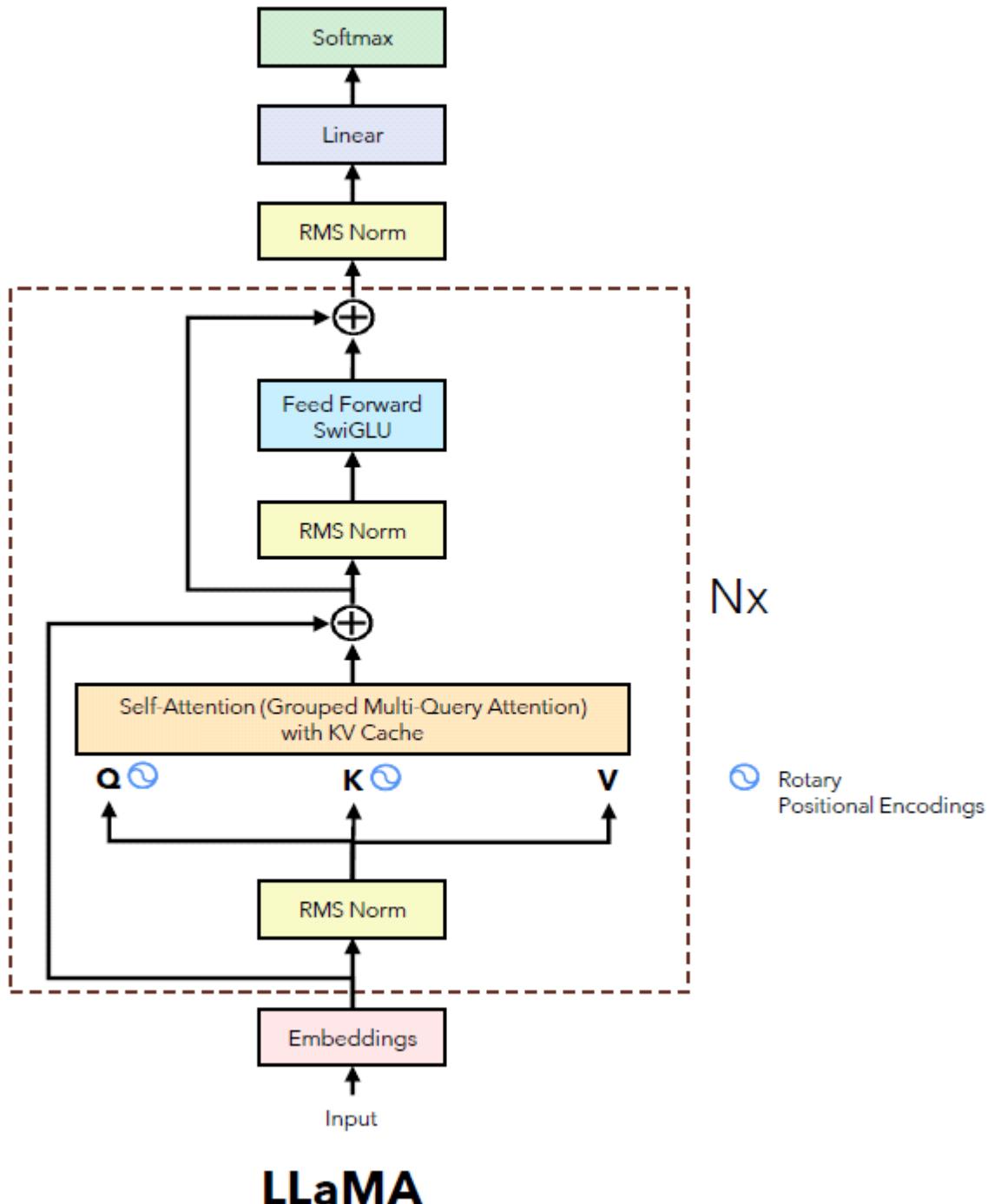


# Llama 2

27 July 2025 22:12



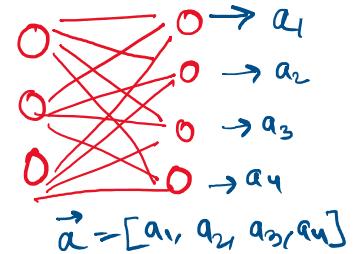
## 4 RMSNorm

A well-known explanation of the success of LayerNorm is its re-centering and re-scaling invariance property. The former enables the model to be insensitive to shift noises on both inputs and weights, and the latter keeps the output representations intact when both inputs and weights are randomly scaled. In this paper, we hypothesize that the re-scaling invariance is the reason for success of LayerNorm, rather than re-centering invariance.

We propose RMSNorm which only focuses on re-scaling invariance and regularizes the summed inputs simply according to the root mean square (RMS) statistic:

$$\bar{a}_i = \frac{a_i}{\text{RMS}(\mathbf{a})} g_i, \quad \text{where } \text{RMS}(\mathbf{a}) = \sqrt{\frac{1}{n} \sum_{i=1}^n a_i^2}. \quad (4)$$

Intuitively, RMSNorm simplifies LayerNorm by totally removing the mean statistic in Eq. (3) at the cost of sacrificing the invariance that mean normalization affords. When the mean of summed inputs is zero, RMSNorm is exactly equal to LayerNorm. Although RMSNorm does not re-center



Just like Layer Normalization, we also have a learnable parameter **gamma** ( $\mathbf{g}$  in the formula on the left) that is multiplied by the normalized values.

element wise multiplication.

```
class RMSNorm(nn.Module):
    Tabnine | Edit | Test | Explain | Document
    def __init__(self, dim: int, eps: float = 1e-6):
        super().__init__()
        self.eps = eps
        ## The gamma parameter
        self.weight = nn.Parameter(torch.ones(dim))

    Tabnine | Edit | Test | Explain | Document
    def _norm(self, x: torch.Tensor):
        # (B, seq_len, dim) * (B, seq_len, 1) -> (B, seq_len, dim)
        # rsqrt: 1/sqrt(x)
        return x * torch.rsqrt(x.pow(2).mean(-1, keepdim=True) + self.eps)

    Tabnine | Edit | Test | Explain | Document
    def forward(self, x: torch.Tensor):
        # (dim) * (B, seq_len, dim) -> (B, seq_len, dim)
        return self.weight * self._norm(x.float()).type_as(x) # This is basically element wise matrix multiplication
```

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

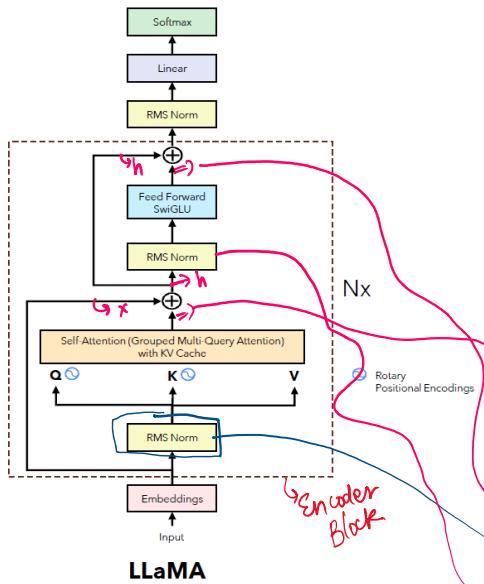
↓

↓

<p style="position: absolute; top:

## Encoder Block

27 July 2025 23:53



```
class EncoderBlock(nn.Module):
    Tabnine | Edit | Test | Fix | Explain | Document
    def __init__(self, args: ModelArgs):
        super().__init__()

        self.n_heads = args.n_heads
        self.dim = args.dim
        self.head_dim = args.dim // args.n_heads

        self.attention = SelfAttention(args)
        self.feed_forward = FeedForward(args)

        # Normalization BEFORE the self.attention
        self.attention_norm = RMSNorm(args.dim, eps=args.norm_eps)

        # Normalization BEFORE the feed forward block
        self.ffn_norm = RMSNorm(args.dim, eps=args.norm_eps)

    Tabnine | Edit | Test | Explain | Document
    def forward(self, x: torch.Tensor, start_pos: int, freqs_complex: torch.Tensor):
        # (N, seq_len, dim) + (b, seq_len, dim) => (b, seq_len, dim)
        h = x + self.attention.forward(self.attention_norm(x), start_pos, freqs_complex)
        out = h + self.feed_forward.forward(self.ffn_norm(h))
        return out
```

$$R_{\Theta, m}^d x = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ \vdots \\ x_{d-1} \\ x_d \end{pmatrix} \otimes \begin{pmatrix} \cos m\theta_1 \\ \cos m\theta_1 \\ \cos m\theta_2 \\ \cos m\theta_2 \\ \vdots \\ \cos m\theta_{d/2} \\ \cos m\theta_{d/2} \end{pmatrix} + \begin{pmatrix} -x_2 \\ x_1 \\ -x_4 \\ x_3 \\ \vdots \\ -x_{d-1} \\ x_d \end{pmatrix} \otimes \begin{pmatrix} \sin m\theta_1 \\ \sin m\theta_1 \\ \sin m\theta_2 \\ \sin m\theta_2 \\ \vdots \\ \sin m\theta_{d/2} \\ \sin m\theta_{d/2} \end{pmatrix}$$

How do we calculate this in code

1)  $\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} \xrightarrow{\text{shape}} \begin{bmatrix} [x_1 \ x_2] \\ [x_3 \ x_4] \end{bmatrix} \xrightarrow{\substack{\text{①} \\ \text{②}}} \begin{bmatrix} x_1 + i x_2 \\ x_3 + i x_4 \end{bmatrix} \xrightarrow{\substack{\text{convert} \\ \text{to polar}}} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} \xrightarrow{\text{③}} \text{u-complex}$

2)  $m = [1, 2, 3, \dots, \text{seq\_len}]$

$$\Theta = [\theta_1, \theta_2, \theta_3, \dots, \theta_{d/2}]$$

$m$       Outer-product       $\Theta \Rightarrow$

$$\begin{bmatrix} m_1 \theta_1 & m_1 \theta_2 & m_1 \theta_3 & \dots & m_1 \theta_{d/2} \\ m_2 \theta_1 & m_2 \theta_2 & m_2 \theta_3 & \dots & m_2 \theta_{d/2} \\ \vdots & & & & \end{bmatrix}$$

freqs-complex

$\Downarrow$  Polar

$$\left[ \cos(m_1 \theta_1) + i \sin(m_1 \theta_1) \right] \left[ \cos(m_1 \theta_2) + i \sin(m_1 \theta_2) \right] \dots$$

element  
wise  
product

(Doing for  $m_1$ )

$\Rightarrow$   $x$ -complex  $\otimes$   $f$ -wise freqs-complex  $\otimes$  (Doing for  $m_1$ )

$$\Rightarrow \begin{bmatrix} x_1 + i x_2 \\ x_3 + i x_4 \end{bmatrix} \otimes \begin{bmatrix} \underbrace{\cos(m_1 \theta_1) + i \sin(m_1 \theta_1)}_{f_1} \\ \underbrace{\cos(m_1 \theta_2) + i \sin(m_1 \theta_2)}_{f_2} \end{bmatrix}$$

$(x_1 + i x_2)(f_1 + i f_2) = x_1 f_1 + i x_1 f_2 + i x_2 f_1 - x_2 f_2$   
 $= (x_1 f_1 - x_2 f_2) + i(x_1 f_2 + x_2 f_1)$

$$\left[ \begin{array}{c} x_1 f_1 - x_2 f_2 + i(x_1 f_2 + x_2 f_1) \\ x_3 f_3 - x_4 f_4 + i(x_3 f_4 + x_4 f_3) \end{array} \right]$$

↓③ Convert to Real

$$\left[ \begin{array}{cc} (x_1 f_1 - x_2 f_2) & x_1 f_2 + x_2 f_1 \\ x_3 f_3 - x_4 f_4 & x_3 f_4 + x_4 f_3 \end{array} \right]$$

↓④ Flatten

$$\left[ \begin{array}{c} x_1 f_1 - x_2 f_2 \\ x_1 f_2 + x_2 f_1 \\ x_3 f_3 - x_4 f_4 \\ \dots \end{array} \right]$$

⇒ This is the matrix that we want.

$$\begin{bmatrix} x_3 f_3 - x_4 f_4 \\ x_3 f_4 + x_4 f_3 \end{bmatrix}$$

1. we want

# Self Attention

August 11, 2025 11:45 PM

Rough  $\Rightarrow$   $n \cdot k \cdot v\text{-heads}$ ,  $q_v\text{-heads}$ ,  $h\text{-rep} = \frac{n \cdot q_v}{n \cdot k \cdot v}$

$$\text{head\_dim} = \text{dim} // n \cdot q_v$$

$w_{qv} \Rightarrow (\text{dim}, n \cdot q_v \times \text{head\_dim})$

$w_k \Rightarrow (\text{dim}, n \cdot k \cdot v \times \text{head\_dim})$

$w_v \Rightarrow (\text{dim}, n \cdot k \cdot v \times \text{head\_dim})$

$w_o \Rightarrow (n \cdot q_v \times \text{head\_dim}, \text{dim})$

$\text{cache\_k} \Rightarrow (\text{max\_bs}, \text{max\_seq\_len}, n \cdot k \cdot v, \text{head\_dim})$

$\text{cache\_v} \Rightarrow (\text{max\_bs}, \text{max\_seq\_len}, n \cdot k \cdot v, \text{head\_dim})$

$x \rightarrow (b, \text{seq\_len}, \text{dim})$

$x_{qv} = w_{qv}(x) \Rightarrow (b, \text{seq\_len}, \text{dim}) * (\text{dim}, \text{dim})$

$\Rightarrow (b, \text{seq\_len}, \text{dim})$

$x_k = w_k(x) \Rightarrow (b, \text{seq\_len}, \text{dim}) * (\text{dim}, n \cdot k \cdot v \times \text{head\_dim})$

$\Rightarrow (b, \text{seq\_len}, n \cdot k \cdot v \times \text{head\_dim})$

$x_v = w_v(x) \Rightarrow (b, \text{seq\_len}, \text{dim}) * (\text{dim}, n \cdot k \cdot v \times \text{head\_dim})$

$\Rightarrow (b, \text{seq\_len}, n \cdot k \cdot v \times \text{head\_dim})$

$x_{qv} \Rightarrow (b, \text{seq\_len}, \text{dim}) \rightarrow (b, \text{seq\_len}, n \cdot q_v, \text{head\_dim})$

$x_k \Rightarrow (b, \text{seq\_len}, n \cdot k \cdot v \times \text{head\_dim}) \rightarrow (b, \text{seq\_len}, n \cdot k \cdot v, \text{head\_dim})$

$x_v \Rightarrow (b, \text{seq\_len}, n \cdot k \cdot v \times \text{head\_dim}) \rightarrow (b, \text{seq\_len}, n \cdot k \cdot v, \text{head\_dim})$

$x_{qv} \Rightarrow \text{RoPE}(x_{qv}) (b, \text{seq\_len}, n \cdot q_v, \text{head\_dim})$

$x_k \Rightarrow \text{RoPE}(x_k) (b, \text{seq\_len}, n \cdot k \cdot v, \text{head\_dim})$

$\text{cache\_k}[:, b, \text{start\_pos}: \text{start\_pos} + \text{seq\_len}] = x_k$

$\text{cache\_v}[:, b, s\text{-pos}: s\text{-pos} + \text{seq\_len}] = x_v$

$\text{key} = \text{cache}[:, b, : s\text{-pos} + \text{seq\_len}]$   
 $\text{value} = \text{cache}[:, b, : s\text{-pos} + \text{seq\_len}]$

$$\text{key} = \text{cache\_k}[:, :b] : \text{seq\_len} \rightarrow [b, \text{seq\_len}]$$

$$\text{value} = \text{cache\_v}[:, :b] : \text{seq\_len} \rightarrow [b, \text{seq\_len}]$$

Repetition of keys and values upto n-key  
but for every  $\forall v$ ,  $k \neq v$  will be same.

$$\therefore \text{key} = \text{repeat\_kv}(\text{key}, \text{seq\_len})$$

$$\Rightarrow (b, \text{seq\_len}, \text{n\_kv}, \text{head\_dim}) \rightarrow (b, \text{seq\_len}, \frac{\text{n\_kv}}{\text{head\_dim}})$$

$$\text{value} = \dots$$

$$\Rightarrow (b, \text{seq\_len}, \text{n\_kv}, \text{head\_dim}) \rightarrow (b, \text{seq\_len}, \frac{\text{n\_kv}}{\text{head\_dim}})$$

$$x_{av} = \text{av. transpose}(1, 2)$$

$$\Rightarrow (b, \text{seq\_len}, \text{n\_av}, \text{head\_dim}) \rightarrow (b, \text{n\_av}, \text{seq\_len}, \text{head\_dim})$$

$$\text{key} = \text{transpose}(1, 2)$$

$$\Rightarrow (b, \text{seq\_len}, \text{n\_av}, \text{head\_dim}) \rightarrow (b, \text{n\_av}, \text{seq\_len}, \text{head\_dim})$$

$$\text{value} = \dots$$

$$\Rightarrow (b, \text{seq\_len}, \text{n\_av}, \text{head\_dim}) \rightarrow (b, \text{n\_av}, \text{seq\_len}, \text{head\_dim})$$

$$x_{av} \leftarrow \text{key. transpose}(2, 3)$$

$$\Rightarrow (b, \text{n\_av}, \text{seq\_len}, \text{head\_dim}) \times (b, \text{n\_av}, \text{head\_dim}, \text{seq\_len})$$

↓

$$(b, \text{n\_av}, \text{seq\_len}, \text{seq\_len})$$

↓ softmax

↓

$$(b, \text{n\_av}, \text{seq\_len}, \text{seq\_len})$$

$$\text{o/p} \Rightarrow \text{sum} \times \text{value}$$

$$\Rightarrow (b, \text{n\_av}, \text{seq\_len}, \text{seq\_len}) \times (b, \text{n\_av}, \text{seq\_len}, \text{head\_dim})$$

$$\times (1 \dots \text{n\_av}, \text{seq\_len}, \text{head\_dim})$$

$\Rightarrow (b, n\_ar, \text{seq\_len}, \text{seq\_len}) \times (n\_ar, \text{head\_dim})$   
 $\Rightarrow (b, n\_ar, \text{seq\_len}, \text{head\_dim})$   
 $\text{op} \Rightarrow (b, n\_ar, \text{seq\_len}, \text{head\_dim})$   
 $\downarrow$   
 $(b, \text{seq\_len}, n\_ar, \text{head\_dim})$   
 $\downarrow$   
 $(b, \text{seq\_len}, \text{dim})$

$\text{op} \Rightarrow \text{Wo}(\text{op})$   
 $\Rightarrow (b, \text{seq\_len}, \text{dim}) \times (\text{dim}, \text{dim})$   
 $\Rightarrow (b, \text{seq\_len}, \text{dim})$