# How Transformer Model Works
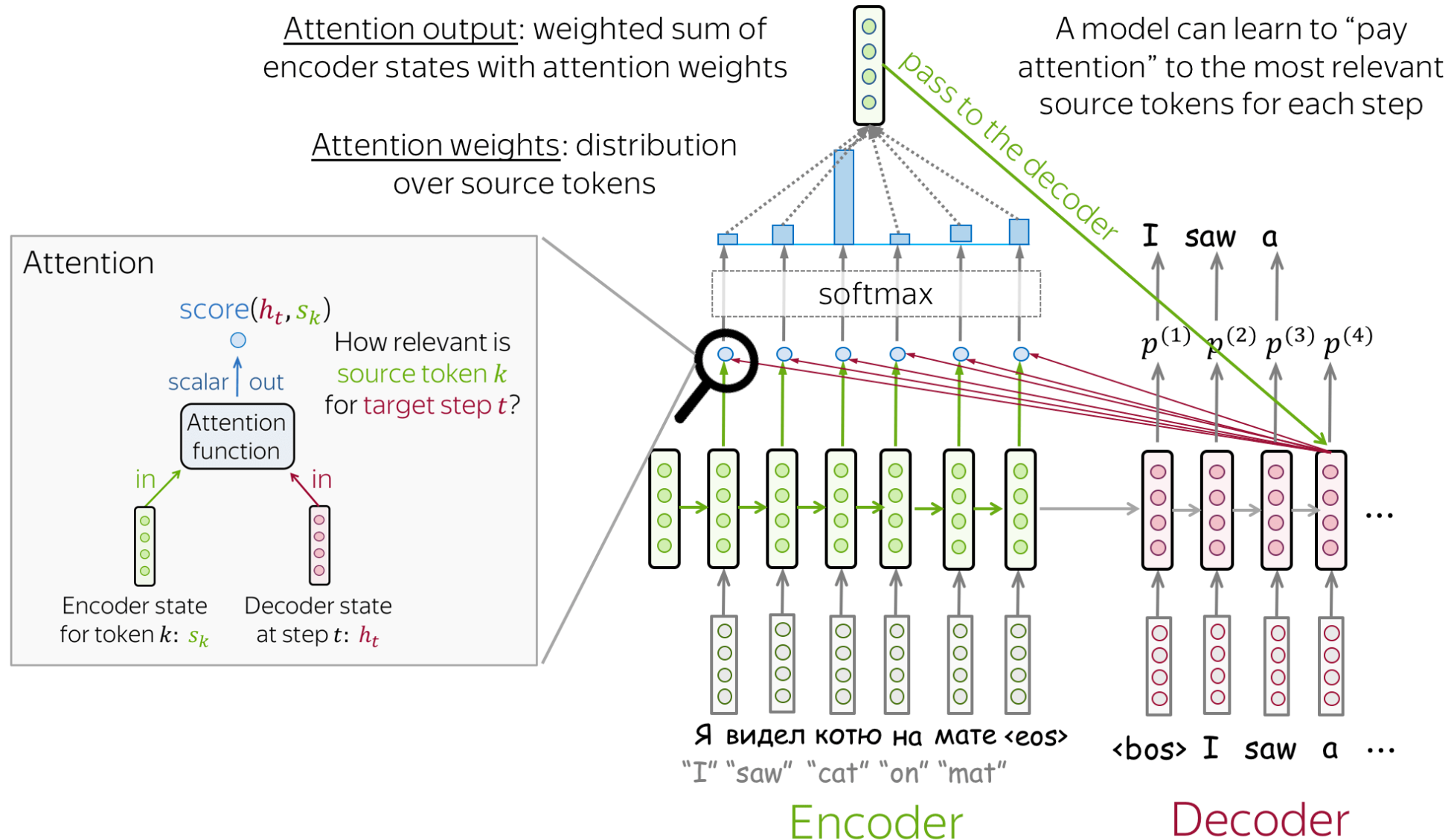
주재걸 교수

KAIST 김재철AI대학원

# Review: Seq2Seq with Attention ← RNN 기반

**Attention output**: weighted sum of encoder states with attention weights

**Attention weights**: distribution over source tokens

A model can learn to "pay attention" to the most relevant source tokens for each step



## Attention

$$\text{score}(h_t, s_k)$$

How relevant is source token $k$ for target step $t$?

scalar | out

Attention function

in | in

Encoder state for token $k$: $s_k$    Decoder state at step $t$: $h_t$

pass to the decoder

softmax

I saw a

$p^{(1)}$ $p^{(2)}$ $p^{(3)}$ $p^{(4)}$

Я видел котю на мате <eos>
"I" "saw" "cat" "on" "mat"

<bos> I saw a …

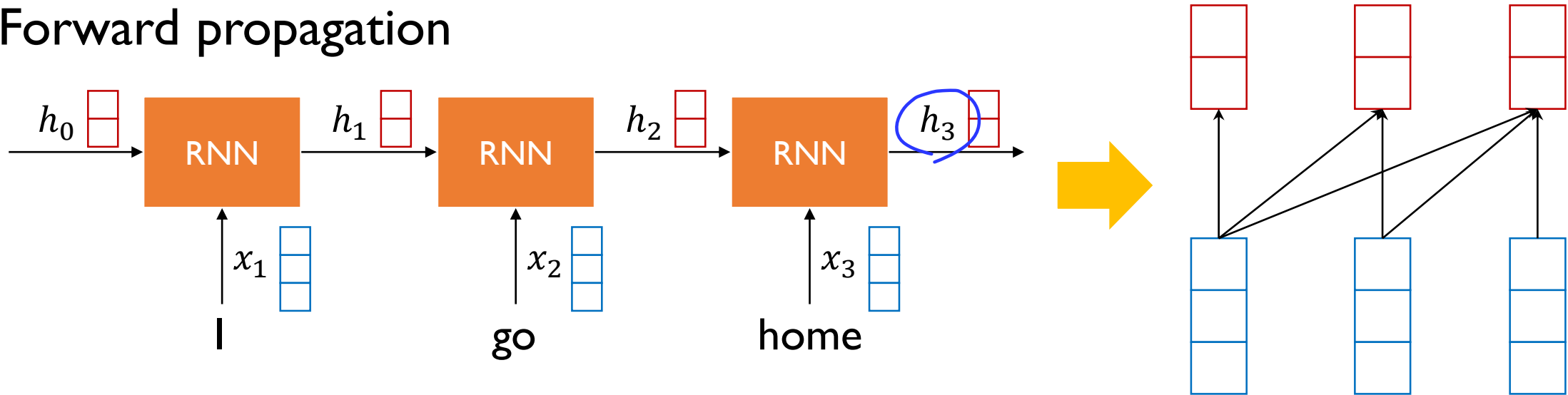**Encoder**          **Decoder**

# Transformer: High-level View

- Attention module can work as both a sequence encoder and a decoder in seq2seq with attention.

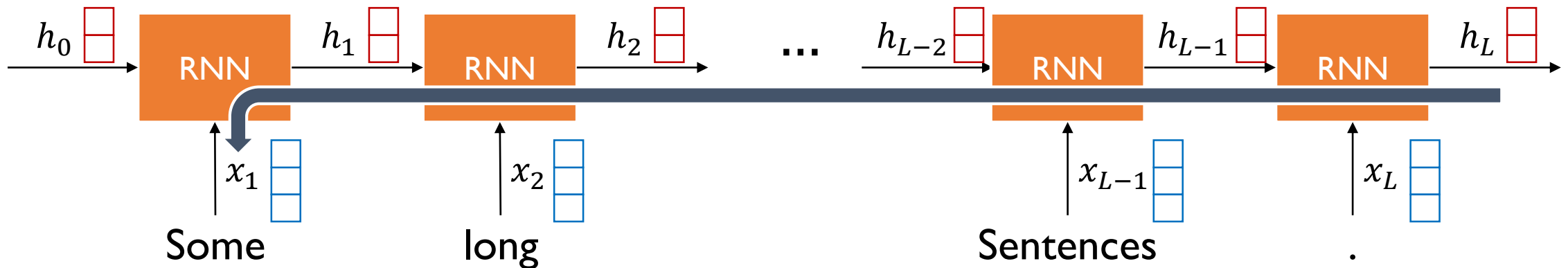- In other words, RNNs or CNNs are no longer necessary, but all we need is attention modules.

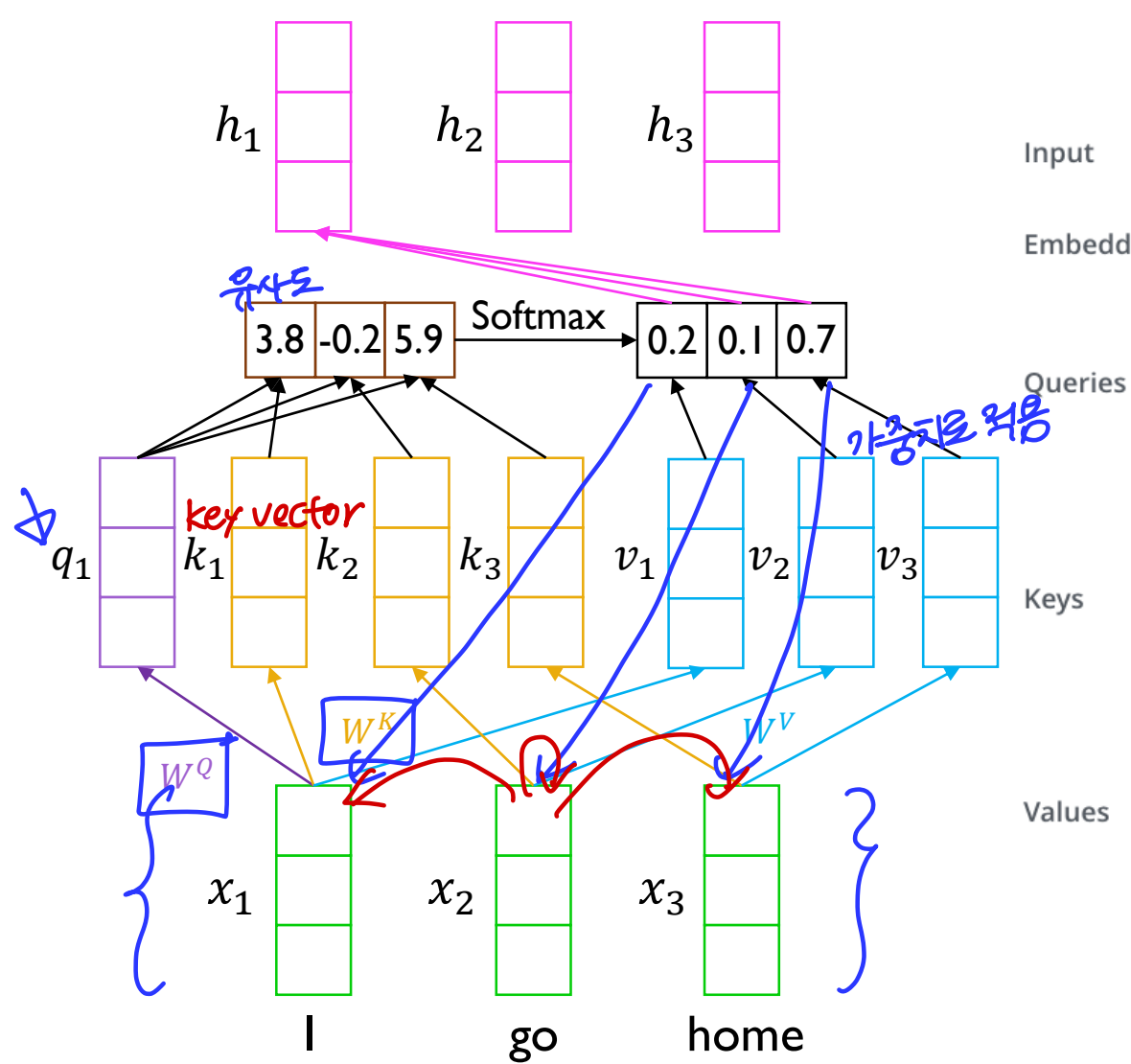# Long-term Dependency Issue of RNN Models

## Forward propagation



## Backpropagation

# Transformer: Solving Long-term Dependency Problem



★ Self-attention

| Input | Thinking | Machines |
| --- | --- | --- |
| Embedding | $X_1$ | $X_2$ |
| Queries | $q_1$ | $q_2$ |
| Keys | $k_1$ | $k_2$ |
| Values | $v_1$ | $v_2$ |

$X$   $W^Q$   $Q$

$X$   $W^K$   $K$   key

$X$   $W^V$   $V$   value

$h_1$  $h_2$  $h_3$

3.8  -0.2  5.9  Softmax → 0.2  0.1  0.7

$q_1$   key vector   $k_1$   $k_2$   $k_3$   $v_1$   $v_2$   $v_3$

$W^Q$   $W^K$   $W^V$

$x_1$   $x_2$   $x_3$

I   go   home
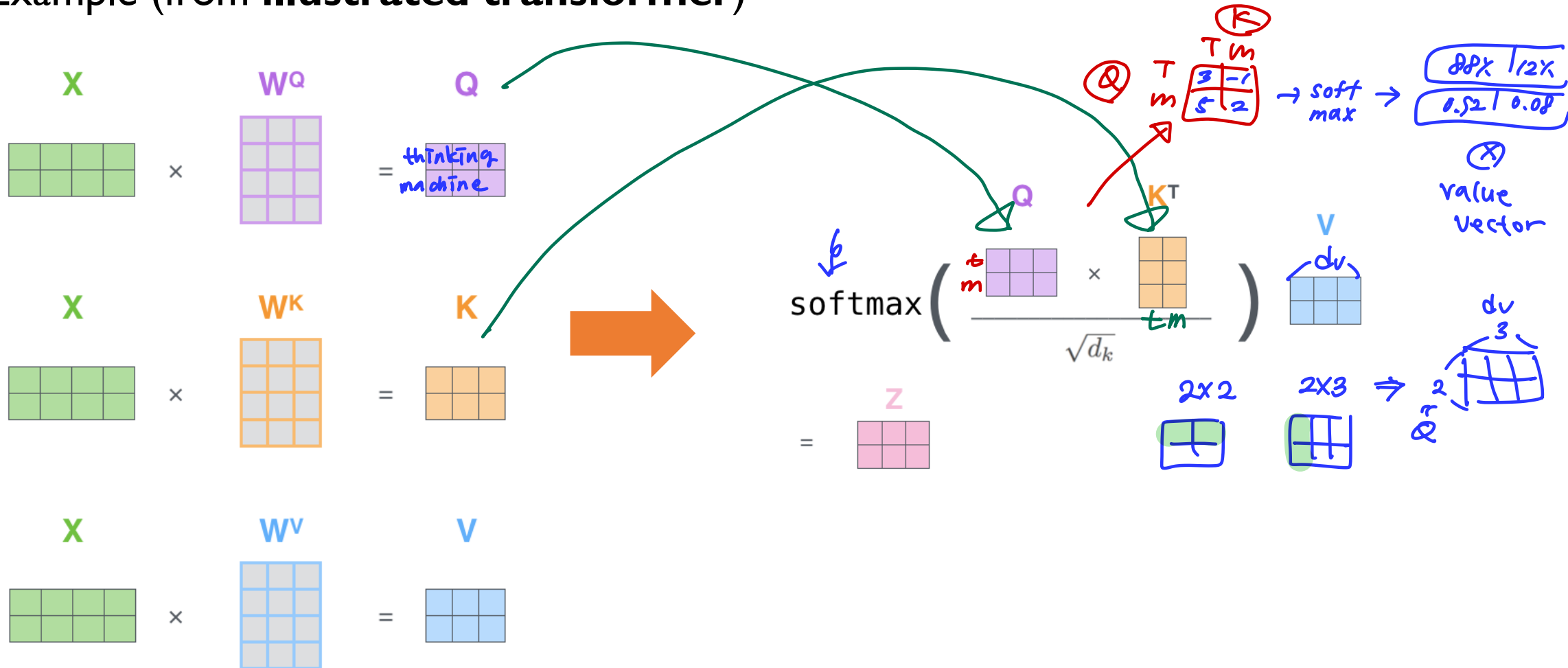
# Transformer: Scaled Dot-product Attention

- Example (from **illustrated transformer**)

# Transformer: Scaled Dot-product Attention

- **Problem**
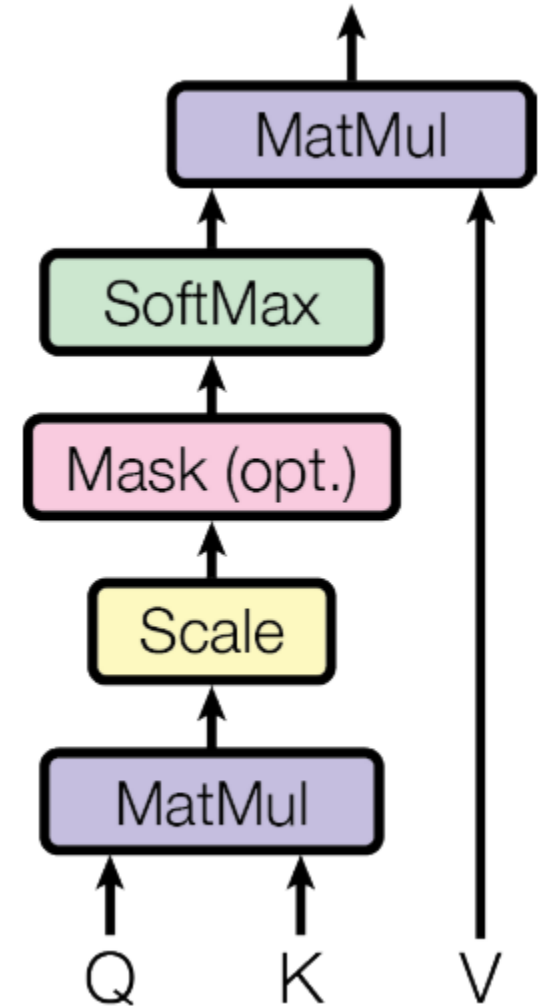
  (0 → (0, 5)
  1000 → (0, 200)

  dimension이 충분히 더 많은 숫자랑 → 연산 중 → softmax 특성상

  - As $d_k$ gets large, the variance of $q^T k$ increases. 한 값이 100% 돌아줌.

    ↓
    gradient 값 x

  - Some values inside the softmax get large.

  - The softmax gets very peaked.

  - Hence its gradient gets smaller.

- **Solution**

  이것정도는
  난순값사
  값도록.

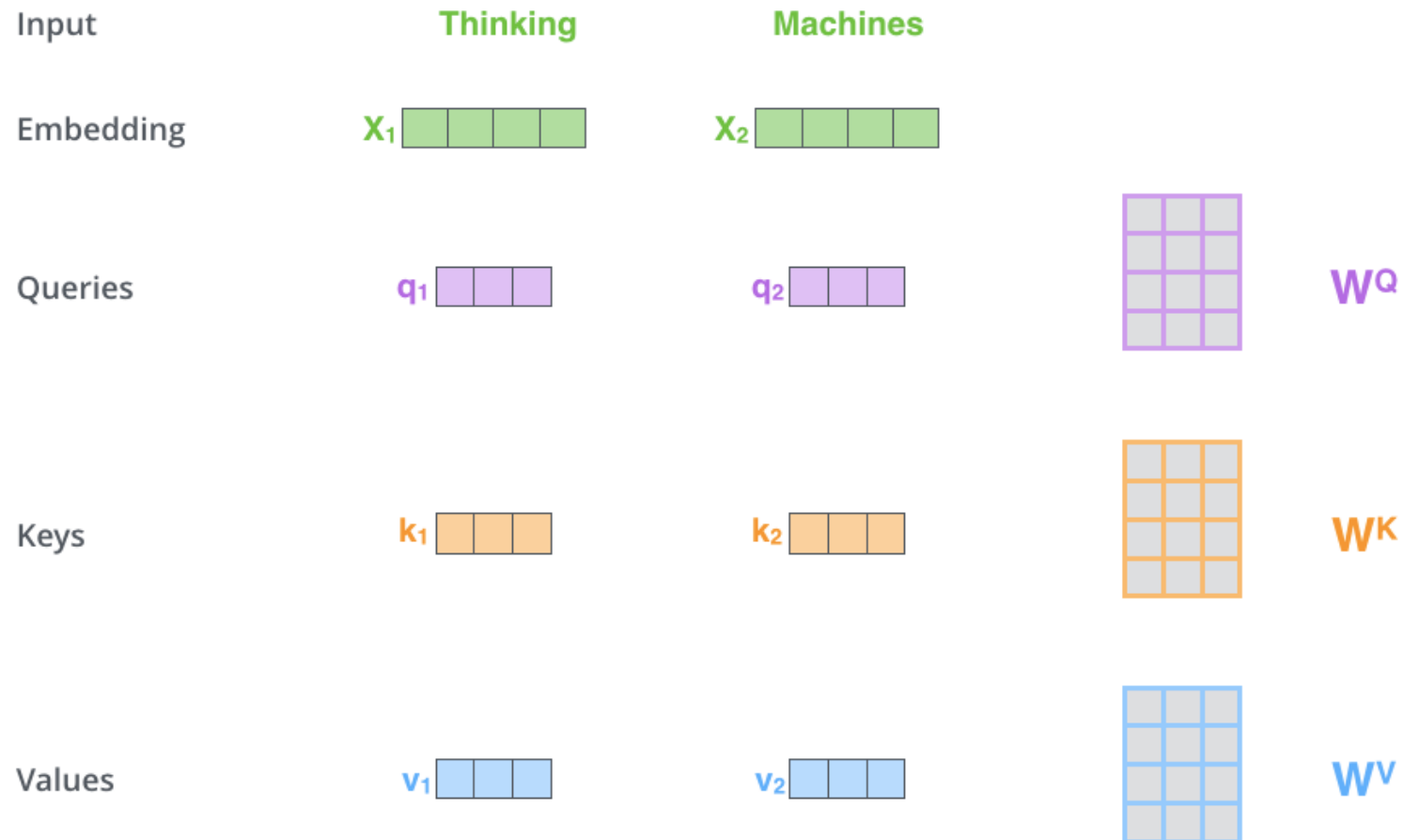  - Scaled by the length of query / key vectors:

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

MatMul

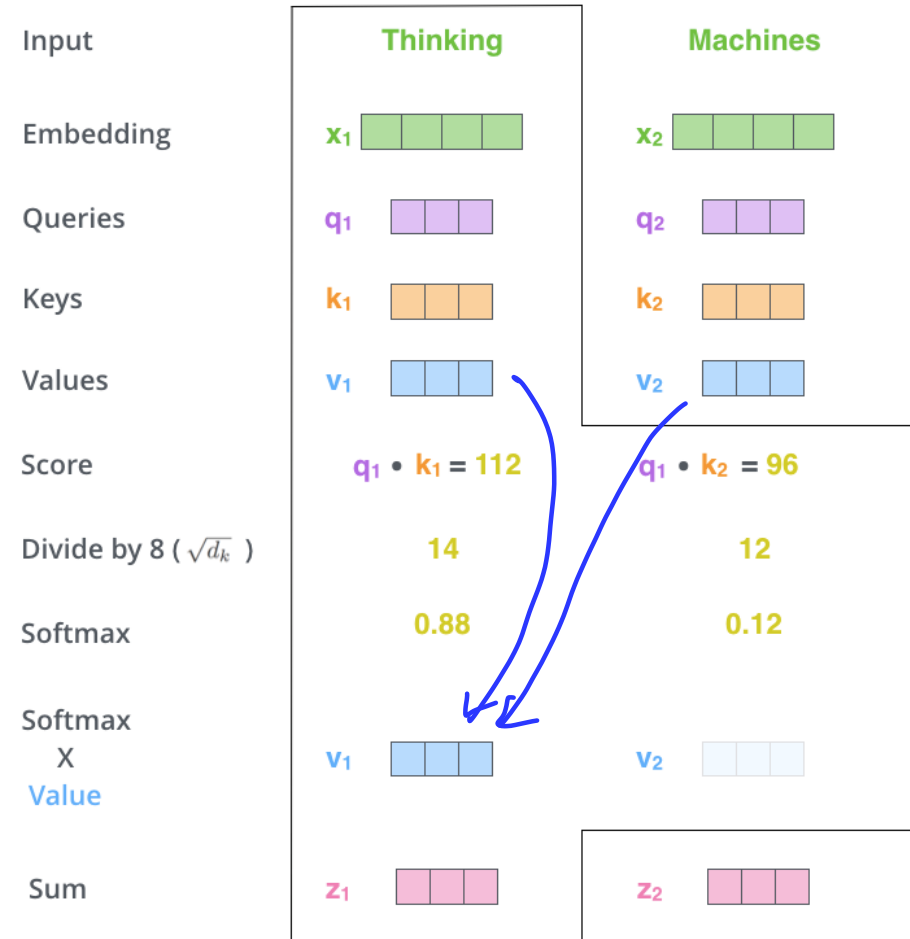SoftMax

Mask (opt.)

Scale

MatMul

Q   K   V

# Transformer: Scaled Dot-product Attention

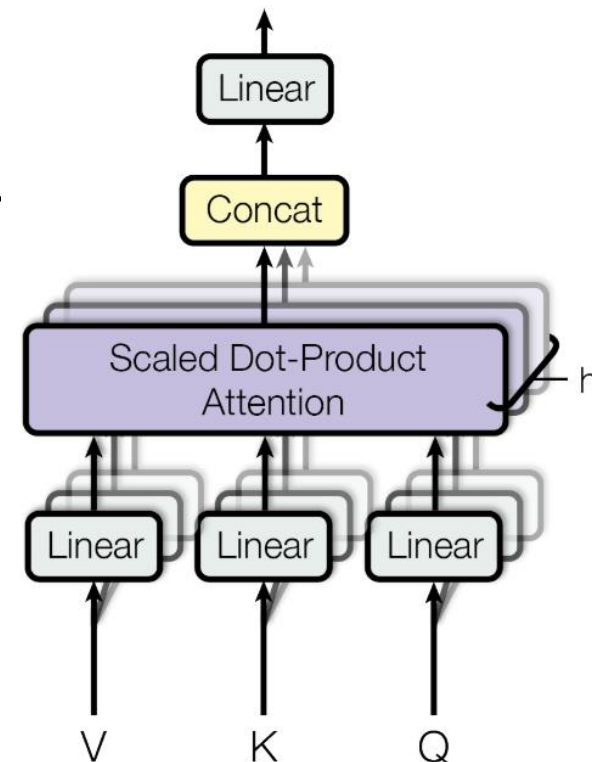- Example (from **illustrated transformer**)

# Transformer: Scaled Dot-product Attention

- Example (from **illustrated transformer**)

# Transformer: Multi-head Attention

- The input word vectors can be the queries, keys and values.

- In other words, the word vectors themselves select one another.

- **Problem:** only one way for words to interact with one another.

- **Solution:** multi-head attention maps $Q, K,$ and $V$ into the $h$ number of lower-dimensional spaces via $W$ matrices.

- Afterwards, apply attention, then concatenate outputs and pipe through linear layer.
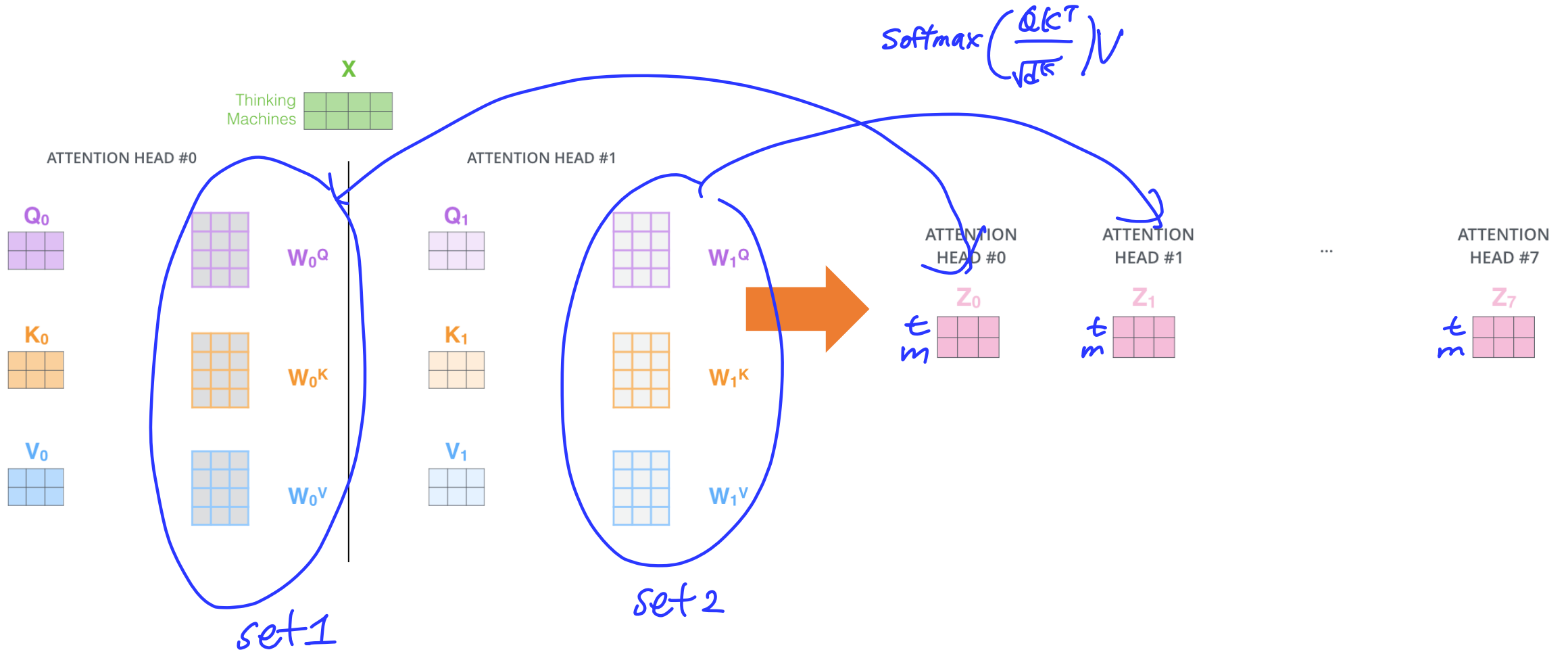


$W^Q \ W^K \ W^V$

$\hookrightarrow$ 차원의 선형변환 × 여러개

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \ldots, \text{head}_h)W^O$$
$$\text{where head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

# Transformer: Multi-head Attention

- Example (from **illustrated transformer**)



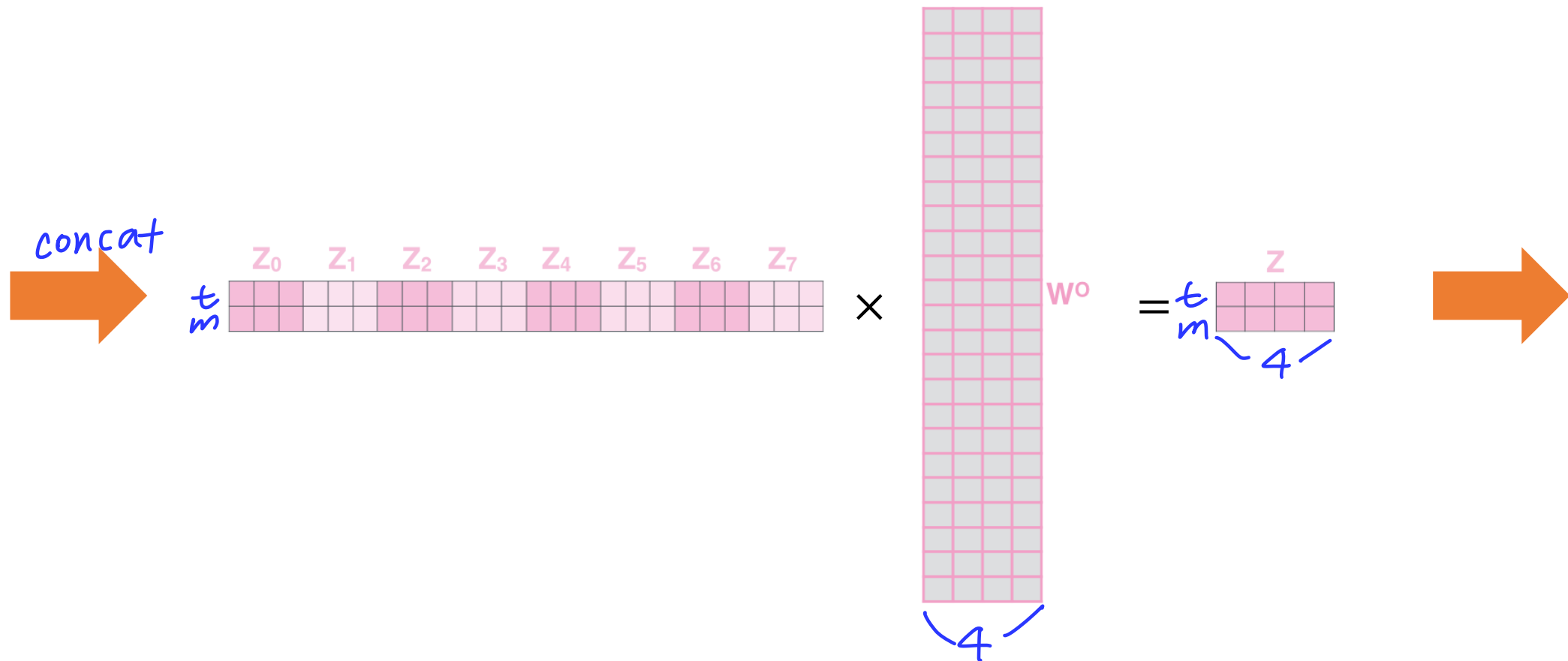$$\text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

# Transformer: Multi-head Attention

- Example (from **illustrated transformer**)

# Transformer: Multi-head Attention
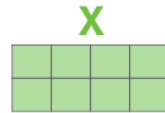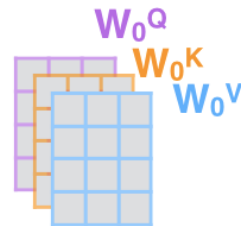
- Example (from **illustrated transformer**)



1) This is our input sentence*
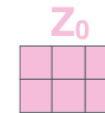
2) We embed each word*

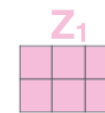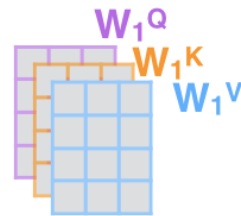3) Split into 8 heads. We multiply X or R with weight matrices

4) Calculate attention using the resulting Q/K/V matrices

5) Concatenate the resulting Z matrices, then multiply with weight matrix $W^O$ to produce the output of the layer
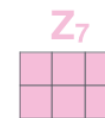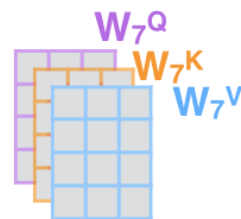
Thinking Machines

X

* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one

R

$W_0^Q$ $W_0^K$ $W_0^V$
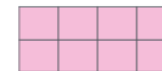
$W_1^Q$ $W_1^K$ $W_1^V$

$W_7^Q$ $W_7^K$ $W_7^V$

$Q_0$ $K_0$ $V_0$

$Q_1$ $K_1$ $V_1$

$Q_7$ $K_7$ $V_7$

$Z_0$

$Z_1$

$Z_7$

$W^O$

$Z$

Table 1: Maximum path lengths, per-layer complexity and minimum number of sequential operations for different layer types. $n$ is the sequence length, $d$ is the representation dimension, $k$ is the kernel size of convolutions and $r$ the size of the neighborhood in restricted self-attention.

| Layer Type | Complexity per Layer | Sequential Operations | Maximum Path Length |
|---|---|---|---|
| Self-Attention | $O(n^2 \cdot d)$ | $O(1)$ | $O(1)$ |
| Recurrent | $O(n \cdot d^2)$ | $O(n)$ | $O(n)$ |
| Convolutional | $O(k \cdot n \cdot d^2)$ | $O(1)$ | $O(log_k(n))$ |
| Self-Attention (restricted) | $O(r \cdot n \cdot d)$ | $O(1)$ | $O(n/r)$ |

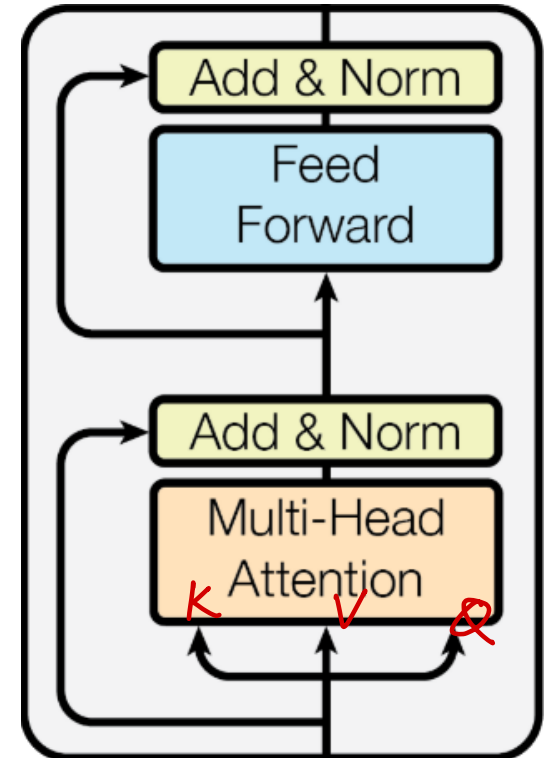*(handwritten annotation near $O(n^2 \cdot d)$: 길이의 제곱에 메모리 사이즈 커짐)*

# Transformer: Block-Based Model

Each block has two **sub-layers**

- Multi-head attention

- Two-layer feed-forward NN (with ReLU)

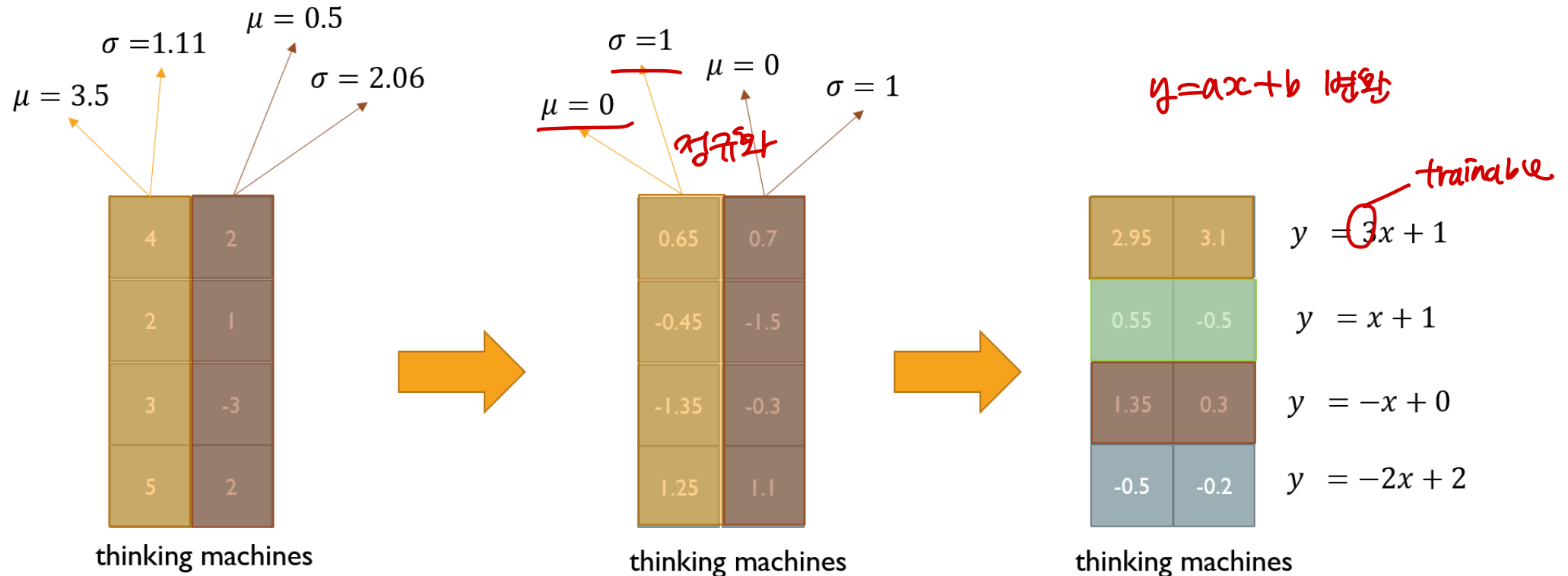Each of these two steps also has

- Residual connection and layer normalization:

    $\text{LayerNorm}(x + \text{sublayer}(x))$

- Layer normalization consists of two steps:
  - Normalization of each word vectors to have zero mean of zero and variance of one.
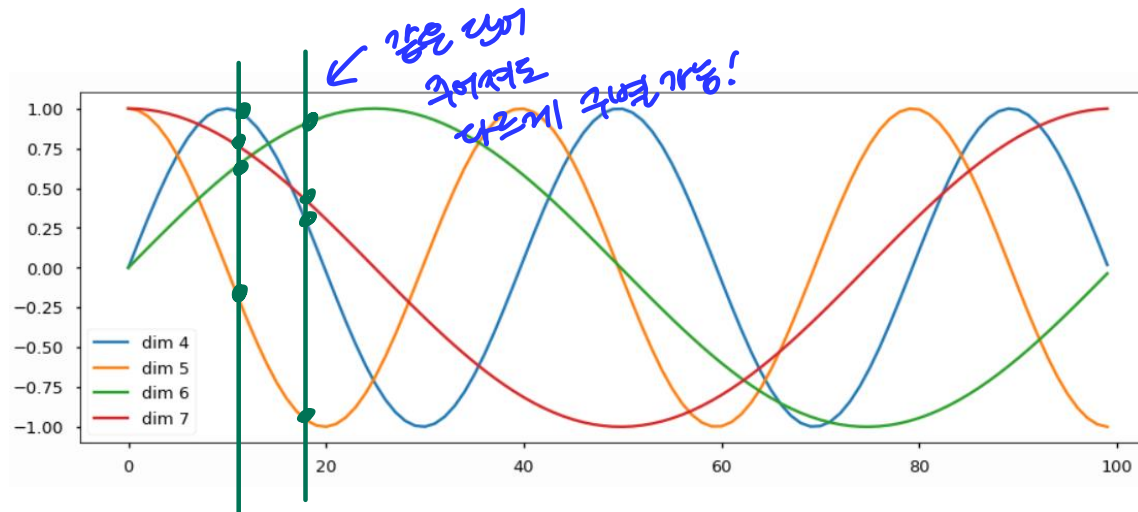  - Affine transformation of each sequence vector with learnable parameters.

# Transformer: Positional Encoding

- Use sinusoidal functions of different frequencies:

$$PE_{(\text{pos},2i)} = \sin(\text{pos}/10000^{2i/d_{\text{model}}})$$

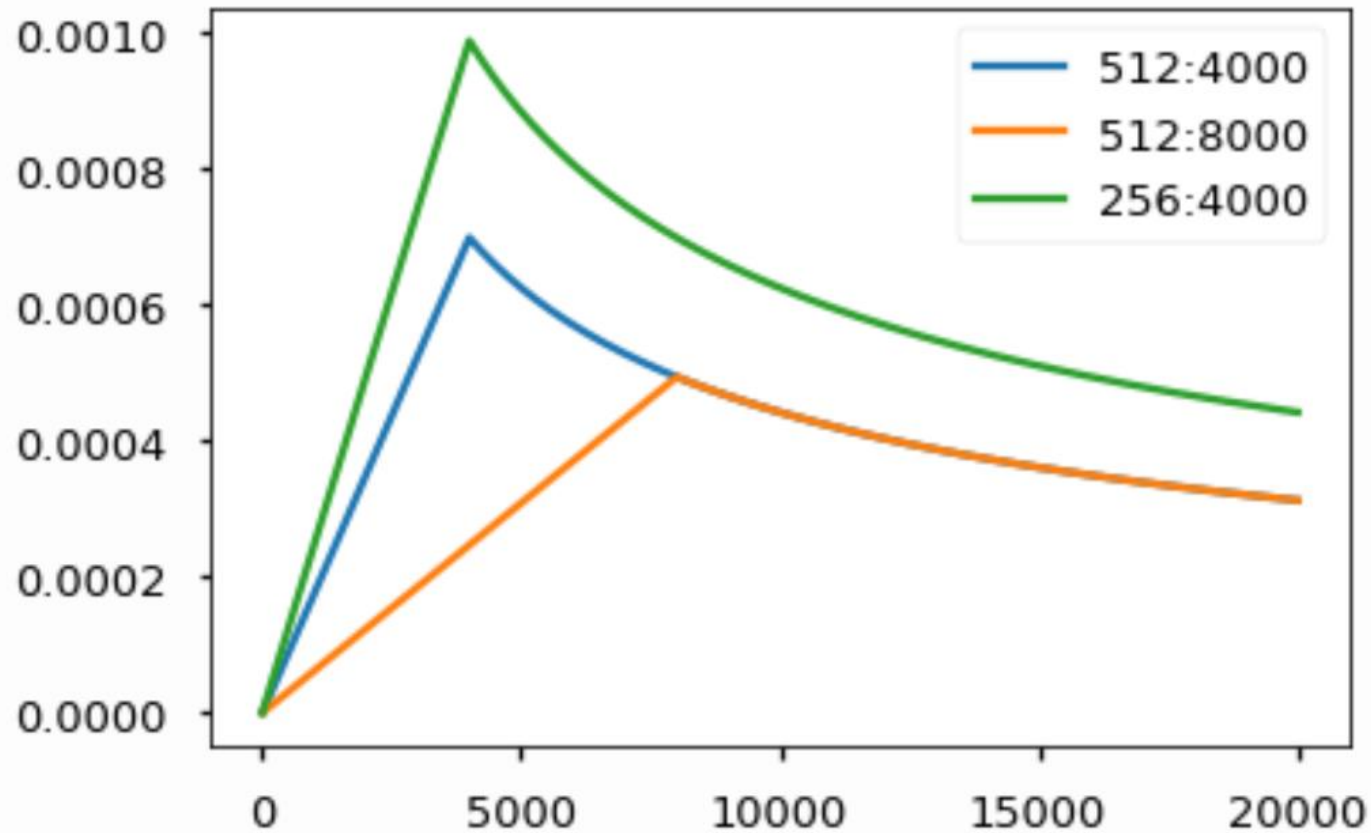$$PE_{(\text{pos},2i+1)} = \cos(\text{pos}/10000^{2i/d_{\text{model}}})$$

*home*

- Easily learn to attend by relative position, since for any fixed offset $k$, $PE_{(\text{pos}+k)}$ can be represented as linear function of $PE_{(\text{pos})}$

- Another positional encoding can also be used (e.g., positional encoding in ConvS2S).
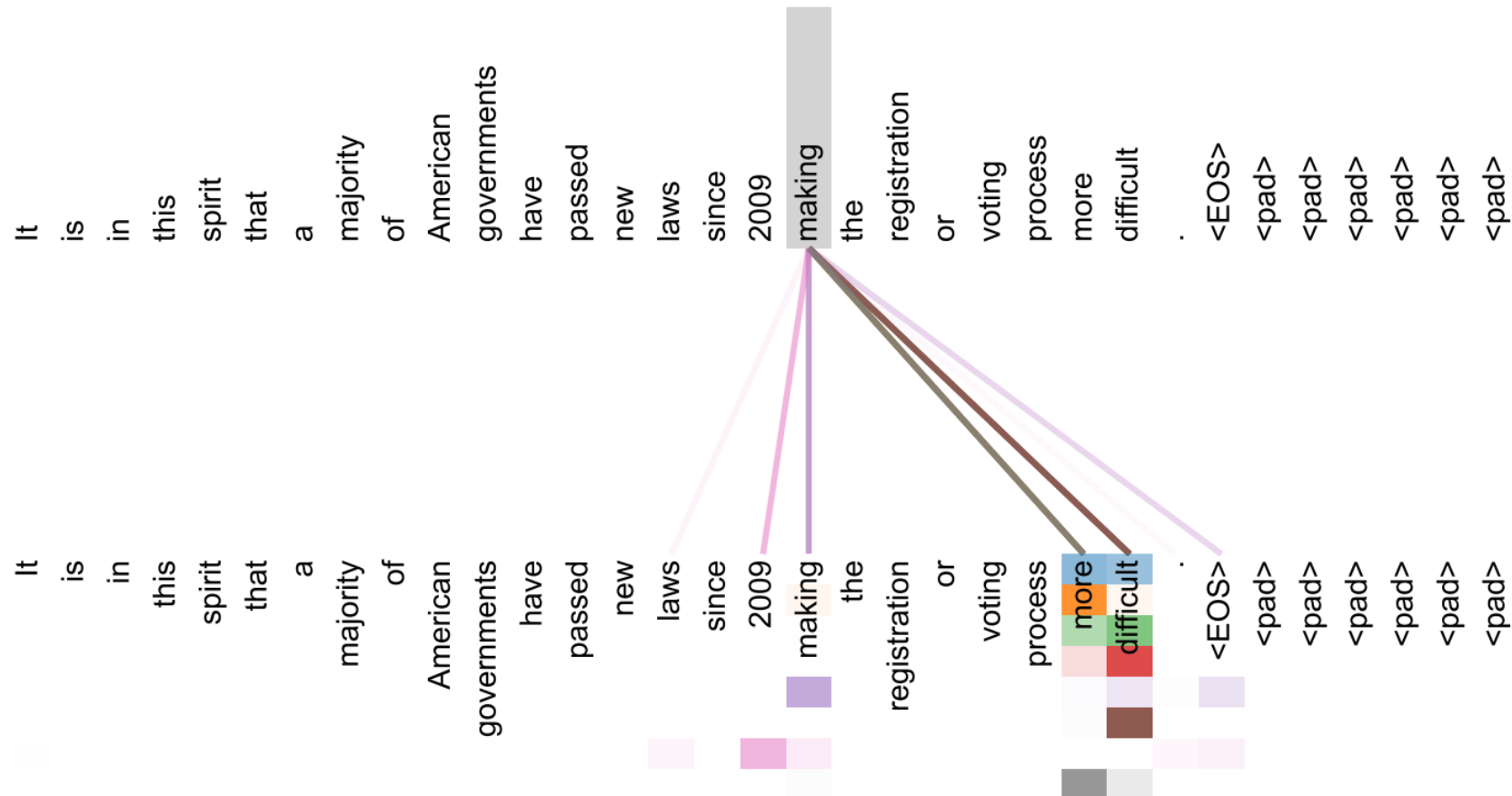
# Transformer: Warm-up Learning Rate Scheduler

- learning rate $= d_{\text{model}}^{-0.5} \cdot \min(\text{\#step}^{-0.5}, \text{\#step} \cdot \text{warmup\_step}^{-1.5})$

- Words start to pay attention to other words in sensible ways.

- Two sub-layer changes in decoder

- Masked decoder self-attention on previously generated outputs:

- Encoder-Decoder attention,
  where queries come from previous decoder layer
  and keys and values come from output of encoder

# Transformer: Masked Self-attention

- Those words not yet generated cannot be accessed during the inference time.
- Renormalization of softmax output prevents the model from accessing not yet generated words.

# Transformer: Masked Self-attention

- Those words not yet generated cannot be accessed during the inference time.
- Renormalization of softmax output prevents the model from accessing not yet generated words.

# Transformer: Masked Self-attention

- Those words not yet generated cannot be accessed during the inference time.
- Renormalization of softmax output prevents the model from accessing not yet generated words.



Query: <SOS>
Key: <SOS>

Query: <SOS>
Key: 집에

Query: 집에
Key: 나는

Query: 나는
Key: 집에

# Transformer: Masked Self-attention
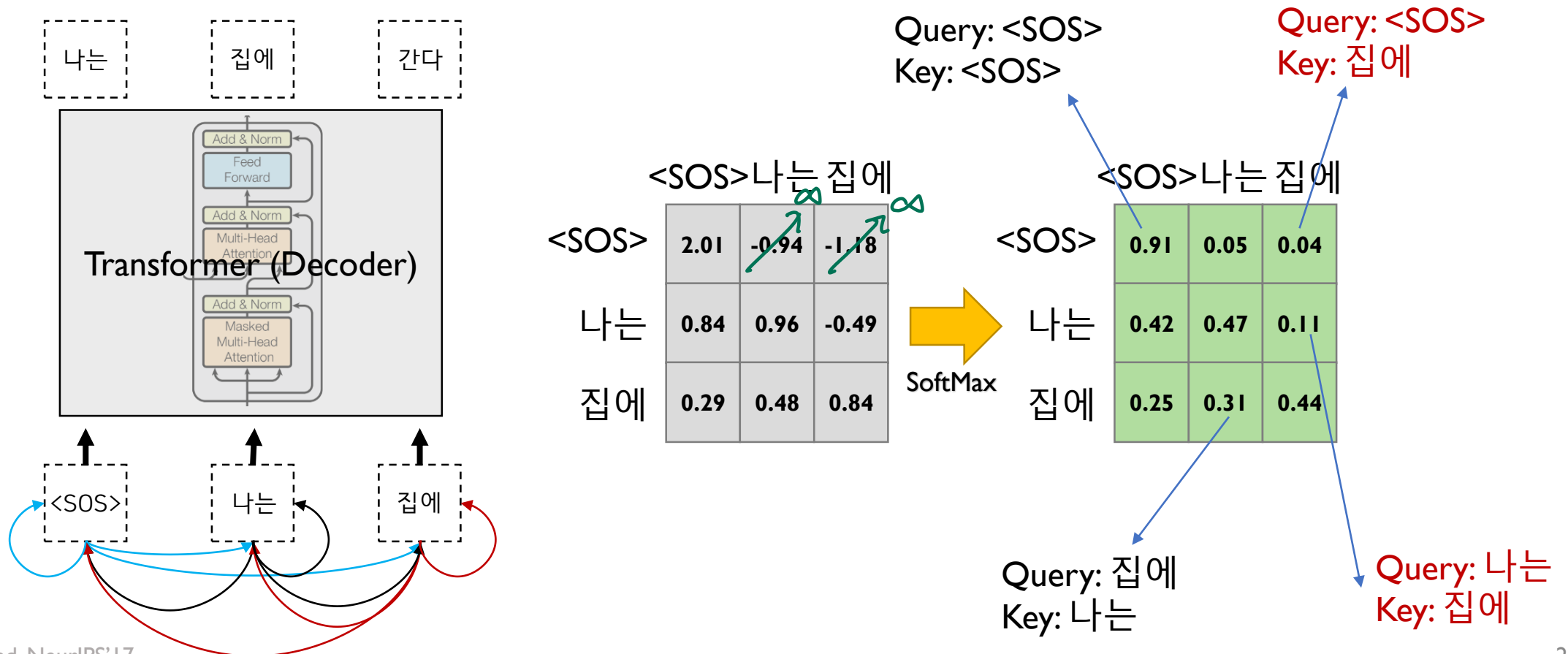
- Those words not yet generated cannot be accessed during the inference time.
- Renormalization of softmax output prevents the model from accessing not yet generated words.

$$\textbf{Softmax} \left( \frac{Query \cdot Key}{\sqrt{d_k}} \right) \cdot Value = \begin{array}{c|c|c|c} & \text{<SOS>} & \text{나는} & \text{집에} \\ \hline \text{<SOS>} & 1.00 & 0 & 0 \\ \hline \text{나는} & 0.47 & 0.53 & 0 \\ \hline \text{집에} & 0.25 & 0.31 & 0.44 \end{array} \cdot Value$$

# Transformer: Experimental Results

- Results on English-German/French translation (newstest2014)

| 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** | $2.3 \cdot 10^{19}$ | |

# Recent Trends

- Transformer model and its self-attention block has become a general-purpose sequence (or set) encoder in recent NLP applications as well as in other areas.

- Training deeply stacked Transformer models via a self-supervised learning framework has significantly advanced various NLP tasks via transfer learning, e.g., BERT, GPT-2, GPT-3, XLNet, ALBERT, RoBERTa, Reformer, T5, …

- Other applications are fast adopting the self-attention architecture and self-supervised learning settings, e.g., computer vision, recommender systems, drug discovery, and so on

- As for natural language generation, self-attention models still require a greedy decoding of words one at a time.

# References

- Harvard NLP - The Annotated Transformer

- Stanford University CS224n – Deep learning for Natural Language Processing

- Fully-parallel text generation for neural machine translation

- Convolution Sequence to Sequence

- The Illustrated Transformer (Eng)

- The Illustrated Transformer (Kor)