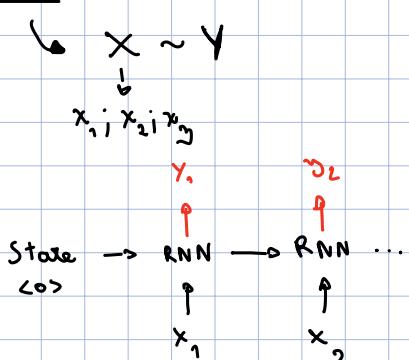


Transformer from scratch

RNN



Problem:

1. Slow computation

2. Vanishing or exploding gradients

3. Difficulty in accessing info from long time ago

$$x \rightarrow f(x,y) = x \cdot y \xrightarrow{\text{square}} g(z) = z^2$$

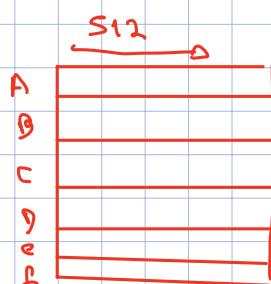
$$\frac{dS}{dx} = \underbrace{\frac{dg}{df}}_{\frac{1}{2}} \cdot \underbrace{\frac{df}{dx}}_{\frac{1}{2}}$$

Transformer

Notations

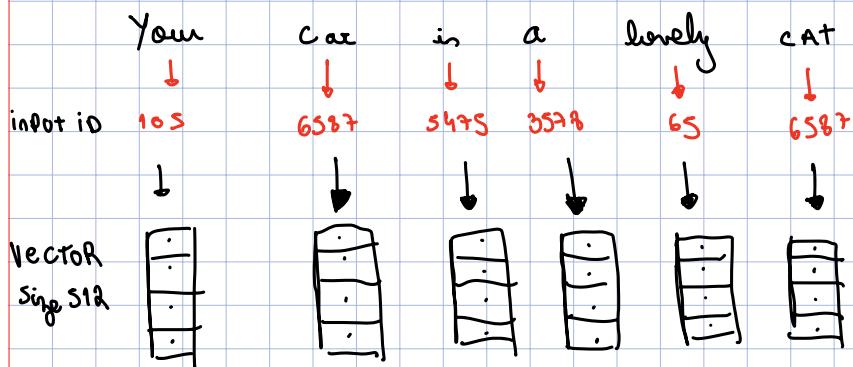
$A \in \mathbb{R}^{6 \times 512}$

$(A)^t = A \in \mathbb{R}^{512 \times 6}$



$$A \otimes C \otimes e^t = B \in \mathbb{R}^{6 \times 6}$$

Encoder

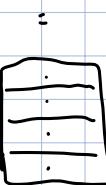
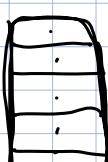


positional encoding → represent a pattern that can learn

$$p_e(p_{en}; i) = \min \frac{p_{en}}{10000^{\frac{i}{n_e}}}$$

$$p_e(p_{en}; i+1) = \cos \frac{i \pi}{10000^{\frac{2}{n_e}}}$$

Positional
embedding
(512)



Encoder
input

- We only need to compute one positional encoding once and then reuse it for every sentence, no matter if it is training or inference.

Multi-head attention

↳ Self attention

$$\hookrightarrow \text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

$\text{req}_f = 6 \wedge d_{\text{model}} = d_k = 512$

$$\begin{array}{c} \text{softmax} \\ \text{refer to } Q \quad K^T \\ \text{dimensions: } (6, 512) \quad (512, 6) \\ \hline \sqrt{512} \end{array}$$

diagonal

	YOUR	CAT	IS	A	LOVELY	CAT
YOUR	0.268	0.119	0.134	0.148	0.179	0.152
CAT	0.124	0.278	0.201	0.128	0.154	0.115
IS	0.147	0.132	0.262	0.097	0.218	0.145
A	0.210	0.128	0.206	0.212	0.119	0.125
LOVELY	0.146	0.158	0.152	0.143	0.227	0.174
CAT	0.195	0.114	0.203	0.103	0.157	0.229

$$X \quad V = \text{Attention} \quad (6, 512)$$

Each row in this matrix captures not only the meaning (given by the embedding) or the position in the sentence (represented by the positional encodings) but also each word's interaction with other words.

Multi-head Attention

Input

$$\begin{aligned} Q \times w^Q &= Q' \rightarrow \text{will see } T \in \mathbb{R} \\ K \times w^K &= K' \\ V \times w^V &= V' \end{aligned}$$

d_{model}
 $\rightarrow n$

$s \times d_k$

$$Q \times w^Q = Q' \rightarrow \text{will see } T \in \mathbb{R}$$

$$\text{Attention } (Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$

$$\text{Head}_i = \text{Attention} (QW_i^Q, KW_i^K; VW_i^V)$$

$$\text{MultiHead}(Q, K, V) = \text{Concat} (\text{Head}_1, \dots, \text{Head}_n) W^O$$

$$\hookrightarrow \begin{matrix} H \\ (\text{req}; \text{Head}_n) \end{matrix} \times \begin{matrix} W^O \\ (\text{dim}_H; \text{dim}_{\text{model}}) \end{matrix} = \boxed{\text{MH-A}}$$

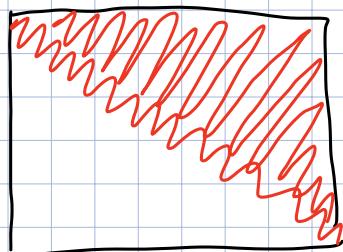
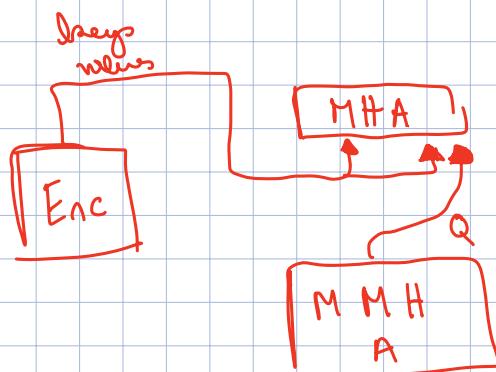
Why $\boxed{Q, K, V}$ \rightarrow amplify noise \boxed{QK}

Layer normalization $\rightarrow \hat{x}_j = \frac{x_j - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$

Decoder

\hookrightarrow Masked Multi-Head
Attention

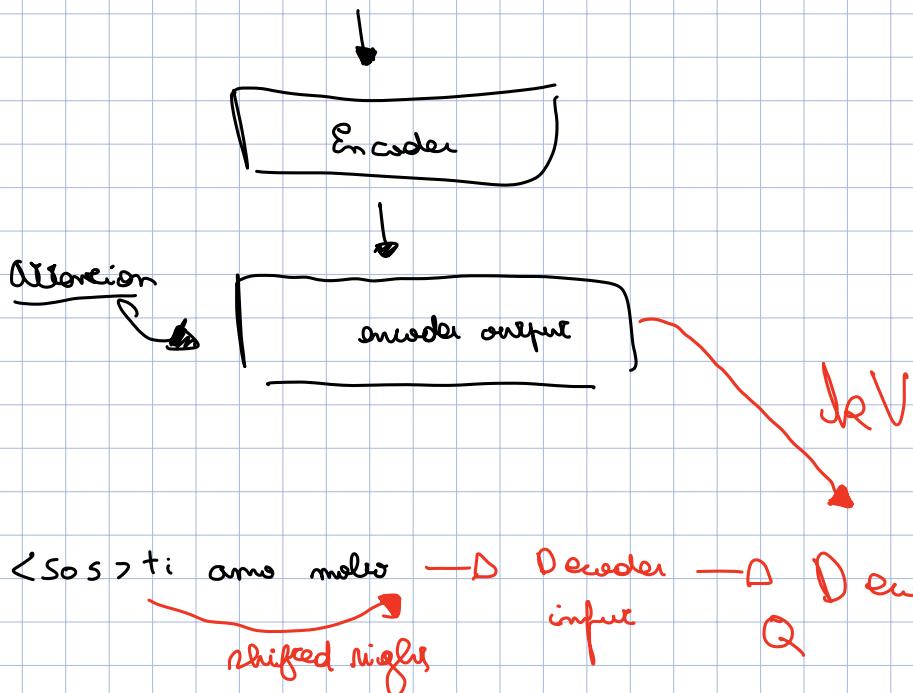
\hookrightarrow the model must not be able to see future words



Inference and training

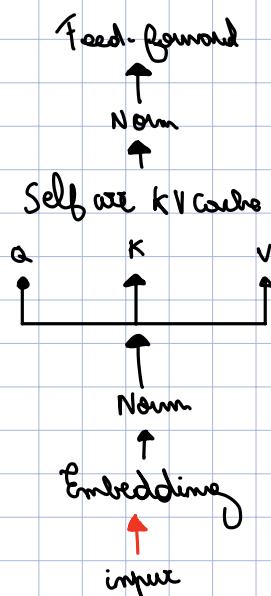
↳ training (translation)

<SOS> I love you very much <EOS>



$$\rightarrow (\text{Seq}; d_{\text{model}}) \xrightarrow{\text{linear}} (\text{Seq}; \text{vocab_size}) \xrightarrow{\text{softmax}} (\text{Seq}; \text{vocab_size})$$

Llama architecture



GQA \Rightarrow grouped query attention

Linear Layer

$$G: x \cdot W^T + b$$

residual weights ↘

- The distribution of internal nodes (neurons) of a neural network along is referred to as **Internal Covariance Shift**
- ↳ Normalization $\xrightarrow{u} \sigma^2$; $y + b \rightarrow$ learnable parameters

!!!

AHAAA!!!!

↓ Gaussian normalization

$$x \sim N(5; 36) \longrightarrow \frac{x - 5}{\sqrt{36}} = z \sim N(0; 1)$$

standard Gaussian

RMS Mean Square Normalization

↑ requires < computation

$$\bar{a}_i = \frac{a_i}{\text{RMS}(a)} g_i ; j \text{ where } \text{RMS}(a) = \sqrt{\frac{1}{m} \sum_{i=1}^m a_i^2}$$

Absolute positional encoding

Given 2 tokens we create a vector that represents a distance

multiple transformations they're applied before

Absolute positional encoding

$$ATT = \left(\frac{QK^T}{\sqrt{d_m}} \right) V$$

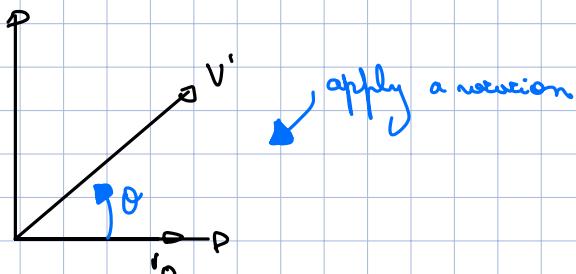
Relative positional encoding

$$Q \left(K + \frac{\alpha_s^K}{\text{distance}} \right)^+$$

Rotary positional encoding

$$\begin{pmatrix} \cos m\theta & -\sin m\theta \\ \sin m\theta & \cos m\theta \end{pmatrix}$$

rotary position encoding



Long-term decay

The intensity of relationships between two tokens encoded with Rotary Positional encoding will be numerically smaller as the distance between them grows

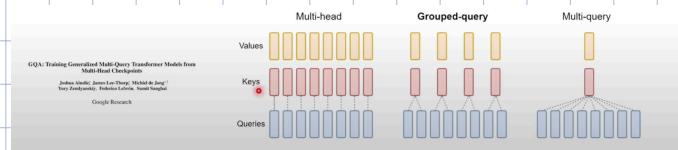
Multi-query att

Self-Attention

$$\text{Attention} = \text{softmax}\left(\frac{QK^T}{\sqrt{dk}}\right)V$$

KV-Cache → less computation
only on during inference

Grouped query Attention
↳ shared heads for KV



feed-forward

SwiGLU Attention function

FP16 → 2 bytes