

Zhang Groups: Automated Structure Prediction by C-I-TASSER and C-QUARK in CASP13

Robin Pearce, Wei Zheng, Chengxin Zhang, Yang Li, S M Mortuza, Yang Zhang

Department of Computational Medicine and Bioinformatics
Department of Biological Chemistry
University of Michigan

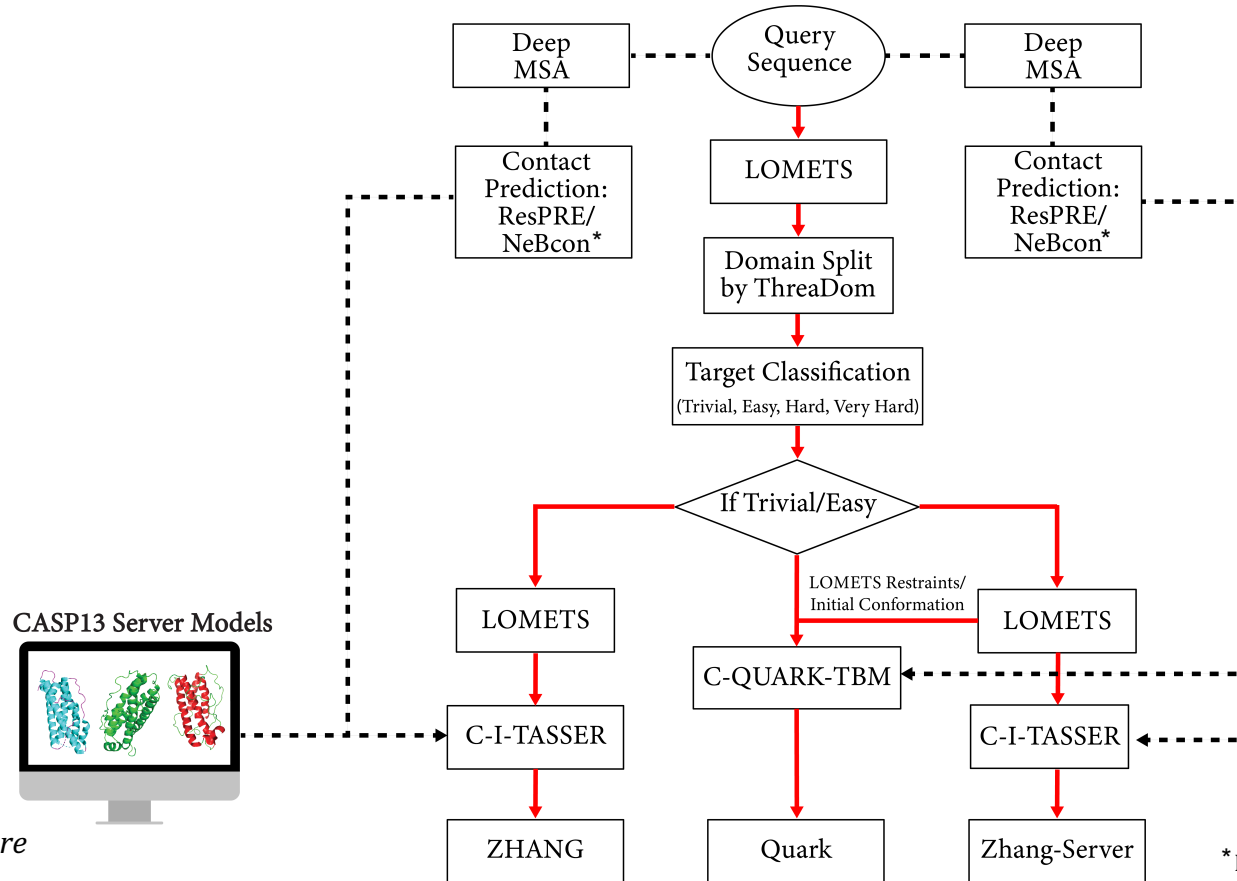


Outline

- **Methods**
 - CASP13 Workflow (Trivial, Easy, Hard, Very Hard Targets)
 - Contact Prediction Method and Potential
 - C-Quark and Zhang-Server Pipelines
- **Results**
 - Summary of FM and TBM Modeling Results
 - Impact of Contact Prediction and Template Quality on Modeling Results
 - Case Studies
 - Problems that We Encountered
- **Summary**

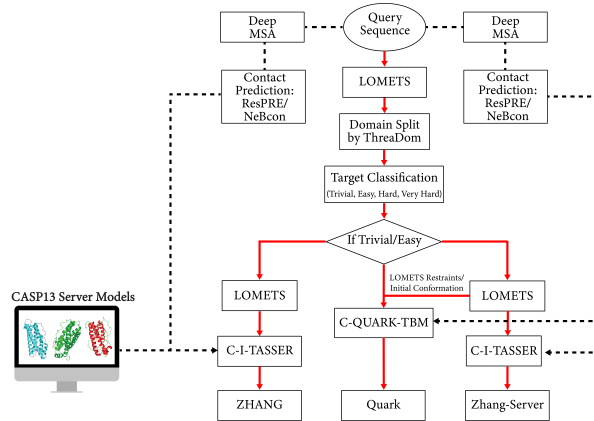
Methods

Pipeline for CASP13 Trivial/Easy Targets

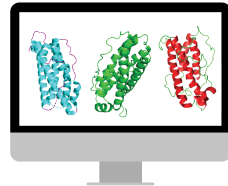


Pipeline for CASP13 Hard/Very Hard Targets

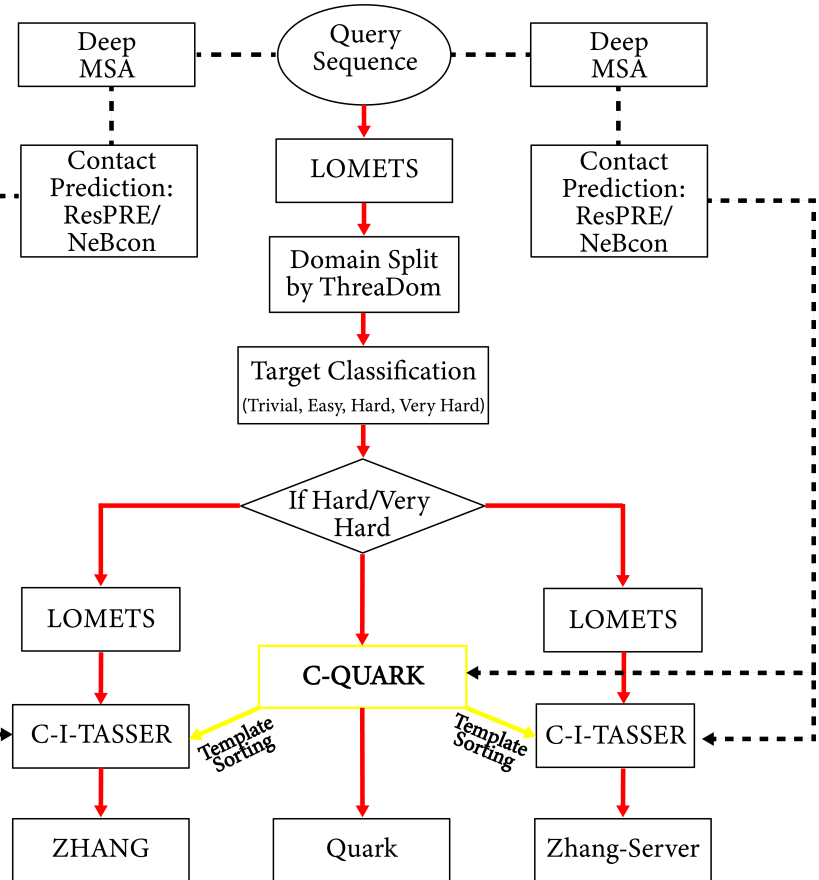
Pipeline for CASP13 Trivial/Easy Targets



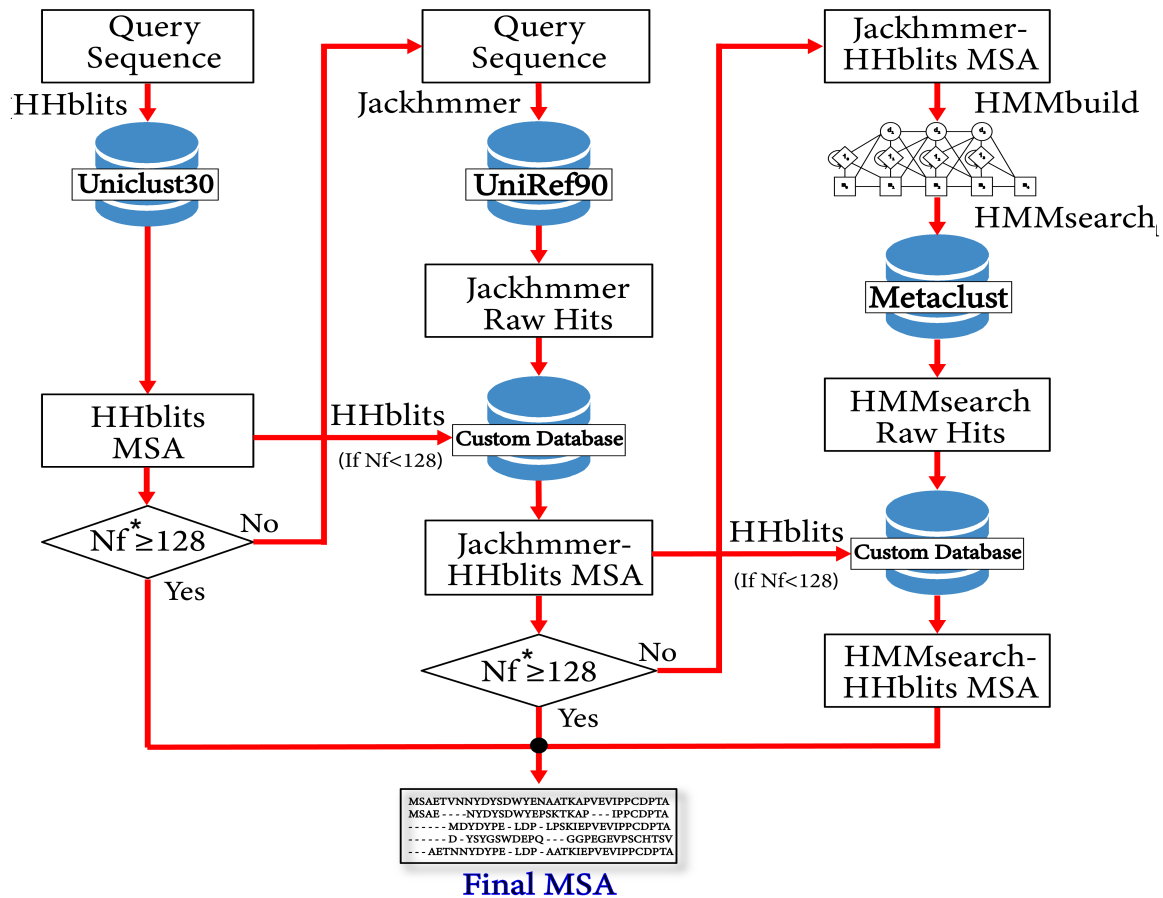
CASP13 Server Models



All Predictions are Automated

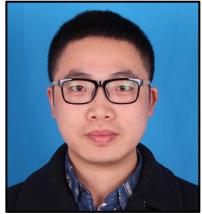


Deep MSA Construction



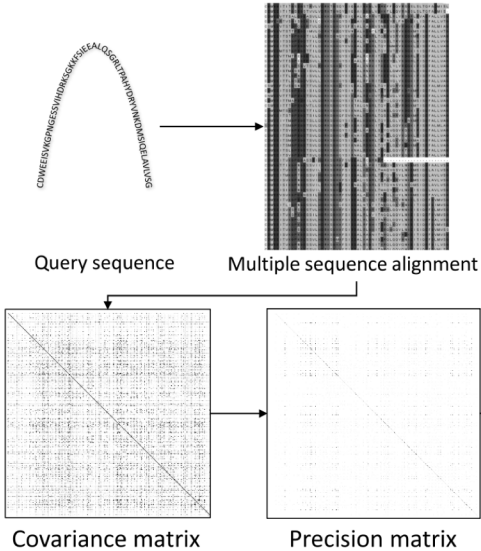
* N_f : Number of Effective Sequences in MSA

Contact Prediction Using ResPRE

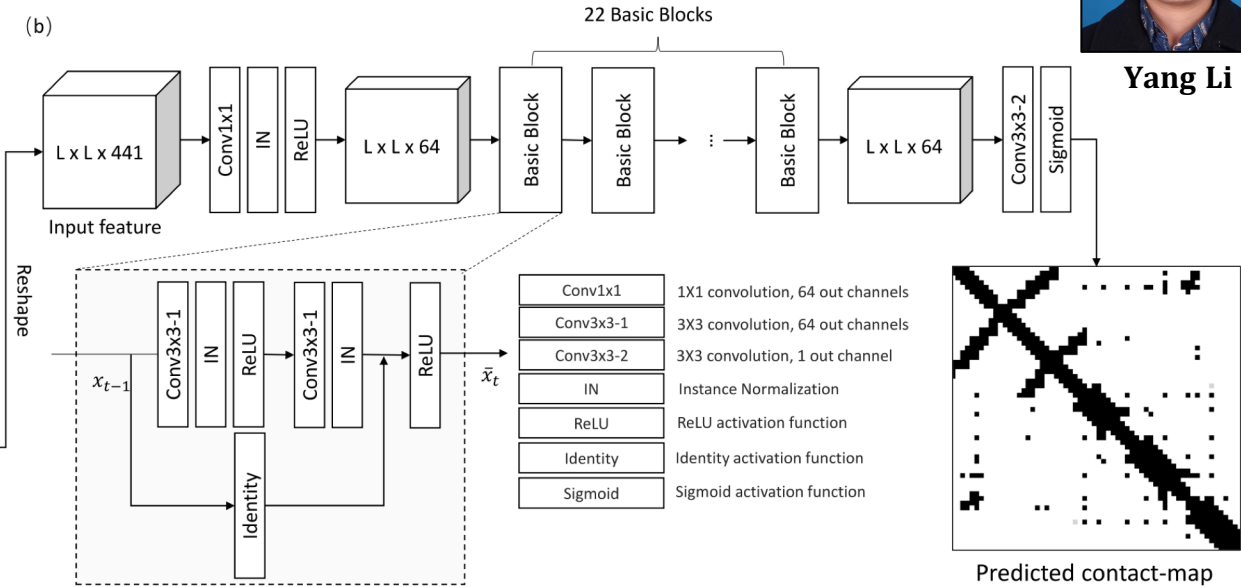


Yang Li

(a)



(b)



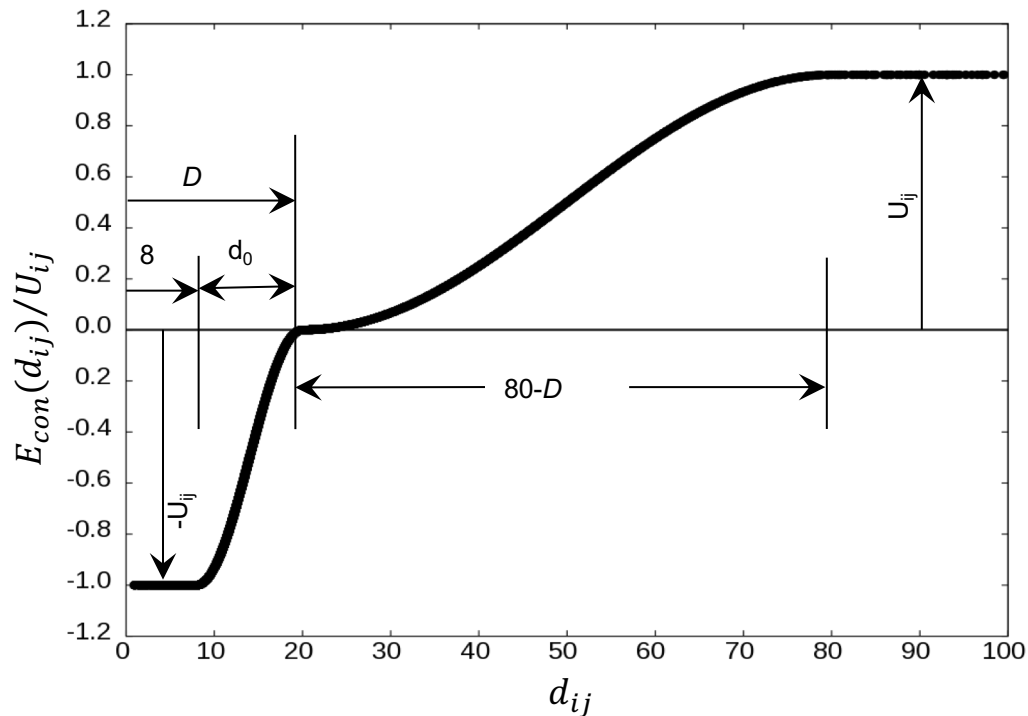
Covariance Matrix (S):

$$S_{i,j}(a,b) = f_{i,j}(a,b) \quad f_i(a) \cdot f_j(b)$$

Precision Matrix (θ):

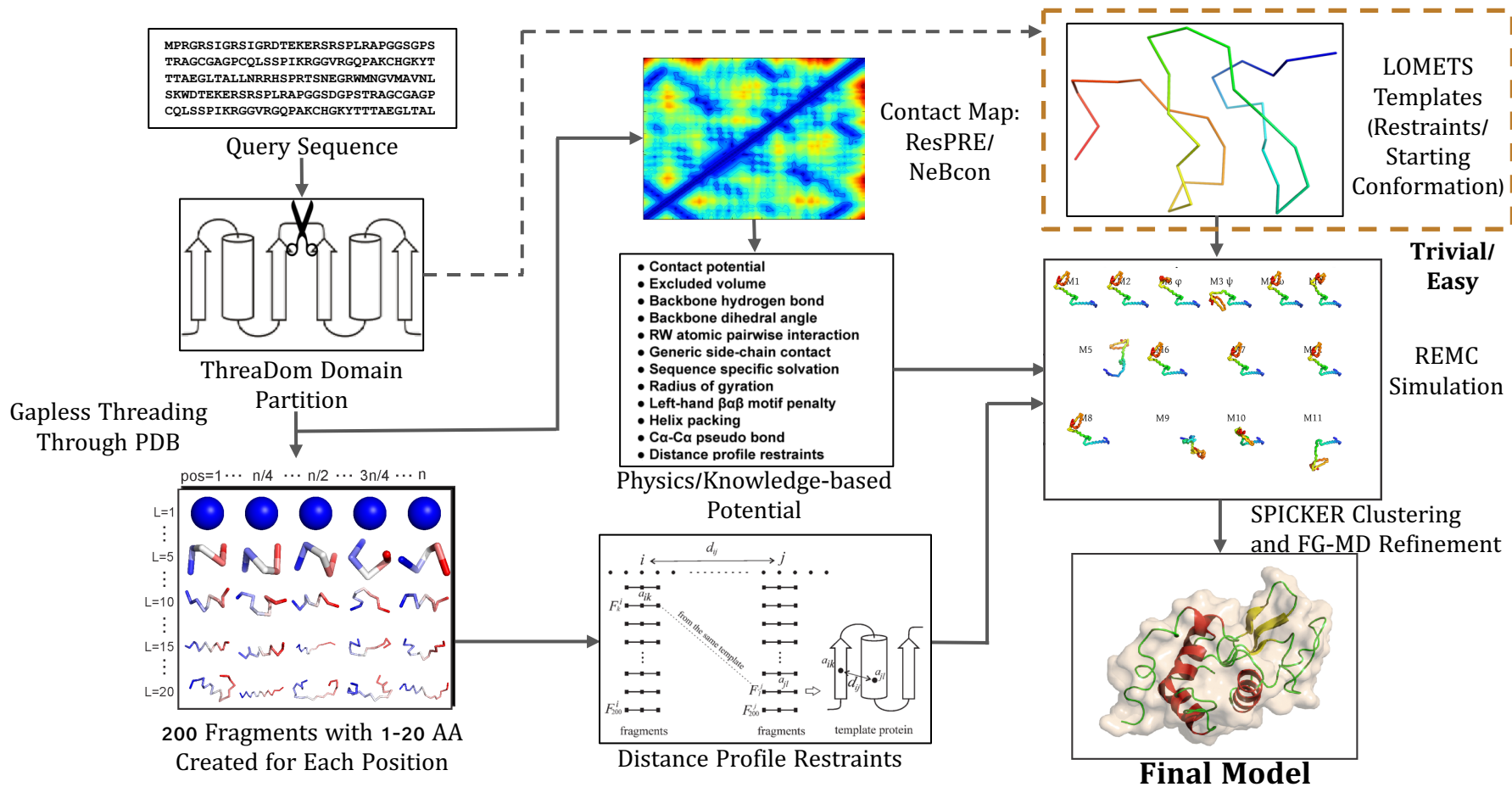
$$\theta = \underset{\theta}{\text{argmin}} (\text{tr}(S \cdot \theta) \quad \log(\det(\theta)) + \rho \cdot \|\theta\|_2^2)$$

Three Gradient Contact Potential in C-I-TASSER/C-QUARK

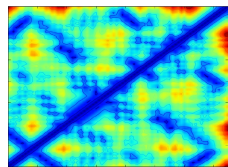


$$E_{con}(d_{ij}) = \begin{cases} -U_{ij}, & d_{ij} < 8\text{\AA} \\ -\frac{1}{2}U_{ij} \left[1 - \sin\left(\frac{d_{ij} - \left(\frac{8+D}{2}\right)\pi}{d_b}\right) \right], & 8\text{\AA} \leq d_{ij} < D \\ \frac{1}{2}U_{ij} \left[1 + \sin\left(\frac{d_{ij} - \left(\frac{D+80}{2}\right)\pi}{(80-D)}\right) \right], & D \leq d_{ij} \leq 80\text{\AA} \\ U_{ij}, & d_{ij} > 80\text{\AA} \end{cases}$$

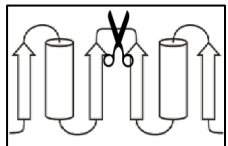
C-QUARK Pipeline with LOMETS



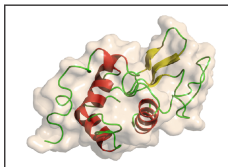
Zhang-Server Pipeline Built on C-I-TASSER



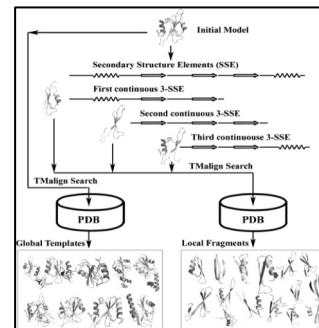
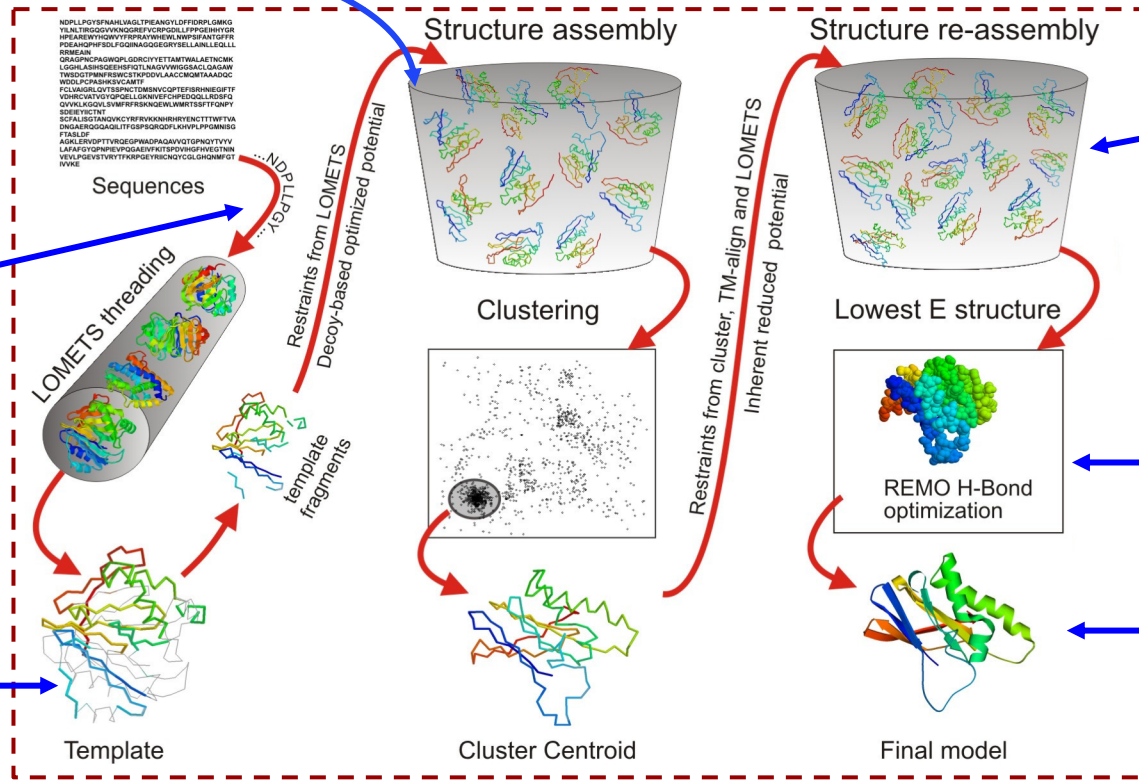
Contact Map:
ResPRE/NeBcon



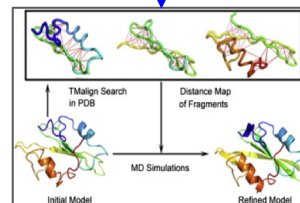
Domain Partition
by ThreaDom



Template Sorting
by C-QUARK
(Hard/Very Hard)



Fragments by TM-align



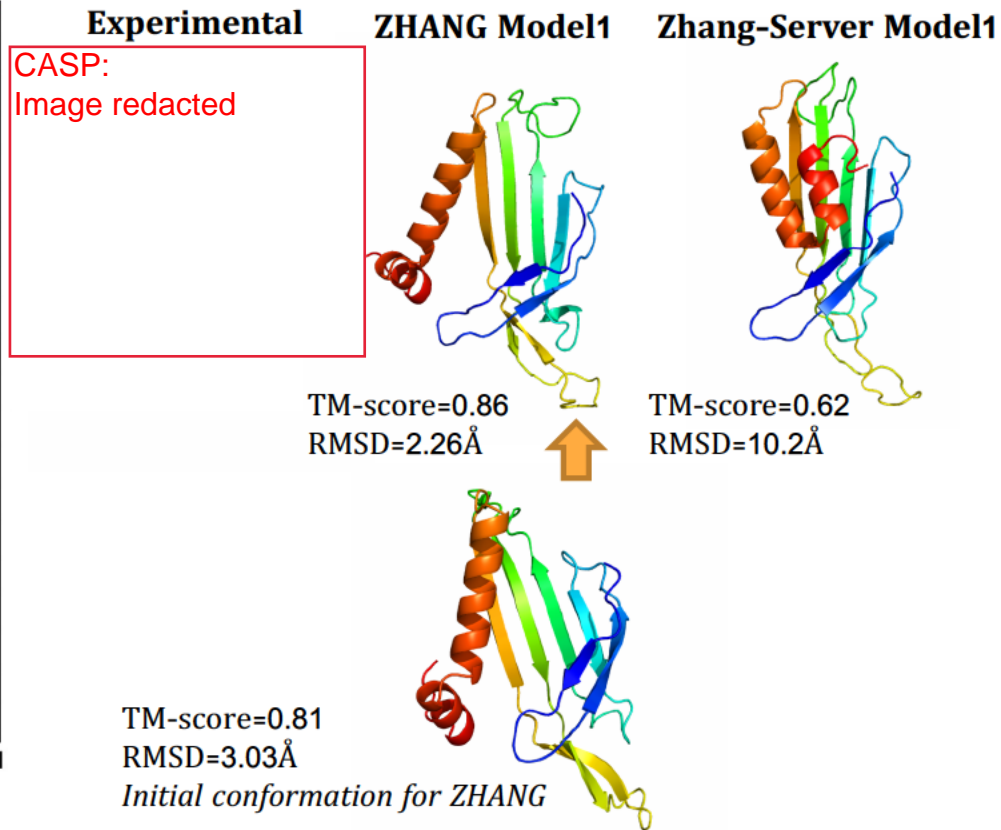
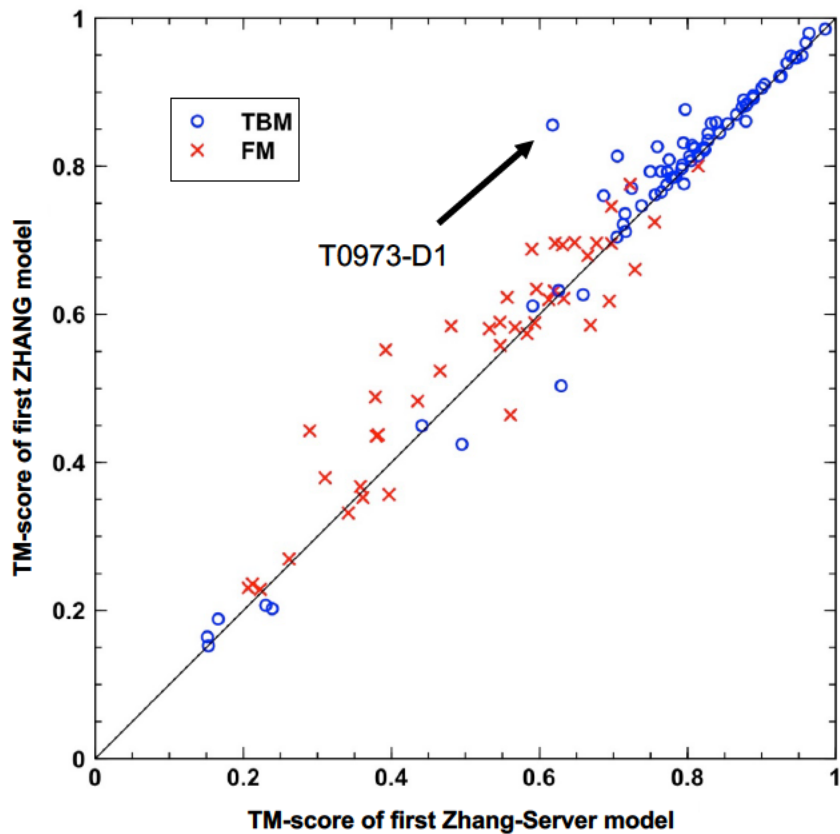
Atomic Refinement
by FG-MD



Residue Quality
Estimation by ResQ

Results

ZHANG vs. Zhang-Server

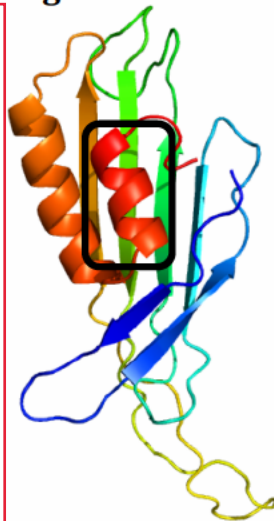


ZHANG vs. Zhang-Server

Experimental

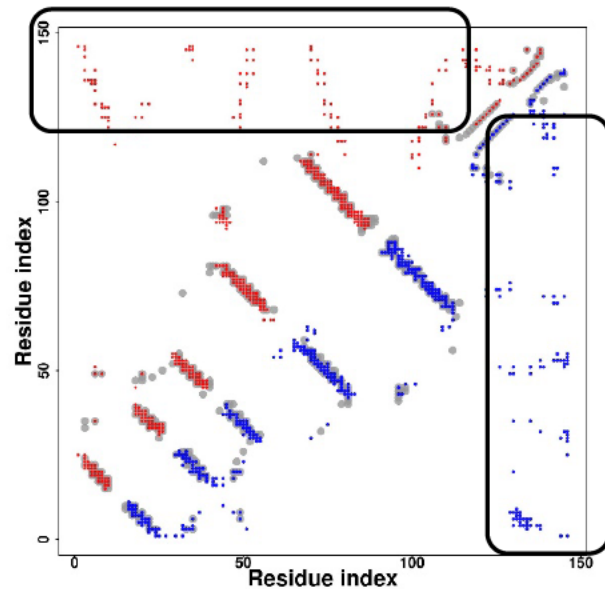
Zhang-Server Model1

CASP:
Images redacted



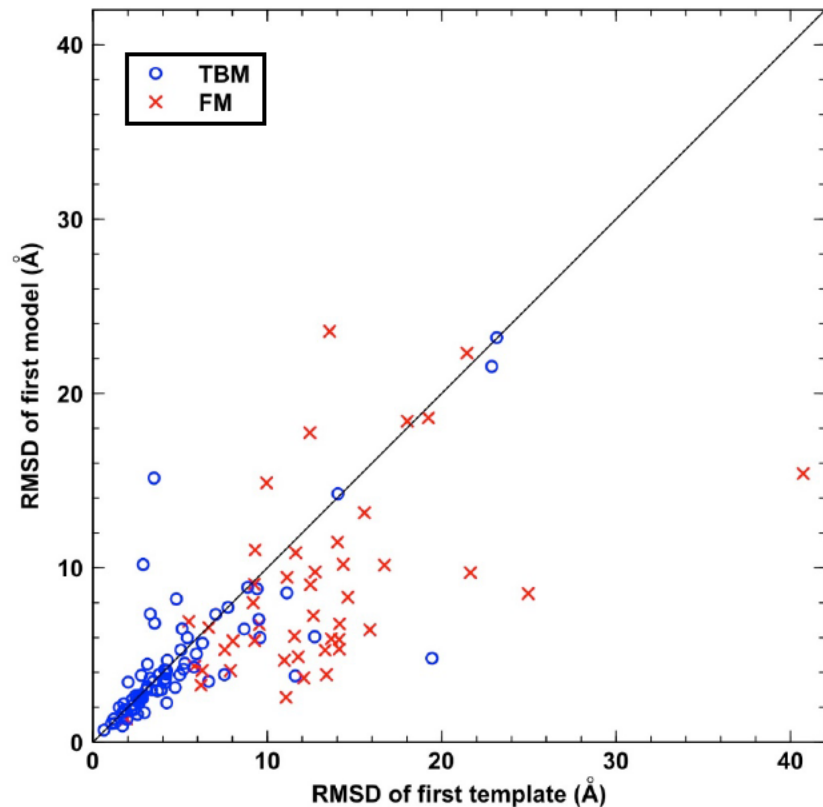
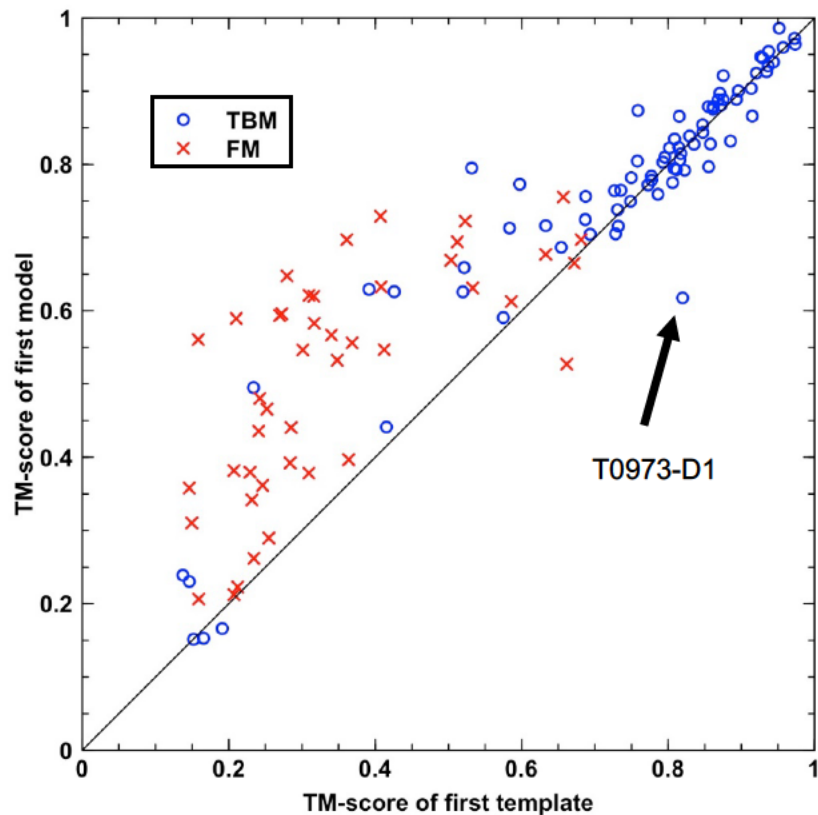
TM-score=0.62
RMSD=10.2Å

Experimental (rainbow)
Template (black)
TM-score=0.82
RMSD=2.88Å

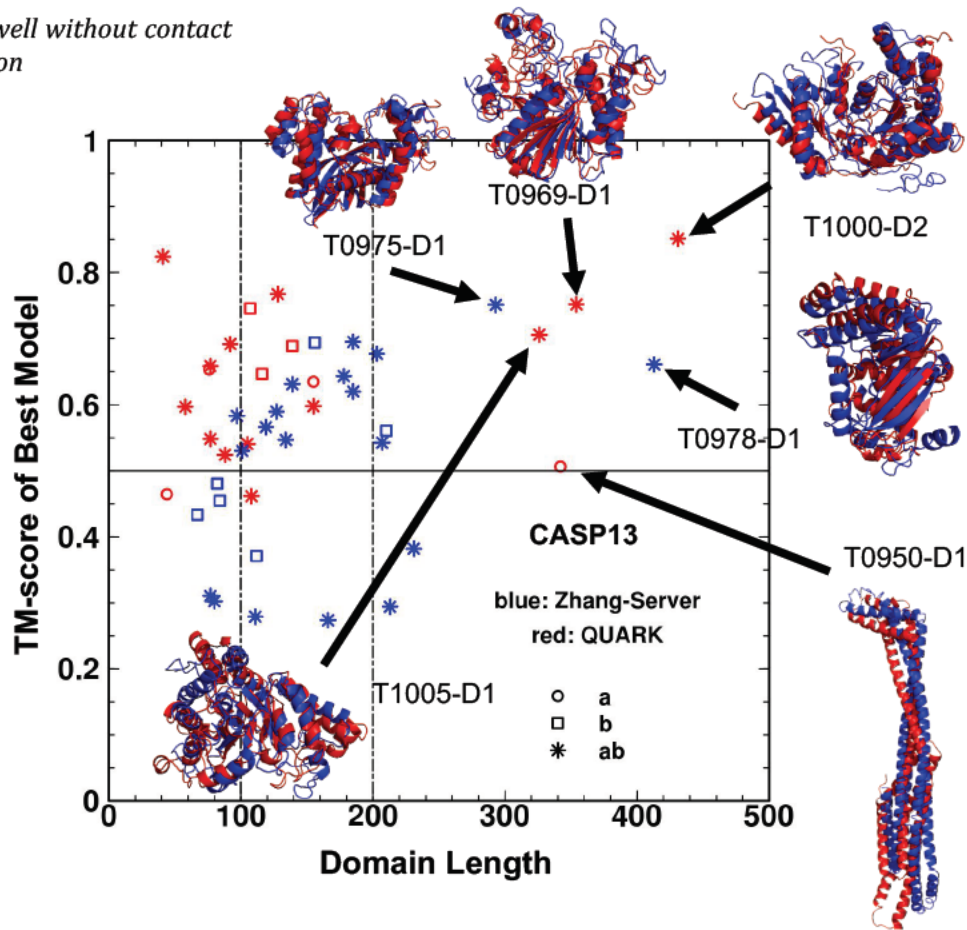
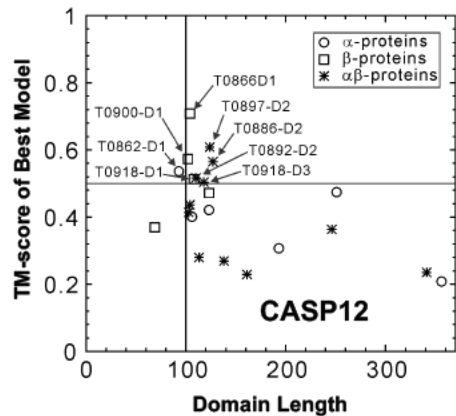
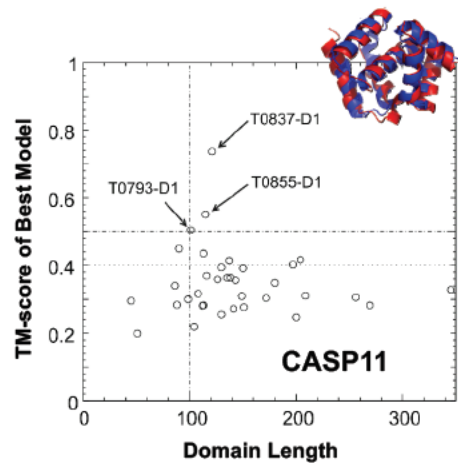


Good template, bad long range contacts!

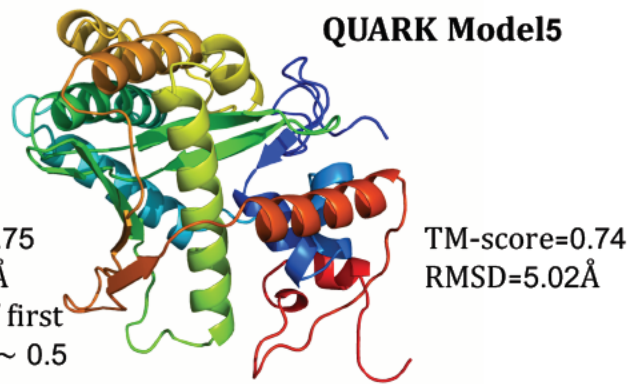
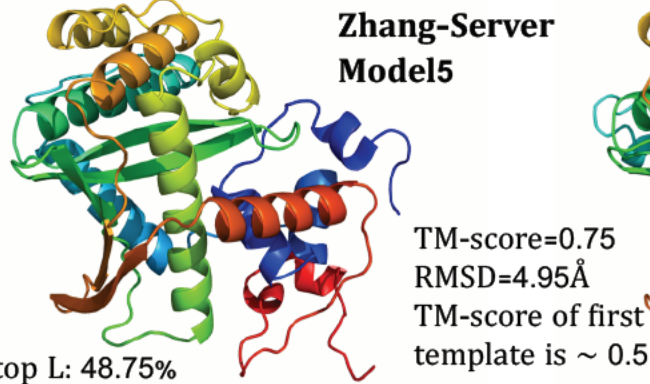
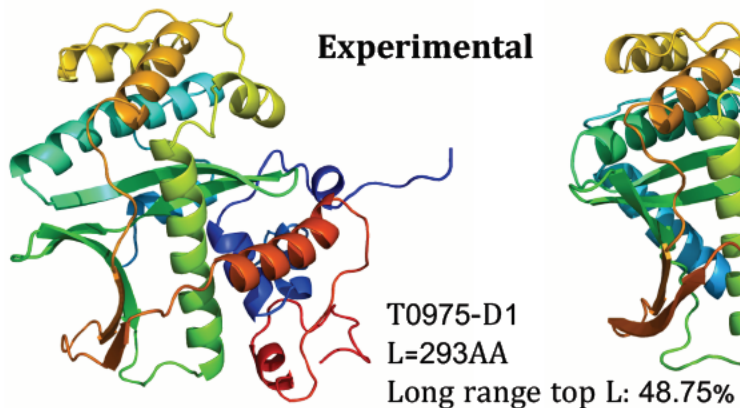
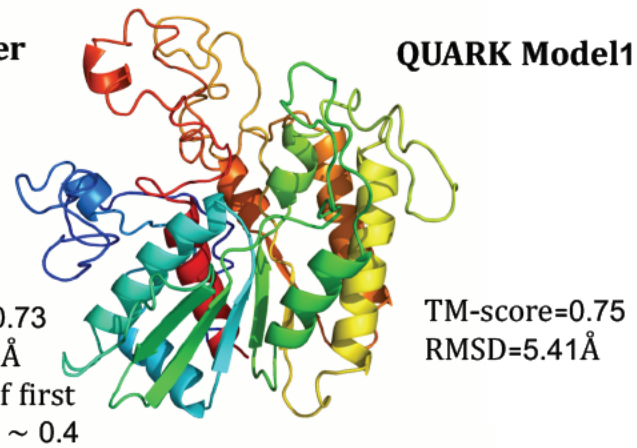
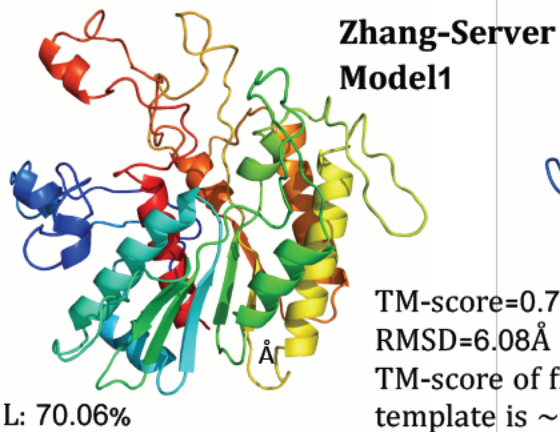
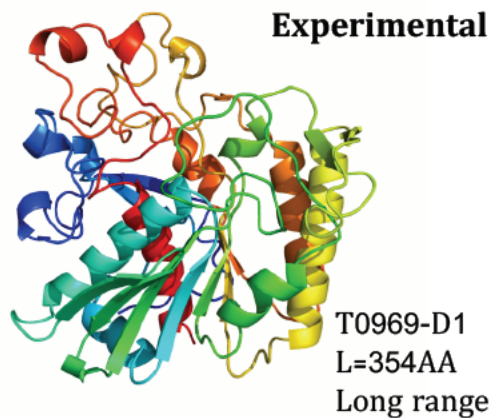
First Zhang-Server Model vs. First LOMETS Template



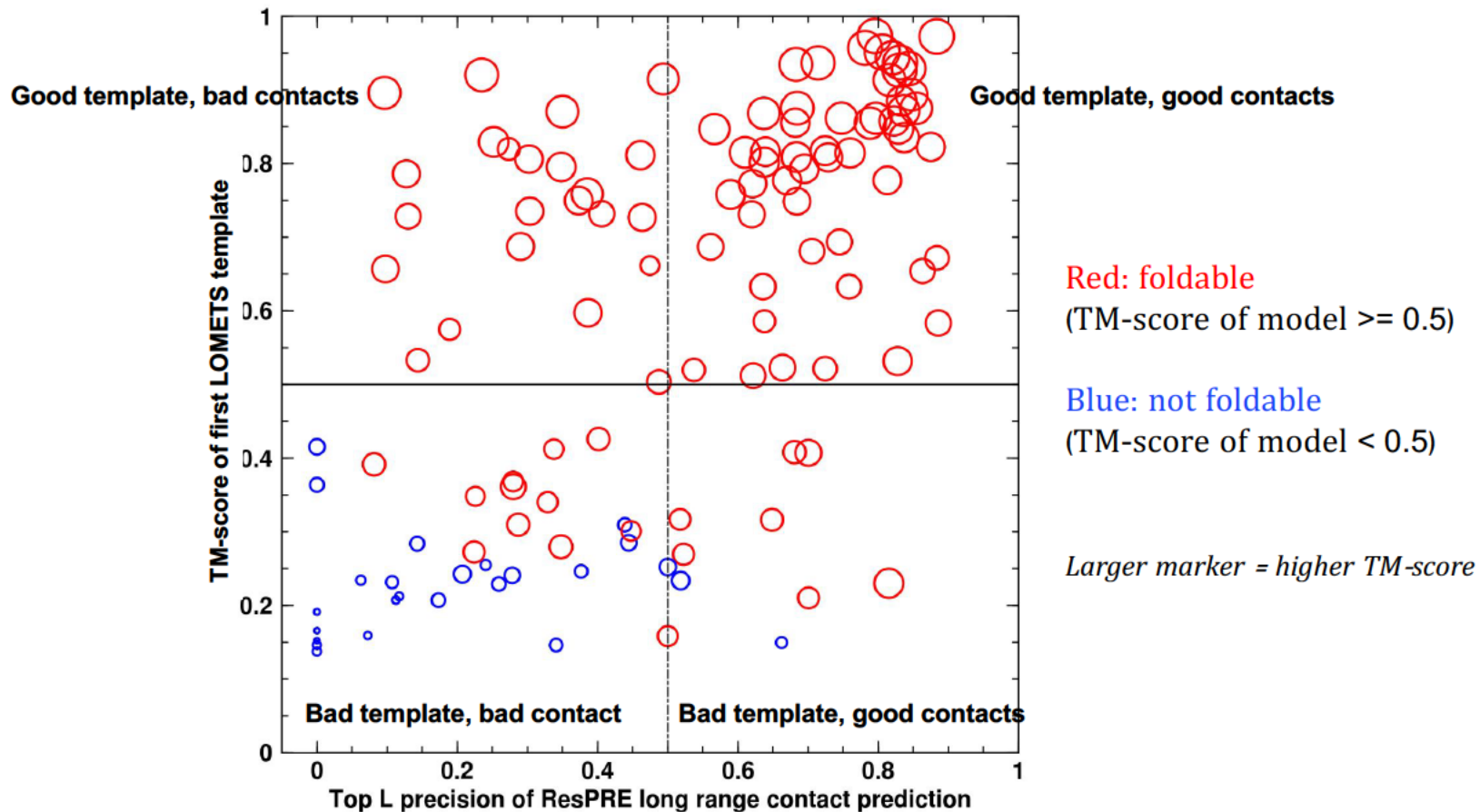
Summary of 45 FM targets (32 FM targets, 13 FM/TBM targets)



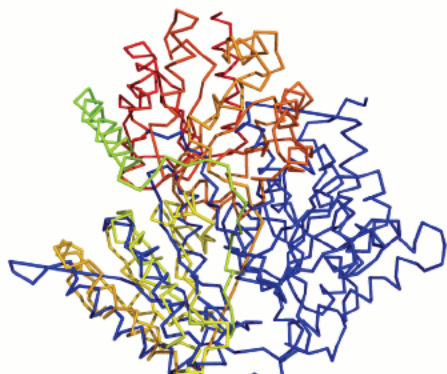
Folding Large Targets



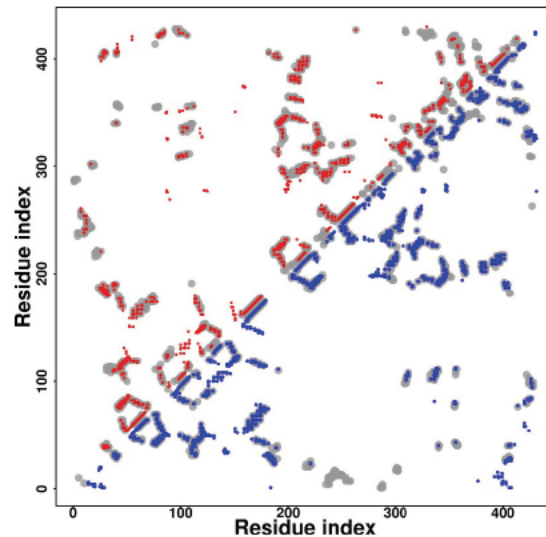
Impact of Contact and Template Quality on Structure Prediction



T1000-D2 (FM Target) (Bad Template, Good Contacts)

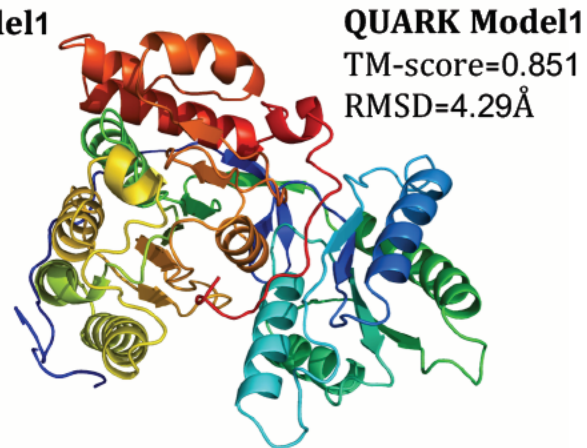
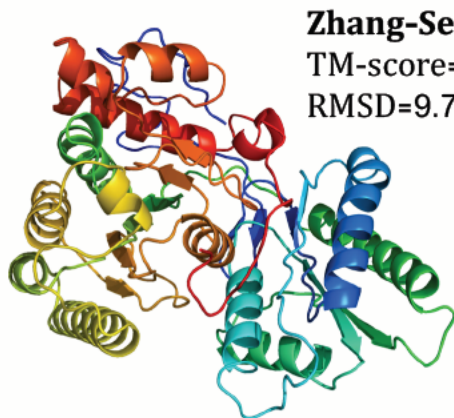
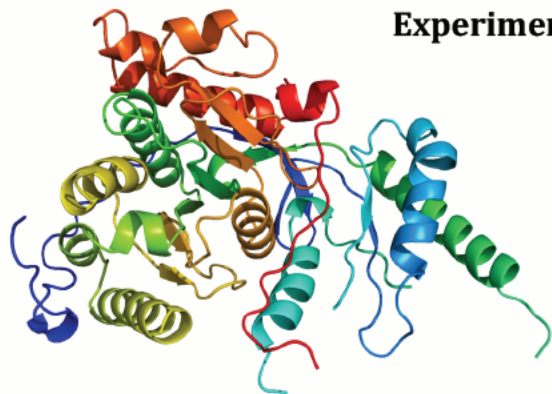


Best template: 1jqkA
TM-score=0.21
RMSD=21.65Å
L=431AA



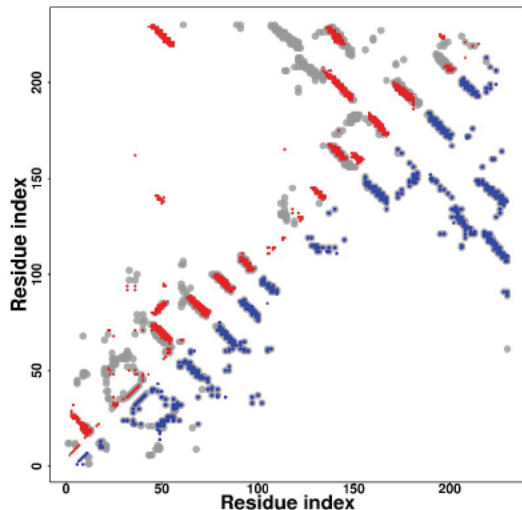
Grey: Native contacts
Red: ResPRE contacts
Blue: Zhang-Server
model1 contacts

ResPRE Long range
top L: 81.52%

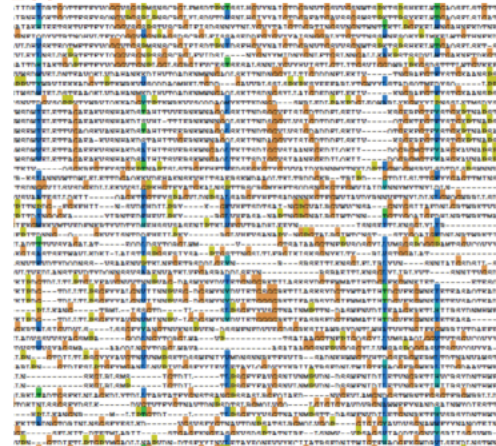


T1004-D3 (Good Template, Bad Contacts)

Sparse MSA for contact prediction but good structure template exists!

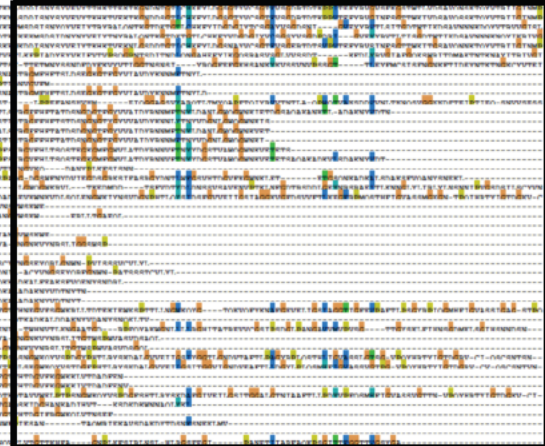


260



MSA

458



CASP:
Image redacted

T1004-D3: 229-458

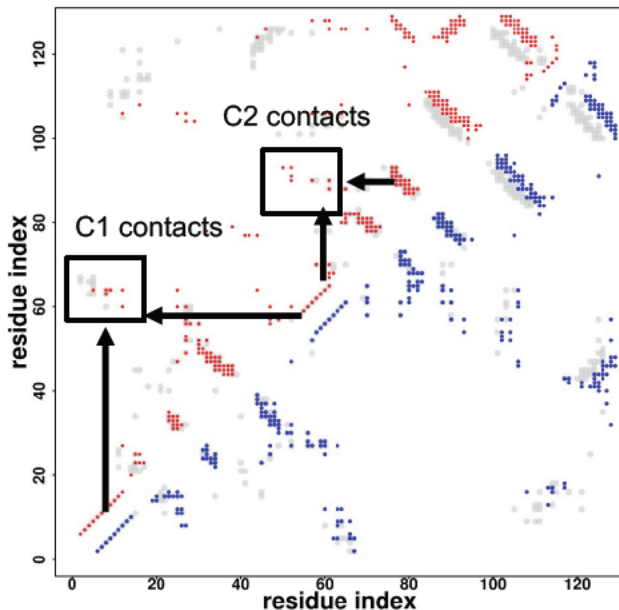
Template: 5m9fA, TM-score=0.92, RMSD=2.02Å

Long range, top L=23.48%

Zhang-Server Model1 TM-score=0.925

T1017s2 (Bad Template, Bad Contacts, Good Fold)

CASP:
Image redacted



CASP:
Image redacted

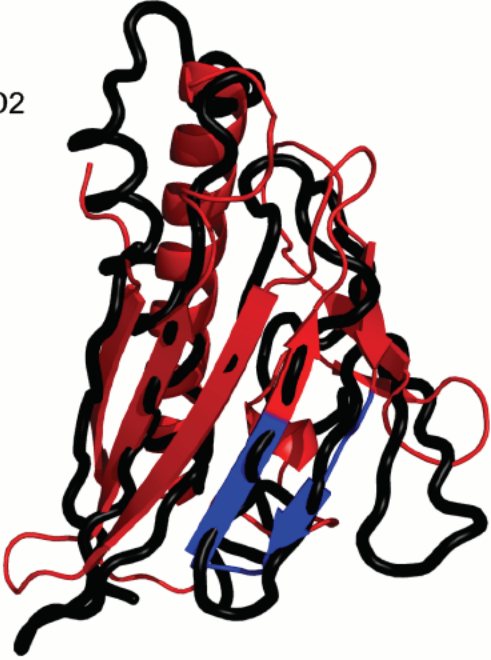
Template: 5m9fA (yellow)
TM-score=0.36
RMSD=6.3Å
ResPRE Long range, top L=28%

Balance template and contacts!

Zhang-Server Model1
TM-score=0.70
RMSD=4.1Å
ResPRE Long range, top L=28%

Secondary Structure Prediction Problem

T0982-D2



```
Conf: 885279857997404452447988989862224757960698778874888877558971
Pred: CCCEEEEECCCCCEEEHHHHHHHHHHHCCEEECCCCEEEEEEEECCCCCCCCCEEEEE
AA: GPGWFFVVDDEGPRRFTLRDWQLDRERALTFAVEIPGARTVTACQVRTEPGERGRTLSVS
      190          200          210          220          230          240
```

Domain Partition Problem

Experimental

T1011-D2



T1011-D1

CASP domain: D1:55-268,433-520 D2: 271-430
Zhang-Server domain: D1:8-267,434-519 D2:268-433

Zhang-Server model1

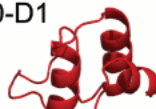
D1: TM-score=0.81

D2: TM-score=0.87



Experimental Domain Partition

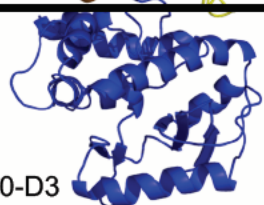
T0990-D1



T0990-D2



T0990-D3



T0990-D1 TM-score=0.57

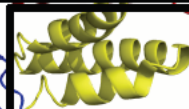
T0990-D2 TM-score=0.38

T0990-D3 TM-score=0.21

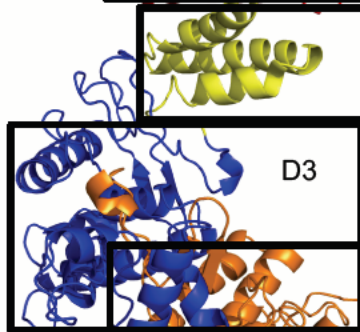
D1



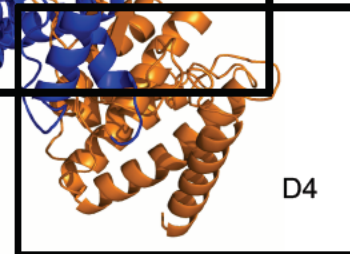
D2



D3

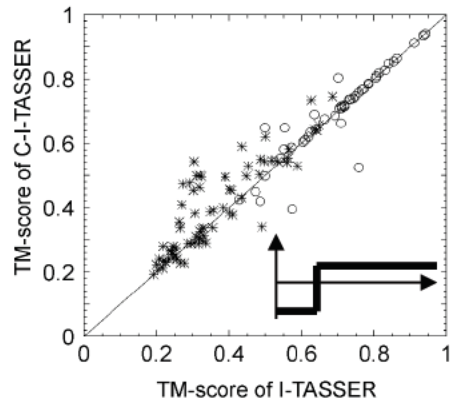
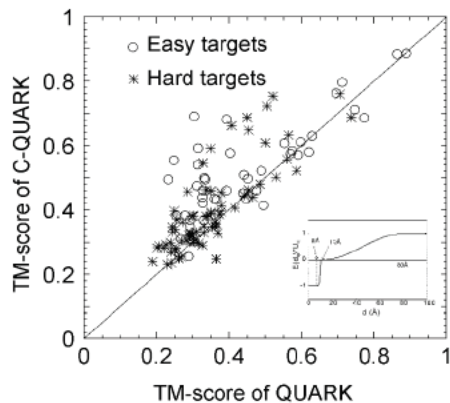


D4

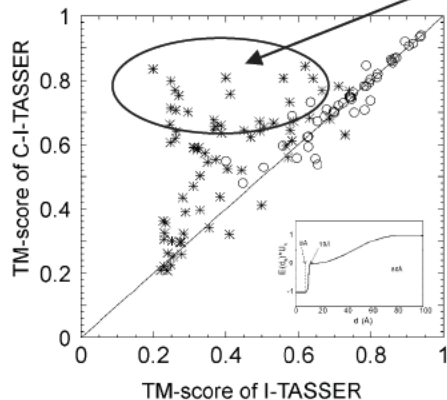
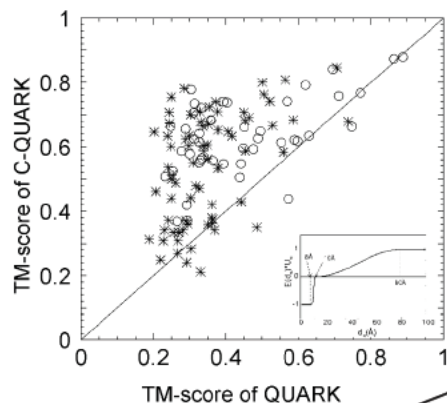


CASP domain: D1:1-76 D2:77-134,348-520 D3: 135-347
Zhang-Server domain: D1:1-128 D2:129-280 D3:281-406 D4: 407-552

Impact of Contact Maps on 3D Structure Prediction



CASP12 (SVMSEQ+NeBcon)



CASP13 (ResPRE+NeBcon)

QUARK vs
C-QUARK

Many non-homology
proteins can now
fold as well as
homology proteins

I-TASSER vs
C-I-TASSER

Summary

- What worked
 - Contact prediction by Deep Neural Networks
 - Integration of contact potential, LOMETS restraints, and inherent potential
- What needs improvement
 - Secondary structure prediction
 - Discontinuous domain partitioning

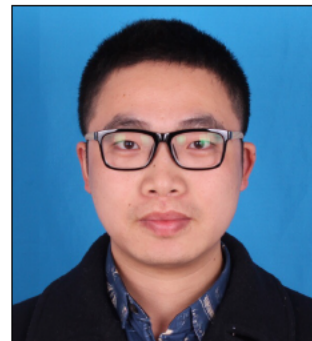
Acknowledgements



Wei Zheng



Chengxin Zhang



Yang Li



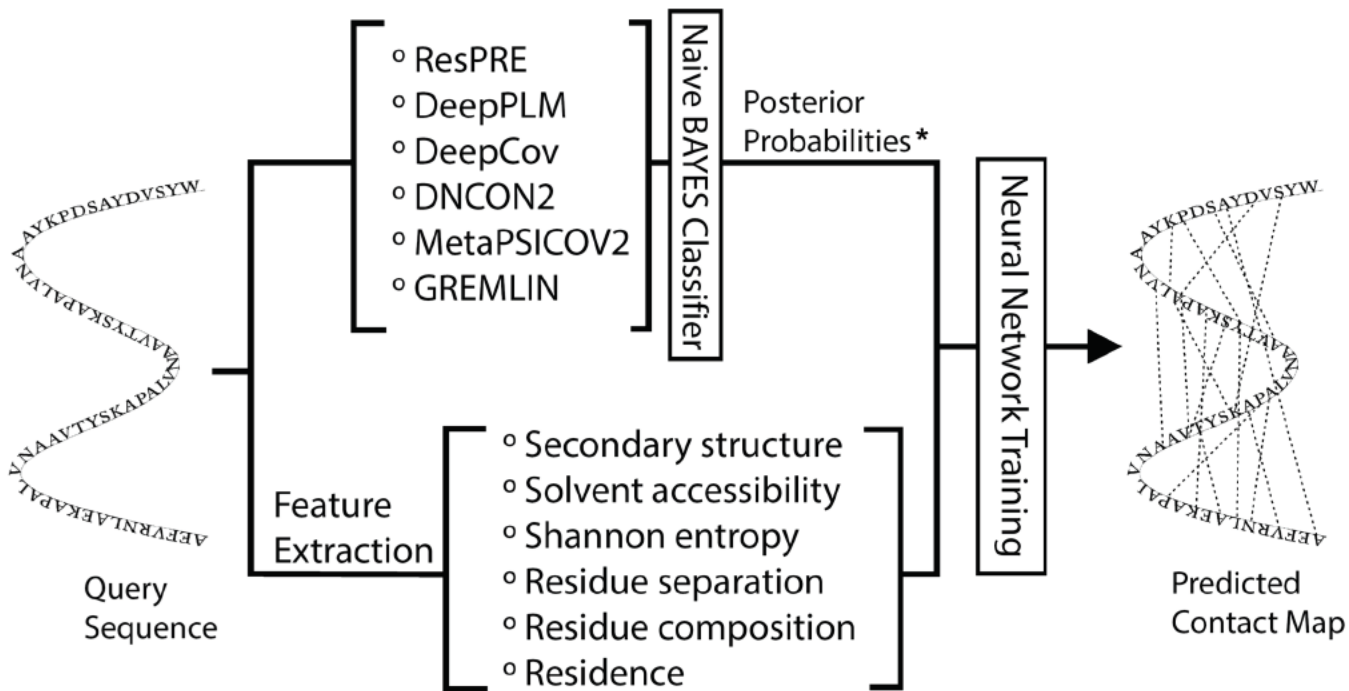
S M Mortuza



Yang Zhang

- **Experimentalists that shared their structures**
- **CASP Organizers**

Contact Prediction using NeBcon



$$*P(C|X_{ij}^1, X_{ij}^2, \dots, X_{ij}^N) = \frac{P(C) \prod_{m=1}^N P(X_{ij}^m|C)}{P(0) \prod_{m=1}^N P(X_{ij}^m|0) + P(1) \prod_{m=1}^N P(X_{ij}^m|1)}$$