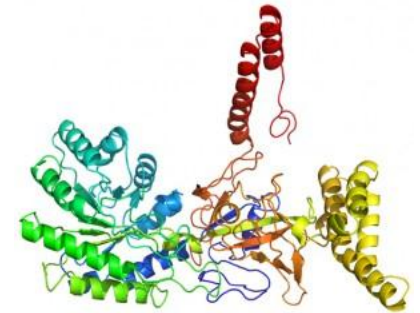# Recursive Protein Structure Modeling

**Jianlin Cheng, PhD**

**Department of Computer Science**

**Informatics Institute**

**University of Missouri, Columbia**

Presented at BIBM Computational Structural Bioinformatics Workshop, Nov. 12, 2011

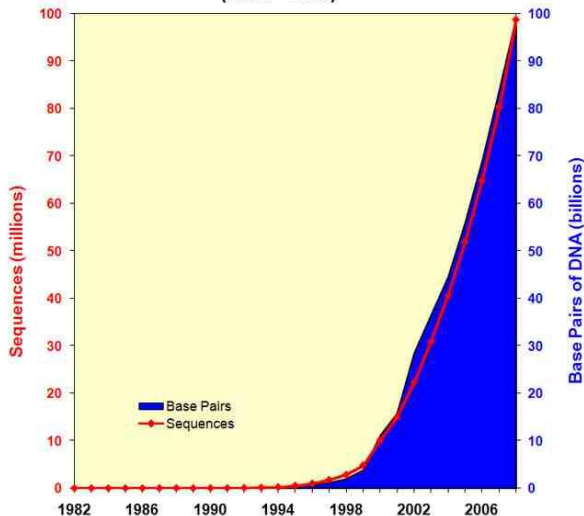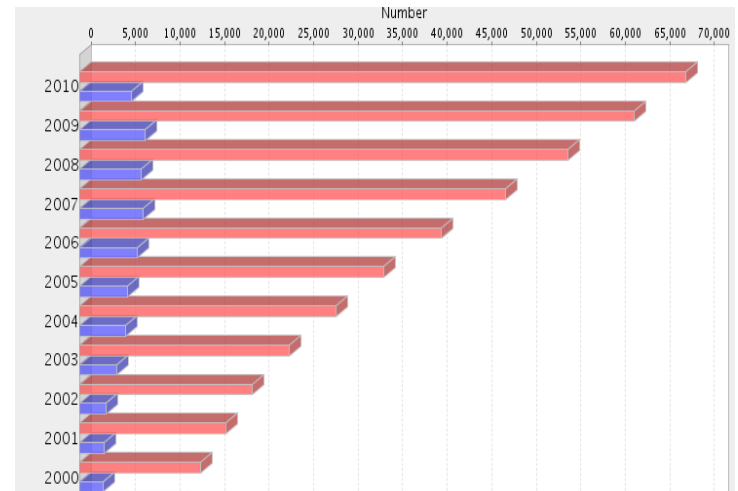# Protein Structure Prediction – A Key Challenge in the Genomic Era

**Genome Sequencing**



Growth of GenBank (1982 - 2008)



**Genome Interpretation**
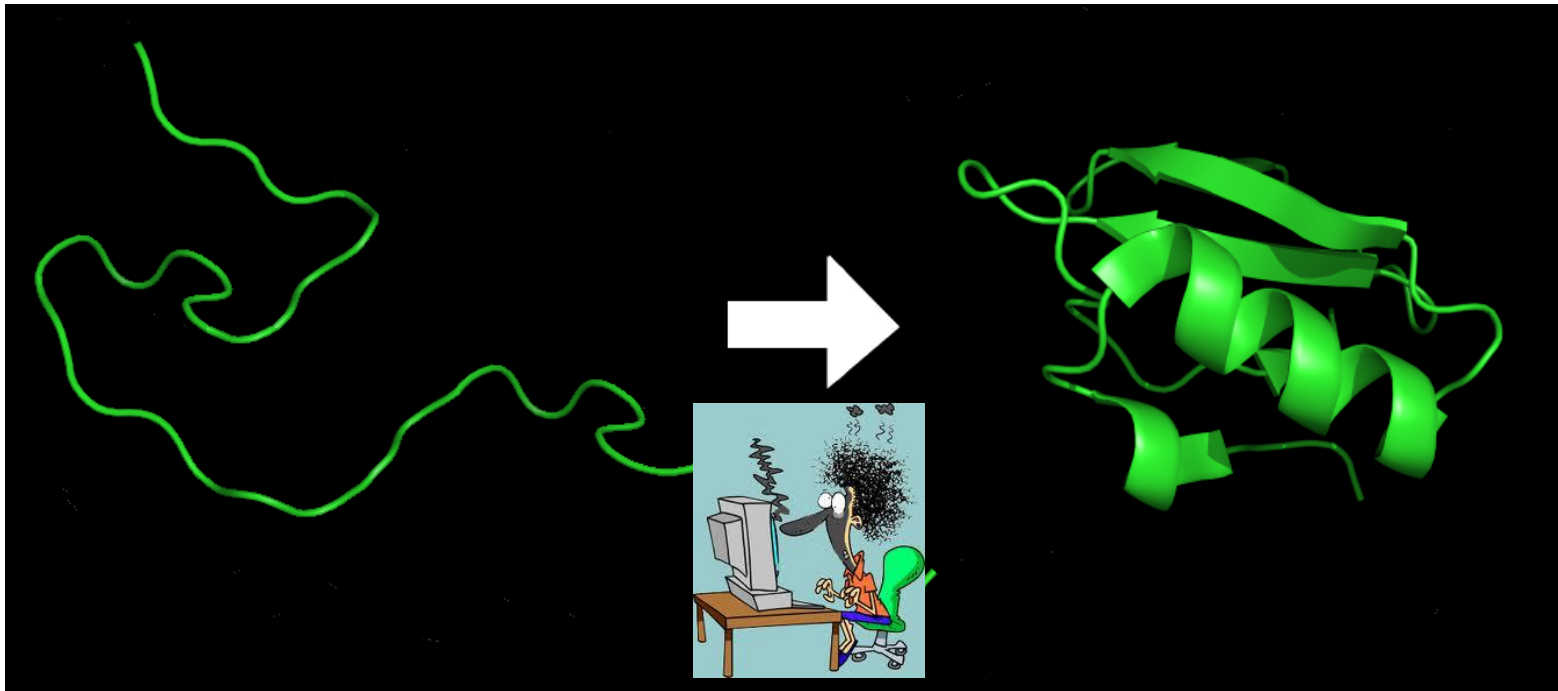


**Growth of PDB Structures**



Images.google.com

# Computational Protein Structure Prediction

**Structure = $f$ ( sequence) ?**    ⬌    **E = MC$^2$**



**Computational Simulation**

# Template-Based Modeling



**Sequence Space**

**Structure Space**

**Target protein**

MWLKKFGINKH...

**Protein Data Bank**

**Fold Recognition**

**Template**

**Alignment**

# Template-Based Modeling



**TARGET**

**TEMPLATE**

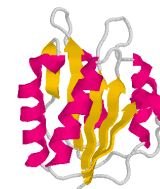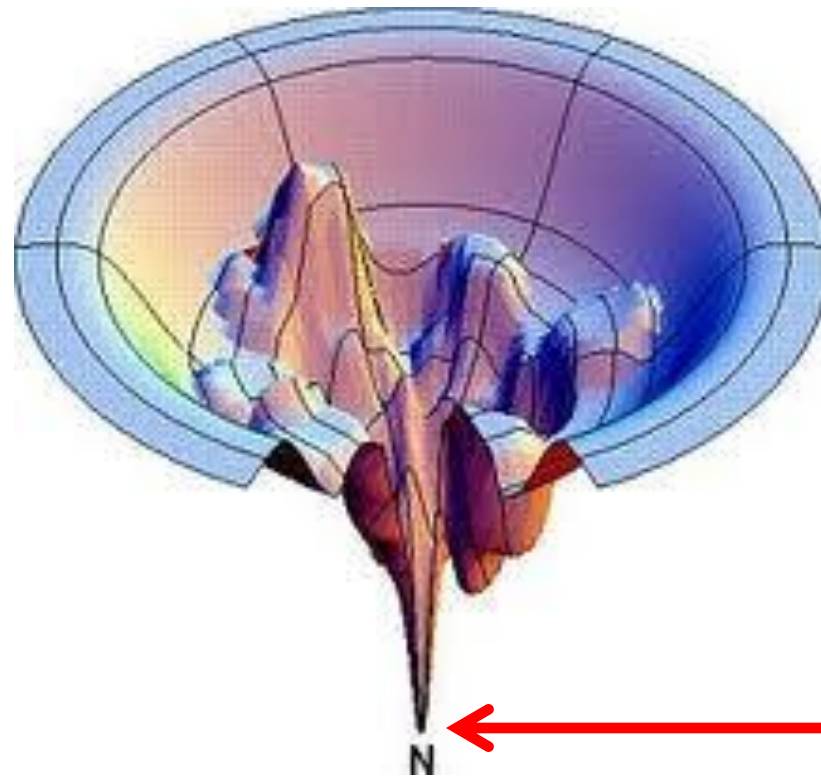ASILPKRLFGNCEQTSDEGLK
IERTPLVPHISAQNVCLKIDD
VPERLIPERASFQWMNDK

ASILPKRLFGNCEQTSDEGLKIERTPLVPHISAQNVCLKIDDVPERLIPE
MSVIPKRLYGNCEQTSEEAIRIEDSPIV---TADLVCLKIDEIPERLVGE

Modeller
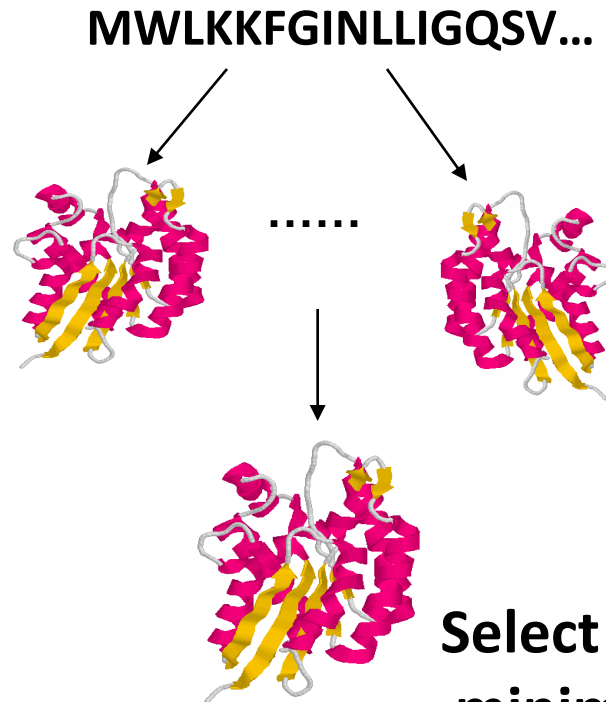
A. Fisher, 2005

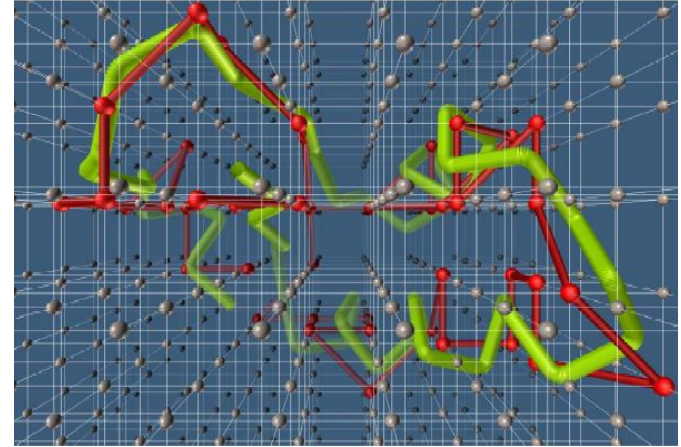# Template-Free (*Ab Initio*) Modeling



**Protein Structure Space**

Native Structure

Dill & Chan, 1997

# Template-Free Modeling

MWLKKFGINLLLIGQSV...

3D Simulation



......

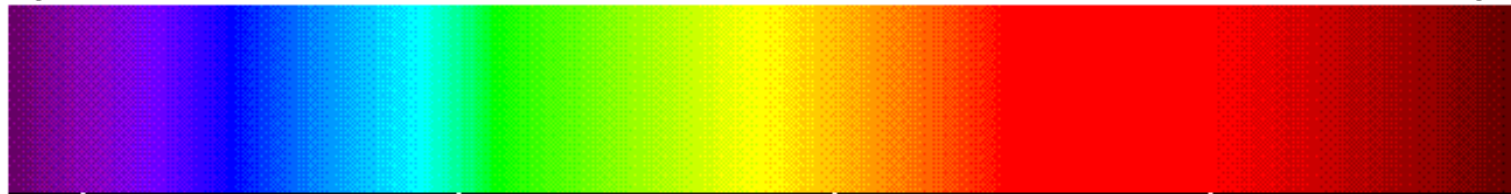Select model with
minimum free energy

**Methods: molecular dynamics, fragment assembly, distance / contact-based modeling**

# Combination of Template-Free and Template-Based Modeling
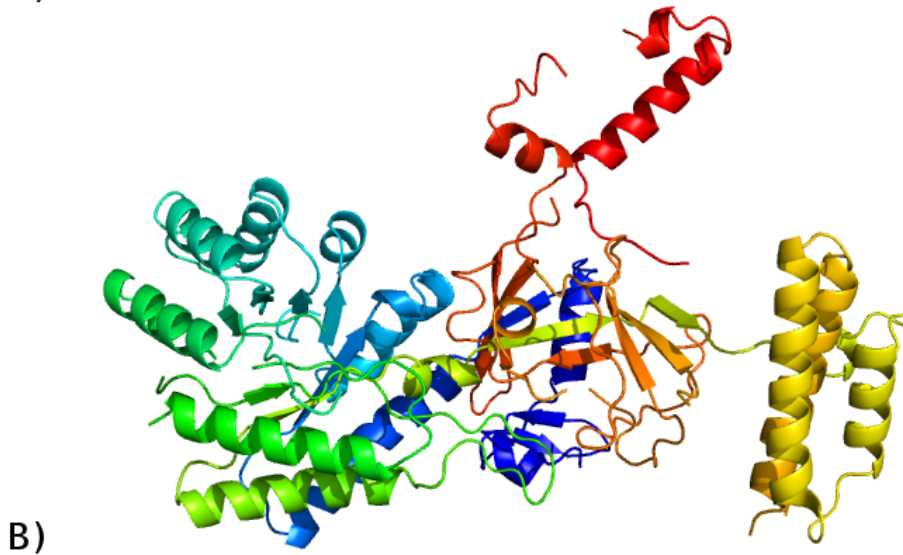
**100% TBM**  **50% TBM+50%FM**  **100% FM**



## Protein Modeling Spectrum

# Region Decomposition from Alignment

# Recursive Protein Modeling – Integrate TBM and FM

# Recursive Modeling Mimics Protein Folding Cascade



Start  100 ns  250 ns  500 ns

1000 ns  2000 ns  2500 ns  Folded

ks.uiuc.edu

# Case 1: Domain-Level Recursive Protein Modeling – CASP9 T0547

| Domain 1 | Domain 2 | Domain 1 | Domain 3 | Domain 1 | Domain 4 |
|----------|----------|----------|----------|----------|----------|

Template-Based

Template-Free

Template-Based

Template-Free

TBM Domain 2
(Insertion)
GDT-TS = 0.74

TBM Domain 1 – Three Discontinuous Fragments
GDT-TS = 0.66

# Case 2: Refine uncertain regions of a largely template-based modeling (T0539)

# Core-Constrained Tail Refinement

**Template protein**

```
>P1;1VI8B
structureX:1VI8B: 1: : 146: : : : :
------------------------SLIWKRKITLEALNAMGEGNMVGFLDIRFEHIGDDTLEATMPVDSRTKQPFGLLHGGASVVL
AESIGSVAGYLCTEGEQKVVGLEINAMHVRSAREGRVRGVCKPLHLGSRHQVWQIEIFDEKGRLCCSSRLTTAILEGGSHHHHH*
```
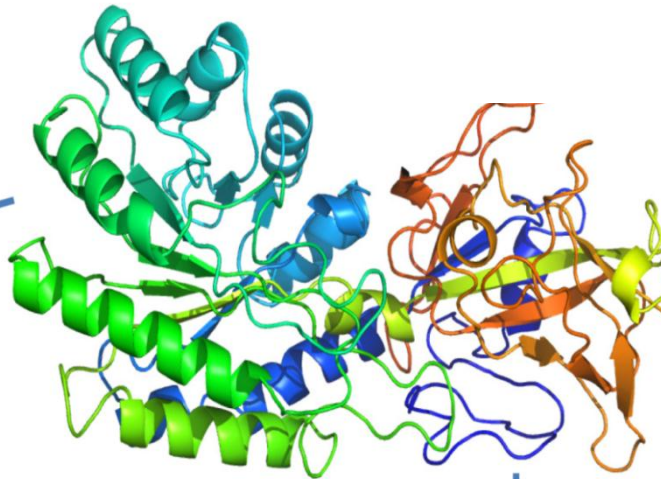
**Query protein**

```
>P1;Query
 : : : : : : : : :
MDKRLQQDRIVDKMERFLSTANEEEKDVLSSIVDGLLAKQERRYATYLASLTQIESQEREDGRFEVRLPIGPLVNNPLNMVHGGITATL
LDTAMGQMVNRQLPDGQSAVTSELNIHYVKPGMGTYLRAVASIVHQGKQRIVVEGKVYTDQGETVAMGTGSFFVLRSRG-----*
```

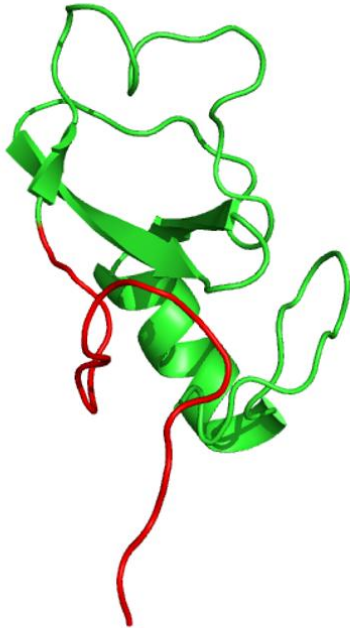# Core-Constrained Tail Refinement



A)

**Before tail refinement**
**GDT-TS = 0.64**

# Core-Constrained Tail Refinement



A)

Before tail refinement
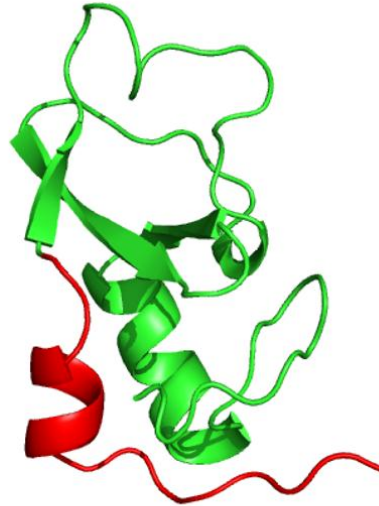GDT-TS = 0.64

B)

After tail refinement
GDT-TS = 0.73

# Core-Constrained Tail Refinement



A)
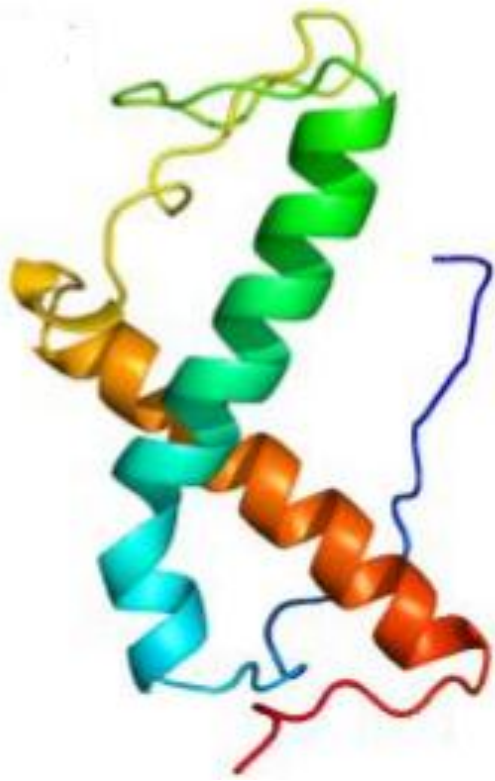
**Before tail refinement**
**GDT-TS = 0.64**

B)

**After tail refinement**
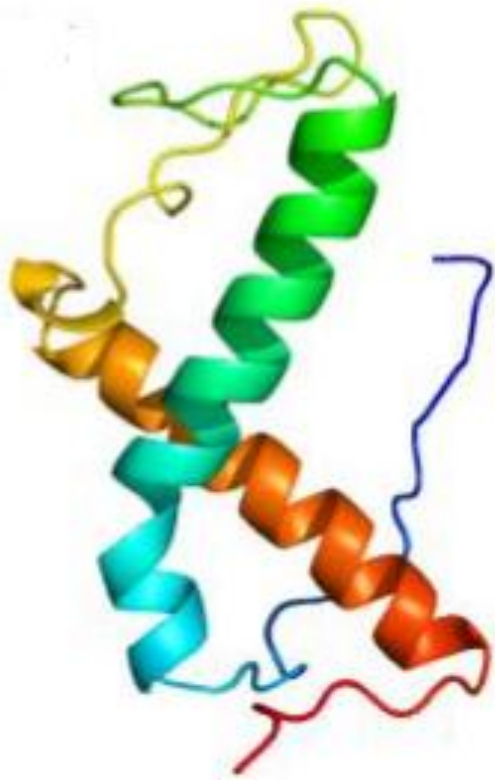**GDT-TS = 0.73**

C)

**Superposition**
**Green: model, Blue: structure**

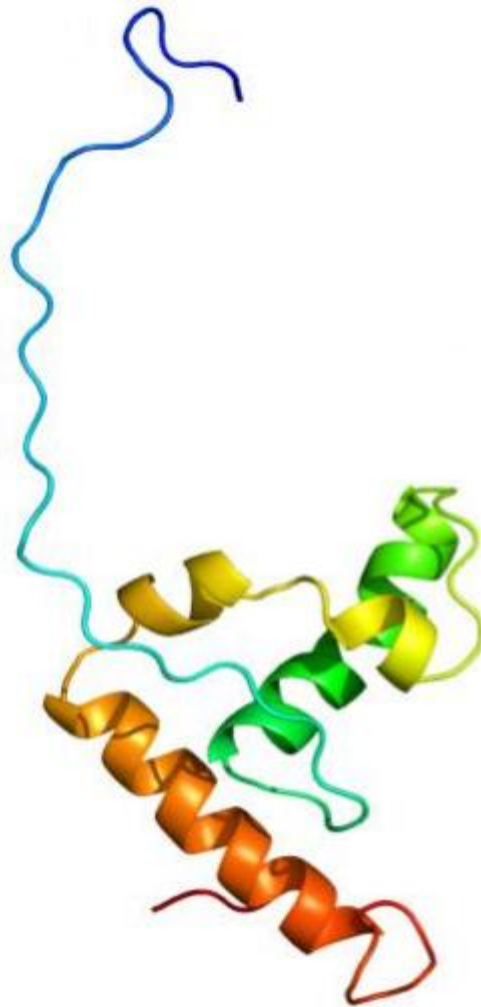# Case 3: Expanding a template-based core into a full structure (T0616)
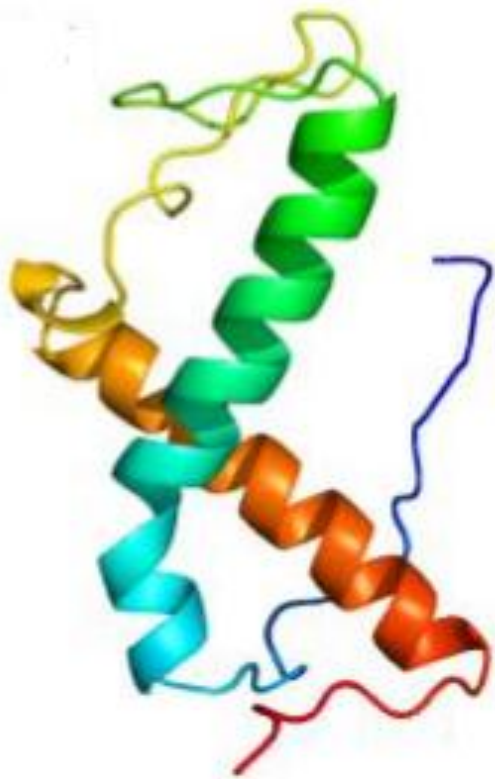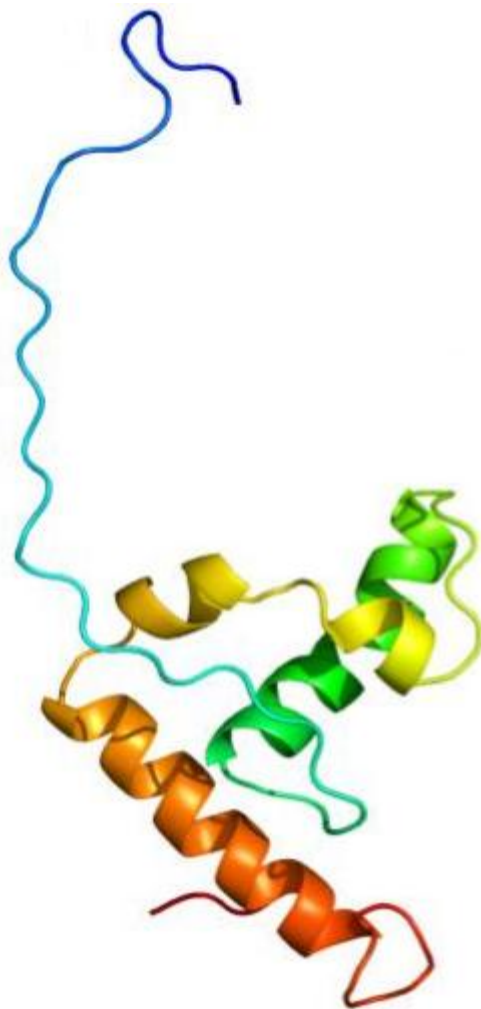
**Native Structure**

**Native Structure**

**Template-based modeling
(GDT-TS = 0.34)**

**Native Structure**

**Template-based modeling (GDT-TS = 0.34)**

**Template-based + Ab Initio (GDT-TS = 0.39)**

# Advantages of Recursive Protein Modeling

- **Avoiding error-prone hard decisions on the classification of a protein target or a region**

- **Combining the strength of template-based modeling and template-free modeling**

- **Improving sampling efficiency by recursively expanding certain regions**
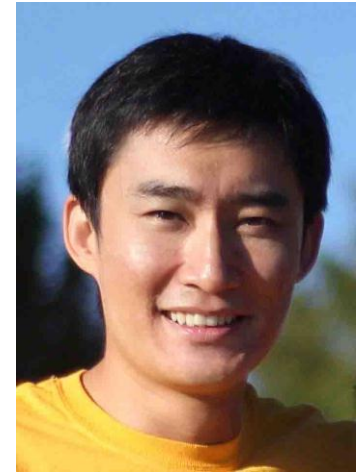
- **Easy to implement and improve**

# Acknowledgements



Xin Deng

Jesse Eickholt

Zheng Wang

# Comparison with Previous Approaches

- **Compared with loop modeling**

  Region of any size: loop, partial domain,
  domain, multiple domains

  Region of any type: helix, strand, loop

- **Compared with  TASSER**

  Common: template-based + template-free

  Different: gap filling VS. alternated,
  recursive certainty expansion