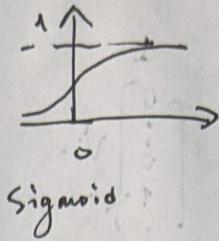
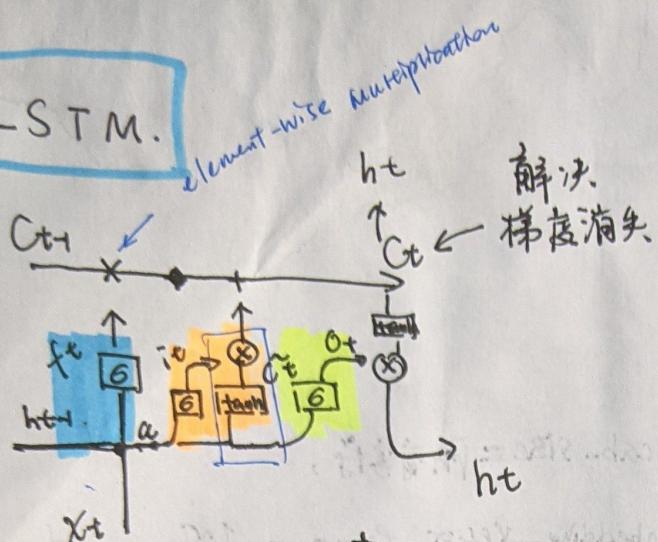


## LSTM.



## Deep Learning & NLP

$$h_t = \left[ \text{tanh} \left( W_f [h_{t-1}]_{x_t} \right) + \text{tanh} \left( W_i [h_{t-1}]_{x_t} \right) * \text{tanh} \left( W_c [h_{t-1}]_{x_t} \right) \right] + \text{tanh} \left( W_o [h_{t-1}]_{x_t} \right)$$

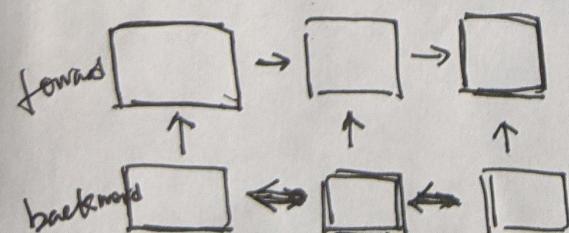
① Forget Gate    ② Input Gate    ③ new Value    ④ Output Gate

$$= h_t$$

$$W_f, W_i, W_c, W_o \rightarrow r^a \left[ \begin{matrix} h_{t-1} \\ x_t \end{matrix} \right]$$

± params:  $4 \times \text{Shape}(h) \times (\text{Shape}(h) + \text{Shape}(x))$

## Bidirectional LSTM.



## Word Embedding

$$d \quad \boxed{\quad} \text{词} \times \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix} v.$$

字典

(Vocab-size = 所有字符数,  $V$ )

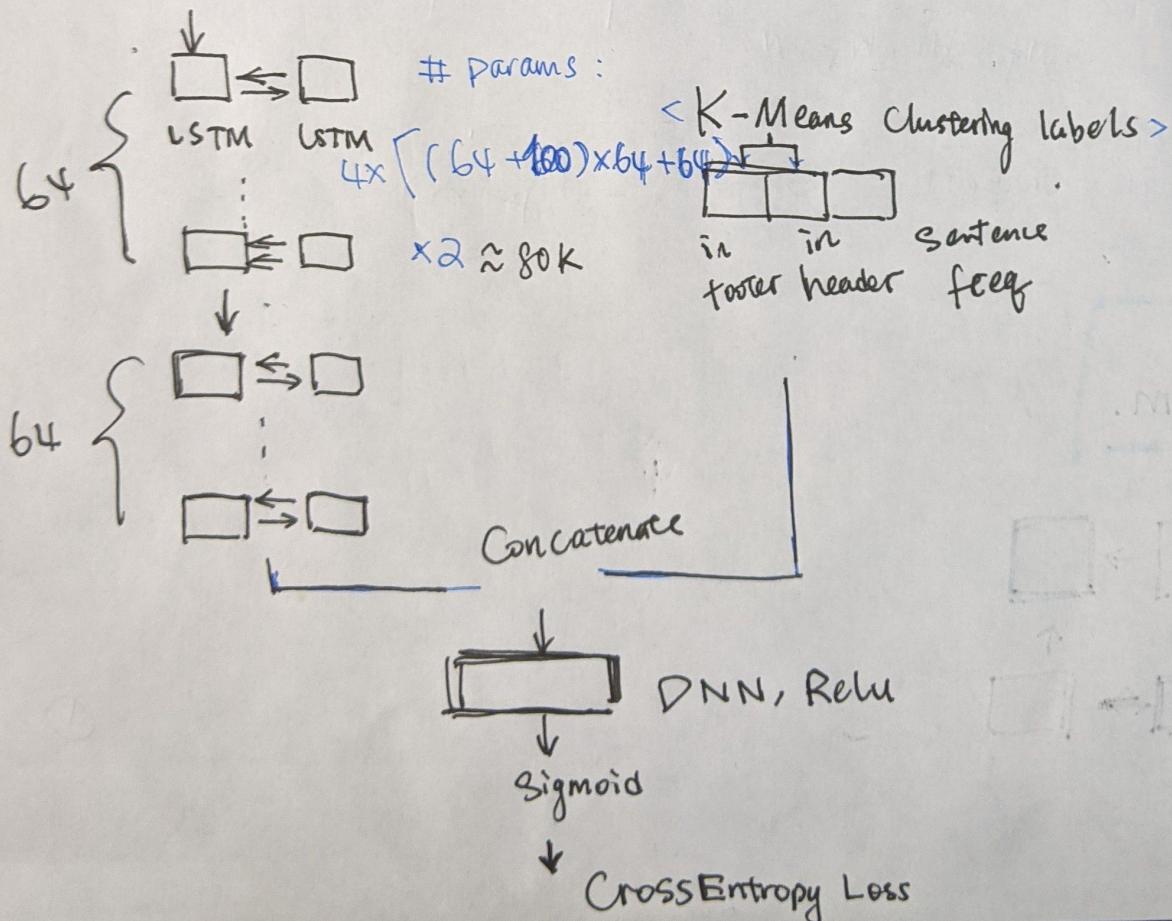
embedding-vector-length = 100

max-length = 50

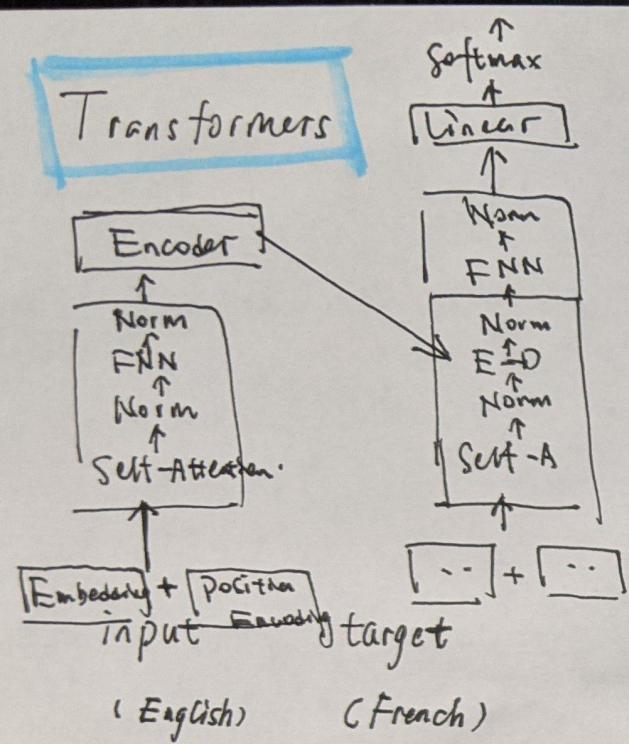
## Our model (DP Intern)

50

$$100 \quad \begin{bmatrix} w_1 & w_2 & w_3 & \dots & w_{50} \\ | & | & | & & | \end{bmatrix}$$



## Transformers



(English)

(French)

Tensorflow :

vs

Pytorch :

✓ Production (Industry)

✓ (Academy)

✓ More complex model.

$$\text{Self-Attention Model: } \begin{bmatrix} w_1 & w_2 & \dots \end{bmatrix} \begin{bmatrix} -w_1 \\ -w_2 \\ \vdots \end{bmatrix}$$

high score if inner prod. is high.

Transformer 模型，

Tensorflow vs. Pytorch.

# BERT

Predict ① masked word, ② next sentence

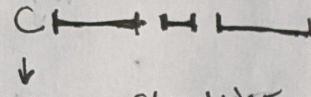
① 15% words  
randomly

Softmax + Cross Entropy Loss

判断是否是 masked word

② 50% true consecutive sentences

50% False (faked)

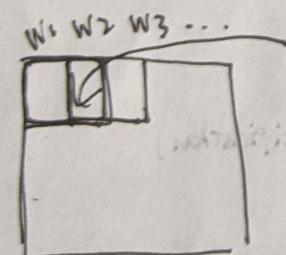


Binary Classifier, Cross Entropy Loss.

• 结构  $\rightarrow$  transformer  $\rightarrow$  Encoder

视觉 classification, 语义 Regression.

## Bert - KPE



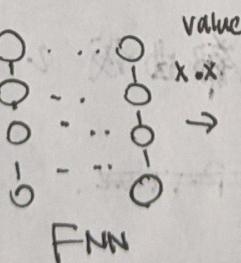
filters

CNN

matrix



$\rightarrow$



max pool +  
for various score

Loss:

Cross Entropy Loss

②