

CHARGE PT \rightarrow GPT

GPT \rightarrow corrective Protective transformers
 \downarrow
80%

CIVIL Engg.

Tokenization

Dataset $\rightarrow \{ \text{NLP is good}, \text{machine learning is good} \dots \}$ $\rightarrow \{ \text{0, 1, 2, } \dots \text{A, B, } \dots \text{1, 2, 3} \}$
[65]

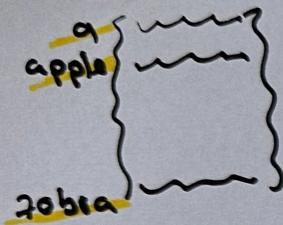
$\{ [39, 11, 15, 63, 8, 18], [63, 6, 14, 14, 3, 63] \}$

Batch size $\rightarrow (BS) = 2$
Block size $\rightarrow BL = 4$
Device = cuda | cpu | mps
embedding size $E = 4$

~~X~~ $\rightarrow [11, 15, 63, 8]$ [2, 4]
 $[63, 6, 14, 14]$

~~-1~~ $\rightarrow [15, 63, 8, 18]$ [2, 4]
 $[6, 14, 14, 3]$

embedding and Positional Encoding



$$\begin{matrix} 0.9 & \{0.1 & 0.2 & 0.3 & 0.4 \\ 1.5 & \{ & \dots & \} \\ \vdots & \dots & \dots & \} \\ 6.5 & \{ \end{matrix} \quad [65, 4]$$

$$x \rightarrow \begin{bmatrix} 11 & 45 & 63 & 8 \\ 63 & 6 & 14 & 14 \end{bmatrix} \rightarrow \left[\begin{bmatrix} 0.5 & 0.2 & -0.1 & 0.7 \\ \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix} \right] \rightarrow x [2, 4, 4]$$

$$x_n = x + p(x)$$

$$\begin{bmatrix}] \\ [2, n, n] \end{bmatrix} \begin{bmatrix}] \\ [n, 1, n] \end{bmatrix} \rightarrow \begin{bmatrix}] \\ [2, n, 1, n] \end{bmatrix}$$

Masked Multihead
 Attention $\frac{n}{2} = 2$
 +local
 $\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$

$$Q \quad x_n W_1^T \rightarrow [2, n, 2] \rightarrow [2, n, 2]$$

~~$$K \quad x_n W_2^T [2, n, 2] \rightarrow [2, n, 2]$$~~

~~$$V \quad x_n W_3^T [2, n, 2] \rightarrow [2, n, 2]$$~~

$$\text{softmax} \left[\frac{QK^T}{\sqrt{d_k}} \right] \rightarrow [2, n, 2] \rightarrow [2, n, 2]$$

$$\begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} \\ [3, 2]$$

$$\begin{bmatrix} \cdot & \cdot \\ \cdot & \cdot \end{bmatrix} \quad \begin{bmatrix} \cdot & \cdot \\ \cdot & \cdot \end{bmatrix} \\ [2, n, 2] \quad [2, n, 2]$$

$$[2, n, 2] \rightarrow O_1$$

$$H \quad \text{Head } 2 \\ Q \quad x_n W_4^T [2, n, 2] \rightarrow [2, n, 2]$$

$$K \quad x_n W_5^T [2, n, 2] \\ V \quad x_n W_6^T [2, n, 2]$$

$$[2, n, 2] \rightarrow O_2$$

$$[2, n, 2] \rightarrow x_0$$

$$x_0 W_7^T [n, 2] \rightarrow [2, n, 2] \rightarrow x_0$$

$$x_0 \rightarrow \text{layer norm}([x_0 + x_0]) \rightarrow [2, n, 2]$$

Forward -> Add -> Norm

$$x_s []_{[2, 4, 4]}$$

Fully connected 1

$$x_s w_1^T \rightarrow [4 \times 4]$$

$$[2, 4, 4] [4, 16]$$

$$\left[\text{ReLU} \left[[2, 4, 16] \right] \right] \rightarrow x_{s2}$$

$$x_{s2} w_2^T [4 \times 2]$$

$$[2, 4, 16] [16, 4]$$

$$[]_{[2, 4, 4]} \rightarrow y_{sco}$$

$$\ln [x_{s2} x_{suv}] \rightarrow x_o$$

$$\text{Attention}(Q, k, v) = \text{softmax}\left(\frac{Q \cdot k^T}{\sqrt{n}}\right)v$$

$$Q \rightarrow x_0 w_{10}^T \xrightarrow{H1} [2, n]$$

$$[2, m, n] \quad [n, 2]$$

$$[] \quad [2, n, 2]$$

$$k \rightarrow x_0 w_{11}^T \quad v \rightarrow x_0 w_{12}^T$$

$$[] \quad [2, n, 2] \rightarrow 0$$

$$Q \rightarrow x_0 w_{13}^T \xrightarrow{H2} []$$

$$K \rightarrow x_0 w_{14}^T \quad []$$

$$Q \rightarrow x_0 w_{15}^T \quad []$$

$$[] \quad [2, n, 2] \rightarrow 0$$

$$[] \quad \rightarrow x_{01}$$

$$[] \quad [2, n, 4]$$

$$x_{01} w_{10}^T [n, 4] \rightarrow []$$

$$[] \quad [?m, n]$$

$$x_g \xrightarrow{LN} [x_0 + x_{01}]$$

ReLU

$$x_g w_{13}^T \xrightarrow{FF1} [h \times h, h]$$

$$[] \quad [2, n, 16] \rightarrow x_{g2}$$

$$x_{g2} w_{14}^T \xrightarrow{FF2} [h, h \times n]$$

$$[] \quad [2, n, h] \rightarrow x_{g30}$$

$$x_0 = \xrightarrow{LN} [x_{g2} + x_{g30}]$$

$$[] \quad [2, n, n]$$

Linear

$$x_0 \begin{bmatrix} \end{bmatrix}$$

$$\begin{bmatrix} 2, 4, 6 \end{bmatrix}$$

$$x_0 w_{1,0}^T \begin{bmatrix} 6, 4 \end{bmatrix}$$

$$\begin{bmatrix} 2, 4, 6 \end{bmatrix} \begin{bmatrix} 6, 4 \end{bmatrix}$$

$$\text{softmax} \begin{bmatrix} \end{bmatrix} \begin{bmatrix} 2, 4, 6 \end{bmatrix}$$

$$\begin{bmatrix} 1, 2, 3 \end{bmatrix} \xrightarrow{\alpha} \begin{bmatrix} 0.09 \\ 0.6652 \\ 0.0466 \end{bmatrix} \xrightarrow{1} \begin{bmatrix} 0.09 \\ 0.6652 \\ 0.0466 \end{bmatrix} \xrightarrow{2} \begin{bmatrix} 0.09 \\ 0.6652 \\ 0.0466 \end{bmatrix} \xrightarrow{1}$$

$$\text{loss} \rightarrow [-\log(0.09)] + [-\log(0.0466)] / 2$$

$$= 2.73$$

→ Back propagation

$$-I = \begin{bmatrix} 1 & 6 & 8 & 14 \\ 6 & 1 & 14 & 3 \end{bmatrix}$$

$$y = [0, 2]$$

$$\text{vocab} \sim \begin{bmatrix} 0 \\ 1 \\ 2 \end{bmatrix} \begin{bmatrix} a \\ b \\ c \end{bmatrix}$$