

Recognition of Ligand Binding Sites with Templates

Thomas Funkhouser
Princeton University
CS597A, Fall 2005

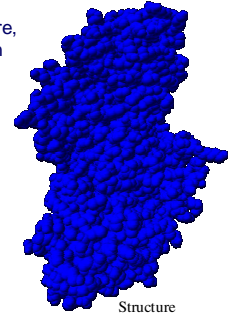
Introduction

Goal:

- Given a sequence & structure, predict its molecular function

```
STAGKVRCKAAVLVIEEKPFSEIEVENAPPEAHEVRKRMVAATGCRSDQ
HVVSGILVPLPLVAGHEAGVGSKEGVTYRPEKRVLPFTQCGRCK
RVCKRPFENCLKNDLSMPKGTMDGTSRFPCKRGRPHHFLGTSFTRQYT
VYREINAKEDIAAPLERVCLERQSTQYGAIVYAKVYRQSTQAVPEL
GGVGLSVMGCKAAGAARHIGVDRNKKDKPAKAEVGAETENPQDYKPKPI
QEVLTETMSNGGIDFSEVHRLDITDRTALKCKQEAIGYVSVVAVPPDSN
LANNPMLLSQRTWKGAHGGFKSDVSKLVAPRMAKRFALDPLDTHL
PFKINIEGFDLRSRGRITLTF
```

Sequence



Structure

1hd

Sequence Motifs

Recognize local patterns (motifs) in sequences indicative of specific functions

Example: Prokaryotic glutathione synthase ATP-binding domain

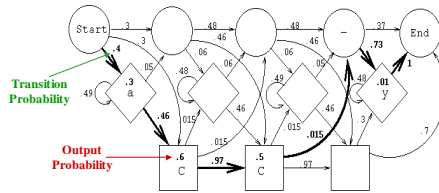
Sequence ID	start	end	weight
1 GSHB_DMSO	123	150	1.00
2 GSHB_NEMO	175	184	0.75
3 GSHB_MTEU	120	147	0.80
4 GSHB_SVNE	3	30	3.00
5 GSHB_JUCAP	4	45	4.00
6 GSHB_WGHR	11	44	11.00
7 GSHB_LJKAP	3	46	3.00
8 GSHB_SVNP	5	47	5.00
9 GSHB_SVNT	2	48	2.00
10 GSHB_PSEK	124	161	1.00
11 GSHB_PSEK	124	161	1.00
12 GSHB_PSEK	124	161	1.00
13 GSHB_PSEK	124	161	1.00
14 GSHB_PSEK	124	161	1.00
15 GSHB_PSEK	124	161	1.00
16 GSHB_PSEK	124	161	1.00
17 GSHB_PSEK	124	161	1.00
18 GSHB_PSEK	124	161	1.00
19 GSHB_PSEK	124	161	1.00
20 GSHB_PSEK	124	161	1.00

ProDom

<http://prodes.toulouse.inra.fr/prodom/>

Sequence Motifs

Recognize local patterns (motifs) in sequences indicative of specific functions



Hidden Markov Model

Sequence Motifs

Many tools match query sequences against sequence motifs

Examples:

- InterPro
- PROSITE
- PRINTS
- PFam-A
- TIGRFAM
- PROFILES
- PRODOM

ProFunc: InterProScan results for 1gsk

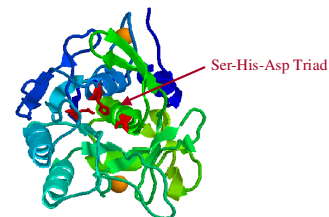
ID	Scan code	Residue range	Residue	Motif name
1	InterPro	1-213	Δ	gshA_Cys glutathione synthase
2	InterPro	1-100	Δ	Prokaryotic glutathione synthase (N-term)
3	InterPro	138-200	Δ	Prokaryotic glutathione synthase (ATP-gr)
4	InterPro	138-200	Δ	ATP_GSHF
5	InterPro	138-200	Δ	OXFOL_SHFL_Q323P1
6	InterPro	138-200	Δ	OXFOL_SHFL_Q323P1
7	InterPro	283-306	Δ	OXFOL_SHFL_Q323P1
8	InterPro	138-200	Δ	OXFOL_SHFL_Q323P1
9	InterPro	138-200	Δ	OXFOL_SHFL_Q323P1
10	InterPro	138-200	Δ	OXFOL_SHFL_Q323P1

InterPro

[Zdobnov01]

Structural Motifs

Recognize local patterns (motifs) in structures indicative of specific functions



Subtilisin (B. amyloliquefaciens)

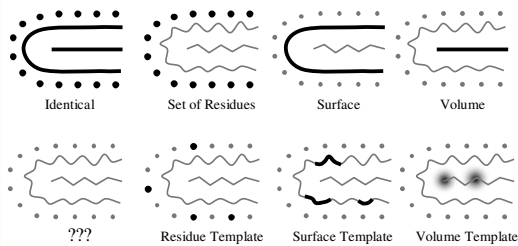
<http://chemistry.umeche.maine.edu/>

Templates



Key idea:

- Encode only the key aspects of the pattern
- Eliminate the noise when matching

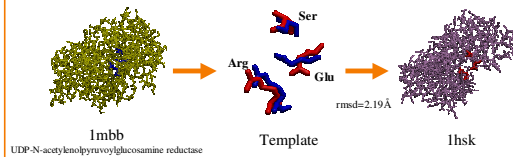


Templates



Methodology:

- Build a structural motif per class
- Search for each structural motif in novel protein
- Report statistically significant "hits"



Slide courtesy of James Watson

Outline



Introduction
 Template construction
 Template search
 Results
 Discussion

Outline



Introduction
Template construction ←
 Template search
 Results
 Discussion

Template Construction



Possible methods:

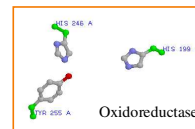
- Human annotation
- Statistical analysis
- All patterns in protein

Template Construction



Possible methods:

- ~~Human annotation~~
- Statistical analysis
- All patterns in protein



The screenshot shows the Catalytic Site Atlas (CSA) interface. It displays a search result for a protein structure, including the protein name, sequence, and a list of residues. The interface is used for searching and analyzing protein structures.

The Catalytic Site Atlas (CSA) contains templates manually curated from scanning the literature

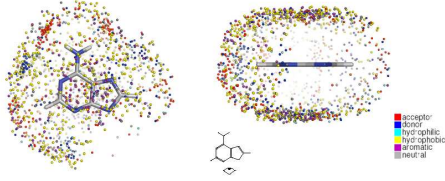
[Porter04]

Template Construction



Possible methods:

- Human annotation
- Ø Statistical analysis
- All patterns in protein



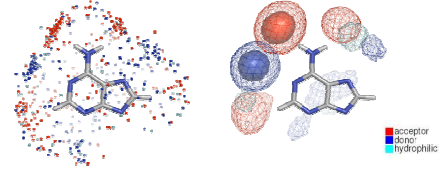
Distribution of atom types contacting adenine rings in PDB [Stockwell05]

Template Construction



Possible methods:

- Human annotation
- Ø Statistical analysis
- All patterns in protein



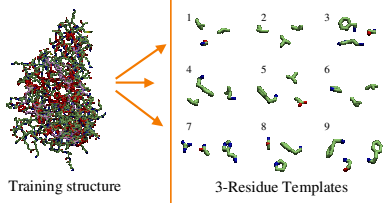
Distribution of polar atoms contacting adenine rings in PDB [Stockwell05]

Template Construction



Possible methods:

- Human annotation
- Statistical analysis
- Ø All patterns in protein



Slide courtesy of James Watson

Outline



Introduction

Template construction

Template search ←

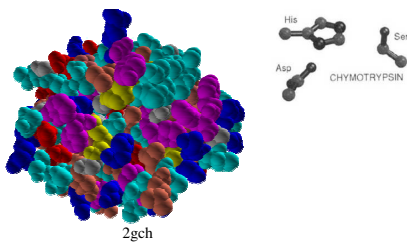
Results

Discussion

Template Search



Does a given pattern appear in the protein?

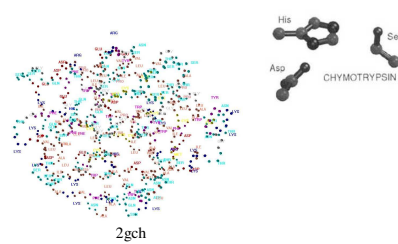


[Wallace97]

Template Search



Does a given pattern appear in the protein?



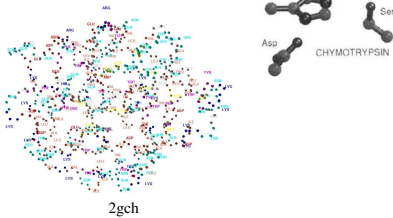
[Wallace97]

Template Search



Challenges:

- Template is subset of structure
- Arbitrary translation
- Arbitrary rotation



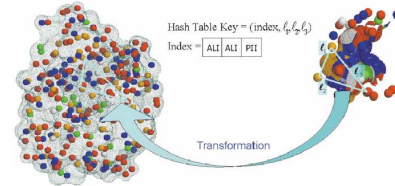
[Wallace97]

Template Search



Methods:

- Geometric hashing
- Association graphs
- Grid correlation



[Shulman-Peleg04]

Outline



Introduction

Template construction

Template search

Results ←

Discussion

Results



Web servers report template hits

b. Probable matches (1.00E-06 < E-value < 0.01)

Hit no.	value	score	[eqqid]	Template id	Matched PDB entry	Longest fitted segment	Seq lengths query/target	Overlap	%-tage	Struc-tural	
1	2.01E-04	257.34	10/1	[21]	ATPc0011	1dV2: The structure of biotin carboxylase, mutant c288k, complexed with asp.	89/186	314/650	378	16.56%	94.1%
2	6.15E-04	242.38	11/1	[20]	POBc0003	1f0v: Complex of d-ala d-ala ligase with asp and a phosphorylphosphonate	138/219	314/506	360	14.71%	99.2%

c. Possible matches (0.01 < E-value < 0.10)
 ... no hits in this category

d. Long shots (0.10 < E-value < 10.0)

Hit no.	value	score	[eqqid]	Template id	Matched PDB entry	Longest fitted segment	Seq lengths query/target	Overlap	%-tage	Struc-tural	
3	1.514	137.69	8/1	[17]	l4Pc0004	1x2b: Inositol 1,3,4-trisphosphate 5b-kinase in complex with mg2-hadpims 1,2,4,6p4	76/176	314/511	375	16.06%	91.4%
4	12.601	109.28	9/1	[14]	PHYc0003	1d4e: D-alanyl-d-lactate ligase	76/166	314/541	390	19.43%	91.3%

[Laskowski05]

Results



Common problems:

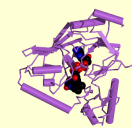
- Too many false positives
- Top hit rarely the correct hit – even in “obvious” cases
- Use of rmsd rarely discriminates true from false positives
 - § Local distortion in structure may give a large rmsd

Need a way to encode more information in templates to avoid false positives

Example Query

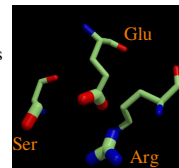


PDB code: **1hsk**
 UDP-N-acetylenolpyruvoylglucosamine reductase (MURB)
 E.C.1.1.1.158



Contains the 3D template that characterises this enzyme class

Sequence identity to template's representative structure (1mbb) is 28%



Slide courtesy of James Watson

Example Query



Hits for 1hsk template:

Hit	E.C number	Rmsd	Enzyme
1.	E.C.1.3.99.2	0.76Å	Acyl-CoA dehydrogenase
2.	E.C.4.2.1.20	0.76Å	Tryptophan synthase α -subunit
3.	E.C.3.2.1.73	1.19Å	Glycosyl hydrolases, family 17
4.	E.C.3.2.1.73	1.21Å	Glycosyl hydrolases, family 16
5.	E.C.4.1.2.13	1.25Å	Fructose-bisphosphate aldolase (class I)
...
102.	E.C.1.1.1.158	2.19Å	UDP-N-acetylmuramate dehydrogenase
...
386.	...	3.94Å	...

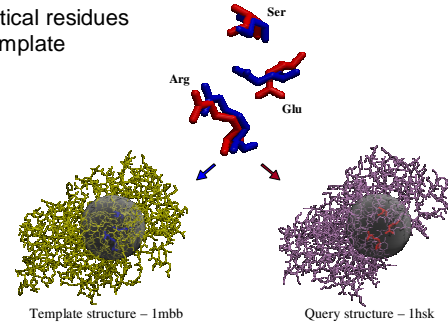


Slide courtesy of James Watson

Example Query



Identical residues in template



Template structure - 1mbb

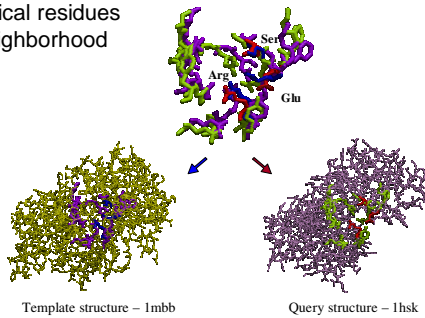
Query structure - 1hsk

Slide courtesy of James Watson

Example Query



Identical residues in neighborhood



Template structure - 1mbb

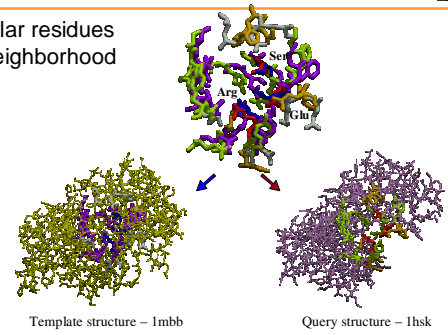
Query structure - 1hsk

Slide courtesy of James Watson

Example Query



Similar residues in neighborhood



Template structure - 1mbb

Query structure - 1hsk

Slide courtesy of James Watson

Example Query



Hits for 1hsk template & neighborhood:

Hit	E.C number	Rmsd	Score	Enzyme
1.	E.C.1.1.1.158	2.08	209.1	UDP-N-acetylmuramate dehydrogenase
2.	E.C.3.2.1.14	2.13	146.0	Chitinase A chitodextrinase 1,4-beta-poly-N-acetylglucosaminidase coly-beta-glucosaminidase
3.	E.C.3.2.1.17	1.92	142.4	Turkey lysozyme
4.	E.C.3.2.1.17	1.89	138.7	Hen lysozyme
5.	E.C.3.5.1.26	1.47	132.3	Aspartylglucosylaminidase
6.	E.C.3.2.1.3	1.54	131.1	Glucan 1,4-alpha-glucosidase

Slide courtesy of James Watson

Discussion



?