# Toward reliability evaluation of computational models of protein molecules and their interactions

**Md Hossain Shuvo, Ph.D.**
Assistant Professor
Department of Computer Science
Prairie View A&M University

**CCSB @PVAMU**

February 19th,2025

# Proteins

A fundamental molecule and workhorse of cells
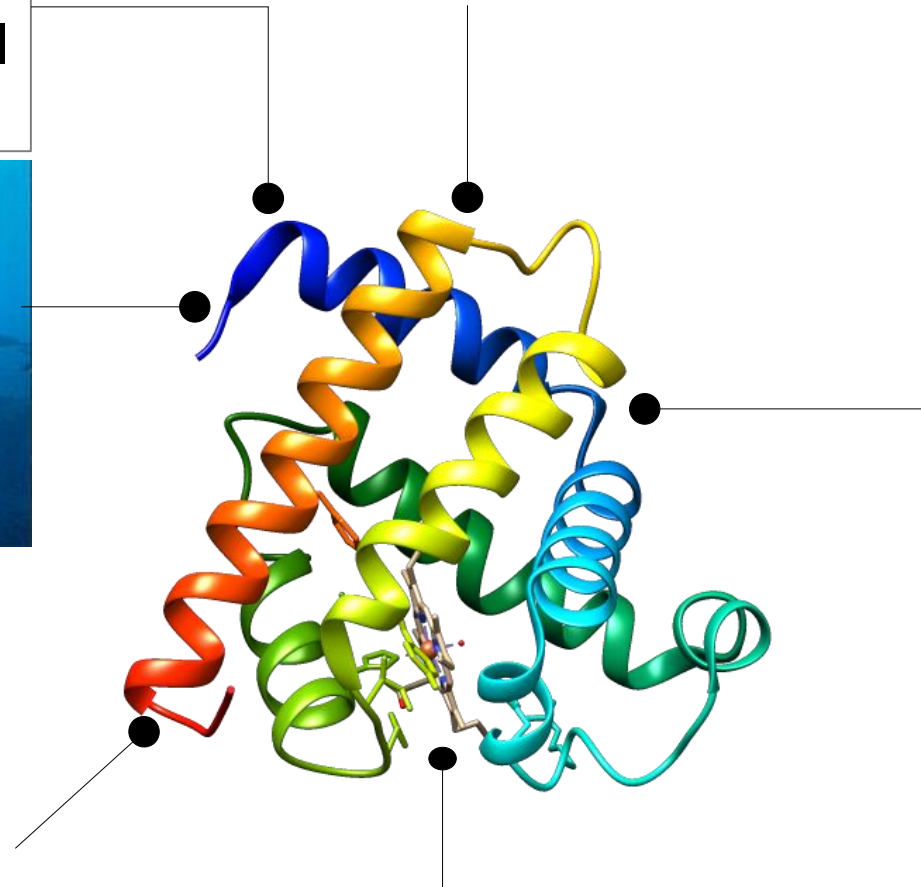
Participates in most biological processes
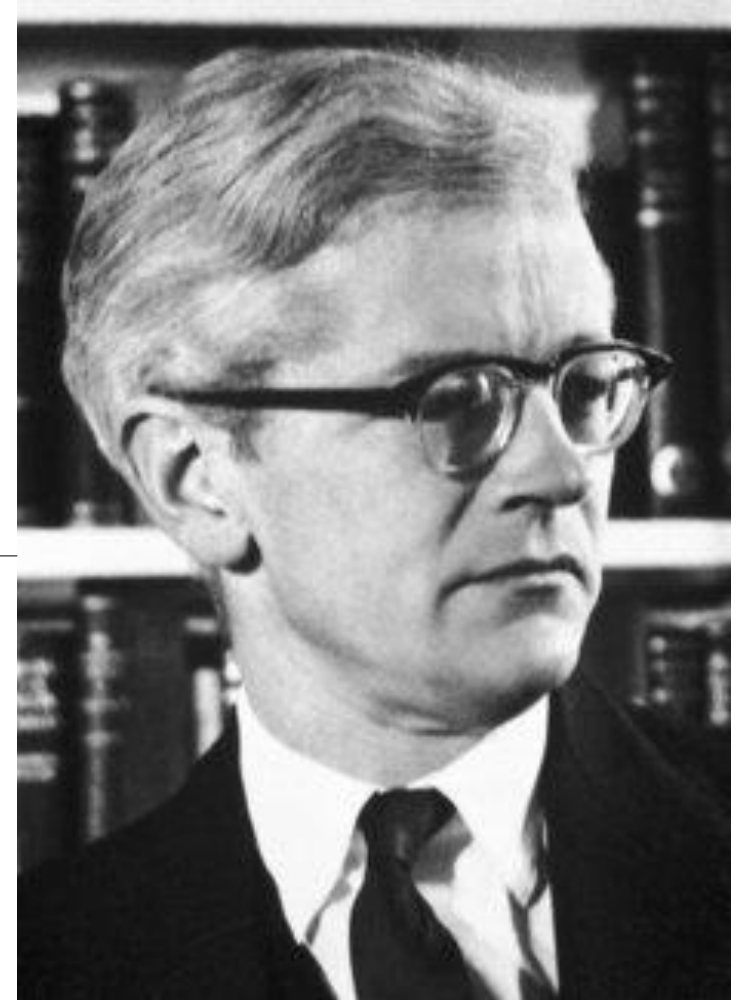

Diving mammals


Skeletal muscle

Myoglobin
PDB ID: 1MBN



Oxygen storage

Dr. John Cowdery Kendrew
Solved the structure in 1957
Noble prize in Chemistry 1962

# How to obtain protein 3D structure?

## Experimental approaches

➢ X-ray crystallography

➢ Nucleic Magnetic Resonance
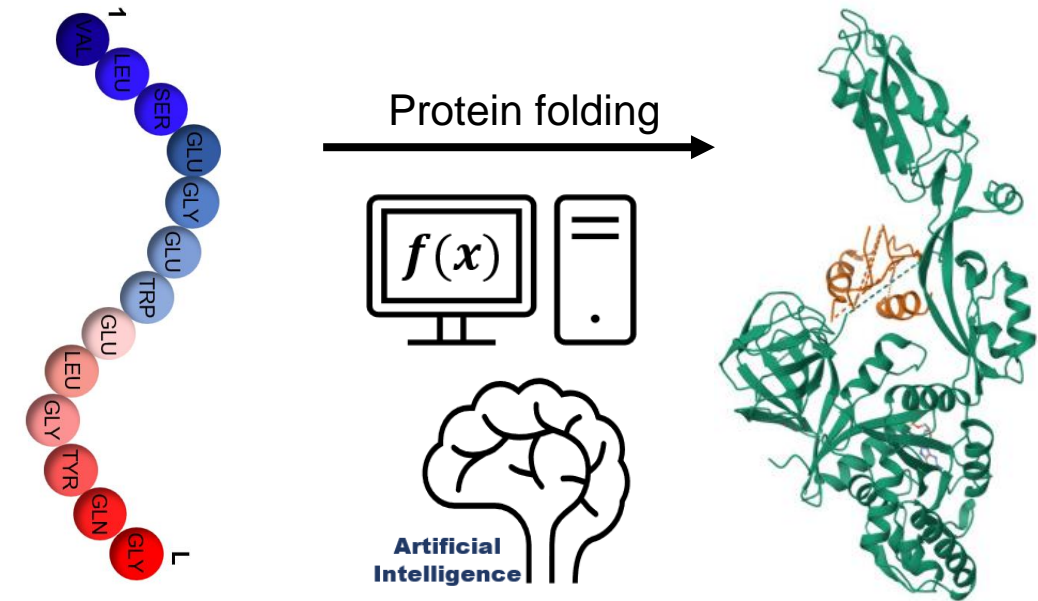
➢ Cryo-Electron Microscopy

**Drawbacks:**

➢ Expensive

➢ Extensive

➢ Leads to gap between sequence and structure



Dr. Christian Anfinsen
Noble prize in 1972

"demonstrated that the amino acid sequence of a protein contained all of the information needed for the protein to reach the native conformation"
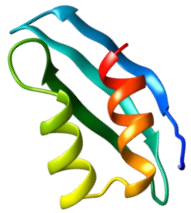


Protein folding

Artificial Intelligence

Computational protein structure prediction can help

Computationally predicted protein models may have error..
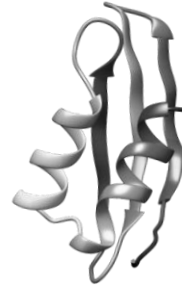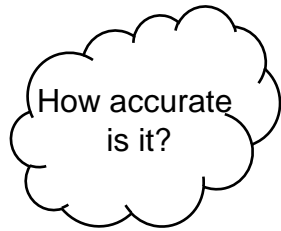
# Protein model quality estimation

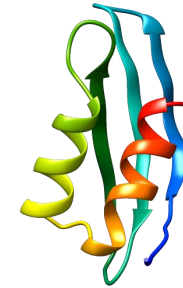# Protein model quality estimation

## Model quality estimation
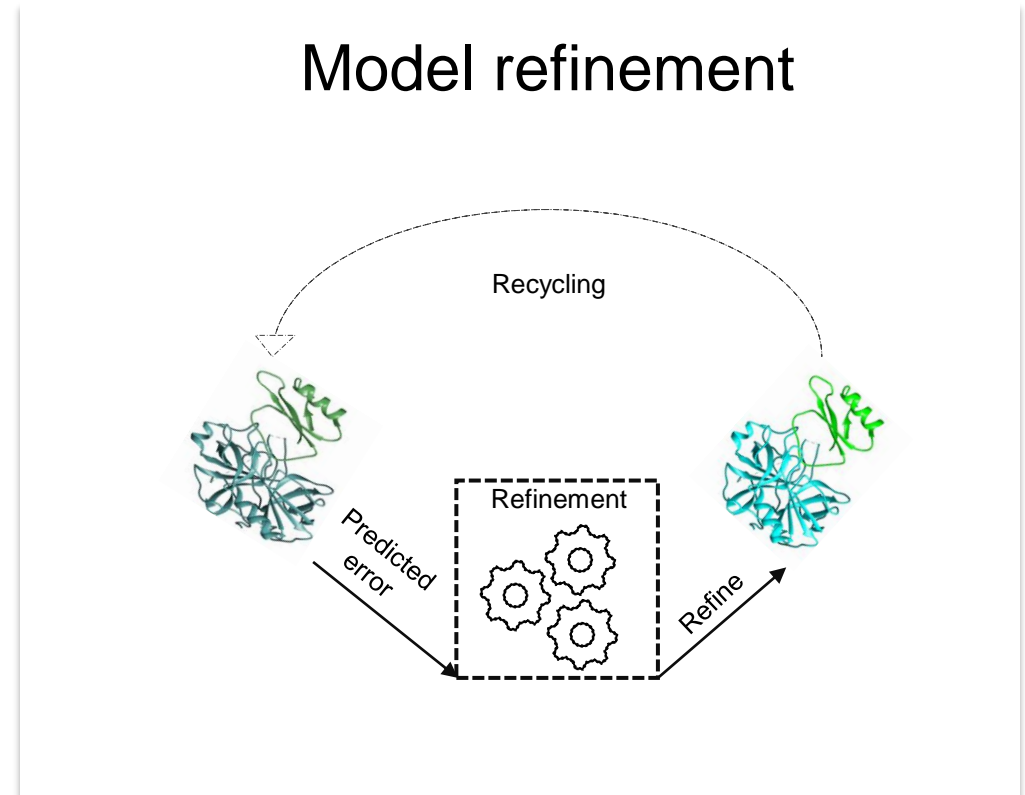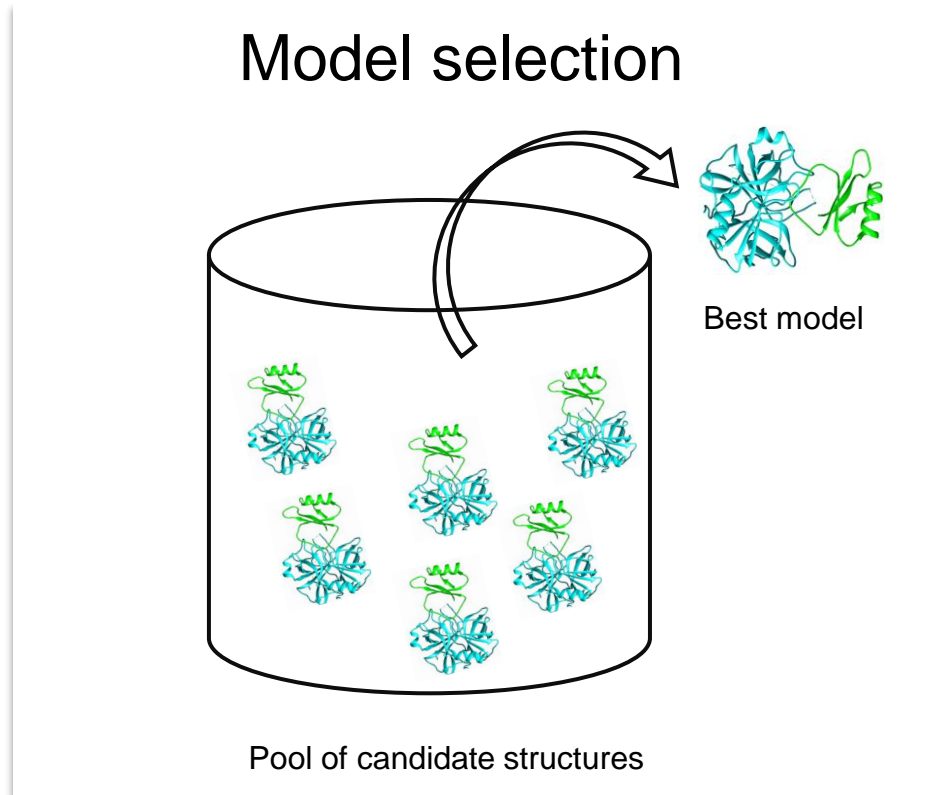


How accurate is it?

Computationally predicted protein model

Model

Experimental/Ground truth

Estimation of protein model quality

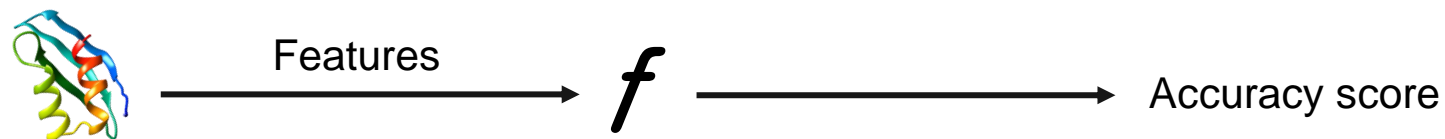➢ When the experimental structure is not known/absent

Model selection

Best model

Pool of candidate structures

Model refinement

Recycling

Predicted error

Refinement

Refine

Helps in accurately guide the process of protein structure prediction

# Application of deep learning in model quality estimation



Features $\xrightarrow{\hspace{2cm}}$ $f$ $\xrightarrow{\hspace{2cm}}$ Accuracy score

Deep learning in protein model quality estimation

- SVM
- MLP
- NN
- RF

Residual Neural Network

Graph Attention Network

Symmetry-aware Graph Neural Network

Protein Language Models

- ProQ3
- ProQ3D
- SVMQA
- RFMQA

- QDeep
- iQDeep
- DeepRefiner

PIQLE

EquiRank

Our methods for protein model quality estimation

Uziela *et al.*, 2016    Manavalan and Lee, 2017    He *et al.*, 2016)    Lin *et al.*, 2023
Uziela *et al.*, 2017    Manavalan *et al.*, 2014    Satorras *et al.*, 2022    **Shuvo *et al.*, 2020, 2021, 2023, 2025**

# Fundamental research questions

**1**

How to estimate the quality of predicted protein models?

1. **CSBJ**, 2025
2. **ISMB** 2020
   also, in **Bioinformatics**
   Oxford Press 2020
3. **JMB** 2023
4. **Bioinformatics advances** 2023
5. **PLOS ONE** 2020
6. **Proteins** 2021

**2**

How to apply quality estimation method to improve quality of predicted protein models?

Improving the quality of less accurate protein models

7. **Nucleic Acids Research** 2021

**JMB**: Journal of Molecular Biology
**ISMB**: International Society for Computational Biology
**CSBJ:** Computational and Structural Biotechnology Journal
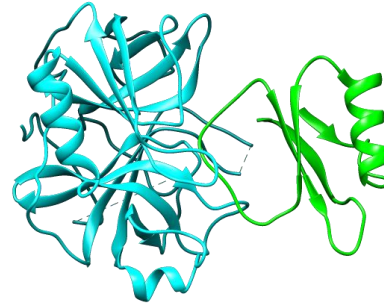
# Key research objectives

## 1

Estimation of **monomeric** protein model quality
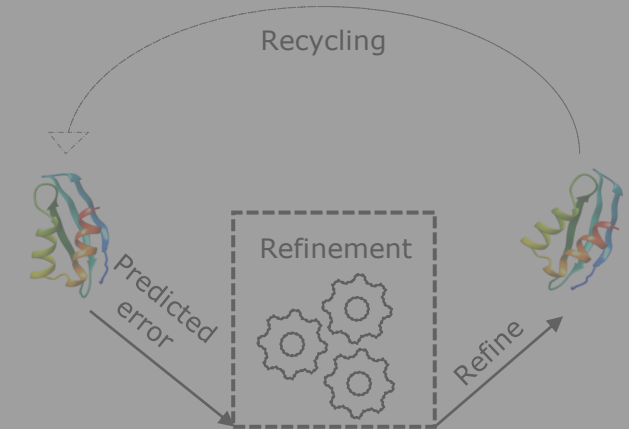


Protein model with single chain

## 2

Estimation of **multimeric** protein model quality



Protein model with multiple chains/subunits

## 3

Application of quality estimation to improve protein model quality



Iterative refinement of predicted protein models

# In this talk…



**View PDF**   Download full issue

## Computational and Structural Biotechnology Journal
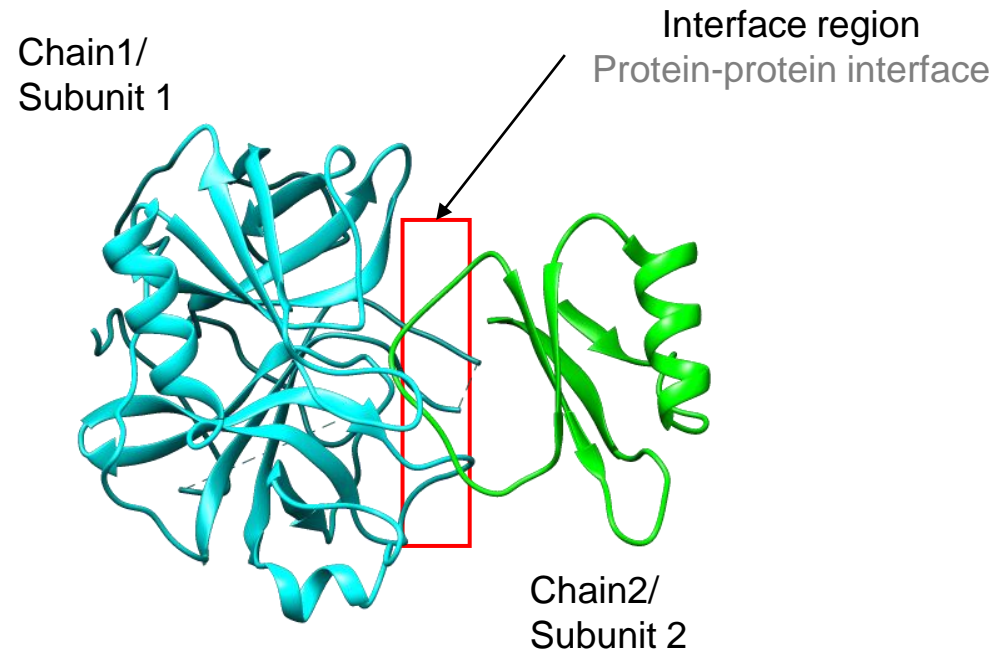Volume 27, 2025, Pages 160-170

ELSEVIER

Research Article

# EquiRank: Improved protein-protein interface quality estimation using protein language-model-informed equivariant graph neural networks

Md Hossain Shuvo [a], Debswapna Bhattacharya [b]

**https://github.com/mhshuvo1/EquiRank**

# Protein-protein interface of multimeric model
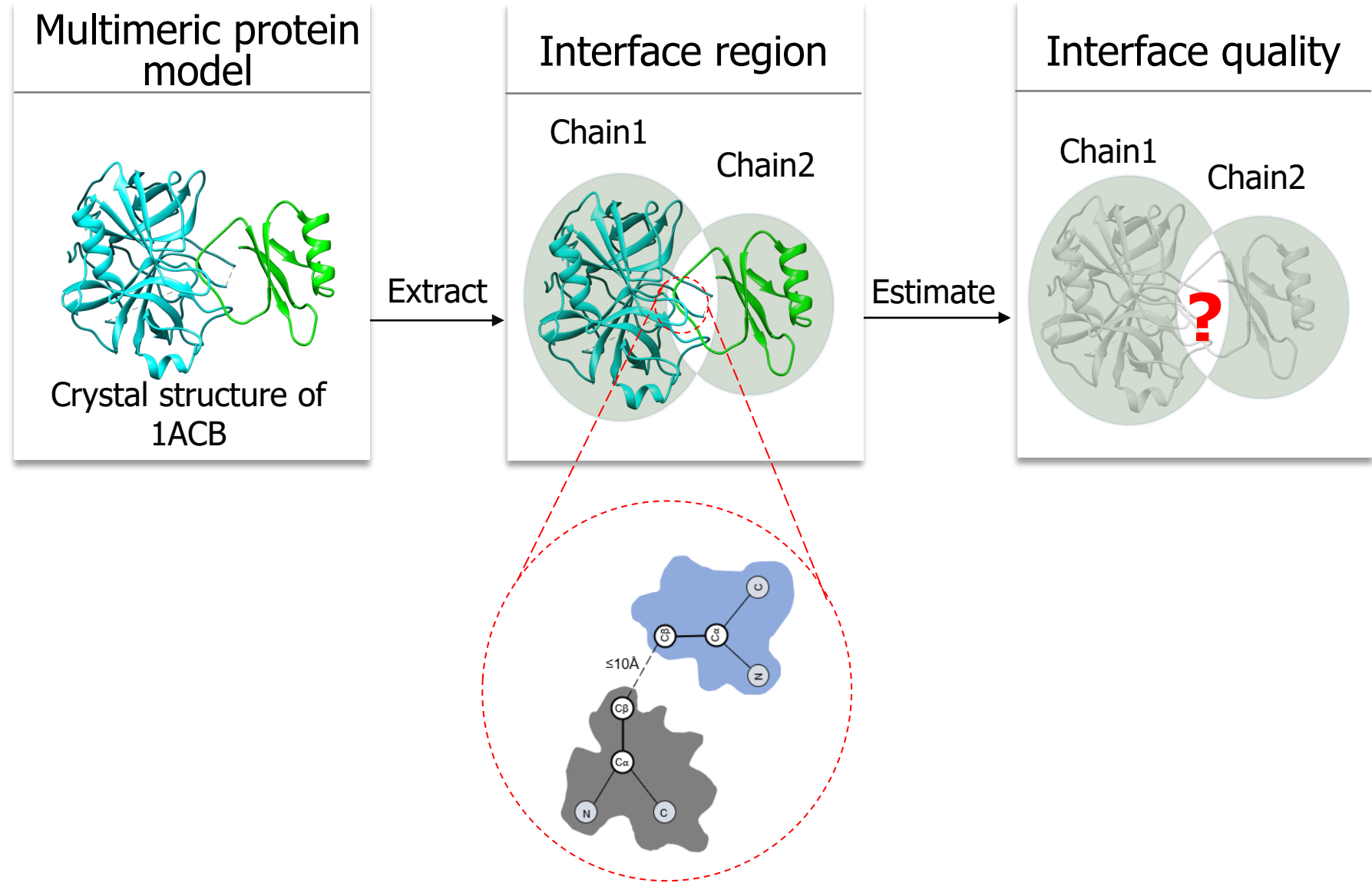
Chain1/
Subunit 1

Interface region
Protein-protein interface



Chain2/
Subunit 2

- ➢ # of chains/subunits > 1

- ➢ A.K.A. protein complex

- ➢ Protein-protein interaction

- ➢ Catalyzes biological processes

- ➢ Interface quality → multimeric protein quality

# Protein-protein interface quality estimation

> Estimation of protein-protein interface quality of computationally predicted protein multimer/protein complexes

> When the experimental structure is not available/absent

**Multimeric protein model**

Crystal structure of 1ACB

Extract →

**Interface region**

Chain1

Chain2

≤10Å

Estimate →

**Interface quality**

Chain1

Chain2

?

# Protein complex structure prediction

## Model selection



Best model

Pool of candidate structures

## Model refinement



Recycling

Predicted error

Refinement

Refine

Helps in accurately guide the process of protein complex structure prediction

# Research questions

**Representation**

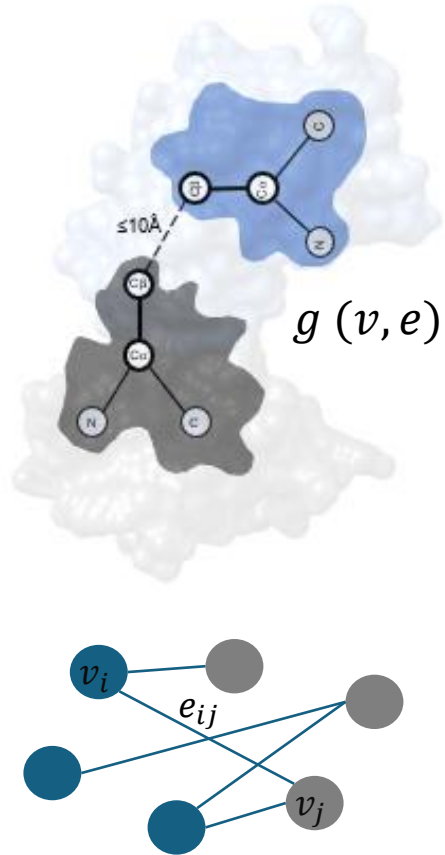**1**

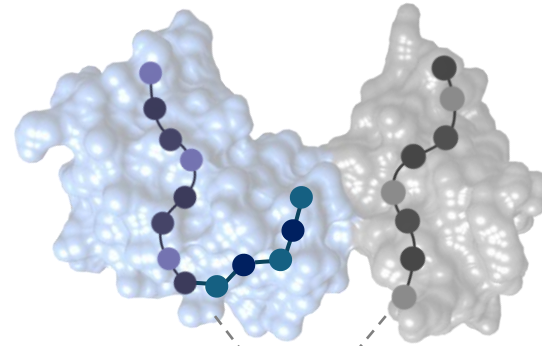How to better represent the interface of a protein complex?

**Learning the representation**

**2**

How to better learn the interface representation?

# 1. Representation: protein-protein interface

## Graph representation



$g\ (v,e)$

## Input representations:
Sequence-and structure- based embeddings



Sequence-based embeddings

Evolutionary information

Protein Language Model

Structure-based embeddings

## Output representations:
Multimeric geometry error



**Model**

$d_{ij}, \vec{a}_{ij}$

$Z_{ij}(d), Z_{ij}(\vec{a})$

$d_{ij}, \vec{a}_{ij}$

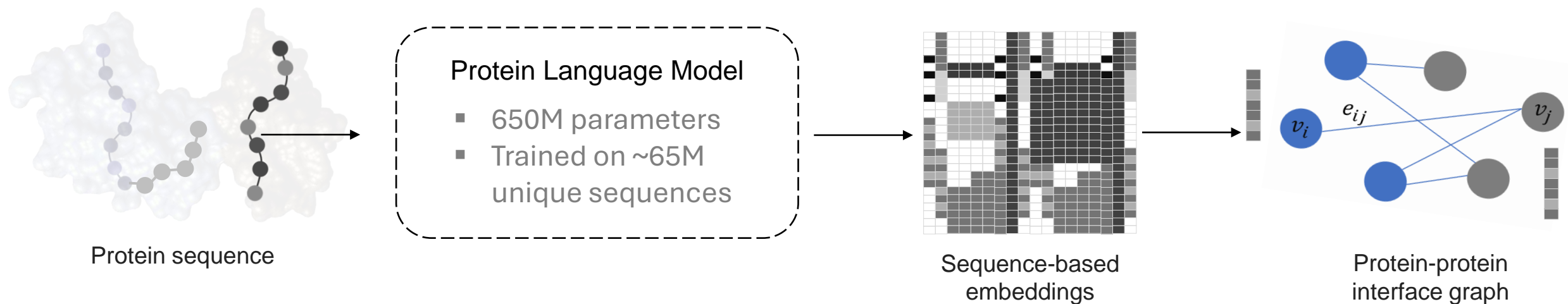**Experimental structure**

**Multimeric distance error:** $Z_{ij}(d)$

**Multimeric orientation errors:**
$Z_{ij}(\vec{a}), \text{where } \vec{a} \in \Omega, \tau_{ij}, \tau_{ji}, \lambda_{ij}, \lambda_{ji}$

Protein sequence

Sequence databases

UniRef100

PDB70

Environmental

Multiple Sequence Alignment (MSA)

$v_i$   $e_{ij}$   $v_j$

Protein-protein interface graph

Protein Language Model

- 650M parameters
- Trained on ~65M unique sequences

Protein sequence

Sequence-based embeddings

$v_i$ $e_{ij}$ $v_j$

Protein-protein interface graph

Captures the evolutionary patterns at scale

- Physiochemical properties
- Structural topology
- Neighborhood information
- Structural geometry

$v_i$  $e_{ij}$  $v_j$

Protein-protein
interface graph

Expressive set of structural embeddings for protein-protein interface

Multimeric distance: $d_{ij}$

Multimeric distance error: $z_{ij}$

$$z_{ij}(d) = \begin{cases} 1 & \text{if } d_{ij}^{\text{model}} < 10\,\text{Å and } d_{ij}^{\text{native}} < 10\,\text{Å} \\[2em] \dfrac{1}{1+\left(\dfrac{\left|d_{ij}^{\text{model}}-d_{ij}^{\text{native}}\right|}{d_0}\right)^2} & \text{otherwise} \end{cases}$$

Multimeric orientations: $\vec{a}_{ij}$



symmetric

asymmetric

**Torsion angles:** $\Omega$, $\tau_{ij}$, $\tau_{ji}$

**Planar angles:** $\lambda_{ij}$, $\lambda_{ji}$

*Symmetric*
$$\Omega_{ij} = c\alpha_i - c\beta_i - c\beta_j - c\alpha_j$$

*Asymmetric*
$$\tau_{ij} = N_i - c\alpha_i - c\beta_i - c\beta_j$$
$$\tau_{ji} = N_j - c\alpha_j - c\beta_j - c\beta_i$$
$$\lambda_{ij} = c\alpha_i - c\beta_i - c\beta_j$$
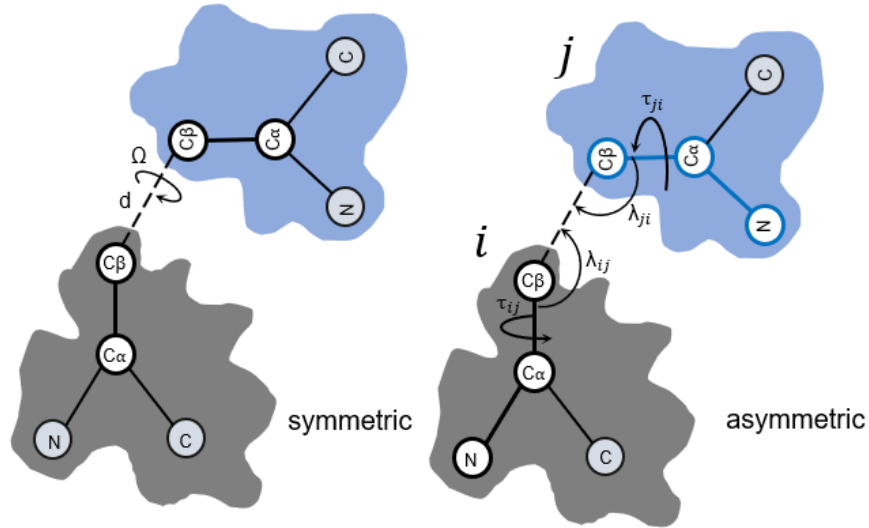$$\lambda_{ji} = c\alpha_j - c\beta_j - c\beta_i$$

Multimeric orientations errors



**Model**

$\vec{a}_{ij}$

$\mathbf{z}_{ij}(\vec{a})$

$\vec{a}_{ij}$

**Native**

*Angular RMSD*

$$z_{ij}(\vec{a}) = \sqrt{\left(\min\left(\left|a_{ij}^{\text{native}} - a_{ij}^{\text{model}}\right|, 2\pi - \left|a_{ij}^{\text{native}} - a_{ij}^{\text{model}}\right|\right)\right)^2}$$

*Where,* $\vec{a} \in \Omega, \tau_{ij}, \tau_{ji}, \lambda_{ij}, \lambda_{ji}$

**Representation**

**1**

How to better represent the interface of a protein complex?

**Learning the representation**

**2**

How to better learn the interface representation?

# How to learn a better mapping between the input and the output representations?

## Graph neural network



$$\mathbf{m}_{ij} = \phi_e(\mathbf{h}_i^l, \mathbf{h}_j^l, a_{ij})$$

$$\mathbf{m}_i = \sum_{j \in \mathcal{N}(i)} \mathbf{m}_{ij}$$

$$\mathbf{h}_i^{l+1} = \phi_h(\mathbf{h}_i^l, \mathbf{m}_i)$$

Permutation invariant

## Graph attention network



PIQLE, published in
Bioinformatics advances

Shuvo *et al.*, 2023

## Symmetry-aware graph neural network



**New output representations**

➢ Multimeric distance error
➢ Multimeric orientation errors

**Equivariance**

Symmetry-aware equivariant neural network

Zhou et al., 2020
Veličković et al., 2018

Equivariance

➢ Let $f: X \to Y$ be a neural network function

➢ $\emptyset^x$ and $\emptyset^y$ are the transformations on $X$ and $Y$, respectively

➢ $f$ is equivariant $iff$

$$f \circ \emptyset_g^x = \emptyset_g^y \circ f$$





Equivariant network

Equivariant network

Satorras *et al.*, 2022

# Research outcome

**?** Can we use graph representation for protein interface and learn the representation using a symmetry-aware graph neural networks?

**Research Outcome**

✓

**Representation**

**1**

Graph (node + edge + output representations)

**Learning the representation**

**2**

Equivariant Graph Neural Network (EGNN)

**EquiRank:** improved protein-protein interface quality estimation using protein-language-model-informed equivariant graph neural networks

# Flowchart of EquiRank



$EquiRank_{score}$ is the probabilistic combination of estimated multimeric distance and orientation error

# Datasets

| | Dataset | Num. targets | Num. decoys | Correct (DockQ ≥ 0.23) | Incorrect (DockQ < 0.23) |
|---|---|---|---|---|---|
| **Training** | VoroIF_GNN_train | 1,097 | 14,400 | 42% | 52% |
| | CASP13 | 20 | 2,386 | 17.24% | 82.76% |
| | CASP14 | 10 | 1,329 | 15.95% | 84.05% |
| **Testing** | VoroIF_GNN_test | 235 | 2,845 | 42% | 58% |
| | CASP15 | 26 | 6,850 | 45.37% | 54.63% |
| | Dockground v1 | 23 | 2,500 | 10.72% | 89.28% |
| **Validation** | VoroIF_GNN_validation | 235 | 2,814 | 40.96% | 59.14% |

**CASP: Critical Assessment of Protein Structure Prediction**

**Dockground v1: Dockground version 1**

# Evaluation metrices

## Ground truth

> DockQ score



**Model**

**Experimental**

> DockQ = [0, 1]

## Ability to Rank

> Reproducibility of ranking w.r.t. ground truth DockQ scores

❑ **Spearman Correlation** between the DockQ and predicted interface quality scores.

❑ Higher correlation indicates better reproducibility

❑ **Top-N Success Rate (N = 1, 5, 10, 15, 20, 25, 30)**

$$SR(N) = \frac{S(N)}{K} \times 100$$

❑ **Top-N Hit Rate (N = 1, 5, 10, 15, 20, 25, 30)**

$$HR(N) = \frac{H(N)}{M} \times 100$$

❑ Higher Rates indicates better performance

## Ability to Distinguish

> High quality protein complex models

❑ DockQ cutoff = 0.80

❑ Receiver Operating Characteristics Area Under the Curve

❑ Higher AUC is better

# Competing methods

➢ **Graph Neural Network-based methods**
1. PIQLE
2. EuDockScore
3. VoroIF_GNN
4. DProQA
5. GDockScore
6. DeepRank-GNN-esm
7. GNN-DOVE

➢ **Transformer-based method**
8. AlphaFold-Multimer

➢ **Convolutional Neural Network-based methods**
9. TRScore
10. DOVE_ATOM20
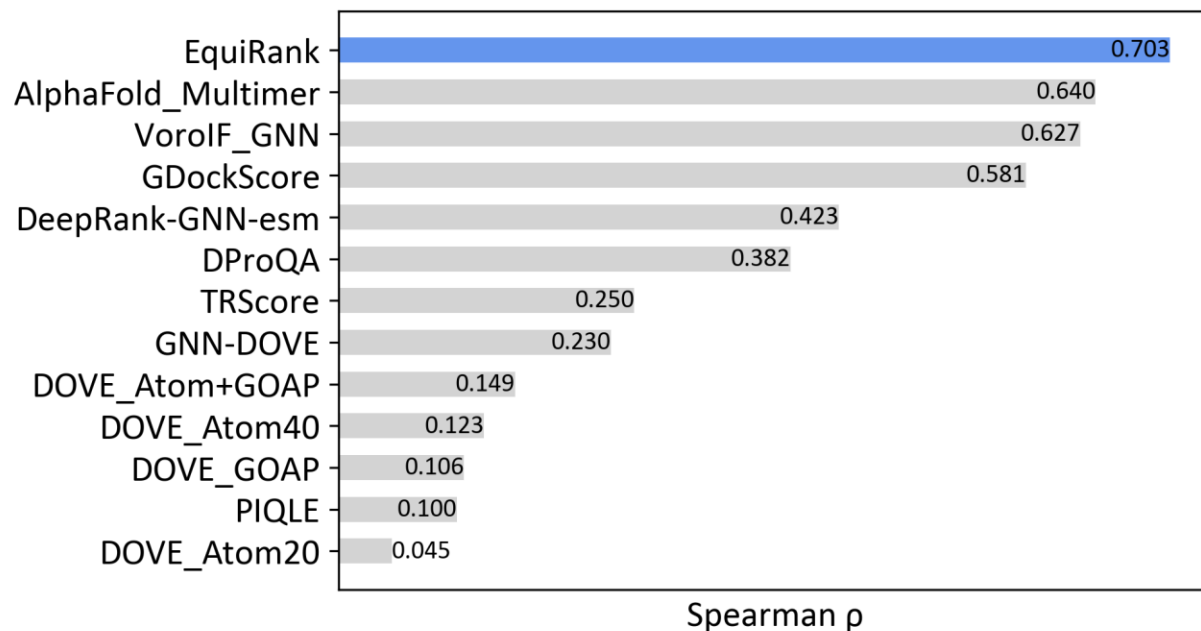11. DOVE_ATOM40
12. DOVE_GOAP
13. DOVE_ATOM_GOAP

Matthew et al., 2024
Shuvo *et al.*, 2023
Olechnovič and Venclovas, 2023
Chen *et al.*, 2023
McFee and Kim, 2023
Xu and Bonvin, 2023
Wang *et al.*, 2021
Evans *et al.*, 2022
Guo *et al.*, 2022
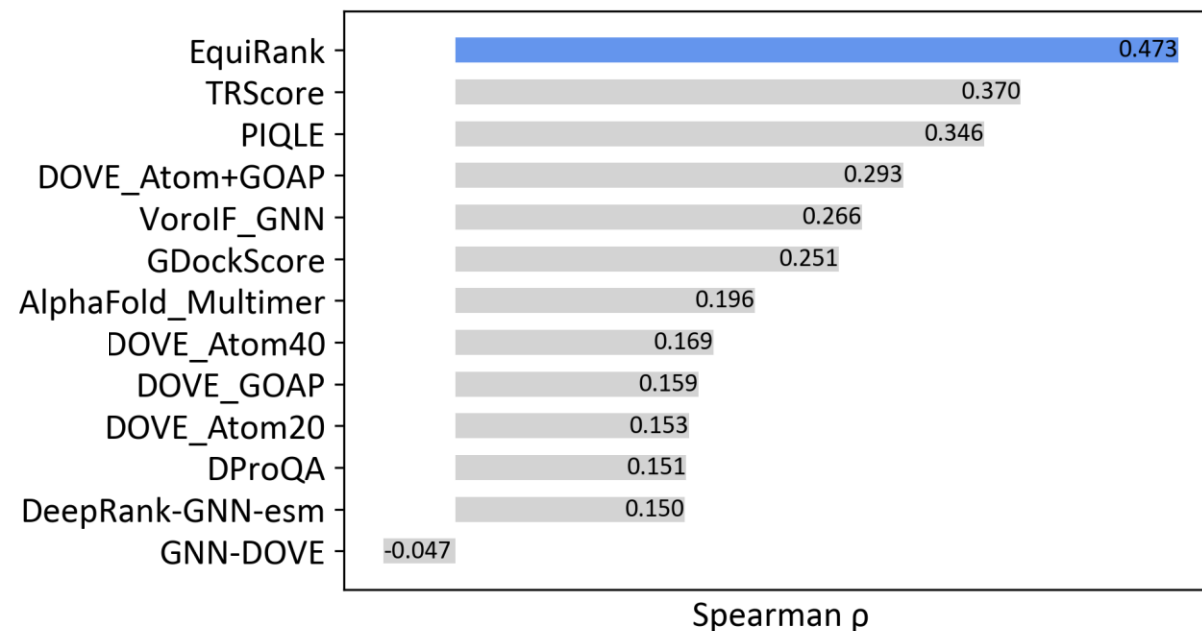Wang *et al.*, 2020

# Results

# **Ability to Rank:** Reproducibility of DockQ ranking
Spearman correlation coefficient between predicted and DockQ scores



VoroIF_GNN_test → total decoys: 2,845 Correct: 42% Incorrect: 58%)

| | Spearman ρ |
|---|---|
| EquiRank | 0.703 |
| AlphaFold_Multimer | 0.640 |
| VoroIF_GNN | 0.627 |
| GDockScore | 0.581 |
| DeepRank-GNN-esm | 0.423 |
| DProQA | 0.382 |
| TRScore | 0.250 |
| GNN-DOVE | 0.230 |
| DOVE_Atom+GOAP | 0.149 |
| DOVE_Atom40 | 0.123 |
| DOVE_GOAP | 0.106 |
| PIQLE | 0.100 |
| DOVE_Atom20 | 0.045 |

Dockground v1→ total decoys: 2,500 Correct: 10.72% Incorrect: 89.28%

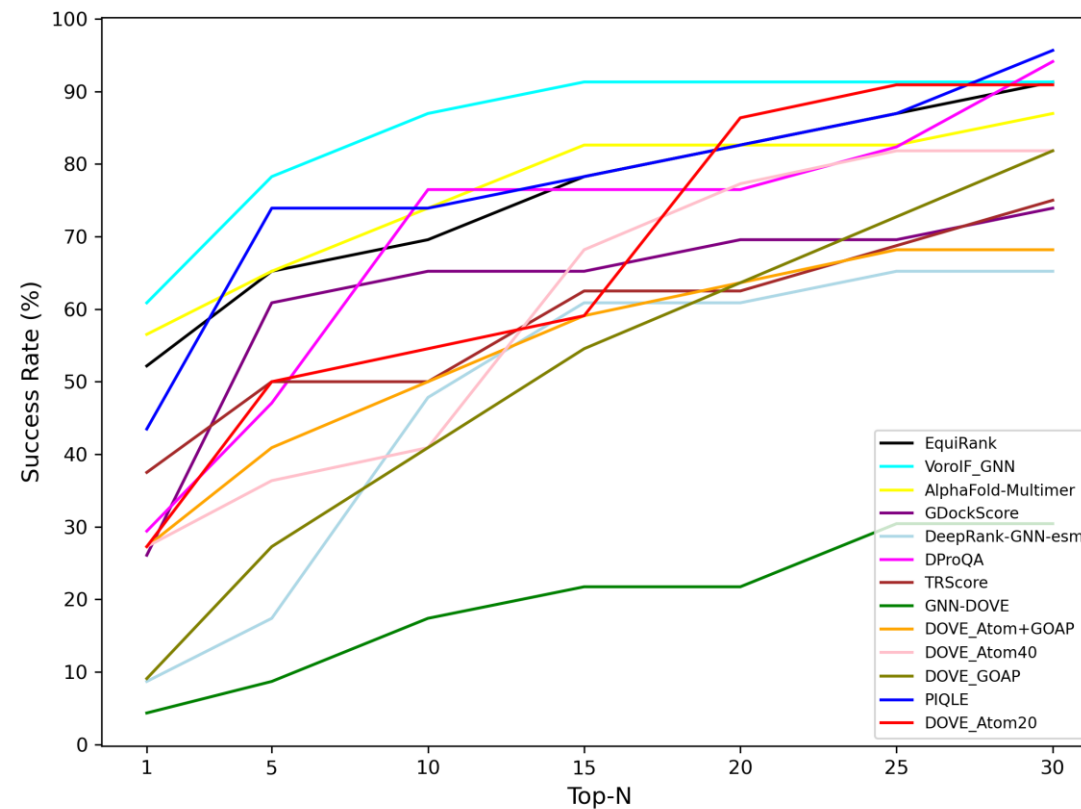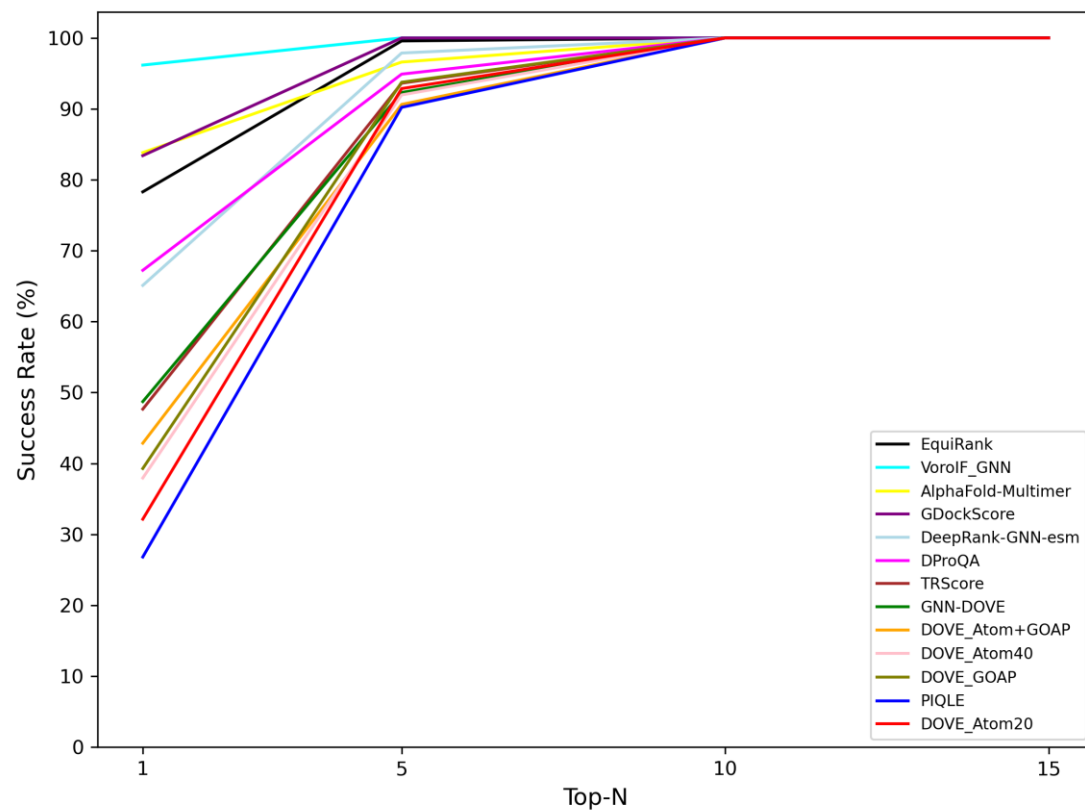| | Spearman ρ |
|---|---|
| EquiRank | 0.473 |
| TRScore | 0.370 |
| PIQLE | 0.346 |
| DOVE_Atom+GOAP | 0.293 |
| VoroIF_GNN | 0.266 |
| GDockScore | 0.251 |
| AlphaFold_Multimer | 0.196 |
| DOVE_Atom40 | 0.169 |
| DOVE_GOAP | 0.159 |
| DOVE_Atom20 | 0.153 |
| DProQA | 0.151 |
| DeepRank-GNN-esm | 0.150 |
| GNN-DOVE | -0.047 |

EquiRank is better in reproducing ground truth ranking

# Ability to Rank: Top-N Success Rate

Percentage of targets with at least one acceptable model among top-N ranked models



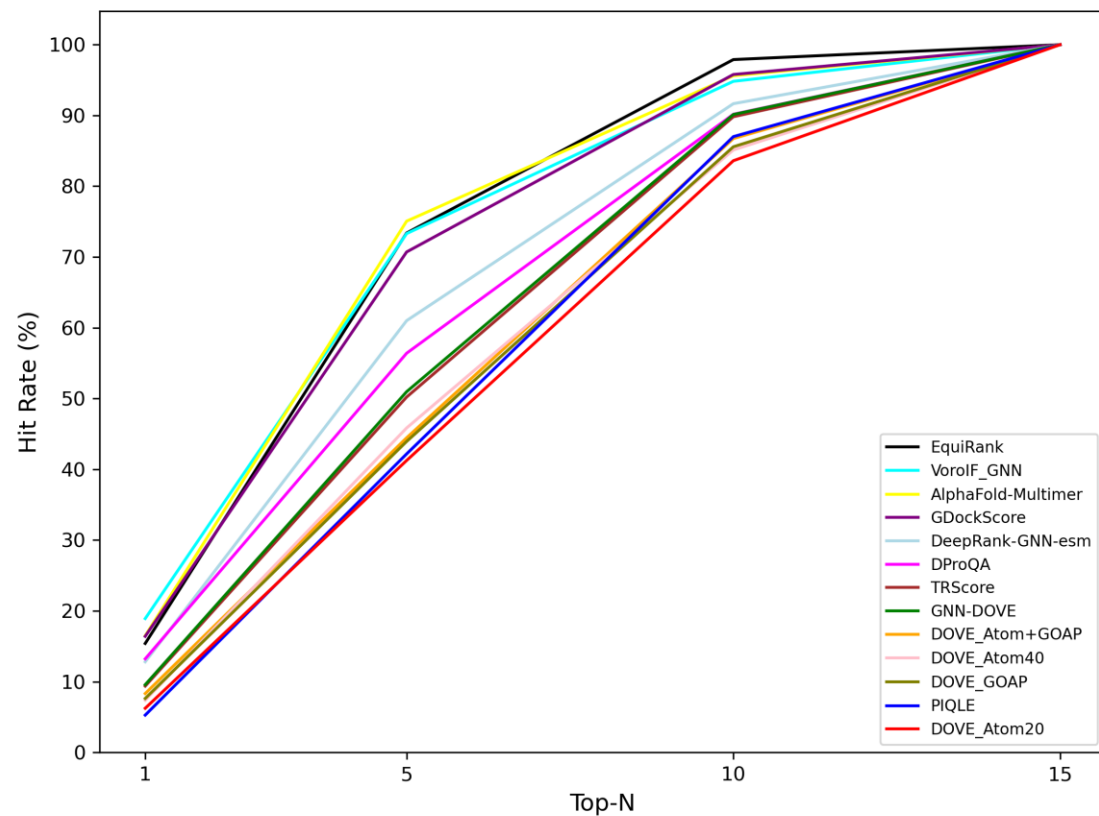VoroIF_GNN_test → total decoys: 2,845 Correct: 42% Incorrect: 58%)

Dockground v1→ total decoys: 2,500 Correct: 10.72% Incorrect: 89.28%
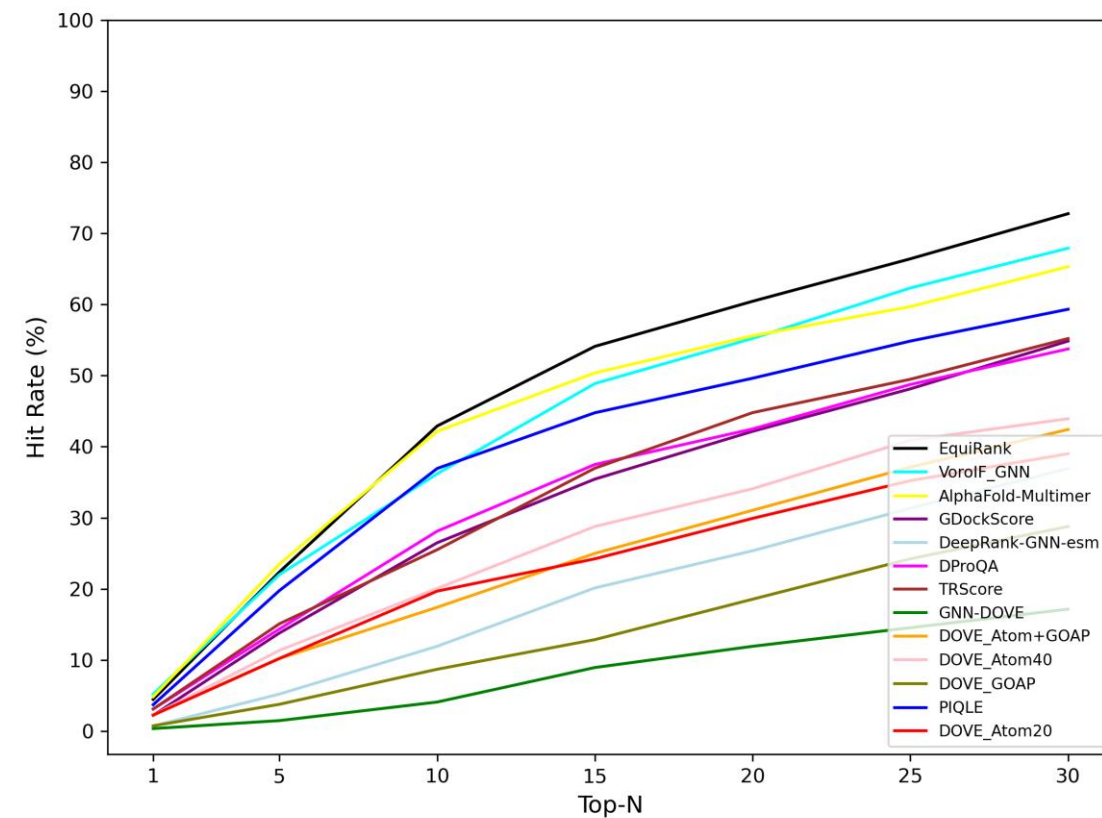
# Ability to Rank: Top-N Hit Rate

Fraction of acceptable models among top-ranked models relative to all acceptable models in the dataset

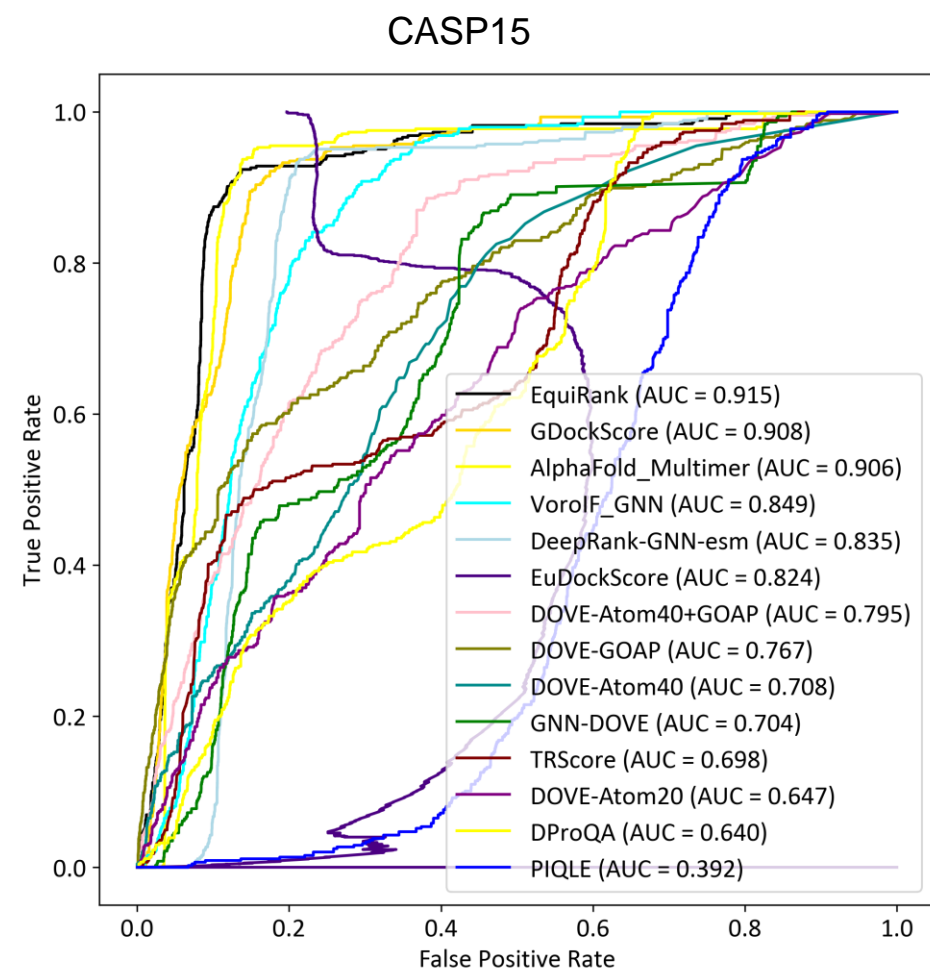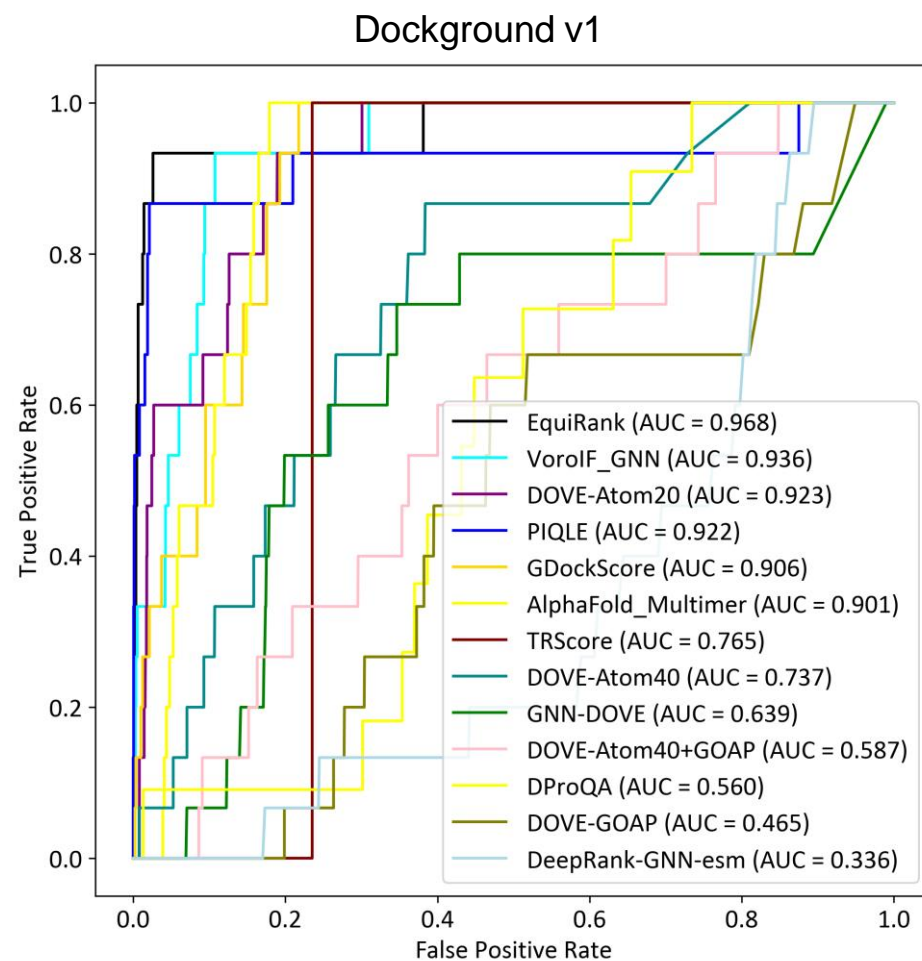VoroIF_GNN_test → total decoys: 2,845 Correct: 42% Incorrect: 58%)

Dockground v1→ total decoys: 2,500 Correct: 10.72% Incorrect: 89.28%
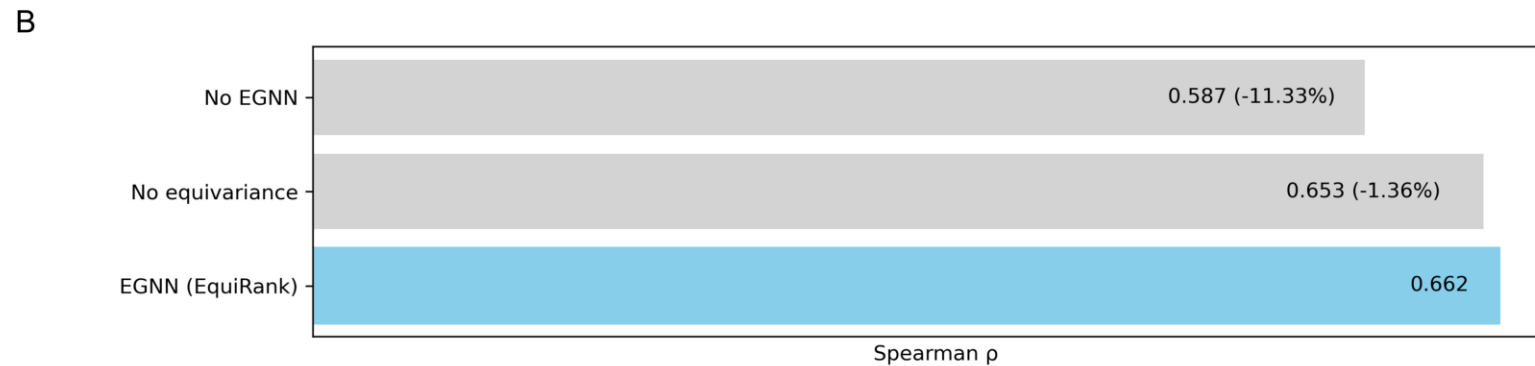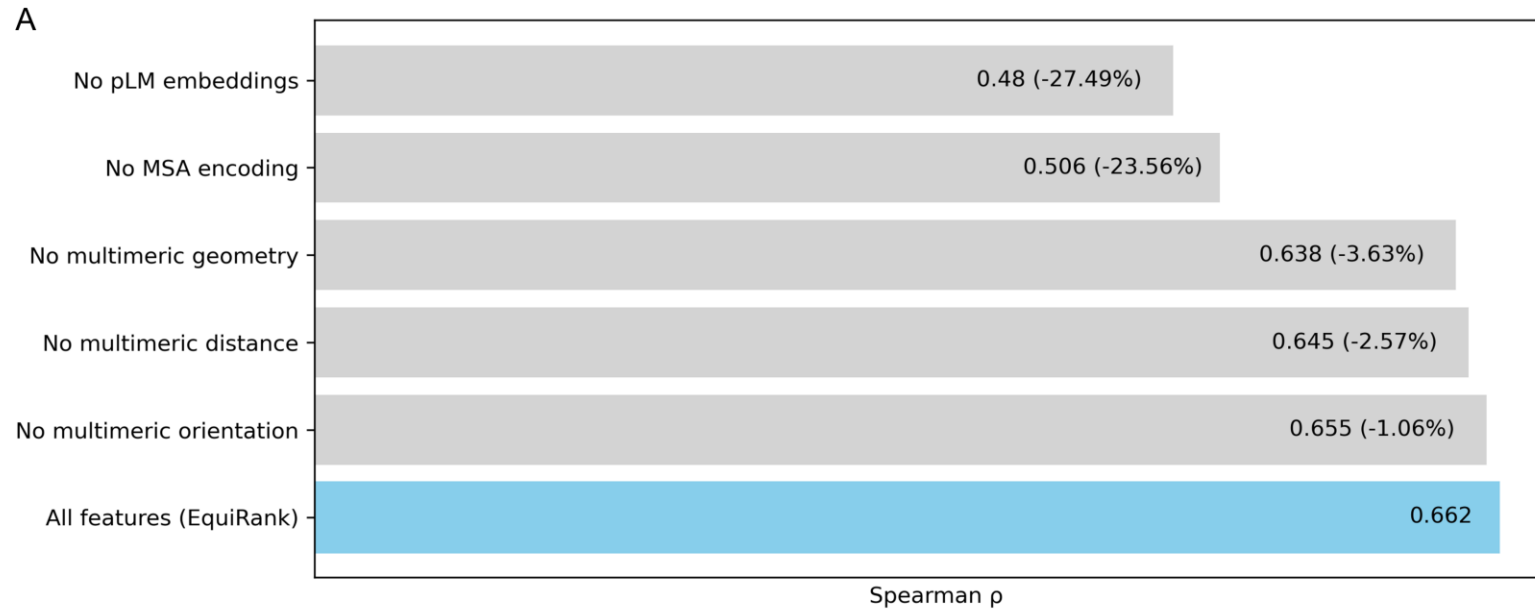
# Ability to Distinguish: Distinguishing high-quality complex models
Area Under the ROC curve with DockQ threshold = 0.80



Dockground v1

CASP15

EquiRank is better in distinguishing high-quality protein complex models

# Ablation studies (on VoroIF_GNN_validation)



A

| | |
|---|---|
| No pLM embeddings | 0.48 (-27.49%) |
| No MSA encoding | 0.506 (-23.56%) |
| No multimeric geometry | 0.638 (-3.63%) |
| No multimeric distance | 0.645 (-2.57%) |
| No multimeric orientation | 0.655 (-1.06%) |
| All features (EquiRank) | 0.662 |

Spearman ρ

B

| | |
|---|---|
| No EGNN | 0.587 (-11.33%) |
| No equivariance | 0.653 (-1.36%) |
| EGNN (EquiRank) | 0.662 |

Spearman ρ

Protein Language Model embeddings and EGNN contributes to the improved model quality estimation performance

# Conclusion and future works

➢ Application of an Ensemble 6 Equivariant Graph Neural Networks

➢ EquiRank is better than other competing methods in terms of reproducibility and distinguishability

➢ EquiRank demonstrates consistent performance on datasets having diverse quality

➢ Protein-language-model-informed equivariant neural network contributes to improved performance

**In the future…**

➢ Improve the generalizability of the multimeric quality estimation method

    ➢ Hyperparameter optimization of the underlying model

➢ Improve the reproducibility of ground truth scores

➢ Application of multimeric model quality estimation to improve the predicted multimeric protein complex models

➢ Development of an integrated framework for multimeric protein model quality estimation

# Thank You!

# Toward reliability evaluation of computational models of protein molecules and their interactions

https://github.com/mhshuvo1/EquiRank

---

**Md Hossain Shuvo, Ph.D.**

mhshuvo@pvamu.edu