

Structural Alignment of Proteins

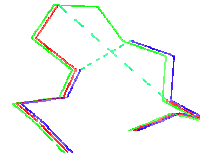
Thomas Funkhouser
Princeton University
CS597A, Fall 2005

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

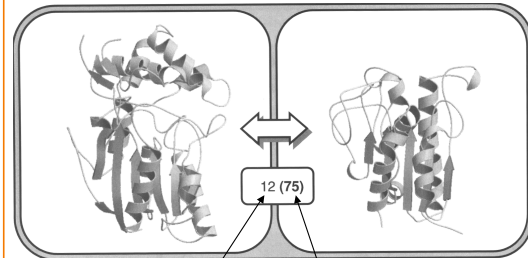
Superposition

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

Alignment

- Find correspondences, and then superpose structures

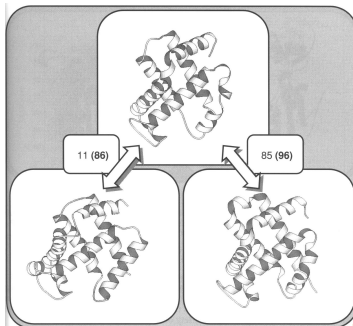
Structure vs. Sequence



Sequence Identity (Structure similarity)

[Orengo04, Fig 6.2]

Structure vs. Sequence



[Orengo04, Fig 6.1]

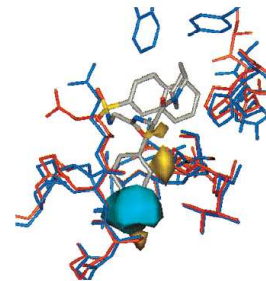
Applications

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

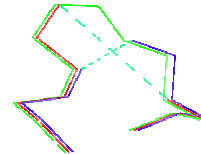
Theoretical Issues



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:

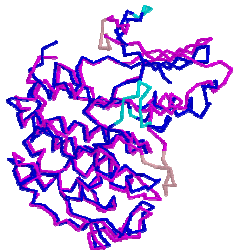
- Residue positions
- Local geometry
- Side chain contacts
- Distance matrices (DALI)
- Properties (COMPARER)
- SSEs (SSM, VAST)
- Geometric invariants

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

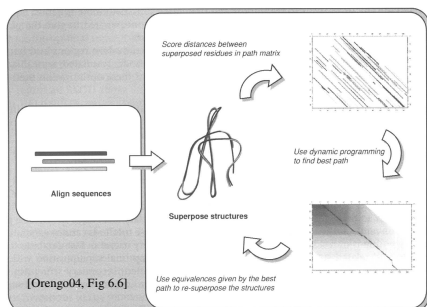


SSAP	Taylor & Orengo, 1989
STRUCTAL	Subbiah, Laurents & Levitt, 1993 Gerstein & Levitt 1998
DALI	Holm & Sander, 1993 Holm & Park, 2000
DEJAVU /LSQMAN	Kleywegt, 1996
CE	Shindyalov & Bourne, 1998
SSM	Krissinel & Henrick, 2003

+ 30 others!

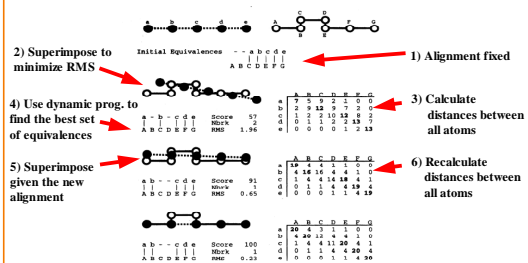
Slide by Rachel Kolodny

STRUCTAL



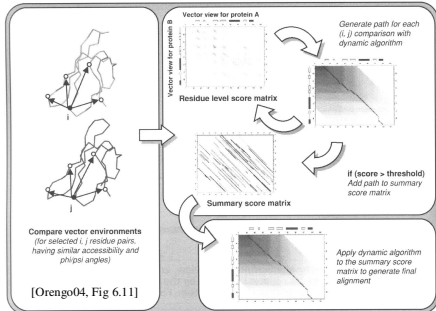
[Subbiah93, Gerstein98]

STRUCTAL



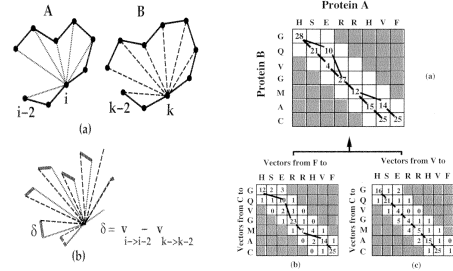
[Subbiah93, Gerstein98]

SSAP

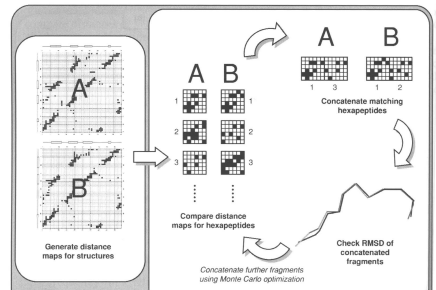


[Orengo96]

SSAP

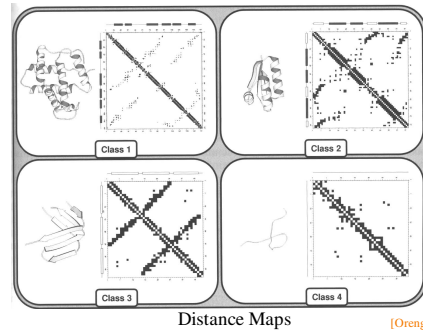


DALI



[Holm93]

DALI

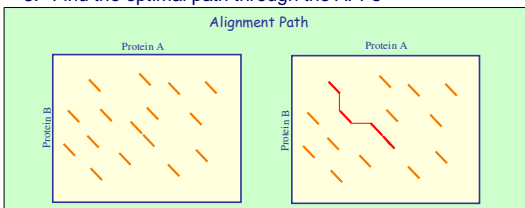


CE



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:

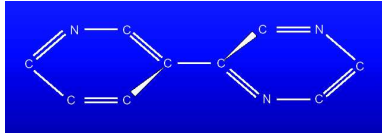
1. Graph representation of structures
2. Graph matching

SSM

Slide by Eugene Krissnel



Graph representation of molecular structures

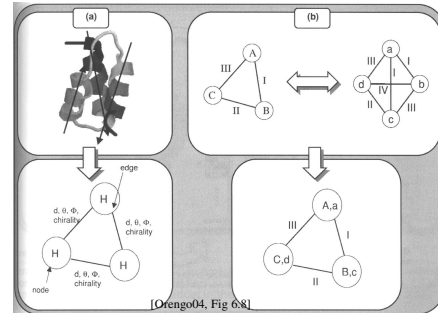


- Simple and intuitive, however results in intractably large graphs for proteins
- Solution: build graphs over stable substructures, such as secondary structure elements (SSEs). Having a correspondence between SSEs, one may use that for the 3D alignment of all core atoms.

Slide by Eugene Krissnel

SSM

Slide by Eugene Krissnel



[Orongo04, Fig 6&8]

Slide by Eugene Krissnel

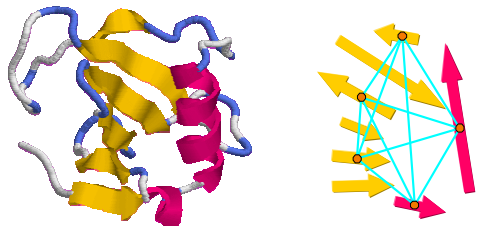
SSM

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. Brutlag (1997) ISMB-97 4:284



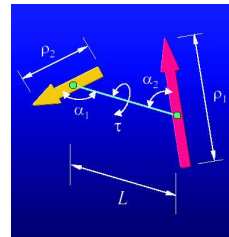
Slide by Eugene Krissnel

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

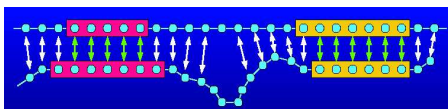
SSM

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 Krissinel



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 Krissinel

SSM

Slide by Eugene Krissinel



List of matches

Structure Alignment Results

Query: pdb entry 1ldc:chan [A:478] residues.
L-LACTATE DEHYDROGENASE, CYTOCHROME C DIOXYGENASE, L-LACTATE DEHYDROGENASE (E.C.1.1.2.3) MUTANT WITH TYR 143 ILLD4 REPLACED BY PHE (Y143F) COMPLEXED WITH PYRUVATE ILLD5

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

#	Scoring			Rmsd	N _{seq}	N _α	% _{seq}	% _{ssr}	Query	Target (PDB entry)		
	Q	P	Z							Match	% _{seq}	N _{seq}
1	1.00	82.6	27.4	0.00	478	0	100	100	1ldc:1A	100	478	L-LACTATE DEHYDROGENASE, CYTOCHROME C DIOXYGENASE, L-LACTATE DEHYDROGENASE (E.C.1.1.2.3) MUTANT WITH TYR 143 ILLD4 REPLACED BY PHE (Y143F) COMPLEXED WITH PYRUVATE ILLD5
2	0.99	62.7	23.8	0.30	478	1	100	100	1ldc:1A	91	480	L-LACTATE DEHYDROGENASE, ILLD4
3	0.98	59.5	23.1	0.41	470	1	100	100	1ldc:1A	89	481	FLAVO-CYTOCHROME B0-(E.C.1.1.2.3) COMPLEXED WITH SULFITE ILLD3
4	0.94	59.1	23.1	0.56	478	1	100	97	1fcb:1A	91	494	CRYSTALLOGRAPHIC STUDY OF THE RECOMBINANT FLAVIN-BINDING DOMAIN OF BAKER'S YEAST FLAVO-CYTOCHROME B0, COMPARISON WITH THE INTERACT WILD-TYPE ENZYME
5	0.91	55.4	22.3	0.51	474	2	98	97	1kb1:1A	86	504	

Slide by Eugene Krissinel

SSM

Slide by Eugene Krissinel



Match details

Match 17 of 22

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

Query PDB 1ldc:A				Alignment			
N _{seq}	% _{seq}	N _{ssr}	% _{ssr}	Q	P	RMSD	N _{gaps}
478	66	31	65	0.514	30.07	1.191	315

L-LACTATE DEHYDROGENASE, CYTOCHROME C DIOXYGENASE, L-LACTATE DEHYDROGENASE (E.C.1.1.2.3) MUTANT WITH TYR 143 ILLD4 REPLACED BY PHE (Y143F) COMPLEXED WITH PYRUVATE ILLD5

view download view superposed view download

Slide by Eugene Krissinel

SSM

Slide by Eugene Krissinel



SSE alignment

Secondary Structure Alignment

Query PDB 1ldc:A		Target PDB thuv:A					
10 IH1	13 IA ASN 124	LEU 136	<=>	1 IH1	13 IA ASN 7	LEU 19	I
11 IH1	10 IA THR 127	SER 146	<=>	2 IH1	10 IA PRO 20	GLY 29	I
12 IH1	10 IA GLU 151	ALA 160	<=>	3 IH1	10 IA ILE 34	VAL 43	I
13 IH1	4 IA PHE 150	VAL 194	<=>	7 IH1	3 IA LEU 78	ILE 77	I
14 IH1	10 IA GLY 208	GLY 217	<=>	9 IH1	14 IA SER 69	GLY 102	I
15 IH1	8 IA ASP 225	ASP 228	<=>	10 IH1	3 IA PHE 105	LEU 107	I
16 IH1	8 IA SER 234	ALA 241	<=>	11 IH1	5 IA SER 114	CTE 122	I
17 IH1	5 IA GLN 249	LEU 253	<=>	12 IH1	5 IA LEU 126	LEU 130	I
18 IH1	15 IA ASN 258	LEU 272	<=>	13 IH1	15 IA SER 134	THR 148	I
19 IH1	4 IA LEU 277	THR 280	<=>	14 IH1	5 IA THR 152	THR 156	I
20 IH1	8 IA ARG 289	LEU 296	<=>	15 IH1	7 IA ARG 165	ASN 171	I
21 IH1	12 IA THR 331	THR 342	<=>	18 IH1	12 IA ASN 213	THR 224	I
22 IH1	6 IA ILE 346	VAL 351	<=>	19 IH1	7 IA LEU 227	LEU 233	I
23 IH1	11 IA ARG 353	ILE 363	<=>	20 IH1	11 IA SER 232	GLY 243	I
24 IH1	4 IA GLY 367	LEU 370	<=>	21 IH1	4 IA GLY 249	LEU 252	I
25 IH1	14 IA ALA 383	ARG 398	<=>	22 IH1	10 IA VAL 269	GLY 278	I
27 IH1	5 IA GLY 405	ASP 409	<=>	24 IH1	3 IA VAL 261	ILE 263	I
28 IH1	11 IA ARG 414	LEU 424	<=>	25 IH1	11 IA ARG 269	LEU 269	I
29 IH1	4 IA GLY 428	LEU 431	<=>	26 IH1	4 IA ALA 303	LEU 306	I
30 IH1	35 IA GLY 432	GLY 446	<=>	27 IH1	35 IA GLY 307	GLY 341	I

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

view download sequence view superposed view download sequence

Slide by Eugene Krissinel

SSM

Slide by Eugene Krissinel



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	96.485
0.524	0.852	-0.020	B	-59.445

9D Structural alignment

PDB 1ldc:A	Dist. (Å)	PDB thuv:A
A1: ASN 203	1.42	A1: ASN 88
A1: ASN 204	0.27	A1: THR 89
A1: PRO 205	2.75	A1: SO 88
A1: LEU 209		A1: LEU 89
A1: VAL 207	1.56	A1: LEU 89
A1: GLY 208	1.66	A1: GLY 90
A1: GLU 209	0.77	A1: ARG 91
A1: SER 210	0.98	A1: LEU 92
A1: ASP 211	0.91	A1: ALA 93
A1: THR 212	0.95	A1: TRP 94
A1: ALA 213	1.57	A1: ALA 95
A1: ARG 214	1.65	A1: ARG 96
A1: GLU 215	1.55	A1: ALA 97
A1: SER 216	1.47	A1: ALA 98
A1: GLY 217	2.01	A1: THR 99
A1: THR 218	2.99	A1: LEU 100
A1: GLY 219		
A1: VAL 220		
A1: PRO 221	2.95	A1: ALA 101
A1: SER 222	1.92	A1: GLY 102

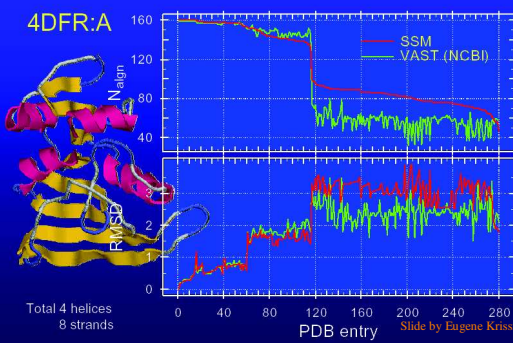
C_α-alignment

Rotational-translation matrix of best superposition

Slide by Eugene Krissinel

SSM Results

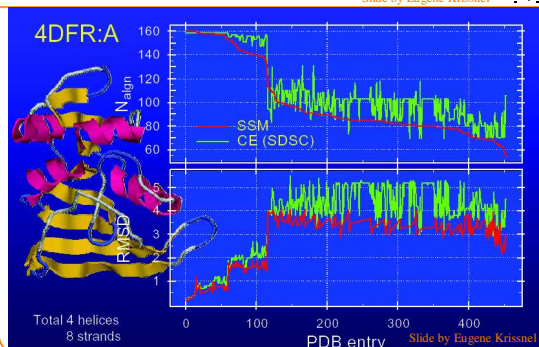
Slide by Eugene Krissinel



Slide by Eugene Krissinel

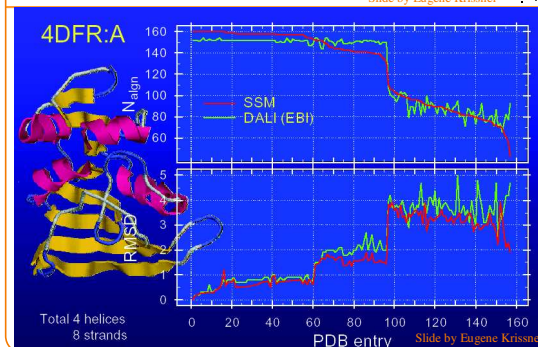
SSM Results

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:

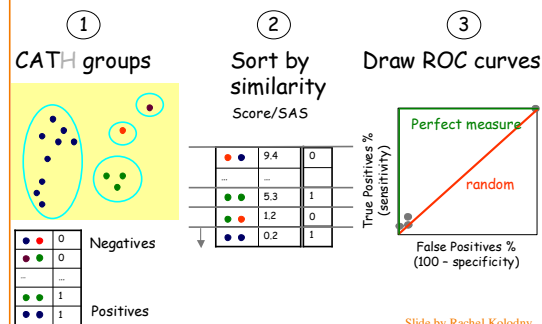
$$SAS = \frac{RMSD \times 100}{N_{mat}}$$

Also penalize gaps:

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

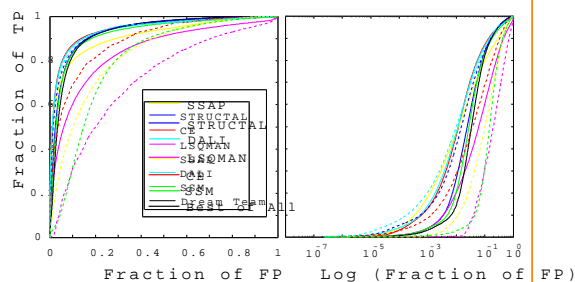
[Kolodny05]

Evaluation Using ROC Curves



Slide by Rachel Kolodny

SAS & Native ROC Curves



Slide by Rachel Kolodny

ROC Curve Issues



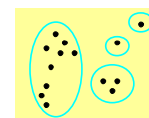
Uses only internal ordering

- Estimation of similarity can be very wrong

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

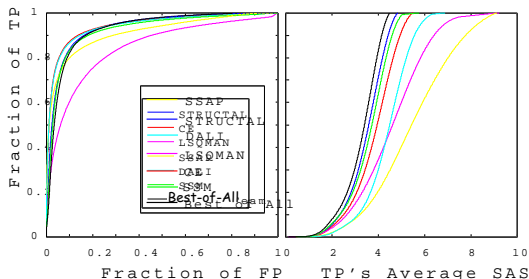
Native scores or SAS

Converts a classification gold standard into binary truth



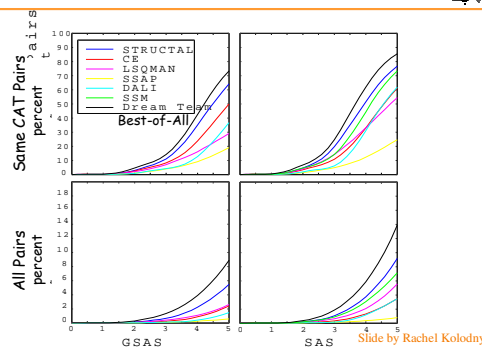
Slide by Rachel Kolodny

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%)	375,542 (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

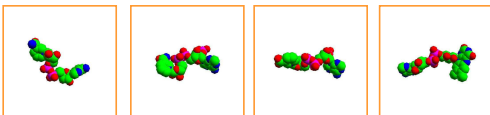
Non-redundant

- 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)

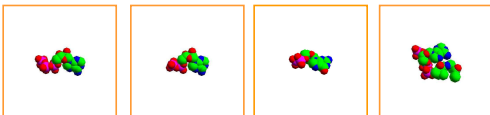


55 NAD (34/9)

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



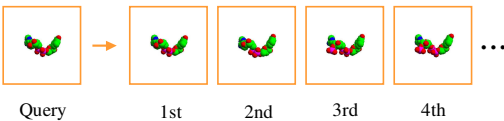
REACTION	NAME	#	REACTION	NAME	#	REACTION	NAME	#
R00145	NAD	2	R00162	ATP	3	R00408	FAD	5
R00214	NAD	2	R00347	ATP	2	R00924	FAD	2
R00342	NAD	7	R00124	ADP	2	R01175	FAD	2
R00538	NAD	3	R00391	ADP	2	MISC	FAD	2
R00623	NAD	5	R00756	ADP	2	R00351	COA	3
R00703	NAD	5	R01512	ADP	2	R00552	COA	2
R01061	NAD	5	R00412	ADP	2	MISC	COA	7
R01463	NAD	2	R00347	AMP	2	R00291	SAM	3
R01778	NAD	2	R00330	GDP	2	MISC	SAM	3
R00112	NAP	2	R01135	GDP	4	R00552	ACO	2
R00343	NAP	2	R01130	IMP	3	R00291	GDU	2
R00625	NAP	2	R00391	IMP	2	R00522	GTT	12
R00553	NAP	2	R02101	UMP	5	R01125	POC	3
R01041	NAP	4	R00965	USP	2	R00190	PRP	2
R01058	NAP	2	R00966	USP	2	R01402	MTA	2
R01155	NAP	2	R01229	SGP	2	R00345	BP	2
R00477	NAP	2	MISC	ATP	16	R00988	CSH	4
R00703	NAI	2	MISC	ADP	19	R01590	ACD	2
R00930	NDP	5	MISC	AMP	10	R00529	ADX	2
R01063	NDP	2	MISC	ASP	5	R00491	SIA	2
R01165	NDP	2	MISC	GTP	2	R00137	MIN	3
MISC	NAD	21	MISC	UDP	4	R00992	MYA	2
MISC	NAP	20	MISC	LMP	1	R00509	13T	2
MISC	NAH	2	MISC	SGP	1	MISC	etc	etc
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)

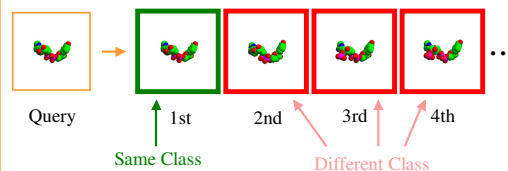


Evaluation Method



"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)

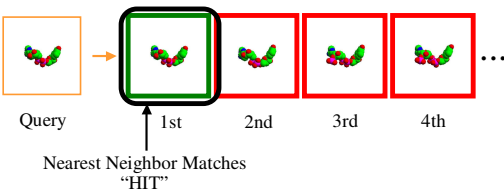


Evaluation Method



"Leave-one-out" classification experiment

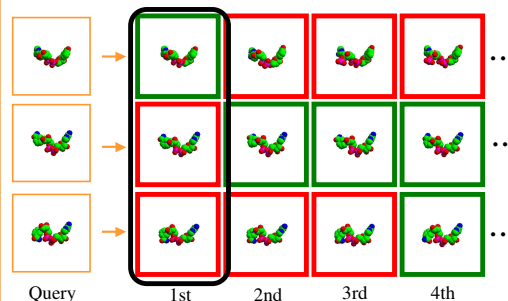
- Match every ligand against all the others in data set
- ∅ Log a "hit" when best match performs same reaction
- Report percentage of hits (correctly classified ligands)



Evaluation Method



Classification rate is 33% in this example



Sequence Alignment Method



Use FASTA to compute Smith-Waterman score for every pair of SCOP domains contacting ligand

```
> fasta34 d1gv0a d1guya
      10      20      30      40      50      60
d1gv0a ACVLDGARFRSFIAMELGVSMQVTEACVLSHGDMVPPVKYTTVAGIPVADLISAERIA
      : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : :
d1guya AGVLDAAARYRTFIAMEAGVSVEDVQAMLMGSHGDEMPLPFPSTISGIPVSEFIAPDRLA
      10      20      30      40      50      60

      70      80      90     100     110     120
d1gv0a ELVERTRTGGAEIVNHLKQSAFYSFSPATSVEMVESIVLDRKRVLTCVSLDQGVGIDGT
      : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : :
d1guya QIVERTKGGGSEIVNLLKTSAYYAPAAATAGMVEAVLKDKKRVMFVAAYLTGQVGLNDI
      70      80      90     100     110     120

      130     140     150     160
d1gv0a FVGVPVKLGKNGVEHIYEIKLDQSDLLDLQSAKIVDENCKML
      : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : :
d1guya YFGVPVILGAGGVKILFLPLNEEMALLNASAKAVRATLDL
      130     140     150     160

54.487% identity
156 out of 163 amino acids overlap
Smith-Waterman score: 588
```

Sequence Alignment Method



Use FASTA to compute Smith-Waterman score for every pair of SCOP domains contacting ligand

```
> fasta34 d1gv0a d1guya
      10      20      30      40      50      60
d1gv0a ACVLDGARFRSFIAMELGVSMQVTEACVLSHGDMVPPVKYTTVAGIPVADLISAERIA
      : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : :
d1guya AGVLDAAARYRTFIAMEAGVSVEDVQAMLMGSHGDEMPLPFPSTISGIPVSEFIAPDRLA
      10      20      30      40      50      60

      70      80      90     100     110     120
d1gv0a ELVERTRTGGAEIVNHLKQSAFYSFSPATSVEMVESIVLDRKRVLTCVSLDQGVGIDGT
      : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : :
d1guya QIVERTKGGGSEIVNLLKTSAYYAPAAATAGMVEAVLKDKKRVMFVAAYLTGQVGLNDI
      70      80      90     100     110     120

      130     140     150     160
d1gv0a FVGVPVKLGKNGVEHIYEIKLDQSDLLDLQSAKIVDENCKML
      : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : :
d1guya YFGVPVILGAGGVKILFLPLNEEMALLNASAKAVRATLDL
      130     140     150     160

54.487% identity
156 out of 163 amino acids overlap
Smith-Waterman score: 588
```

Sequence Alignment Method



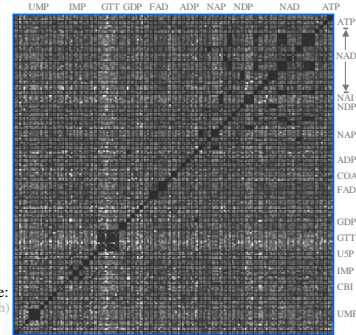
Use FASTA to compute Smith-Waterman score for every pair of SCOP domains contacting ligand

$$D(A, B) = 1 / \max_{A_i, B_j} \text{SmithWaterman}(A_i, B_j)$$

Sequence Alignment Results



Similarity matrix:

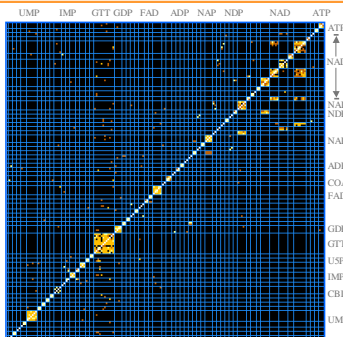


1/SmithWaterman Score:
(Darker means better match)

Sequence Alignment Results



Tier matrix:



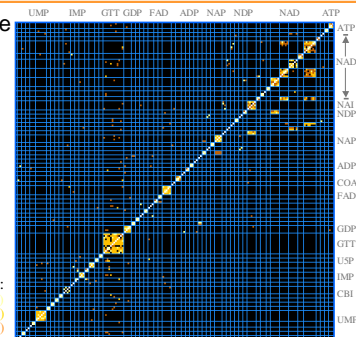
Best Matches:
(Beige = 1st tier match)
(Yellow = 1st tier match)
(Orange = 2nd tier match)

Sequence Alignment Results



Classification rate

FASTA = 68%
Random = <1%



Best Matches:
(Beige = 1st tier match)
(Yellow = 1st tier match)
(Orange = 2nd tier match)

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 Aligned: 1000 (100.0%)
Align: 1000 (100.0%)
Rmsd = 2.28Å
Z-Score = 6.8
Gaps = 30 (3.0%)
Seq1:
Seq2:
```

```
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)
```

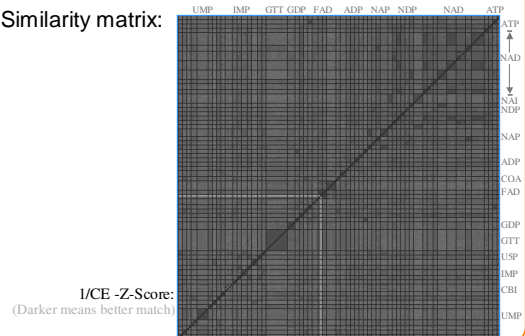


Image from Shindyalov and Bourne (1998)

Structure Alignment Results



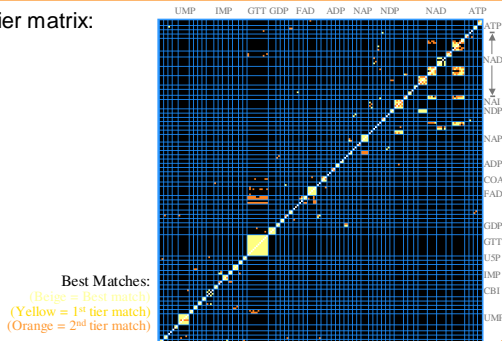
Similarity matrix:



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%

CATH Matching Method

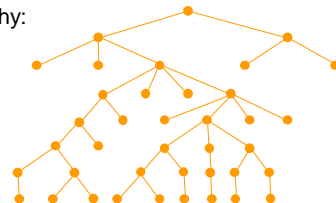


Distance measure is proximity in CATH hierarchy

- $D(A,B)$ = least #levels to common ancestor in hierarchy for any pair of contacting chains

CATH hierarchy:

- Class
- Architecture
- Topology
- Homology
- S35 (Family)
- S95
- S100



CATH Matching Method

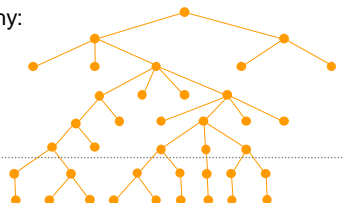


Distance measure is proximity in CATH hierarchy

- $D(A,B)$ = least #levels to common ancestor in hierarchy for any pair of contacting chains

CATH hierarchy:

- Class
- Architecture
- Topology
- Homology
- S35 (Family)
- S95
- S100



CATH Matching Method

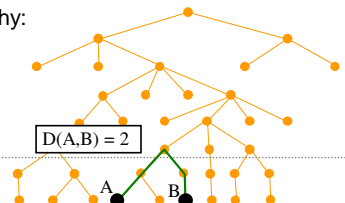


Distance measure is proximity in CATH hierarchy

- $D(A,B)$ = least #levels to common ancestor in hierarchy for any pair of contacting chains

CATH hierarchy:

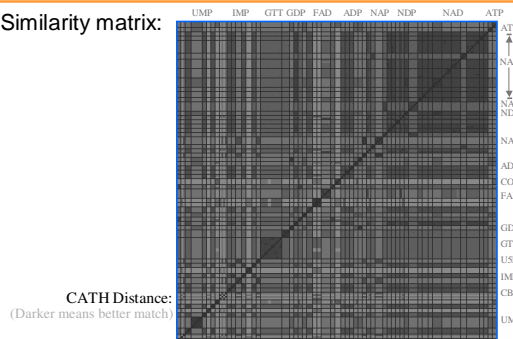
- Class
- Architecture
- Topology
- Homology
- S35 (Family)
- S95
- S100



CATH Matching Results



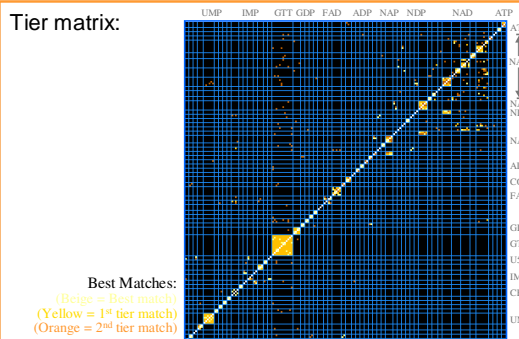
Similarity matrix:



CATH Matching Results



Tier matrix:



CATH Matching Results



Classification rate: When Smith-Waterman \geq 500:

FASTA = 68%	FASTA = 80%
CE = 65%	CE = 72%
CATH = 58%	CATH = 65%
Random = <1%	Random = <1%

When Smith-Waterman < 500:

CE = 53%
CATH = 44%
FASTA = 44%
Random = <1%

SCOP Matching Results



Classification rate: When Smith-Waterman \geq 500:

FASTA = 68%	FASTA = 80%
CE = 65%	CE = 72%
SCOP = 64%	SCOP = 72%
CATH = 58%	CATH = 65%
Random = <1%	Random = <1%

When Smith-Waterman < 500:

CE = 53%
SCOP = 47%
CATH = 44%
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