# Accurate Contact/Distance Prediction by tFold
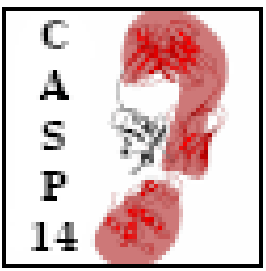
**tFold develop team:**

**Sheng Wang***#, Haidong Lan*, Tao Shen*, Jiaxiang Wu*,
Liangzhen Zheng*, Jianguo Pei*, Yuyi Liu, Junhong Huang,
Ningqiao Huang, Zhenlei Xu, Wei Liu#, and Junzhou Huang#

**CASP14 Conference**

2020.12.03

# Acknowledgement to CASP14 Organizers
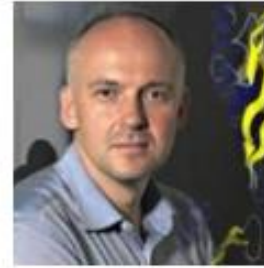
## CASP Organizers

**John MOULT**
President,
University of Maryland,
USA

**Krzysztof FIDELIS**
University of California,
Davis,
USA

**Andriy KRYSHTAFOVYCH**
University of California, Davis,
USA

**Torsten SCHWEDE**
University of Basel
SIB Swiss Institute of
Bioinformatics
Switzerland

**Maya TOPF**
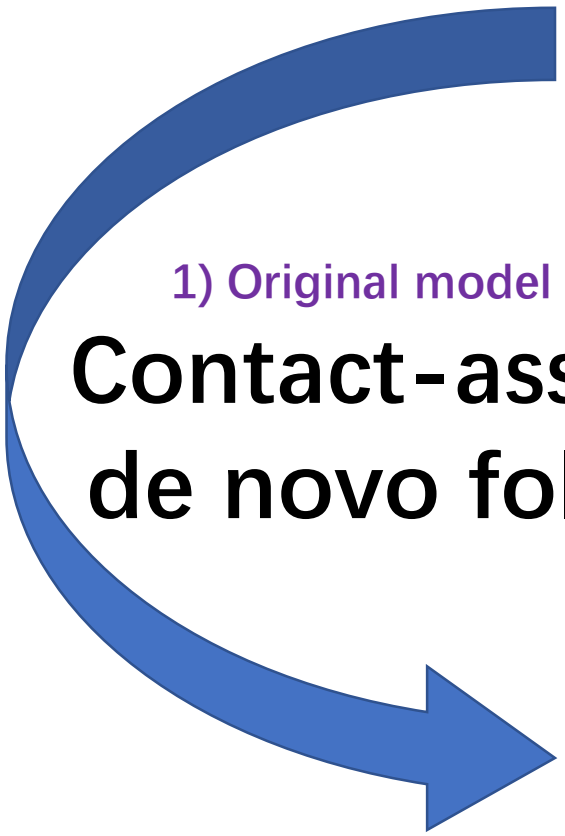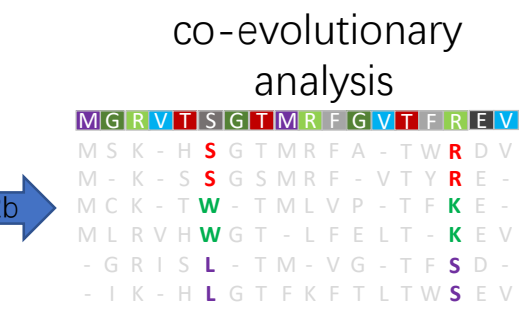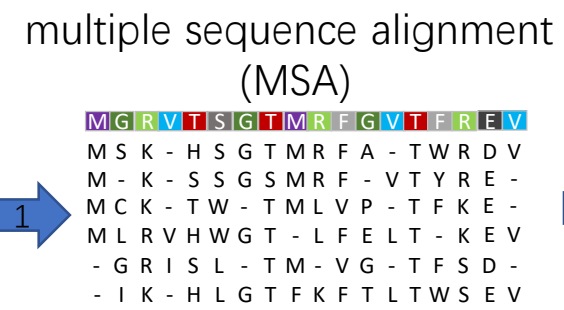Birkbeck, University of London, UK
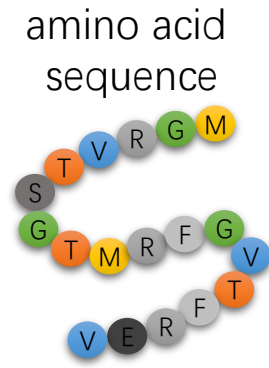CSSB (HPI and UKE) Hamburg,
Germany

## Contact-prediction Section Chair

Jinbo Xu

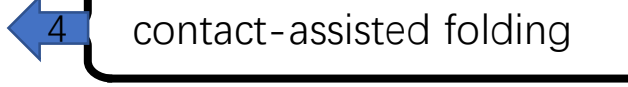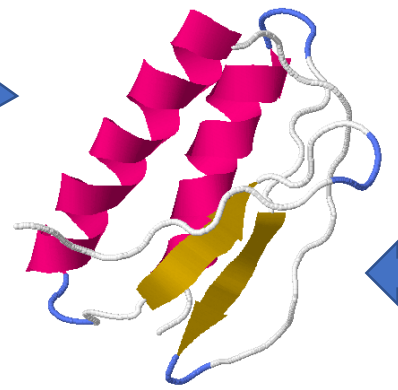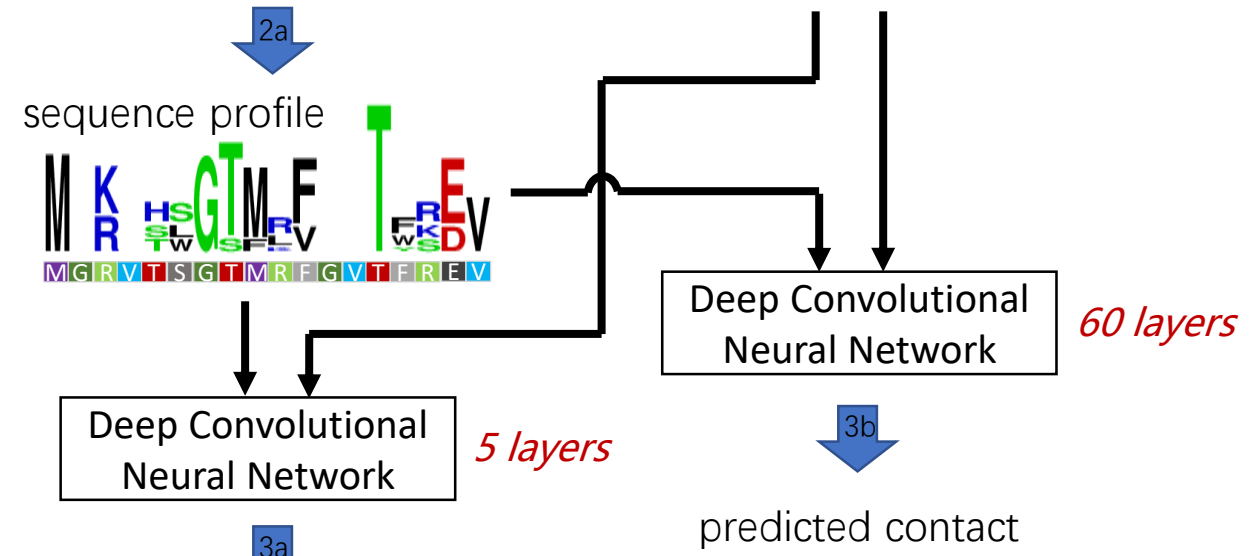**Affiliation:** Toyota Technological Institute at Chicago

amino acid sequence

multiple sequence alignment (MSA)

co-evolutionary analysis

MGRVTSGTMRFGVTFREV

MSK-HSGTMRFA-TWRDV
M-K-SSGSMRF-VTYRE-
MCK-TW-TMLVP-TFKE-
MLRVHWGT-LFELT-KEV
-GRISL-TM-VG-TFSD-
-IK-HLGTFKFTLTWSEV

1

2b

MGRVTSGTMRFGVTFREV

MSK-HSGTMRFA-TWRDV
M-K-SSGSMRF-VTYRE-
MCK-TW-TMLVP-TFKE-
MLRVHWGT-LFELT-KEV
-GRISL-TM-VG-TFSD-
-IK-HLGTFKFTLTWSEV

2a

sequence profile

MGRVTSGTMRFGVTFREV

**1) Original model**

# Contact-assisted de novo folding

Deep Convolutional Neural Network *60 layers*

Deep Convolutional Neural Network *5 layers*

3b

predicted contact

3a

predicted local structure

4 contact-assisted folding

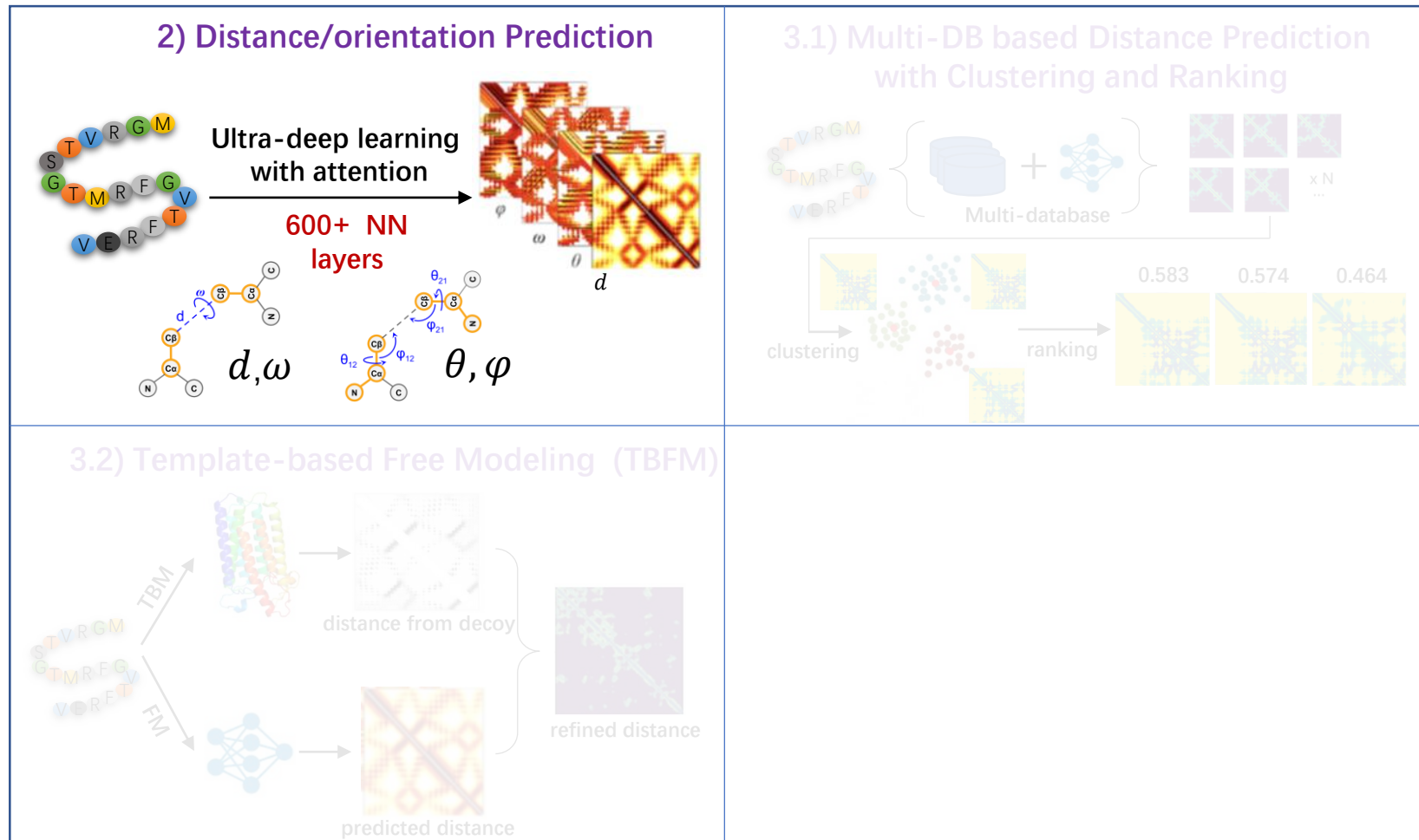**S Wang**, ···, J Xu[#]. *PLoS Computational Biology* **13**(1), (2017)

# Candidate issues of my previous work

- Contact V.S. distance/orientation

- Shallow network architecture

- Insufficient data usage
  a.  More input MSAs
  b.  More input decoys
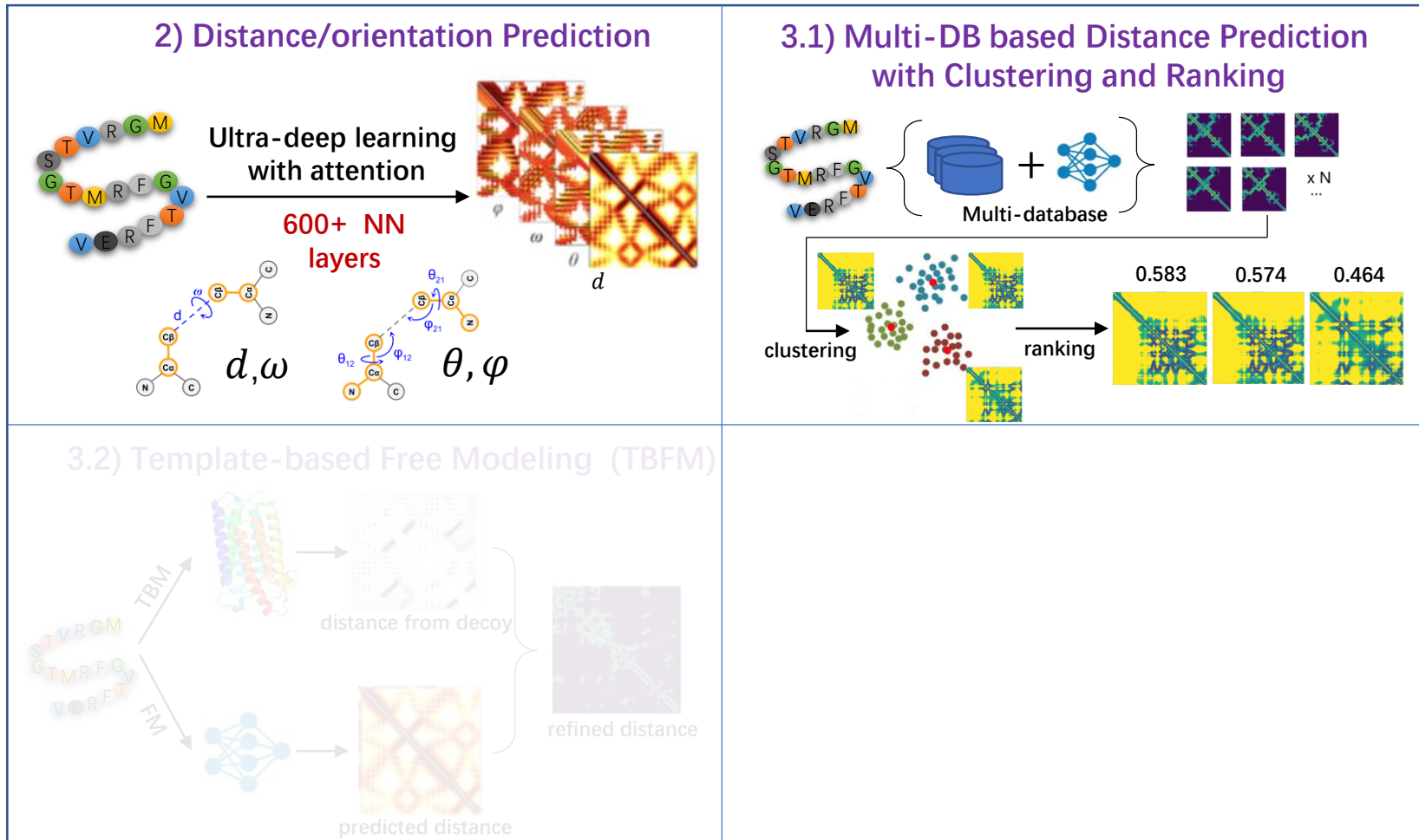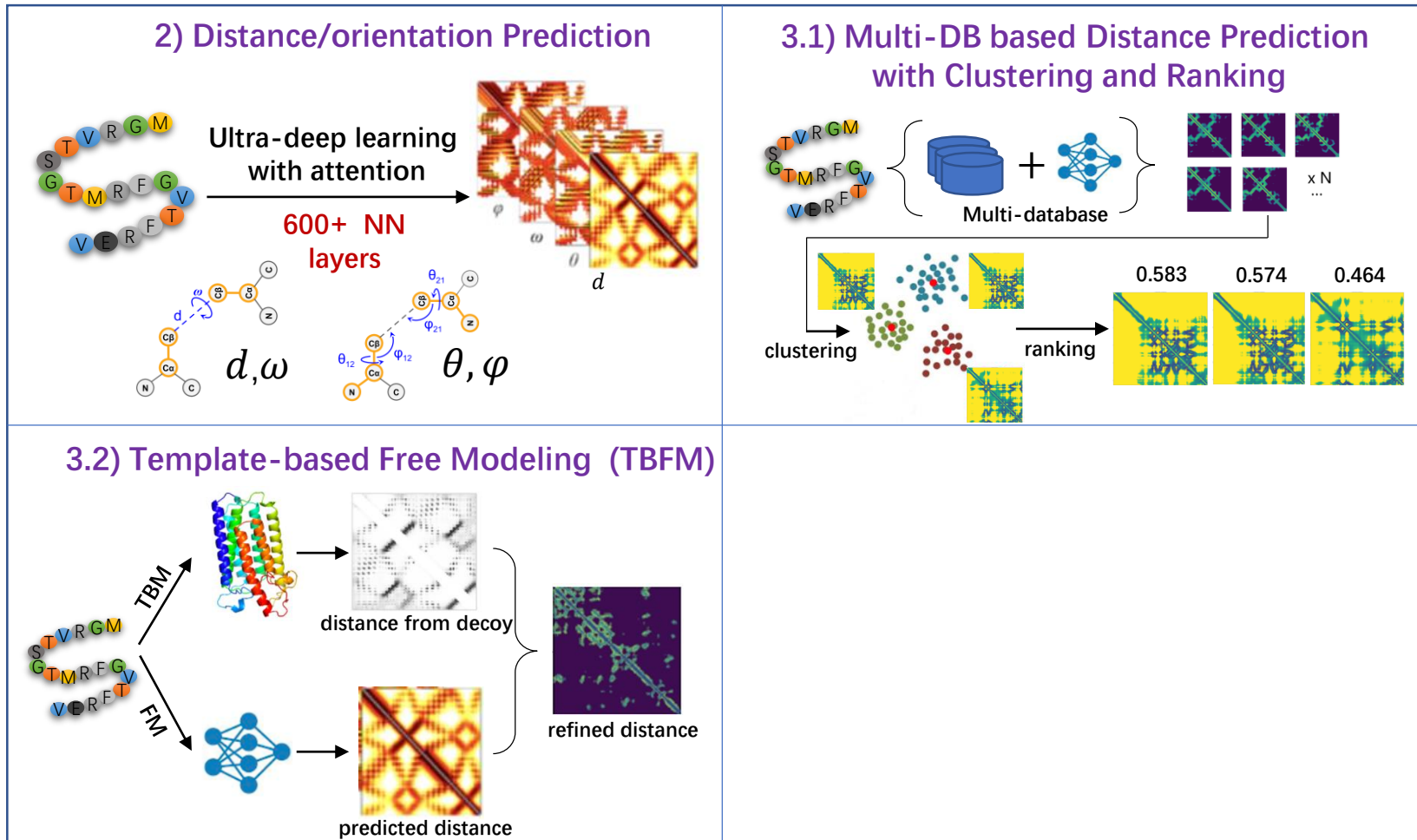
# New developments in tFold contact prediction
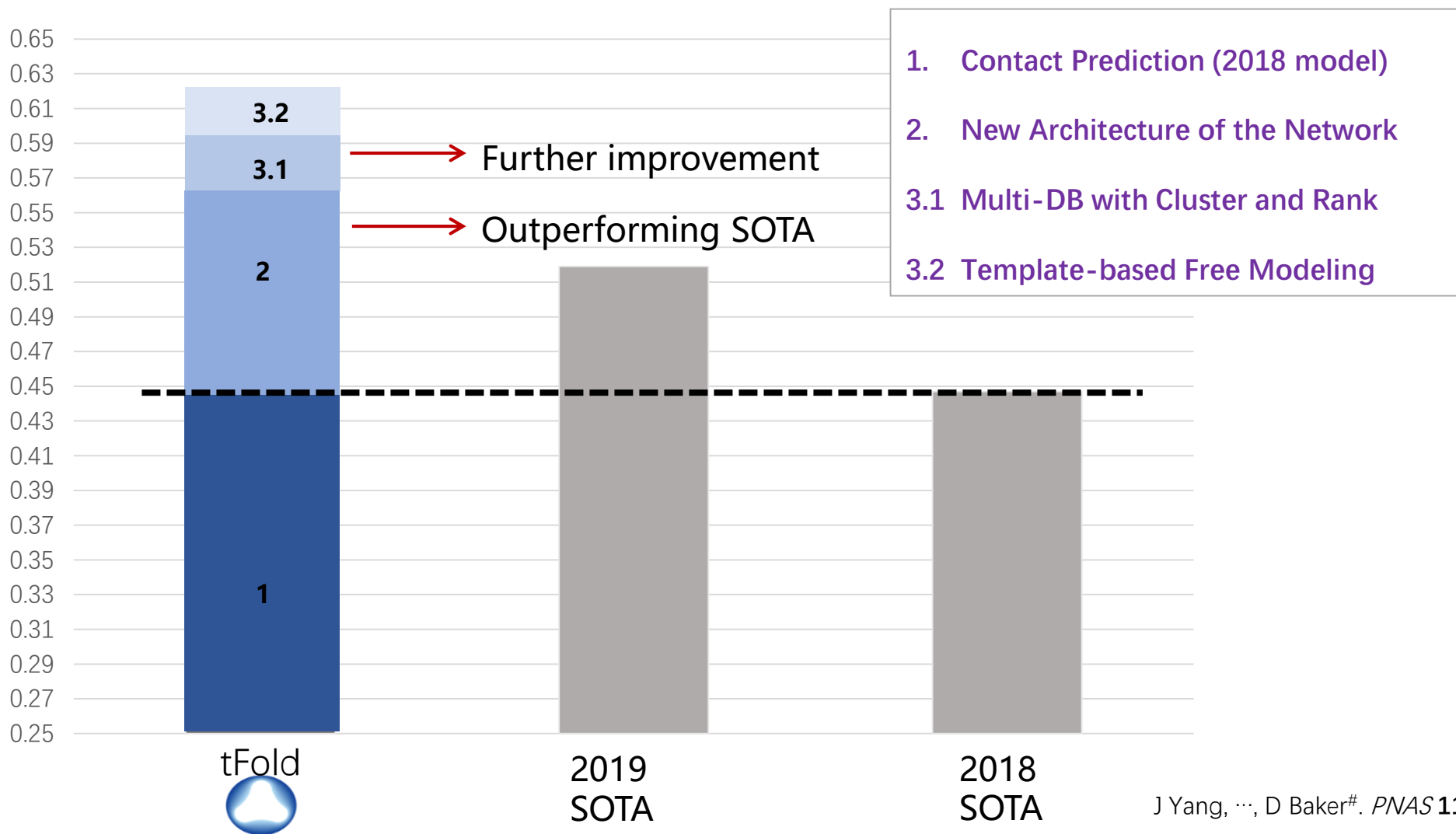


2) Distance/orientation Prediction

Ultra-deep learning with attention

600+ NN layers

$d, \omega$

$\theta, \varphi$

3.1) Multi-DB based Distance Prediction with Clustering and Ranking

Multi-database

x N

clustering          ranking

0.583      0.574      0.464

3.2) Template-based Free Modeling (TBFM)

TBM

distance from decoy

FM

predicted distance

refined distance

Distance -> J Zhu, **S Wang**,···, J Xu[#]. *Bioinformatics* **34**(13), (2018)     Orientation -> J Yang, ···, D Baker[#]. *PNAS* **117**(3), (2020)

# New developments in tFold contact prediction



2) Distance/orientation Prediction

Ultra-deep learning with attention

600+ NN layers

$d,\omega$     $\theta,\varphi$

3.1) Multi-DB based Distance Prediction with Clustering and Ranking

Multi-database

x N

clustering     ranking

0.583     0.574     0.464

3.2) Template-based Free Modeling (TBFM)

TBM

distance from decoy

FM

predicted distance

refined distance

# New developments in tFold contact prediction



### 2) Distance/orientation Prediction

Ultra-deep learning with attention

600+ NN layers

$d, \omega$      $\theta, \varphi$

### 3.1) Multi-DB based Distance Prediction with Clustering and Ranking

Multi-database

x N

clustering     ranking

0.583    0.574    0.464

### 3.2) Template-based Free Modeling (TBFM)

TBM

FM

distance from decoy

predicted distance

refined distance

[note]: the similar idea of TBFM first appears in J Xu[#], **S Wang**. *Proteins* **87**(12), (2019)

# TopL Long Range Contact Precision



*Data draw using CASP13 targets*

J Yang, ···, D Baker[#]. *PNAS* **117**(3), (2020)

© 2020 Tencent AI Lab

# tFold: Distance/Orientation Prediction with SEResNet+2DAttention with a Multi-input Multi-task Scheme



H Wang, et. al. *ECCV* (2020)

J Hu, et. al. *CVPR* (2018)

## Comparison of trRosetta and tFold

|  | trRosetta | tFold |
| --- | --- | --- |
| Architecture | ResNet | SEResNet+2DAttention |
| Size | 120 layers | 600+ layers |
| Time | 9 days (1 * RTX Titan) | 2 days (16 * V100) |

# Ablation study of the Deep Learning model

a) Data construction：
- Construct MSAs from multi databases
...

b) Network architecture and loss design：
- SE ResNet module
- 2D Attention module
- Multi-task learning
...

c) Training strategy:
- Progressive training strategy
- 600+ layer ultra-deep network
...

| Model | CASP13 TopL long range contact |
|---|---|
| Baseline model | 51.32%  ~ 2019 SOTA |
| Baseline + a) | 53.67% |
| Baseline + a) + b) | 55.15% |
| Baseline + a) + b) + c) | 56.37% |

# tFold-CaT: Multi-DB based Distance Prediction with Clustering and Ranking



Amino Acid Sequence

M G R V T S G T

A variety of MSAs

Ultra-deep CNN with 2D attention

*tFold*

1      2      3

hierarchical clustering

Quality Assessment upon Probability Distribution

Rank1      Rank2

Probability distribution      Probability distribution

Coverage ratio      0.68      Coverage ratio      0.48

Standard deviation      Standard deviation

TopK merge

# tFold-CaT: Multi-DB based Distance Prediction with Clustering and Ranking

# tFold-IDT: Template-based Free Modeling  (TBFM)



[note]: the similar idea first appears in J Xu[#], **S Wang**. *Proteins* **87**(12), (2019)

# tFold-IDT: Template-based Free Modeling (TBFM)

We can incorporate a variety of decoys

FM

Others

TBM

Template based modeling

Distance from initial model

Deviation prediction

Distance prediction model

Predicted distance

2D branch

Amino Acid Sequence

SPSSQGQHKHKYHFQK...

tFold

xN Blocks

Refined distance

Co-evolution features

1D branch

Sequential Features (One-hot, PSSM, ..)

Local LDDT score

0.1 0.3 0.4 0.5 0.3 0.1 ...

# The relationship between the decoy quality and the distance prediction enhancement

- High quality decoy

  significant enhance

- Low quality decoy

  won't influence much

The robustness of our algorithm with respect of the decoy quality



*Data draw using CAMEO targets from 2020-02-01 to 2020-05-02*

# What goes right and why?

# What goes right and why?

# What goes wrong and why?

We didn't use…

☹

Modern metagenomics databases



| ID | TripleRes | tFold-orig | tFold-BFD | diff |
|---|---|---|---|---|
| T1027-D1 | 0.4343 | 0.3939 | 0.5152 | 0.1213 |
| T1029-D1 | 0.0560 | 0.0400 | 0.0320 | (0.0080) |
| T1031-D1 | 0.3789 | 0.0632 | 0.2526 | 0.1894 |
| T1033-D1 | 0.1200 | 0.1400 | 0.1600 | 0.0200 |
| T1037-D1 | 0.5767 | 0.4455 | 0.5347 | 0.0892 |
| T1038-D1 | 0.2719 | 0.3070 | 0.3333 | 0.0263 |
| T1039-D1 | 0.4596 | 0.0994 | 0.6149 | 0.5155 |
| T1040-D1 | 0.4000 | 0.0769 | 0.2000 | 0.1231 |
| T1041-D1 | 0.6901 | 0.6322 | 0.7107 | 0.0785 |
| T1042-D1 | 0.4891 | 0.3225 | 0.5290 | 0.2065 |
| T1043-D1 | 0.0473 | 0.2568 | 0.2568 | 0.0000 |
| T1047s1-D1 | 0.4834 | 0.6256 | 0.6445 | 0.0189 |
| T1049-D1 | 0.7388 | 0.8284 | 0.8284 | 0.0000 |
| T1061-D2 | 0.7232 | 0.5867 | 0.7048 | 0.1181 |
| T1064-D1 | 0.0652 | 0.0326 | 0.0543 | 0.0217 |
| T1074-D1 | 0.3864 | 0.2879 | 0.5758 | 0.2879 |
| T1090-D1 | 0.6138 | 0.7831 | 0.7696 | (0.0135) |
| T1093-D1 | 0.0284 | 0.1631 | 0.2553 | 0.0922 |
| T1093-D3 | 0.0566 | 0.6321 | 0.5566 | (0.0755) |
| T1094-D2 | 0.6473 | 0.4541 | 0.5411 | 0.0870 |
| T1096-D1 | 0.5569 | 0.3412 | 0.3765 | 0.0353 |
| T1096-D2 | 0.4971 | 0.2164 | 0.4035 | 0.1871 |
| Average | 0. 3964 | 0.3513 | **0.4477** | 0.0964 |

| # | Gr.# | Gr. Name | No. Domains | | F1 | | Prec | |
|---|---|---|---|---|---|---|---|---|
| | | | No Submit. | No Total | Submit. | Total | Submit. | Total |
| 1. | 368 | tFold-CaT_human | 22 | 22 | 41.158 | 41.158 | 41.783 | 41.783 |
| 2. | 488 | tFold-IDT_human | 22 | 22 | 39.374 | 39.374 | 40.504 | 40.504 |
| 3. | 010 | TripletRes | 22 | 22 | 39.282 | 39.282 | 39.641 | 39.641 |
| 4. | 125 | PreferredFold | 22 | 22 | 38.696 | 38.696 | 39.440 | 39.440 |
| 5. | 024 | DeepPotential | 22 | 22 | 38.286 | 38.286 | 38.586 | 38.586 |
| 6. | 009 | tFold_human | 22 | 22 | 36.821 | 36.821 | 38.056 | 38.056 |
| 7. | 183 | tFold-CaT | 22 | 22 | 35.465 | 35.465 | 37.107 | 37.107 |
| 8. | 351 | tFold-IDT | 22 | 22 | 34.774 | 34.774 | 36.516 | 36.516 |
| 9. | 238 | tFold | 22 | 22 | 33.548 | 33.548 | 35.130 | 35.130 |

tFold-orig only uses metaclust50 (year 2018) as the metagenomics databases.
tFold-BFD adds BFD (year 2019) as the additional metagenomics databases.

# What goes wrong and why?



T1039 (orig)

T1039 (BFD)

T1042 (orig)

T1042 (BFD)

| ID | tFold-orig | tFold-BFD | diff |
|---|---|---|---|
| T1027-D1 | 0.3939 | 0.5152 | 0.1213 |
| T1029-D1 | 0.0400 | 0.0320 | (0.0080) |
| T1031-D1 | 0.0632 | 0.2526 | 0.1894 |
| T1033-D1 | 0.1400 | 0.1600 | 0.0200 |
| T1037-D1 | 0.4455 | 0.5347 | 0.0892 |
| T1038-D1 | 0.3070 | 0.3333 | 0.0263 |
| T1039-D1 | 0.0994 | 0.6149 | 0.5155 |
| T1040-D1 | 0.0769 | 0.2000 | 0.1231 |
| T1041-D1 | 0.6322 | 0.7107 | 0.0785 |
| T1042-D1 | 0.3225 | 0.5290 | 0.2065 |
| T1043-D1 | 0.2568 | 0.2568 | 0.0000 |
| T1047s1-D1 | 0.6256 | 0.6445 | 0.0189 |
| T1049-D1 | 0.8284 | 0.8284 | 0.0000 |
| T1061-D2 | 0.5867 | 0.7048 | 0.1181 |
| T1064-D1 | 0.0326 | 0.0543 | 0.0217 |
| T1074-D1 | 0.2879 | 0.5758 | 0.2879 |
| T1090-D1 | 0.7831 | 0.7696 | (0.0135) |
| T1093-D1 | 0.1631 | 0.2553 | 0.0922 |
| T1093-D3 | 0.6321 | 0.5566 | (0.0755) |
| T1094-D2 | 0.4541 | 0.5411 | 0.0870 |
| T1096-D1 | 0.3412 | 0.3765 | 0.0353 |
| T1096-D2 | 0.2164 | 0.4035 | 0.1871 |
| | | | |
| Average | 0.3513 | 0.4477 | 0.0964 |

tFold-orig only uses metaclust50 (year 2018) as the metagenomics databases.
tFold-BFD add BFD (year 2019) as the additional metagenomics databases.

# **Take home messages**

- Distance/orientation matters.

- Deeper and attention-based network works.

- Sufficient data usage will increase the robustness:
  a. More input MSAs
  b. More input decoys

# Accurate De Novo Protein Structure Prediction by tFold Server

https://drug.ai.tencent.com/console/en/tfold

# From MRF to Distance



# From Distance to 3D model

# Excellent performance of tFold Server on CAMEO



**Performance on hard targets**

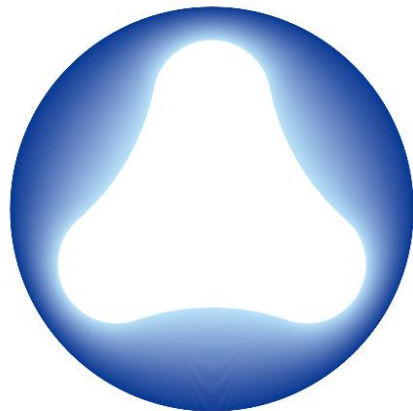Junhong Huang  Tao Shen  Junzhou Huang  Wei Liu  Jiaxiang Wu  Jianguo Pei

Ningqiao Huang  Haidong Lan  Liangzhen Zheng  Yuyi Liu

Zhenlei Xu  Sheng Wang

# Thank you

Tencent AI Lab