# EnsembleFold: Alternative conformation prediction using multi-MSA strategy and structural clustering

12/04/2024

**Speaker: Wei Zheng**

Wei Zheng, Qiqige Wuyun, Quancheng Liu, Chunxiang Peng,
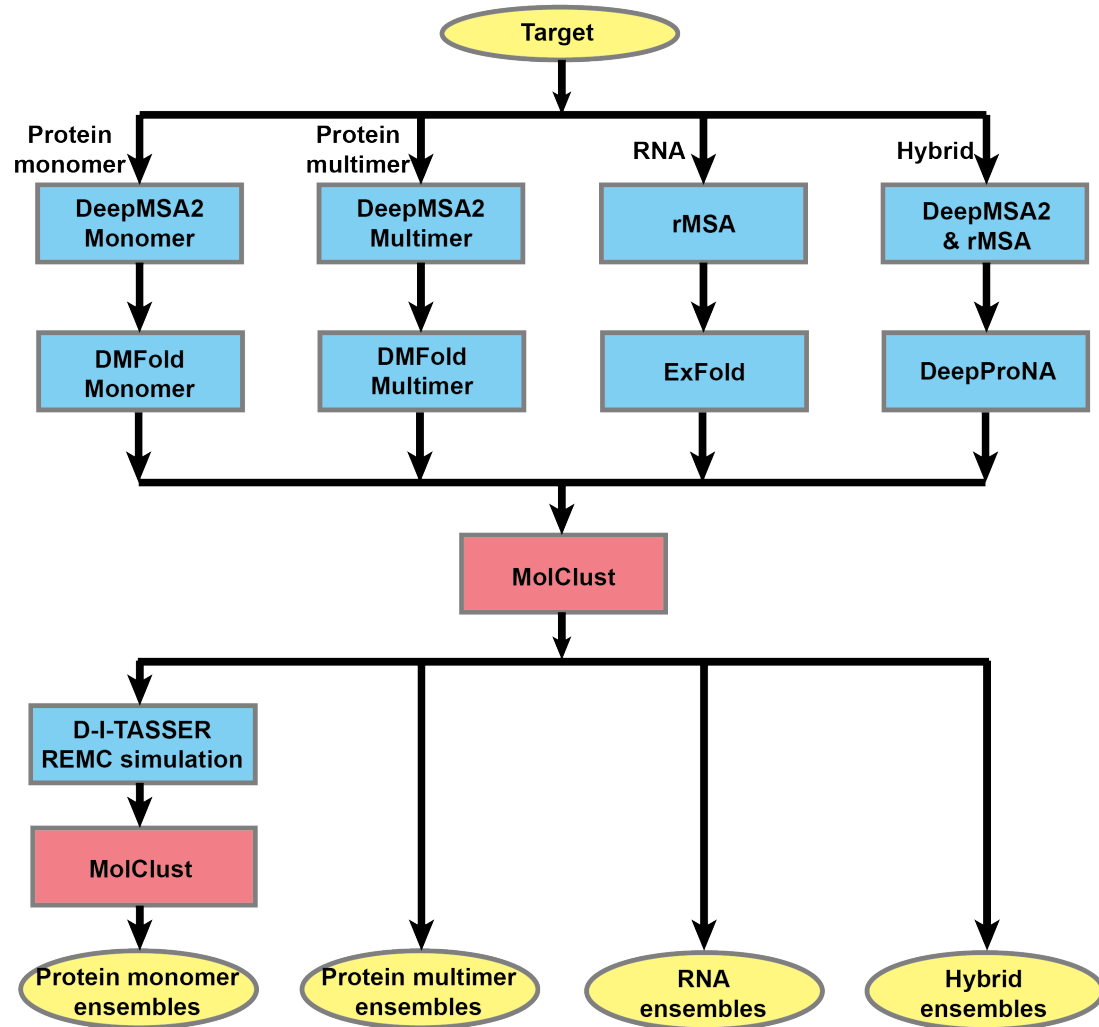Wentao Ni, Gang Hu, Ziying Zhang, Xiaogen Zhou, Lydia Freddolino
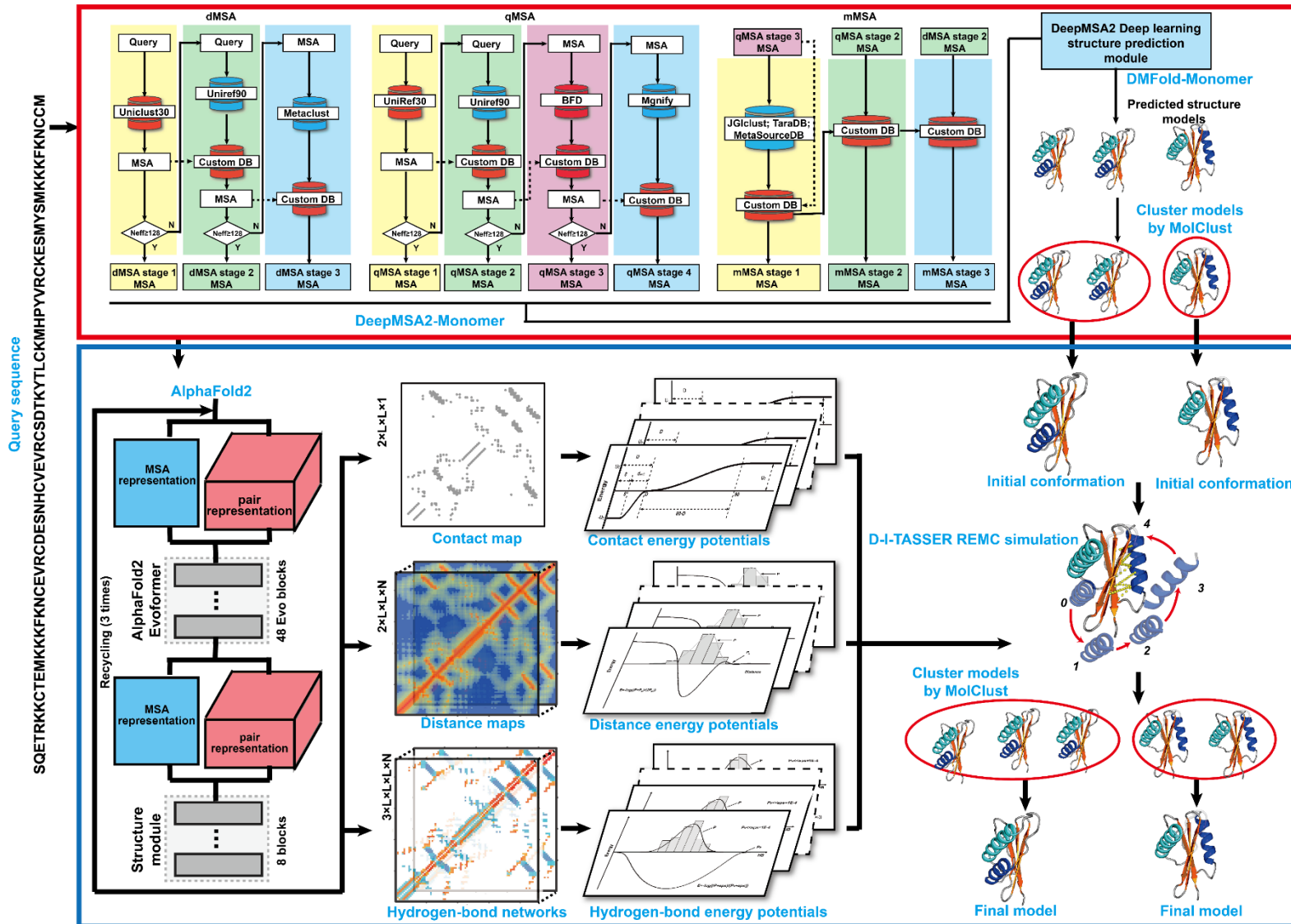
- Zheng-Server, Zheng-Multimer, Zheng, MIEnsembles-Server, and NKRNAs participated in CASP16
- MIEnsembles-Server (server group) and Zheng (human group) focus on ensemble targets
- Same pipeline, Zheng has longer running time and more combinations of MSAs
- Four different pipelines for handling protein monomer, protein complex, RNA, and hybrid targets
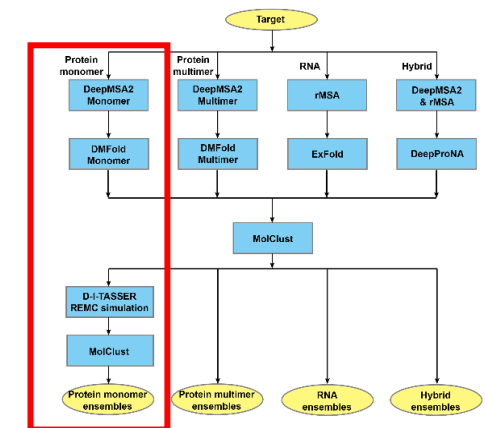
**DeepMSA2-Monomer & DMFold-Monomer**

**Key points:**
1. Using all MSAs from DeepMSA2
2. Using DMFold models and spatial restraints from all representative models in replica exchange Monte Carlo (REMC) simulation
3. Clustering decoys from REMC simulation

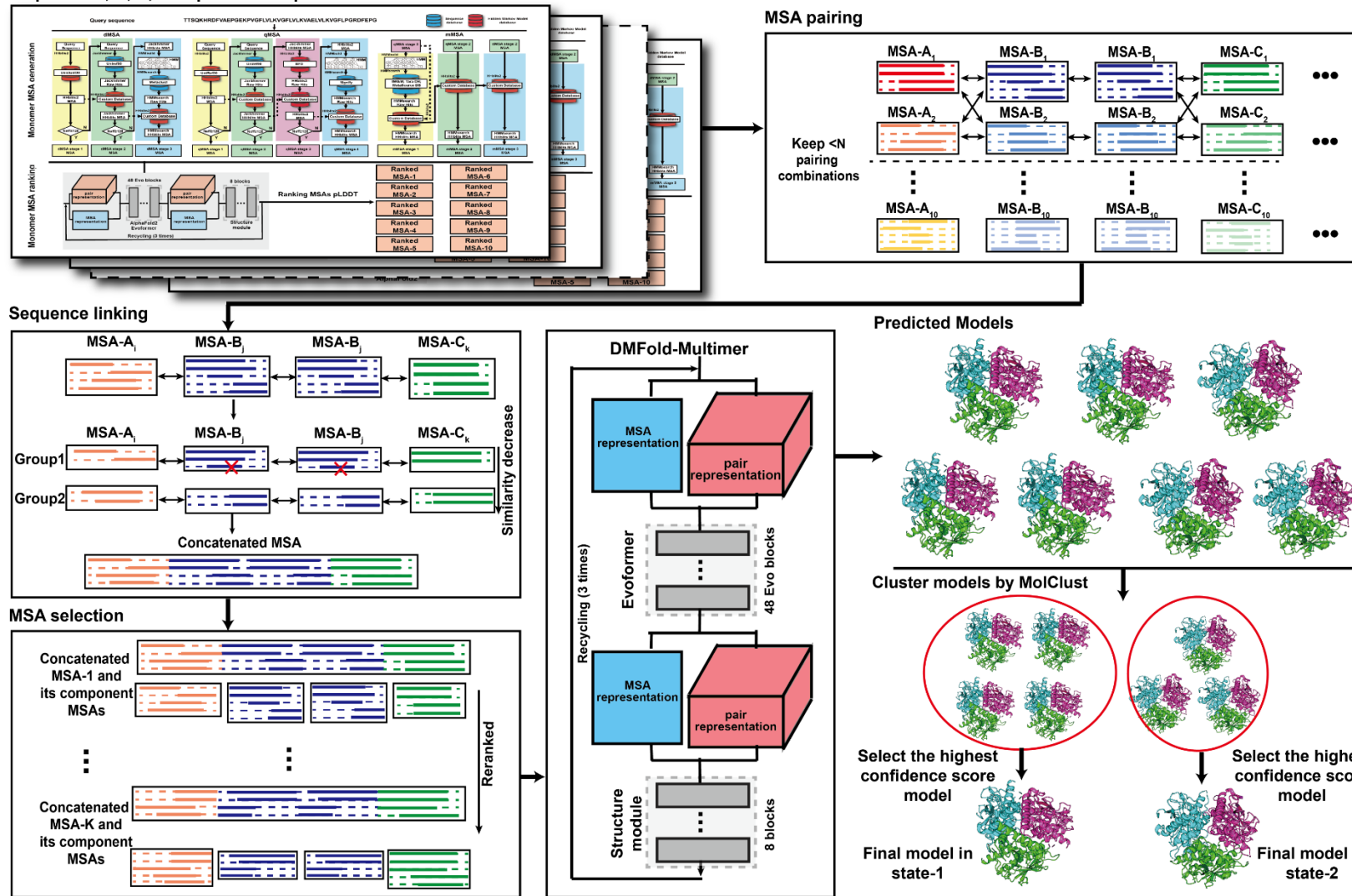**D-I-TASSER REMC simulation**

**Structural clustering by MolClust**
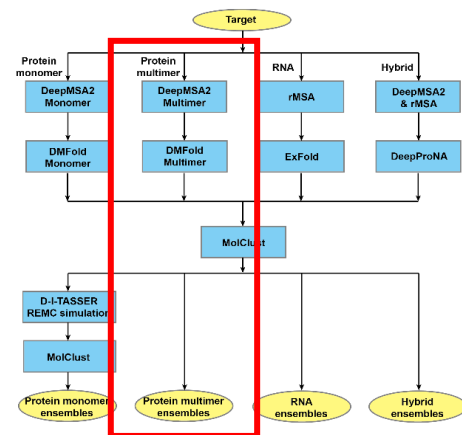1. SPICKER
2. US-align

**For targets: T1214, T1200 and T1300.**

4

**Key points:**
1. Larger metagenomes than CASP15 version
2. More combinations of MSA pairing
3. Sampling strategy in modeling stage: using the template or not, opening the drop up rate or not, and using different alphafold2 pre-trained parameters (v1 v2 v3).
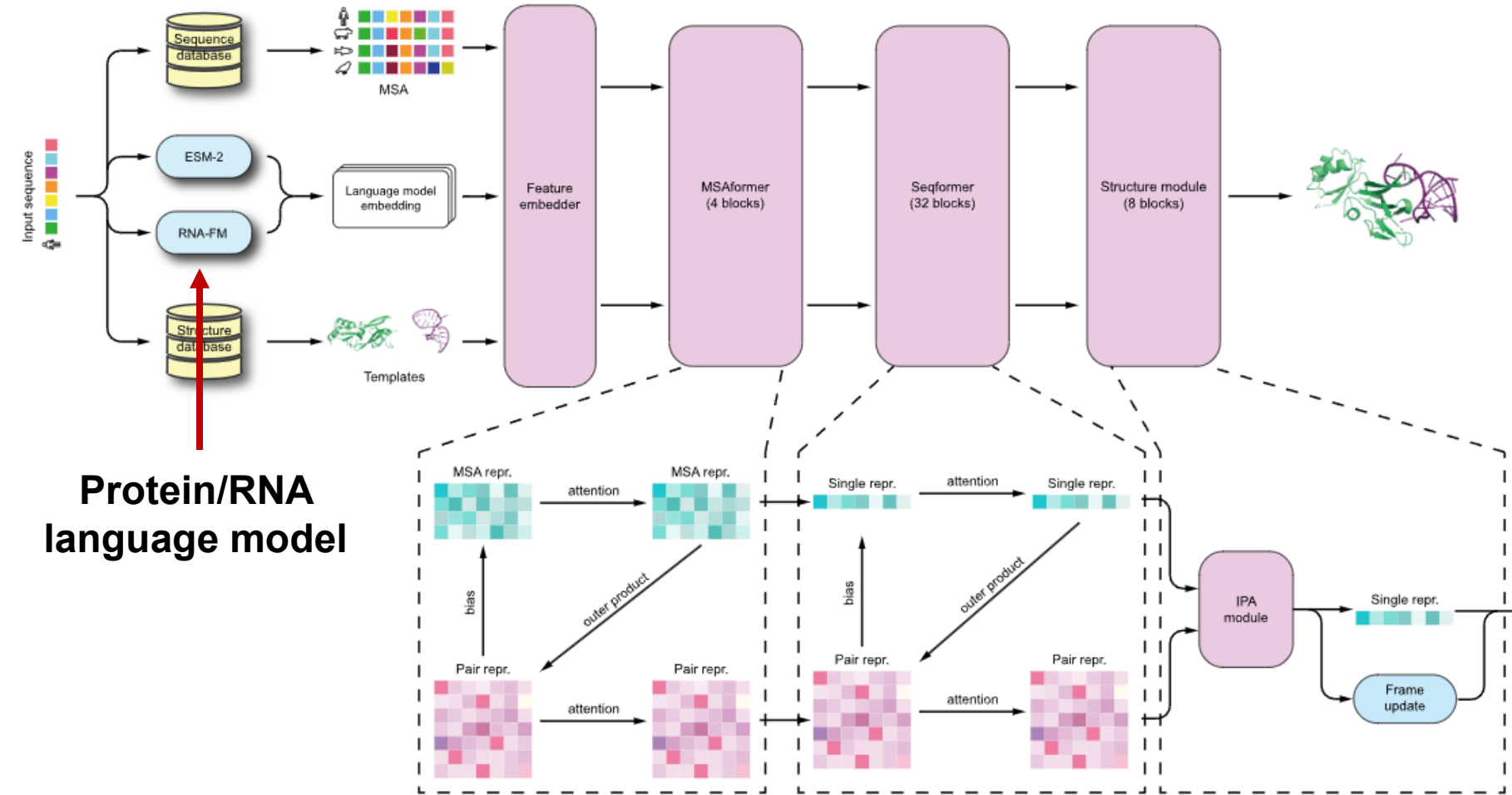4. Clustering models by structural similarity, rank by highest confidence score of the members

**For targets: T1249v1/v2 and T1294v1/v2.**

5

Wentao Ni

**Protein/RNA language model**

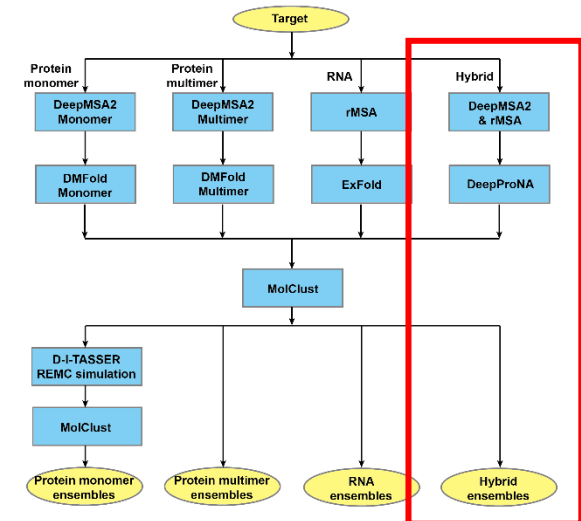**Key points:**
1. Modified from AlphaFold2 pipeline
2. Using Protein/RNA language model
3. Using multiple sets of MSAs as input
4. Clustering the models

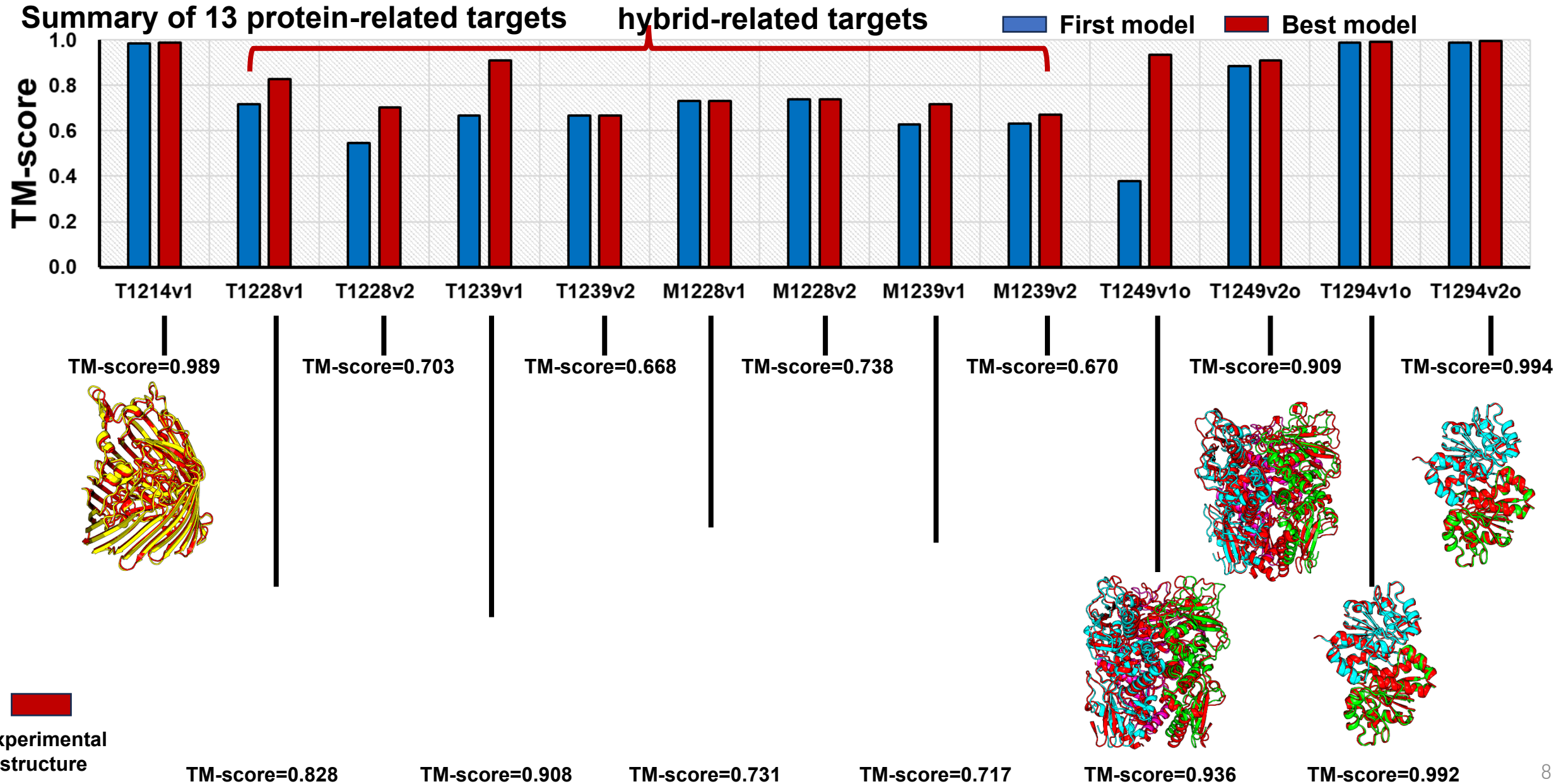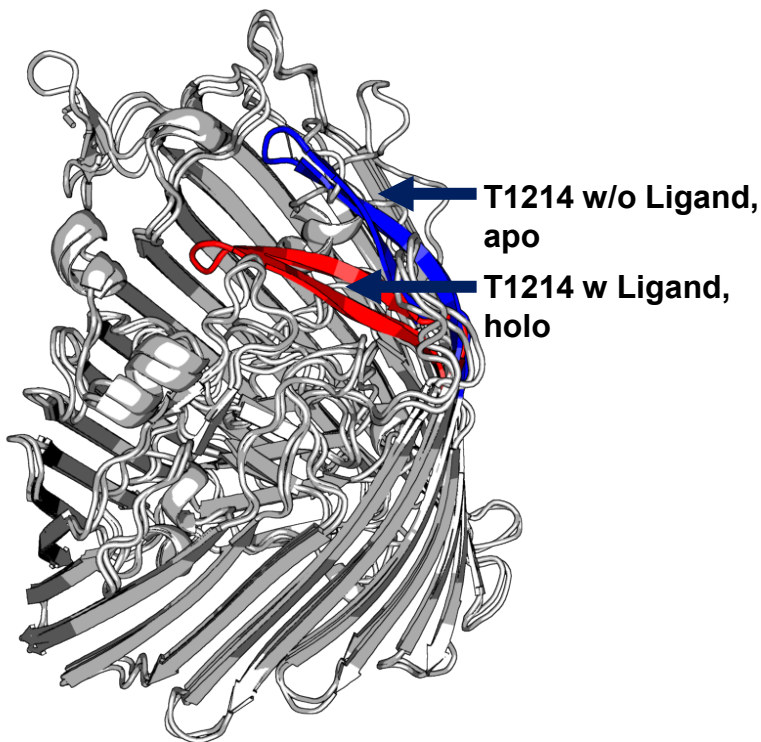For targets: M1228v1/v2, T1228v1/v2, M1239v1/v2, and T1239v1/v2.

| 1 | Methods |
| --- | --- |
| **2** | **Results** |

Summary of 13 protein-related targets

hybrid-related targets

First model   Best model

TM-score

1.0
0.8
0.6
0.4
0.2
0.0

T1214v1  T1228v1  T1228v2  T1239v1  T1239v2  M1228v1  M1228v2  M1239v1  M1239v2  T1249v1o  T1249v2o  T1294v1o  T1294v2o

TM-score=0.989   TM-score=0.703   TM-score=0.668   TM-score=0.738   TM-score=0.670   TM-score=0.909   TM-score=0.994

Experimental structure

TM-score=0.828   TM-score=0.908   TM-score=0.731   TM-score=0.717   TM-score=0.936   TM-score=0.992

**MolClust result for DMFold-Monomer decoys**

X-axis: TM-score of decoy to T1214 w/o ligand
Y-axis: TM-score of decoy to T1214 with ligand

**DMFold-Monomer models**

1 — TM-score=0.947 qMSA stage 3 MSA *Neff*=390.3 — Model

2 — Model — TM-score=0.973 qMSA stage 2 MSA *Neff*=106.9 — Model

3 — TM-score=0.976 dMSA stage 1 MSA *Neff*=860.5

T1214 w/o Ligand, apo

T1214 w Ligand, holo

**MolClust result for EnsembleFold decoys**

X-axis: TM-score of decoy to T1214 w/o ligand
Y-axis: TM-score of decoy to T1214 with ligand

New cluster (after REMC simulation)

**EnsembleFold models**

First model TM-score=0.986

Best model TM-score=0.989

Superpose ligand to the model

Prediction without ligand but ligand fits well!
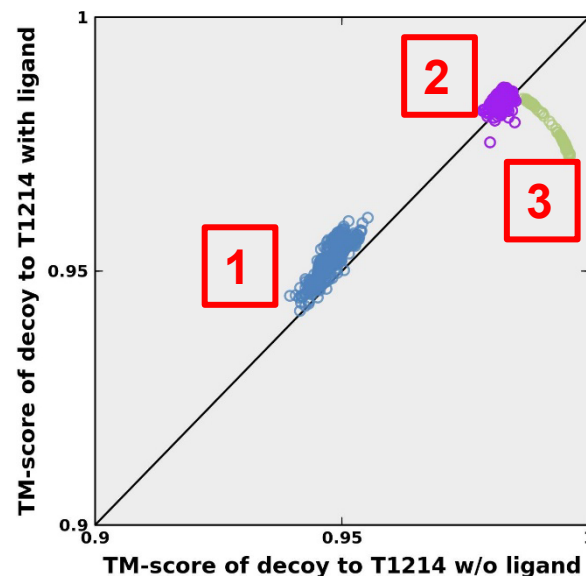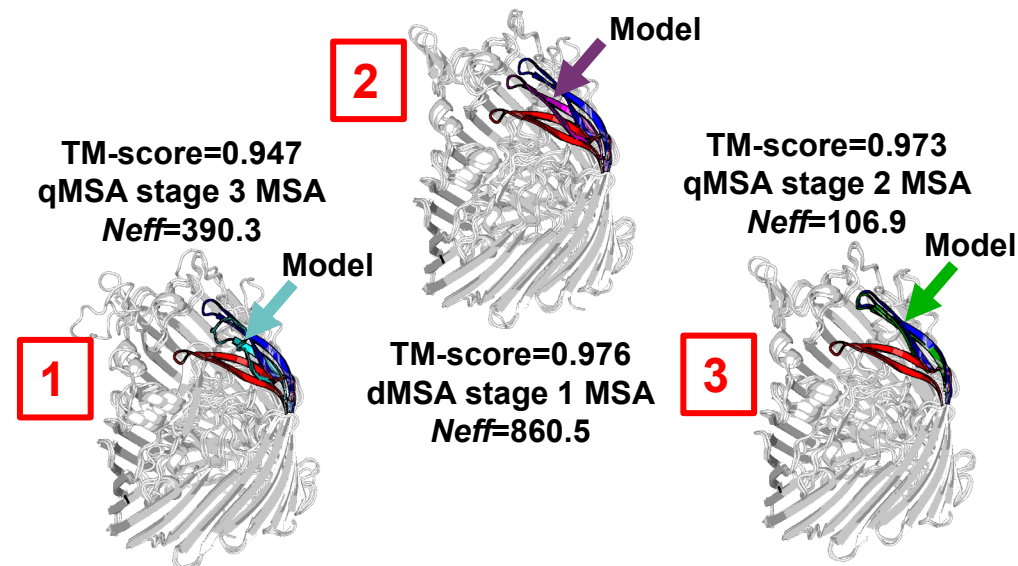
**Experimental structures**

**Model this target without ligand**

- The model from each MSA corresponds to one 'state'
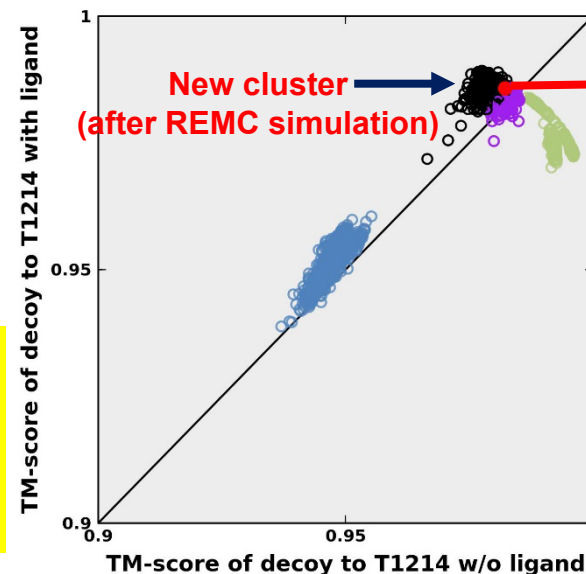- REMC simulation helps create a diverse set of models

9

**Experimental structures**

M1228v1

M1228v2

**Best decoy for M1228v1**
**TM-score=0.780**
**QA=0.524**

**MolClust result for EnsembleFold decoys**

**Best decoy for M1228v2**
**TM-score=0.798**
**QA=0.612**

**Model1 for T1228v1**
**TM-score=0.718**

**Model1 for T1228v2**
**TM-score=0.545**

Structural clustering works well for picking the correct model of each state in this case

protein with both states shown

**Experimental structures**

M1239v2

M1239v1

**Best decoy for M1239v2**
**TM-score=0.796**
**QA=0.456**



**MolClust result for EnsembleFold decoys**



**Best decoy for M1239v1**
**TM-score=0.816**
**QA=0.544**



Model1 for T1239v2
TM-score=0.668

Model1 for T1239v1
TM-score=0.668

**Structural clustering works well for picking the correct model of each state in this case**

protein with both states shown

**Experimental structure**
**Trimer**



TM-score=0.83

**Experimental structures**

- 🟥 **Open state**
- 🟦 **Closed state**



TM-score of decoy to T1249 open state (y-axis)
TM-score of decoy to T1249 closed state (x-axis)

1
2
3

**Predicted open state**

**Open state** | **Closed state**

**Cluster1**
Correct open state — TM-score=0.954
Predicted closed state — TM-score=0.884
QA=0.566

**Cluster2**
Predicted open state — TM-score=0.377
TM-score=0.405
QA=0.528

**Cluster3**
TM-score=0.891
Correct closed state — TM-score=0.929
QA=0.512

12

**Dimer**

**Minor difference**

**Experimental structures**

🟥 T1294v1

🟦 T1294v2

**TM-score=0.999 between two states**

**MolClust result for EnsembleFold decoys**

All decoys are almost identical, and fall into a single cluster

TM-score of decoy to T1294v1

TM-score of decoy to T1294v2

**Almost identical in the selected region**

**TM-score=0.988 to T1294v1**
**TM-score=0.989 to T1294v2**

**Predicting ensemble structures with minor variations remains highly challenging**

Future direction: dynamic selection of clustering regions and thresholds may be necessary to emphasize sampling of important regions in candidate clusters

# Summary

## What went right by EnsembleFold?

- **Diverse sets of MSAs** help create models with multiple states for ensemble targets
- **Knowledge-based REMC simulation** helps create diverse set of models
- **Structural clustering** works well for picking the correct model of each state in most cases

## What went wrong by EnsembleFold?

- Current confidence scores are not sensitive enough for selecting correct state model
- Predicting ensemble structures with minor variations remains highly challenging

# Team FZZH

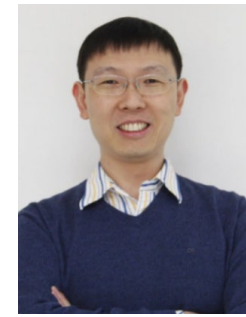## Freddolino & Zheng & Zhou & Hu Team



Dr. Lydia Freddolino
Umich

Dr. Wei Zheng
Umich -> NK

Dr. Xiaogen Zhou
ZJUT

Dr. Gang Hu
NK
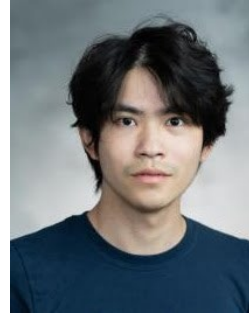
Dr. Pengshuo Yang
SDFNU
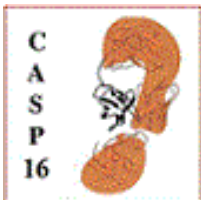
Dr. Qiqige Wuyun
MSU

Dr. Chunxiang Peng
Umich

Quancheng Liu
Umich

Wentao Ni
NK

Ziying Zhang
ZJUT

# Thank you!

# Q&A