



KOREA UNIVERSITY

**DSBA** Transformer survey paper study

# A Survey of Transformers

## Appendix: Complexity, Parameters, and Scaling

arXiv preprint

---



고려대학교 산업경영공학과

**Data Science & Business Analytics Lab**

이유경, 김명섭, 윤훈상, 김지나, 허재혁, 김수빈

발표자 : 김명섭

# Complexity and Parameters

Matrix, Complexity, and Parameters

- 행렬 연산과 Complexity, Parameter의 관계를 알아보면...

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \times \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} = \begin{bmatrix} 1 \times 1 + 2 \times 2 + 3 \times 3 \\ 4 \times 1 + 5 \times 2 + 6 \times 3 \end{bmatrix} = \begin{bmatrix} 14 \\ 32 \end{bmatrix}$$

Diagram illustrating matrix multiplication with dimensions and arrows:

- Matrix 1:  $(2 \times 3)$  (Blue arrows)
- Matrix 2:  $(3 \times 1)$  (Orange arrows)
- Result:  $(2 \times 1)$  (Green arrows)

$$\text{Computational Complexity (Mult)} = 2 \times 3 \times 1 = 6$$

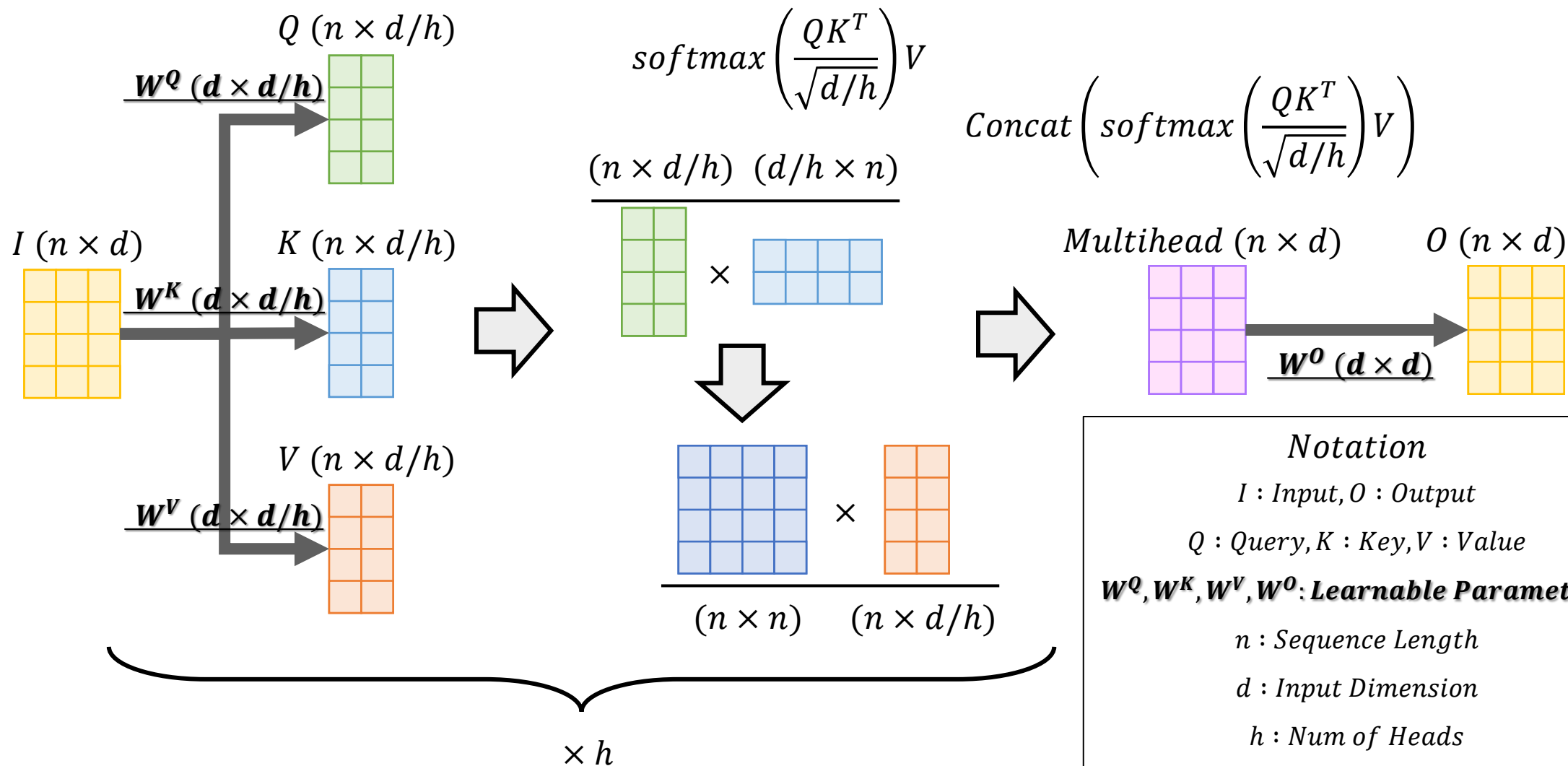
$$\text{Parameters} = 3 \times 1 = 3$$

$(3 \times 1)$  size matrix is **learnable**

# Complexity and Parameters

## Computational Complexity of Self-Attention

Module	Complexity	#Parameters
Self-attention	$O(n^2 \cdot d)$	$4d^2$
Position-wise FFN	$O(n \cdot d^2)$	$8d^2$



### Notation

$I$  : Input,  $O$  : Output

$Q$  : Query,  $K$  : Key,  $V$  : Value

$W^Q, W^K, W^V, W^O$  : **Learnable Parameters**

$n$  : Sequence Length

$d$  : Input Dimension

$h$  : Num of Heads

Multihead : Concatenated Attention Value

# Complexity and Parameters

## Computational Complexity of Self-Attention

Module	Complexity	#Parameters
Self-attention	$O(n^2 \cdot d)$	$4d^2$
Position-wise FFN	$O(n \cdot d^2)$	$8d^2$

$$\text{Query Projection : } I\mathbf{W}^Q \rightarrow (n \times d) \times (\mathbf{d} \times \mathbf{d}/\mathbf{h}) \rightarrow \frac{nd^2}{h}$$

$$\text{Key Projection : } I\mathbf{W}^K \rightarrow (n \times d) \times (\mathbf{d} \times \mathbf{d}/\mathbf{h}) \rightarrow \frac{nd^2}{h}$$

$$\text{Value Projection : } I\mathbf{W}^V \rightarrow (n \times d) \times (\mathbf{d} \times \mathbf{d}/\mathbf{h}) \rightarrow \frac{nd^2}{h}$$

$$\text{Self-Attention : } \text{softmax}\left(\frac{QK^T}{\sqrt{d/h}}\right)V \rightarrow \frac{n^2d}{h} + \frac{n^2d}{h} = \frac{2n^2d}{h}$$

$$\text{Attention Score : } \frac{QK^T}{\sqrt{d/h}} = A \rightarrow (n \times d/h) \times (d/h \times n) \rightarrow \frac{n^2d}{h}$$

$$\text{Attention Value : } \text{softmax}(A)V \rightarrow (n \times n) \times (n \times d/h) \rightarrow \frac{n^2d}{h}$$

$$\text{Output Projection : } \text{Multihead}\mathbf{W}^O \rightarrow (n \times d) \times (\mathbf{d} \times \mathbf{d}) \rightarrow nd^2$$

$\times h$

### Notation

$I$  : Input,  $O$  : Output

$Q$  : Query,  $K$  : Key,  $V$  : Value

$\mathbf{W}^Q, \mathbf{W}^K, \mathbf{W}^V, \mathbf{W}^O$  : **Learnable Parameters**

$n$  : Sequence Length

$d$  : Input Dimension

$h$  : Num of Heads

Multihead : Concatenated Attention Value

$$\text{Multi-Head Self-Attention Computational Complexity} = 4nd^2 + 2n^2d \rightarrow O(nd^2 + n^2d)$$

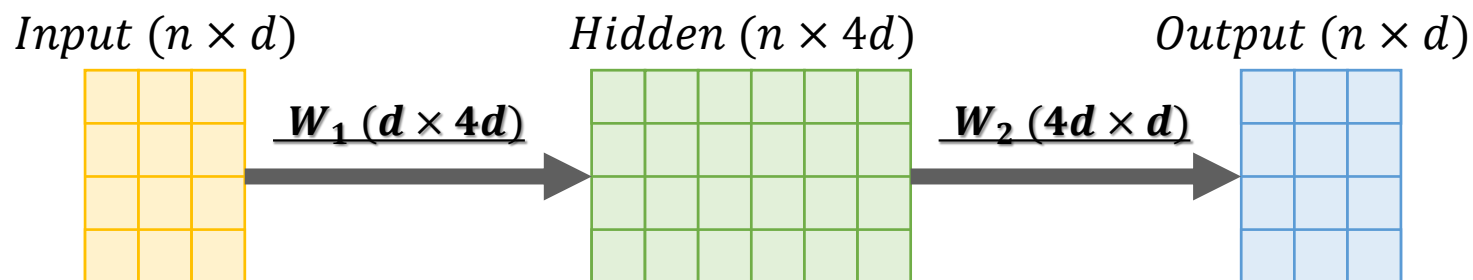
$$\text{Self-Attention Computational Complexity} = 2n^2d \rightarrow O(n^2d)$$

$$\text{Num of Parameters} = \underline{4d^2}$$

# Complexity and Parameters

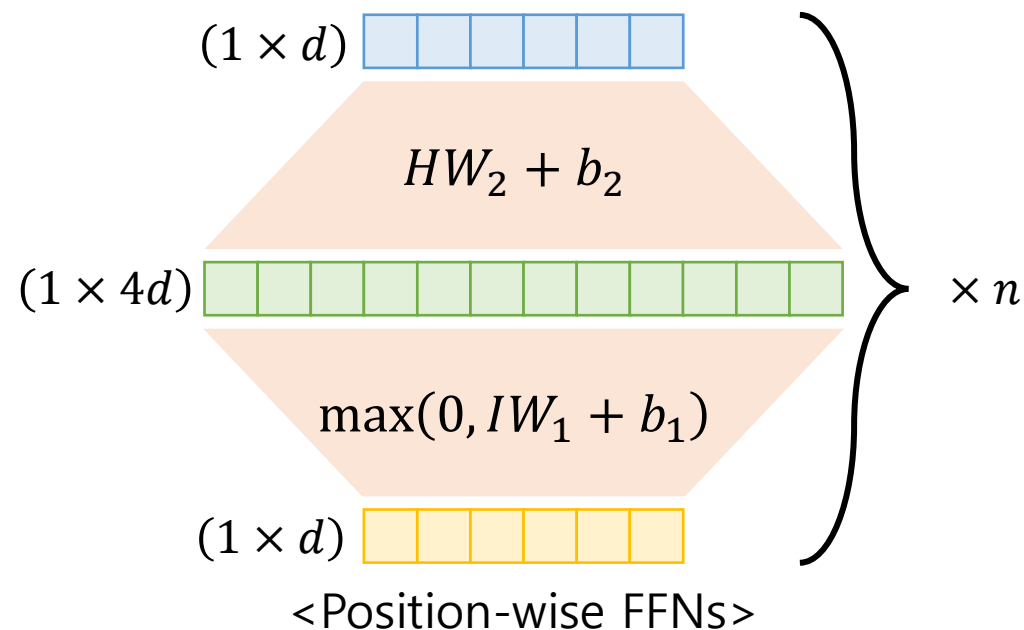
## Computational Complexity of Position-wise FFN

Module	Complexity	#Parameters
Self-attention	$O(n^2 \cdot d)$	$4d^2$
Position-wise FFN	$O(n \cdot d^2)$	$8d^2$



Input to Hidden :  $I\mathbf{W}_1 = H \rightarrow (n \times d) \times (d \times 4d) \rightarrow 4nd^2$

Hidden to Output :  $H\mathbf{W}_2 \rightarrow (n \times 4d) \times (4d \times d) \rightarrow 4nd^2$



### Notation

$I$  : Input,  $H$  : Hidden,  $O$  : Output

$\mathbf{W}_1, \mathbf{W}_2$ : **Learnable Parameters**

$n$  : Sequence Length

$d$  : Input Dimension

**Position-wise FFN Computational Complexity =  $8nd^2 \rightarrow O(nd^2)$**

**Num of Parameters =  $\underline{8d^2}$**

**Question:**

Therefore, the total complexity of the layer is  $O(n^2 d + n d^2)$ , which is worse than that of a traditional RNN layer. I obtained the same result for multi-headed attention too, on considering the appropriate intermediate representation dimensionalities ( $d_k$ ,  $d_v$ ) and finally multiplying by the number of heads  $h$ .

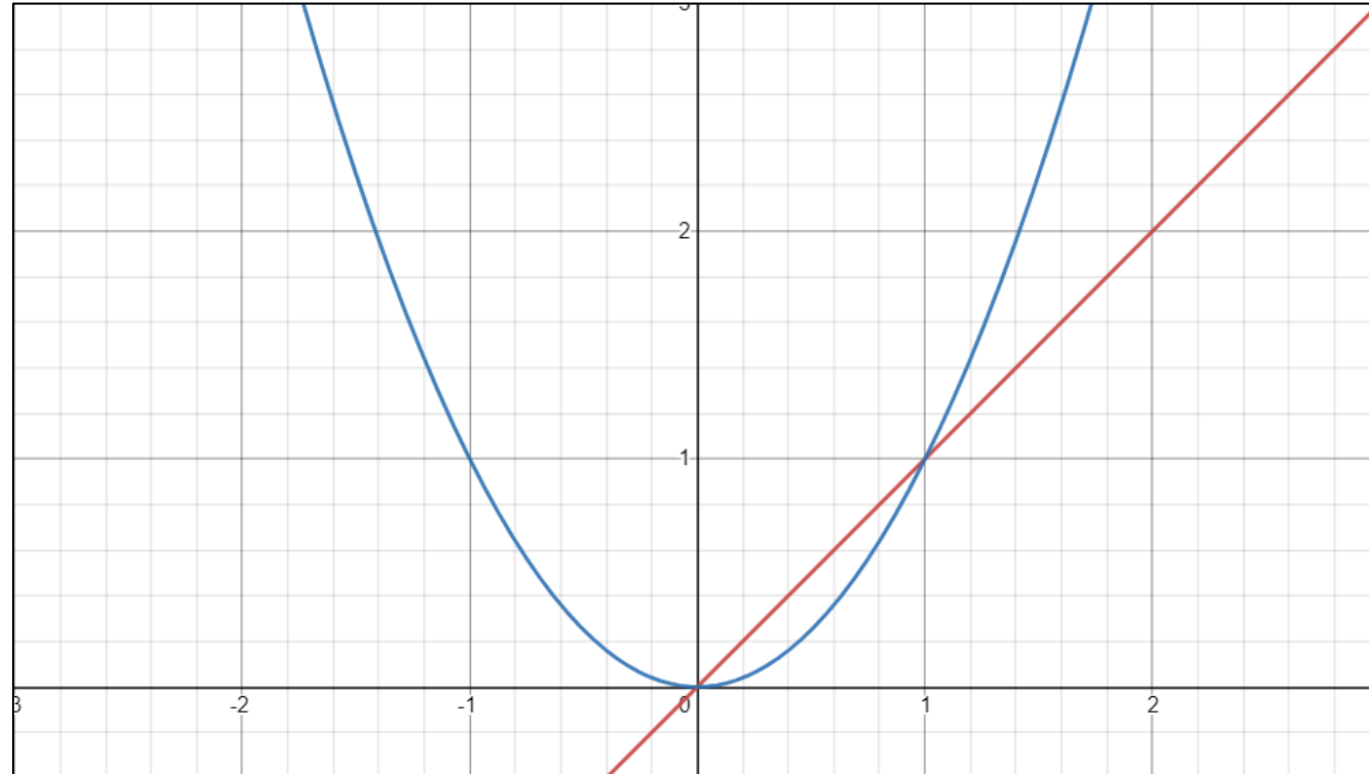
**Answer:**

When the original [Attention paper](#) was first introduced, it didn't require to calculate  $q$ ,  $v$  and  $k$  matrices, as the values were taken directly from the hidden states of the RNNs, and thus the complexity of Attention layer is  $O(n^2 \cdot d)$ .

Now, to understand what Table 1 contains please keep in mind how most people scan papers: they read title, abstract, then look at figures and tables. Only then if the results were interesting, they read the paper more thoroughly. So, the main idea of the [Attention is all you need](#) paper was to replace the RNN layers completely with attention mechanism in seq2seq setting because RNNs were really slow to train. If you look at the Table 1 in this context, you see that it compares RNN, CNN and Attention and highlights the motivation for the paper: using Attention should have been beneficial over RNNs and CNNs. It should have been advantageous in 3 aspects: constant amount of calculation steps, constant amount of operations and lower computational complexity for usual Google setting, where  $n \approx 100$  and

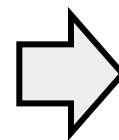
# Complexity and Parameters

## Quadratic Computation Complexity



$$Q (n \times d_k) \times K (d_k \times n) = \text{Attention Score Matrix} (n \times n)$$

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



$$q (1 \times d_k) \times k (d_k \times 1) = \text{Attention Score (scalar)}$$

$$\text{Attention Score} = q \cdot k = \sum_{i=1}^{d_k} q_i k_i$$

$$\text{if } q_i \sim E[q_i] = 0, \text{Var}[q_i] = 1,$$

$$k_i \sim E[k_i] = 0, \text{Var}[k_i] = 1, q_i \perp k_i,$$

$$\rightarrow E[q \cdot k] = 0, \text{Var}[q \cdot k] = d_k$$

❖ if all random variables are independent,

$$\checkmark E[\sum_{i=1}^n X_i] = \sum_{i=1}^n E[X_i]$$

$$\checkmark E[\prod_{i=1}^n X_i] = \prod_{i=1}^n E[X_i]$$

$$\checkmark \text{Var}[\sum_{i=1}^n X_i] = \sum_{i=1}^n \text{Var}[X_i]$$

$$\checkmark \text{Var}[\prod_{i=1}^n X_i] = \prod_{i=1}^n (\text{Var}[X_i] + (E[X_i])^2) - \prod_{i=1}^n (E[X_i])^2$$



**Algorithm 1** Disentangled Attention

**Input:** Hidden state  $H$ , relative distance embedding  $P$ , relative distance matrix  $\delta$ . Content projection matrix  $W_{k,c}$ ,  $W_{q,c}$ ,  $W_{v,c}$ , position projection matrix  $W_{k,r}$ ,  $W_{q,r}$ .

```

1:  $K_c = HW_{k,c}$ ,  $Q_c = HW_{q,c}$ ,  $V_c = HW_{v,c}$ ,  $K_r = PW_{k,r}$ ,  $Q_r = PW_{q,r}$ 
2:  $A_{c \rightarrow c} = Q_c K_c^\top$ 
3: for  $i = 0, \dots, N - 1$  do
4:    $\tilde{A}_{c \rightarrow p}[i, :] = Q_c[i, :] K_r^\top$ 
5: end for
6: for  $i = 0, \dots, N - 1$  do
7:   for  $j = 0, \dots, N - 1$  do
8:      $A_{c \rightarrow p}[i, j] = \tilde{A}_{c \rightarrow p}[i, \delta[i, j]]$ 
9:   end for
10: end for
11: for  $j = 0, \dots, N - 1$  do
12:    $\tilde{A}_{p \rightarrow c}[:, j] = K_c[:, j] Q_r^\top$ 
13: end for
14: for  $j = 0, \dots, N - 1$  do
15:   for  $i = 0, \dots, N - 1$  do
16:      $A_{p \rightarrow c}[i, j] = \tilde{A}_{p \rightarrow c}[\delta[j, i], j]$ 
17:   end for
18: end for
19:  $\tilde{A} = A_{c \rightarrow c} + A_{c \rightarrow p} + A_{p \rightarrow c}$ 
20:  $H_o = \text{softmax}(\frac{\tilde{A}}{\sqrt{3d}}) V_c$ 

```

**Output:**  $H_o$