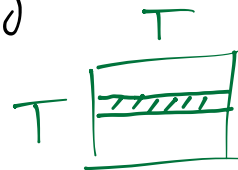
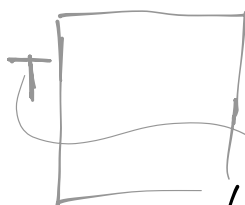


Attention

$$\frac{\sum_j e^{\langle q_i, k_j \rangle}}{z}$$

$$v_j = \sum_j \frac{\text{softmax}(\langle q_i, k_j \rangle)}{\sum_j s_j} v_j$$



Self-Attention

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

$$W^Q \in \mathbb{R}^{d_x \times d_k}$$

$$W^K \in \mathbb{R}^{d_x \times d_k}$$

$$W^V \in \mathbb{R}^{d_x \times d_v}$$

$$\text{softmax}_{\text{axis}=-1} \left(X W^Q \cdot (X W^K)^T \right) \cdot X W^V \quad (\text{SAT})$$

$T \times d_k \quad d_k \times T \quad T \times d_v$

cost $O(d_k \cdot T^2)$

Multihood

$$\left(\text{SAT}_1(X), \dots, \text{SAT}_H(X) \right), \quad M \sim H \cdot d_v \times d_x$$

$T \times H \cdot d_v \quad (\dots) \cdot M$

$+ O(T^2 d_v)$

→ quadratic cost severely constrains the feasible context length more on this soon

Positional encoding

Note that (AT) is invariant to permuting the key/value pairs.

(SAT) is equivariant to permuting X :

$$Y_{\pm i} := X_{\sigma(t)i} \quad \text{some } \sigma \in S_T$$

$$\Rightarrow \text{SAT}(Y)_{\pm i} = \text{SAT}(X)_{\sigma(t)i}$$

For many applications one wants to break this invariance/equivariance.

One wants to encode the positions of things.

\leadsto instead of $q_m = f_q(x_m)$ above $x_m \cdot W^Q$
 \vdots
 $v_m = f_v(x_m)$ ~~x_m~~ $\cdot W^V$

take

$$q_m = f_q(x_m, m)$$

$$\vdots$$

$$v_m = f_v(x_m, m)$$

sin/cos (additive)

Already in [Vaswani2017attention](#) the following positional encoding is suggested.

$$pe_{2i}^{\sin/\cos} = \left(\sin\left(\frac{pos}{10000^{2i/d_{model}}}\right) \right)$$

$$pe_{2i+1}^{\sin/\cos} = \left(\cos\left(\frac{pos}{10000^{2i/d_{model}}}\right) \right), \quad i = 0, \dots, d_{model}/2 - 1, pos = 1, \dots, d_{model}.$$

$$f_\eta(x, m) = \text{~~multiply~~ } (x + pe_m) \cdot W^\eta$$

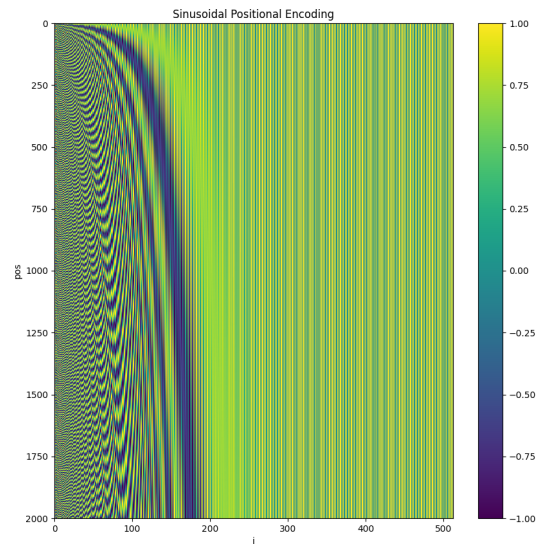
sin/cos (multiplicative)

$$f_\eta(x, m) = R_m W_\eta x \text{~~re~~}$$

where

$$R_m = \begin{pmatrix} \cos(m\theta_1) & \sin(m\theta_1) & 0 & \dots & 0 & 0 \\ -\sin(m\theta_1) & \cos(m\theta_1) & 0 & \dots & 0 & 0 \\ 0 & 0 & \cos(m\theta_2) & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \cos(m\theta_{d/2}) & \sin(m\theta_{d/2}) \\ \dots & \dots & \dots & \dots & -\sin(m\theta_{d/2}) & \cos(m\theta_{d/2}) \end{pmatrix},$$

$$\theta_i := 10000^{-2(i-1)/d}, i = 1, \dots, d/2.$$



Transformer architecture

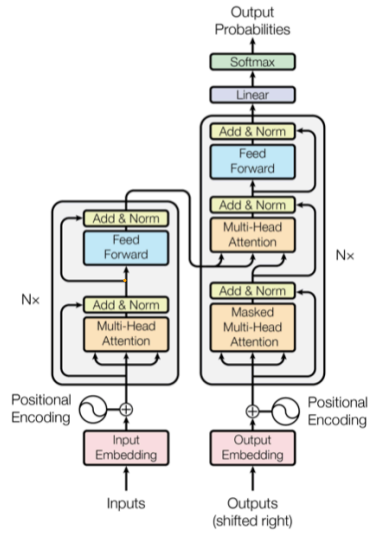


Figure 1: The Transformer - model architecture.

GPT, Llama .. modify this original Transformer architecture in that they are "decoder-only", they only have the right track (and discard the second Multi-Head-Attention, which is also called Cross-Attention)

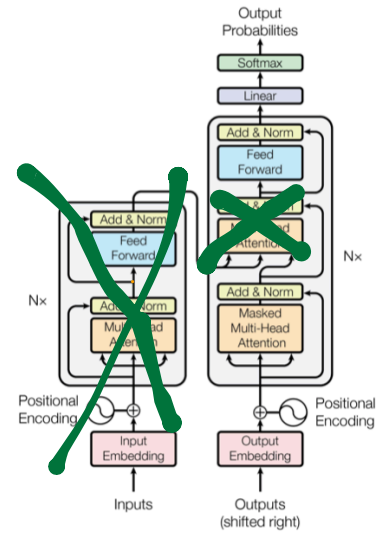


Figure 1: The Transformer - model architecture.

Vaswani et al '17

→ bdk at makemore