



上海科技大学
ShanghaiTech University

Simulating 500 Million Years of Evolution with a language Model

Yifan Qin



立志成才报国裕民

Intro: ESM2



上海科技大学
ShanghaiTech University

- Only trained on the sequence information
- Transformer based structure
- BERT-like Masked Language Task



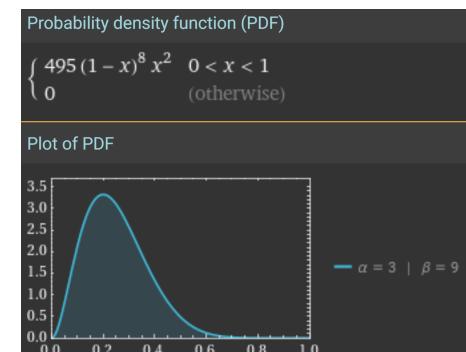
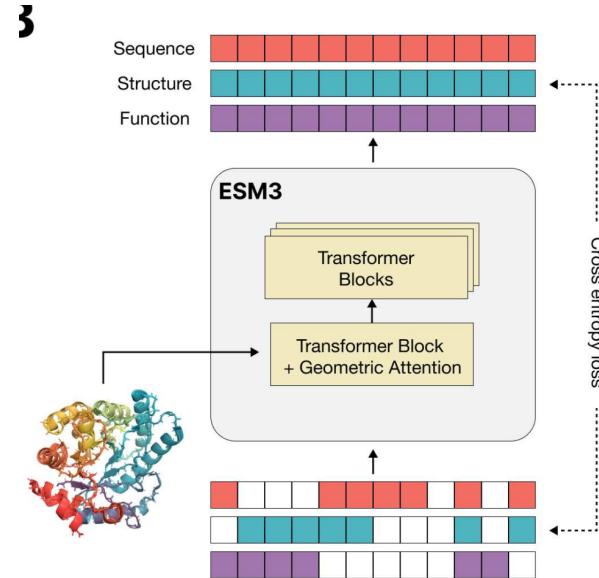
立志成才报国裕民



- Trained on structure, function annotation, and sequence etc.
- Transformer blocks with attention tailored
- Masked language task with special mask strategy



- Geometric attention to encode coordinates
- Special sequence mask strategy
 - 80% time, ratio~ $\beta(3, 9)$
 - 20% time, ratio~ Uniform(0, 1)



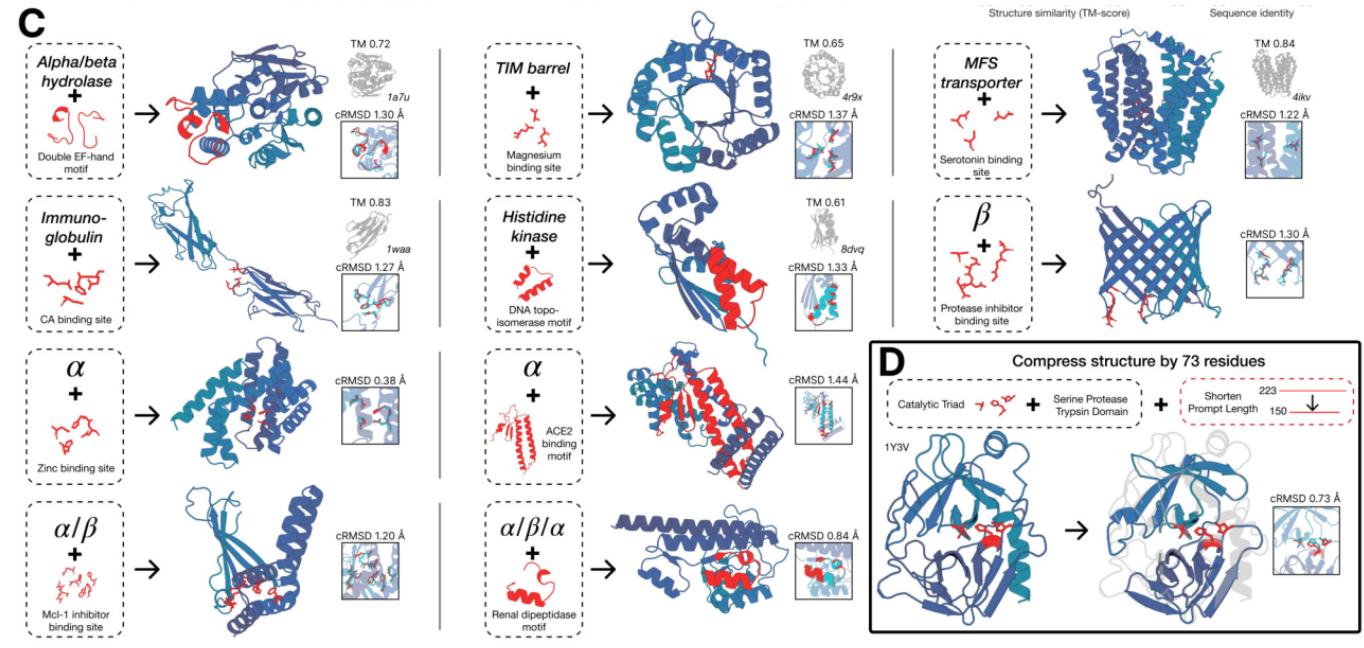
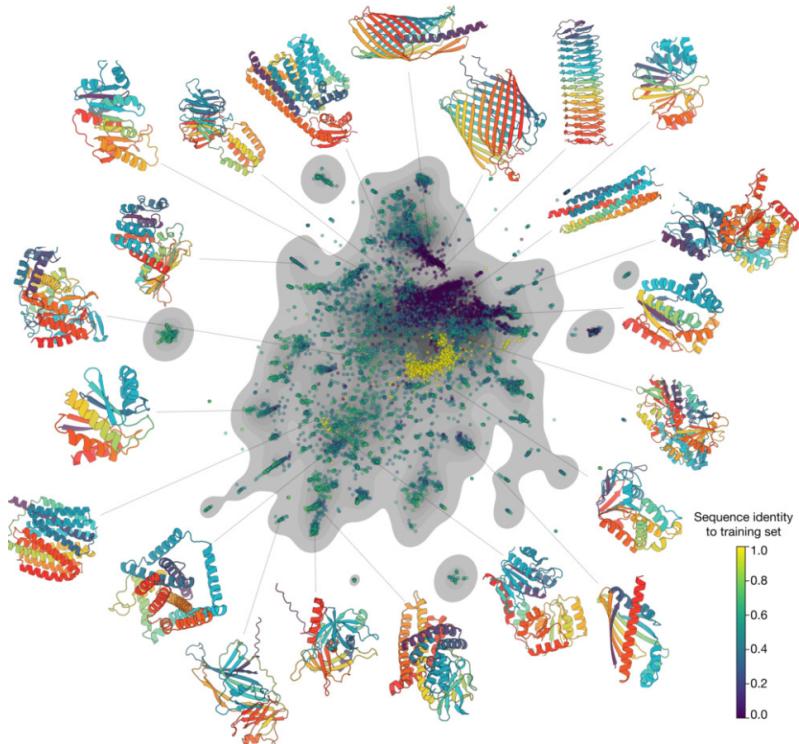
立志成才报国裕民

Insights & Results



上海科技大学
ShanghaiTech University

- Network architecture
- Unconditional design & Conditioned design



立志成才报国裕民



- Raw Input

- Sequence info
- Structure coordinates
- Secondary structure(SS) annotation
- Solvent Accessible Surface Area(SASA)
- Functional annotation
- Residue annotation



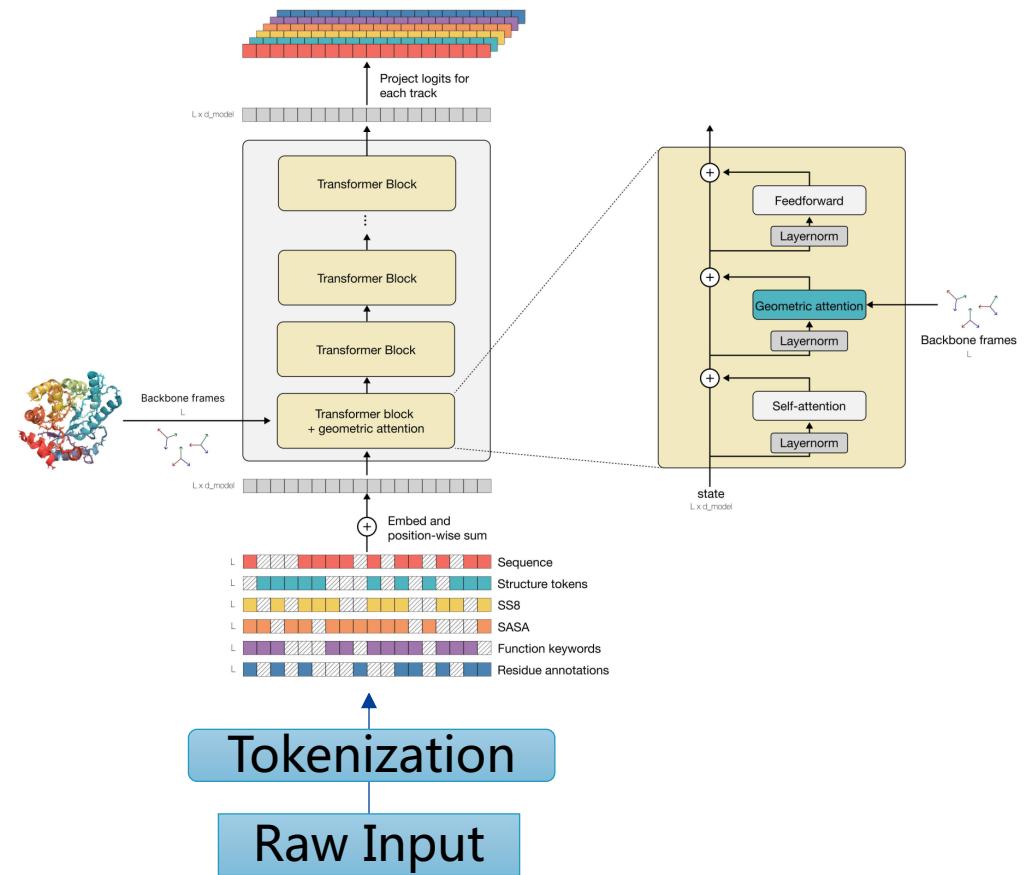
立志成才报国裕民

ESM3: Architecture



上海科技大学
ShanghaiTech University

- Tokenization: raw input -> tokens
- Transformer Trunk: transformer directly apply on tokens
- Decoder: tokens -> each track of prediction



立志成才报国裕民

ESM3: Architecture: Tokenization



上海科技大学
ShanghaiTech University

- Sequence info: 20 standard aa + 4 non-standard aa + BOS, EOS, mask, pad, unknown = 29 types
- SS annotation: 8 predefined SS type + mask, unknown = 10 types



立志成才报国裕民

ESM3: Architecture: Tokenization



上海科技大学
ShanghaiTech University

- SASA: continuous area value -> 16 fixed bin + mask,
unknown = 18 types
- Residue annotation: InterPro labels -> 1478 multi-hot
- Tokens:
 - Sequence: $\{0 \dots 28\}^L$, SS: $\{0 \dots 10\}^L$
 - SASA: $\{0 \dots 17\}^L$, Residue annotation: $\{0,1\}^{L \times 1478}$



立志成才报国裕民

ESM3: Structure Tokenization



上海科技大学
ShanghaiTech University

- Target: 3D coordinates to tokens for transformer input
- Use a VQ-VAE model to learn encode and decode of the coordinates
- Use encode part of VQ-VAE to get tokens for transformer input



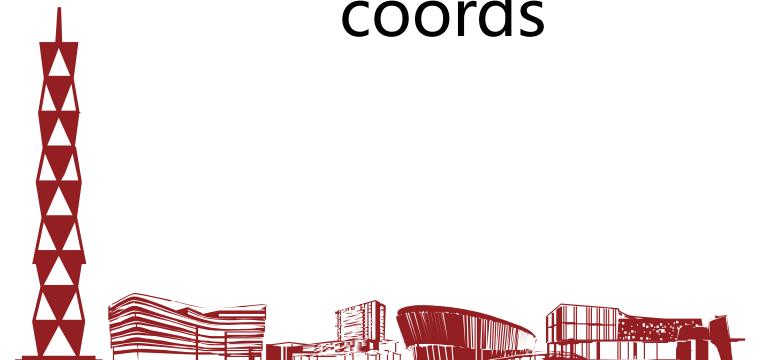
立志成才报国裕民

ESM3: Structure Tokenization



上海科技大学
ShanghaiTech University

- Frame & Geometric attention
- Frame $T \in SE(3)$
 - Given N, CA, C coords, get rotation and translation by gram_schmidt
 - Use rotation & translation to change between local and global coords



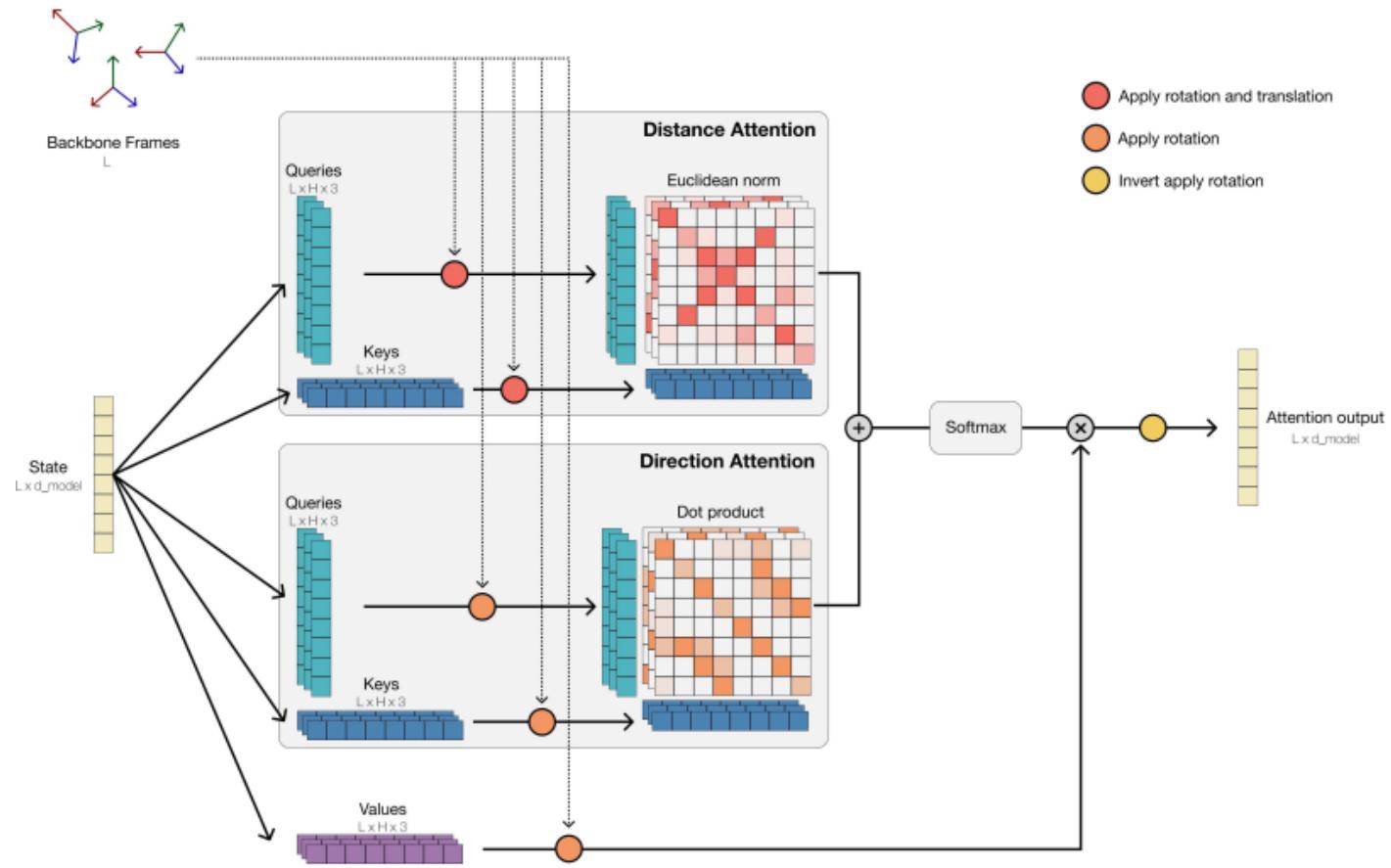
立志成才报国裕民

ESM3: Geometric Attention



上海科技大学
ShanghaiTech University

- Overview



立志成才报国裕民

ESM3: Geometric Attention



上海科技大学
ShanghaiTech University

- Given latent & frames
- Map latent to coords
- Convert coords to global distance frame and rotation frame , seen as Q, K, V

Algorithm 6 geometric_mha

Input: $X \in \mathbb{R}^{L \times d}, T \in SE(3)^L$

- 1: $Q_r, K_r, Q_d, K_d, V = \text{Linear}(X)$ $\triangleright (\mathbb{R}^{L \times h \times 3})_{\times 5}$
 - 2: $(\mathbf{R}_i, \mathbf{t}_i) = T_i$ $\triangleright (SO(3)^L, \mathbb{R}^{L \times 3})$
 - 3: $[Q_r]_{i,h,:} = \mathbf{R}_i([Q_r]_{i,h,:})$ $\triangleright \mathbb{R}^{L \times h \times 3}$
 - 4: $[K_r]_{i,h,:} = \mathbf{R}_i([K_r]_{i,h,:})$ $\triangleright \mathbb{R}^{L \times h \times 3}$
 - 5: $[Q_d]_{i,h,:} = T_i([Q_d]_{i,h,:})$ $\triangleright \mathbb{R}^{L \times h \times 3}$
 - 6: $[K_d]_{i,h,:} = T_i([K_d]_{i,h,:})$ $\triangleright \mathbb{R}^{L \times h \times 3}$
 - 7: $[R]_{i,j,h} = \frac{1}{\sqrt{3}}[q_r]_{i,h,:} \cdot [k_r]_{j,h,:}$ $\triangleright \mathbb{R}^{L \times L \times h}$
 - 8: $[D]_{i,j,h} = \frac{1}{\sqrt{3}}\|[q_r]_{i,h,:} - [k_r]_{j,h,:}\|_2$ $\triangleright \mathbb{R}^{L \times L \times h}$
 - 9: $A = \text{softplus}(\bar{w}_r)R - \text{softplus}(\bar{w}_d)D$ $\triangleright \mathbb{R}^{L \times L \times h}$
 - 10: $A = \text{softmax}_{\textcolor{brown}{j}}(A)$
 - 11: $[V]_{i,h,:} = \mathbf{R}_i([V]_{i,h,:})$
 - 12: $\textcolor{red}{O} = A \cdot V$ $\triangleright \mathbb{R}^{L \times h \times 3}$
 - 13: $[O]_{i,h,:} = \mathbf{R}_i^{-1}([O]_{i,h,:})$
 - 14: $X = X + \text{Linear}(\textcolor{yellow}{O})$ $\triangleright \mathbb{R}^{L \times d}$
-



立志成才报国裕民

ESM3: Geometric Attention



上海科技大学
ShanghaiTech University

- Attention within distance frame and rotation frame
- Combine attention and change V under rotation frame
- Map back to latent

Algorithm 6 geometric_mha

Input: $X \in \mathbb{R}^{L \times d}, T \in SE(3)^L$

- 1: $Q_r, K_r, Q_d, K_d, V = \text{Linear}(X)$ $\triangleright (\mathbb{R}^{L \times h \times 3})_{\times 5}$
 - 2: $(\mathbf{R}_i, \mathbf{t}_i) = T_i$ $\triangleright (SO(3)^L, \mathbb{R}^{L \times 3})$
 - 3: $[Q_r]_{i,h,:} = \mathbf{R}_i([Q_r]_{i,h,:})$ $\triangleright \mathbb{R}^{L \times h \times 3}$
 - 4: $[K_r]_{i,h,:} = \mathbf{R}_i([K_r]_{i,h,:})$ $\triangleright \mathbb{R}^{L \times h \times 3}$
 - 5: $[Q_d]_{i,h,:} = T_i([Q_d]_{i,h,:})$ $\triangleright \mathbb{R}^{L \times h \times 3}$
 - 6: $[K_d]_{i,h,:} = T_i([K_d]_{i,h,:})$ $\triangleright \mathbb{R}^{L \times h \times 3}$
 - 7: $[R]_{i,j,h} = \frac{1}{\sqrt{3}}[q_r]_{i,h,:} \cdot [k_r]_{j,h,:}$ $\triangleright \mathbb{R}^{L \times L \times h}$
 - 8: $[D]_{i,j,h} = \frac{1}{\sqrt{3}}\|[q_r]_{i,h,:} - [k_r]_{j,h,:}\|_2$ $\triangleright \mathbb{R}^{L \times L \times h}$
 - 9: $A = \text{softplus}(\bar{w}_r)R - \text{softplus}(\bar{w}_d)D$ $\triangleright \mathbb{R}^{L \times L \times h}$
 - 10: $A = \text{softmax}_{\textcolor{brown}{j}}(A)$
 - 11: $[V]_{i,h,:} = \mathbf{R}_i([V]_{i,h,:})$
 - 12: $\textcolor{red}{O} = A \cdot V$ $\triangleright \mathbb{R}^{L \times h \times 3}$
 - 13: $[O]_{i,h,:} = \mathbf{R}_i^{-1}([O]_{i,h,:})$
 - 14: $X = X + \text{Linear}(\textcolor{yellow}{O})$ $\triangleright \mathbb{R}^{L \times d}$
-



立志成才报国裕民

ESM3: Structure Tokenization Encode



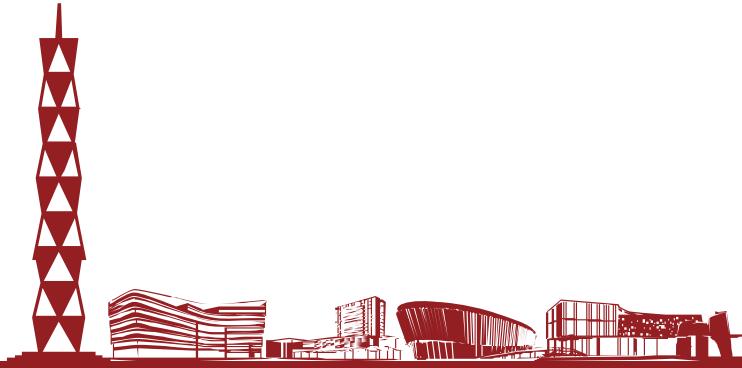
上海科技大学
ShanghaiTech University

- Find nearest neighbor in 3D (structure distance)
- Get latent from sequence distance
- Geometric attn f_{enc}
- Quantize by VQVAE codebook

Algorithm 7 structure_encode

Input: $x_{C_\alpha} \in \mathbb{R}^{L \times 3}, T \in SE(3)^L$

- | | |
|-----------------------------------------------------------|----------------------------------------------------|
| 1: $N_{\text{idx}} = \text{knn}(x_{C_\alpha})$ | $\triangleright \{0..L-1\}^{L \times 16}$ |
| 2: $T_{\text{knn}} = T[N_{\text{idx}}]$ | $\triangleright SE(3)^{L \times 16}$ |
| 3: $\Delta i = \text{clamp}(N_{\text{idx}} - i, -32, 32)$ | |
| 4: $N = \text{embed}(\Delta i)$ | $\triangleright \mathbb{R}^{L \times 16 \times d}$ |
| 5: $N = f_{enc}(N, T_{\text{knn}})$ | $\triangleright \mathbb{R}^{L \times 16 \times d}$ |
| 6: $z = \text{Linear}(N_{:,0,:})$ | $\triangleright \mathbb{R}^{L \times d'}$ |
| 7: $z = \text{quantize}(z)$ | $\triangleright \{0..4095\}^L$ |
-



立志成才报国裕民

ESM3: Structure Tokenization Decode



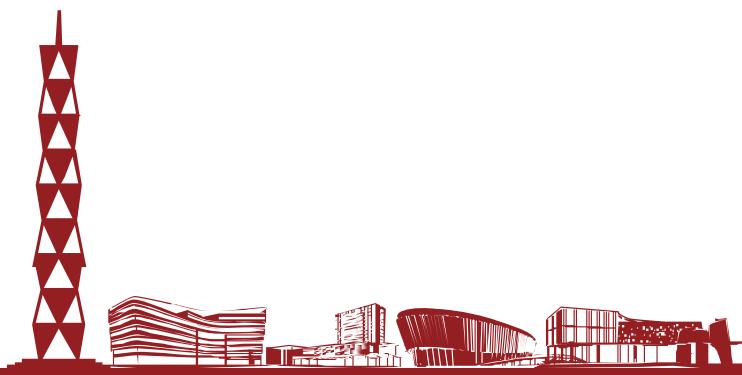
上海科技大学
ShanghaiTech University

- Transformer on the embedding
- Direct regress backbone coords and angles for calculating torsions

Algorithm 8 structure_decode

Input: $z \in \{0..4099\}^{L \times 16}$

- 1: $z = \text{embed}(z)$ $\triangleright \mathbb{R}^{L \times d}$
 - 2: $z = f_{dec}(z)$ $\triangleright \mathbb{R}^{L \times d}$
 - 3: $\vec{t}, \vec{x}, \vec{y}, \sin \theta, \cos \theta = \text{proj}(z)$ $\triangleright (\mathbb{R}^{L \times 3})_{\times 3}, (\mathbb{R}^{L \times 7})_{\times 2}$
 - 4: $T = \text{gram_schmidt}(\vec{t}, -\vec{x}, \vec{y})$ $\triangleright SE(3)^L$
 - 5: $\sin \theta = \frac{\sin \theta}{\sqrt{\sin \theta^2 + \cos \theta^2}}$ $\triangleright [-1, 1]^{L \times 7}$
 - 6: $\cos \theta = \frac{\cos \theta}{\sqrt{\sin \theta^2 + \cos \theta^2}}$ $\triangleright [-1, 1]^{L \times 7}$
 - 7: $T_{\text{local}} = \text{rot_frames}(\sin \theta, \cos \theta)$ $\triangleright SE(3)^{L \times 7}$
 - 8: $T_{\text{global}} = \text{compose}(T_{\text{local}}, T)$ $\triangleright SE(3)^{L \times 14}$
 - 9: $\vec{X} = T_{\text{global}}(\vec{X}_{ref})$ $\triangleright \mathbb{R}^{L \times 14 \times 3}$
-



立志成才报国裕民

ESM3: Structure Tokenization Decode



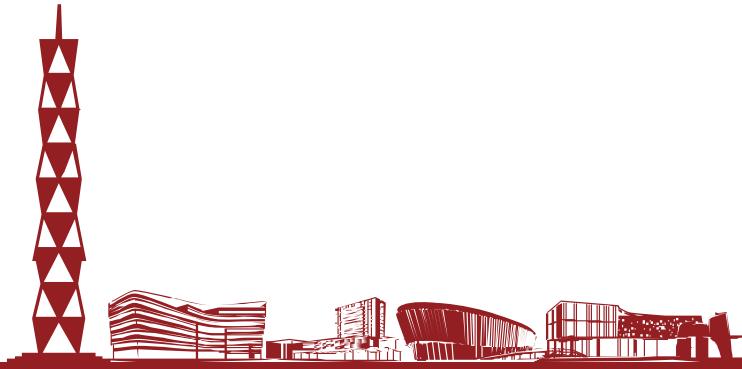
上海科技大学
ShanghaiTech University

- Backbone coords to frame
- Torsion angle to local frame
- Use backbone frame to global frame and rotate ref positions

Algorithm 8 structure_decode

Input: $z \in \{0..4099\}^{L \times 16}$

- 1: $z = \text{embed}(z)$ $\triangleright \mathbb{R}^{L \times d}$
 - 2: $z = f_{dec}(z)$ $\triangleright \mathbb{R}^{L \times d}$
 - 3: $\vec{t}, \vec{x}, \vec{y}, \sin \theta, \cos \theta = \text{proj}(z)$ $\triangleright (\mathbb{R}^{L \times 3})_{\times 3}, (\mathbb{R}^{L \times 7})_{\times 2}$
 - 4: $T = \text{gram_schmidt}(\vec{t}, -\vec{x}, \vec{y})$ $\triangleright SE(3)^L$
 - 5: $\sin \theta = \frac{\sin \theta}{\sqrt{\sin \theta^2 + \cos \theta^2}}$ $\triangleright [-1, 1]^{L \times 7}$
 - 6: $\cos \theta = \frac{\cos \theta}{\sqrt{\sin \theta^2 + \cos \theta^2}}$ $\triangleright [-1, 1]^{L \times 7}$
 - 7: $T_{\text{local}} = \text{rot_frames}(\sin \theta, \cos \theta)$ $\triangleright SE(3)^{L \times 7}$
 - 8: $T_{\text{global}} = \text{compose}(T_{\text{local}}, T)$ $\triangleright SE(3)^{L \times 14}$
 - 9: $\vec{X} = T_{\text{global}}(\vec{X}_{ref})$ $\triangleright \mathbb{R}^{L \times 14 \times 3}$
-



立志成才报国裕民

ESM3: Structure Tokenization Train



上海科技大学
ShanghaiTech University

- Two stage
- Stage I
 - Encoder and codebook are learnt with small decoder
 - 5 losses
 - Backbone distance loss, Backbone direction loss
 - Binned dir classification loss, Distogram loss, Inverse folding loss



立志成才报国裕民

ESM3: Structure Tokenization Train



上海科技大学
ShanghaiTech University

- Stage II
 - Encoder and codebook are frozen
 - Larger decoder applied
 - Sequence embedding added at input to decoder
 - All atom distance loss and All atom direction loss
 - pLDDT head and pAE head



立志成才报国裕民

ESM3: Function Tokenization



上海科技大学
ShanghaiTech University

- Find frequency of keywords for annotations -> TF-IDF vector
- Max pool vector to create vector per residue
- Locality sensitive hashes quantize vector per residue -> 8 tokens, each in range 0,...,255
- Still cannot figure out, code inspection



立志成才报国裕民

ESM3: Tokenization



上海科技大学
ShanghaiTech University

- Tokens:
 - Sequence: $\{0 \dots 28\}^L$, SS: $\{0 \dots 10\}^L$
 - SASA: $\{0 \dots 17\}^L$, Residue annotation: $\{0,1\}^{L \times 1478}$
 - Structure: $\{0 \dots 4099\}^L$, func: $\{0 \dots 258\}^{L \times 8}$
- Before transformer, they are embedded to $[L, d_{model}]$ and fused by summing with each other



立志成才报国裕民

ESM3: Transformer



上海科技大学
ShanghaiTech University

- Transformer applies to the embedded and summed tokens
- Attention in first block is changed to Geometric attention



立志成才报国裕民

ESM3: Decoder



上海科技大学
ShanghaiTech University

- Apply regression head to get the predicted logit and token for each track
- For sequence, ss, sasa, residue annotation, no further decode is needed



立志成才报国裕民

ESM3: Decoder



上海科技大学
ShanghaiTech University

- For structure, take the token after regression head and use the decoder trained in VQVAE
- For function annotation, learn an offline 3-layer transformer as the inverse of tokenization process



立志成才报国裕民

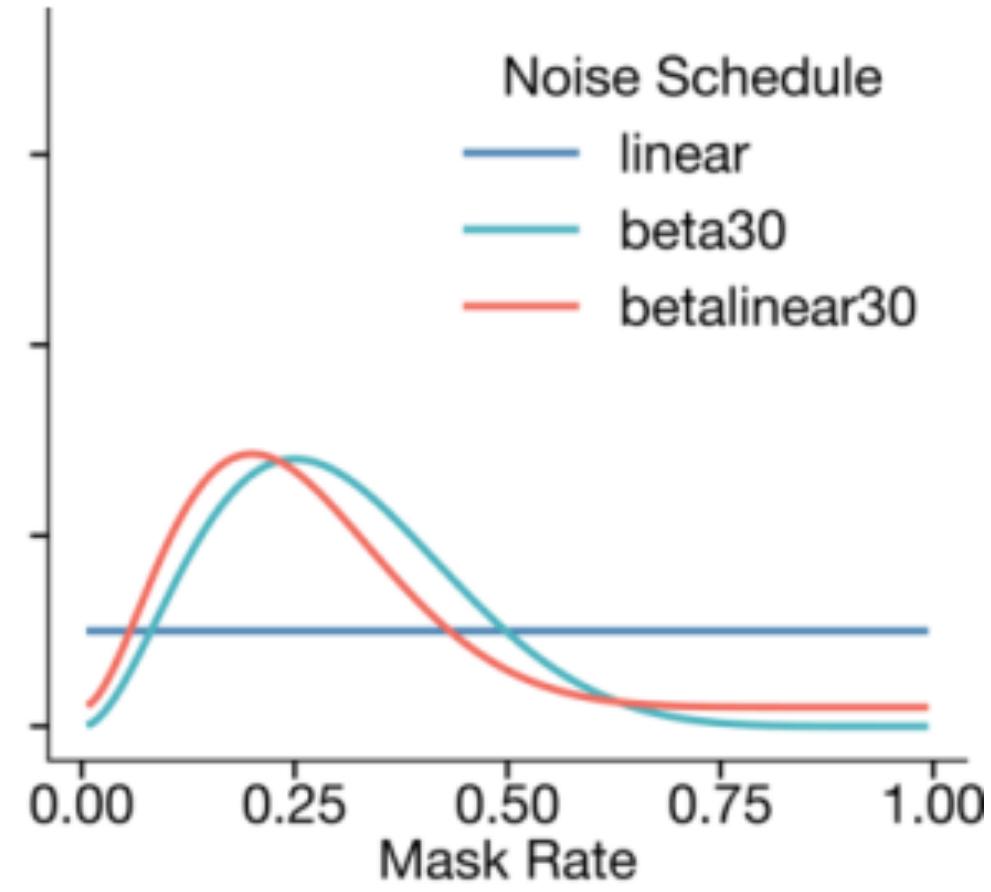
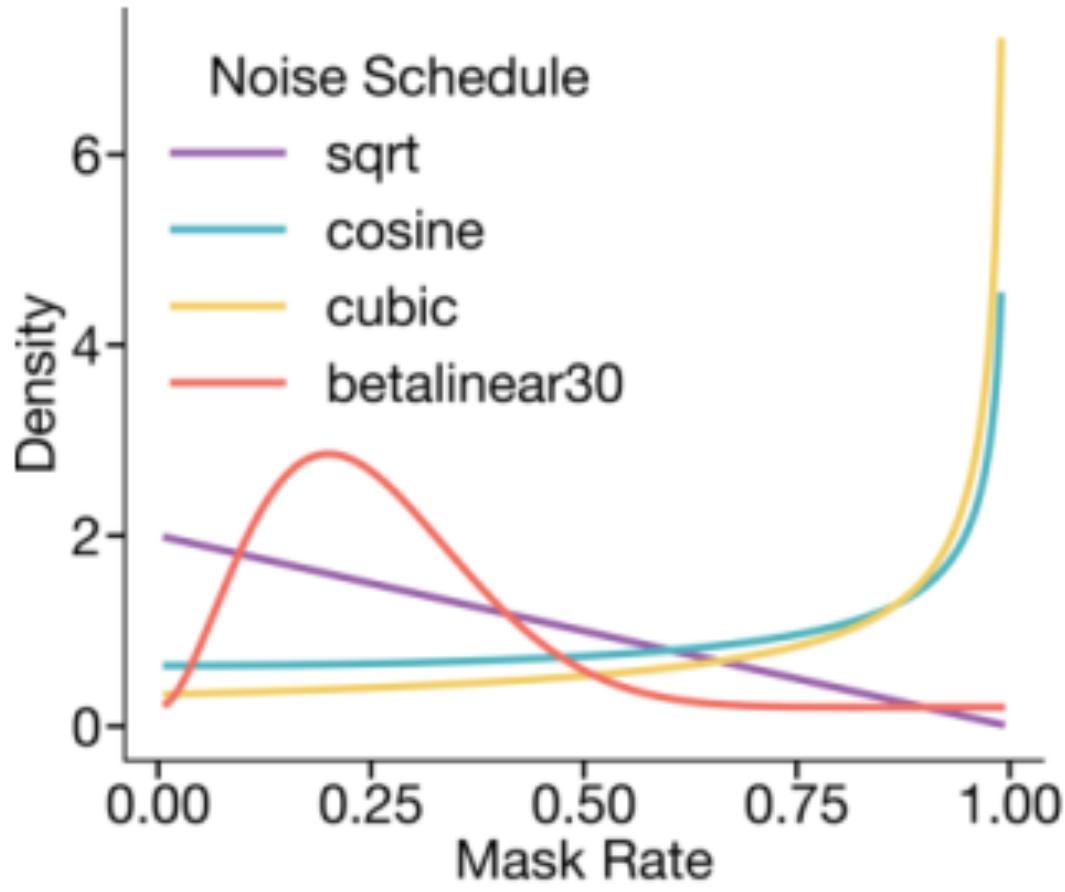


上海科技大学
ShanghaiTech University

Thanks for listening



立志成才报国裕民





Algorithm 22 Invariant point attention (IPA)

```
def InvariantPointAttention({ $\mathbf{s}_i$ }, { $\mathbf{z}_{ij}$ }, { $T_i$ },  $N_{\text{head}} = 12, c = 16, N_{\text{query points}} = 4, N_{\text{point values}} = 8$ ) :  
    1:  $\mathbf{q}_i^h, \mathbf{k}_i^h, \mathbf{v}_i^h = \text{LinearNoBias}(\mathbf{s}_i)$   $\mathbf{q}_i^h, \mathbf{k}_i^h, \mathbf{v}_i^h \in \mathbb{R}^c, h \in \{1, \dots, N_{\text{head}}\}$   
    2:  $\vec{\mathbf{q}}_i^{hp}, \vec{\mathbf{k}}_i^{hp} \in \mathbb{R}^3, p \in \{1, \dots, N_{\text{query points}}\}$ , units: nanometres  
    3:  $\vec{\mathbf{v}}_i^{hp} = \text{LinearNoBias}(\mathbf{s}_i)$   $\vec{\mathbf{v}}_i^{hp} \in \mathbb{R}^3, p \in \{1, \dots, N_{\text{point values}}\}$ , units: nanometres  
    4:  $b_{ij}^h = \text{LinearNoBias}(\mathbf{z}_{ij})$   
    5:  $w_C = \sqrt{\frac{2}{9N_{\text{query points}}}},$   
    6:  $w_L = \sqrt{\frac{1}{3}}$   
    7:  $a_{ij}^h = \text{softmax}_j \left( w_L \left( \frac{1}{\sqrt{c}} \mathbf{q}_i^{h\top} \mathbf{k}_j^h + b_{ij}^h - \frac{\gamma^h w_C}{2} \sum_p \|T_i \circ \vec{\mathbf{q}}_i^{hp} - T_j \circ \vec{\mathbf{k}}_j^{hp}\|^2 \right) \right)$   
    8:  $\tilde{\mathbf{o}}_i^h = \sum_j a_{ij}^h \mathbf{z}_{ij}$   
    9:  $\mathbf{o}_i^h = \sum_j a_{ij}^h \mathbf{v}_j^h$   
    10:  $\vec{\mathbf{o}}_i^{hp} = T_i^{-1} \circ \sum_j a_{ij}^h (T_j \circ \vec{\mathbf{v}}_j^{hp})$   
    11:  $\tilde{\mathbf{s}}_i = \text{Linear} \left( \text{concat}_{h,p}(\tilde{\mathbf{o}}_i^h, \mathbf{o}_i^h, \vec{\mathbf{o}}_i^{hp}, \|\vec{\mathbf{o}}_i^{hp}\|) \right)$   
    12: return { $\tilde{\mathbf{s}}_i$ }
```

Algorithm 22 Invariant point attention (IPA)

```
def InvariantPointAttention({ $\mathbf{s}_i$ }, { $T_i$ },  $N_{\text{head}} = 12, c = 16, N_{\text{query points}} = 4, N_{\text{point values}} = 8$ ) :  
    1:  $\mathbf{q}_i^h, \mathbf{k}_i^h, \mathbf{v}_i^h = \text{LinearNoBias}(\mathbf{s}_i)$   $\mathbf{q}_i^h, \mathbf{k}_i^h, \mathbf{v}_i^h \in \mathbb{R}^c, h \in \{1, \dots, N_{\text{head}}\}$   
    2:  $\vec{\mathbf{q}}_i^{hp}, \vec{\mathbf{k}}_i^{hp} \in \mathbb{R}^3, p \in \{1, \dots, N_{\text{query points}}\}$ , units: nanometres  
    3:  $\vec{\mathbf{v}}_i^{hp} = \text{LinearNoBias}(\mathbf{s}_i)$   $\vec{\mathbf{v}}_i^{hp} \in \mathbb{R}^3, p \in \{1, \dots, N_{\text{point values}}\}$ , units: nanometres  
    5:  $w_C = \sqrt{\frac{2}{9N_{\text{query points}}}},$   
    6:  $w_L = \sqrt{\frac{1}{3}}$   
    7:  $a_{ij}^h = \text{softmax}_j \left( w_L \left( \frac{1}{\sqrt{c}} \mathbf{q}_i^{h\top} \mathbf{k}_j^h - \frac{\gamma^h w_C}{2} \sum_p \|T_i \circ \vec{\mathbf{q}}_i^{hp} - T_j \circ \vec{\mathbf{k}}_j^{hp}\|^2 \right) \right)$   
    9:  $\mathbf{o}_i^h = \sum_j a_{ij}^h \mathbf{v}_j^h$   
    10:  $\vec{\mathbf{o}}_i^{hp} = T_i^{-1} \circ \sum_j a_{ij}^h (T_j \circ \vec{\mathbf{v}}_j^{hp})$   
    11:  $\tilde{\mathbf{s}}_i = \text{Linear} \left( \text{concat}_{h,p}(-\mathbf{o}_i^h, \vec{\mathbf{o}}_i^{hp}, \|\vec{\mathbf{o}}_i^{hp}\|) \right)$   
    12: return { $\tilde{\mathbf{s}}_i$ }
```

