



Boulder

# Computational Biology

**Maziar Raissi**

**Assistant Professor**

Department of Applied Mathematics

University of Colorado Boulder

[maziar.raissi@colorado.edu](mailto:maziar.raissi@colorado.edu)

# Improved Protein Structure Prediction using Potentials from Deep Learning

**Protein folding problem:** determine the three-dimensional shape of a protein from its amino acid sequence (21 amino acids)

SQETRKKCTEMKKFKNCEVRCDESNCHEVRCSDTKYTL



Protein Data Bank (PDB)

<https://www.rcsb.org/3d-view/5W9F>

Critical Assessment of Protein Structure Prediction (CASP13)

$S = (s_1, \dots, s_L) \rightarrow$  amino acid sequence of a protein

$s_i \rightarrow i\text{-th residue}$

$\text{MSA}(S) \rightarrow$  multiple sequence alignment features (HHblits & PSI-BLAST)

The input to the network consists of a two-dimensional array of features in which each  $i, j$  feature is the concatenation of the one-dimensional features for both  $i$  and  $j$  as well as the two-dimensional features for  $i, j$ .

$P(\varphi_i, \psi_i | S, \text{MSA}(S)) \rightarrow$  discrete probability distributions of backbone torsion angles

Neural Network

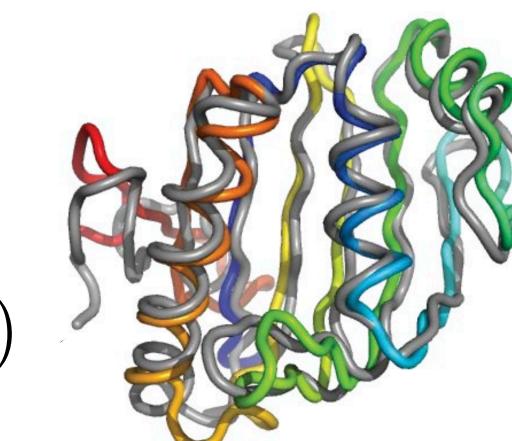
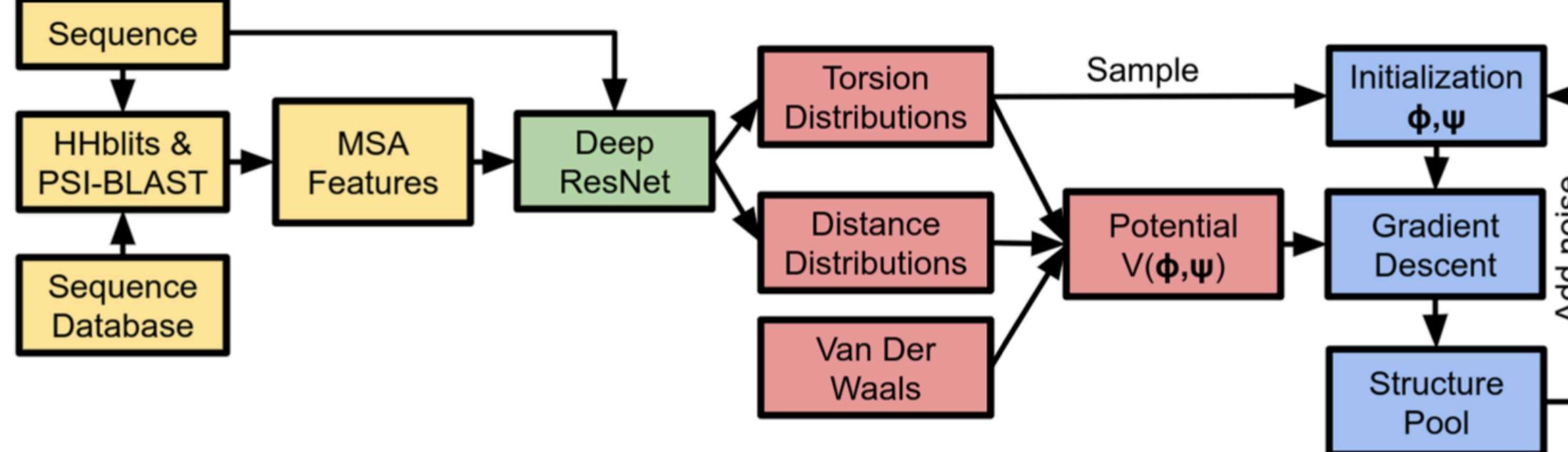
$x_i \rightarrow$  coordinates for residue  $i$

$x = G(\varphi, \psi) \rightarrow$  build a differentiable model  $G$  (Neural Network) of protein geometry

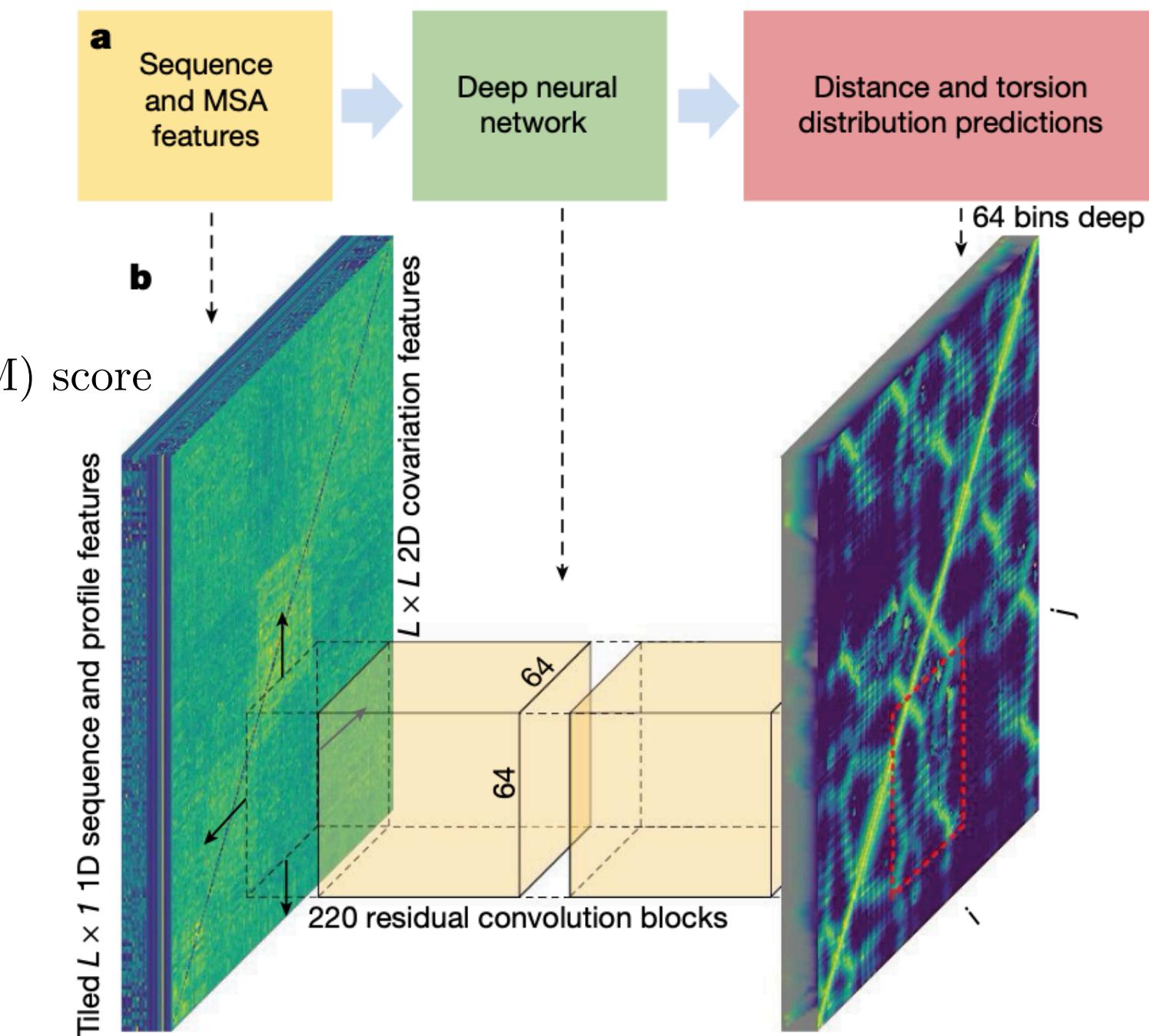
$d_{ij} = \|x_i - x_j\| \rightarrow$  inter-residue distances

$P(d_{ij} | S, \text{MSA}(S)) \rightarrow$  discrete probability distribution for every  $ij$  pair

Neural Network



Template Modelling (TM) score



Potentials

$$\begin{aligned}
 V_{\text{distance}}(\mathbf{x}) &= - \sum_{i,j, i \neq j} \log P(d_{ij} | \mathcal{S}, \text{MSA}(\mathcal{S})) && \text{fitting a spline} \\
 V_{\text{distance}}(\mathbf{x}) &= - \sum_{i,j, i \neq j} \log P(d_{ij} | \mathcal{S}, \text{MSA}(\mathcal{S})) - \log P(d_{ij} | \text{length}, \delta_{\alpha\beta}) && \text{reference} \\
 V_{\text{torsion}}(\boldsymbol{\phi}, \boldsymbol{\psi}) &= - \sum_i \log p_{\text{vonMises}}(\phi_i, \psi_i | \mathcal{S}, \text{MSA}(\mathcal{S})) && \text{glycine } (\text{C}_\alpha \text{ atom}) \text{ or not } (\text{C}_\beta) \\
 V_{\text{total}}(\boldsymbol{\phi}, \boldsymbol{\psi}) &= V_{\text{distance}}(G(\boldsymbol{\phi}, \boldsymbol{\psi})) + V_{\text{torsion}}(\boldsymbol{\phi}, \boldsymbol{\psi}) + V_{\text{score2_smooth}}(G(\boldsymbol{\phi}, \boldsymbol{\psi})) && \text{a van der Waals term}
 \end{aligned}$$



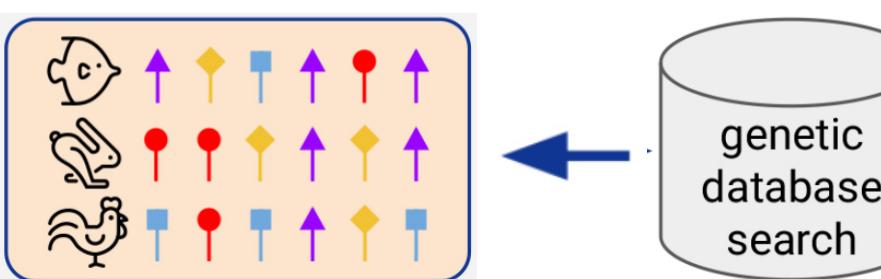
Boulder

# Highly Accurate Protein Structure Prediction with AlphaFold

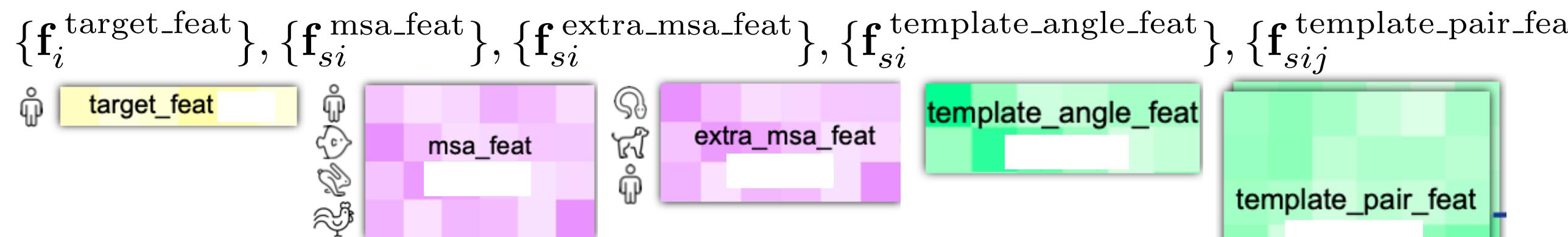
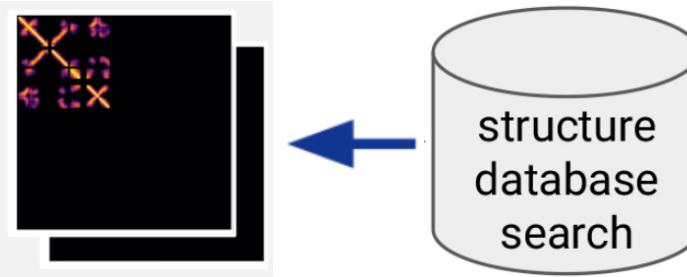
AlphaFold receives input features derived from the amino-acid sequence, MSA, and templates and outputs features including atom coordinates, the distogram, and per-residue confidence scores.



MSA (multiple sequence alignment): sequences of evolutionary related proteins



Templates: 3D atom coordinates of a small number of homologous structures



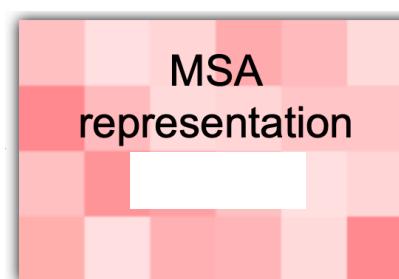
$i, j, k \rightarrow$  operate on the residue dimension

$s, t \rightarrow$  operate on the sequence dimension

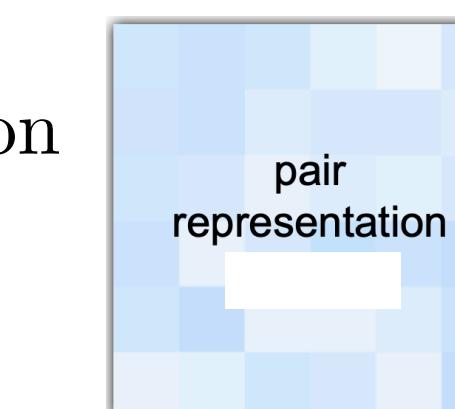
## Input Embedder

$\{f_i^{\text{target\_feat}}\}, \{f_{si}^{\text{msa\_feat}}\} \mapsto \{m_{si}\}, \{z_{ij}\}$

$\{m_{si}\} \rightarrow$  MSA representation



$\{z_{ij}\} \rightarrow$  pair representation



## Embed Templates

$\{f_{si}^{\text{template\_angle\_feat}}\}, \{f_{sij}^{\text{template\_pair\_feat}}\}, \{z_{ij}\} \mapsto \{z_{ij}\}$

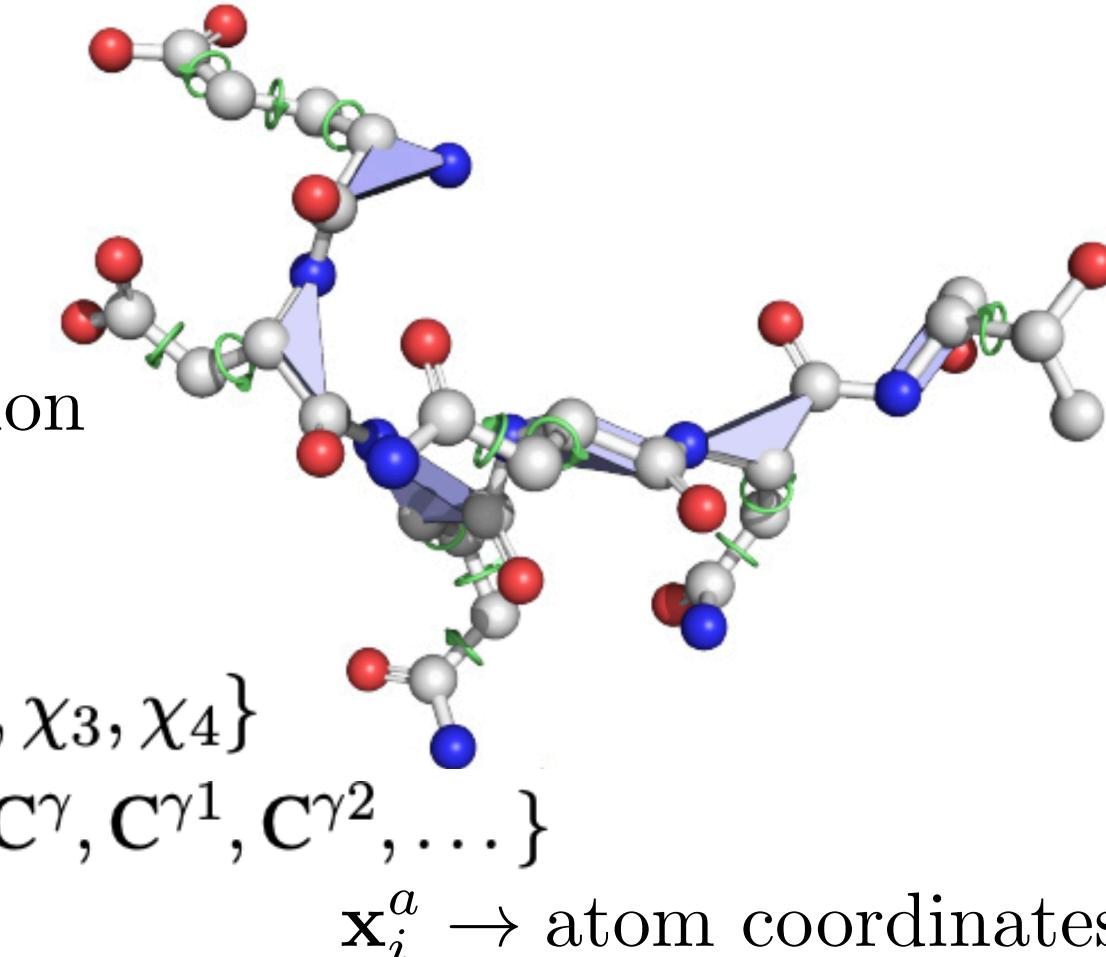
## Embed extra MSA features

$\{f_{si}^{\text{extra\_msa\_feat}}\}, \{z_{ij}\} \mapsto \{z_{ij}\}$

## Evoformer

$\{m_{si}\}, \{z_{ij}\} \mapsto \{m_{si}\}, \{z_{ij}\}, \{s_i\}$

$\{s_i\} \rightarrow$  abstract single representation



## Structure Module

$\{z_{ij}\}, \{s_i\} \mapsto \{T_i^f\}, \{x_i^a\}$

$f \in \mathcal{S}_{\text{torsion names}} = \{\omega, \phi, \psi, \chi_1, \chi_2, \chi_3, \chi_4\}$

$a \in \mathcal{S}_{\text{atom names}} = \{N, C^\alpha, C, O, C^\beta, C^\gamma, C^{\gamma 1}, C^{\gamma 2}, \dots\}$

$T_i = (R_i, t_i) \rightarrow$  backbone frames

$\vec{x}_{\text{global}} = T_i \circ \vec{x}_{\text{local}} = R_i \vec{x}_{\text{local}} + \vec{t}_i$

## Frame Aligned Point Error (FAPE)

```
def computeFAPE({T_i}, {x_j}, {T_i^true}, {x_j^true}, Z = 10Å, d_clamp = 10Å, ε = 10⁻⁴Å²):
```

$T_i, T_i^{\text{true}} \in (\mathbb{R}^{3 \times 3}, \mathbb{R}^3)$

$\vec{x}_j, \vec{x}_j^{\text{true}} \in \mathbb{R}^3$

$\vec{x}_{ij} \in \mathbb{R}^3$

$\vec{x}_{ij}^{\text{true}} \in \mathbb{R}^3$

$d_{ij} \in \mathbb{R}$

Chirality property!

$$1: \vec{x}_{ij} = T_i^{-1} \circ \vec{x}_j$$

$$2: \vec{x}_{ij}^{\text{true}} = T_i^{\text{true}-1} \circ \vec{x}_j^{\text{true}}$$

$$3: d_{ij} = \sqrt{\|\vec{x}_{ij} - \vec{x}_{ij}^{\text{true}}\|^2 + \epsilon}$$

$$4: \mathcal{L}_{\text{FAPE}} = \frac{1}{Z} \text{mean}_{i,j}(\min(d_{\text{clamp}}, d_{ij}))$$

Self-distillation procedure, (BERT)-style objective, iterative refinement (“recycling”), predict model confidence, and structural violations penalty.



Boulder

# Questions?

---