

RMS Norm:

$$y = \frac{x}{\sqrt{\frac{1}{n} \sum x^2 + \epsilon}} * \gamma$$

layer norm:

$$y = \frac{x - E(x)}{\sqrt{\text{Var}[x] + \epsilon}} * \gamma + \beta$$

ROPE

$$\bar{g}_m = f_g(x_m, m) \quad \bar{k}_n = f_k(x_n, n)$$

$$\langle f_g(x_m, m), f_k(x_n, n) \rangle = g(x_m, x_n, m-n)$$

$$f_g(x_m, m) = (W_g x_m) e^{im\theta}, \quad f_k(x_n, n) = (W_k x_n) e^{in\theta}$$

$$g(x_m, x_n, m-n) = \text{Re}[(W_g x_m)(W_k x_n)^* e^{i(m-n)\theta}]$$

$$\begin{bmatrix} x_m^{i'} & x_m^{j'} \end{bmatrix} = \begin{bmatrix} x_m^i & x_m^j \end{bmatrix} \begin{bmatrix} \cos(m\theta_m) & -\sin(m\theta_m) \\ \sin(m\theta_m) & \cos(m\theta_m) \end{bmatrix}$$

KV Cache

Q: (seq len, dim)

K: (dim, seq len)

V: (seq len, dim)

Z: (seq len, dim)

q_1
q_2
q_3
\vdots
q_i

x

k_1	k_2	\dots	k_i
-------	-------	---------	-------

Q

K^T

QK^T

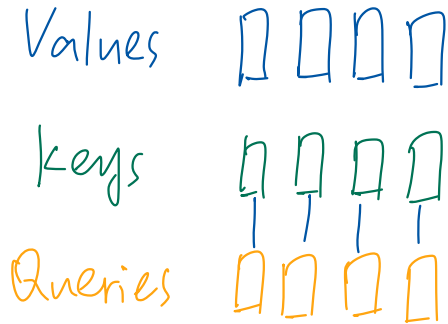
Token i

T_1	T_2	T_3	T_4
-------	-------	-------	-------

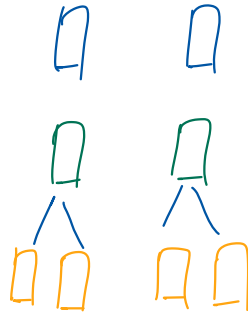
--	--	--	--

Grouped Multi-Query Attention.

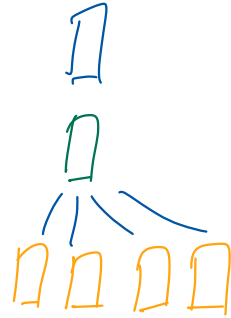
MHA



Grouped-Query



Multi-query.



Feedforward

$$\text{out} = W_3 \sigma(W_1 x * W_2 x)$$

$$\text{SiLU}(x) = \frac{x}{1 + e^{-x}}$$

Transformer:

$$\text{out} = \sigma(W_1 x) W_2$$

$$hb = x_b^T W_1$$

$$hb_2 = x_b^T W_3$$

$$hb = \text{SiLU}(hb)$$

$$hb = hb * hb_2$$

$$xb = hb^T W_2$$

$$x = xb + x$$