

Lec 5 Transformer

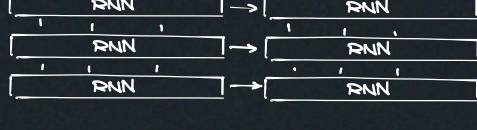
Transformer: self-attention or self-attention + cross-attention

Query keys to get attn. scores, softmax to get weight

Attention is All You Need

seq2seq model entirely based on attention

Before:



Transformer: - fast!
- everything matrix

Types: - Encoder-decoder - TS, MBART
- Decoder-only - GPT, LLaMa

→ See Transformer diagrams

Input: split using subwords

Multi-head attention - attend to different parts of sentence to disambiguate

attend multiple places for multiple functions

I run a small business

I run a mile

The robber made a run for it

The stocking had a run

MultiHead(Q, k, v) = Concat(head...n)

head_i = Attention(W^QQ_i, W^kK_i, W^vV_i)

Query Q = $\begin{bmatrix} \square & \square & \square \\ \square & \square & \square \\ \square & \square & \square \end{bmatrix}$ Key K = $\begin{bmatrix} \square & \square & \square & \square \\ \square & \square & \square & \square \\ \square & \square & \square & \square \end{bmatrix}$ Val V = $\begin{bmatrix} \square & \square & \square \\ \square & \square & \square \\ \square & \square & \square \end{bmatrix}$

$\downarrow * W^Q$ $\downarrow * W^K$ $\downarrow * W^V$

split into multi-heads

do attention on each head

prefix LM: do both masked and unmasked attention

concat

linear

Positional Emb - deal with attention ignoring order

A big dog and a big cat

$W_{big} + W_{pos2}$

$W_{big} + W_{pos6}$

→ Sinusoidal func $P_t^{(i)} = f(t)^{(i)} = \begin{cases} \sin(W_k \cdot t) & i = 2k \\ \cos(W_k \cdot t) & i = 2k + 1 \end{cases}$

$$W_k = \frac{1}{10000^{2k/d}}$$

motivation: dot prod of two embs higher if words closer

→ Or just learn the positional embedding

↳ simpler, but won't extrapolate to larger input

→ Abs vs. Rel encoding

↳ positional, and hope relative info is recovered

↳ explicitly encode relative location e.g. whether key is -5 of query

→ Rotary Positional Encoding - dot prod of emb and rel position

$$f_q(x_m, m) \cdot f_k(x_n, n) = g(x_m, x_n, m - n)$$

{ RoPE uses imaginary number stuff involved

Training stability

Layer normalisation

$$\text{LayerNorm}(x; g, b) = \frac{g}{\sigma(x)} \odot (x - \mu(x)) + b$$

vec stddev vec mean



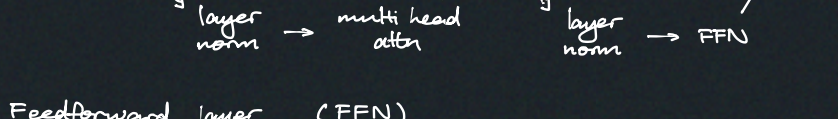
RMSNorm - simplified LayerNorm, no mean subtraction & no bias

Residual Connections

from input to Add & Norm

→ attention learns to attend elsewhere, since self is free through residual

Pre/post layer norm

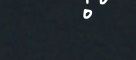


Feedforward layer (FFN)

$$\text{FFN}(x; W_1, b_1, W_2, b_2) = \underbrace{f(xW_1 + b_1)}_{\text{sometimes just disable}} \underbrace{W_2 + b_2}_{\text{non linearity}}$$

Choosing f

- ReLU



↳ LLaMa

- Swish

↳ SiLU



Model Optimisation

Optimisers

- SGD
- Adam
- Vaswani + 2017 lr schedule



- AdamW - weight decay to regularise Adam

Low-precision Training

- 16 bit training
 - ↳ IEEE half precision, 5-bit exponent
 - ↳ bfloat16, 8-bit precision ← usually more stable

Checkpoint & Restart

- Monitor gradient norm
- Check your code

Comparing Transformers

Original	LLaMa
post norm	pre norm
LayerNorm	RMSNorm
ReLU	SiLU
Sinusoidal	RoPE