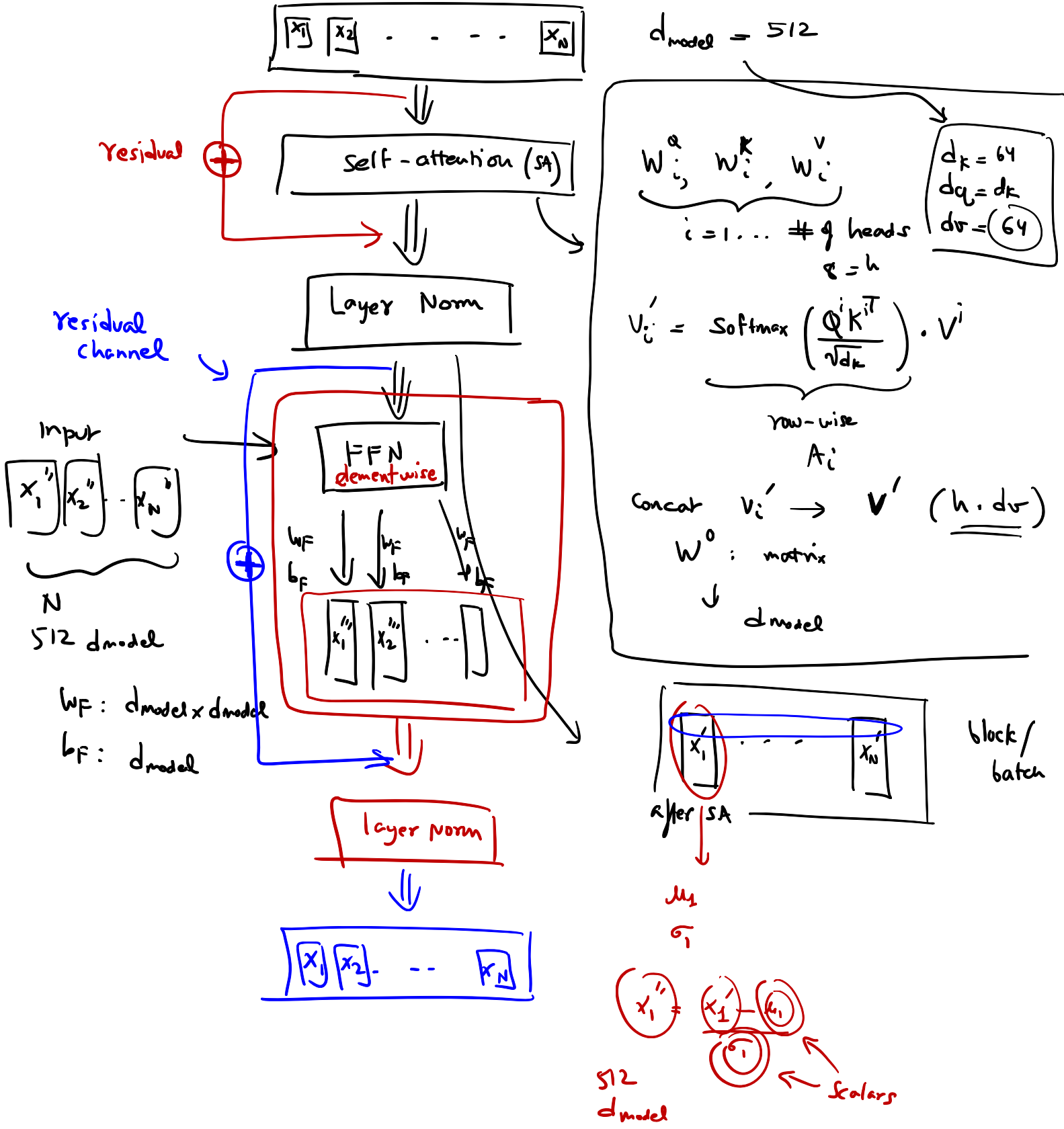
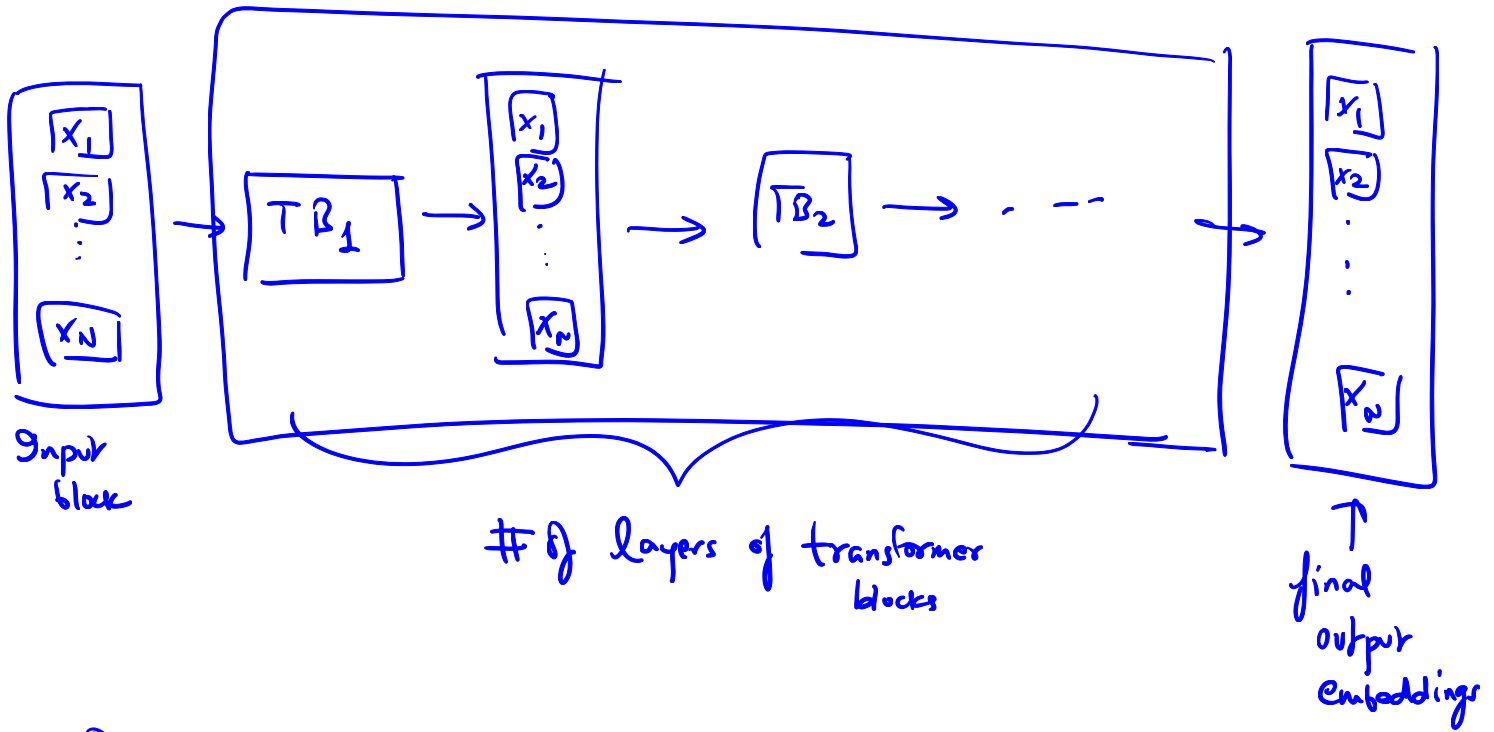
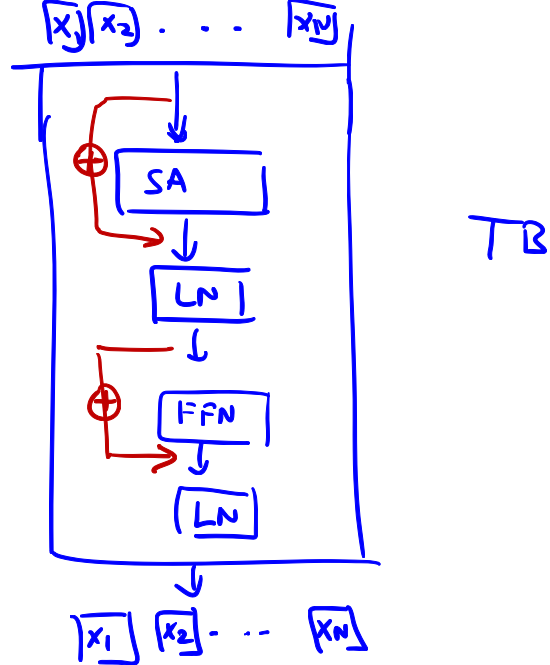


Transformer

self-attention : multi-headed



Transformer block
 $d_{model} = 512$



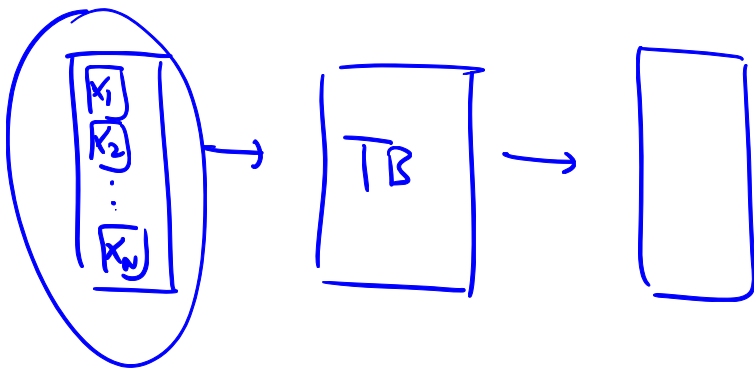
Pretraining

→ Word2vec : 1 token/word → 1 embedding

I was sitting by the river bank →
 I went to the bank and sat on the bench →

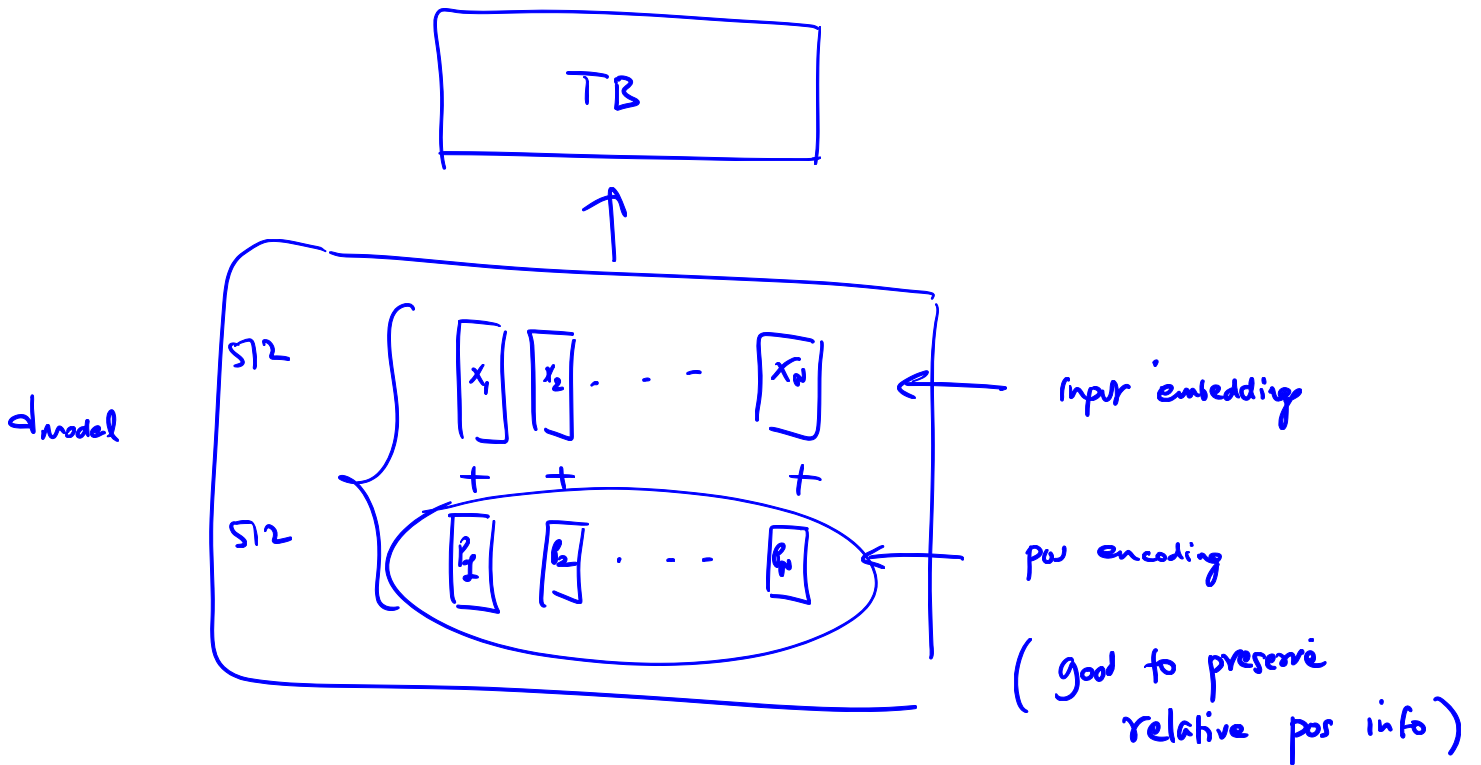
→ transformer :

1 token / word → multiple embeddings
 (as many blocks / occurrences)



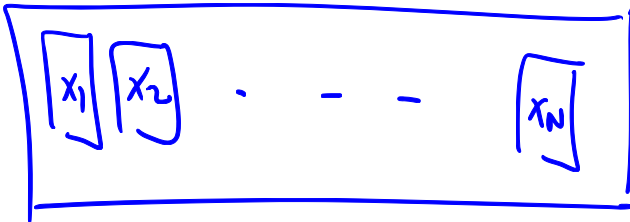
no order info

↪ Position Encoding

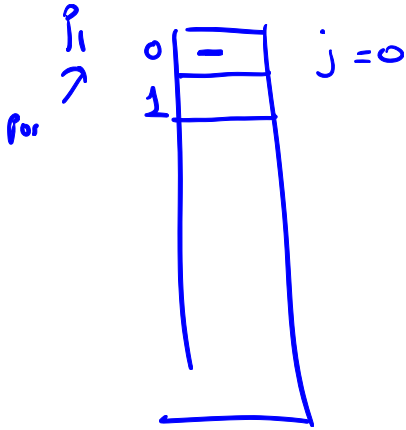


- Pos embedding:
- 1) Cos / sin based
 - 2) preserves relative pos info

Input block



size dim

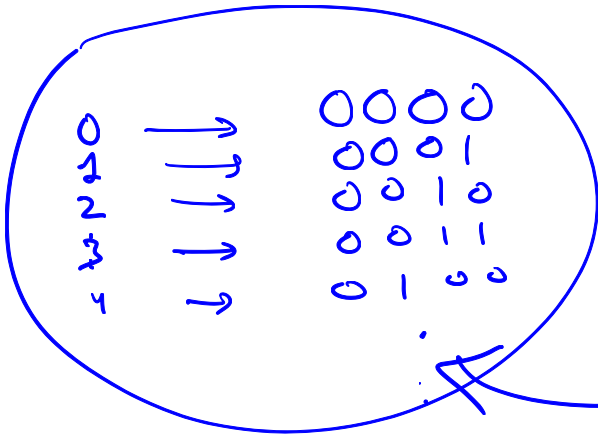


$$P_{pos, z_j} = \sin \left(\frac{pos}{(10000)^{z_j/d_{model}}} \right)$$

Labels: 'dim' with a downward arrow pointing to z_j ; 'pos=1' with an upward arrow pointing to P_{pos, z_j} .

$$P_{pos, z_{j+1}} = \cos \left(\frac{pos}{(10000)^{z_{j+1}/d_{model}}} \right)$$

Label: 'dim' with an upward arrow pointing to z_{j+1} .



binary encoding

Continuous version of binary encoding

BERT: pre-trained model to generate contextual embeddings

↳ base model

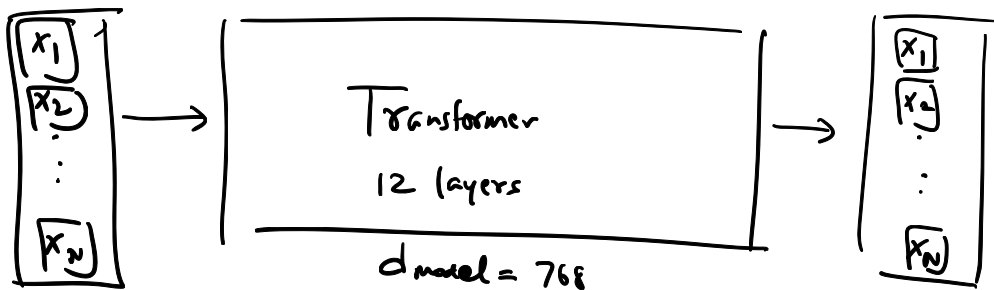
$$d_{\text{model}} = 768$$

$$\text{block size} = 512 = N \quad x_1 \dots x_N$$

$$\# \text{ of layers} = 12 \quad (\text{transformer layers})$$

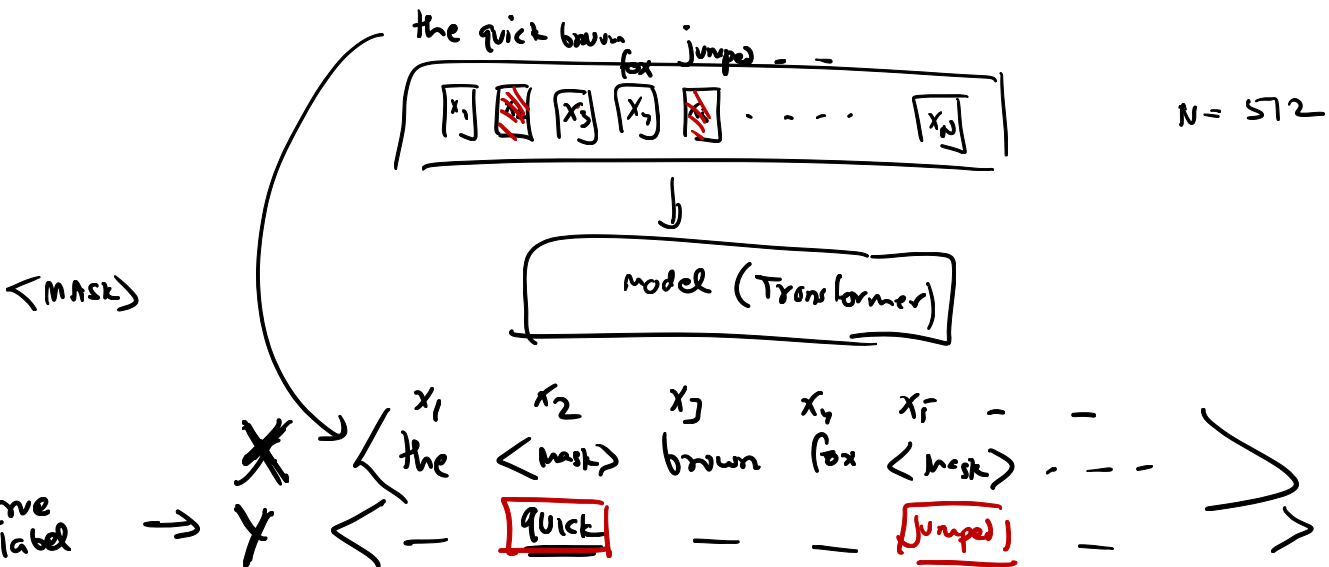
$$\# \text{ of heads} = \underline{\underline{8/12}} ?$$

$$\underbrace{d_k = d_q, d_v}_{64}$$

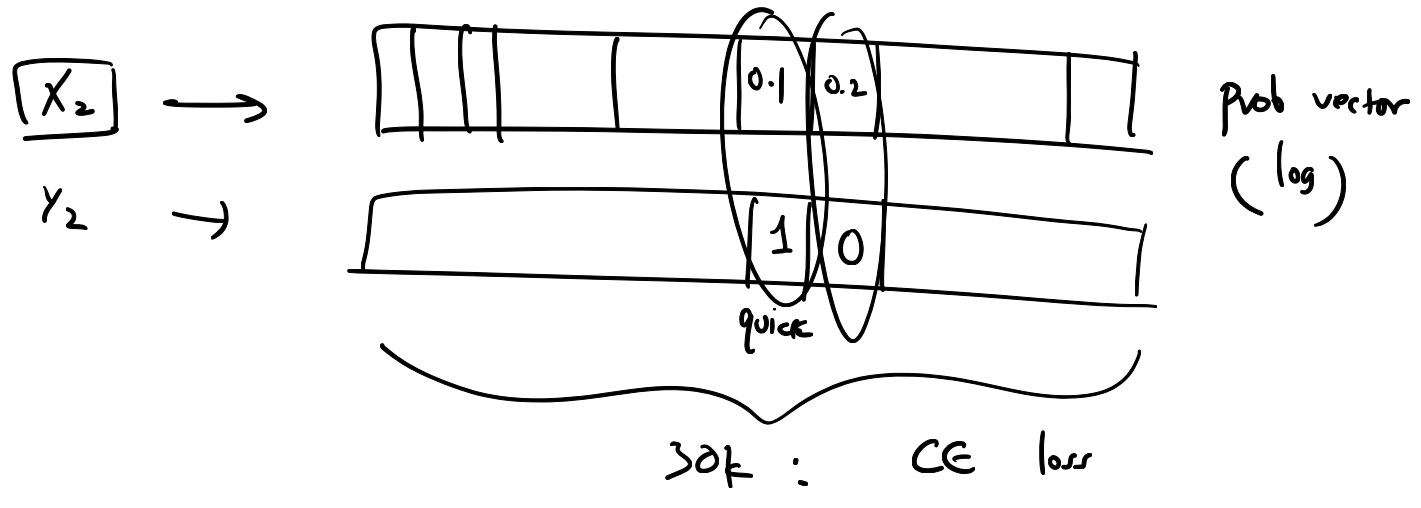
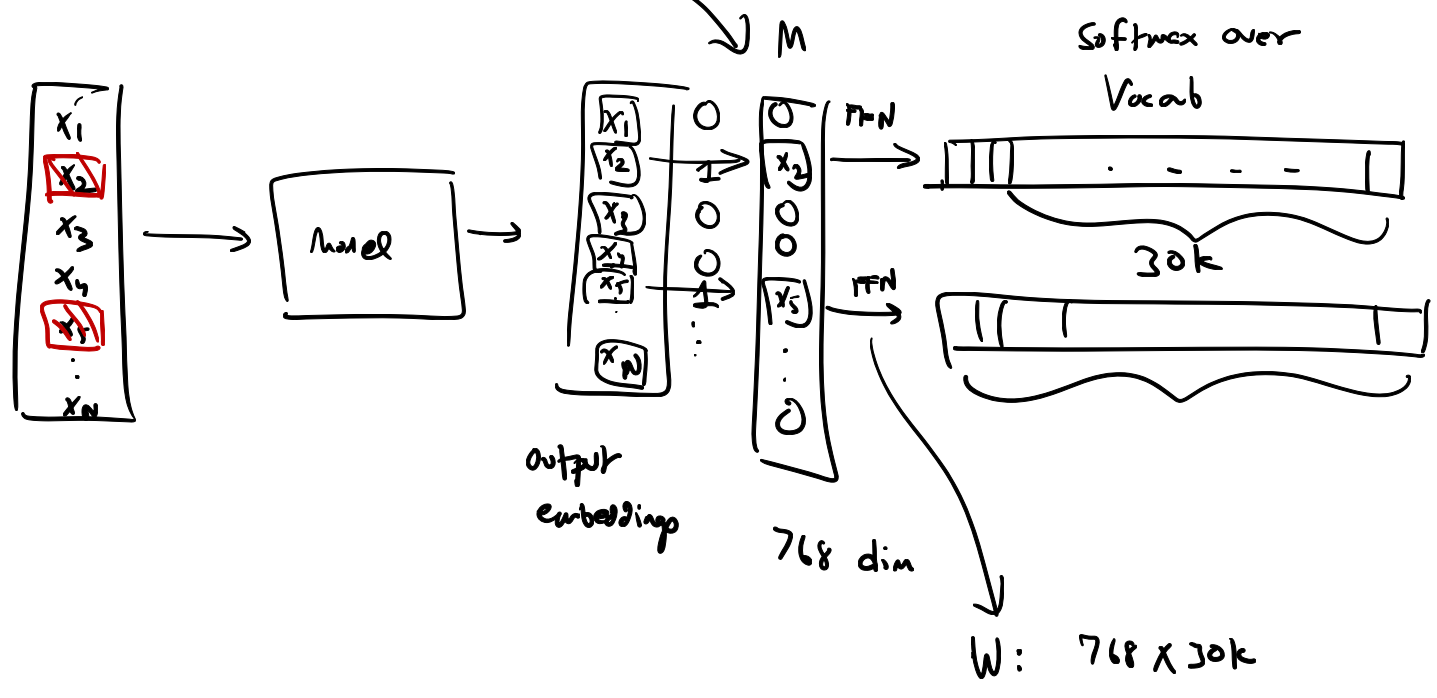


30k vocab
of tokens
(word-piece)

MLM: Masked Language Model
(self-supervised)



Mask (loss) $\rightarrow M < 0 \ 1 \ 0 \ 0 \ 1 \ 0 \dots >$



BERT :

15% of the positions are masked randomly

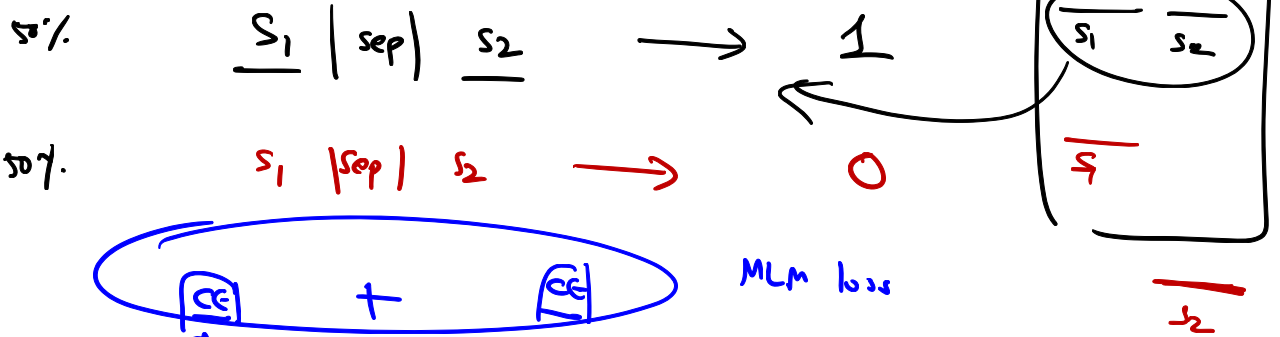
a) 80% replace the actual token with $\langle \text{mask} \rangle$

b) 10% replace with a random token
 the quick brown fox
 the bank brown fox

c) 10%. leave as is
 the quick brown fox

I) MLM \leftarrow loss CE over 30k vocab per masked position

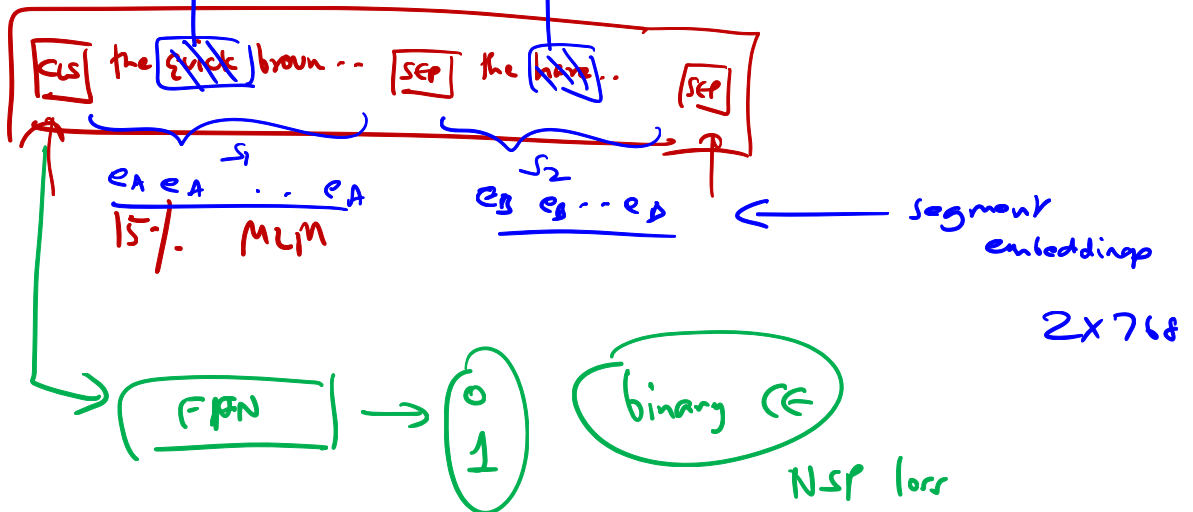
II) NSP \leftarrow next sentence prediction



MLM loss

CE + CE

BERT

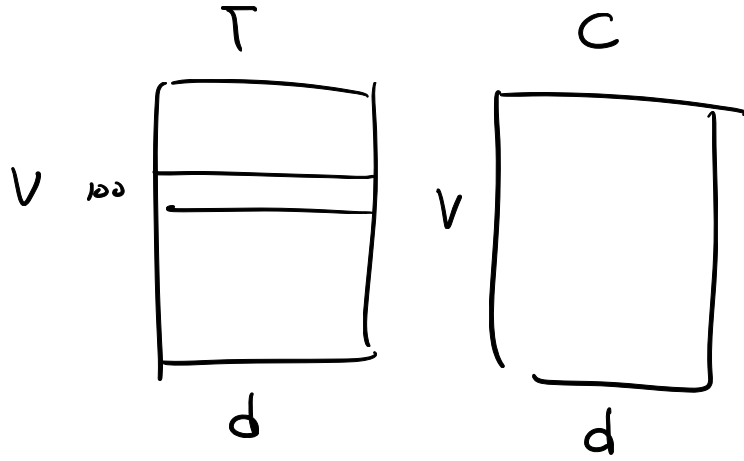


total loss $L = \text{MLM} + \text{NSP loss}$

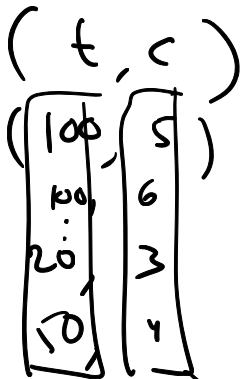
$X_i + P_i + e_A/e_B$

Model

nn. Embedding (linear / Dense)



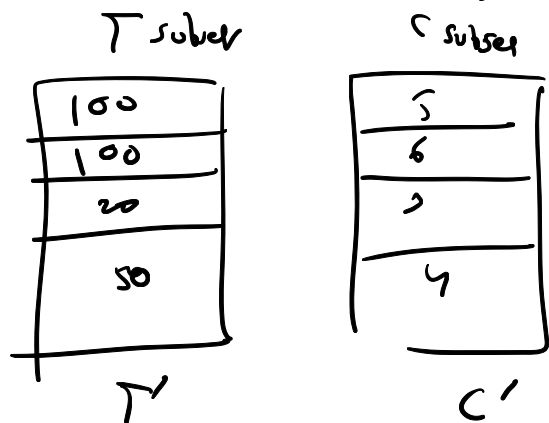
$$T = \text{nn. Embedding}(V, d)$$



$$T(100)$$

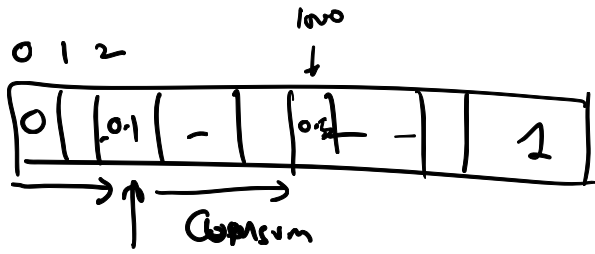
$$T([100, 100, 20, 50])$$

$$C([5, 6, 3, 4])$$

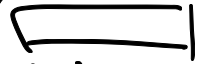


logits = torch.sum($T' * C'$, dim = 1)

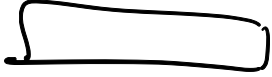
BCE with logits



$\gamma_{\text{random}} = 0.1 \text{ idx } 2$
 $0.5 \text{ idx } 1000$

NP. Random. Choices ( , 1000)
 prob vector

NP. Search sorted ()

$(\frac{0.1, 0.5, \dots}{1024})$ 
 Cumsum
 V