

A Structure Encoder Module

```

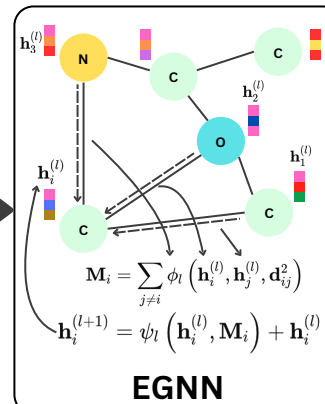
ATOM 1 N ASP A 2
-28.777 1.946 35.269
1.00197.02 N
ATOM 2 CA ASP A 2
-30.146 1.801 35.856
1.00197.61 C
ATOM 3 C ASP A 2
-31.246 1.547 34.797
...
  
```

PDB File

Graphin Library

**Chain-Specific
Filtering**

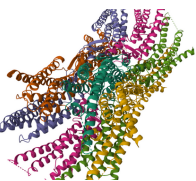
Protein Atom Graph



$e_{\text{pdb}}(1)$
 $e_{\text{pdb}}(2)$
 ...
 $e_{\text{pdb}}(512)$

**PDB
Embeddings**

B



Structure Encoder

C

**Structure
Projection Layer**

**Sequence
Projection Layer**

**Language
Projection Layer**

Loss function

$$\mathcal{L}_{\text{pdb, desc}} = \text{CE} \left(\text{softmax} \left(\frac{\mathbf{S}_{\text{pdb, desc}}}{\tau} \right), \mathbf{I} \right)$$

$$\mathcal{L}_{\text{desc, aas}} = \text{CE} \left(\text{softmax} \left(\frac{\mathbf{S}_{\text{desc, aas}}}{\tau} \right), \mathbf{I} \right)$$

$$\mathcal{L}_{\text{pdb, aas}} = \text{CE} \left(\text{softmax} \left(\frac{\mathbf{S}_{\text{pdb, aas}}}{\tau} \right), \mathbf{I} \right)$$

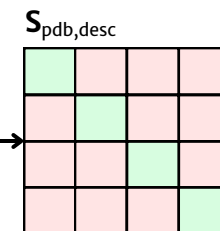
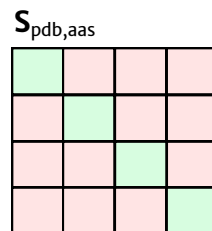
$$\mathcal{L}_{\text{total}} = \frac{1}{3} (\mathcal{L}_{\text{pdb, desc}} + \mathcal{L}_{\text{desc, aas}} + \mathcal{L}_{\text{pdb, aas}}) + \lambda \|\theta_{\text{enc}}\|_2^2$$

Training time

Alignment Module

Projected Structure Embeddings

$x(1)$ $x(2)$ $x(3)$... $x(512)$

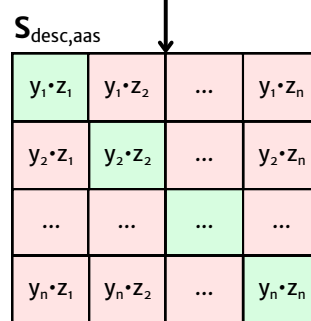


Projected Sequence Embeddings

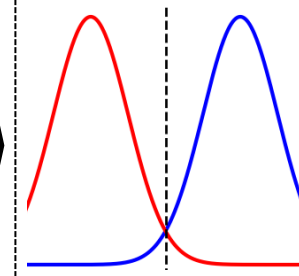
$y(1)$ $y(2)$ $y(3)$... $y(512)$

Projected Language Embeddings

$z(1)$ $z(2)$ $z(3)$... $z(512)$



Thresholding



Inference time

PROTEIN_NAMES: Gap junction beta-2 protein, Connexin-26.
 PROTEIN_DESCRIPTOR: Structural component of gap junctions. Gap junctions are dodecameric channels that connect the cytoplasm of adjoining cells ...

Language

BioGPT