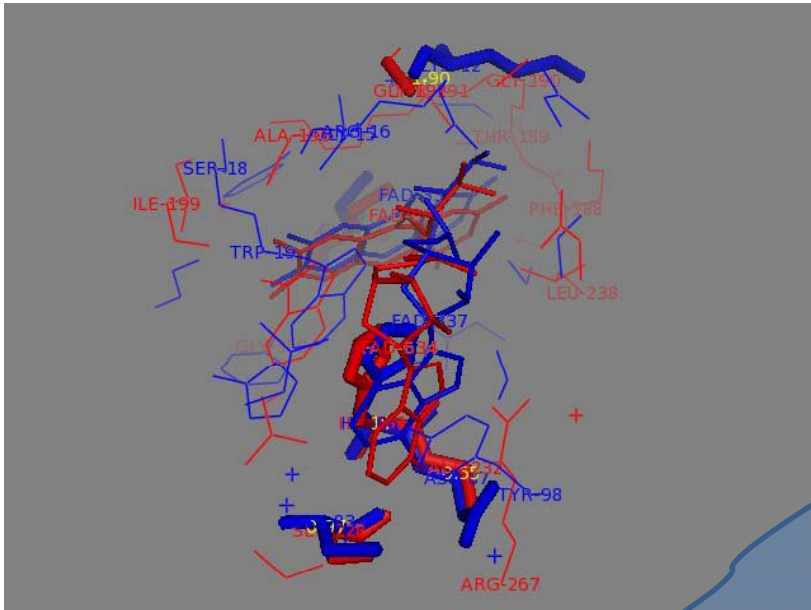
## Detection of part-similarities by PocketMatch.

Examples illustrating binding of different ligands in essentially the same binding pocket, but with different orientations. The part-similarities in these were identified correctly by PocketMatch. Binding of different trypsin inhibitors (stick models) complexed to trypsin variants (wire) as in PDB entries (a) 1GJC and 1V2Q and (b) 1GJC and 2AYW.

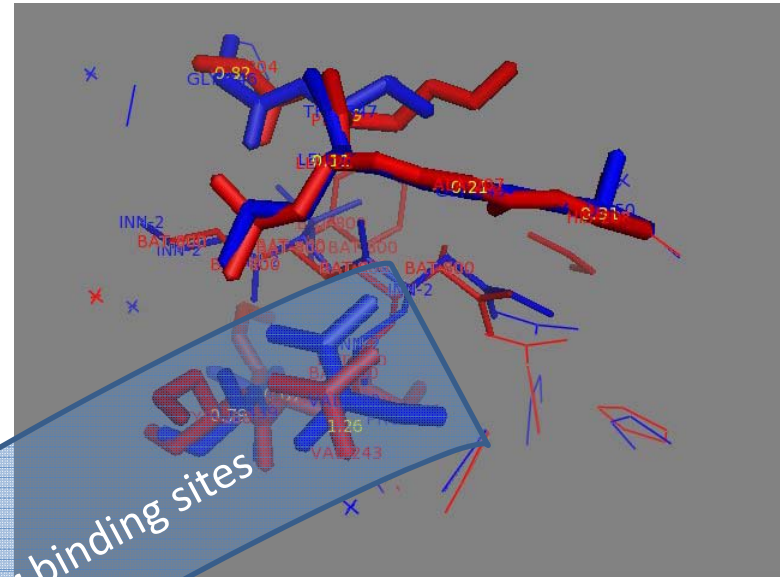


# Need for site alignment

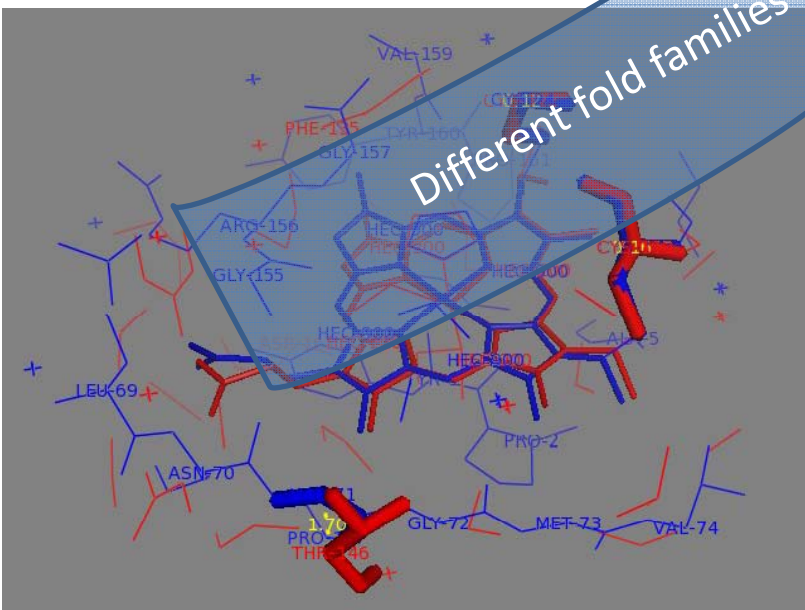
- Different folds exhibit similar binding sites
  - Ex – cofactor binding sites HEM, NAD, FAD
- Difficult to detect local similarities by human – error prone
- Structural motifs – determining function



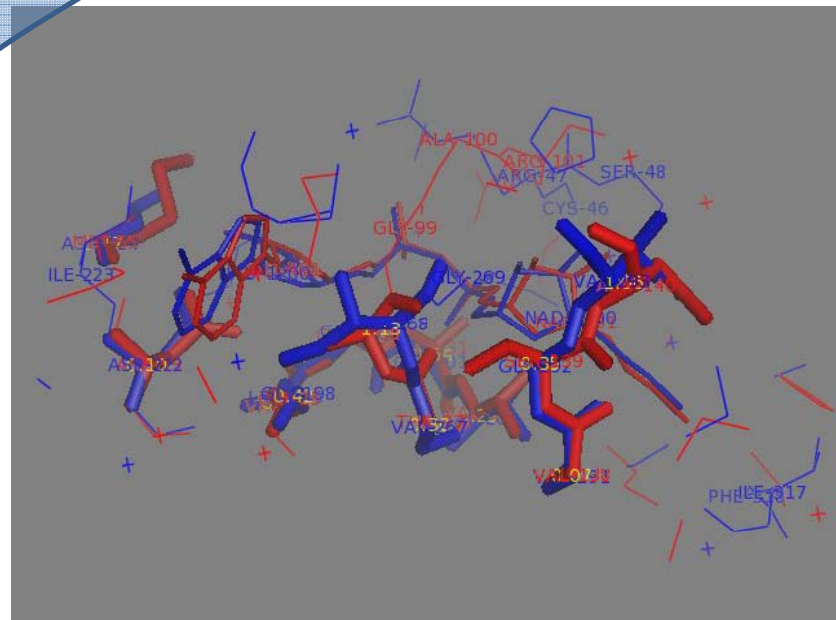
pdb1rp4.ent\_FAD-pdb1jra.ent\_FAD-pa



pdb1rm8.ent\_BAT-pdb1bkc.ent\_INN-pa



pdb1m6z.ent\_HEC-pdb1e2z.ent\_HEC-pa



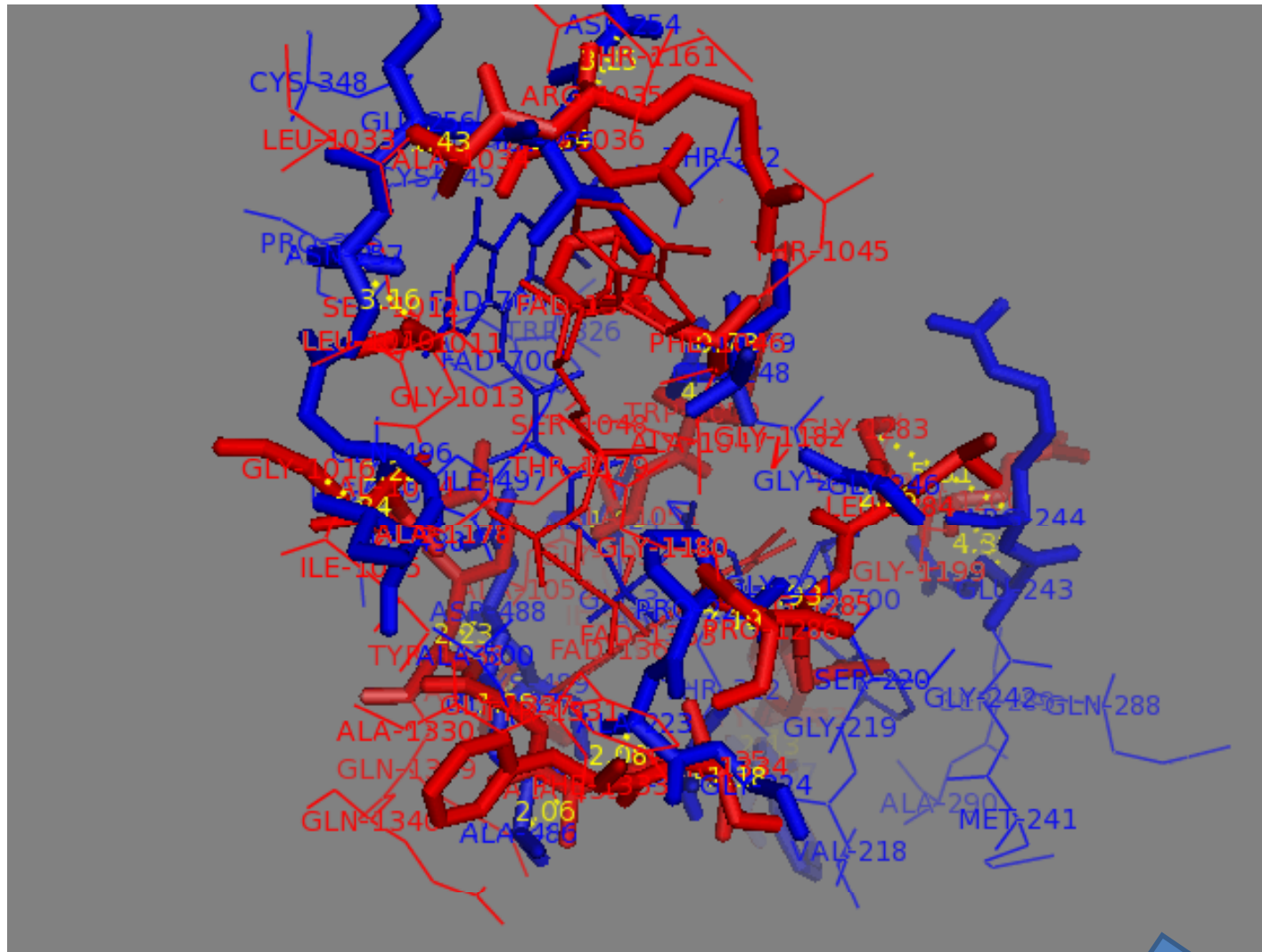
pdb9ldt.ent\_NAD-pdb1ee2.ent\_NAD-pa

Different fold families — similar binding sites

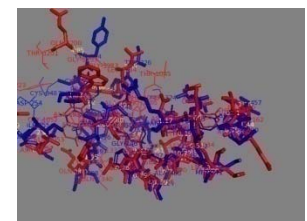
## *Challenges in alignment*

- Many possible local similarities exist
  - Exhaustive enumeration is impractical
- Finding out the best is tough
  - Quantify when does an expert call superposition 'good'
- What level to consider structural match
  - Atomic, Residue, C-alpha





Extract



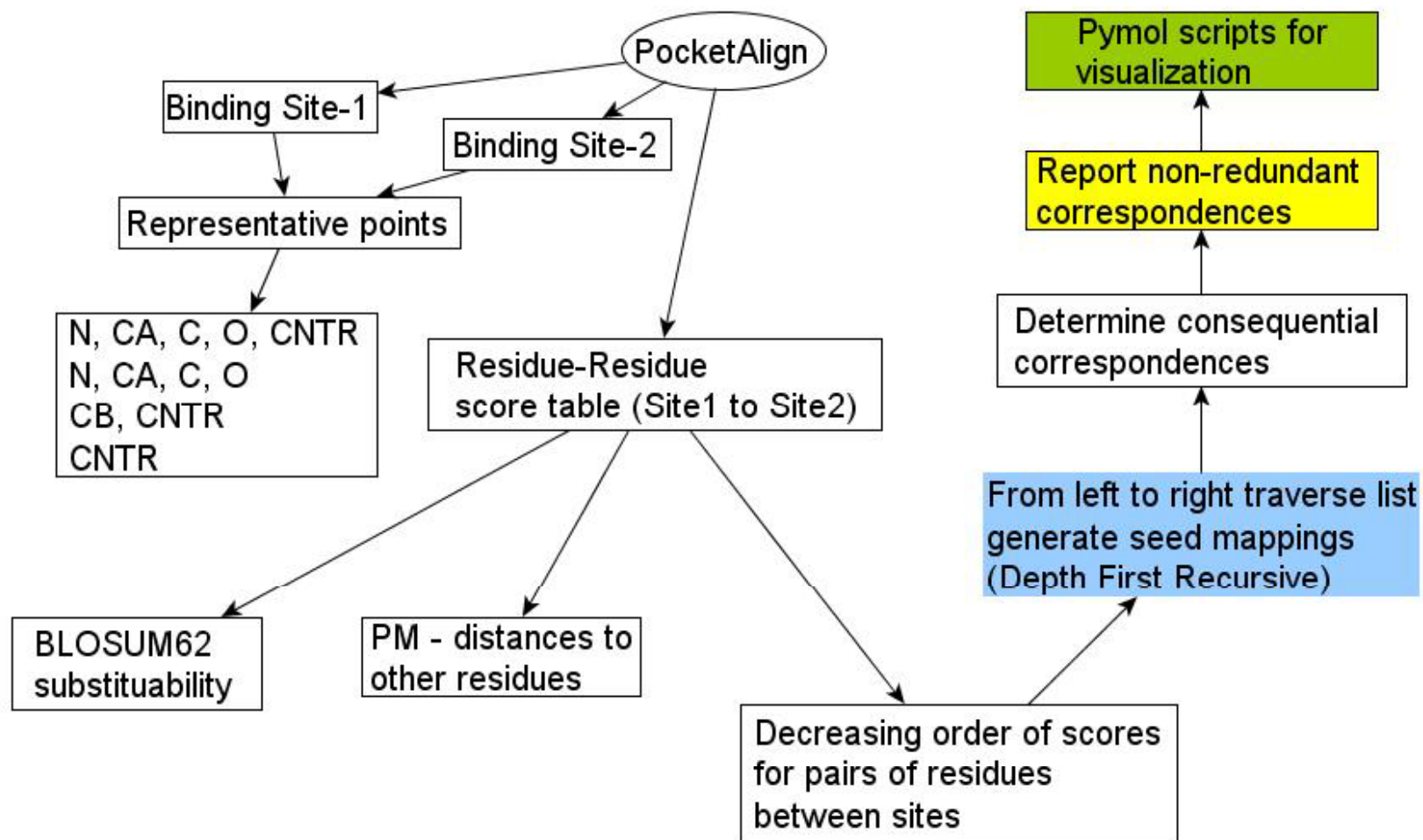
Superpositions of FAD binding sites  
1COP and 1HYU



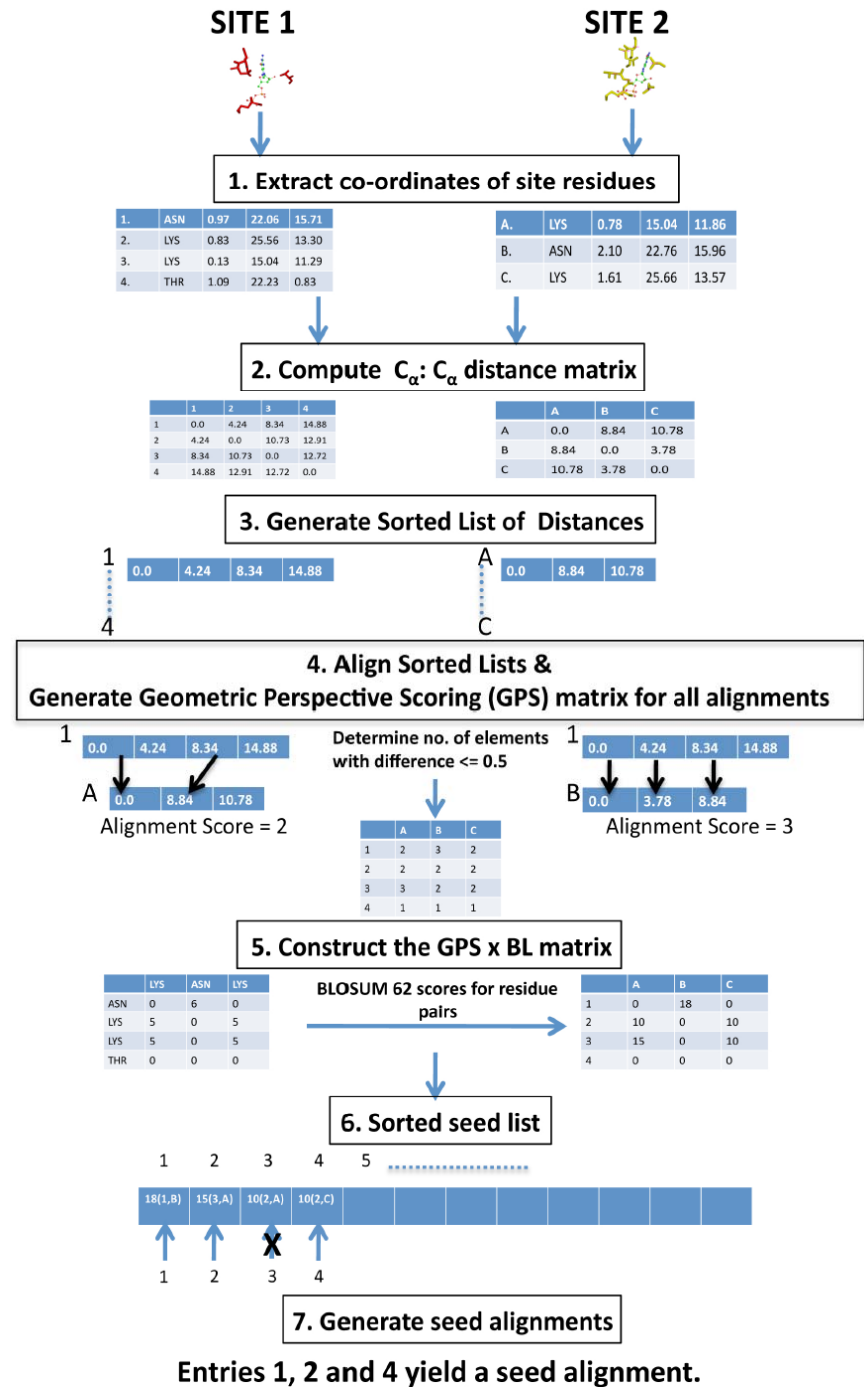
**POCKETALIGN**

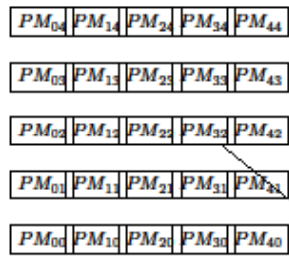
# PocketAlign

## Alignment of binding sites

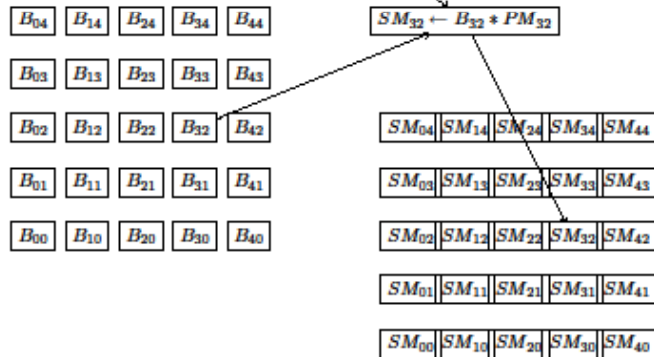


# Schematic





- Scores between residue pairs
- Descending order sorted pair-scores
- Selection of top pair from left moving right on the string




---

### Algorithm 1 Generation of seed alignments

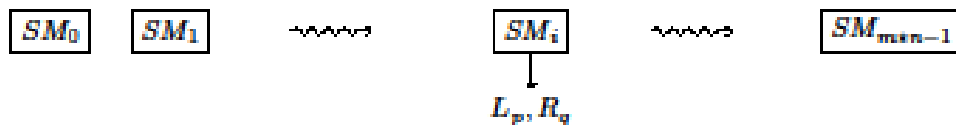
---

```

for  $i = 1$  to  $m * n$  do
  for  $j = i$  to  $m * n$  do
     $p \leftarrow SMM[j].residue[0]$ 
     $q \leftarrow SMM[j].residue[1]$ 
    if  $p$  and  $q$  are not already mapped then
      if RMSD criteria is met for  $map \cup \langle p, q \rangle$  then
        {Above check considers current alignment type}
         $map \leftarrow map \cup \langle p, q \rangle$ 
      end if
    end if
  end for {#j}
  Update database with  $map$ 
end for {#i}

```

---



Two binding sites are represented as sets of residues

$S = \{R_1 \dots R_m\}$  where  $R_i$  is  $i^{th}$  residue of first site

Each residue defines a partitioning of the set of atoms,  $A$

$R = \{a \in A\} \subset A$

$R_i \cap R_j = \emptyset (\forall i \neq j \in |S|)$

Where  $|S|$  denotes cardinality of the set,  $S$

Similarly second site is represented by  $S' = \{R'_1 \dots R'_n\}$  on set of atoms,  $A'$

Chemical similarities are denoted by a function  $BL : S \times S' \rightarrow N$

Geometric similarities (GPS) are denoted by  $GPS : S \times S' \rightarrow N$

A combination scoring scheme is defined  $GPS \times BL_{ij} \rightarrow GPS_{ij} * BL_{ij}$

A linearization of  $GPS \times BL$  is performed

A one-to-one function is defined  $L : [1 \dots m] \times [1 \dots n] \rightarrow [1 \dots m * n]$

SeedList is created by obtaining values from  $GPS \times BL$

$SeedList_{L(i,j)}^V \leftarrow GPS \times BL_{ij}$  for storing the values

$SeedList_{L(i,j)}^P \leftarrow (i, j)$  for storing the residue pairs

SeedList is sorted such that  $(\forall p \leq q) SeedList_p^V \geq SeedList_q^V$

A mapping is defined as residuewise correspondences between the two sites

A one-to-one function, for a mapping  $M : [1 \dots m] \rightarrow [1 \dots n]$

Seed mapping or alignment  $B$  is derived by traversal of SeedList

$B \leftarrow \{(p, q)\} \subset SeedList^P$

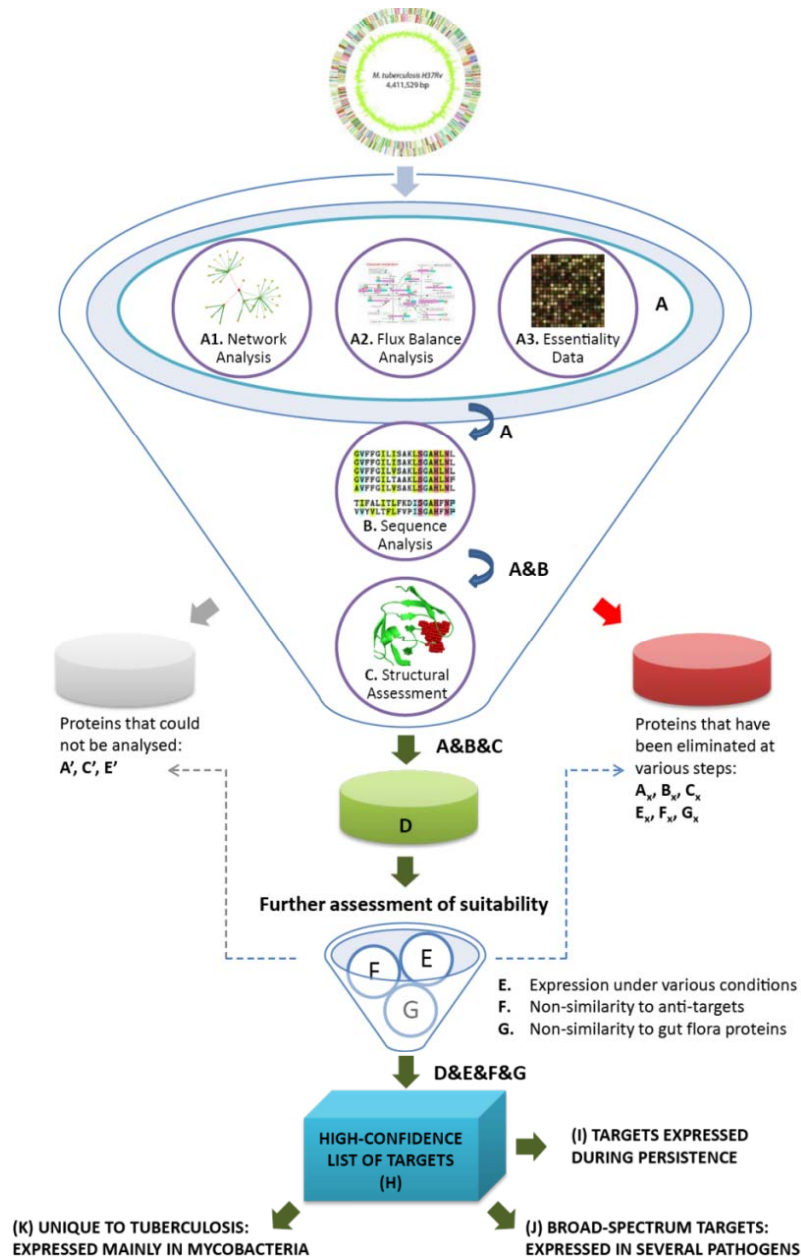
# PocketAlign (Validation & New Results)

- Ran for a set of 34 pairs of sites known to be similar
- Encouraging results obtained from a set of 143 pairs of histamines, 29 pairs of lectins, 209 pairs of sites of carbohydrate (GAL, GLC and MAN) sites and ATP binding sites





# targetTB – Target Identification Pipeline



## (A&B) Systems and Sequence Level Filters

- A1 Node deletions on STRING + Metabolic Influences network
- A2 Essential genes from *Mtb* *inj661*, GSMN-TB
- A3 High-throughput Transposon Site Hybridisation (TraSH) Mutagenesis study
- B Eliminated proteins with close homologues in human proteome

## (C) Structural Assessment of Targetability

- Binding site prediction and comparison – *Mtb* vs. *Hsa*
- Structural models obtained from ModBase
- Binding sites identified using PocketDepth and compared using PocketMatch (cut-off: 0.80) (A&B&C ⇒ D)

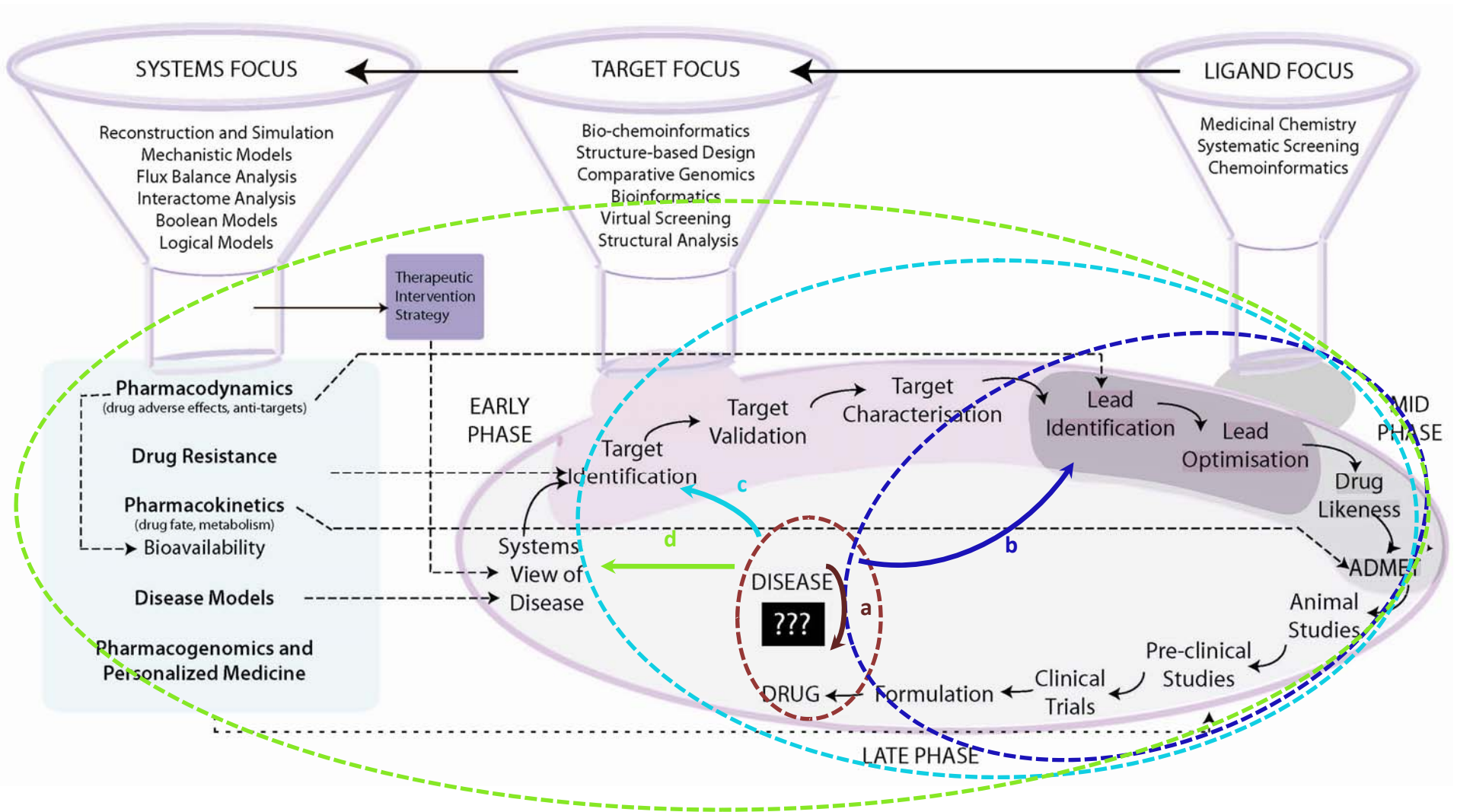
## Other Filters (applied to (D))

- E Expression of target (Microarray data)
- F Non-similarity to human 'anti-targets'
- G Non-similarity to gut flora proteins
- Paths to resistance mechanisms

## Multiple Lists of Targets

- H Passing filters A–G
- I H-List targets upregulated in persistence
- J H-List targets that can serve as broad-spectrum targets
- K H-List targets, unique to *Mtb*.

# THE NEW DRUG DISCOVERY PIPELINE







Thank You