

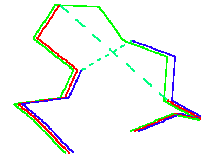
Structural Alignment of Proteins

Thomas Funkhouser
Princeton University
CS597A, Fall 2007

Goal

Align protein structures

```
1 2 3 4 5 6 7 8 9 10 11 12 13 14
PHE ASP ILE CYS ARG LEU PRO GLY SER ALA GLU ALA VAL CYS
PHE ASN VAL CYS ARG THR PRO --- --- --- GLU ALA ILE CYS
PHE ASN VAL CYS ARG --- --- --- THR PRO GLU ALA ILE CYS
```



[Marian Novotny]

Terminology

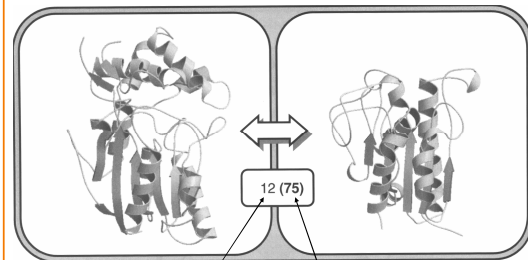
Superposition

- Given correspondences, compute optimal alignment transformation, and compute alignment score

Alignment

- Find correspondences, and then superpose structures

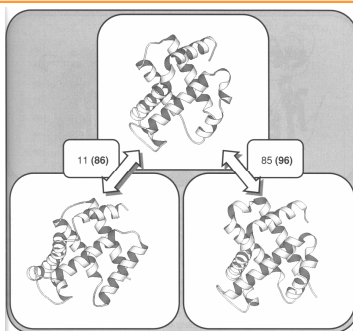
Structure vs. Sequence



Sequence Identity (Structure similarity)

[Orengo04, Fig 6.2]

Structure vs. Sequence



[Orengo04, Fig 6.1]

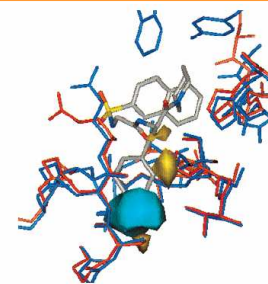
Applications

Fundamental step in:

- Analysis
- Visualization
- Comparison
- Design

Useful for:

- Structure classification
- Structure prediction
- Function prediction
- Drug discovery



Comparison of S1 binding pockets of thrombin (blue) and trypsin (red).

[Katzenholtz00]

Goals



Desirable properties:

- Automatic
- Discriminating
- Fast

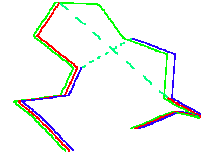
Theoretical Issues



NP-complete problem

- Arbitrary gap lengths
- Global scoring function

```
1 2 3 4 5 6 7 8 9 10 11 12 13 14
PHE ASP ILE CYS ARG LEU PRO GLY SER ALA GLU ALA VAL CYS
PHE ASN VAL CYS ARG THR PRO --- --- --- GLU ALA ILE CYS
PHE ASN VAL CYS ARG --- --- --- THR PRO GLU ALA ILE CYS
```



Methodological Issues



Choices:

- Representation
- Scoring function
- Search algorithm

Methodological Issues



Factors governing choices:

?

Methodological Issues



Factors governing choices:

- Application: homology detection, drug design, etc.
- Granularity: atom, residue, fragment, SSE
- Representation: inter-molecular, intra-molecular
- Scoring: geometric, gaps, chemical, structural, etc.
- Correspondences: sequential, non-sequential
- Gap penalty: expect gaps near loops, etc.
- Flexibility: rigid, flexible
- Target: single protein, representative proteins, PDB

Methodological Issues



Representations:

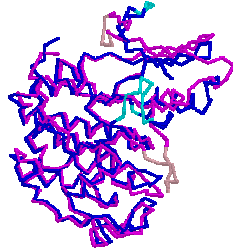
- Residue positions
- Local geometry
- Side chain contacts
- Distance matrices (DALI)
- Properties (COMPARER)
- SSEs (SSM, VAST)
- Geometric invariants

Methodological Issues



Scoring functions:

- Distances (RMSD)
- Substitutions
- Gaps



Methodological Issues



Search algorithms:

- Heuristics (CE)
- Monte Carlo (DALI, VAST)
- Dynamic programming (STRUCTAL, SSAP)
- Graph matching (SSM)

Outline



Alignment issues

Example alignment methods ←

Fold prediction experiment

Function prediction experiment

Example Methods

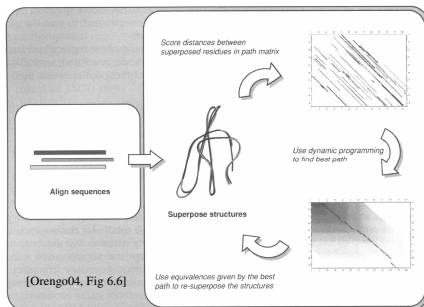


SSAP	Taylor & Orengo, 1989
STRUCTAL	Subbiah, Laurents & Levitt, 1993 Gerstein & Levitt 1998
DALI	Holm & Sander, 1993 Holm & Park, 2000
DEJAVU /LSQMAN	Kleywegt, 1996
CE	Shindyalov & Bourne, 1998
SSM	Krissinel & Henrick, 2003

+ 30 others!

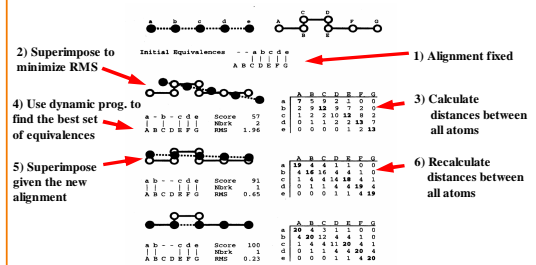
Slide by Rachel Kolodny

STRUCTAL

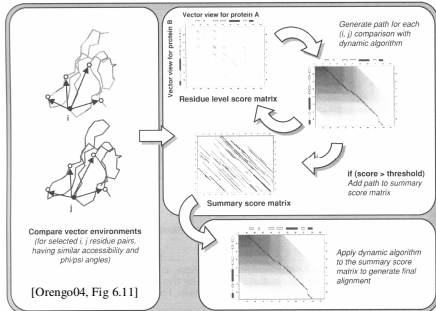


[Subbiah93, Gerstein98]

STRUCTAL

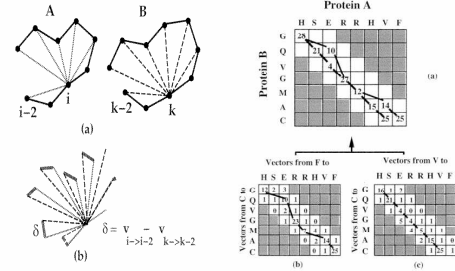


SSAP



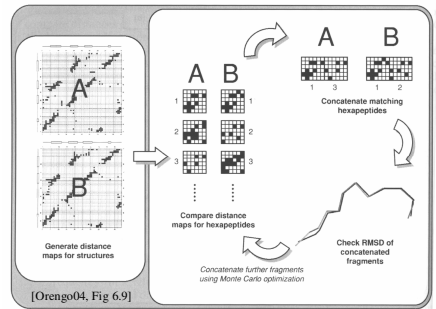
[Orengo96]

SSAP



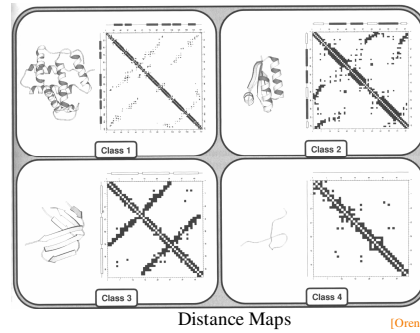
[Orengo96]

DALI



[Holm93]

DALI

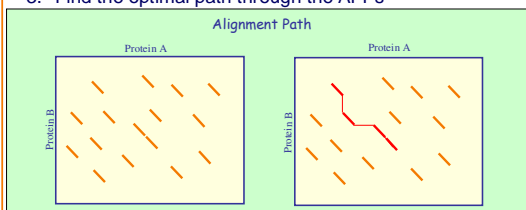


[Orengo04, Fig 6.7]

CE

Basic steps:

1. Compare octameric fragments to create candidate aligned fragment pairs (AFP)
2. Stitch together AFPs according to heuristics
3. Find the optimal path through the AFPs



SSM

Two-step solution:

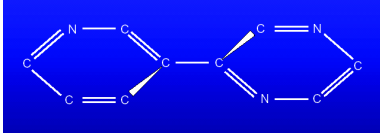
1. Graph representation of structures
2. Graph matching

SSM

Slide by Eugene Krissnel



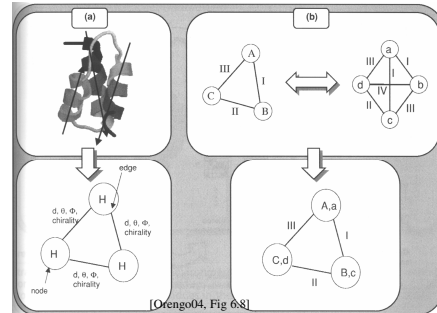
Graph representation of molecular structures



- Simple and intuitive, however results in intractably large graphs for proteins
- Solution: build graphs over stable substructures, such as secondary structure elements (SSEs). Having a correspondence between SSEs, one may use that for the 3D alignment of all core atoms.

Slide by Eugene Krissnel

SSM



[Orongo04, Fig 688]

Slide by Eugene Krissnel

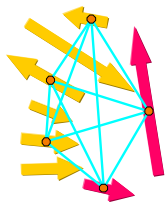
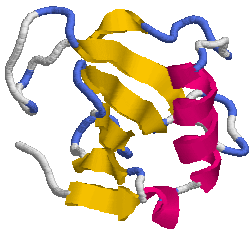
SSM

Slide by Eugene Krissnel



Graph representation of protein SSEs

E. M. Mitchell et al. (1990) J. Mol. Biol. 212:151
A. P. Singh and D. L. Bruttig (1997) ISMB-97 4:284



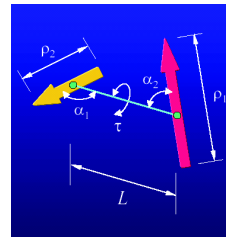
Slide by Eugene Krissnel

SSM

Slide by Eugene Krissnel



Protein graph labeling



Composite label of a vertex

- type - helix or strand
- length r

Composite label of an edge

- length L (directed if connects vertices from the same chain)
- vertex orientation angles α_1 and α_2
- torsion angle τ

Vertex and edge labels are matched with thresholds on particular quantities

Slide by Eugene Krissnel

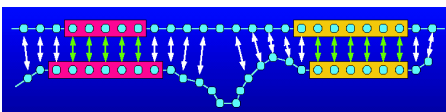
SSM

Slide by Eugene Krissnel



C_α alignment

- SSE-alignment is used as an initial guess for C_α -alignment
- C_α -alignment is an iterative procedure based on the expansion of shortest contacts at best superposition of structures



- C_α -alignment is a compromise between the alignment length N_a and $r.m.s.d.$. The optimised quantity is

$$Q = \frac{N_a^2}{[1 + (r.m.s.d./R_0)^2] N_1 N_2}$$

Slide by Eugene Krissnel

SSM

Slide by Eugene Krissnel



Statistical significance of match

- The overall probability of getting a particular match score by chance is the measure of the statistical significance of the match

$$P_{value} = 1 - \left(1 - P(S_{N_a}) P(S_{r.m.s.d.}) \prod_{SSE} P(S_{SSE}) \right)^{N_{combinations}}$$

- P_M is traditionally expressed through so-called Z-characteristics

$$P_0 = P(S_{N_a}) P(S_{r.m.s.d.}) \prod_{SSE} P(S_{SSE})$$

$$P_0 = \int_Z^\infty \omega(y) dy$$

$$\omega(y) = \sqrt{\frac{2}{\pi}} \exp\left(-\frac{1}{2}y^2\right)$$

Slide by Eugene Krissnel

SSM

Slide by Eugene Krissnel



SSM output

- Table of matched Secondary Structure Elements (SSE alignment)
- Table of matched core atoms (C_{α} -alignment) with dists between them
- Rotational-translation matrix of best structure superposition
- *R.m.s.d.* of C_{α} -alignment
- Length of C_{α} -alignment N_{α}
- Number of gaps in C_{α} -alignment N_g
- Quality score Q
- Probability estimate for the match P_M
- Z-characteristics
- Sequence identity

Slide by Eugene Krissnel

SSM

Slide by Eugene Krissnel



List of matches

Structure Alignment Results

Query: pdb entry 1ldc:chan [A, 478 residues]
 L-LACTATE DEHYDROGENASE, CYTOCHROME C DIOXIDOREDUCTASE, L-LACTATE DEHYDROGENASE (E.C.1.1.2.3) MUTANT WITH THR 143 ILEU4 REPLACED BY PHE (Y143F) COMPLEXED WITH PYRUVATE 1LD C5

Examined 19295 entries (39511 chains).
 Matches 1-14 of 14.

#	Scoring	Rmsd	N_{seq}	N_{ss}	N_{ca}	$\%_{\text{seq}}$	Query	Target (PDB entry)	Title
	Q P Z								
1	1.00 82.6 27.4	0.00	478	0	100	100	1ldc:chan	100 478	L-LACTATE DEHYDROGENASE, CYTOCHROME C DIOXIDOREDUCTASE, L-LACTATE DEHYDROGENASE (E.C.1.1.2.3) MUTANT WITH THR 143 ILEU4 REPLACED BY PHE (Y143F) COMPLEXED WITH PYRUVATE 1LD C5
2	0.99 62.7 23.8	0.30	478	1	100	100	1ldc:a	91 480	L-LACTATE DEHYDROGENASE, ILEU4
3	0.98 59.5 23.1	0.41	470	1	100	100	1ldc:a	89 481	FLAVOCYTOCHROME B2-(E.C.1.1.2.3) COMPLEXED WITH SULFITE 1LD3
4	0.94 59.1 23.1	0.56	478	1	100	97	1fcb:a	91 494	CRYSTALLOGRAPHIC STUDY OF THE RECOMBINANT FLAVIN-BINDING DOMAIN OF BAKER'S YEAST FLAVOCYTOCHROME B2 COMPARISON WITH THE INTACT WILD-TYPE ENZYME
5	0.91 55.4 22.3	0.51	474	2	98	97	1kb1:a	86 504	

Slide by Eugene Krissnel

SSM

Slide by Eugene Krissnel



Match details

Match 17 of 22

Back to match list | first match | << | >> | last match

Query PDB 1ldc:A		Alignment	
N_{seq}	$\%_{\text{seq}}$	N_{ss}	$\%_{\text{ss}}$
478	66	31	65

Q P RMSD N_{seq} N_{ss} N_{ca} $\%_{\text{seq}}$ $\%_{\text{ss}}$ $\%_{\text{ca}}$

0.514 30.07 1.191 315 6

Z N_{seq} N_{ss} N_{ca} $\%_{\text{seq}}$ $\%_{\text{ss}}$ $\%_{\text{ca}}$

33.3 17.30 20 6

L-LACTATE DEHYDROGENASE, CYTOCHROME C DIOXIDOREDUCTASE, L-LACTATE DEHYDROGENASE (E.C.1.1.2.3) MUTANT WITH THR 143 ILEU4 REPLACED BY PHE (Y143F) COMPLEXED WITH PYRUVATE 1LD C5

view | download | view superposed | view | download

Slide by Eugene Krissnel

SSM

Slide by Eugene Krissnel



SSE alignment

Secondary Structure Alignment

Query PDB 1ldc:A	Target PDB thuv:A
10 IHI 13 IAI ASN 124 ILEU 136 I	<=> 1 IHI 13 IAI ASN 7 ILEU 19 I
11 IHI 10 IAI THR 127 ISEP 146 I	<=> 2 IHI 10 IAI PRO 20 IGLU 29 I
12 IHI 10 IAI GLU 151 IALA 160 I	<=> 3 IHI 10 IAI ILEU 34 IVAL 43 I
13 IHD 10 IAI PHE 156 IVAL 194 I	<=> 7 IHD 10 IAI ILEU 79 ILEU 77 I
14 IHI 10 IAI GLU 208 IGLU 217 I	<=> 9 IHI 14 IAI SER 69 ILEU 102 I
15 IHD 4 IAI GLN 225 ISEP 228 I	<=> 10 IHD 3 IAI PHE 105 ILEU 107 I
16 IHI 8 IAI SER 234 IALA 241 I	<=> 11 IHI 9 IAI SER 114 ITHR 122 I
17 IHD 5 IAI GLN 249 ILEU 253 I	<=> 12 IHD 5 IAI ILEU 126 ILEU 130 I
18 IHI 15 IAI ASN 258 ILEU 272 I	<=> 13 IHI 15 IAI SER 134 ITHR 148 I
19 IHD 4 IAI ILEU 277 ITHR 280 I	<=> 14 IHD 5 IAI THR 152 ITHR 156 I
20 IHI 8 IAI ASN 289 ILEU 296 I	<=> 15 IHI 7 IAI ASN 165 IASN 171 I
21 IHI 12 IAI THR 331 ITHR 342 I	<=> 18 IHI 12 IAI ASN 213 ITHR 224 I
22 IHD 6 IAI ILE 346 IVAL 351 I	<=> 19 IHD 7 IAI ILE 227 ILEU 233 I
23 IHI 11 IAI ASN 353 ILEU 363 I	<=> 20 IHI 11 IAI SER 232 IGLU 243 I
24 IHD 4 IAI GLY 367 ILEU 370 I	<=> 21 IHD 4 IAI GLY 249 ILEU 252 I
25 IHI 14 IAI ALA 383 IASP 398 I	<=> 23 IHI 10 IAI VAL 269 IGLU 278 I
27 IHD 5 IAI GLY 405 IASP 409 I	<=> 24 IHD 3 IAI VAL 261 ILEU 263 I
28 IHI 11 IAI ASN 414 ILEU 424 I	<=> 25 IHI 11 IAI ASN 269 ILEU 299 I
29 IHD 4 IAI GLY 428 ILEU 431 I	<=> 26 IHD 4 IAI ALA 303 ILEU 306 I
30 IHI 35 IAI GLY 432 IGLU 446 I	<=> 27 IHI 35 IAI GLY 307 IGLU 341 I

OCA | SCOP domain | SCOP family
 GeneCensus | ESSP | SDET | GATH | PDBsum
 SWISS-PROT | TrEMBL | GOX | SPIDER
 Protomap | MDL | PSEUD | GOX | SPIDER

view | download sequence | view superposed | view | download sequence

Slide by Eugene Krissnel

SSM

Slide by Eugene Krissnel



Rotational-translation matrix
 (to be applied to the query)

-0.710	0.423	-0.563	X	98.778
-0.471	0.309	0.824	X Y +	96.485
0.524	0.852	-0.020	Z	-59.445

9D Structural alignment

PDB 1ldc:A	Dist. (Å)	PDB thuv:A
AI: ILE 203	1.42	AI: ILE 19
AI: ASN 204	0.27	AI: THR 29
AI: PRO 205	2.75	AI: PRO 88
AI: ILE 209		AI: ILE 89
AI: VAL 208	1.56	AI: VAL 90
AI: ILE 209	0.77	AI: ASN 91
AI: ILE 210	0.48	AI: ILE 92
AI: ASP 211	0.91	AI: ALA 93
AI: THR 212	0.93	AI: THR 94
AI: ALA 213	1.57	AI: ALA 95
AI: ARG 214	1.65	AI: ARG 96
AI: ILE 215	1.55	AI: ALA 97
AI: SER 216	1.47	AI: ALA 98
AI: GLY 217	2.01	AI: THR 99
AI: THR 218	2.99	AI: ILE 100
AI: GLY 219		
AI: VAL 220		
AI: PHE 221	2.95	AI: ALA 101
AI: ILE 222	1.92	AI: GLY 102

Rotational-translation matrix of best superposition

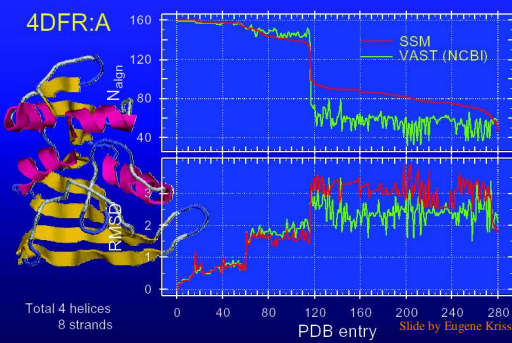
C α -alignment

Total 4 helices
8 strands

Slide by Eugene Krissnel

SSM Results

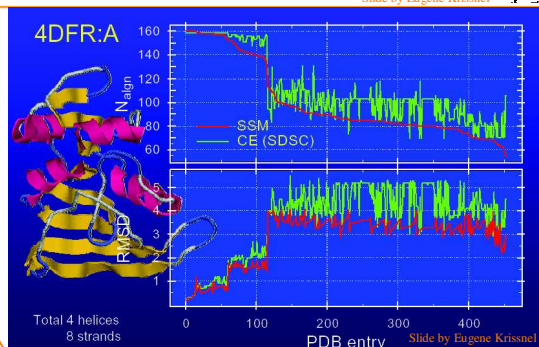
Slide by Eugene Krissnel



Slide by Eugene Krissnel

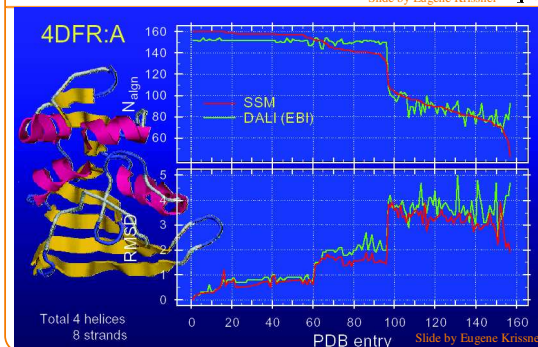
SSM Results

Slide by Eugene Krissinel



SSM Results

Slide by Eugene Krissinel



Outline



Alignment issues

Example alignment methods

Fold prediction experiment ←

Function prediction experiment

Fold Prediction Experiments



Evaluate how useful alignment algorithms are for predicting a protein's fold

How?

Fold Prediction Experiments



Kolodny, Koehl, & Levitt [2005]

- ROC curves and geometric measures using CATH

Sierk & Pearson [2004]

- ROC curves using CATH

Novotny et al. [2004]

- Checked a few dozen cases using CATH

Lepplae & Hubbard [2002]

- ROC curves using SCOP

Fold Prediction Experiments



Kolodny, Koehl, & Levitt [2005] ←

- ROC curves and geometric measures using CATH

Sierk & Pearson [2004]

- ROC curves using CATH

Novotny et al. [2004]

- Checked a few dozen cases using CATH

Lepplae & Hubbard [2002]

- ROC curves using SCOP

Kolodny, Koehl, & Levitt [2005]



Large scale alignment study

- 2,930 structures (all pairs)
- 6 structural alignment algorithms
- 4 geometric scoring functions
- Evaluation with respect to CATH topology level
- 20,000 hours of compute time

Tested Methods



SSAP	Taylor & Orengo, 1989
STRUCTAL	Subbiah, Laurents & Levitt, 1993 Gerstein & Levitt 1998
DALI	Holm & Sander, 1993 Holm & Park, 2000
DEJAVU /LSQMAN	Kleywegt, 1996
CE	Shindyalov & Bourne, 1998
SSM	Krissinel & Henrick, 2003
Best-of-All	Best of above methods

Slide by Rachel Kolodny

Scoring Functions



Consider # aligned residues & geometric similarity:

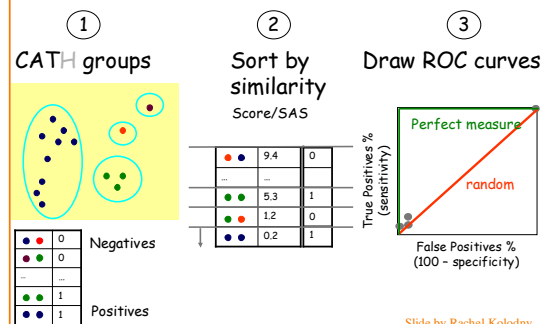
$$SAS = \frac{RMSD \times 100}{N_{mat}}$$

Also penalize gaps:

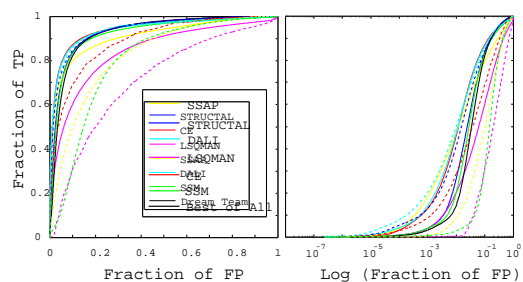
$$GSAS = \begin{cases} \text{if } (N_{mat} > N_{gap}) & \frac{RMSD \times 100}{N_{mat} - N_{gap}} \\ \text{else} & 99.9 \end{cases}$$

[Kolodny05]

Evaluation Using ROC Curves



SAS & Native ROC Curves



ROC Curve Issues



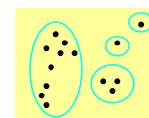
Uses only internal ordering

- Estimation of similarity can be very wrong

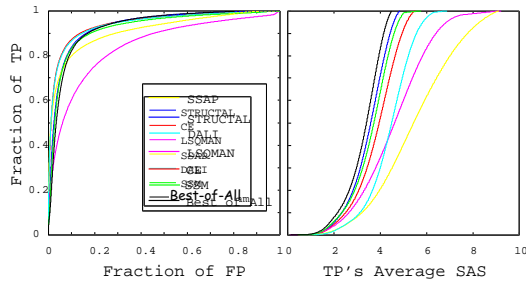
●●	9.4	●●	9400
—	—	—	—
●●	5.3	●●	5300
●●	1.2	●●	1200
●●	0.2	●●	200

Native scores or SAS

Converts a classification gold standard into binary truth

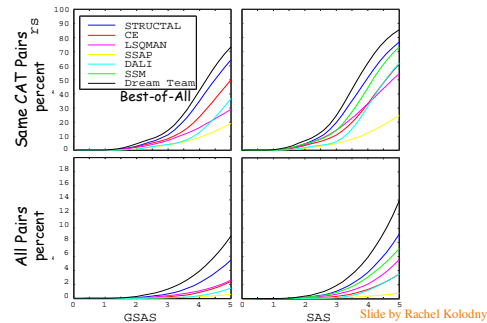


Comparing SAS Values Directly



Slide by Rachel Kolodny

GSAS & SAS Distributions



Slide by Rachel Kolodny

Contributions to "Best-of-All"



	Total	SSAP	STRUCTAL	DALI	LSQMAN	CE	SSM
GSAS ≤ 5 Å (100%)	275,547	832 (0.3%)	189,871 (69%)	5868 (2.1%)	54,606 (20%)	24,370 (8.8%)	-
SAS ≤ 5 Å (100%)	539,755	498 (0.09%)	286,972 (53%)	15,648 (2.9%)	103,408 (19.2%)	15,844 (2.9%)	117,385 (21.8%)
SI ≤ 5 Å (100%)	978,531	3745 (0.4%)	497,330 (51%)	24,767 (2.5%)	201,202 (21%)	17,142 (1.8%)	234,345 (24%)
MI ≤ 0.8 (100%)	880,503	4579 (0.5%)	575,542 (65%)	31,402 (3.6%)	63,088 (7.2%)	72,974 (8.3%)	134,918 (15.3%)

The absolute number of alignments contributed by each method is listed and the percentage of alignments is given in parentheses. The largest contributor is shown in bold.

[Kolodny05]

Outline



Alignment issues

Example alignment methods

Fold prediction experiment

Function prediction experiment ←

Function Prediction Experiment



Evaluate how useful alignment methods are for predicting a protein's molecular function

How?

Data Set



Proteins crystallized with bound ligands

- PDB file must have resolution ≤ 3 Angstroms
- Ligands must have ≥ 20 HETATOMS

Classified by reaction/reactant

- PDB file must have an EC number (enzymes only)
- EC number must have a KEGG reaction with a reactant whose graph closely matches ligand in PDB file

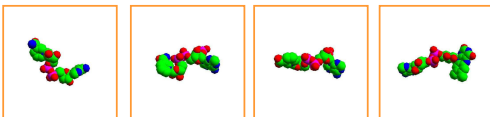
Non-redundant

- No two ligands contacting domains with same CATH S95
- No two ligands contacting domains with same SCOP SP
- No two ligands from same PDB file

Data Set



351 proteins / 58 Reactions (189 outliers)

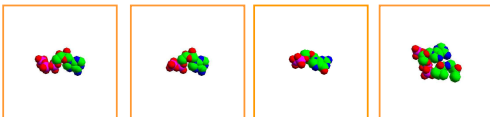


55 NAD (34/9)

25 NDP (9/3)

38 NAP (18/8)

11 FAD (9/3)



21 ATP (5/2)

29 ADP (10/5)

6 GDP (6/2)

12 COA (5/2)

Data Set



REACTION	NAME	#	REACTION	NAME	#	REACTION	NAME	#
R00145	NAD	2	R00162	ATP	3	R00408	FAD	5
R00214	NAD	2	R00847	ATP	2	R00924	FAD	2
R00342	NAD	7	R00124	ADP	2	R01175	FAD	2
R00338	NAD	3	R00897	ADP	2	MISC	FAD	2
R00623	NAD	5	R00758	ADP	2	R00351	COA	3
R00703	NAD	5	R01512	ADP	2	R03552	COA	2
R01061	NAD	5	R02412	ADP	2	MISC	COA	7
R01463	NAD	2	R03847	AMP	2	R02261	SAM	3
R01778	NAD	2	R00330	GDP	2	MISC	SAM	3
R00112	NAP	2	R01135	GDP	4	R03552	ACO	2
R00343	NAP	2	R01130	IMP	3	R02911	GDU	2
R00365	NAP	2	R03394	TMP	2	R03622	GT1	12
R00359	NAP	2	R02101	UMP	6	R01146	PCC	3
R01041	NAP	4	R00965	USP	2	R01090	PRP	2
R01068	NAP	2	R00966	USP	2	R01402	MTA	2
R01195	NAP	2	R01229	SSP	2	R03435	BP	2
R02477	NAP	2	MISC	ATP	19	R02928	PSI	4
R00703	NAI	2	MISC	ADP	19	R01590	ACD	2
R00939	NDP	5	MISC	AMP	10	R00529	ADX	2
R01063	NDP	2	MISC	ASP	5	R03491	SIA	2
R01195	NDP	2	MISC	QMP	2	R01137	MIN	3
MISC	NAD	21	MISC	UDP	4	R03992	MYA	2
MISC	NAP	20	MISC	UMP	1	R03509	13T	2
MISC	NAH	2	MISC	5CP	1	MISC	etc	etc
MISC	NAI	2						
MISC	NDP	16						

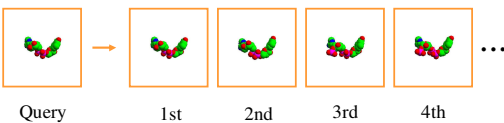
Evaluation Method



“Leave-one-out” classification experiment

∅Match every ligand against all the others in data set

- Log a “hit” when best match performs same reaction
- Report percentage of hits (correctly classified ligands)



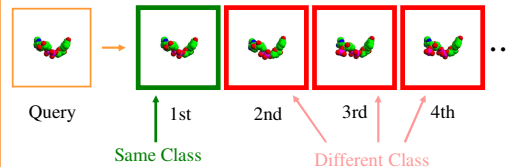
Evaluation Method



“Leave-one-out” classification experiment

∅Match every ligand against all the others in data set

- Log a “hit” when best match performs same reaction
- Report percentage of hits (correctly classified ligands)



Evaluation Method

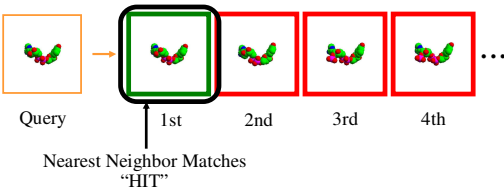


“Leave-one-out” classification experiment

• Match every ligand against all the others in data set

∅Log a “hit” when best match performs same reaction

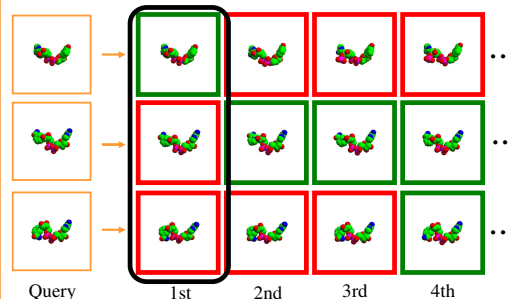
- Report percentage of hits (correctly classified ligands)



Evaluation Method



Classification rate is 33% in this example



Structure Alignment Method



Use CE to compute similarity of protein structures

```
CE - /ebi/data/pdbe/1jsu.pdb A -/ebi/data/pdbe/1hcl.pdb _ scratch
Structure Alignment Calculator, version 1.02, last modified: Jun 15, 2001.
```

```
CE Alignment: 1.000000
Align:
Rmsd = 2.28A
Z-Score = 6.8
Gaps = 30 (11.5%)
Seqs:
```

```
X2 = ( 0.597420)*X1 + ( 0.071548)*Y1 +
      ( 0.005923)*Z1 + (-93.687386)
Y2 = ( 0.059473)*X1 + (-0.777232)*Y1 +
      (-0.608397)*Z1 + ( 119.695427)
Z2 = (-0.040214)*X1 + ( 0.625133)*Y1 +
      (-0.779462)*Z1 + ( 84.334138)
```

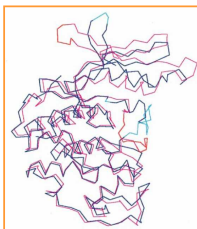
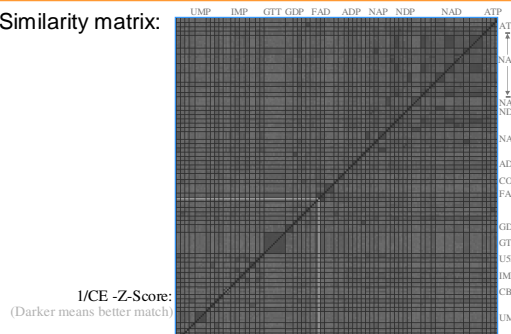


Image from Shindyalov and Bourne (1998)

Structure Alignment Results



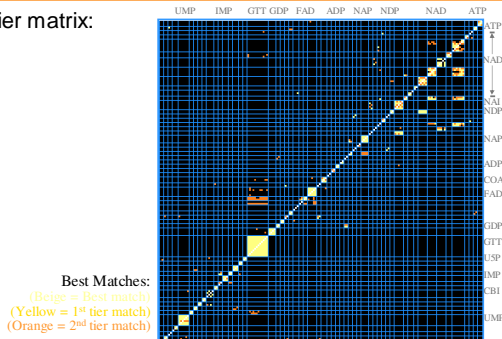
Similarity matrix:



Structure Alignment Results



Tier matrix:



Structure Alignment Results



Classification rate:

FASTA = 68%
CE = 65%
Random = <1%

Structure Alignment Results



Classification rate: When Smith-Waterman \geq 500:
FASTA = 68% Sequence = 80%
CE = 65% CE = 72%
Random = <1% Random = <1%

When Smith-Waterman < 500:
CE = 53%
FASTA = 44%
Random = <1%

Conclusion



Many algorithms for structural alignment, differing according to

- Application: homology detection, drug design, etc.
- Granularity: atom, residue, fragment, SSE
- Representation: inter-molecular, intra-molecular
- Scoring: geometric, gaps, chemical, structural, etc.
- Correspondences: sequential, non-sequential
- Gap penalty: expect gaps near loops, etc.
- Flexibility: rigid, flexible
- Target: single protein, representative proteins, PDB

None seems best for all situations

All probably provide some benefit over sequence