# D-I-TASSER: Integrating Deep Learning with multi-MSAs and threading alignments for protein structure prediction

## (Groups "UM-TBM" and "Zheng")

Wei Zheng[1,2], Qiqige Wuyun[3], and Peter L. Freddolino[1,2]

[1]Department of Computational Medicine and Bioinformatics , University of Michigan

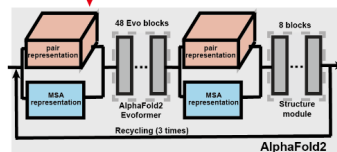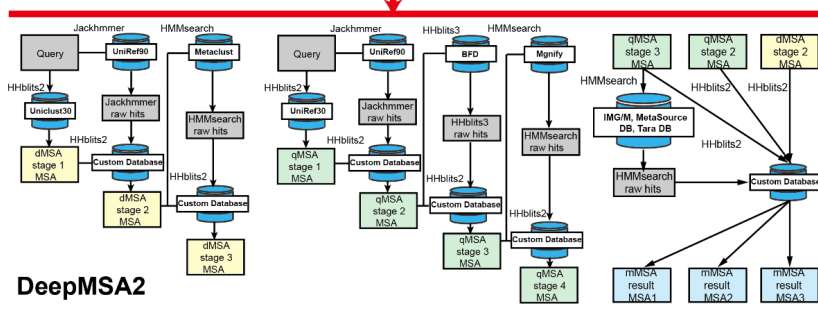[2]Department of Biological Chemistry, University of Michigan

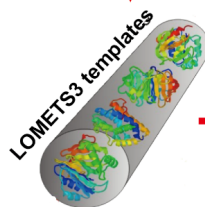[3]Department of Computer Science and Engineering, Michigan State University

# Methods

# UM-TBM server built from D-I-TASSER for single-chain protein modeling

# UM-TBM server built from D-I-TASSER for single-chain protein modeling

# DeepMSA2 for monomeric MSA construction

# UM-TBM server built from D-I-TASSER for single-chain protein modeling



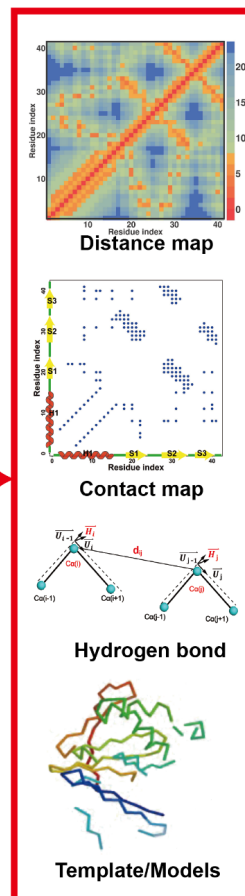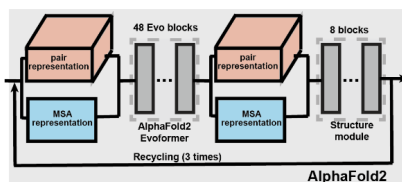**Query sequence** TTSQKHRDFVAEPGEKPVGFLVLKVGFLVLKVAELVLKVGFLPGRDFEPG

$$E=E_{knowledge}+E_{template}+E_{distance}+E_{contact}+E_{HB}$$

DeepMSA2

Ranking MSAs by the predicted models' pLDDT

Final MSA

LOMETS3 templates

DeepPotential AttentionPotential AlphaFold2

Distance map
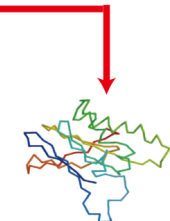
Contact map

Hydrogen bond

Template/Models

Distance/Contact/HB-guided simulation

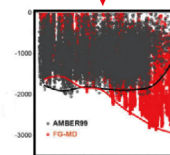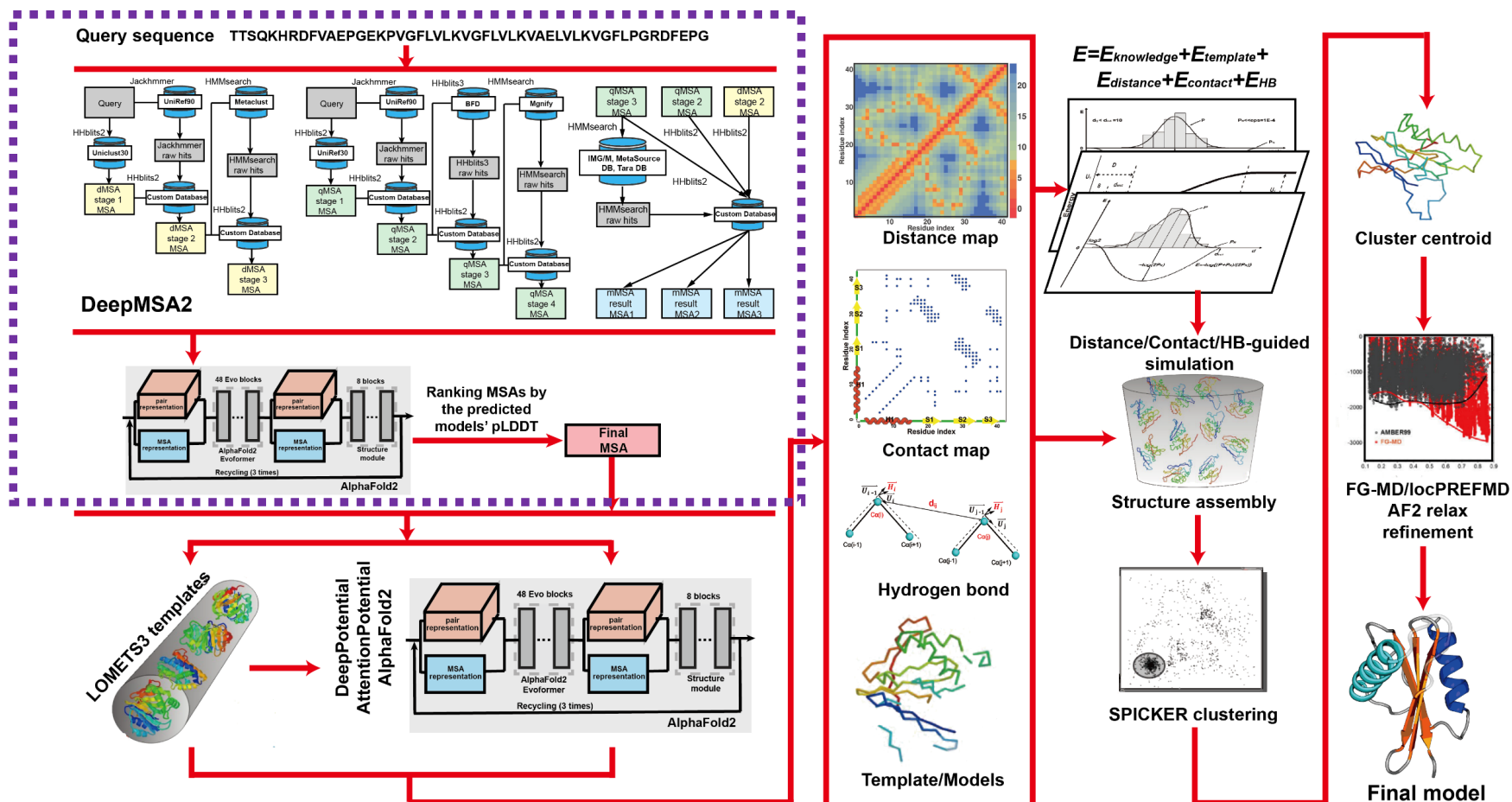Structure assembly

SPICKER clustering
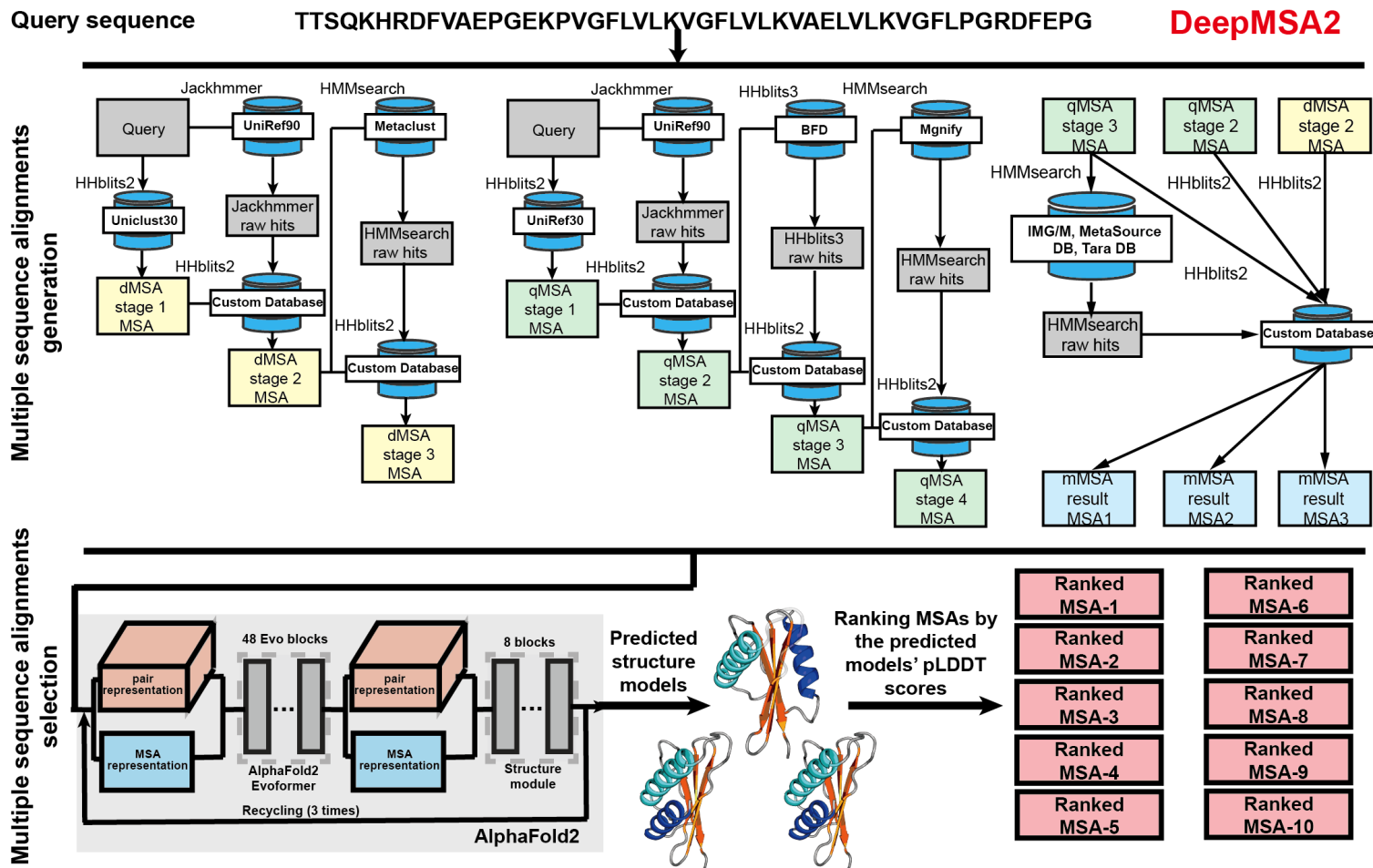
Cluster centroid

FG-MD/locPREFMD AF2 relax refinement

Final model

# LOMETS3 for threading template detection



Query sequence

MIKFLSALILLLVTTAAQAERIRDLTSVQGVRQNTQTLNNMLSQLGITVPTGTNMQLKNVAAVMVTAS

Full-length MSA

Full-length

Contact-map:
**FUpred** for domain boundary prediction for non-homologous targets

Threading templates:
**ThreaDom** for domain boundary prediction for homologous targets

AlphaFold2

Template

Domain-level sequences

Domain-level MSAs

GQTIDVVVSSMGNAKSL

YGLVVGLQTFTQTLNNM

MIKFLSALILLLVTT

AlphaFold2 etc

Contacts/Distance/HB    Template    LOMETS threading

Reference full-length model by **DeepFold**

PDB database

Domain-level templates

Domain-level Distance maps

DEMO2 assembly    Assembled templates
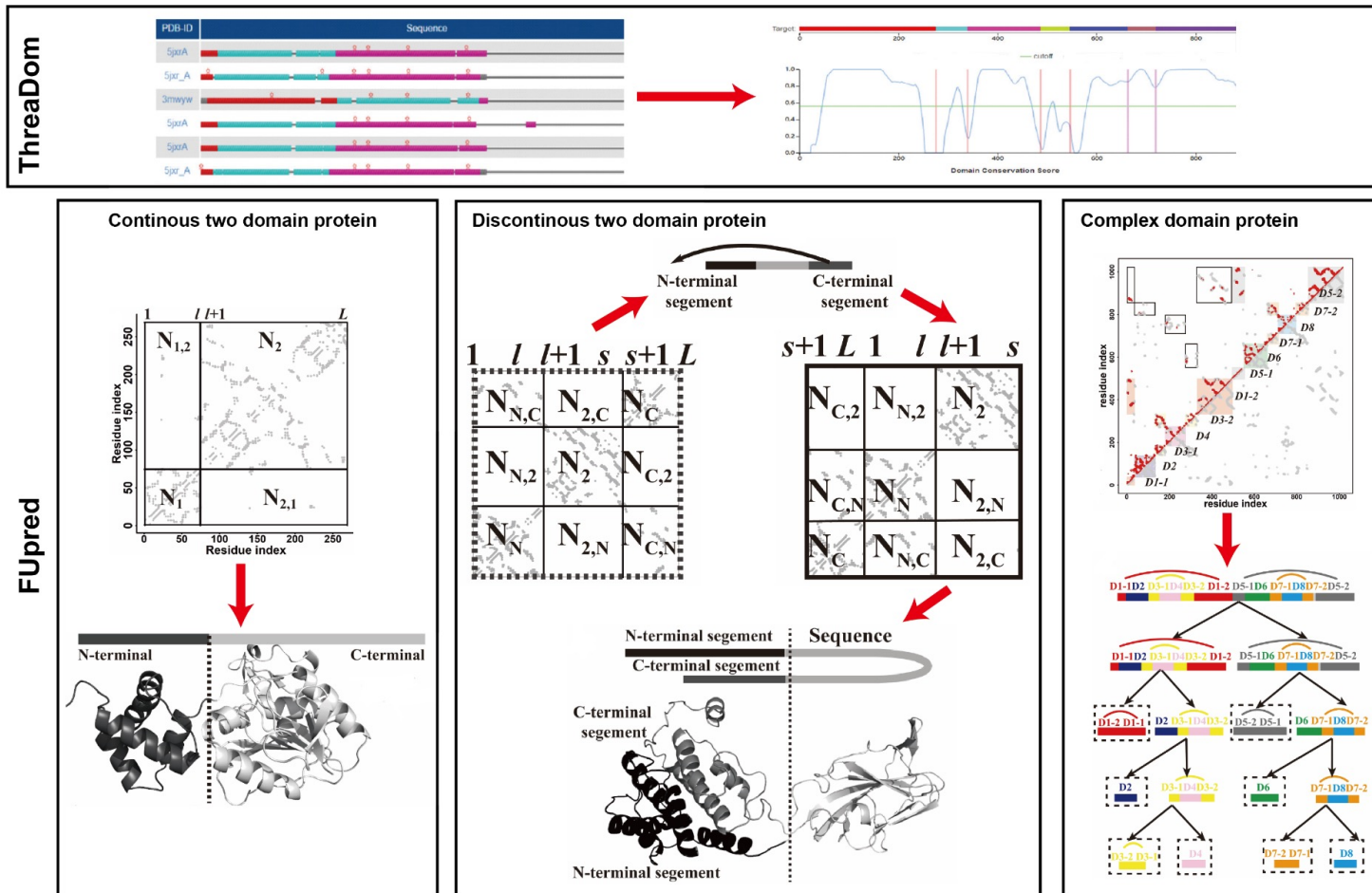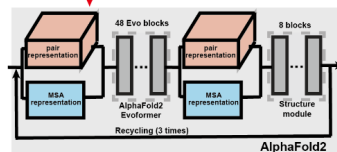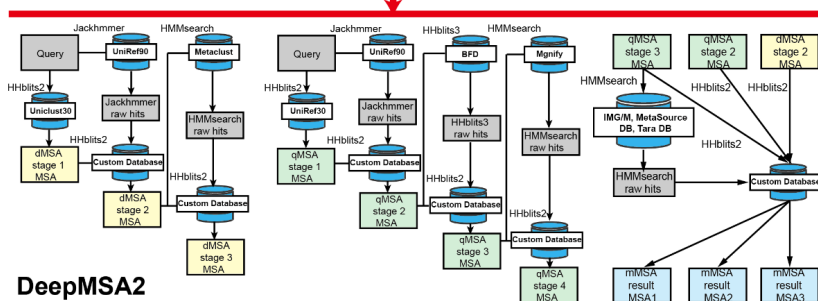
**The domain-level distance maps and the whole chain-level distance map will be combined to guide the D-I-TASSER folding simulation**
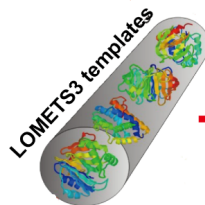
# Domain boundary prediction in UM-TBM server

# UM-TBM server built from D-I-TASSER for single-chain protein modeling



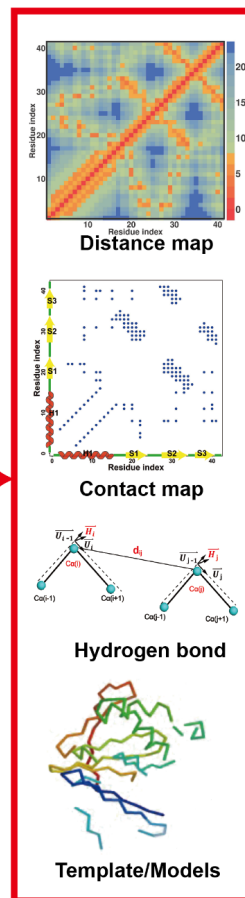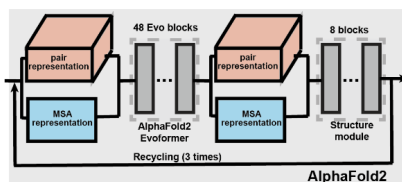**Query sequence** TTSQKHRDFVAEPGEKPVGFLVLKVGFLVLKVAELVLKVGFLPGRDFEPG
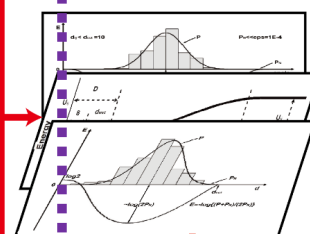
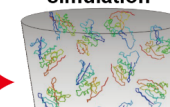**DeepMSA2**

Ranking MSAs by the predicted models' pLDDT

**Final MSA**

**AlphaFold2**

**LOMETS3 templates**

**DeepPotential AttentionPotential AlphaFold2**

**AlphaFold2**

**Distance map**

**Contact map**

**Hydrogen bond**

**Template/Models**

$E = E_{knowledge} + E_{template} + E_{distance} + E_{contact} + E_{HB}$
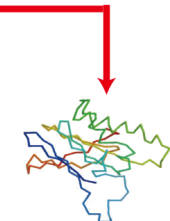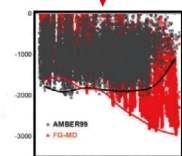
**Distance/Contact/HB-guided simulation**

**Structure assembly**

**SPICKER clustering**

**Cluster centroid**

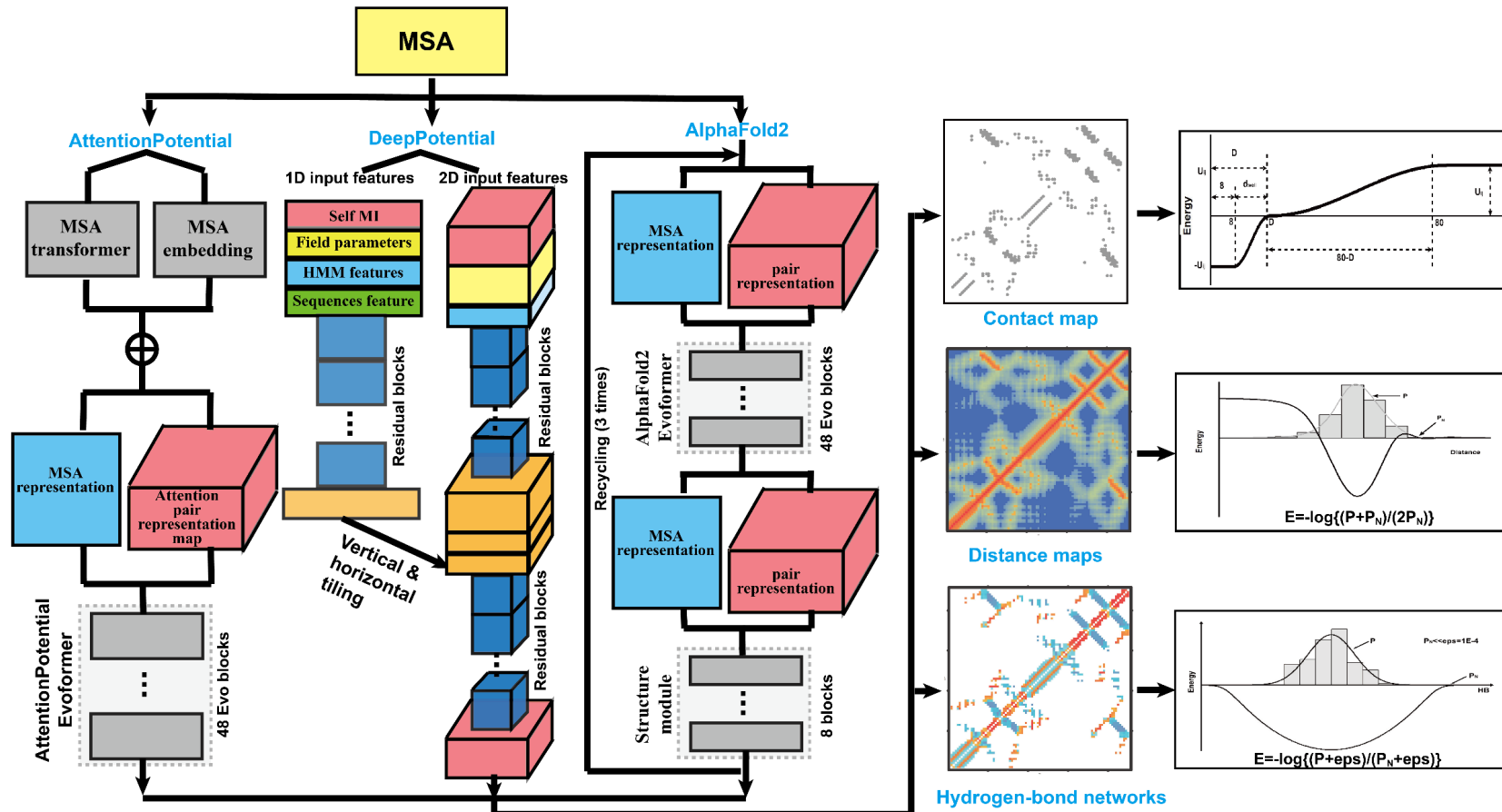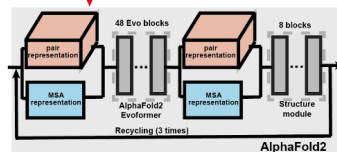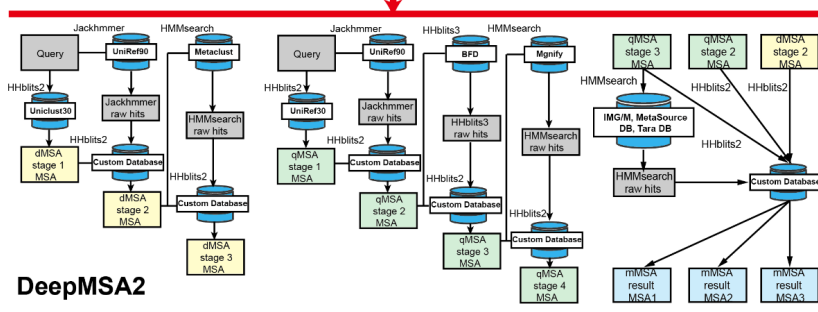**FG-MD/locPREFMD AF2 relax refinement**

**Final model**
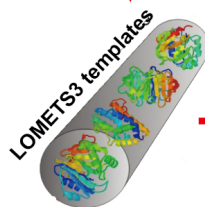
# Deep learning-based spatial restraints prediction

# UM-TBM server built from D-I-TASSER for single-chain protein modeling

# 'Zheng' built on DeepMSAFold-Multimer for protein complex modeling

# Zheng built on DeepMSAFold-Multimer for protein complex modeling

# Zheng built on DeepMSAFold-Multimer for protein complex modeling

# Homo-oligomer MSA paring and sequence linking

# Heteromer MSA pairing and sequence linking



For a heteromer A2B2C1

$$P = N^M \leq 100$$

# Zheng is based on DeepMSAFold-Multimer for protein complex modeling

# Complex model generation based on Multi-MSAs

# Results for UM-TBM (D-I-TASSER) server

# Summary of FM targets folded by I-TASSER series algorithm

# D-I-TASSER vs standard AlphaFold2 on CASP15 domains

**T1169-D1**

AF2 TM=0.15  UM-TBM TM=0.84
FM target, 345AA

**T1169-D3**

AF2 TM=0.77  UM-TBM TM=0.96
TBM-hard target, 401AA

**T1169-D2**

AF2 TM=0.93  UM-TBM TM=0.94
FM target, 1434AA

**T1169-D4**
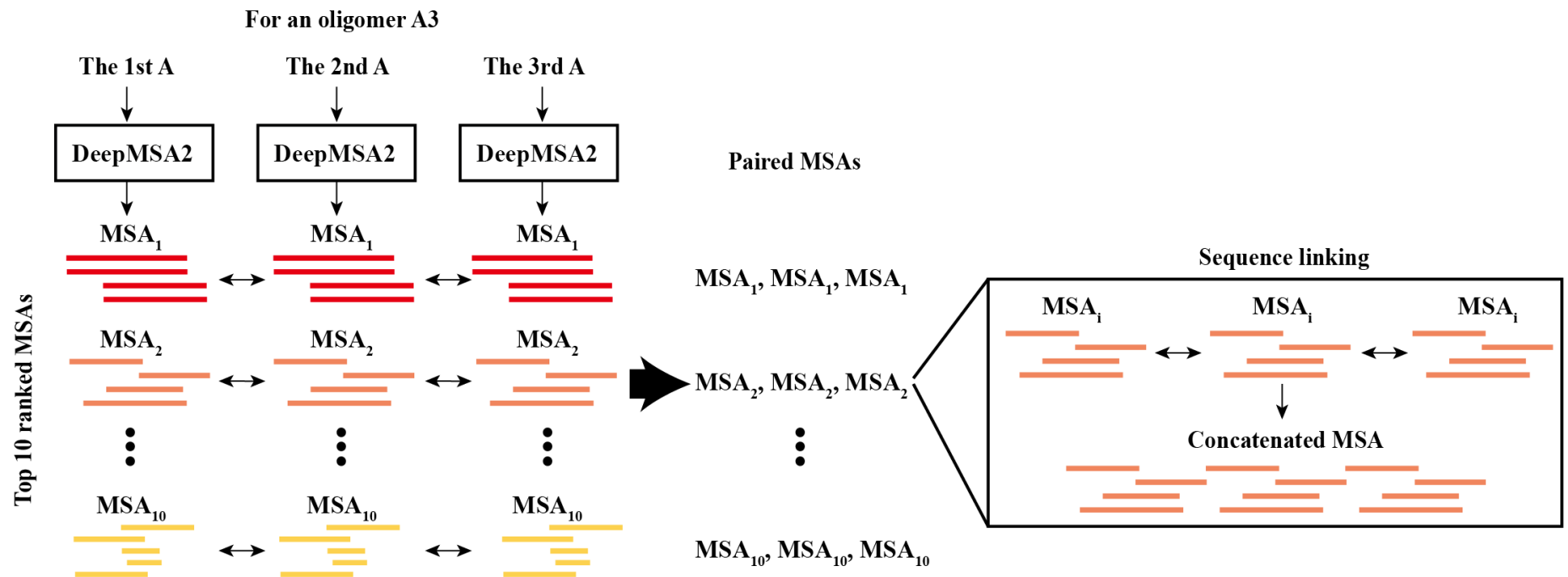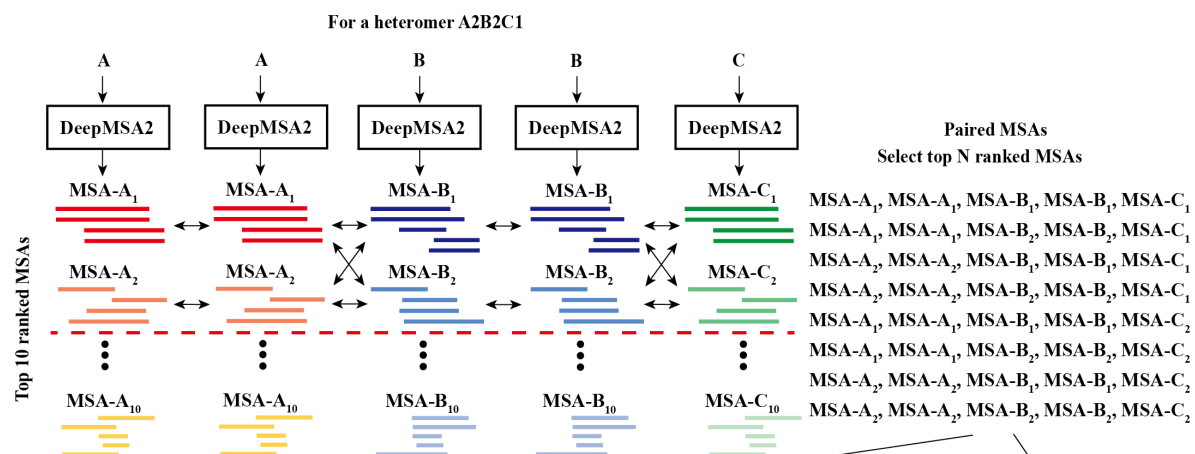
AF2 TM=0.64  UM-TBM TM=0.94
FM target, 523AA

TM-score of AlphaFold2 model

TM-score of UM-TBM (D-I-TASSER) model

○ TBM-easy & TBM-hard
✕ FM/TBM & FM

■ Experimental structure  ■ AlphaFold2 model  ■ UM-TBM model

Standard AlphaFold2 data are taken from CASP15 NBIS-AF2-standard server.

# Impact of different components in D-I-TASSER

Alphafold2 (with LOMETS3 and DeepMSA2) models are used as the initial conformation of D-I-TASSER simulations, thus the further improvement represents the contribution of D-I-TASSER folding simulation and other modules used in the UM-TBM pipeline.

# A case study of T1125-D2 to highlight the advance of D-I-TASSER folding



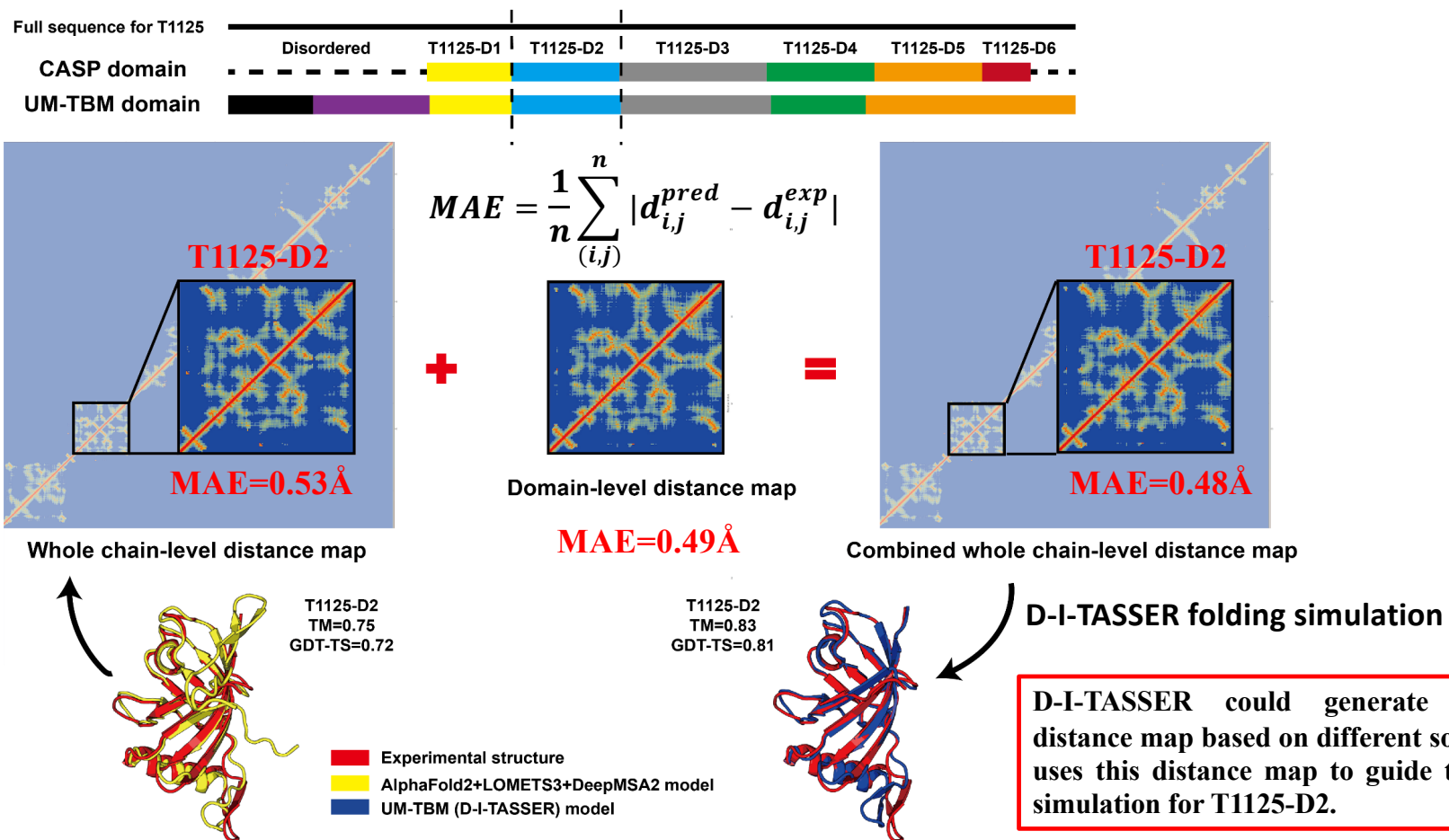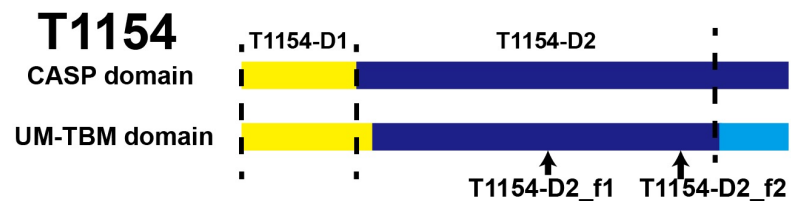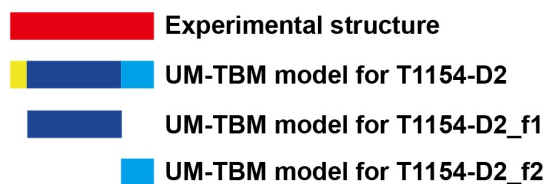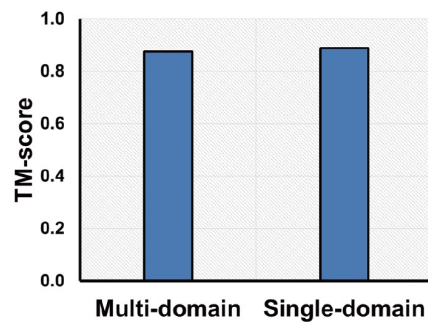$$MAE = \frac{1}{n}\sum_{(i,j)}^{n} |d_{i,j}^{pred} - d_{i,j}^{exp}|$$

Whole chain-level distance map

Domain-level distance map
MAE=0.49Å

Combined whole chain-level distance map

D-I-TASSER folding simulation

T1125-D2
TM=0.75
GDT-TS=0.72

T1125-D2
TM=0.83
GDT-TS=0.81

■ Experimental structure
■ AlphaFold2+LOMETS3+DeepMSA2 model
■ UM-TBM (D-I-TASSER) model

D-I-TASSER could generate a better distance map based on different sources, and uses this distance map to guide the folding simulation for T1125-D2.

Domain partition problem in T1154-D2

# Results for 'Zheng' (DeepMSAFold-Multimer) group

# DeepMSAFold-Multimer vs standard AlphaFold2-Multimer

**Global fold modeling quality**
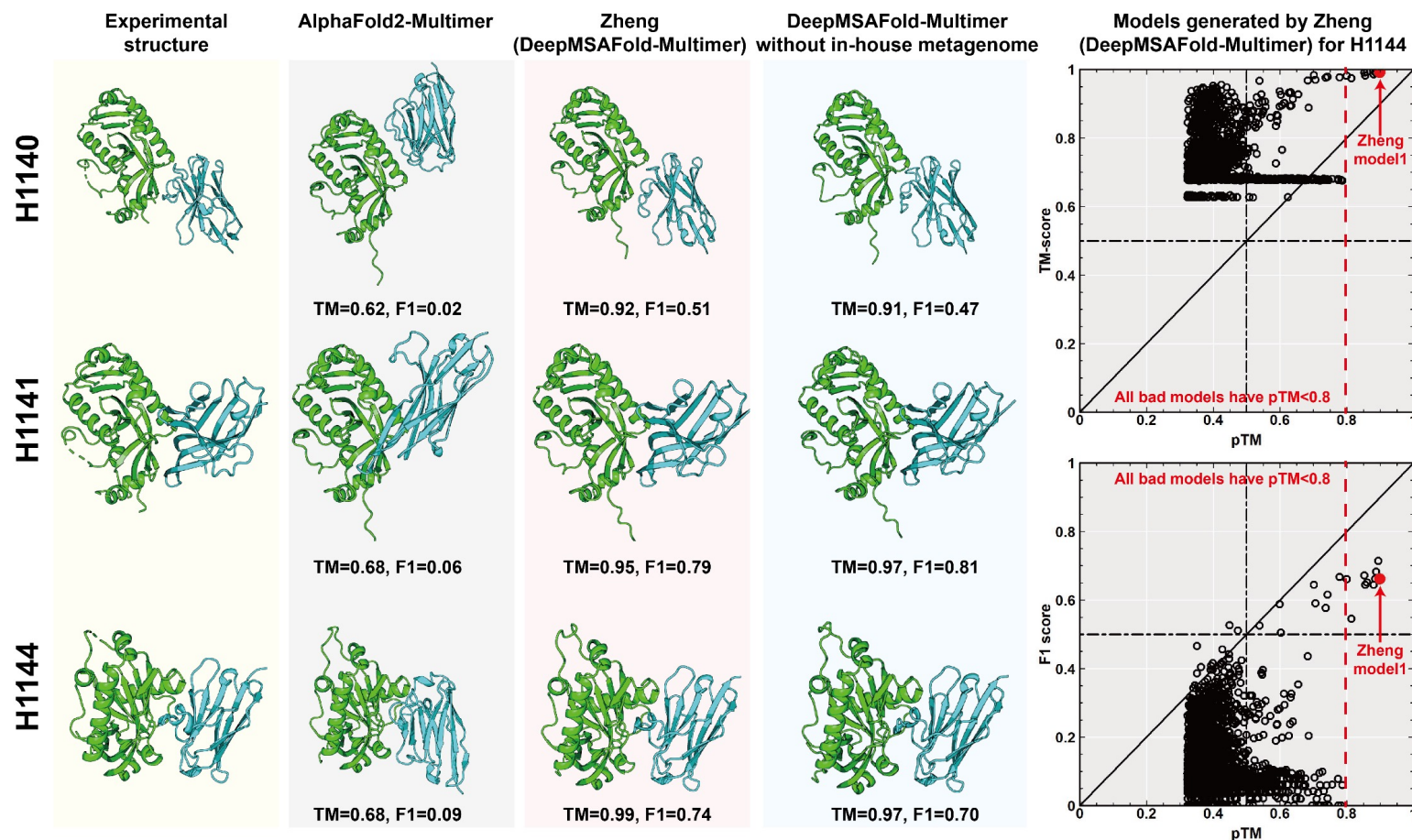
**Interface modeling quality**



Standard AlphaFold2 data are taken from CASP15 NBIS-AF2-multimer server.

# Impact of the MSA generation and database for complex modeling

# High quality model generation for Nanobody-antigen by pairing MSAs



DeepMSAFold-Multimer works well for producing high quality models by paring the different MSAs and ranking the model with pTM.
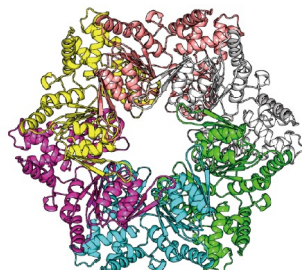
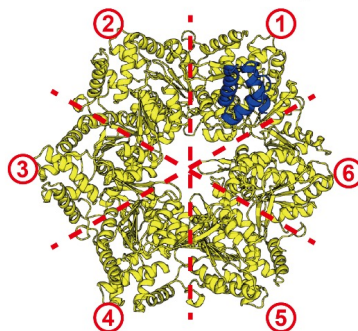# Model ranking problem with pTM score in H1172



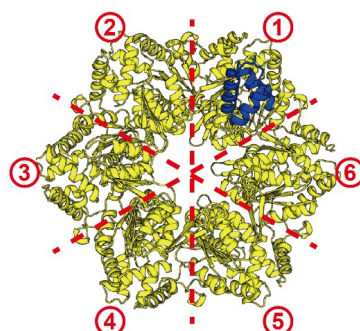**T1170o, A6, 1908AA**

Experimental structure

Zheng model1
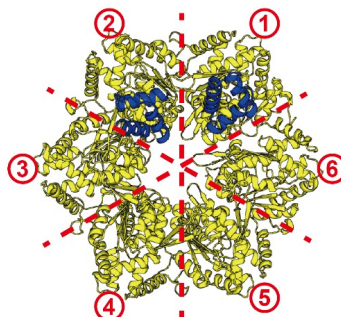pTM=0.793
TM=0.93, F1=0.58

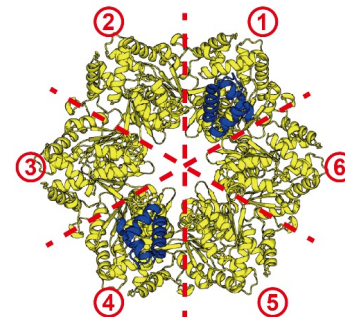**H1171, A6B1, 1956AA**

Experimental structure

Zheng model1
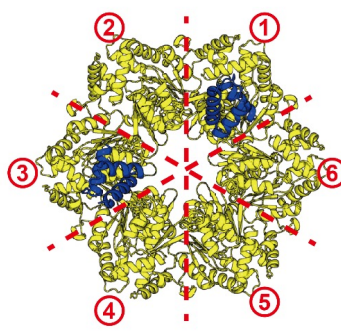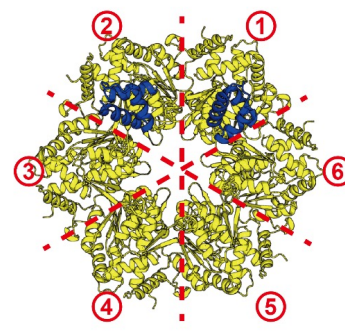pTM=0.764
TM=0.93, F1=0.51

**H1172, A6B2, 2004AA**

Experimental structure

Zheng model1
pTM=0.746
TM=0.91, F1=0.54

Zheng model2
pTM=0.738
TM=0.91, F1=0.50

Zheng model3
pTM=0.733
TM=0.93, F1=0.54

The pTM scores are not sufficient for identifying the best model among models that have similar pTM scores.

# Summary

- What went right
  - D-I-TASSER algorithm that integrates threading templates, structure model-ranked MSAs, deep learning-based spatial restraints, and an optimized folding system with comprehensive force field works well for protein monomer structure prediction.
  - DeepMSA2 with high-quality MSA generating, ranking, and paring system helps improve the model quality for both protein monomer and protein complex.

- What went wrong
  - AlphaFold2's QA score (pLDDT and pTM) has the ability to distinguish models that are significantly better, but it is not sensitive enough to rank high-quality models.
  - Domain partition is still meaningful and important for large multi-domain protein modeling, however the accuracy of domain boundary prediction is still not yet satisfactory.

# Acknowledgements
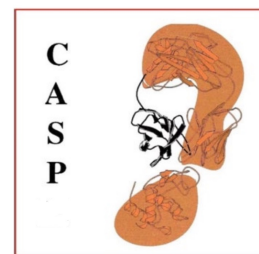


Yang Zhang

Yang Li

Robin Pearce

Jonathan Poisson

# Thank you
# Q&A