

15-780: Graduate AI

Lecture 10: transformers continued

Aditi Raghunathan

Recap of self-attention

$$\underline{X} = \begin{bmatrix} \leftarrow x_1 \rightarrow \\ \leftarrow x_2 \rightarrow \\ \vdots \\ \leftarrow x_t \rightarrow \end{bmatrix}$$

$\underline{X} \underline{W}$: treat rows independently

"batching"

$$\underline{A} \underline{X}$$

$$\underline{A}(\underline{X}) \underline{X}$$

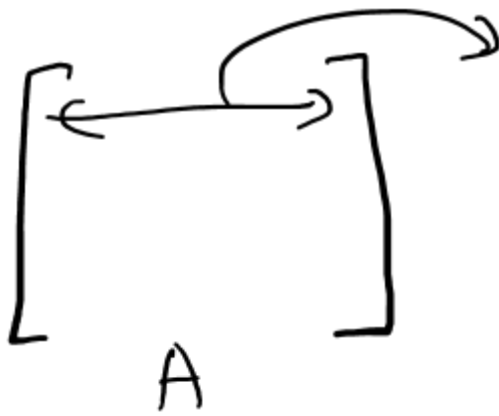
Recap of self-attention

$$A = \text{softmax} \left(\frac{\overbrace{X W_Q}^Q \overbrace{W_K^T X^T}^K}{\sqrt{d}} \right)$$

$Y \in \mathbb{R}^{T \times d}$

$X \in \mathbb{R}^{T \times d}$

$$Y = (A) X \underbrace{W V}_{\text{one set of linear coefficients}}$$



one set of linear coefficients

colored quantities are
"learned" (by SGD)

$$X \Leftarrow \textcircled{e_{\text{the}}} W_Q W_K^T \textcircled{e_{\text{quick}}} - \textcircled{\text{the}} \textcircled{\text{quick}} \text{ brown fox}$$

Properties of self-attention - $\textcircled{\text{the}} \text{ fox brown } \textcircled{\text{quick}}$

• Order invariant: permutations don't change the mixing

• Multi-head attention ✓

• Full mixing: every word can depend on everything else, *masking*

Position embeddings

$$T \times d \begin{bmatrix} X \end{bmatrix} + \begin{bmatrix} E \end{bmatrix} \quad E_{\text{pos}} \text{ depends on pos}$$

$$A(X) \rightarrow A(X + E) \neq A(X) + E$$

$$(X + E) W_Q W_K^T (X + E)^T$$

Position embeddings

→ E_{pos} : $u_{pos} \in \mathbb{R}^d$ that is learnt

→ "Sine embeddings"

$$E_{pos} = \begin{bmatrix} \sin c_1 \cdot pos \\ \cos c_1 \cdot pos \\ \sin c_2 \cdot pos \\ \cos c_2 \cdot pos \\ \vdots \end{bmatrix}$$

$$c_i = \left(\frac{1}{10,000} \right)^{i/d}$$

for $pos = 1 \dots T$

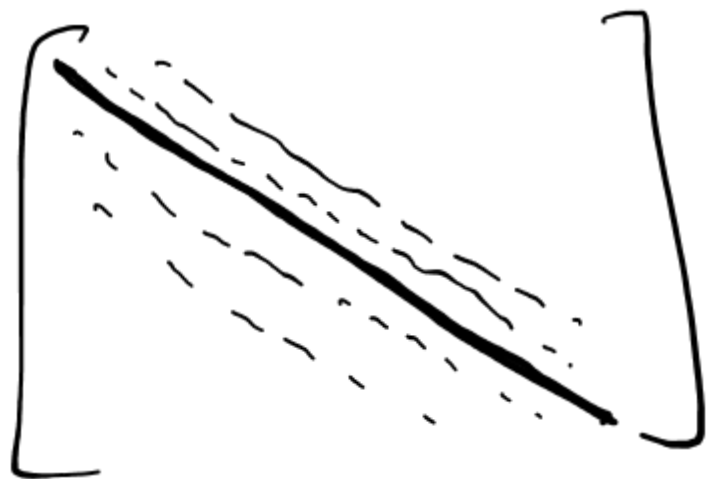
Interpretation of position embedding

For now, $w_Q, w_K \equiv \text{identity}$

$$A(X) = (X + E)(X + E)^T = \underbrace{XX^T}_{\substack{\text{original} \\ \text{attention term}}} + \cancel{XE^T} + \cancel{EX^T} + \underbrace{EE^T}_{\substack{\text{newly} \\ \text{added}}}$$

can ignore

$$EE^T =$$



large at diagonal and slowly decay



Relative embeddings

- $A(x) = x x^T + E E^T$ ↑ write the dot product directly

$\rightarrow (E_i)^T (E_j) \equiv \alpha_{ij} \rightarrow$ learnt α depends on $(j-i)$

- "rotary embeddings": E matrix that encodes relative position

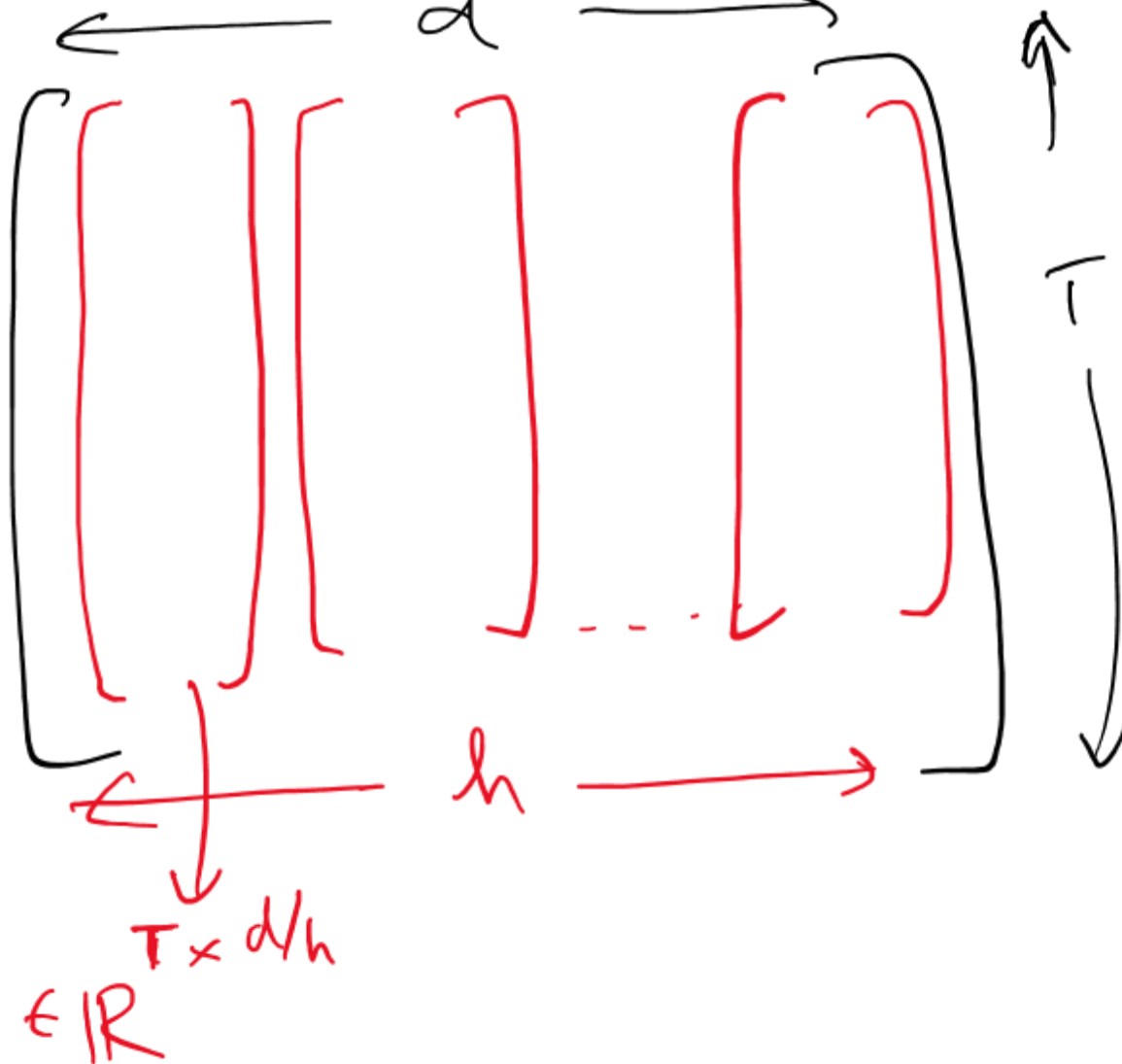
$x E$ (multiplied)

Back to self-attention

$$Y = AX$$

$$A = \text{softmax} \left(\frac{XW_Q W_K^T X^T}{\sqrt{d}} \right)$$

X :



h diff blocks

Back to self-attention

$$Y = A \tilde{X} = \begin{bmatrix} A \end{bmatrix} \begin{bmatrix} \begin{bmatrix} \end{bmatrix} \begin{bmatrix} \end{bmatrix} \dots \end{bmatrix}$$

Diagram illustrating the matrix multiplication $Y = A \tilde{X}$. The matrix A is shown in black. The matrix \tilde{X} is represented by a large bracket containing several smaller red brackets, each representing a column vector. These columns are labeled $\tilde{X}^{(1)}$, $\tilde{X}^{(2)}$, and $\tilde{X}^{(h)}$ in red. A red 'X' is drawn over the label $\tilde{X}^{(1)}$. A dashed red line indicates the continuation of the sequence of columns.

$$Y = A \tilde{X} \text{ mixes the rows}$$

Same mixing vs different mixing

$$\begin{bmatrix} \lambda & (1-\lambda) \\ \beta & (1-\beta) \end{bmatrix} \begin{bmatrix} \hat{a} & \hat{b} \\ c & d \end{bmatrix} = \begin{bmatrix} \lambda \hat{a} + (1-\lambda)c & \lambda \hat{b} + (1-\lambda)d \\ \beta \hat{a} + (1-\beta)c & \beta \hat{b} + (1-\beta)d \end{bmatrix}$$

Piazza poll

$$X = \begin{bmatrix} \text{[red box]} & \text{[red box]} & \text{[red box]} & \dots \end{bmatrix} \quad h \text{ diff blocks}$$

Select the true statements under the standard self-attention formulation $Y = A X$:

- ~~(A)~~ All blocks share the same linear combination or mixing
- (B) Different blocks can have different mixing

we want this in practice

Multi-head attention


$$\begin{aligned}
 & X = \left[\begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right] \quad \text{with head width } h \\
 & W_K, W_Q, W_V = \left[\begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right] \quad \text{with head width } h \\
 & K: W_K X \\
 & Q: W_Q X \\
 & Y = \left[\begin{array}{cc} \tilde{A}^{(1)} & \tilde{X}^{(1)} \\ \tilde{A}^{(2)} & \tilde{X}^{(2)} \end{array} \right] \dots
 \end{aligned}$$

can be different!

Multi-head attention

W_V : could be of higher dimensions

W_O : combines the different heads in some way

$$Y = \underbrace{A(x) \times W_V}_{\text{multi-head}} W_O$$


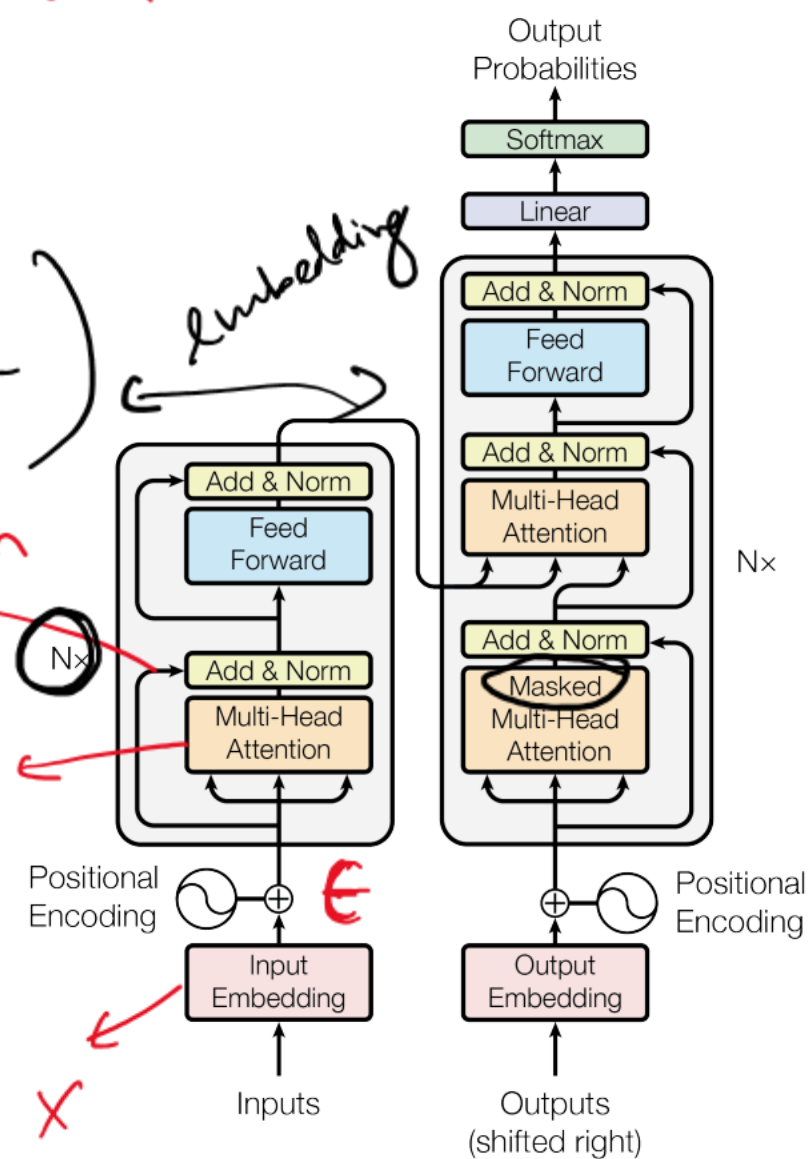
Transformer block

$$\begin{cases} Y = \text{Layer Norm} (X + \text{M.Self A} (X)) \\ Z = \text{Layer Norm} (Y + \sigma(X W_1) W_2) \end{cases}$$

$Y \in \mathbb{R}^{T \times d_{\text{model}}}$
 $W_1 \in \mathbb{R}^{d_{\text{model}} \times d_{\text{ff}}}$
 $W_2 \in \mathbb{R}^{d_{\text{ff}} \times d_{\text{model}}}$



W_K, W_Q, W_V, W_O



decoder part

GPT-3

- Dimension of hidden state: $d_{\text{model}}=12288$
- Dimension of the intermediate feed-forward layer: $d_{\text{ff}}=4d_{\text{model}}$
- Number of heads: $n_{\text{heads}}=96$
- Context length: ~~l~~ $=2048$

Masking

- A way to set certain attention scores to zero

A: completely determined by w_Q, w_K, w_V, E, x

→

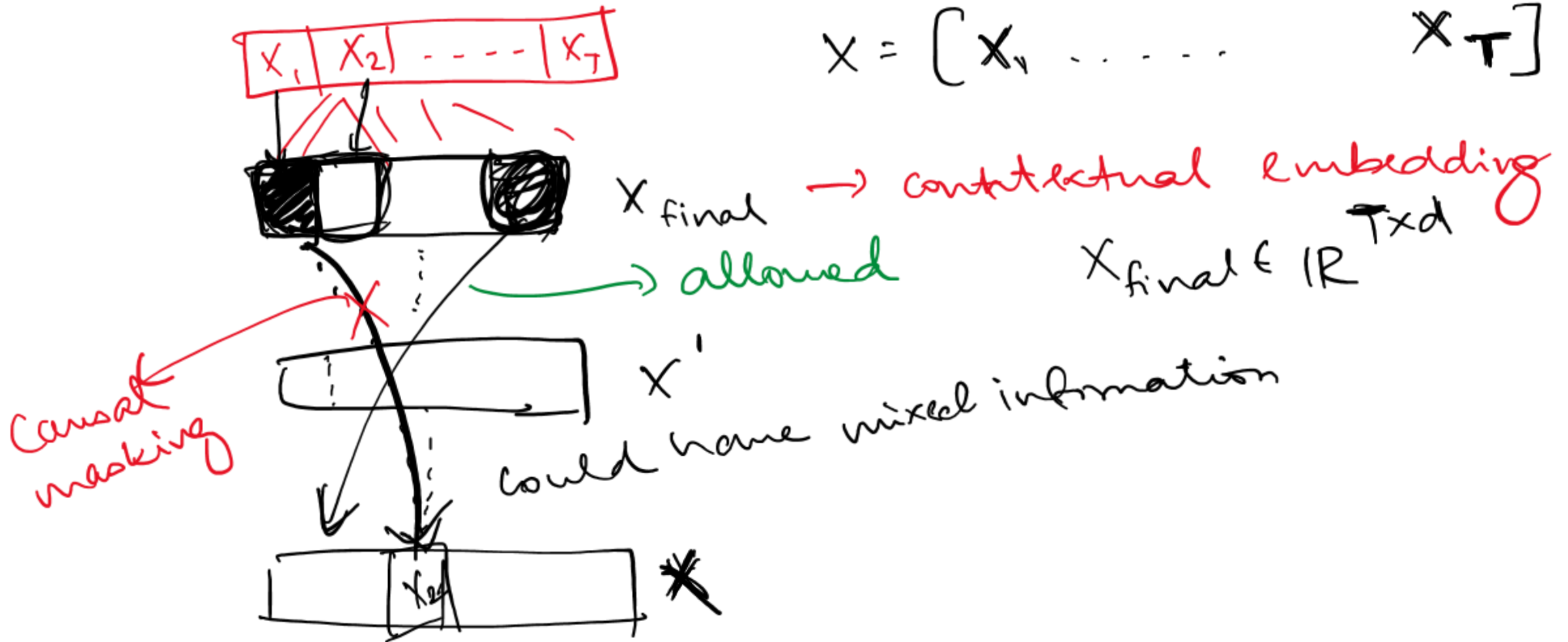
Masking

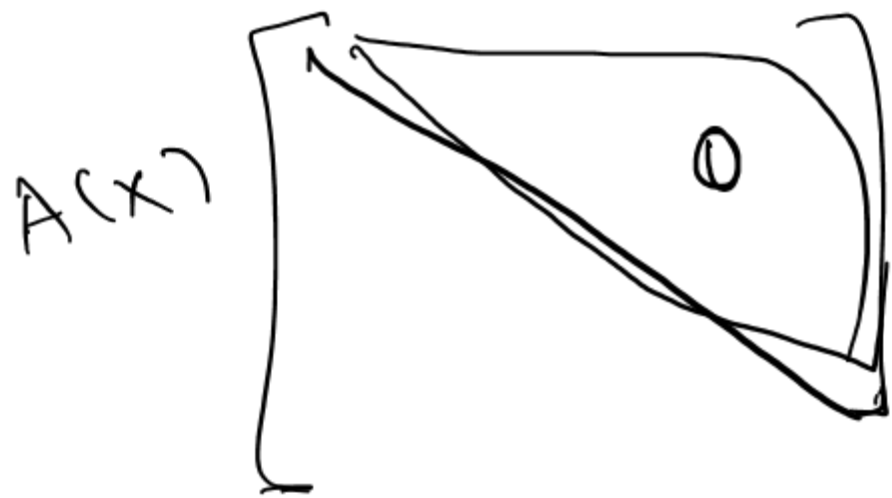
$$P(x_i | x_1, \dots, x_{i-1})$$

"Causal masking"

$$Y = [x_1, x_2, \dots, x_T]$$

$$X = [x_1, \dots, x_T]$$





$$A_{ij} = 0 \quad \text{if } j > i$$

$$A(x) = \text{Softmax} \left(K Q^T + \begin{bmatrix} & -\delta \\ 0 & \end{bmatrix} \right)$$

Summary

train by minimizing
loss via Adam

- You know how to build GPT3!

→ MLPs

→ Add Norm

→ Self Attention

multihhead

masking

→ position embeddings