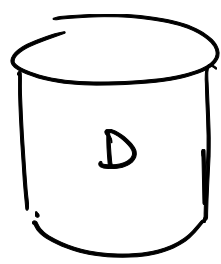
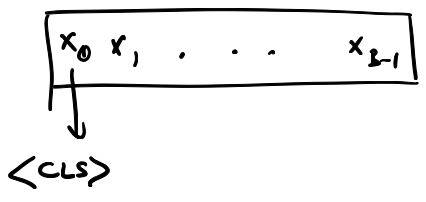


# BERT Contextual Model



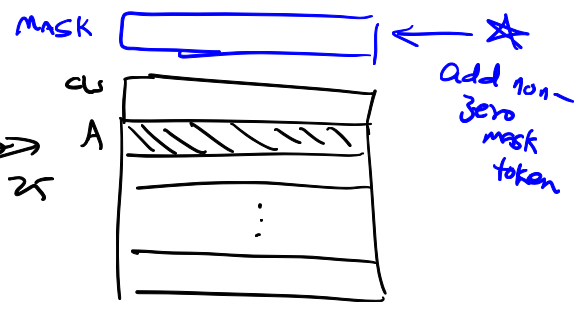
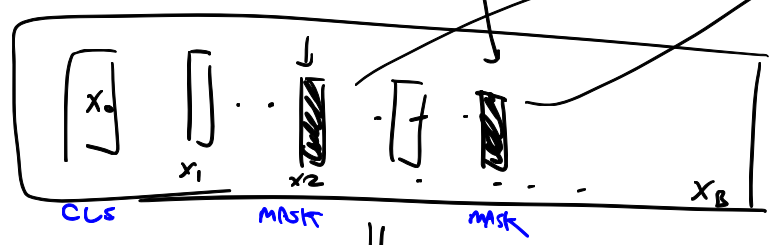
$S_1$   
 $S_2$   
 $\vdots$   
 $S_N$

$N \approx 500k$   
 $|S_i| \leq 1000$   
 $k\text{-mer} = 1$



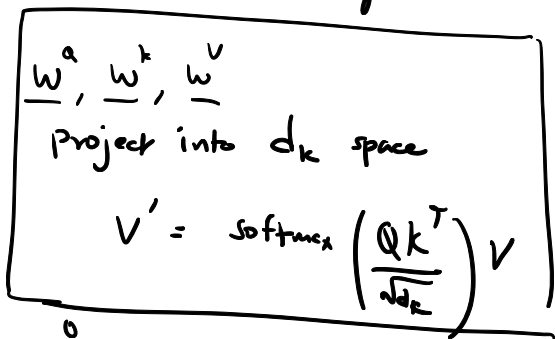
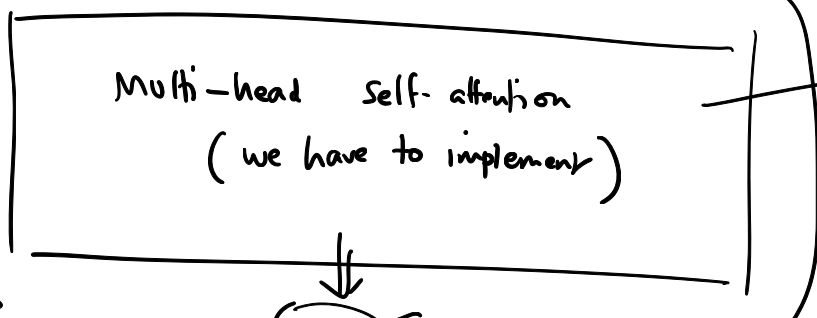
$B = \text{block size}$   
 $B = 1000$

M A I G L A . . .



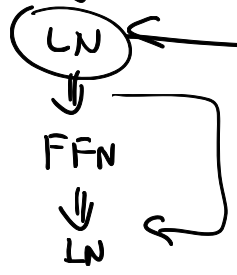
$d_{\text{model}} = ? 128$

initial token embeddings



$\underline{w^0}$  ← project back into  $d_{\text{model}}$

$h$ -heads  
Concat  
Use  $w^0$



$d_{\text{model}}$

$$d_k = d_q = d_v$$

Transformer Block:

```

--init-- ()
  mha = MHSA_net ( )
  ln1 = LN ( )
  ffn = FFN ( . . . )
  ln2 = LN ( )
--forward-- ( x )
  x' = mha ( x )
  x'' = ln1 ( x + x' )
  x''' = ffn ( x'' )
  x'''' = ln2 ( x''' + x'' )
  return x''''

```

MHSA\_net ( )

```

--init--
  for h = 0 ... 7
    A[h] = Attn_net ( )
  W0 = linear ( h x d_k,
               d_model )
--forward-- ( x )
  for h = 0 ... 7
    V[h] = Attn_net ( x,
                     A[h] )
  V = Concat ( V[0] ... V[7] )
  x = W0 ( V )

```

Transformer ( #layers )

$$\# \text{layers} = 1$$

```

for l = 0 ... #layers
  TB[l] = Transformer Block ( . . . )

```

Training Code

One sequence

```

X = < x1, x2, ... xL >
M = random binary mask
Y = Transformer ( X )

```



```

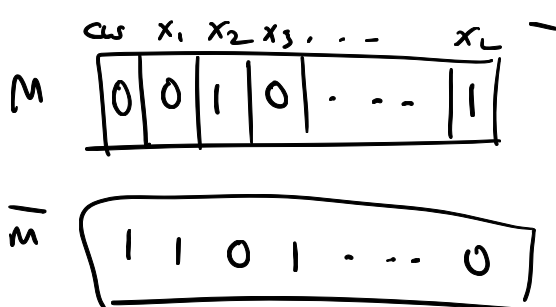
Define ProbBERT:
--init--
  Transformer ( )
  Linear ( 128, 25 )
  ...

```



Contextualized embeddings

p = prob of masking = 0.15



$$x' = X \odot \bar{M}$$

input  $\rightarrow$

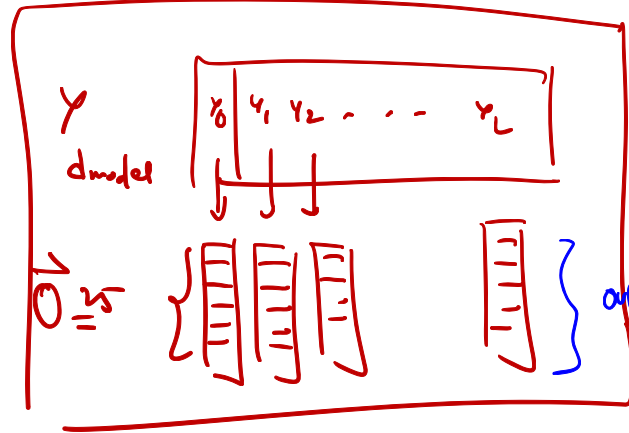
Instead of zeroing out masked positions, replace with <mask> token  
 $x_2, x_L$

$$\vec{O} = \text{Transformer}(X')$$

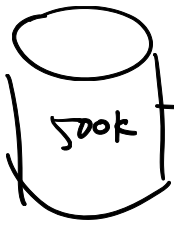
$(25 \times L)$  matrix

$\vec{O}' =$  remove the  $\phi$  positions from  $\vec{O}$   
 $25 \times \underline{L}$

loss = CE - with - logits ( $\vec{O}'$ , True label)  
 $[x_2, x_L]$



output should always be # of Amino acids (25)  
 do not include <CLS>, <MASK>

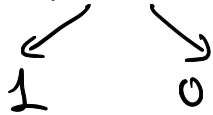


train BERT

pre-trained

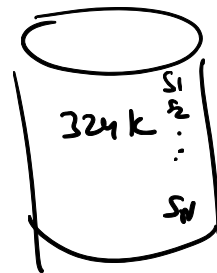
down stream task

protein family classification



binary logistic regression  
 MLP

SOS ribosomal GTPase  
 3084 instances  
 (+1)

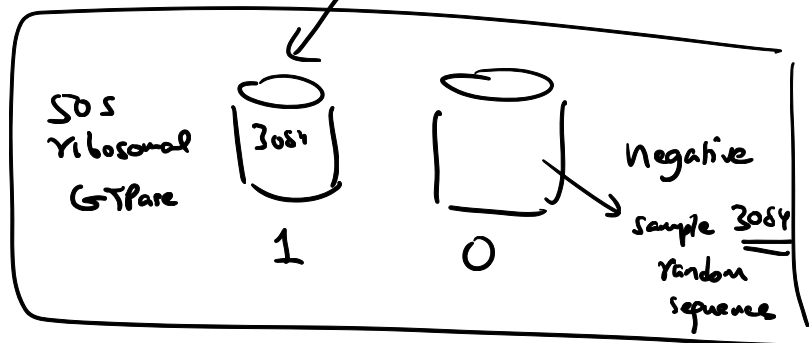


metadata



family name

SOS DB



$X =$  SOS DB  $\rightarrow$  3084 x 2 sequences

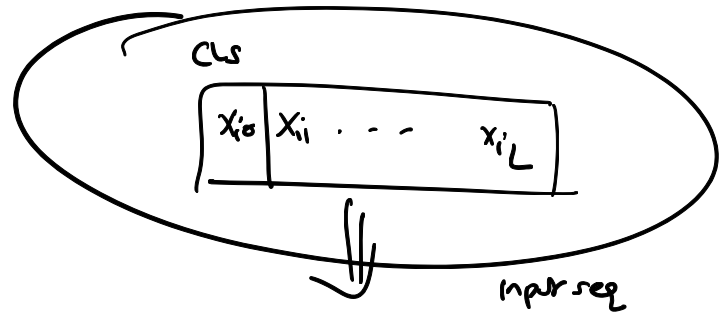
0/1 labels

for  $i = 1 \dots 6178$

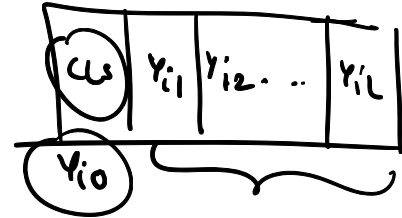
$$Y_i = \text{Transformer}(X_i)$$

$$Z_i = \text{MLP}(Y_i)$$

$$\text{BCE loss}(Z_i, \text{true})$$



Transformer

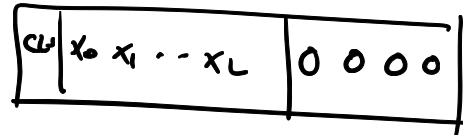


$d_{\text{model}} = 128 \text{ dim}$

Zero-pad all blocks for transformer

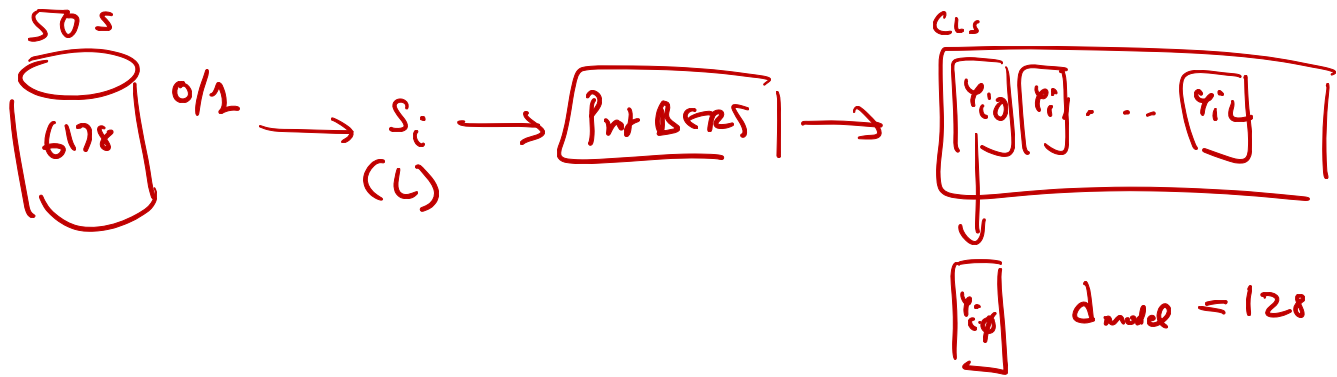
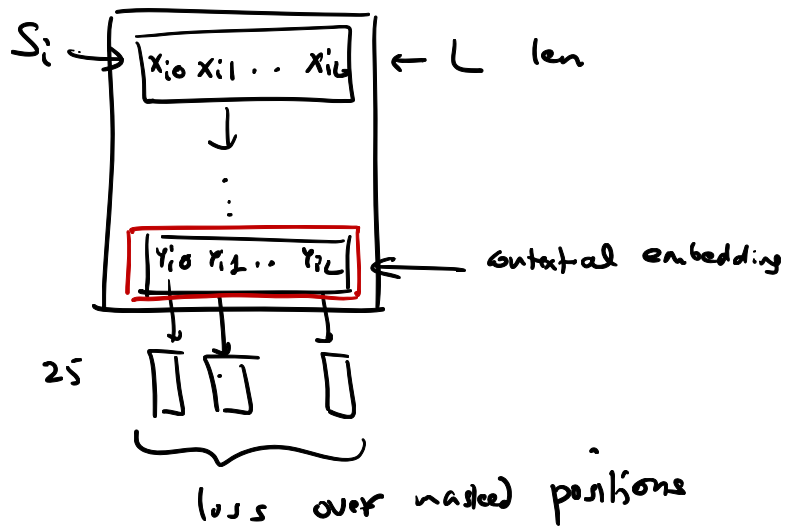
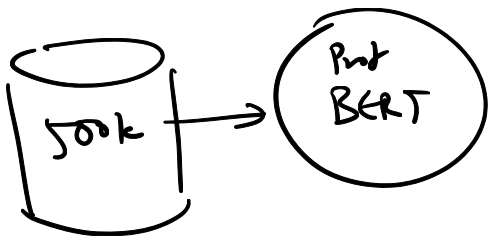
SOS DB

if a seq > 999  
look at 1<sup>st</sup> 999  
elements



$L = 300$

You can use <PAD> token and init <PAD> to zero vector in embedding layer

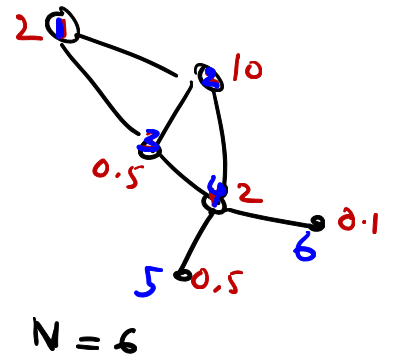


## Graph Neural Networks

$$L = D - A \quad \text{matrix}$$

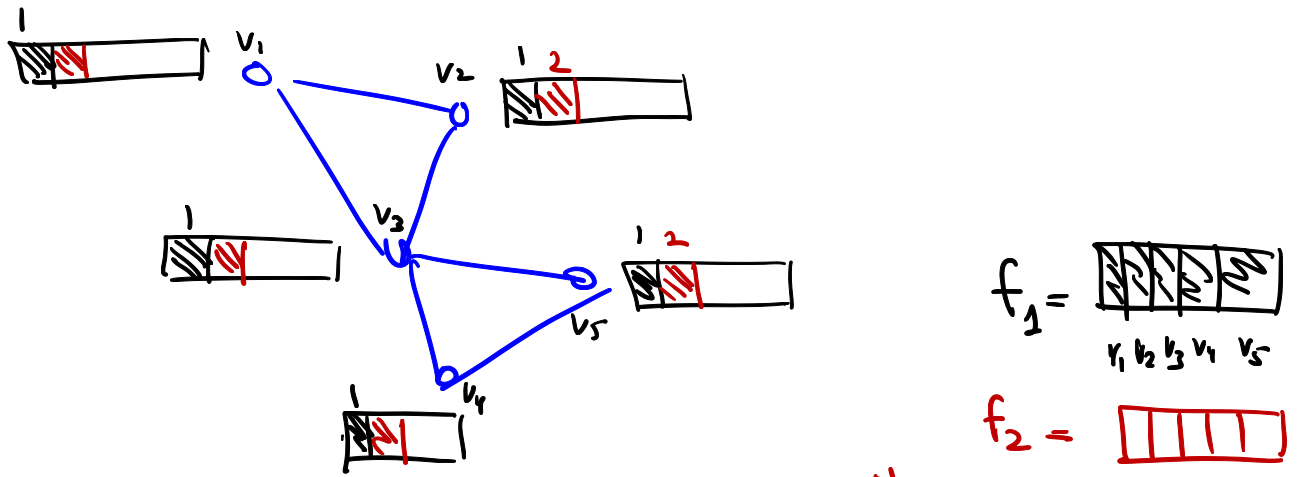
$$f_1: \text{feature vector} \in \mathbb{R}^N$$

$$f_1 = [2, 10, 0.5, 2, 0.5, 0.1]^T$$

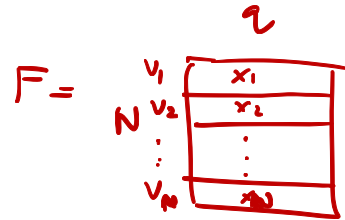
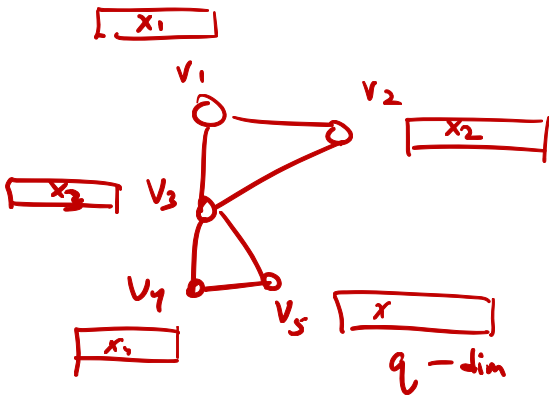
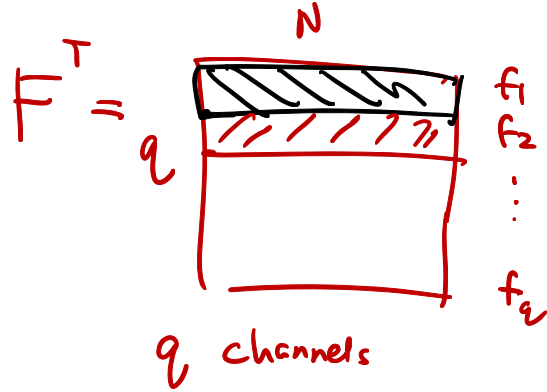


$$\text{Smoothness} = f_1^T L f_1 = \sum_i \sum_{j \in N(i)} g_{ij} (f_1^{[i]} - f_1^{[j]})^2$$

↑  
filter we need to learn



GNN



$x_i$ : feature vector for node  $v_i$   
 $q$ -dim

Each layer update the node feature

- a) spatial update / neighborhood-based (1-hop)
- b) spectral update

GCN: graph convolutional network

$$x'_i = \sum_{j \in N(x_i)} \frac{1}{\sqrt{d_i} \sqrt{d_j}} \vec{x}_j \quad W_{ij} \quad N(x_i) = N(x_i) \cup \{x_i\}$$

