# Attention is all you need

**Author:** Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L.,
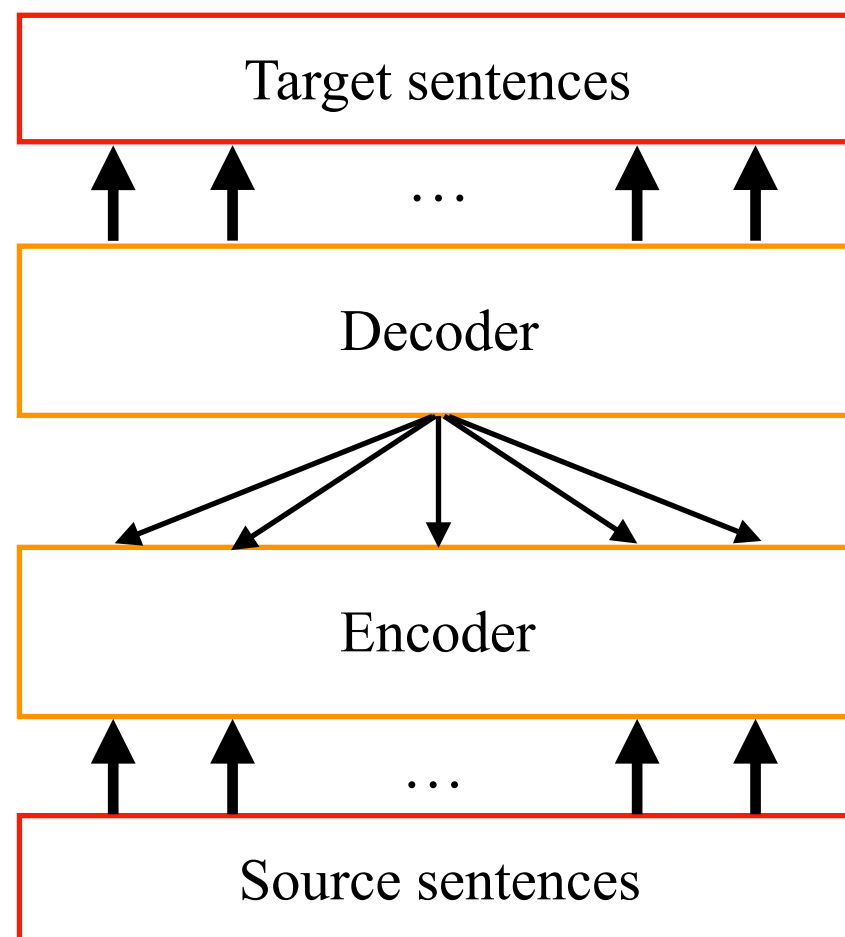Gomez, A. N., ... & Polosukhin, I..

Presenter: WENWEI KANG

- Introduction

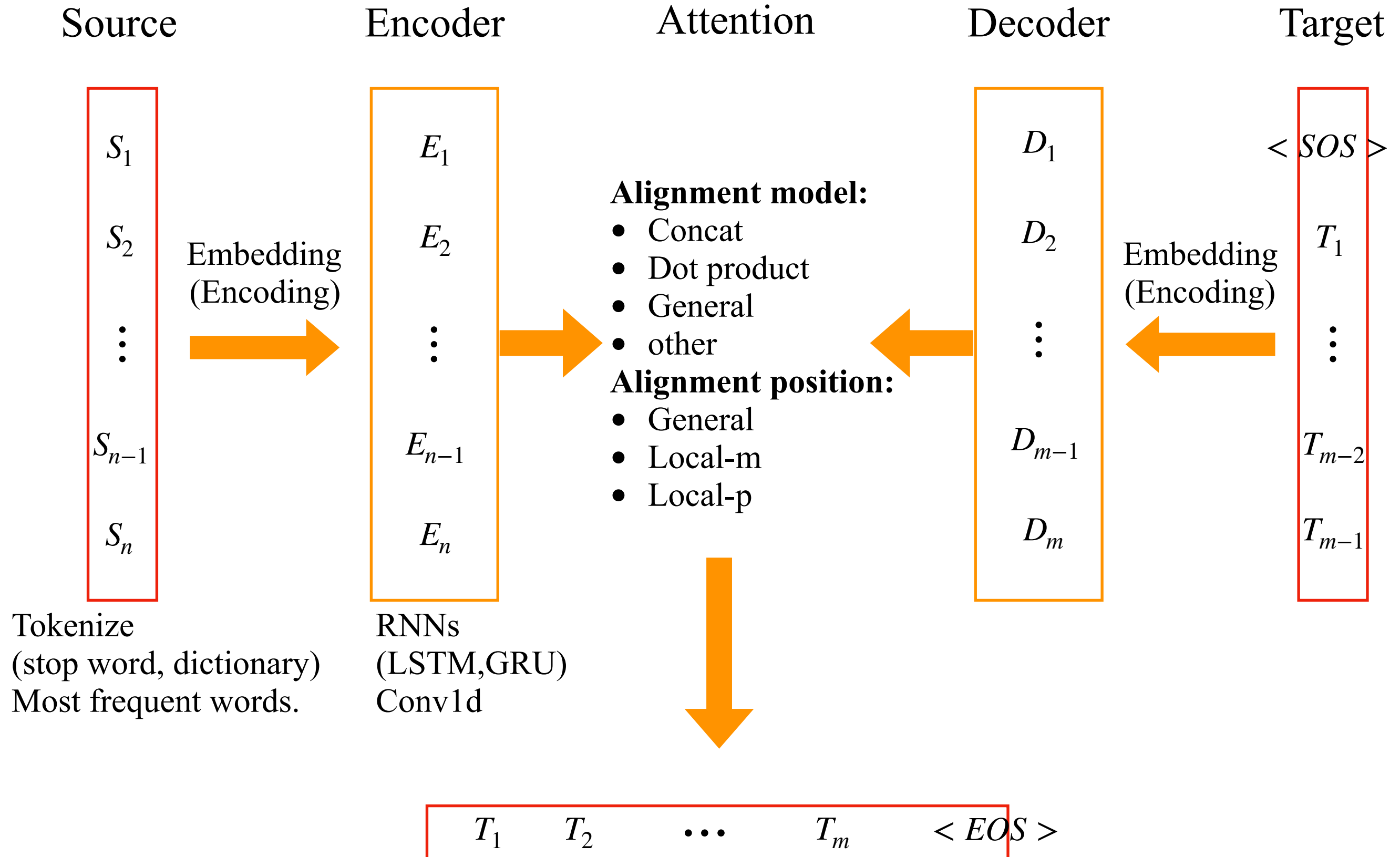- Encoder - Decoder

- Transformer

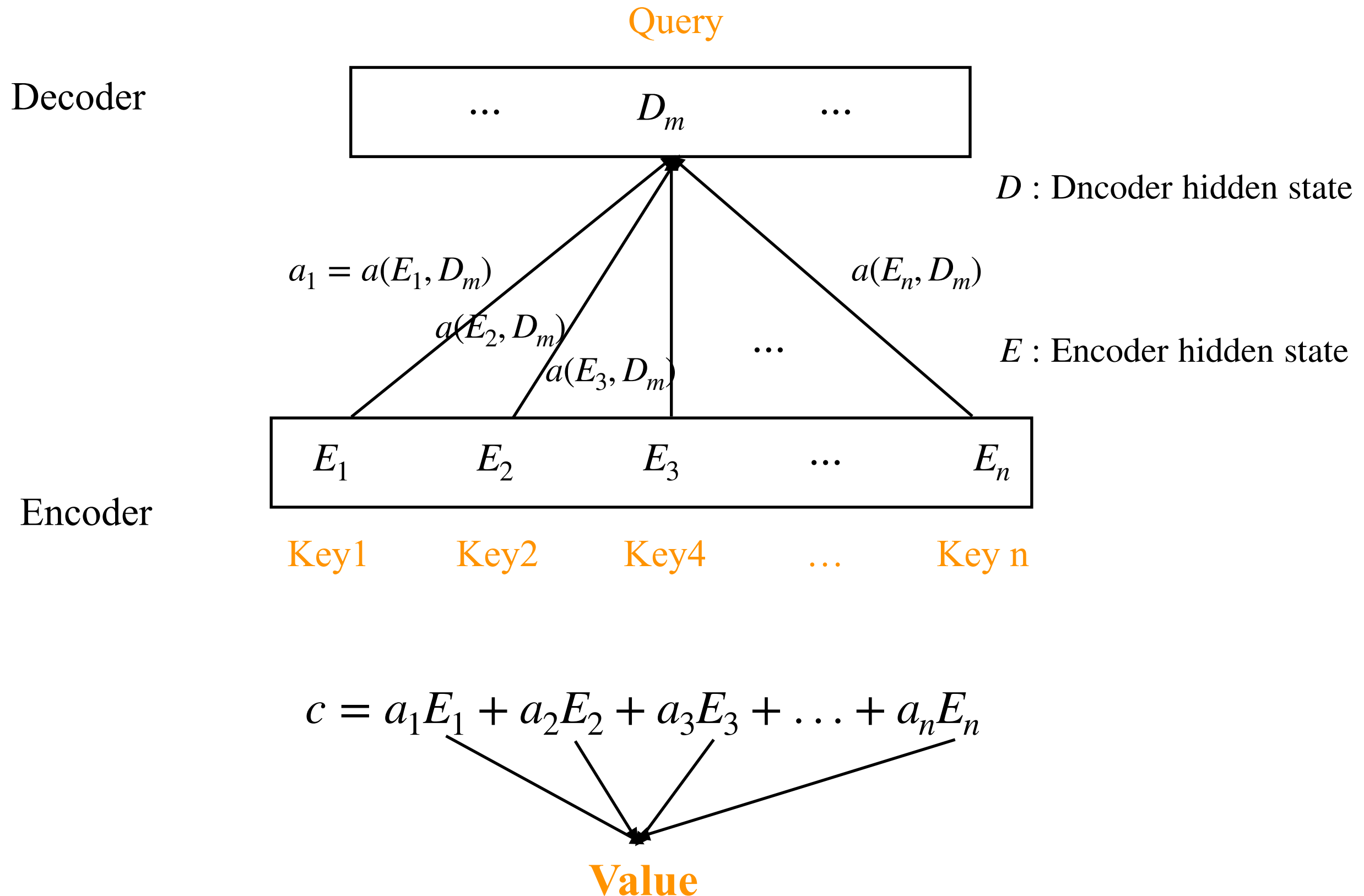- Evaluation

# Introduction

Neural Machine Translation(NMT):

- **Statistical based:** Phrase-based + large LM (Moses)

- **NN based:** Encoder - Decoder (Seq2seq, ConvS2S, ensemble …)

# Encoder - Decoder

| Source | Encoder | Attention | Decoder | Target |
|---|---|---|---|---|

$S_1$

$S_2$

$\vdots$

$S_{n-1}$

$S_n$

Tokenize
(stop word, dictionary)
Most frequent words.

Embedding
(Encoding)

$E_1$

$E_2$

$\vdots$

$E_{n-1}$

$E_n$

RNNs
(LSTM,GRU)
Conv1d

**Alignment model:**
- Concat
- Dot product
- General
- other

**Alignment position:**
- General
- Local-m
- Local-p

$D_1$

$D_2$

$\vdots$

$D_{m-1}$

$D_m$

Embedding
(Encoding)

$< SOS >$

$T_1$

$\vdots$

$T_{m-2}$

$T_{m-1}$

| $T_1$ | $T_2$ | $\cdots$ | $T_m$ | $< EOS >$ |

# Encoder - Decoder

Query

Decoder

| ... | $D_m$ | ... |

$D$ : Dncoder hidden state

$a_1 = a(E_1, D_m)$

$a(E_2, D_m)$

$a(E_3, D_m)$

$a(E_n, D_m)$

...

$E$ : Encoder hidden state

| $E_1$ | $E_2$ | $E_3$ | ... | $E_n$ |

Encoder

Key1        Key2        Key4        …        Key n

$$c = a_1 E_1 + a_2 E_2 + a_3 E_3 + \ldots + a_n E_n$$

**Value**

# Transformer



**Encoder**

**Decoder**

# Transformer

$$Q \quad K \quad V$$

(Query)    (Key)    (Value)



Positional
Encoding

Input
Embedding

Inputs

**Positional Encoding**

$$PE = \{p_0, p_1, p_2, \ldots, p_{n-1}\}$$

$$
\begin{cases}
PE_{(pos,2i)} = \sin(pos/10000^{\frac{2i}{d_{model}}}) \\
PE_{(pos,2i+1)} = \cos(pos/10000^{\frac{2i}{d_{model}}})
\end{cases}
$$

**Word Embedding**

$$W = \{w_1, w_2, \ldots, w_n\}, w_i \in d_{model}$$
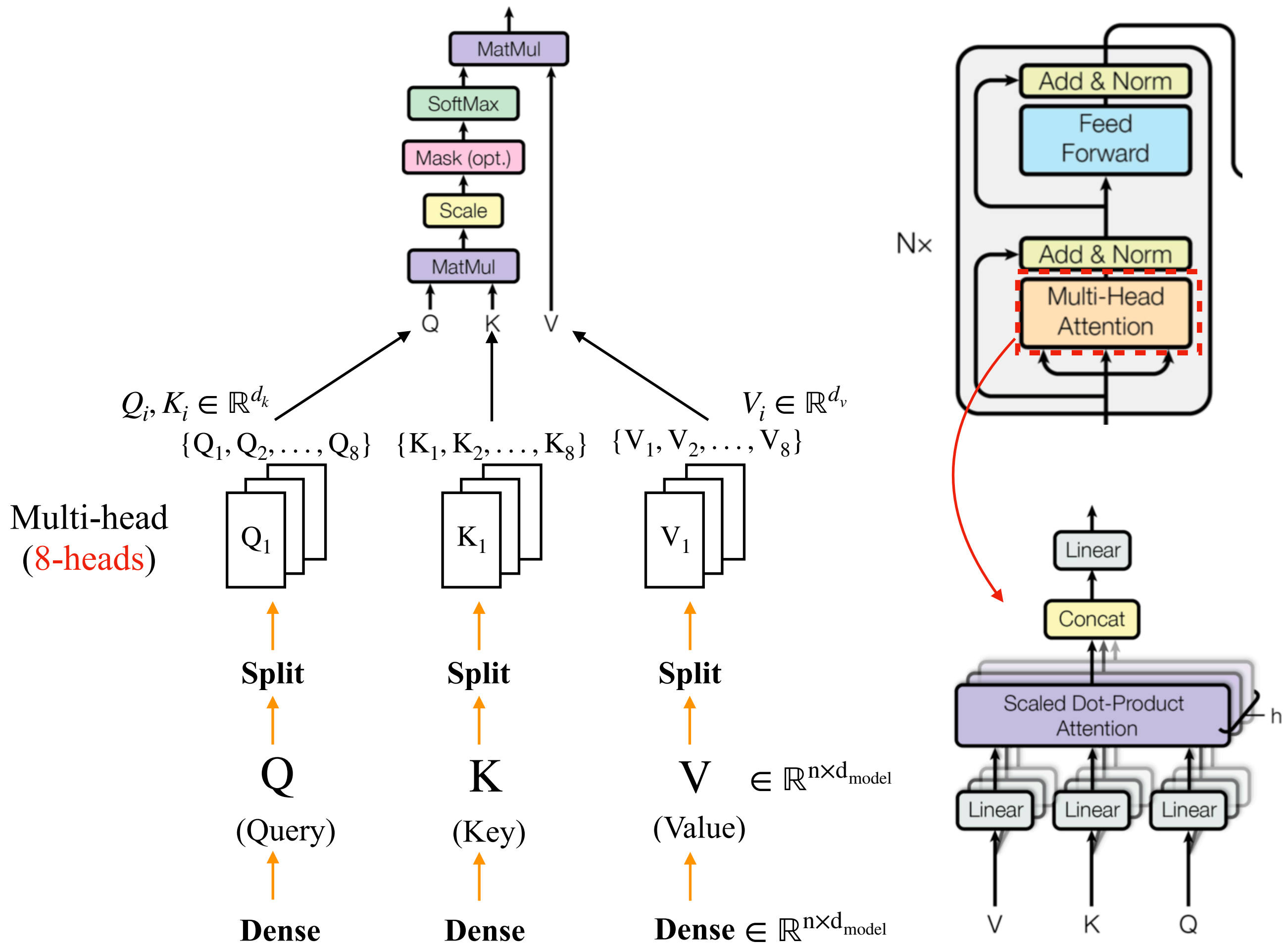
# Positional Encoding

Positional Encoding: $\{p_0, p_1, \ldots, p_{49}\}$, $p_i \in d_{model}$

1. Calculates the inner product $p_{24}$ to the others $p_0$ $to$ $p_{49}$.
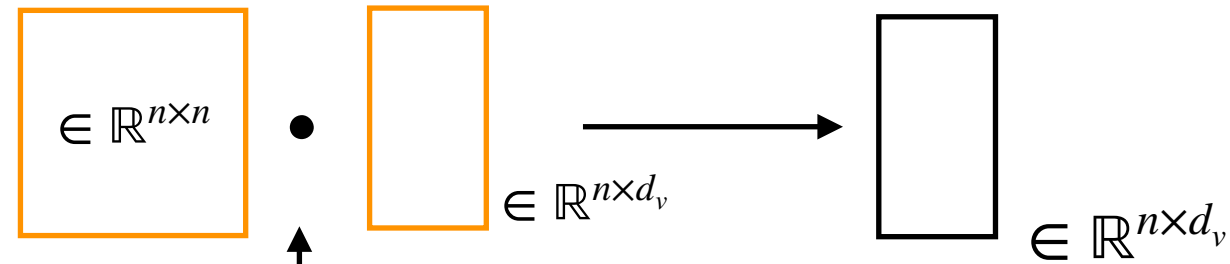2. Get product $V_0, V_1, \ldots, V_{49}$.



$V$

MatMul

SoftMax

Mask (opt.)

Scale

MatMul

Q        K        V

$Q_i, K_i \in \mathbb{R}^{d_k}$

$V_i \in \mathbb{R}^{d_v}$

$\{Q_1, Q_2, \ldots, Q_8\}$   $\{K_1, K_2, \ldots, K_8\}$   $\{V_1, V_2, \ldots, V_8\}$

Multi-head
(8-heads)

$Q_1$        $K_1$        $V_1$

**Split**        **Split**        **Split**

Q        K        V $\in \mathbb{R}^{n \times d_{model}}$

(Query)        (Key)        (Value)

**Dense**        **Dense**        **Dense** $\in \mathbb{R}^{n \times d_{model}}$

Add & Norm

Feed
Forward

N×

Add & Norm

Multi-Head
Attention

Linear

Concat

Scaled Dot-Product
Attention

h

Linear        Linear        Linear

V        K        Q

9

# Transformer

Multi-heads attention

$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

$\in \mathbb{R}^{n \times n}$ • $\in \mathbb{R}^{n \times d_v}$ → $\in \mathbb{R}^{n \times d_v}$

Softmax

Scale    Divided by $\sqrt{d_k}$
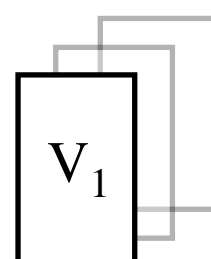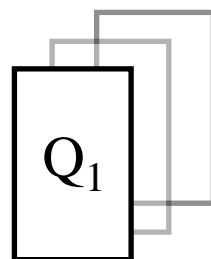
MatMul    $\in \mathbb{R}^{n \times n}$    *Self-attention*

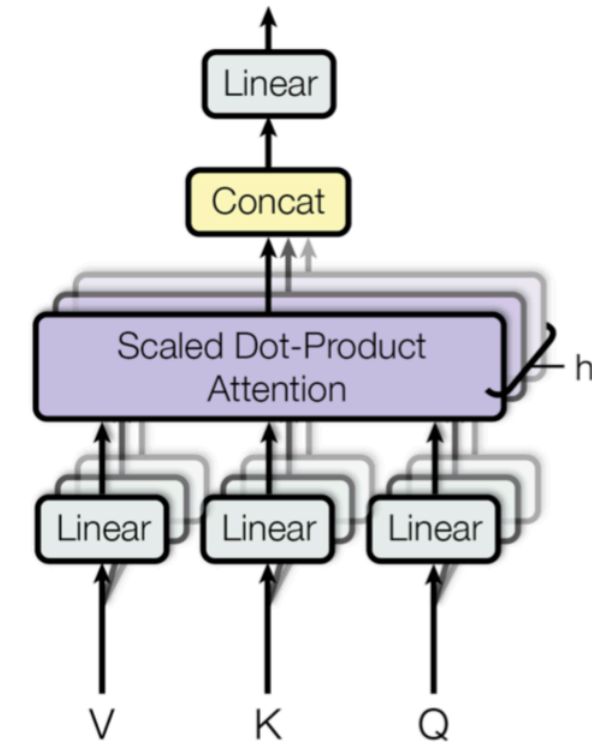{Q₁}    {K₁}    {V₁}

One-head    Q₁    K₁    V₁

Scaled Dot-Product Attention
(Self-Attention)

# Transformer



**Residual connection**

$$\text{Normalize}\left( \quad \cdots \quad + \begin{array}{c} Q \\ \text{(Query)} \end{array} \right)$$

$\in \mathbb{R}^{n \times d_{model}}$ $\qquad$ $\in \mathbb{R}^{n \times d_{model}}$

$$= \text{Normalize}\left( \quad \in \mathbb{R}^{n \times d_{model}} \quad \right)$$

Layer Normalize: $LN(z; \alpha, \beta) = \dfrac{(z - \mu)}{\sigma} \odot \alpha + \beta$

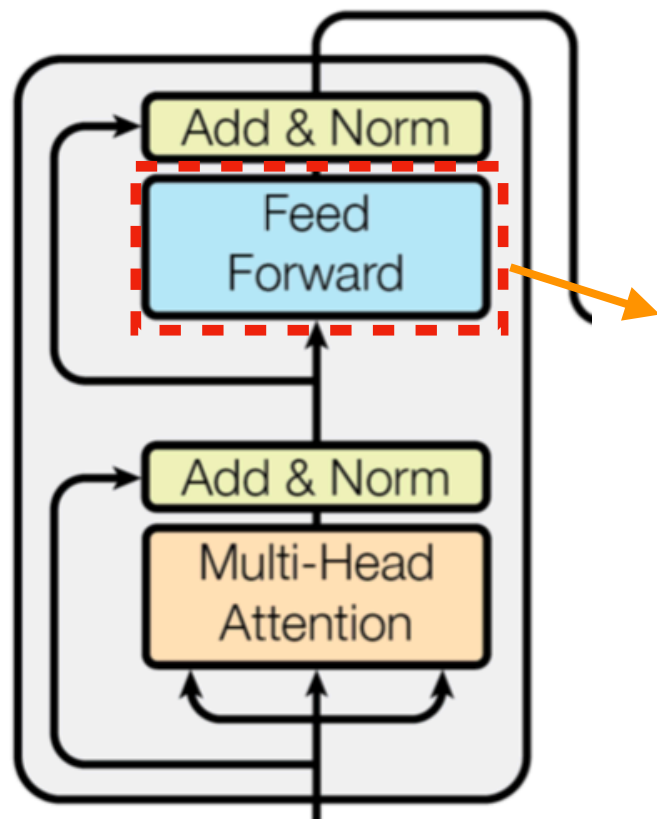Mean : $\mu^l = \dfrac{1}{D} \sum\limits_{i=1}^{D} z_i^l$ $\qquad$ Gains : $\alpha$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ Biases : $\beta$

Standard Deviation : $\sigma^l = \sqrt{\dfrac{1}{D} \sum\limits_{i=1}^{D} (z_i^l - \mu^l)^2}$

# Transformer

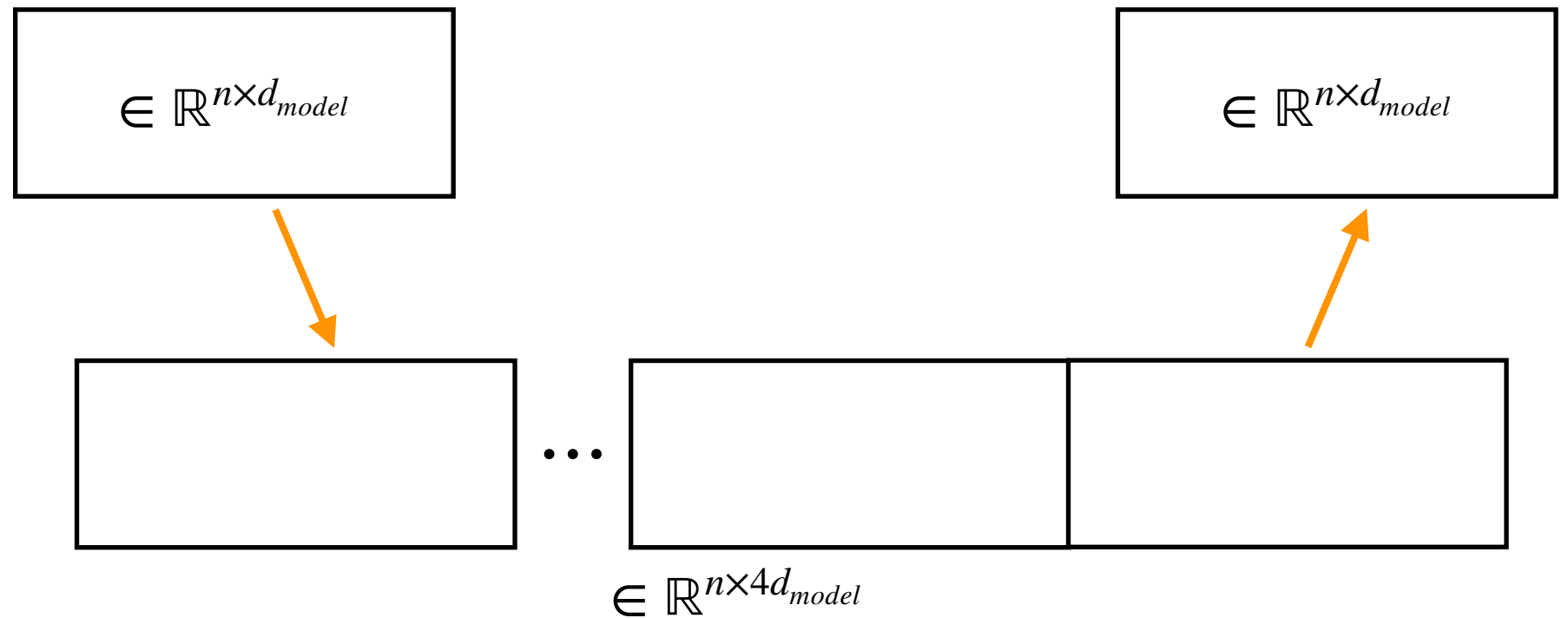**Feed Forward**
**(Dense、Conv1d)**
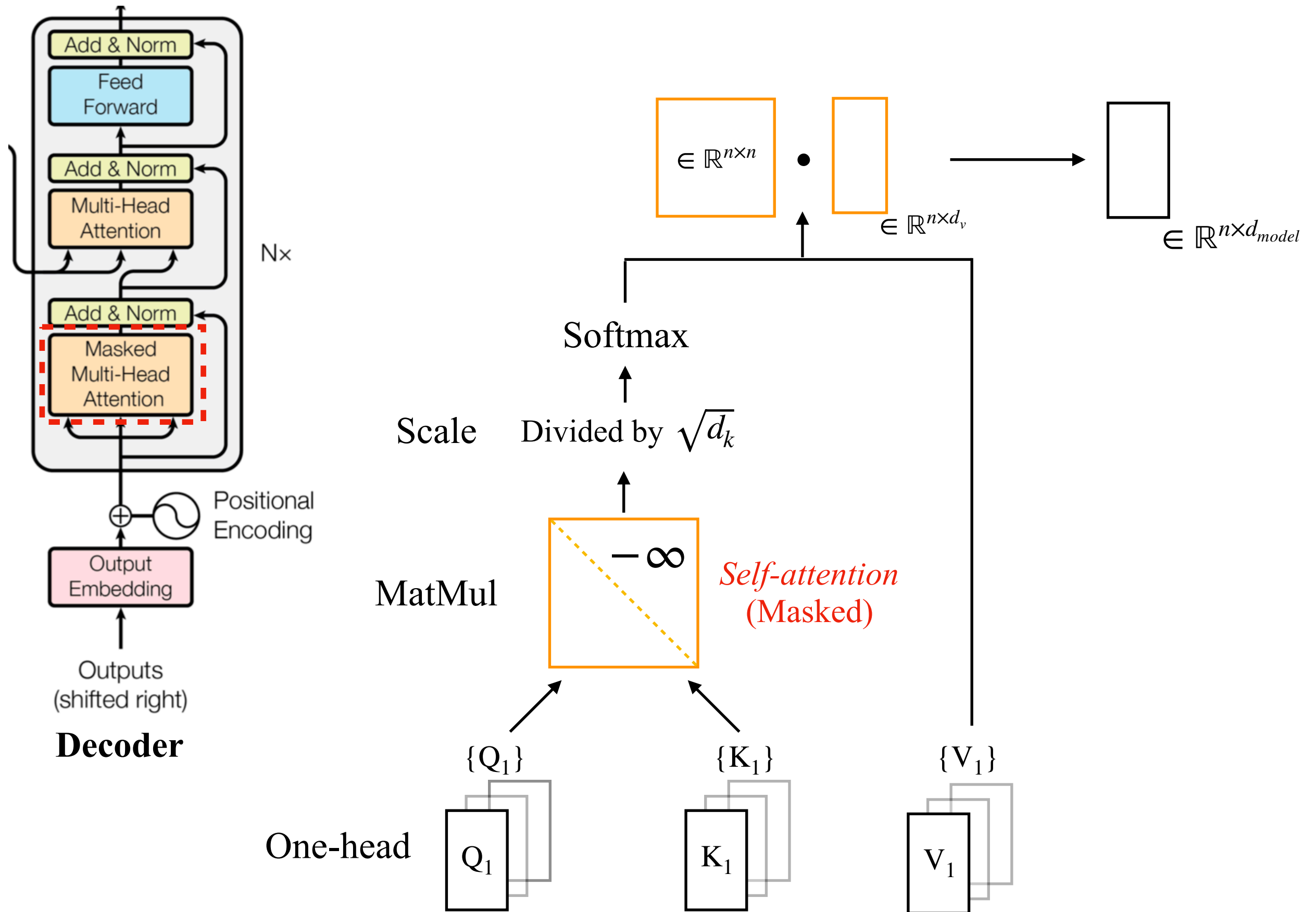
$$\text{FFN}(x) = max(0, xW_1 + b_1)W_2 + b_2$$

$$\text{Dense}(d_{model}) \rightarrow \text{Dense}(4 \times d_{model}) \rightarrow \text{Dense}(d_{model})$$



$\in \mathbb{R}^{n \times d_{model}}$

$\in \mathbb{R}^{n \times d_{model}}$

...

$\in \mathbb{R}^{n \times 4d_{model}}$

**Encoder**

# Transformer



**Decoder**

Scale — Divided by $\sqrt{d_k}$

*Self-attention* (Masked)

MatMul

One-head

$\{Q_1\}$  $\{K_1\}$  $\{V_1\}$

$Q_1$  $K_1$  $V_1$

$\in \mathbb{R}^{n \times n}$  $\in \mathbb{R}^{n \times d_v}$  $\in \mathbb{R}^{n \times d_{model}}$

Softmax

$-\infty$

# Transformer



Add & Norm
Feed Forward
Add & Norm
Multi-Head Attention
Add & Norm
Masked Multi-Head Attention
N×
Positional Encoding
Output Embedding
Outputs (shifted right)

**Decoder**

Softmax $\longrightarrow$ $0$

Scale  Divided by $\sqrt{d_k}$

We need to prevent leftward information flow in the decoder

$M_{1,1}$: The similarity between $Q_1$ and $K_1$

$$\begin{matrix} M_{1,1} & & & & \\ M_{2,1} & M_{2,2} & -\infty & & \\ \vdots & M_{3,2} & & & \\ \vdots & \vdots & \ddots & & \\ \vdots & \vdots & \ddots & \ddots & \\ M_{n,1} & M_{n,2} & \cdots & \cdots & M_{n,n} \end{matrix}$$

*Self-attention* (Masked)

$w_1$ $w_2$ $w_3$ $w_4$ $w_5$ $\ldots$ $w_n$

$w_1$ $w_2$ $w_3$ $w_4$ $w_5$ $\ldots$ $w_n$

$w_1$ $w_2$ $w_3$ $w_4$ $w_5$ $\ldots$ $w_n$

$\vdots$

$w_1$ $w_2$ $w_3$ $w_4$ $w_5$ $\ldots$ $w_n$

$\{Q_1\}$          $\{K_1\}$
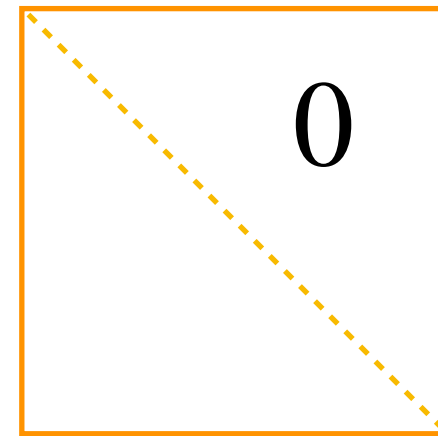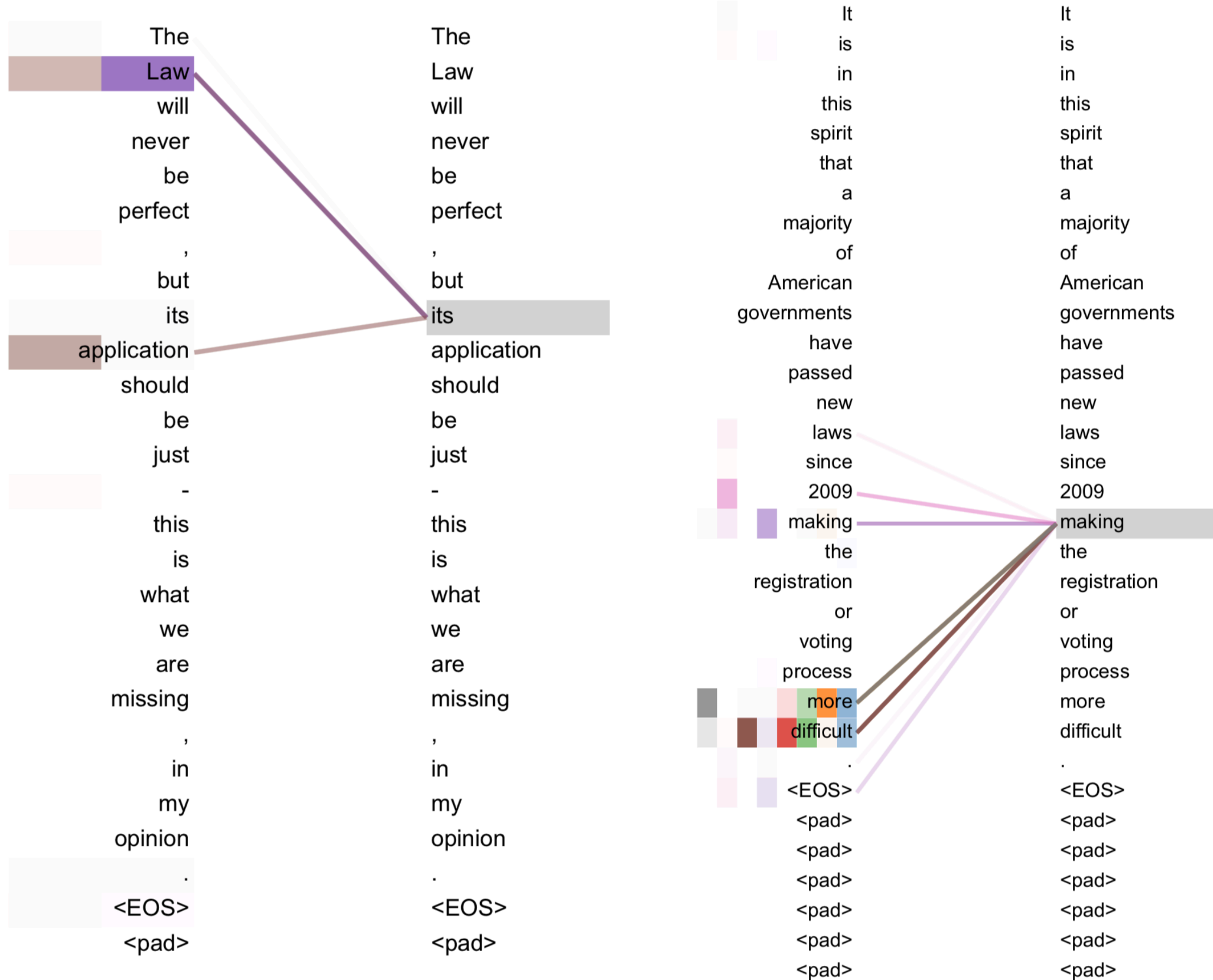
One-head      $Q_1$          $K_1$

# Evaluation

Table 2: The Transformer achieves better BLEU scores than previous state-of-the-art models on the English-to-German and English-to-French newstest2014 tests at a fraction of the training cost.

| Model | BLEU | | Training Cost (FLOPs) | |
|---|---|---|---|---|
| | EN-DE | EN-FR | EN-DE | EN-FR |
| ByteNet [18] | 23.75 | | | |
| Deep-Att + PosUnk [39] | | 39.2 | | $1.0 \cdot 10^{20}$ |
| GNMT + RL [38] | 24.6 | 39.92 | $2.3 \cdot 10^{19}$ | $1.4 \cdot 10^{20}$ |
| ConvS2S [9] | 25.16 | 40.46 | $9.6 \cdot 10^{18}$ | $1.5 \cdot 10^{20}$ |
| MoE [32] | 26.03 | 40.56 | $2.0 \cdot 10^{19}$ | $1.2 \cdot 10^{20}$ |
| Deep-Att + PosUnk Ensemble [39] | | 40.4 | | $8.0 \cdot 10^{20}$ |
| GNMT + RL Ensemble [38] | 26.30 | 41.16 | $1.8 \cdot 10^{20}$ | $1.1 \cdot 10^{21}$ |
| ConvS2S Ensemble [9] | 26.36 | **41.29** | $7.7 \cdot 10^{19}$ | $1.2 \cdot 10^{21}$ |
| Transformer (base model) | 27.3 | 38.1 | $\mathbf{3.3 \cdot 10^{18}}$ | |
| Transformer (big) | **28.4** | **41.8** | $2.3 \cdot 10^{19}$ | |

- **ByteNet:** 2 convolution layers
- **Deep-Att + PosUnk:** 2 Bi-LSTM layers(Encoder) + 1 LSTM layer (Decoder)
- **GNMT + RL:** 7 LSTM layers + 1 Bi-LSTM layer(Encoder) + 8 LSTM layers(Decoder)
- **Transformer(base):** training 100,000 steps(12 hours), 0.4 seconds per steps
- **Transformer(big):** training 300,000 steps(3.5 days), 1.0 seconds per steps

# Evaluation

Source sentence:
 Taiwan is a beautiful country.

# Evaluation

Predict sentence:
 台湾是一个美丽的国家。