

Jean Soler, Olivier Serris, Nicolas Castanet  
Sorbonne Université - M2 DAC


$$Attention(Q, K, V) = Softmax\left(\frac{QK^T}{\sqrt{d}}\right)V$$

Dans l'attention, on a les matrices **Querys**, **Keys** et **Values** :

$$Q \in \mathbb{R}^{Ld}, K^T \in \mathbb{R}^{dL} \text{ et } V \in \mathbb{R}^{Ld}$$

The diagram illustrates the dimensions and complexity of two matrix multiplication sequences:

- Left Sequence:  $(QK^T)V$** 
  - Matrix  $Q$  has dimensions  $L \times L$ .
  - Matrix  $K^T$  has dimensions  $L \times d$ .
  - Matrix  $V$  has dimensions  $L \times d$ .
  - The intermediate product  $QK^T$  has dimensions  $(L \times L)(L \times d) = (L \times d)$ .
  - The final product  $(QK^T)V$  has dimensions  $(L \times d)(L \times d) = (L \times d)$ .
  - The time complexity is  $\mathcal{O}(L^2d)$ .
- Right Sequence:  $Q(K^TV)$** 
  - Matrix  $K^T$  has dimensions  $L \times d$ .
  - Matrix  $V$  has dimensions  $d \times L$ .
  - Matrix  $Q$  has dimensions  $L \times L$ .
  - The intermediate product  $K^TV$  has dimensions  $(d \times L)(L \times d) = (d \times d)$ .
  - The final product  $Q(K^TV)$  has dimensions  $(L \times L)(d \times d) = (L \times d)$ .
  - The time complexity is  $\mathcal{O}(d^2L)$ .

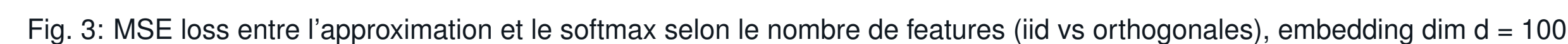
⇒ Le softmax force la complexité de calcul à être quadratique en L :  $\mathcal{O}(L^2d)$

$$Q'K'^T \approx \text{Softmax}\left(\frac{QK^T}{\sqrt{d}}\right)$$
$$Attention(Q, K, V) \approx Q'(K'^T V)$$
$$\begin{cases} A = \exp(QK^T/\sqrt{d}) \\ D = \text{diag}(A1_L) \end{cases}$$

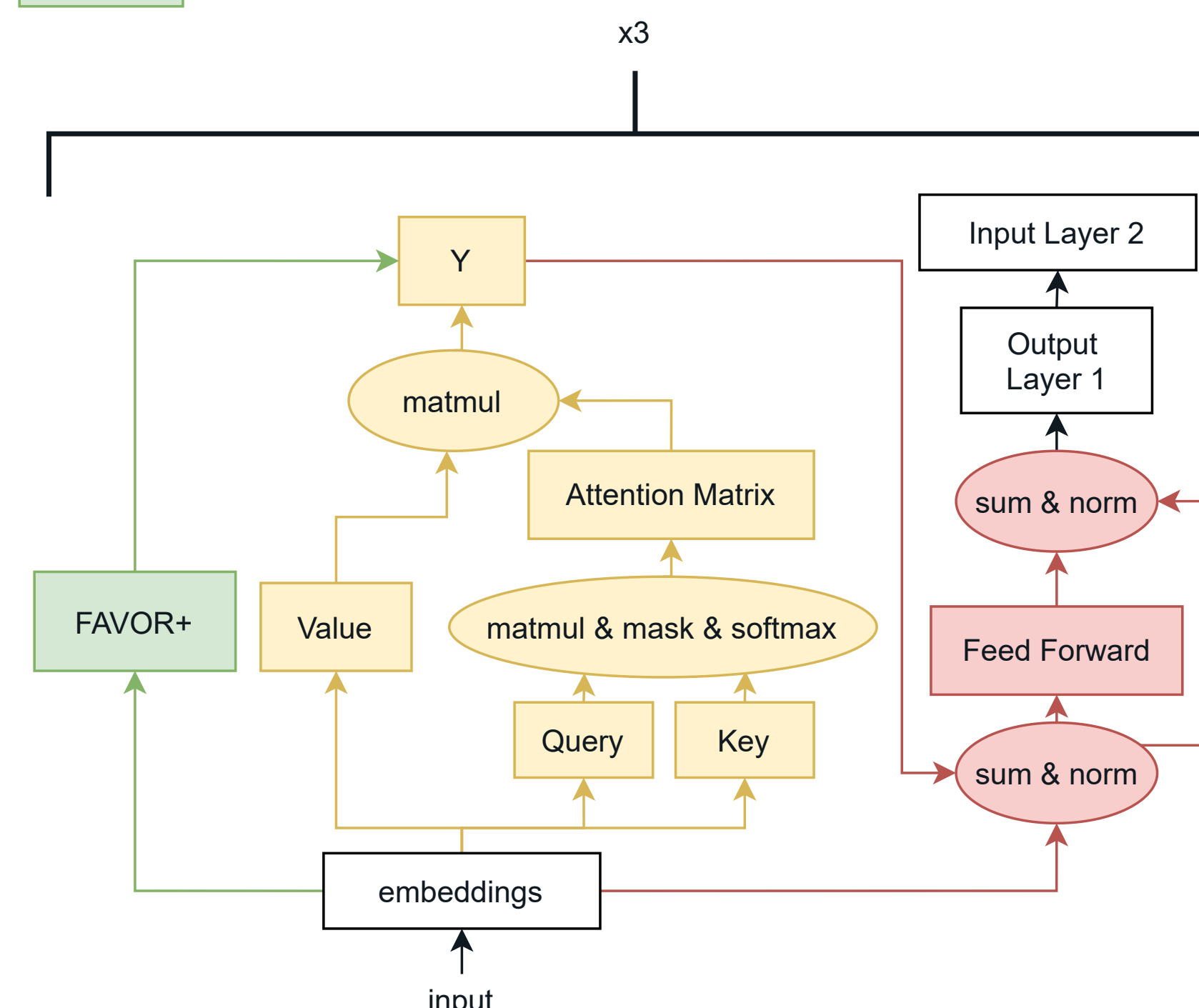
The diagram illustrates the proposed parallelized attention mechanism. It shows the flow from input matrices  $A$  and  $V$  to the final output  $V$ . Matrix  $A$  (blue,  $L \times L$ ) is processed by an attention mechanism (indicated by a red dashed box) with complexity  $O(L^2d)$  to produce  $Q'$  (red,  $L \times r$ ). Matrix  $V$  (yellow,  $L \times d$ ) is also processed by the attention mechanism. The resulting  $Q'$  is then multiplied by the transpose of the key matrix  $(K')^\top$  (green,  $r \times L$ ) to produce the final output  $V$  (yellow,  $L \times d$ ). The complexity of this step is  $O(Lrd)$ . The overall complexity of the mechanism is  $O(Lrd)$ .

$$K_{SM}(x, y) = \exp(xy^T) = \phi(x)^T \phi(y)$$
$$K_{SM}(x, y) = \mathbb{E}_{w \sim \mathcal{N}(0, I_d)} [\exp(-\frac{\|x\|^2}{2}) \exp(-\frac{\|y\|^2}{2}) \exp(w^T x) \exp(w^T y)]$$

$$\Rightarrow \phi(x) = \frac{1}{\sqrt{R}} \exp\left(-\frac{\|x\|^2}{2}\right) (\exp(w_1^T x), \dots, \exp(w_R^T x))$$



The diagram illustrates the components of the Transformer block. It consists of three colored boxes: a yellow box labeled 'Self-Attention', a red box labeled 'Normalisation & Skip-connection', and a green box labeled 'Approximation of Attention'.



Sur des séquences de grandes tailles le performer permet d'améliorer la vitesse de calcul et/ou d'augmenter la taille du batch.

nb iteration en batch	transformer accuracy	performer accuracy
0	0.50	0.52
250	0.75	0.50
500	0.78	0.75
750	0.78	0.78
1000	0.80	0.80
1250	0.78	0.80
1500	0.80	0.82
1750	0.81	0.83
2000	0.81	0.84

Les performances d'un transformer VS un performer ne sont pas différentiables sur une base simple comme IMDB.

Pour le même nombre d'époch l'apprentissage est 50% plus rapide chez le performer.

[1] Krzysztof et al. *Rethinking Attention with Performers* - ICLR 2021  
 [2] Vaswani et al. *Attention Is All You Need* - 2017  
 [2] Yu et al. *Orthogonal Random Features* - 2016  
 [3] Rahimi et al. *Random Features for Large-scale Kernel Machines*