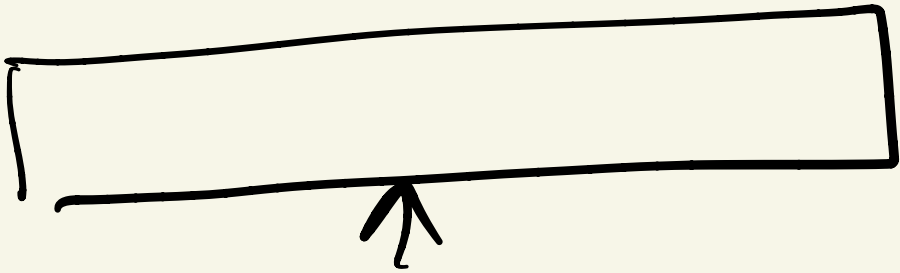


Suppose we have a
pre-trained decoder-only
model.

We have sum data of
the form (news art., sum)
How do we fine-tune
the model on this

(x, y)



$x_1 \dots x_N$

News
Article

$y_1 \dots y_M$

Summary

$$S(x, y) \prod_{i=1}^k$$

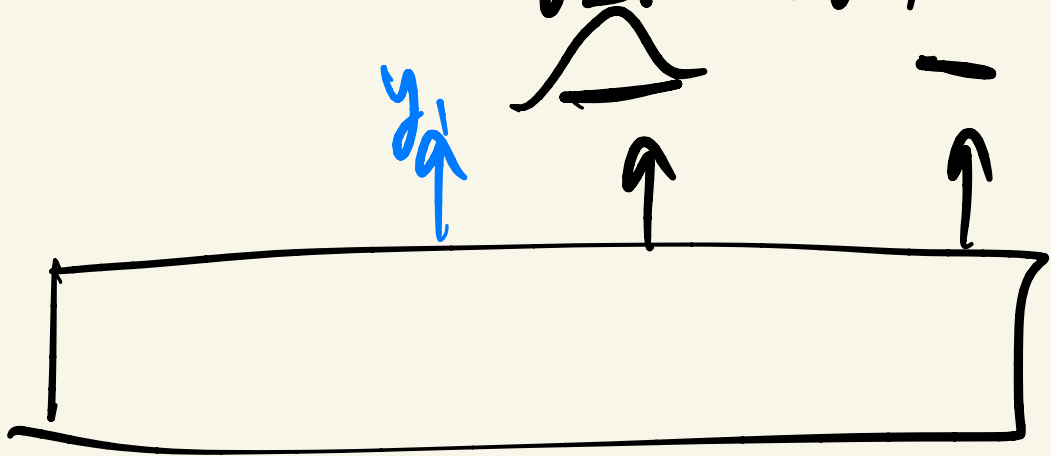
$$\max_{\theta} \prod_{i=1}^k P(y_i | x)$$

$$\max_{\theta} \sum_{i=1}^k \log P(y_i | x)$$

$$\min - \sum \log P(y_i | x)$$

loss

$$\sum_{i=1}^N \log P(y_i | x_{1:i-1}, y_{1:i-1})$$



$x_1 \dots x_N$ $y_1 \dots y_N$



$$\text{Loss} = - \sum_{j=1}^N \log P(y_j | x_{1:j-1}, y_{1:j-1})$$

Pre-training

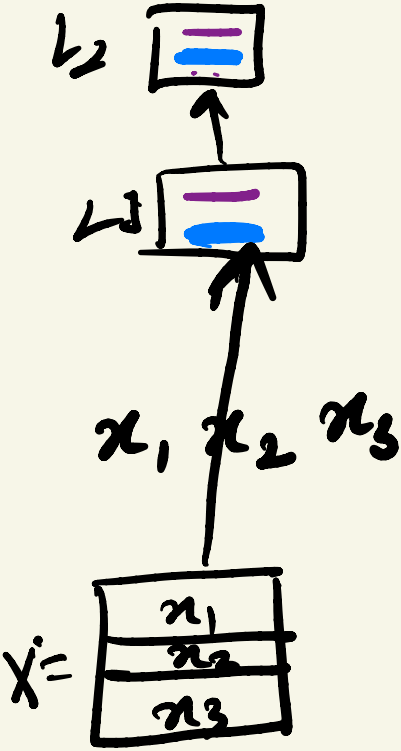
$$\max_{\theta} P(x)$$

$$= \max_{\theta} \prod P(x_i | \dots x_{i-1})$$



Transformer \rightarrow decoder-only
2 Layers

\rightarrow Attention



$$L1 \rightarrow X = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

→ Attention
 $W^Q \quad W^V \quad W^K$

$$Q = \begin{bmatrix} q_1 \\ q_2 \\ q_3 \end{bmatrix}$$

$$K = \begin{bmatrix} k_1 \\ k_2 \\ k_3 \end{bmatrix}$$

$$V = \begin{bmatrix} v_1 \\ v_2 \\ v_3 \end{bmatrix}$$

$$O = \text{Softmax}\left(\frac{QK^T}{\sqrt{d.}} + \text{mask}\right)$$

$$a = O \cdot V$$

$$\begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix}$$



Output of
Attention

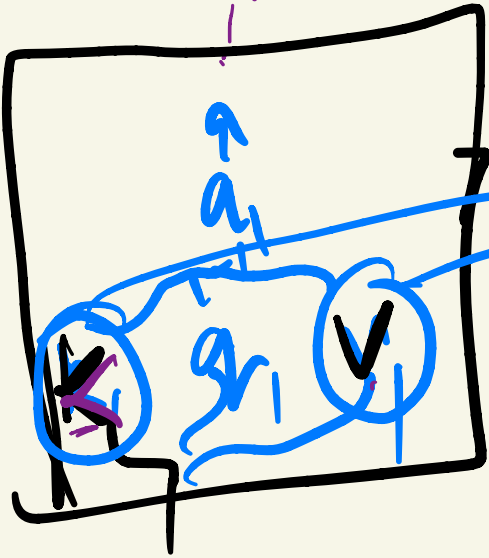
→ LN + Residual + FFN

Generation

x_1

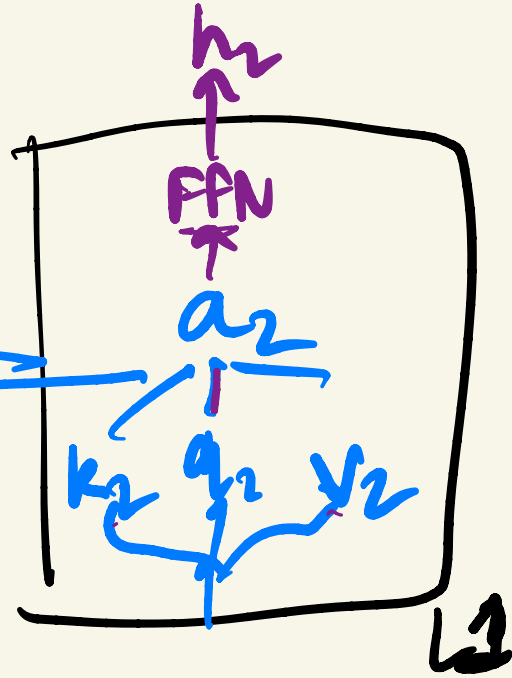


h_1



x_1

Generate
Cond. on x_1



x_2

KV - cache

k_1 v_1

for each

k_2 v_2

layer

,

for each

,

prev. timestep

,

