

Challenges and Solutions to Protein Structure Prediction in the Genomic Era



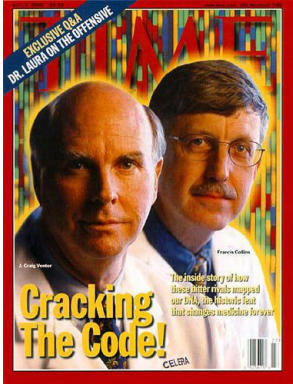
Jianlin Jack Cheng, PhD

Computer Science Department

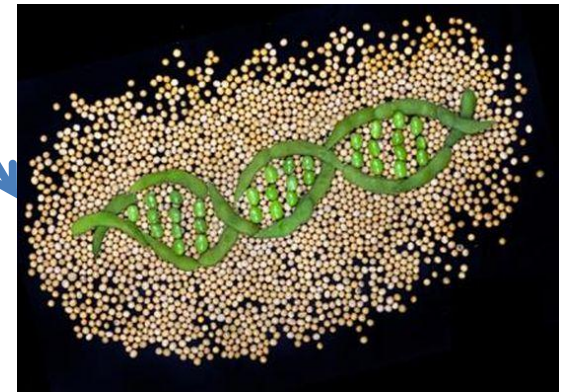
University of Missouri, Columbia

Oct. 5, 2011

The Genomic Era



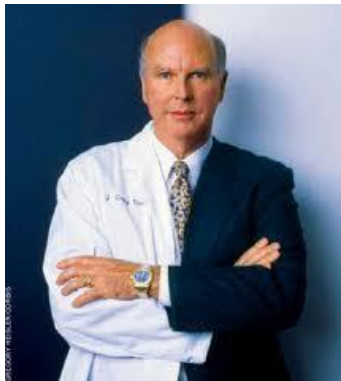
Collins, Venter, Human Genome, 2000



NIH \$1,000 Genome Project

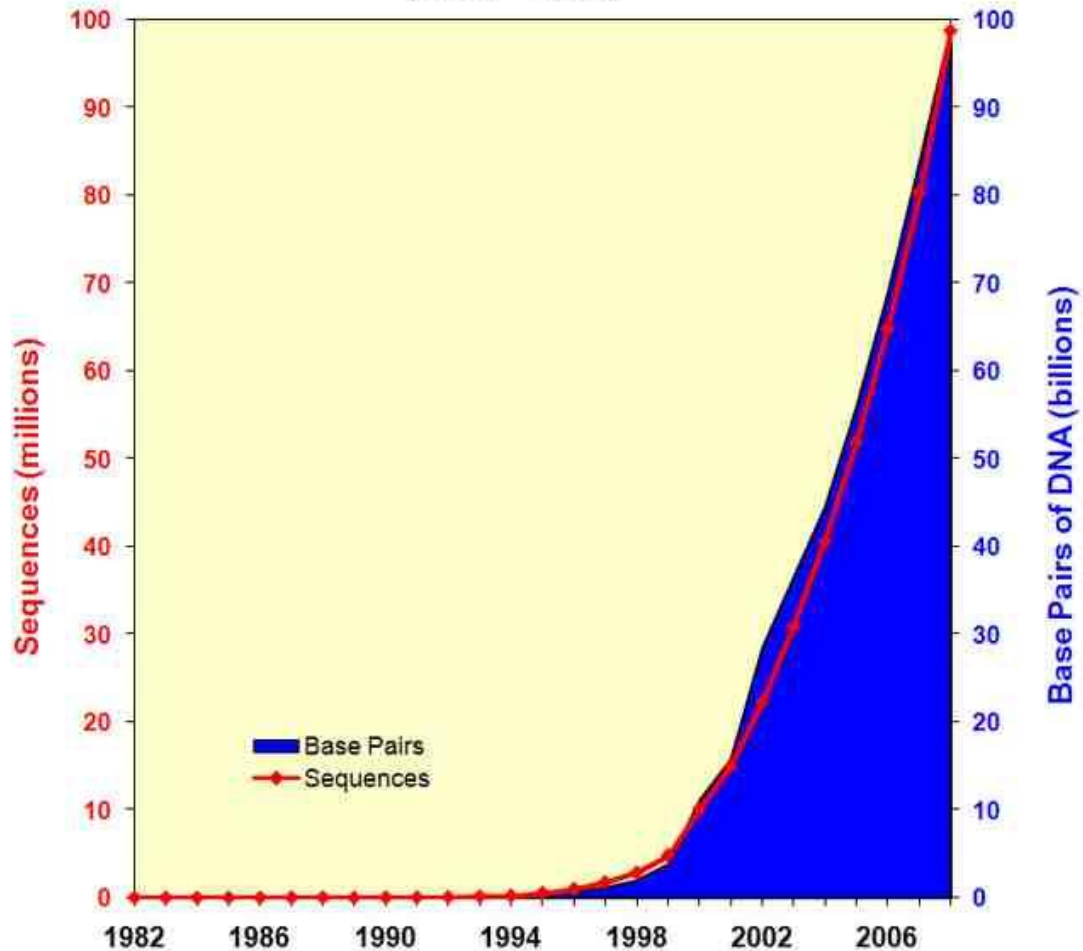


- **Personal Computer in 1970s**
- **Personal Genome in 2010s**

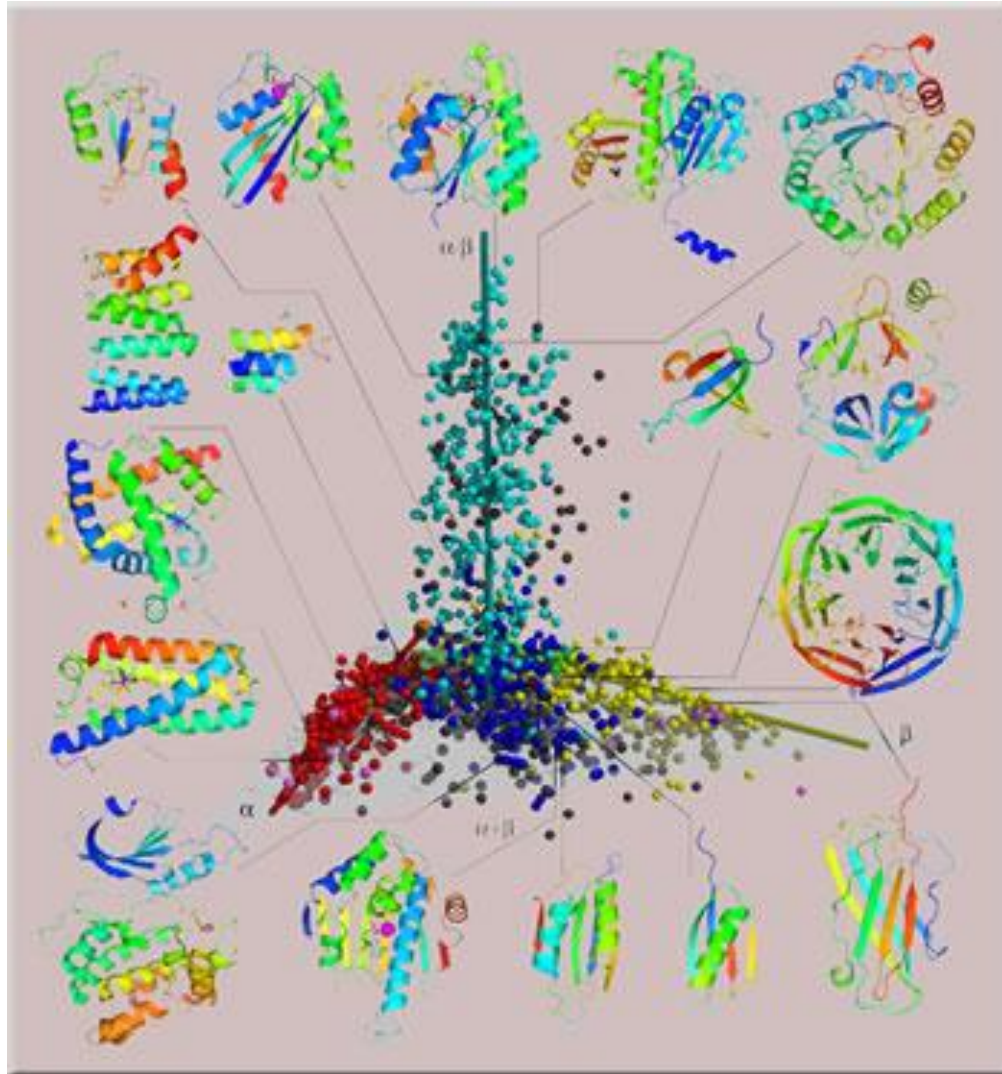


Growth of Protein Sequences

Growth of GenBank
(1982 - 2008)



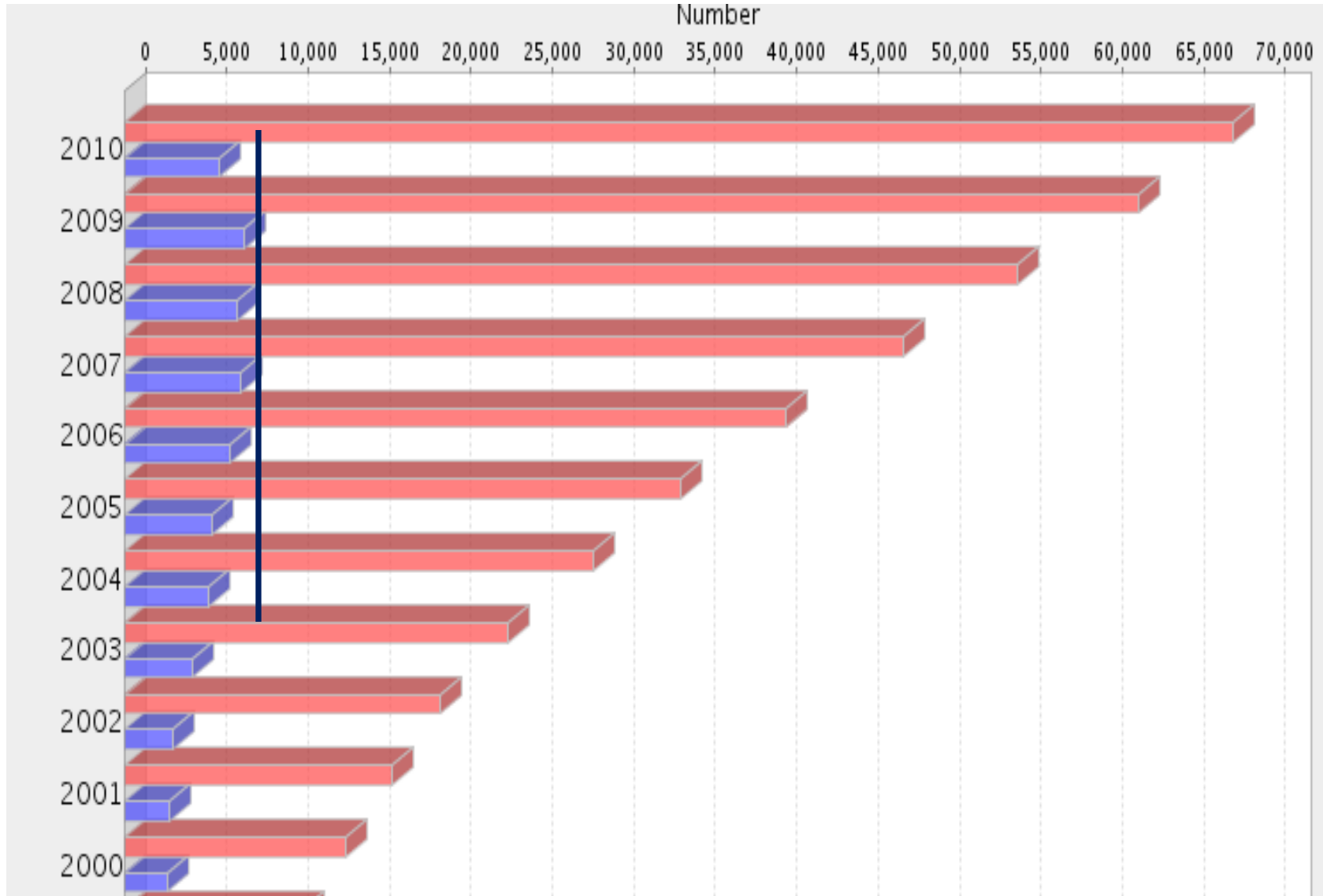
NIH Structural Genomics Project



Chothia, Nature, 1992

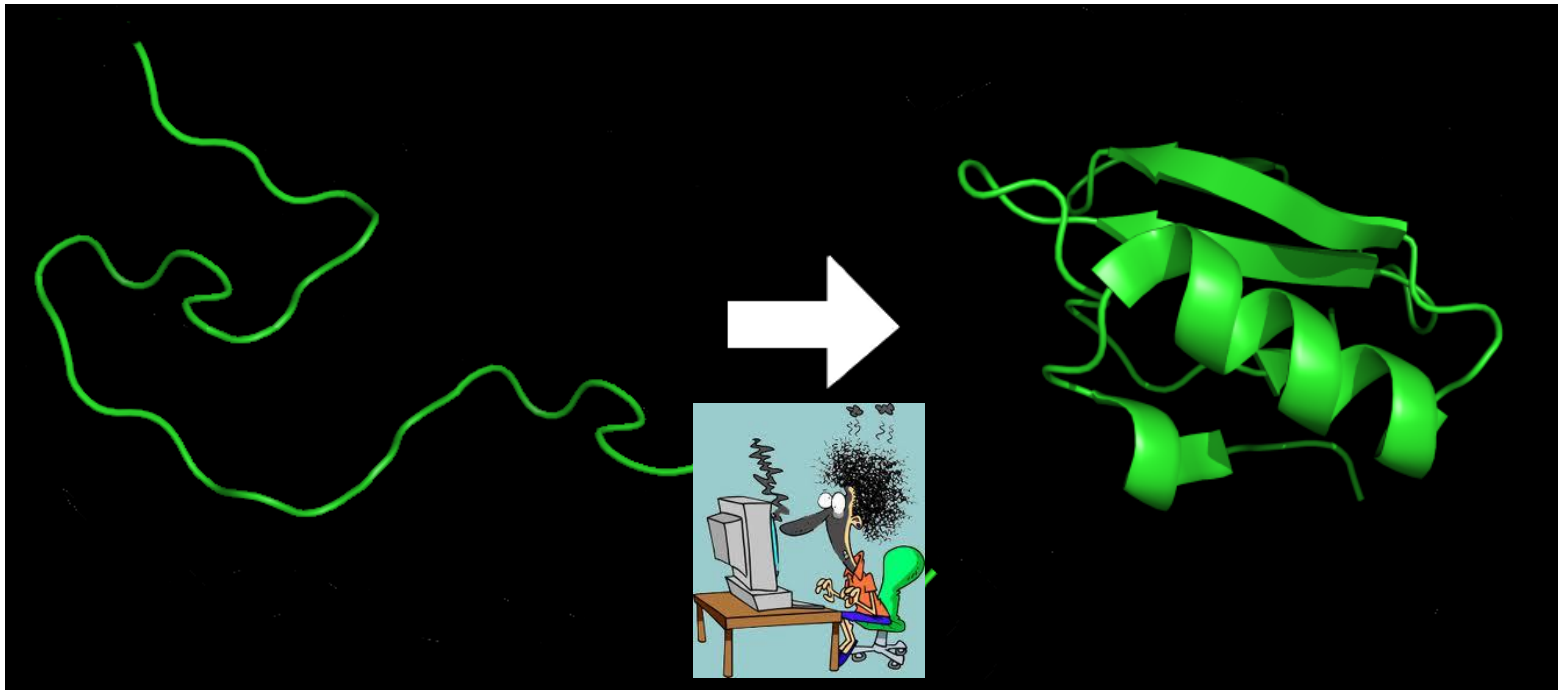
Protein Sequence and Fold Space

Growth of Protein Structures in PDB

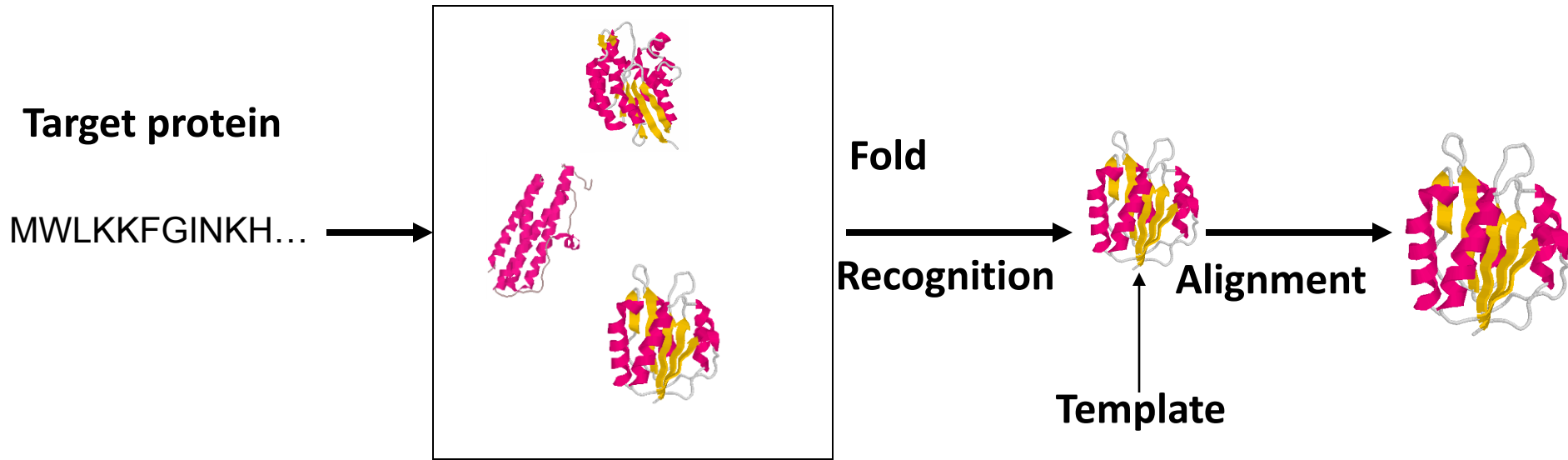


Computational Protein Structure Folding / Prediction

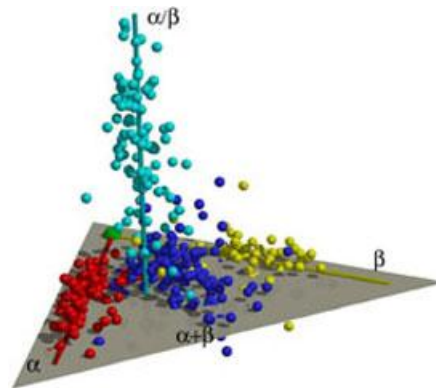
Structure = $f(\text{sequence})$? \leftrightarrow $E = MC^2$



Template-Based Approach



Protein Data Bank



Protein Structure Space

TARGET

ASILPKRLFGNCEQTSDEGLK
IERTPLVPHISAQNVCLKIDD
VPERLIPERASFQWMNDK

TEMPLATE



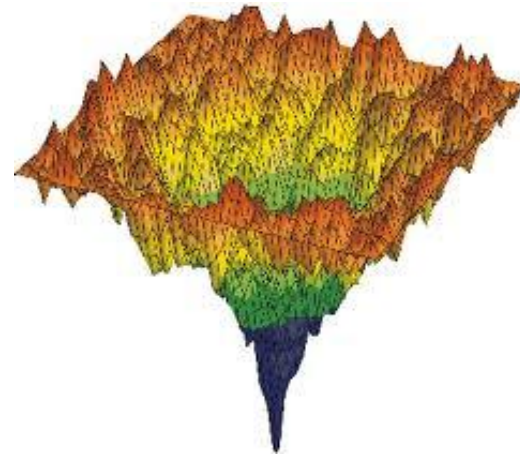
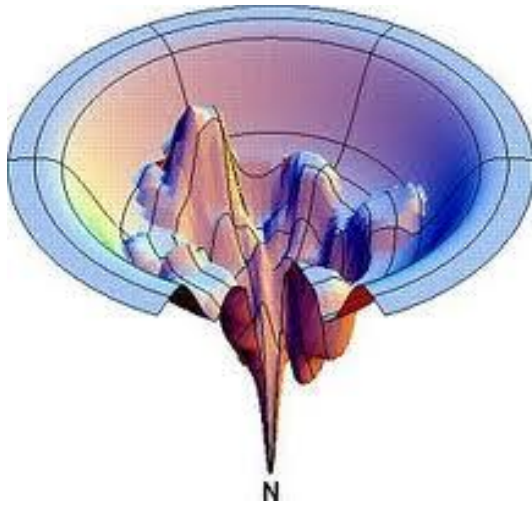
ASILPKRLFGNCEQTSDEGLK IERTPLVPHISAQNVCLKIDD VPERLIPE
MSVIPKRLYGNCEQTSEEAIRIEDSPIV---TADLVCLKIDEIPERLVGE



Modeller

A. Fisher, 2005

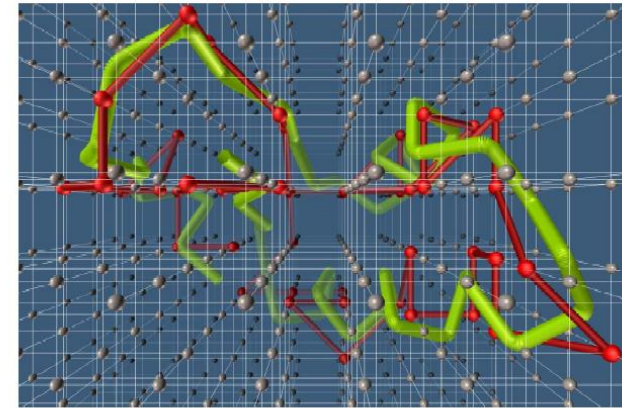
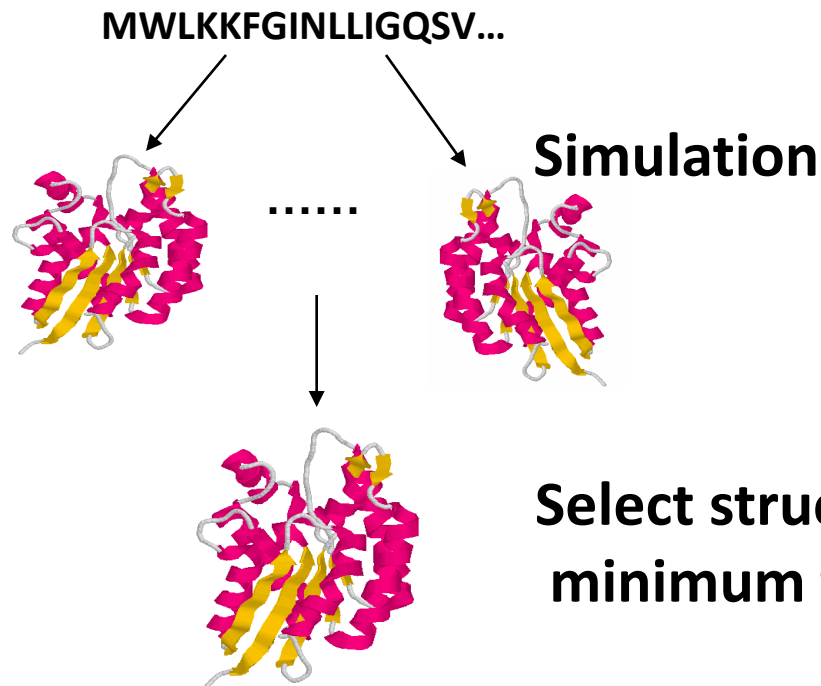
Template-Free Protein Structure Prediction



Dill & Chan, 1997

Template-Free Approach

Methods: molecular dynamics, fragment assembly, contact-based modeling



Major Challenges in Protein Structure Prediction

- **Select best templates?**
- **Generate best alignments?**
- **Generate best models?**
- **Select best models?**
- **Model refinement?**

Major Challenges in Protein Structure Prediction

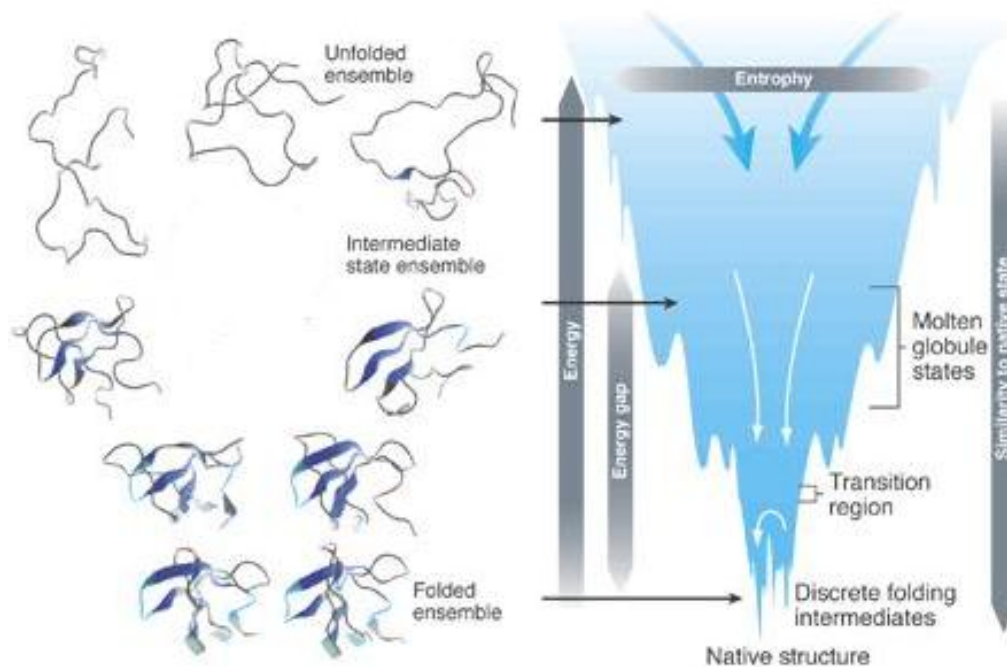
- Select best template
- Generate

models?

model refinement?

Protein Conformation Ensemble

- $P(\text{conformation}) \propto P(-\text{energy})$
- Conformation Distribution
- Maximum Likelihood & Maximum a Posterior

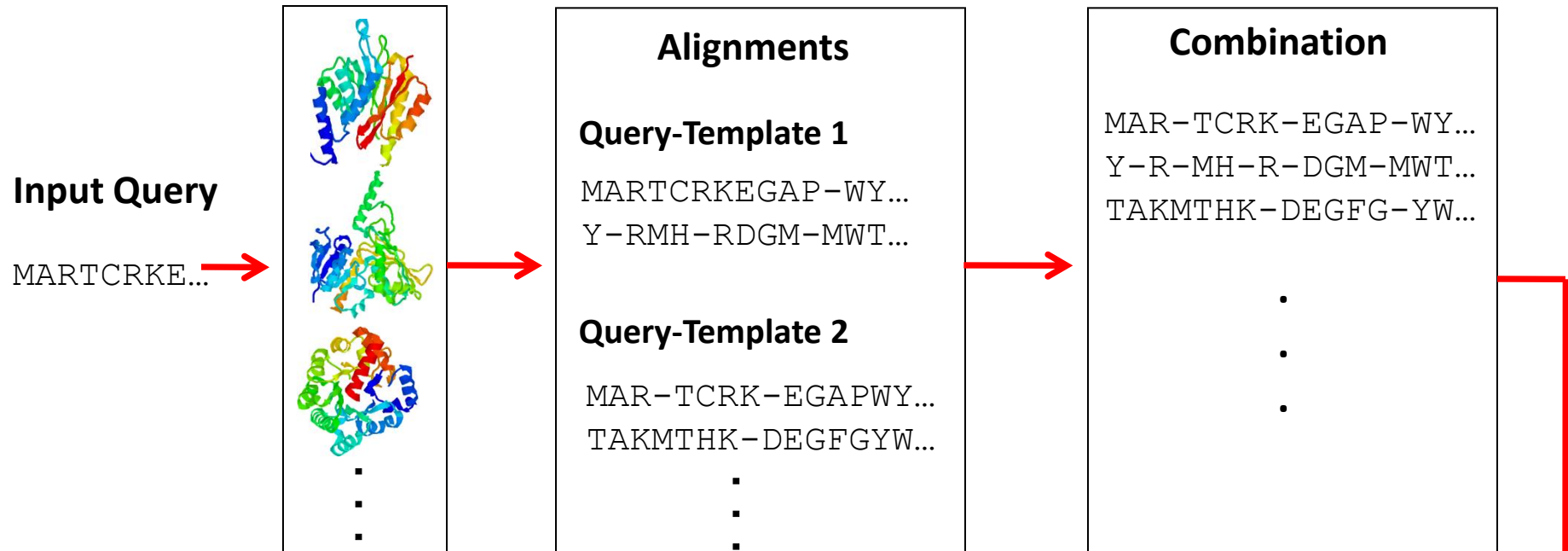


Brooks et al., 2001

A Unified Protein Structure Prediction Pipeline

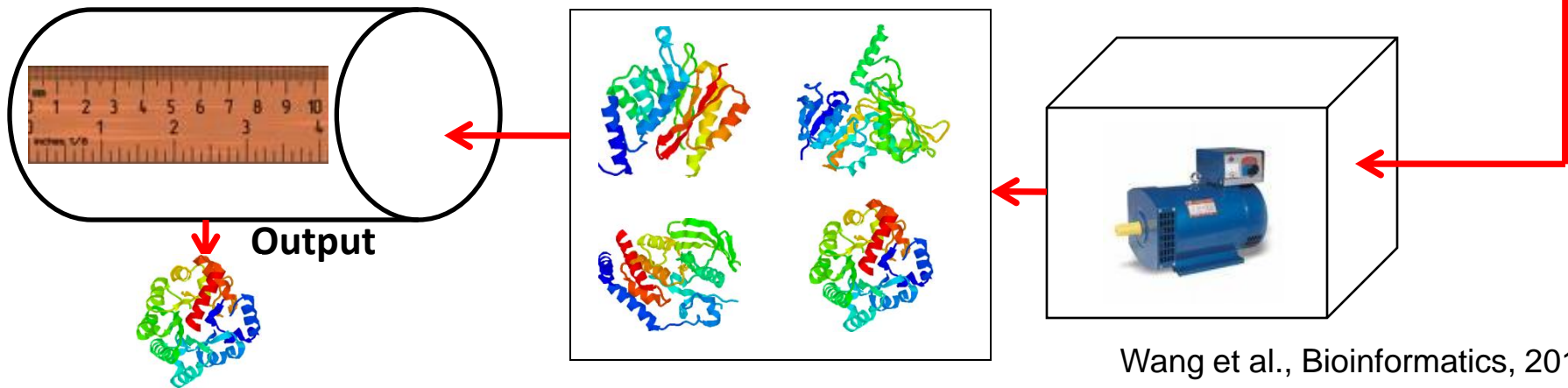
1. Template Ranking

2. Multiple-Template Combination



4. Evaluation & Refinement

3. Model Generation



Template Ranking / Fold Recognition

- **PSI-BLAST (sequence – profile)**
- **SAM (sequence – HMM)**
- **HMMer (sequence – HMM)**
- **Compass (profile – profile)**
- **HHSearch (HMM - HMM)**
- **PRC (HMM-HMM)**
- **FOLDpro (machine learning)**
- **MSACompro (profile-profile)**

Multi-Template Combination in Template and Alignment Space

Query VR-RNNMGMP LI E S S Y H D A L F T L G Y A G D R I S Q M L G M R Y A N N L H D L F L A E G Y Y E A S Q R K R

Temp1 I A H I Y A N N L H D L F L A E G Y Y E A S Q R L F E I E L F G L M G N L S S W V G A
(10^{-80})

Temp2 L L A Q - G R L S E M A G A D A L D V N I Y I D S N G
(10^{-70})

Temp3 Q G T A R D R A W Q L E V E R H R A Q G T S A S F L
(10^{-10})

Temp4 A A N Q L D A M R A L G Y A Q E R Y F E M D L M R R A P A G E L S E L F G A K A V D L K
(10^{-5})



Multi-Template Combination in Template and Alignment Space

Query VR-RNNMGMP LI ESSSYHDA LFTLGYAGDRISQMLGMRYANNLHDLFLAEGYYEASQRKR

Temp1 IAHIYANNLHDLFLAEGYYEASQRLFEIEL-----FGLMGN-----LSSWVGA-----
(10^{-80})

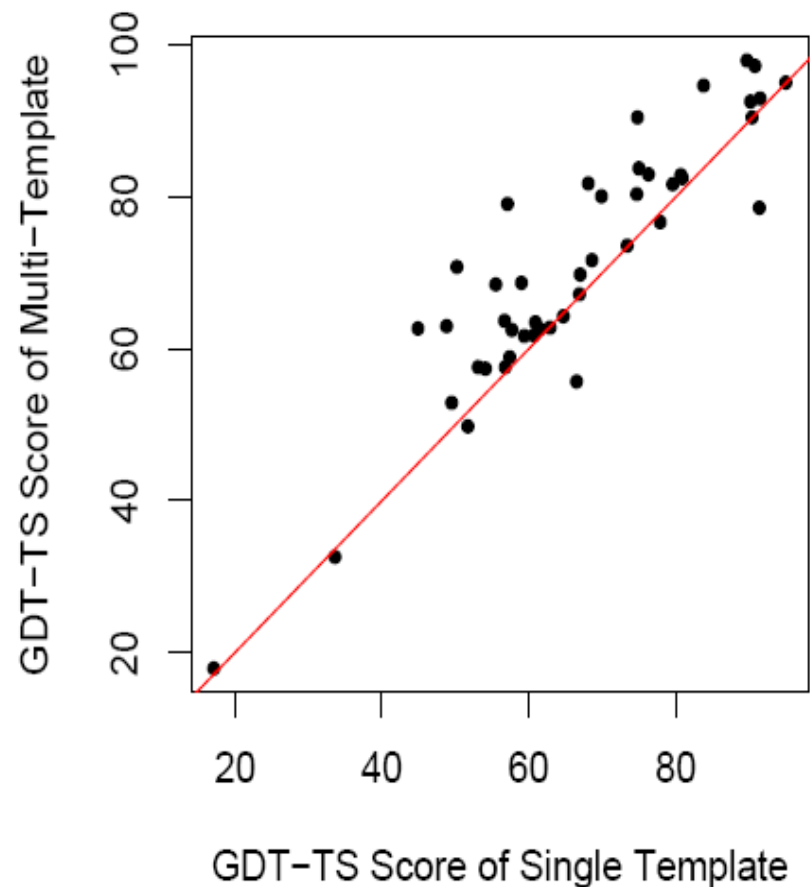
Temp2 LLAQ-GRLSEMAGADALDVNIYIDSNG-----
(10^{-70})

Temp3 -----ARDRAWQLEVERHRAQGTSASFL-----
(10^{-10})

Temp4 -----GAKAVDLK-----
(10^{-5})



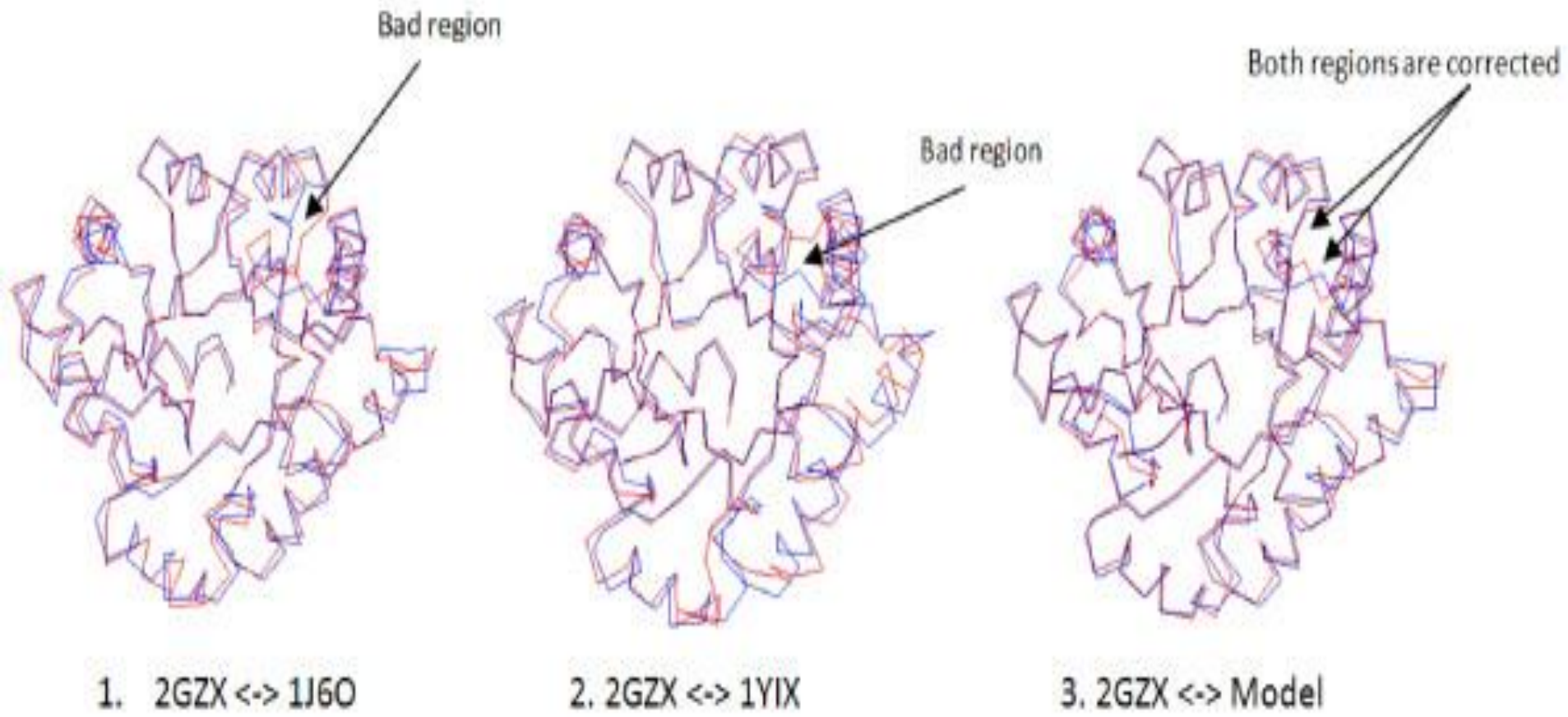
Multi-Template VS Single-Top-Template



Improve 38 / 45 targets

Improvement by 6.8%

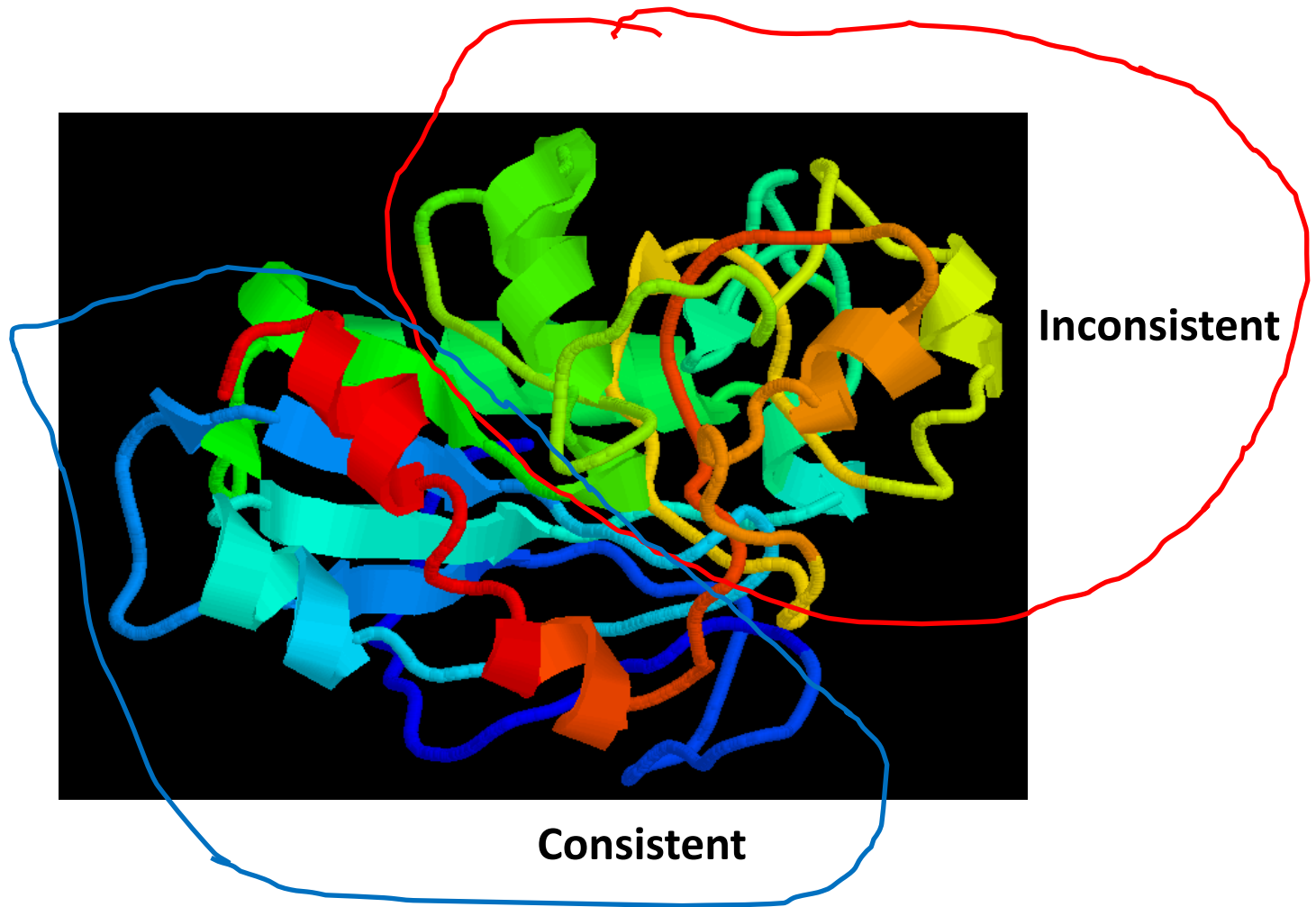
P-value < 10^{-4}



Cheng, BMC Structure Biology, 2008

Advantage: reduce variance of modeling

Problem – Atom Clashes



Combine 49 templates – CASP8 target T0413

Structure Comparison-Guided Multi-Template Combination (CASP9)

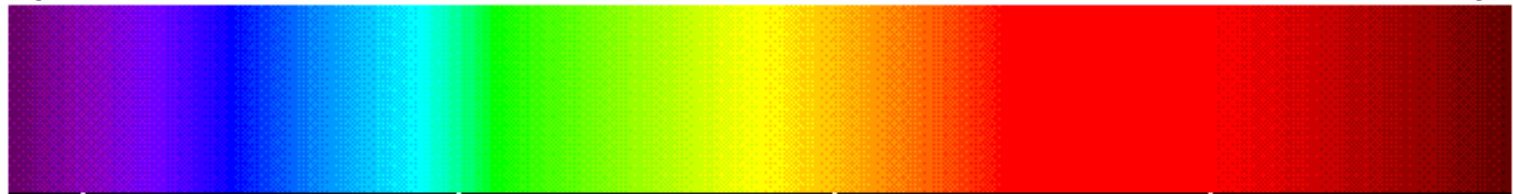
- **Check the structural consistency between query-template alignments**
- **Only combine query-template alignments that are structurally consistent**
- **Completely removed atom clashes in CASP9**

Combination of Template-Free and Template-Based Modeling for Model Generation

100% TBM

50% TBM+50%FM

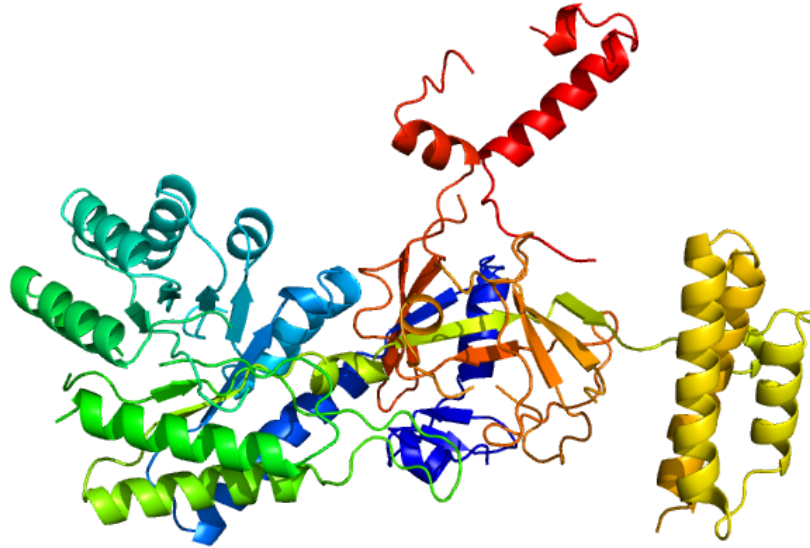
100% FM



Protein Modeling Spectrum

Region Decomposition

A)



B)

Template: 1TWIA

```

-----MLGNDTVEIK-DGRFFI---DGYDA-IELA EK-----FGTPLYVMSEEQIKINYNRYI
EAFKRWEEETG--KEFIVAYAYKANANLAITRLLAKLGC---GADVVS GGEL YIAKLSNVPSK----K
IVFNGNC-KTKEEIIIMGIE---ANIRA-FNVDSISELILINETAKE-LGETANVAFRINPNVNPKTHPK
ISTGLKKNKFGLDVESGIAMKAIKMALE--MEYV-NVVG VHC HIGSQLTDISPFI EETRKMDFVVELK
E-----E-GI-EIEDVNLGGGLGTPYYKDKOT---PTOKDLADATINTMIKYKD--KVEMPNI TLEPG
RSLVATAGYLLGKVHHIKETPVT-----
-----KWVMIDAGMNDMMRP-AMYE
AY-HHIINCK----VKN----EKEVVS IAGGLCESSDVFGRDR-----ELD-KVEVGD---VLAIFD
VGAYGISMAN-NYNARGRPRMVLTS--KKG-V--FLIRERETYADLIAK-----
-----
    
```

Query: T0547

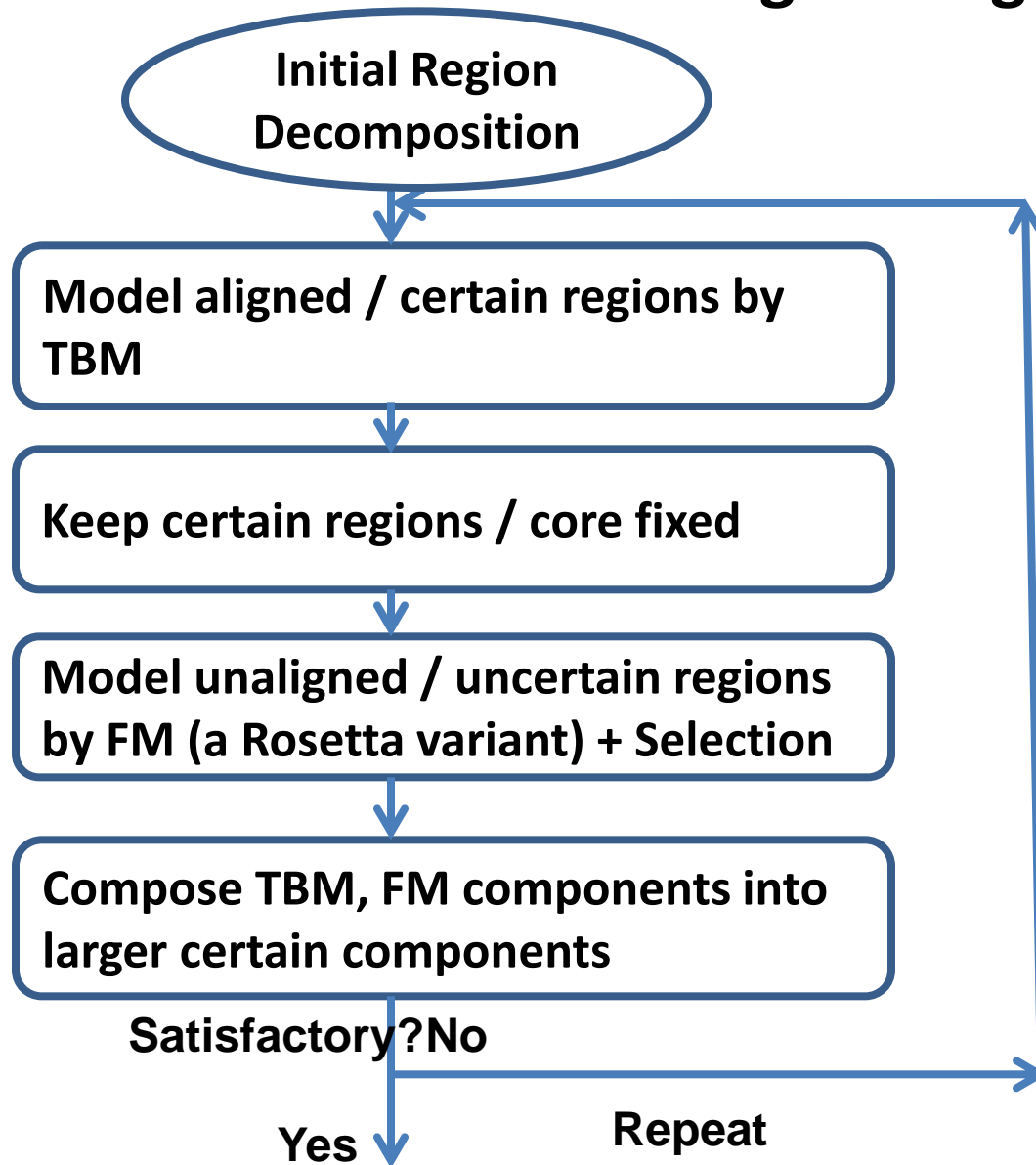
```

MMDYGIDIWGNENFIK-NGKVCINYEKKPAI-IDIVLELR----DDGYKG PLLLRFPHLIQKQIE NIY
GNFNKARKEFGYKGGFNAVYPLKVNQYPGFVKNLVKLGKDNYGLEAGSKAELLAMAYNNEGA---P
ITVNG-F-KDRELINIGFIAAEMGHNITLTIEGLNEVEAIIIDIAKERFKPKPNIGLRVRLHSAGVGI-W
AKSGGINSKFGLTSTE--LIEAVNLLKE--NKLLEOETMIHEHLGSOTTETHPLKKALNEAGNTITELR
K-----M-GAKNIKATNLGGGLAVEYSOFKNEKSRNYTLREYANDVVFETLKNIAEOKKDL E PDI FIE SG
REVAANHAVL TAPVLELFSOEYAENKLT LKKONPKLID-ELYDL YKSTI-KPSNALEYLHDSIHLESI
LTLFDLGYVDLQDRSNAEILTHLITKKAILLGDKQN PADL LAIQDEVQERYLVNFSLFQSMPOFWGLE
QN-FPIMPLD----RLD--EEPTRASIWIDITCDS DGEISYSKD---KPLFLH-DVDVEKENFLGFFL
VGAYOEVLGM-KHNL FTHPT EATISINEKG-YEVEGII EAQSILDTLEDL DYDIHAIMDILNERISNSK
LVNDKQKKHILGELYLFLNDNGYLKSIGV
    
```

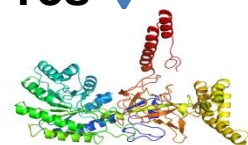
Domain 3

Domain 4

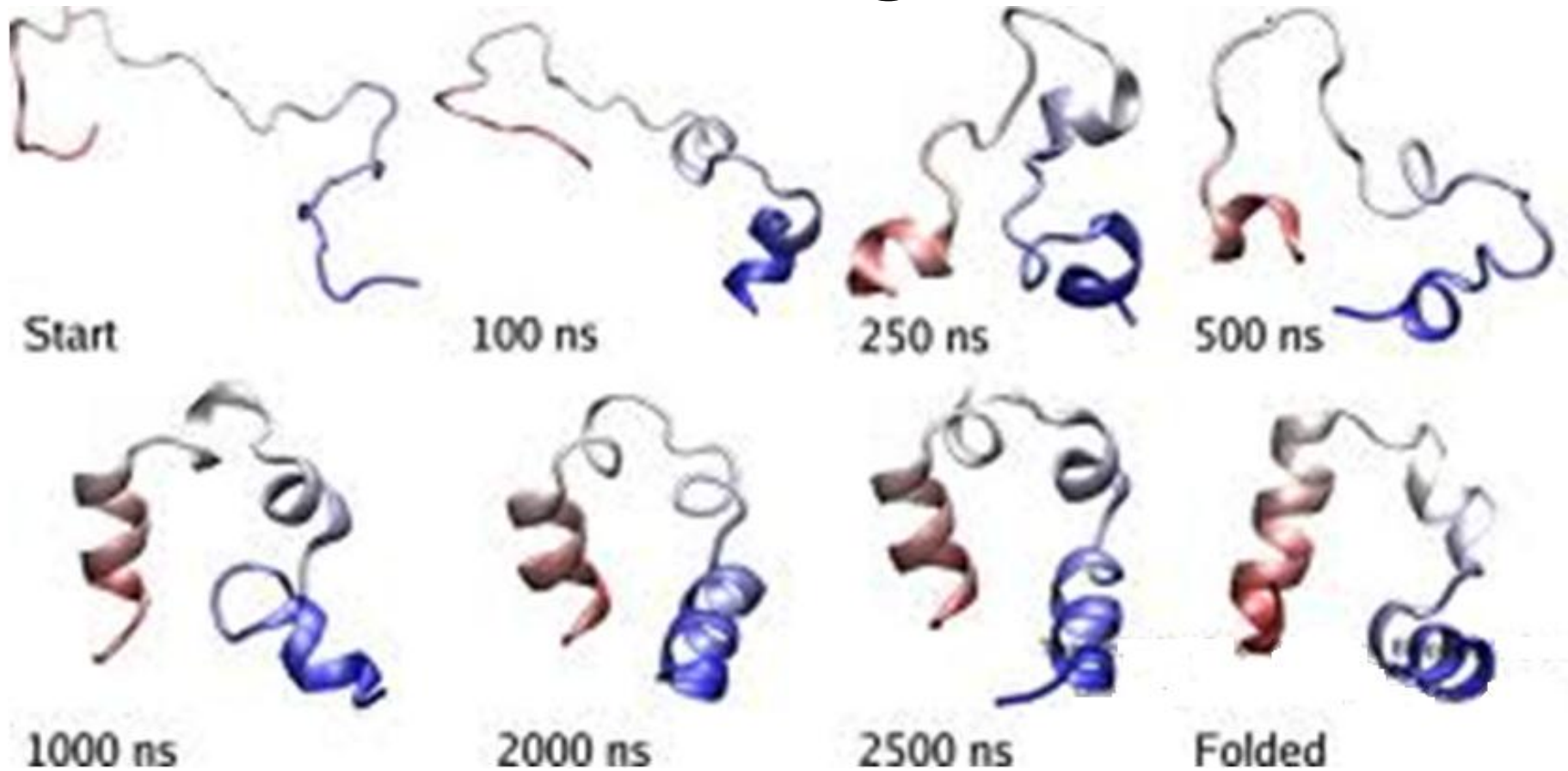
Recursive Protein Modeling – Integrate TBM and FM



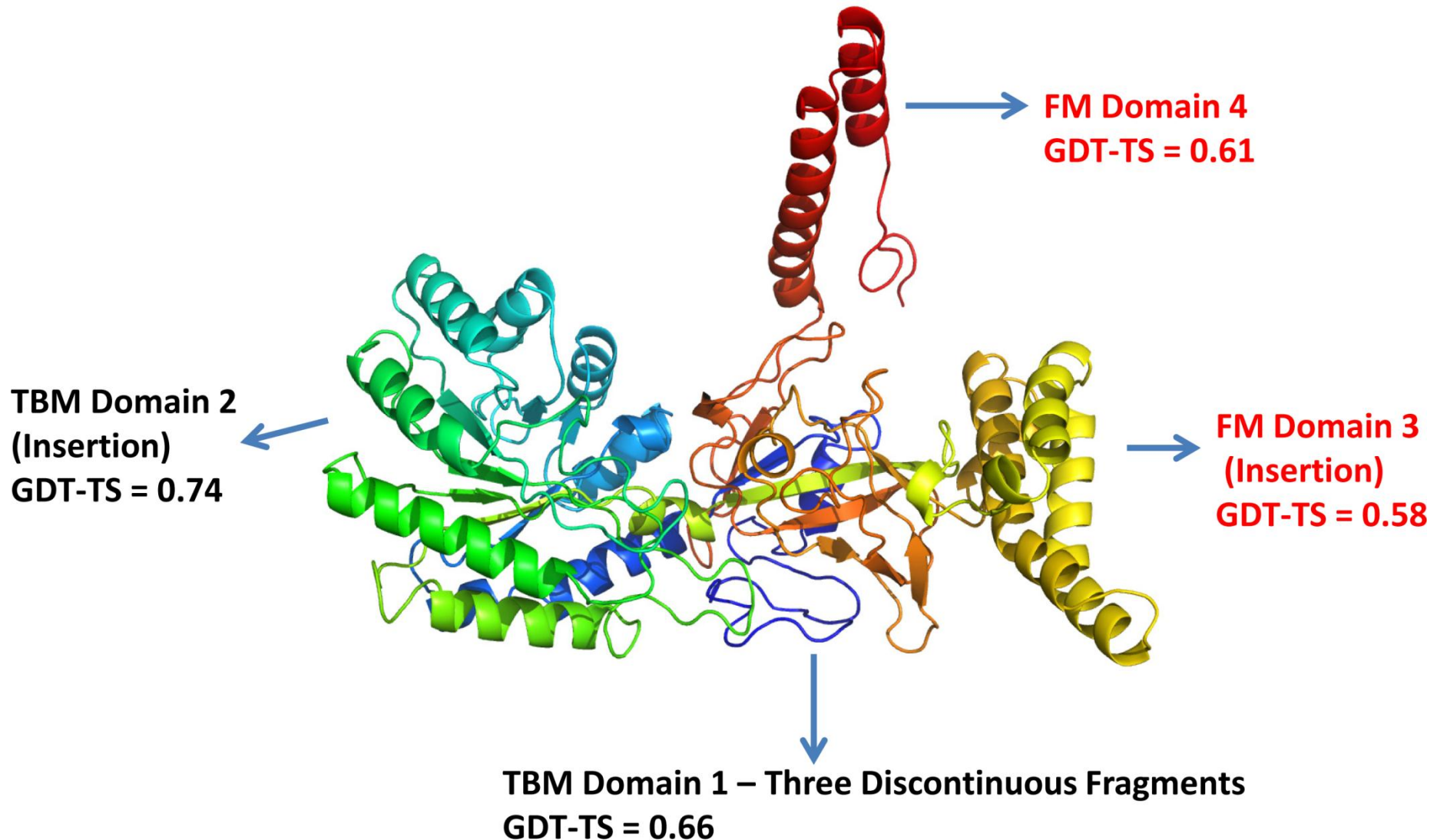
Divide & Conquer



Recursive Modeling Mimics Protein Folding Cascade



Domain Level Integration (CASP9)



Core-Constrained Tail Refinement

>P1;1VI8B

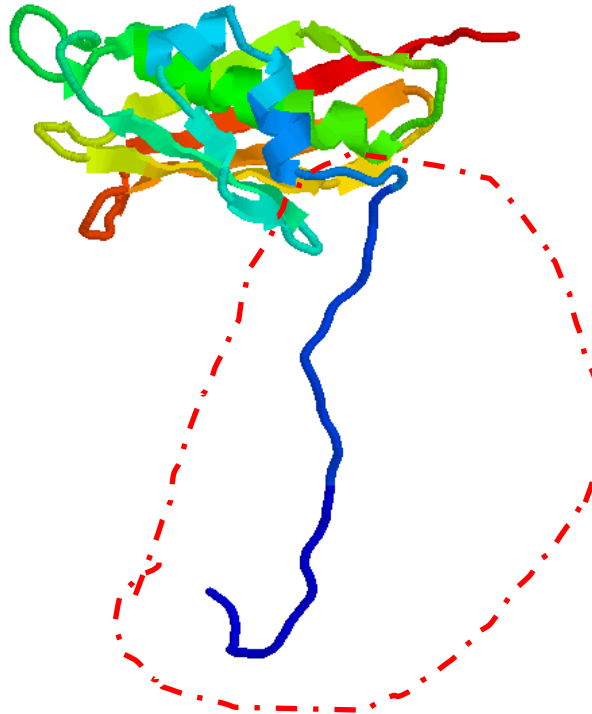
structureX:1VI8B: 1: : 146 : : : :

-----+SLIWKRKITLEALNAMGEGNMVGFDIRFEHIGDDTLEATMPVDSRTKQPFGLLHGGASVVL
AESIGSVAGYLCTEGEQKVVGLEINAMHVRSAAREGRVVRGVCKPLHLGSRHQVWQIEIFDEKGRLLCCSSRLTTAILEGGSHHHH*

>P1;Query

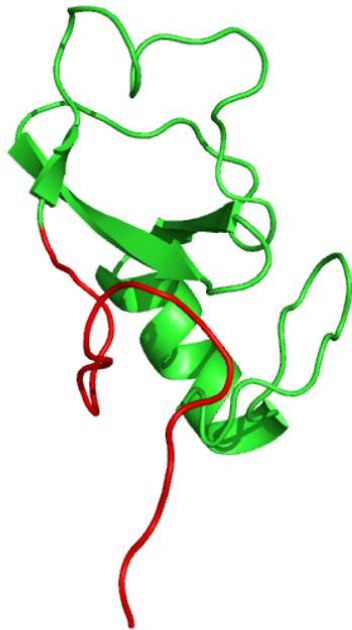
: : : : : : : :

~~MDKRLQQDRIVDKMERFLSTANEEKDV~~LSSIVDGLLAKQERRYATYLASLTQIESQEREDGRFEVRLPIGPLVNNPLNMVHGGITATL
LDTAMGQMVNRQLPDGQSAVTSELNIHYVKPGMGTYLRAVASIVHQKQRIVVEGKVYTDQGETVAMGTGSFFVLRSG-----*



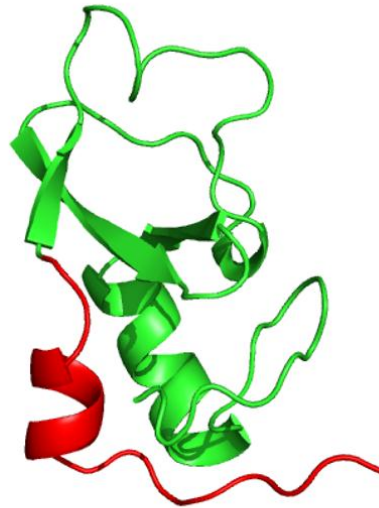
Template-Based + Template-Free & Recursive Modeling (CASP9)

A)



Before tail refinement
GDT-TS = 0.64

B)



After tail refinement
GDT-TS = 0.73

C)



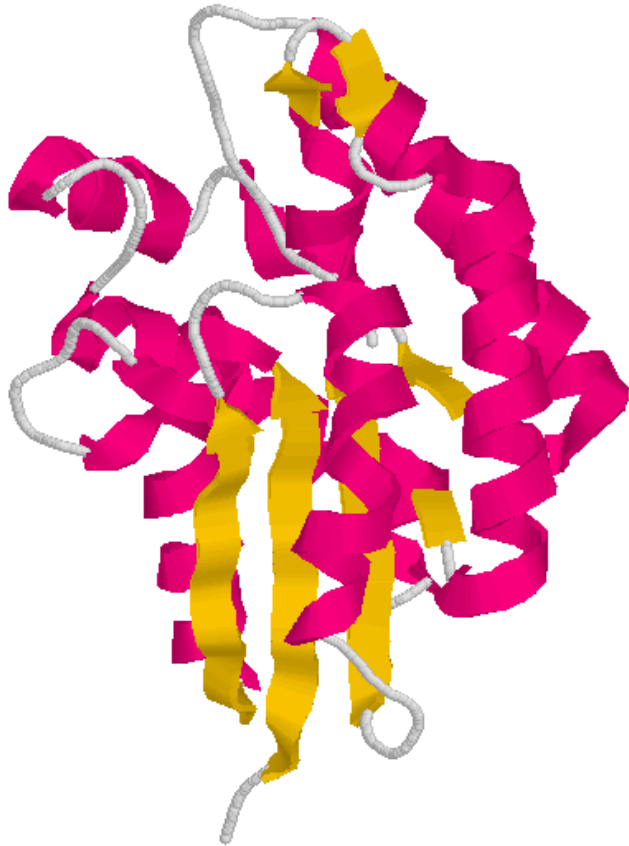
Superposition
Green: model, Blue: structure

Improve Template-Free Modeling Using Contact Restraints

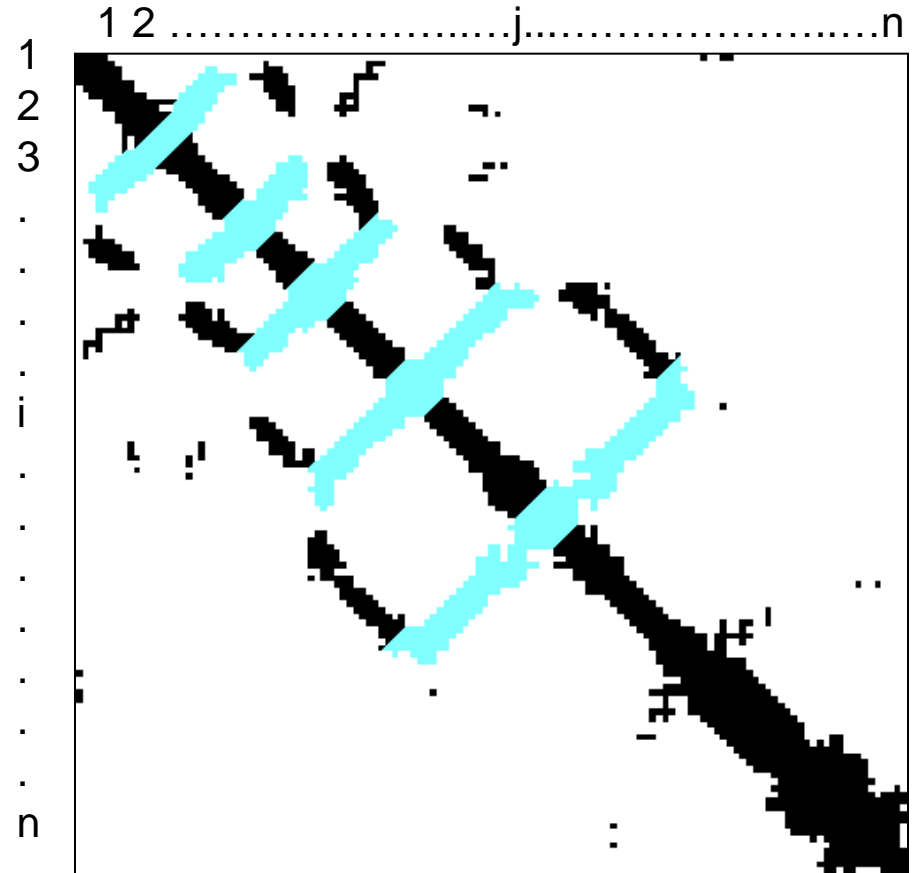
- **Fragment assembly (e.g. Rosetta)**
- **Conformations with good secondary structures and compactness, but often wrong topology**
- **Huge bottleneck (30% proteins)**

Contact Map Prediction

3D Structure



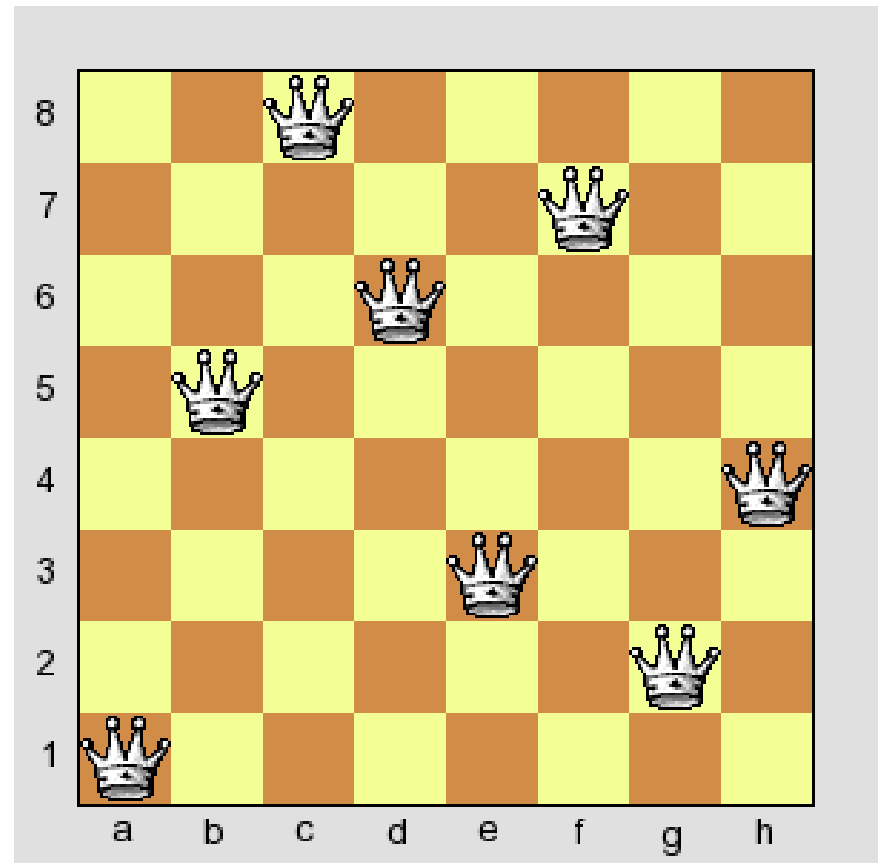
2D Contact Map



Cheng, Randall, Sweredoski, Baldi, 2005
Cheng and Baldi, 2007.
Tegge et al., 2009

Contact-Based Protein Structure Reconstruction

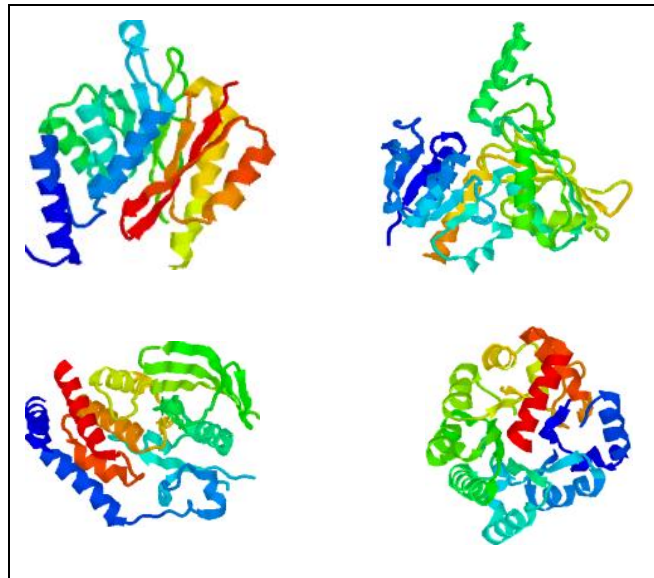
- A constrained optimization problem
- A small number of contacts $L / 10$
- 30% accuracy



Eight Queens Problem

Contact Prediction

- **2D Neural Networks (CASP7, 8, 9)**
- **Support Vector Machine (CASP7, 8, 9)**
- **A Conformation Ensemble Approach (CASP9)**



Contact Voting

CASP9 Contact Predictions

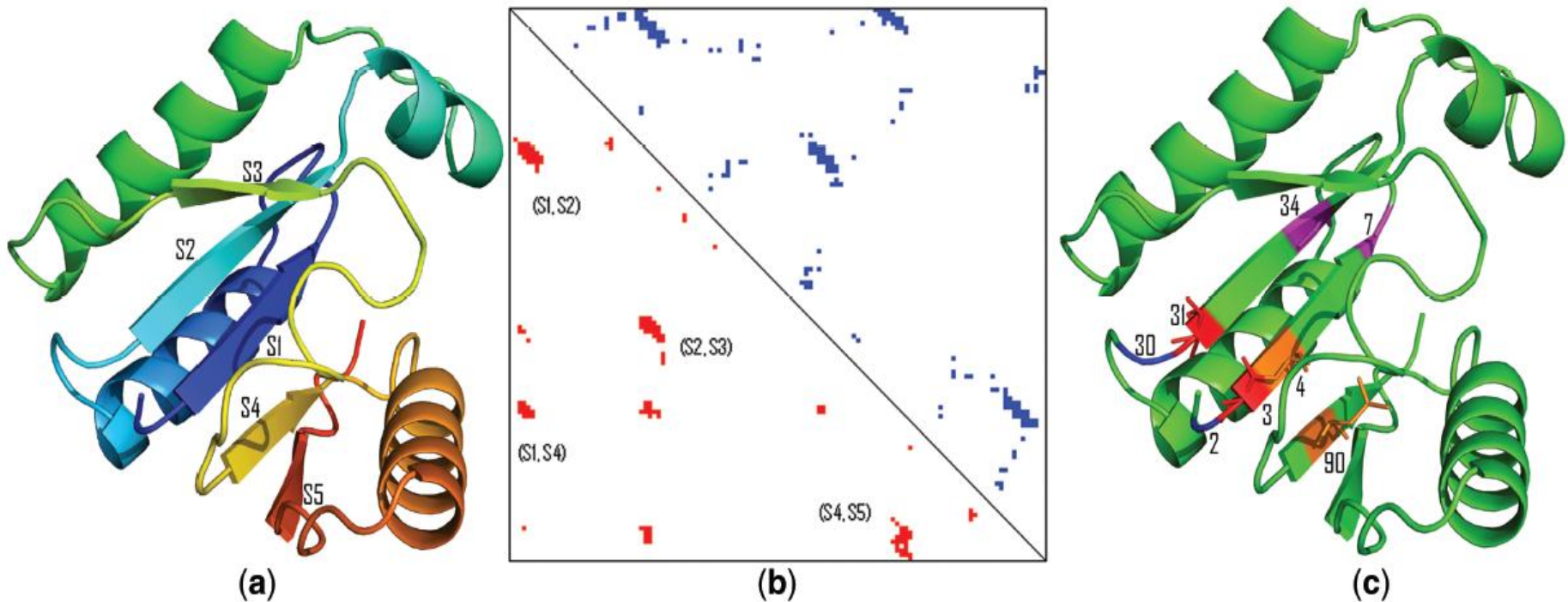
C
A
S
P
9



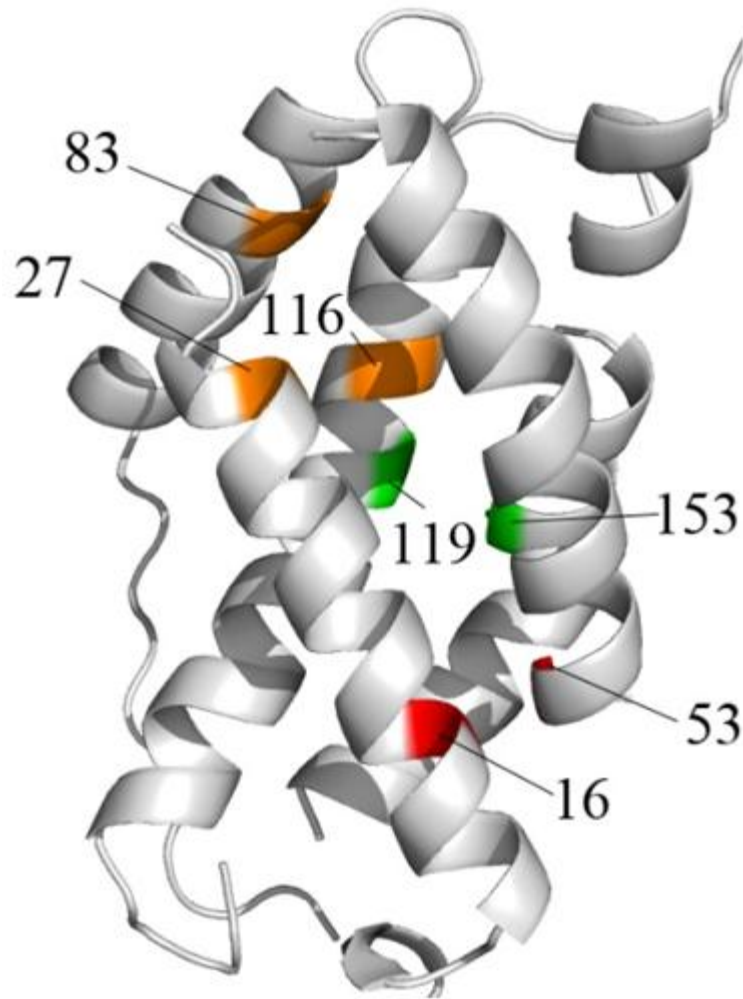
Contact Predictors	≥ 12	≥ 24
Ensemble	.34	.30
Ensemble (± 1)	.55	.48
Zhang-Server	.28	.23
RosettaVote	.27	.20

Top L/5 contacts, on 26 CASP9 *ab initio* domains

A CASP8 Example: T0507



A CASP9 Example

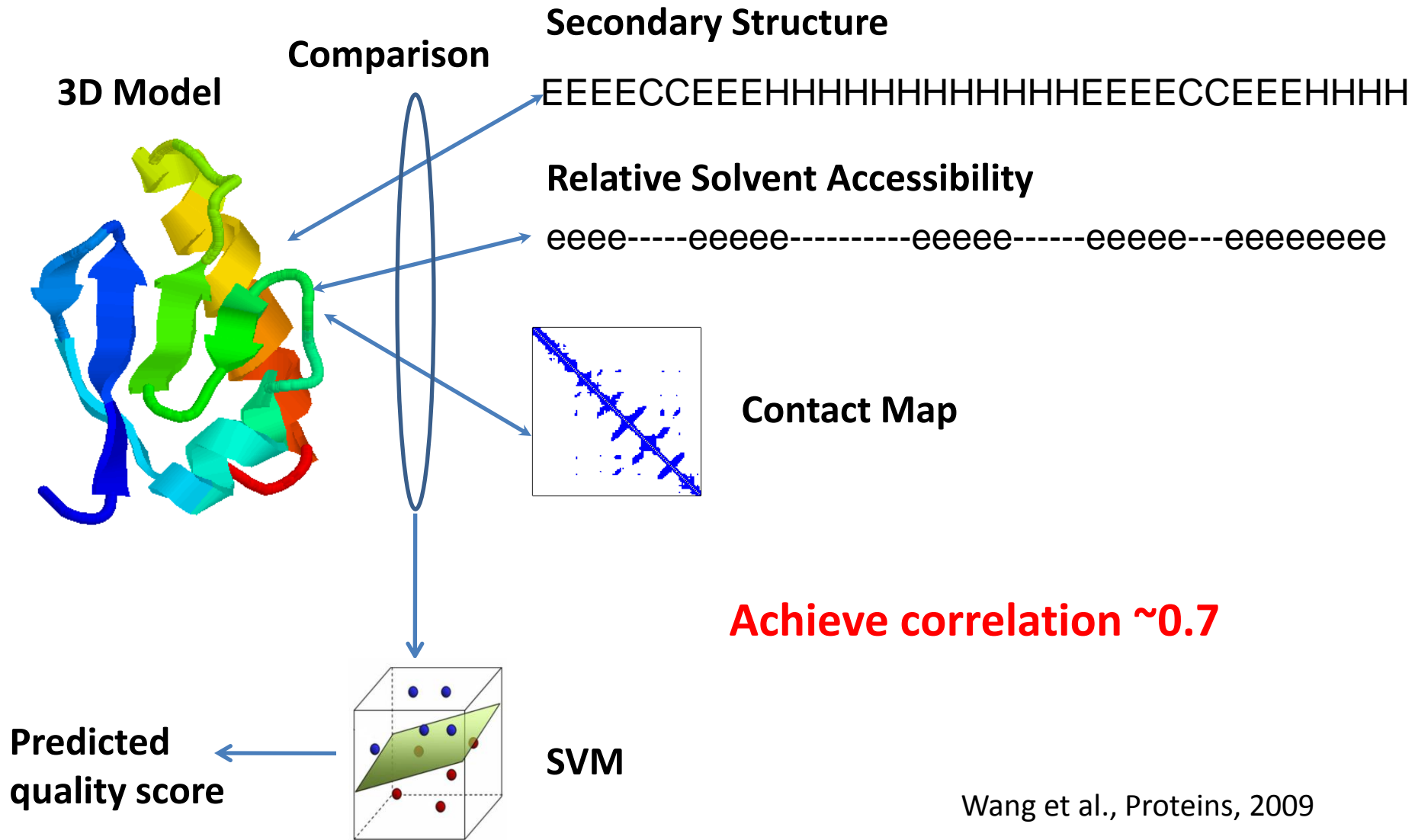


Model Evaluation

- **Single model approach**
- **Pair-wise model assessment**
- **Hybrid approach**

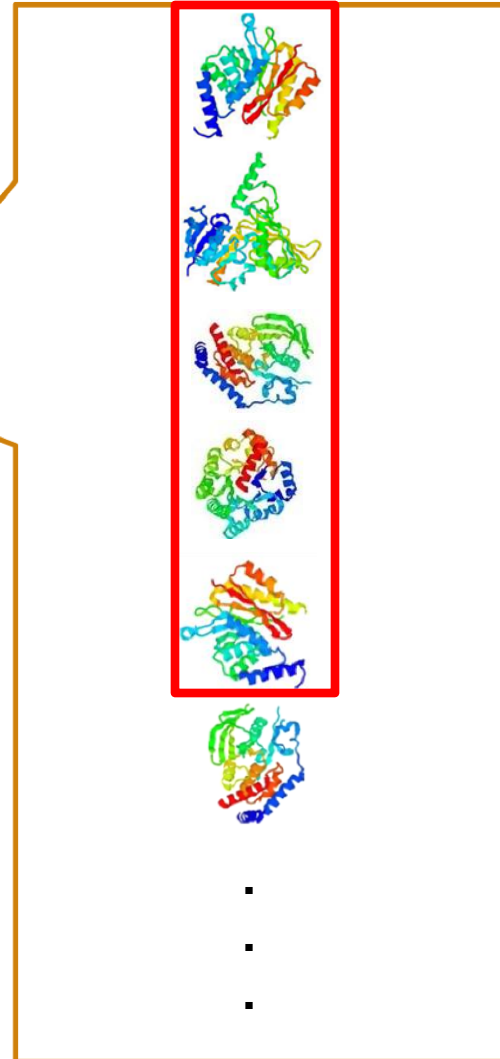
Wang et al., Proteins, 2009
Cheng et al., Proteins, 2009
Wang et al., Bioinformatics, 2011

ModelEvaluator

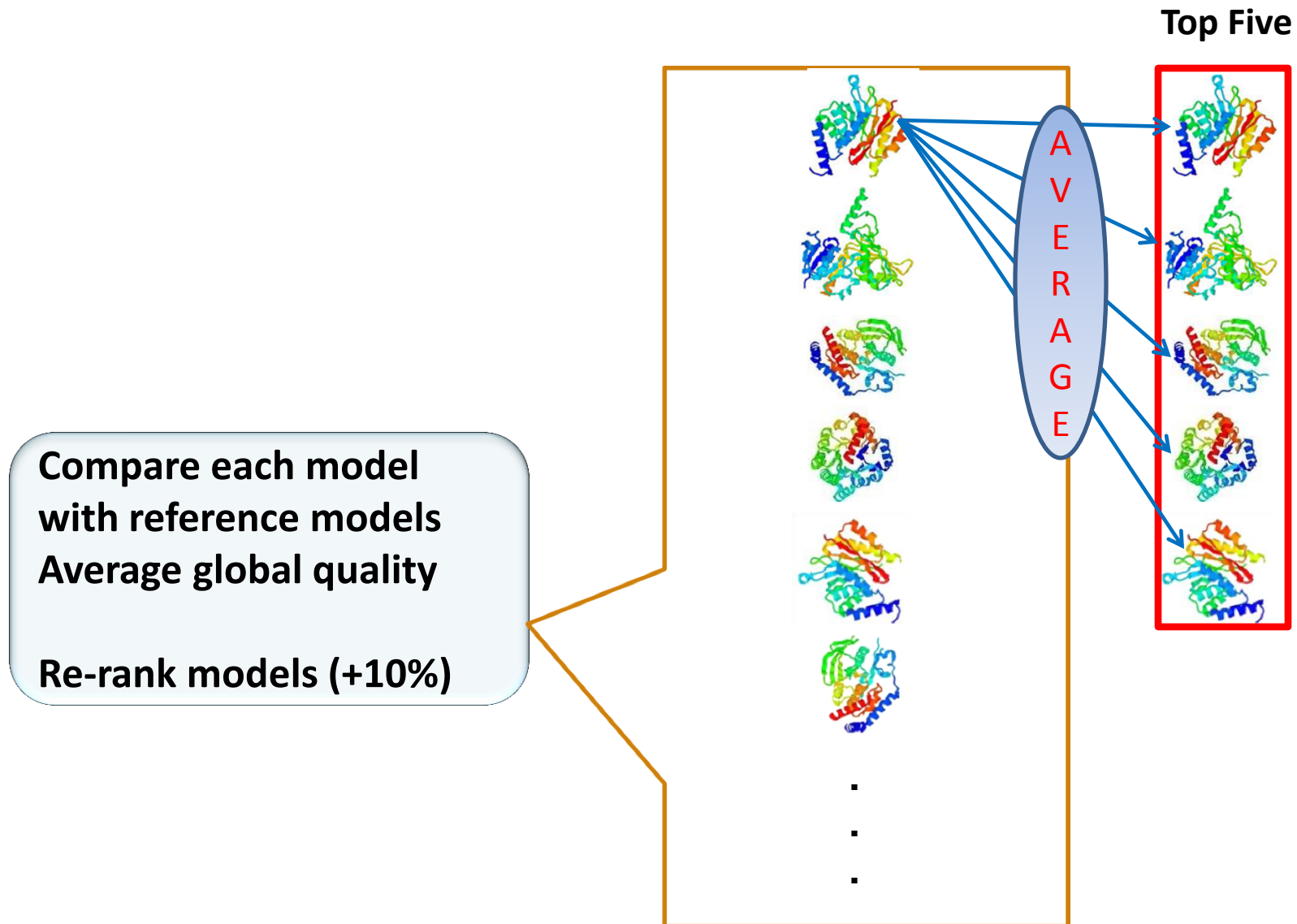


Model Quality Evaluation

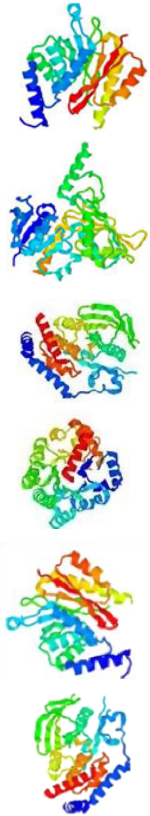
Select top 5 ranked models as references



Model Quality Assessment

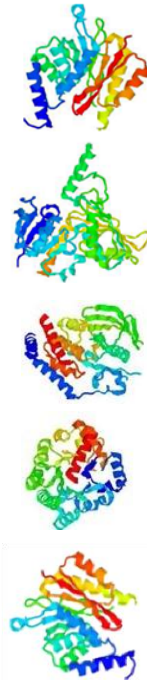


Global-Local Model Combination



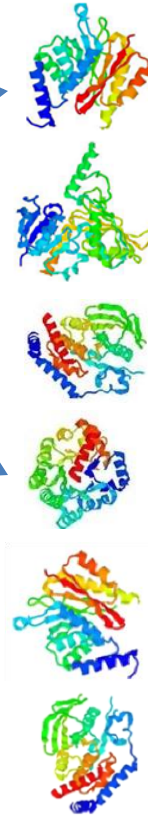
·
·
·

Model ranking



Select top 5 models
as seed models

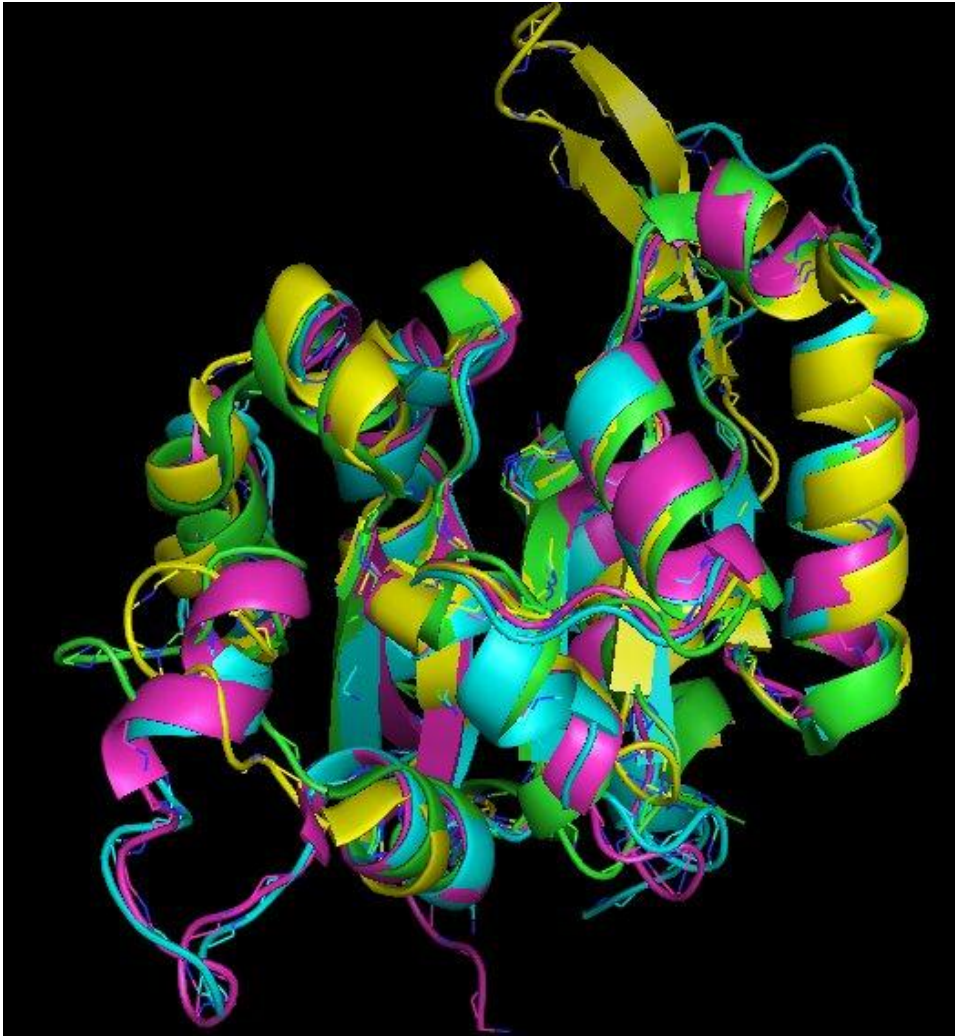
Structure comparison



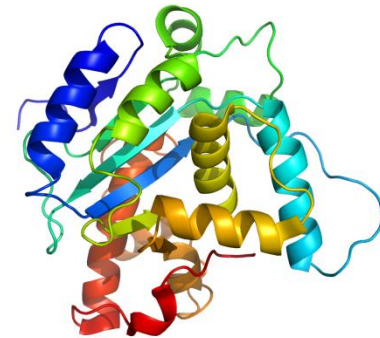
·
·
·

Identify similar models
or fragments

Model Combination and Averaging

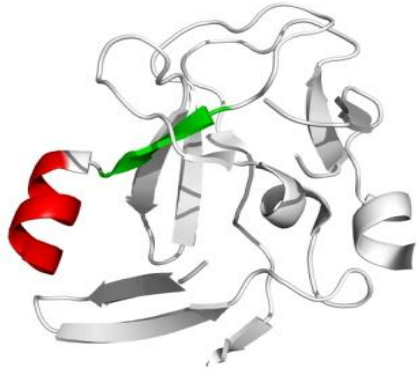


Average

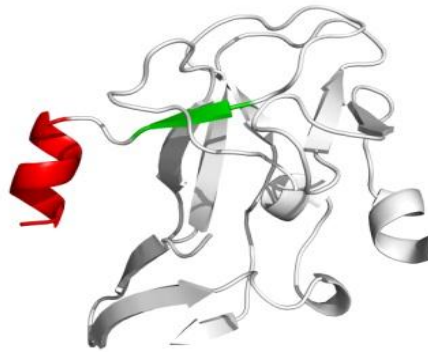


Advantage: reduce variance of modeling while increasing fitness

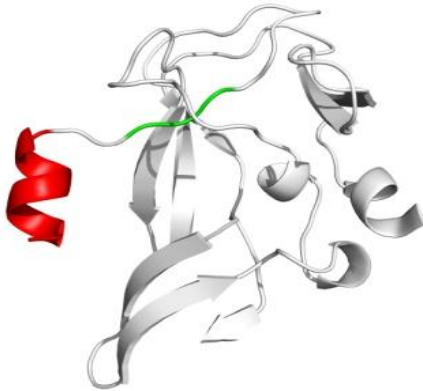
(a) Native



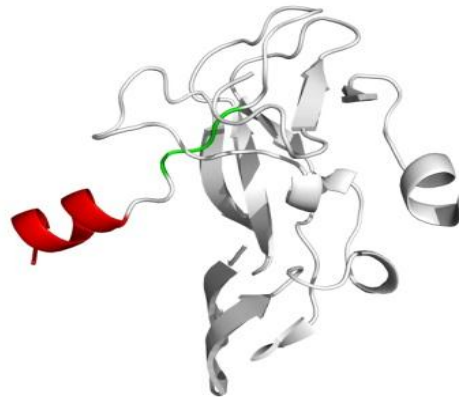
(b) MULTICOM_1 (83.26)



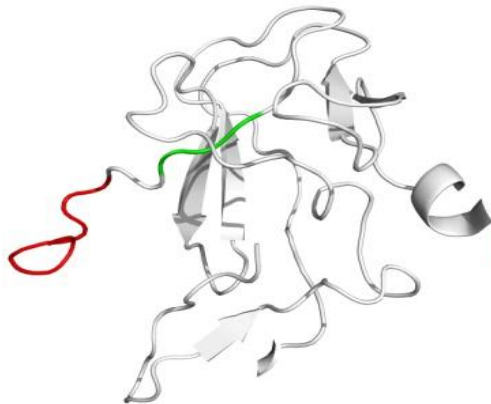
(c) SAM-T08-server_1 (82.20)



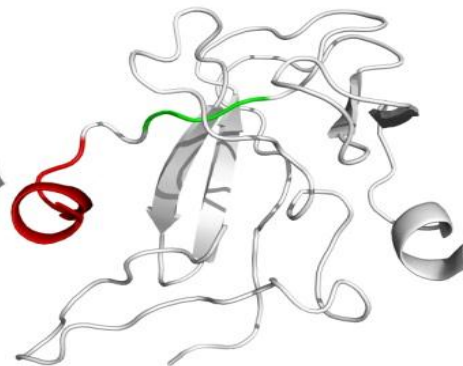
(d) SAM-T06-server_1 (79.03)



(e) pro-sp3-TASSER_5 (80.08)



(f) pro-sp3-TASSER_4 (80.30)



C
A
S
P
8



Wang et al., Bioinformatics, 2010

CASP9 Top 20 Servers

#	GR #	GR name	Domains Count	SUM Z-score (GDT_TS)	AVG Z-score (GDT_TS)	AVG GDT_TS
1.	380	QUARK	147	115.788	0.788	62.675
2.	428	Zhang-Server	147	113.242	0.770	62.765
3.	077	RaptorX-MSA	147	103.270	0.703	61.774
4.	286	RaptorX	147	103.010	0.701	61.731
5.	276	RaptorX-Boost	147	99.845	0.679	61.453
6.	449	HHpredB	147	93.104	0.633	59.528
7.	453	HHpredA	147	93.104	0.633	59.528
8.	346	HHpredC	147	91.821	0.625	59.361
9.	452	Seok-server	147	89.542	0.609	60.158
10.	002	MULTICOM-CLUSTER	147	88.944	0.605	59.987
11.	321	BAKER-ROSETTASERVER	145	87.240	0.602	58.768
12.	119	MULTICOM-REFINE	147	86.441	0.588	59.519
13.	215	MULTICOM-NOVEL	147	82.825	0.563	59.371
14.	236	gws	145	82.645	0.570	58.931
15.	457	chunk-TASSER	147	82.609	0.562	58.846
16.	174	Phyre2	147	78.792	0.536	58.823
17.	080	MULTICOM-CONSTRUCT	147	76.446	0.520	58.703
18.	253	pro-sp3-TASSER	147	75.358	0.513	58.117
19.	481	MUFOLD-Server	147	68.676	0.467	56.260
20.	127	FAMSD	147	68.669	0.467	57.295

CASP9 Top 20 Servers on AB Initio Targets

#	GR #	GR name	Domains Count	SUM Z-score (GDT_TS)	AVG Z-score (GDT_TS)	AVG GDT_TS
1.	380	QUARK	29	31.632	1.091	31.851
2.	428	Zhang-Server	29	26.506	0.914	30.539
3.	119	MULTICOM-REFINE	29	22.423	0.773	28.971
4.	457	chunk-TASSER	29	20.697	0.714	28.626
5.	002	MULTICOM-CLUSTER	29	19.954	0.688	28.445
6.	286	RaptorX	29	19.754	0.681	27.730
7.	077	RaptorX-MSA	29	19.343	0.667	27.494
8.	321	BAKER-ROSETTASERVER	29	18.973	0.654	27.118
9.	253	pro-sp3-TASSER	29	18.892	0.651	27.996
10.	276	RaptorX-Boost	29	18.286	0.631	27.400
11.	215	MULTICOM-NOVEL	29	18.266	0.630	27.813
12.	055	MUFOLD-MD	28	16.900	0.604	24.986
13.	080	MULTICOM-CONSTRUCT	29	16.582	0.572	27.219
14.	063	Jiang_Assembly	29	14.721	0.508	26.249
15.	236	gws	29	13.931	0.480	26.002
16.	047	BioSerf	29	13.590	0.469	24.716
17.	103	SAM-T08-server	29	12.686	0.437	24.618
18.	452	Seok-server	29	12.552	0.433	24.762
19.	481	MUFOLD-Server	29	10.590	0.365	23.567
20.	174	Phyre2	29	10.385	0.358	24.465

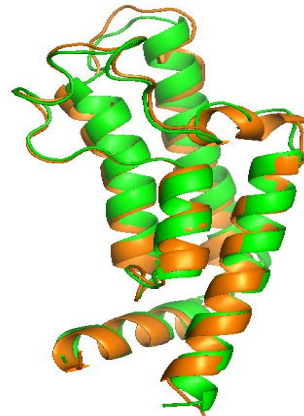
Some High-Quality CASP Predictions



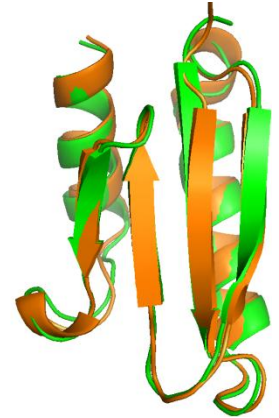
T0390
GDT=0.90



T0426
GDT=0.97



T0432
GDT=0.92



T0458
GDT=0.97

Orange: structure; Green: model

50 of 120 CASP8 targets are in high-accuracy, $\text{RMSD} < 2 \text{ \AA}$

Genome-Wide Annotation of Protein Structure



Schmutz et al., Nature, 2010

DATABASE

Open Access

SoyDB: a knowledge database of soybean transcription factors

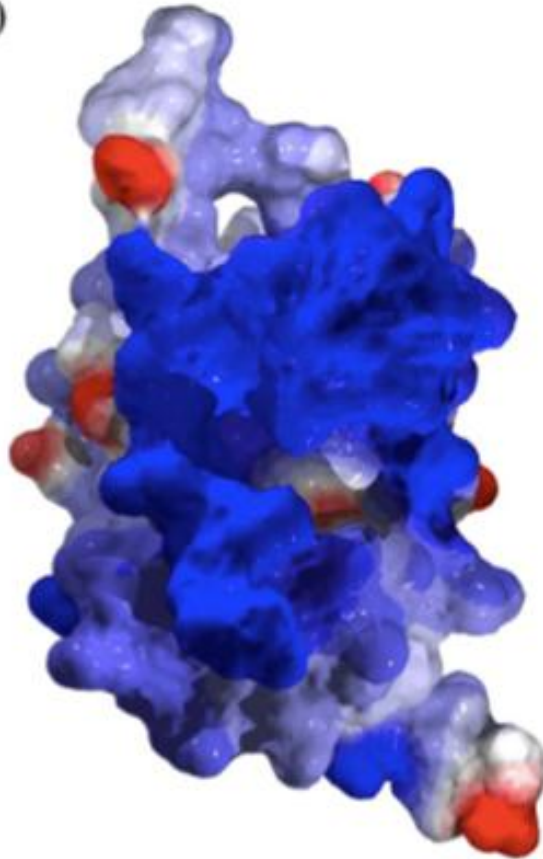
BMC Plant Biology

Zheng Wang¹, Marc Libault^{2,3}, Trupti Joshi^{1,2}, Babu Valliyodan^{2,3}, Henry T Nguyen^{2,3}, Dong Xu^{1,2,4}, Gary Stacey^{2,3}, Jianlin Cheng^{1,2,4*}

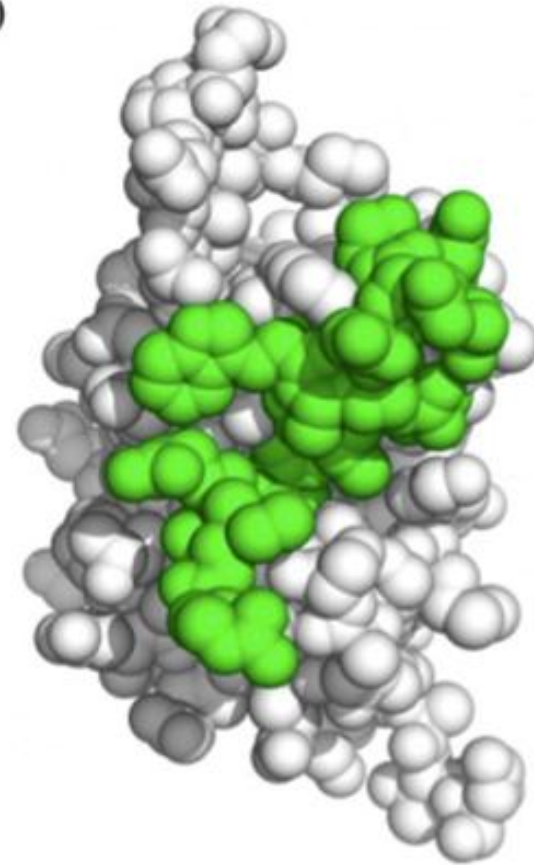
Highly accessed

Transcription Factor Glyma01g32810

(a)



(b)



Acknowledgements

Group Members

- Xin Deng
- Jesse Eickholt
- Jianfeng He
- Jilong Li
- Cuong Nguyen
- Allison Tegge
- Zheng Wang

