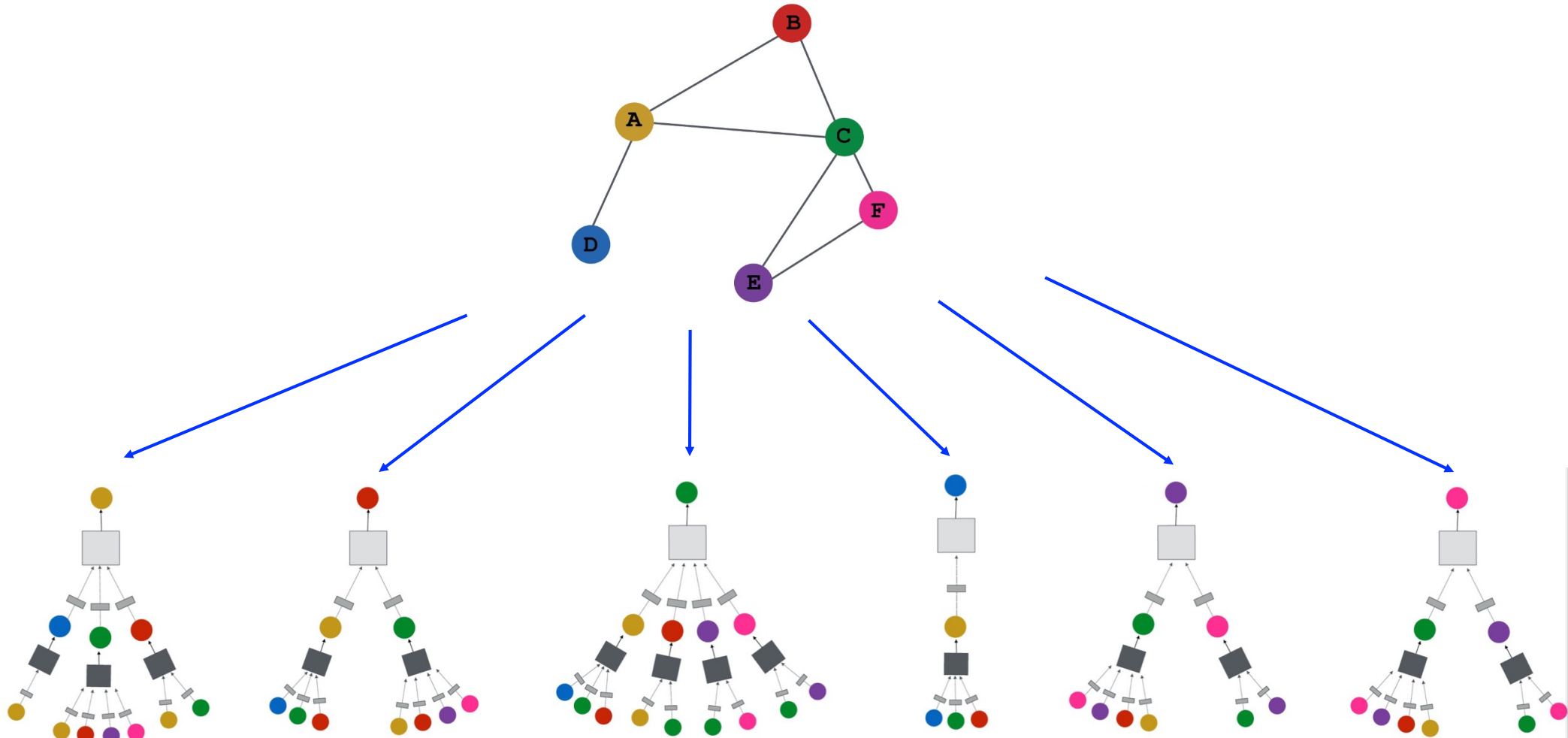


RNN vs GNN

$$\mathbf{h}_u^{(k)} = \sigma \left(\mathbf{W}_{\text{self}}^{(k)} \mathbf{h}_u^{(k-1)} + \mathbf{W}_{\text{neigh}}^{(k)} \sum_{v \in \mathcal{N}(u)} \mathbf{h}_v^{(k-1)} + \mathbf{b}^{(k)} \right)$$

$$\mathbf{h}^{(t)} = \sigma \left(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_e \mathbf{e}^{(t)} + \mathbf{b}_1 \right)$$

Graph convolutional network

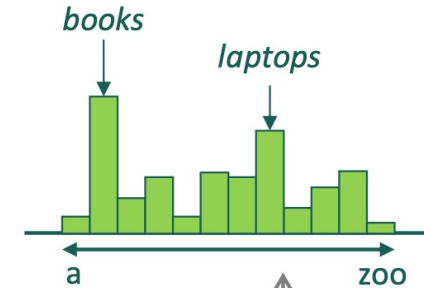


RNN language model

output distribution

$$\hat{y}^{(t)} = \text{softmax}(U h^{(t)} + b_2) \in \mathbb{R}^{|V|}$$

$$\hat{y}^{(4)} = P(x^{(5)} | \text{the students opened their})$$



hidden states

$$h^{(t)} = \sigma(W_h h^{(t-1)} + W_e e^{(t)} + b_1)$$

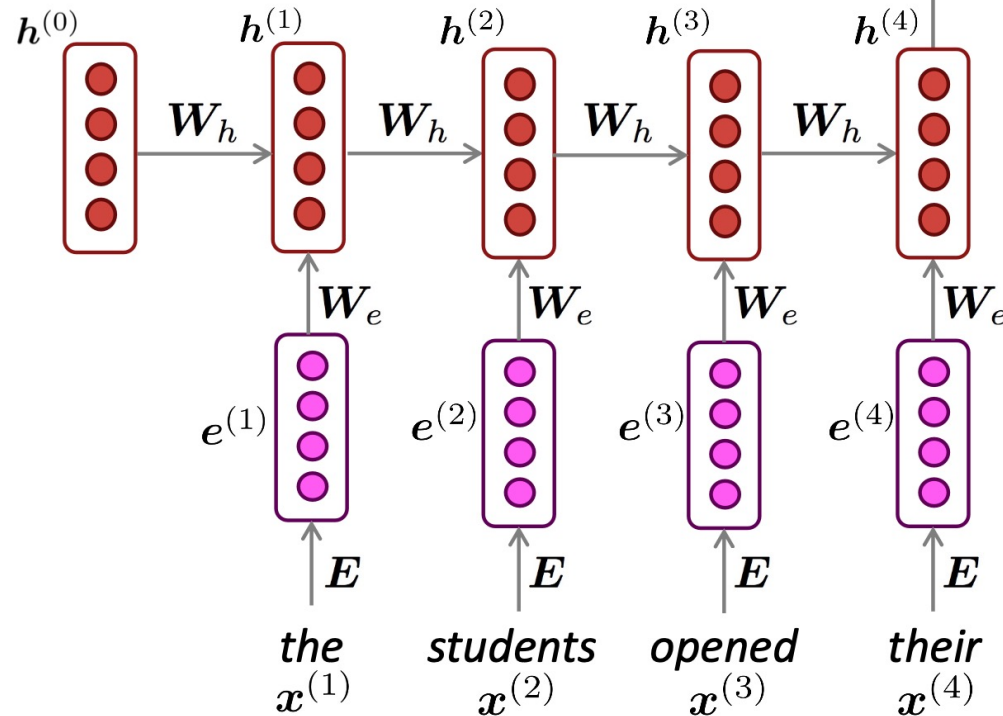
$h^{(0)}$ is the initial hidden state

word embeddings

$$e^{(t)} = E x^{(t)}$$

words / one-hot vectors

$$x^{(t)} \in \mathbb{R}^{|V|}$$



RNN with attention

$$\mathbf{a}_t = \sum_{i=1}^N \alpha_i^t \mathbf{h}_i \in \mathbb{R}^h$$

Weighted sum of encoder hidden states based on the attention distribution

Softmax

$$\alpha^t = \text{softmax}(\mathbf{e}^t) \in \mathbb{R}^N$$

Dot product

$$\mathbf{e}^t = [s_t^T \mathbf{h}_1, \dots, s_t^T \mathbf{h}_N] \in \mathbb{R}^N$$

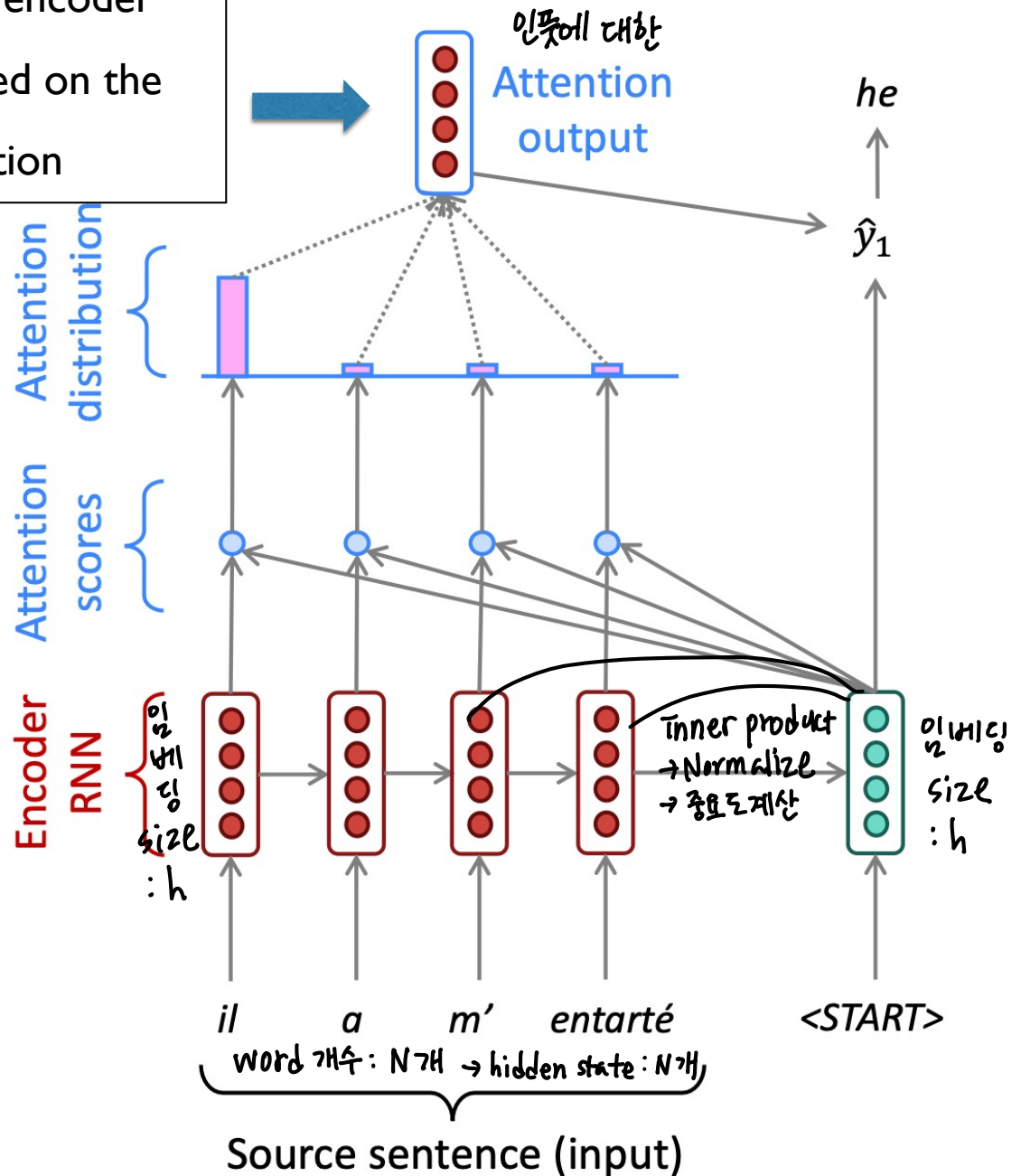
→ Vector size : 1×1

Encoder hidden states

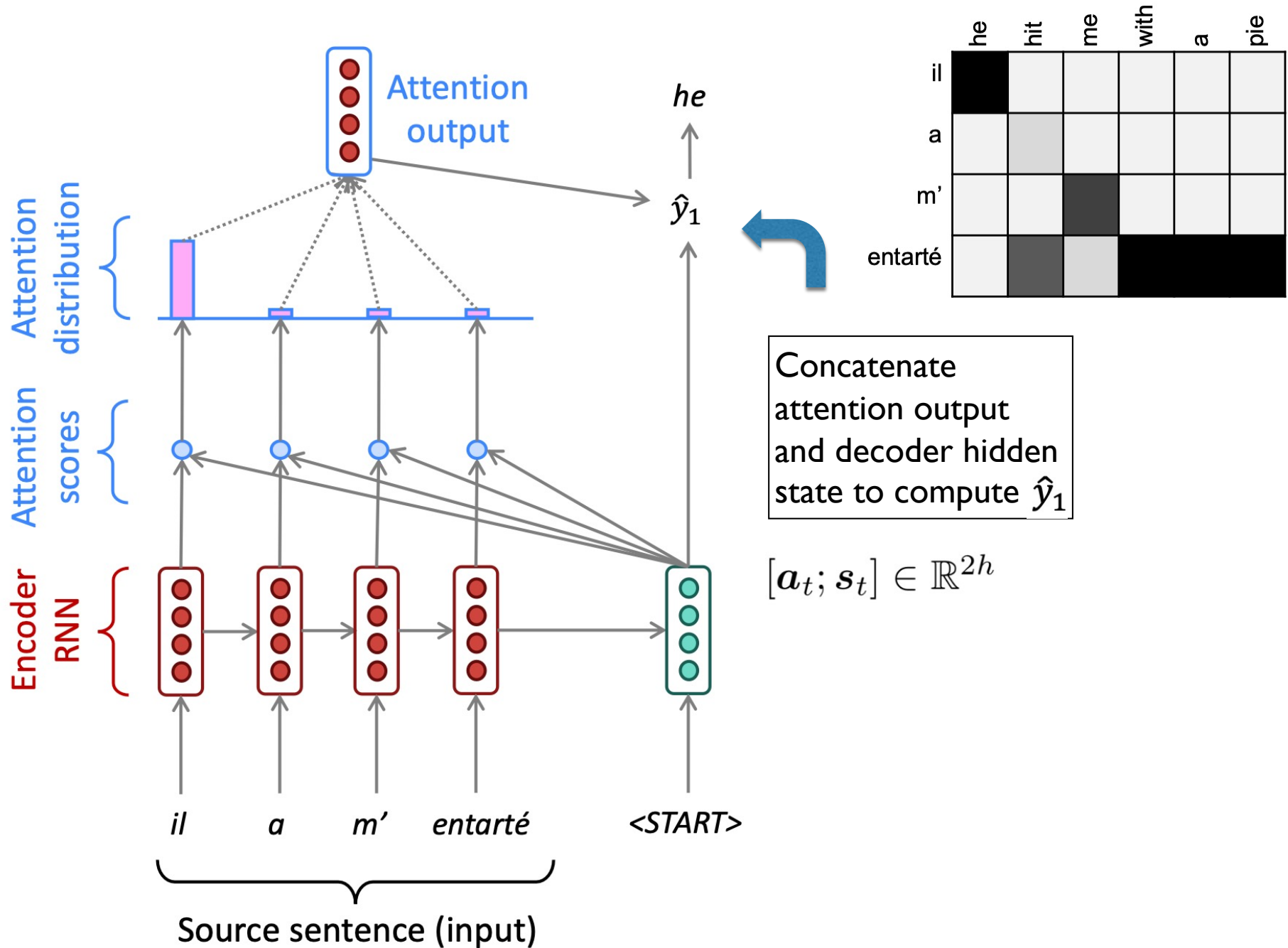
$$\mathbf{h}_1, \dots, \mathbf{h}_N \in \mathbb{R}^h$$

Decoder hidden states

$$\mathbf{s}_t \in \mathbb{R}^h$$



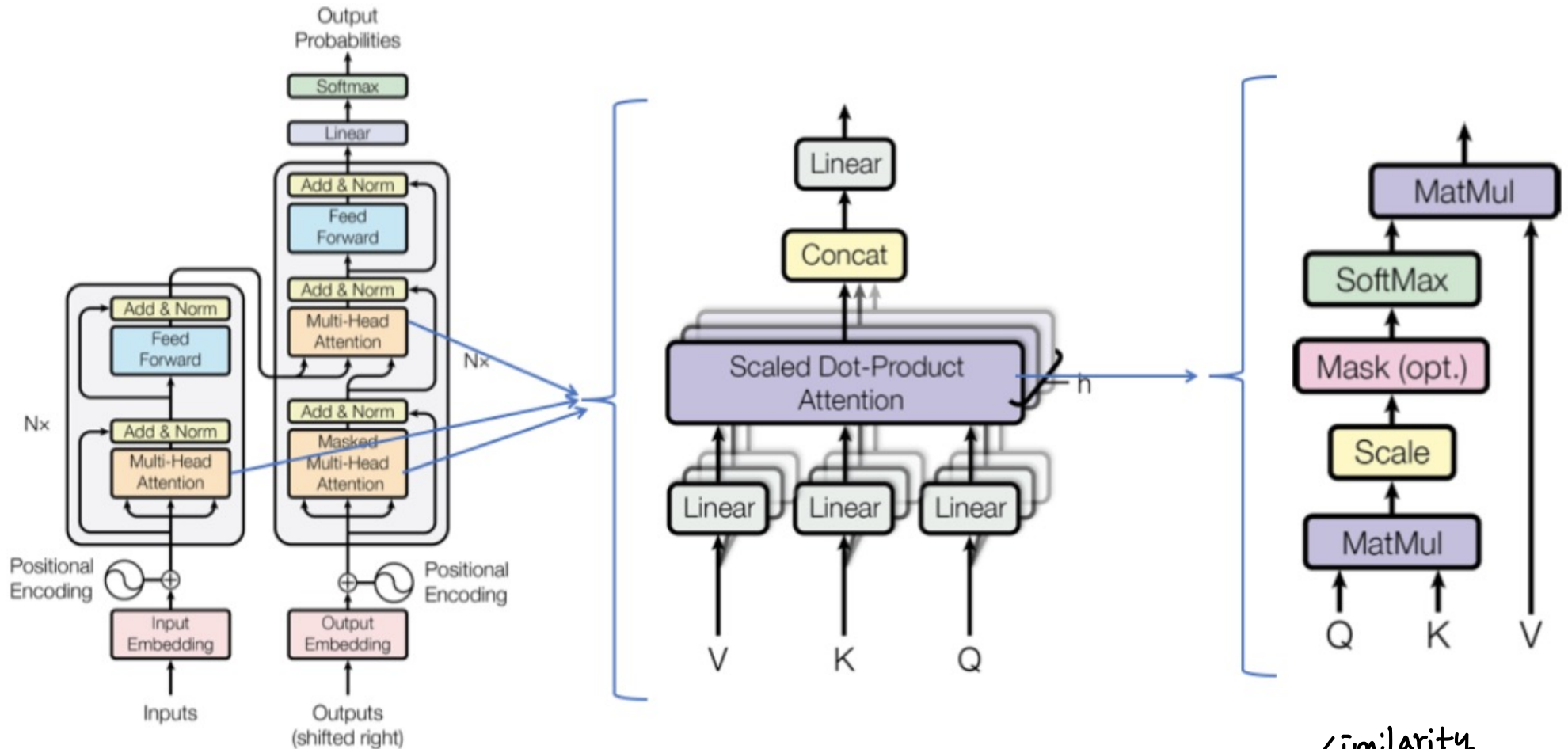
RNN with attention



RNN with attention

- **Attention** allows the model to focus on the relevant parts of the input sequence as needed
- **Self-attention** is an attention mechanism relating different positions of a single sequence in order to compute a representation of the sequence

Self-attention

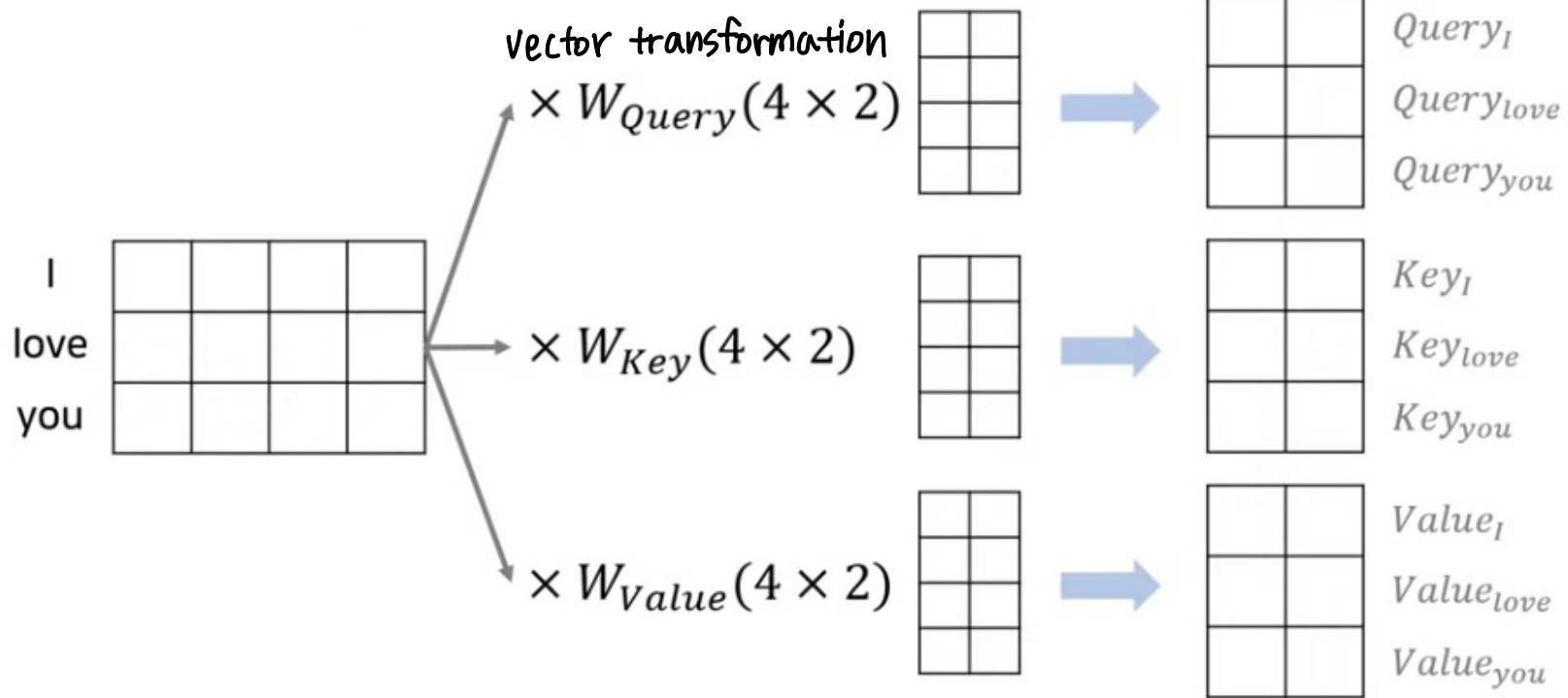
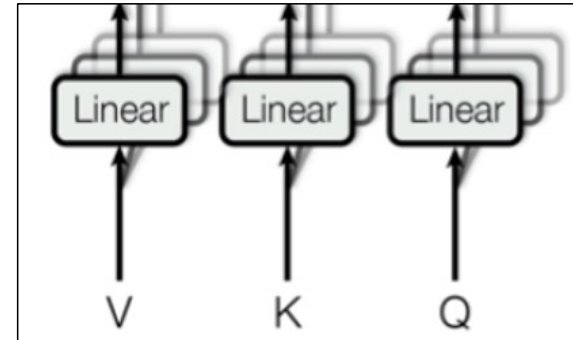


$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Handwritten annotations: "similarity" with an arrow pointing to QK^T , and "weight 계산" (weight calculation) with an arrow pointing to the softmax operation.

Self-attention

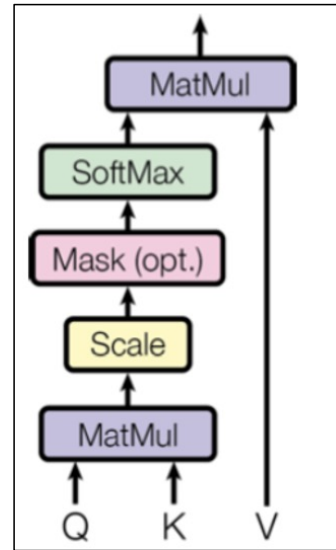
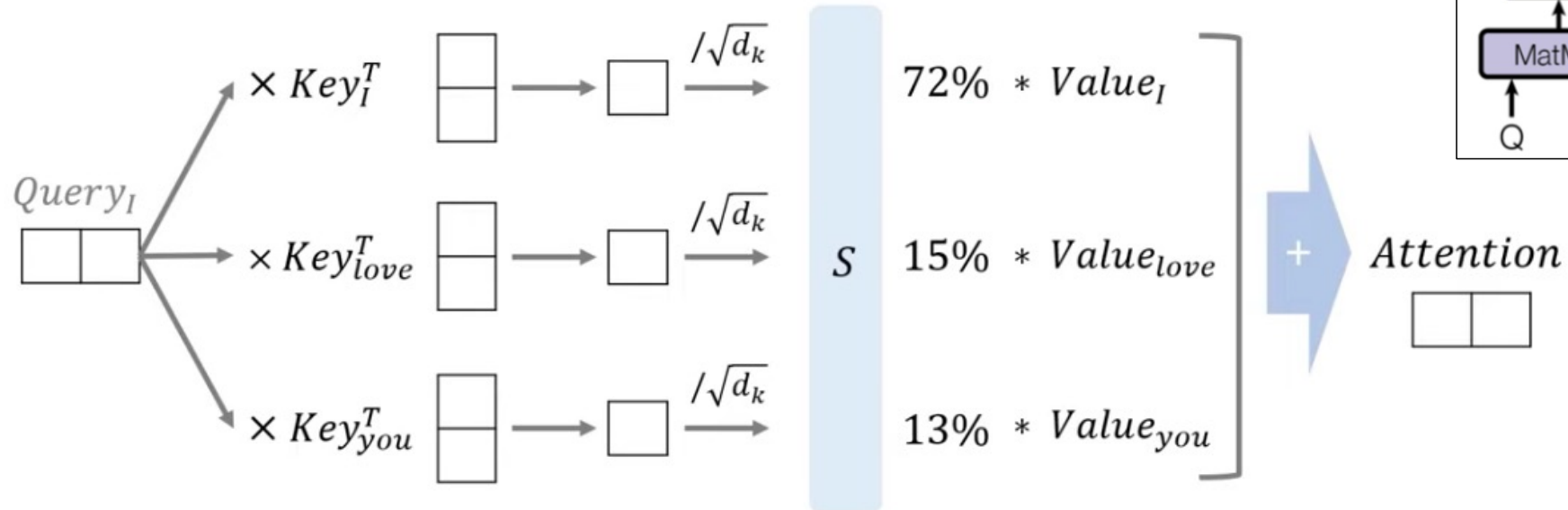
- Transform the embedding to the embedding for query, key, and value



Self-attention

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

element by element similarity
 Normalization



$$e_{ij} = q_i^T k_j$$

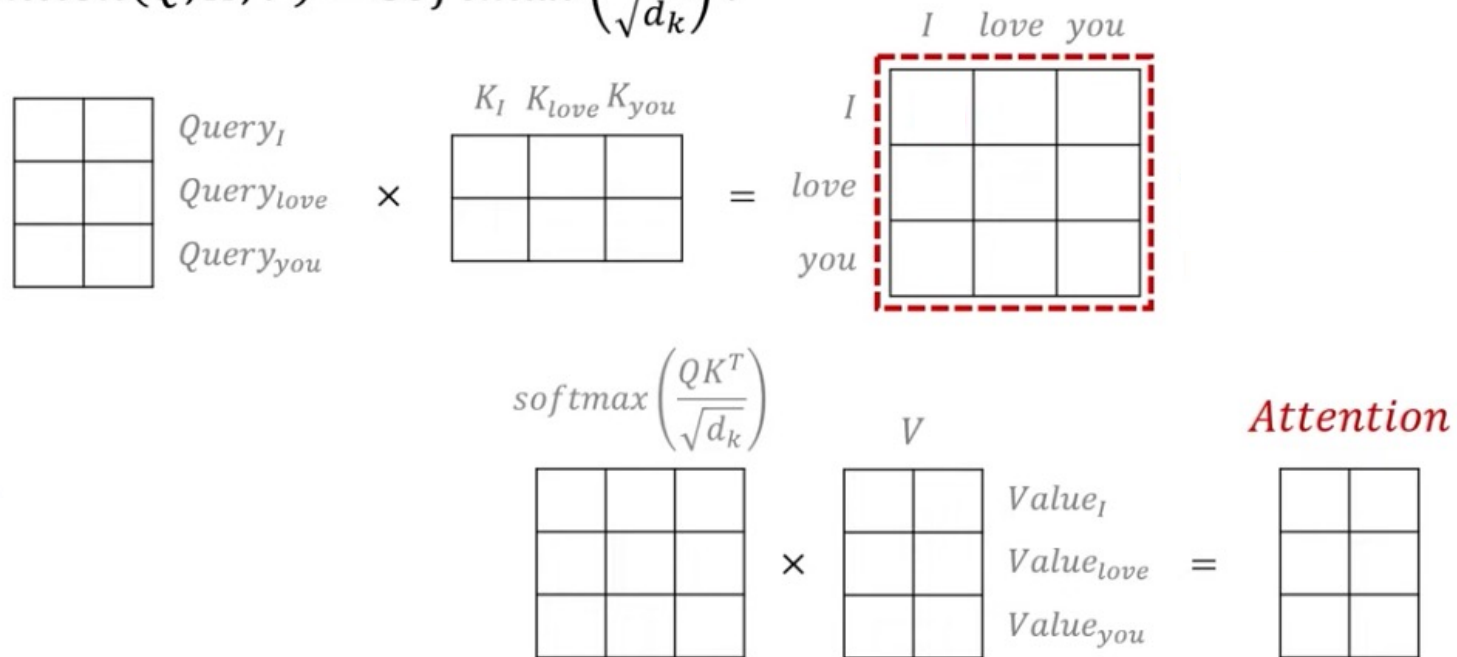
$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{j'} \exp(e_{ij'})}$$

$$\text{output}_i = \sum_j \alpha_{ij} v_j$$

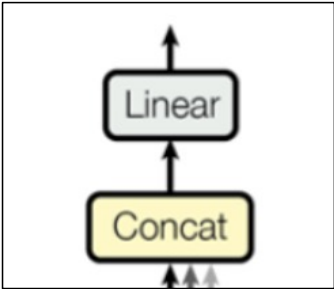
Normalize \rightarrow coefficient

Self-attention

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



Self-attention



$Concat(head_1, \dots, head_h) =$

$head_1$

$head_2$

$head_3$

\dots

$head_h$

$d_{model} = d_v \times h$

$MultiHead(Q, K, V) =$

$d_{model} = d_v \times h$

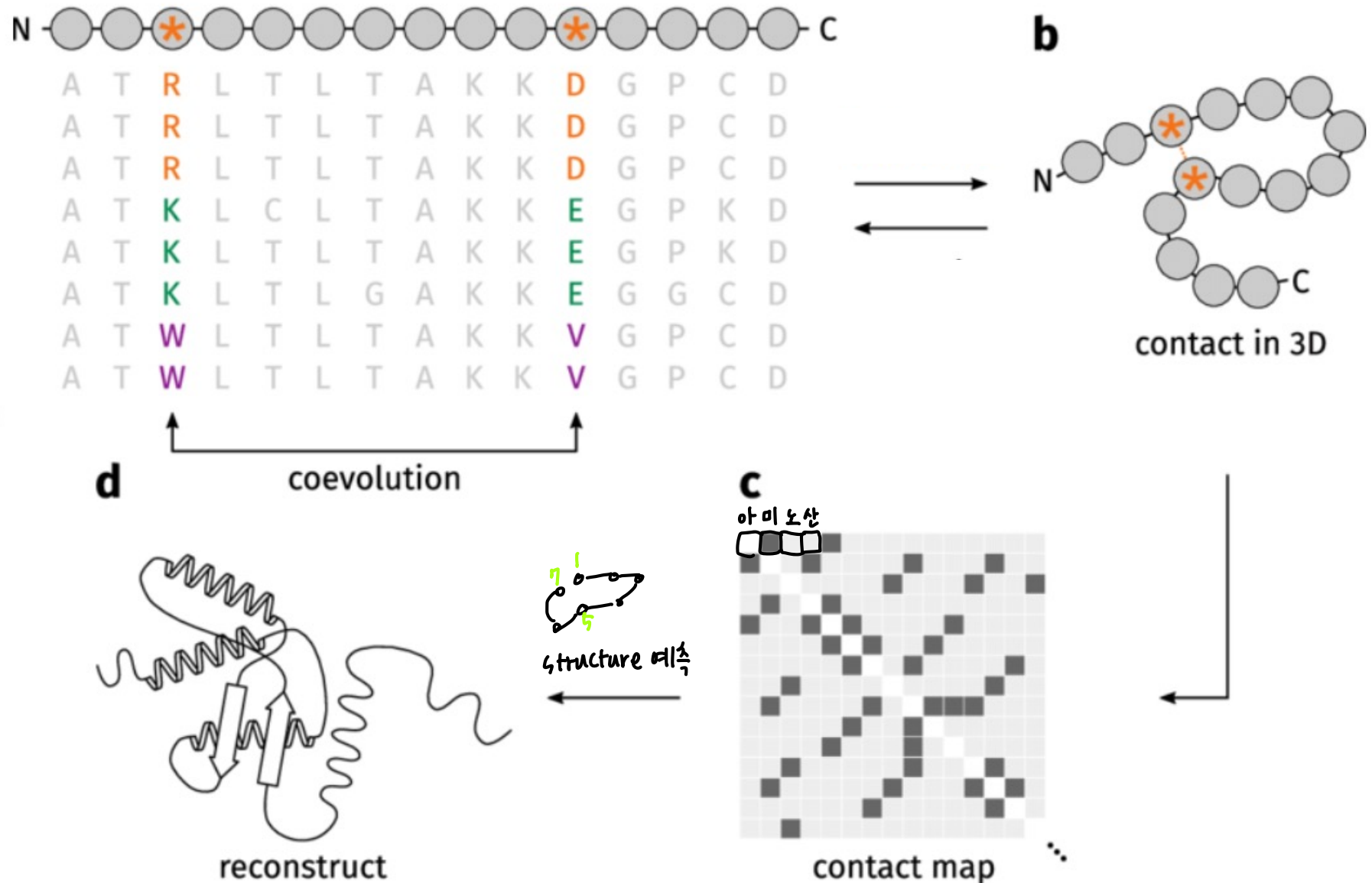
\times

d_{model}

$seq_len \times$

d_{model}

Self-attention applied to genomic sequences



Masked language model

(Bert)

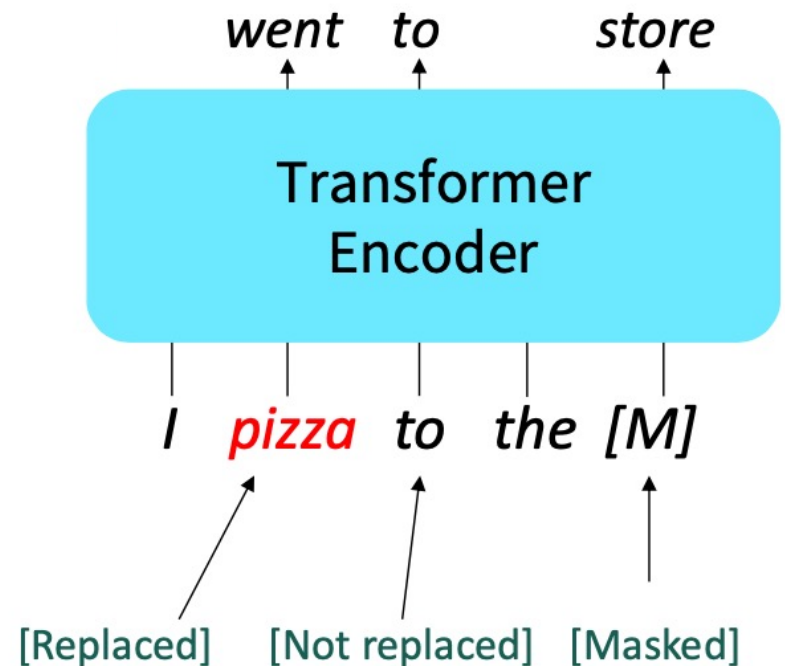
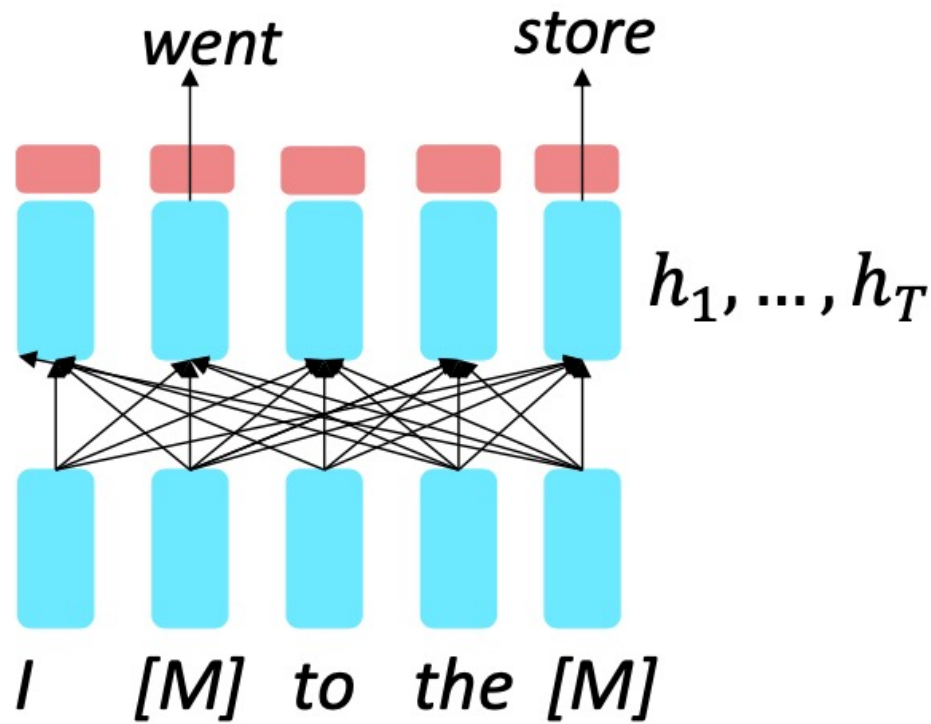
Self-supervised learning

→ 가장 큰 목적: 레이블링된 데이터가 적어도

Input으로 좋은 representation 만들기

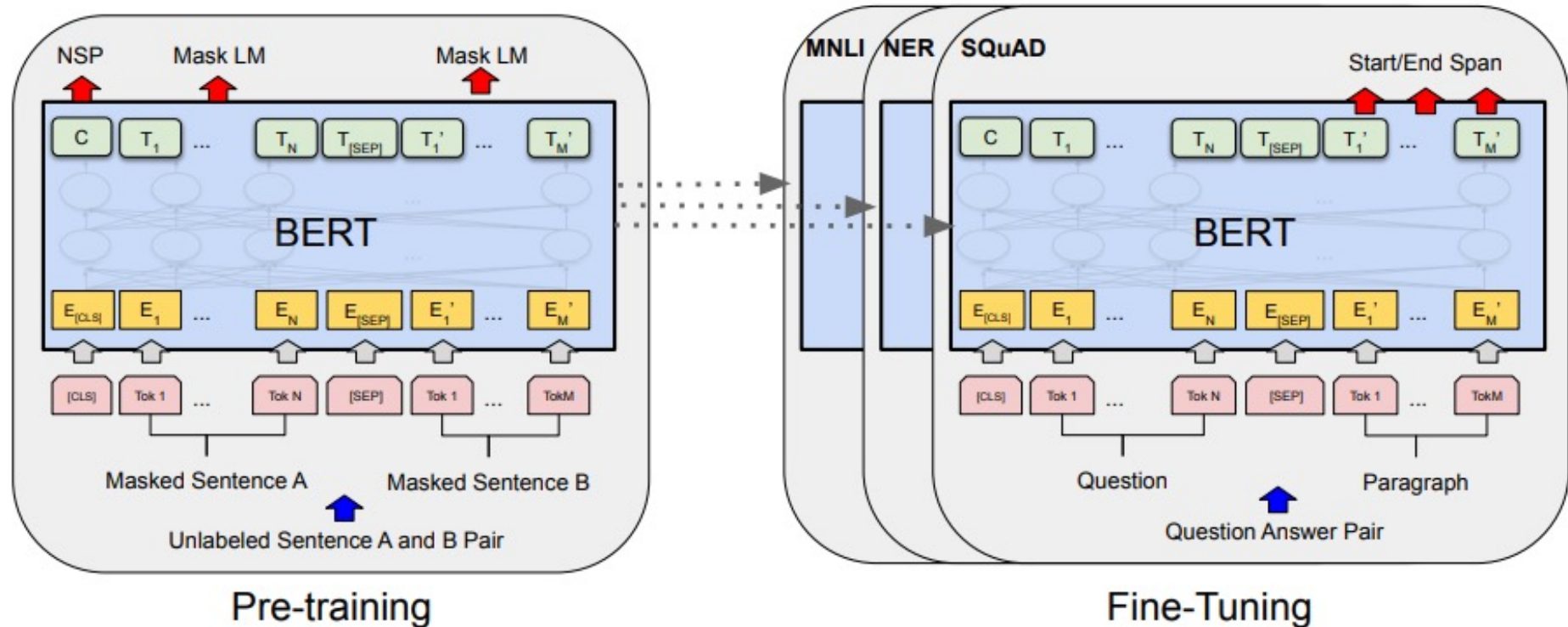
\tilde{x} is the masked version of x , we're learning $p(x|\tilde{x})$

- replace some fraction of words in the input with a special [MASK] token
- predict these words.



Bidirectional Encoder Representation

Input: word piece



- Apart from output layers, the same architectures are used in both pre-training and fine-tuning
- The same pre-trained model parameters are used to initialize models for different down-stream tasks
- During fine-tuning, all parameters are fine-tuned.

BERT applied to genomic sequences

