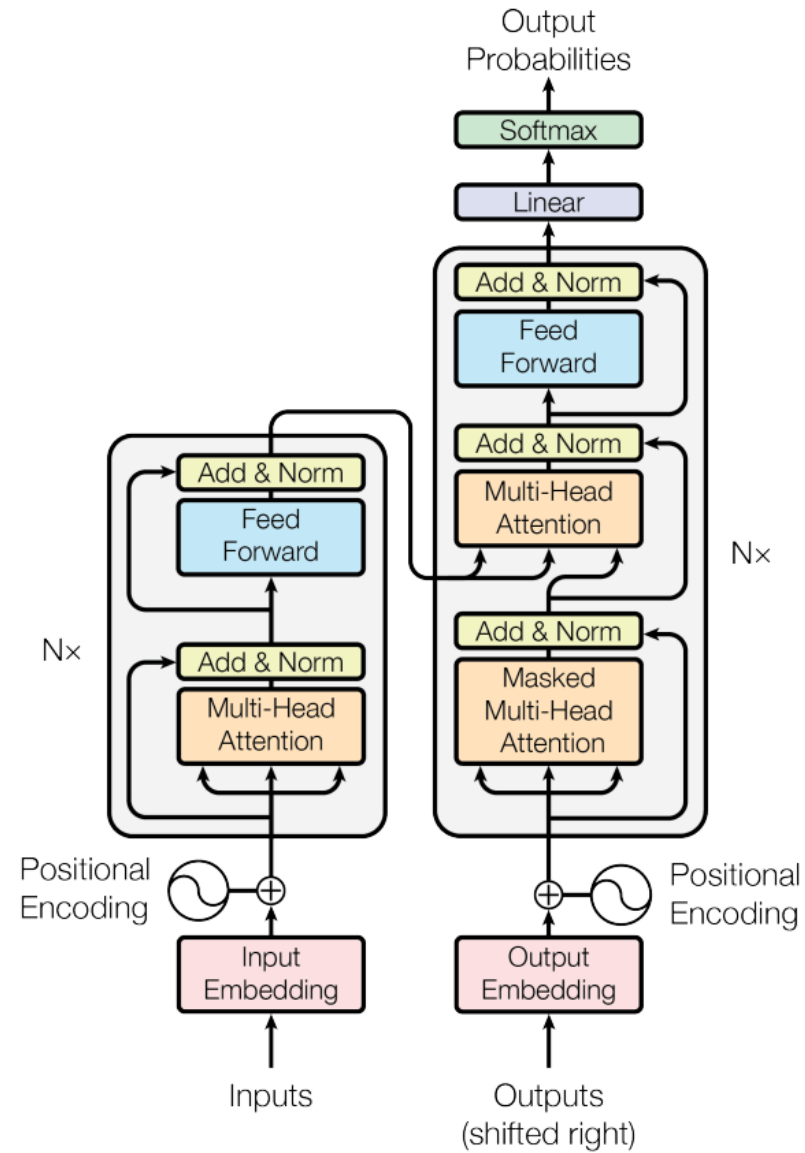


# 인공지능 심화교육 -자연어처리 실습 2

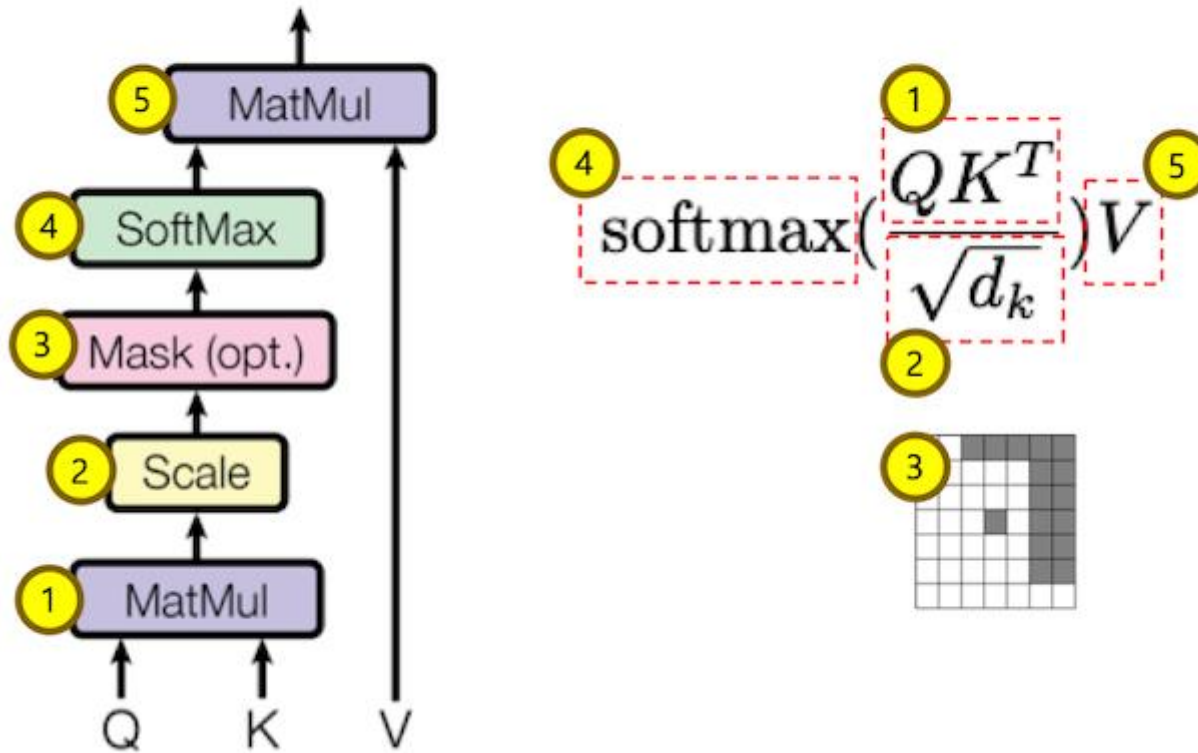
전북대학교 이성민

# Transformer



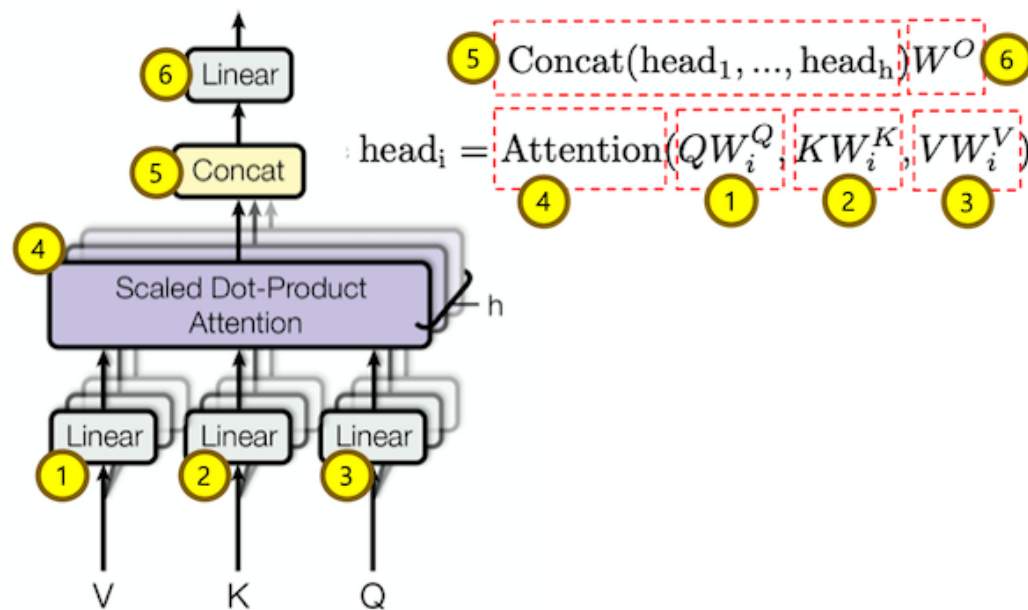
# Transformer

## Scaled Dot Product Attention



# Transformer

Multi-Head Attention

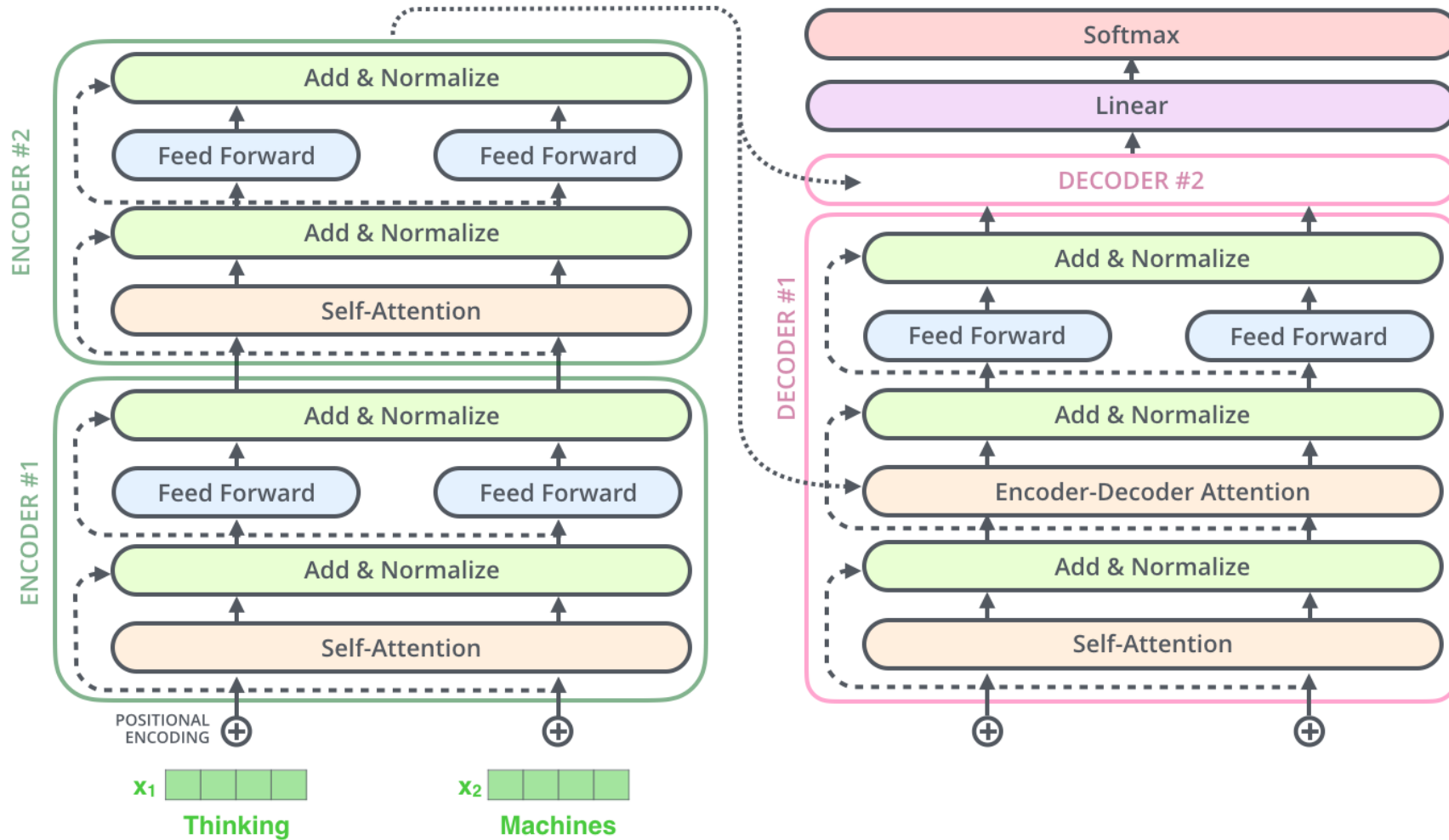


$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W^O$$

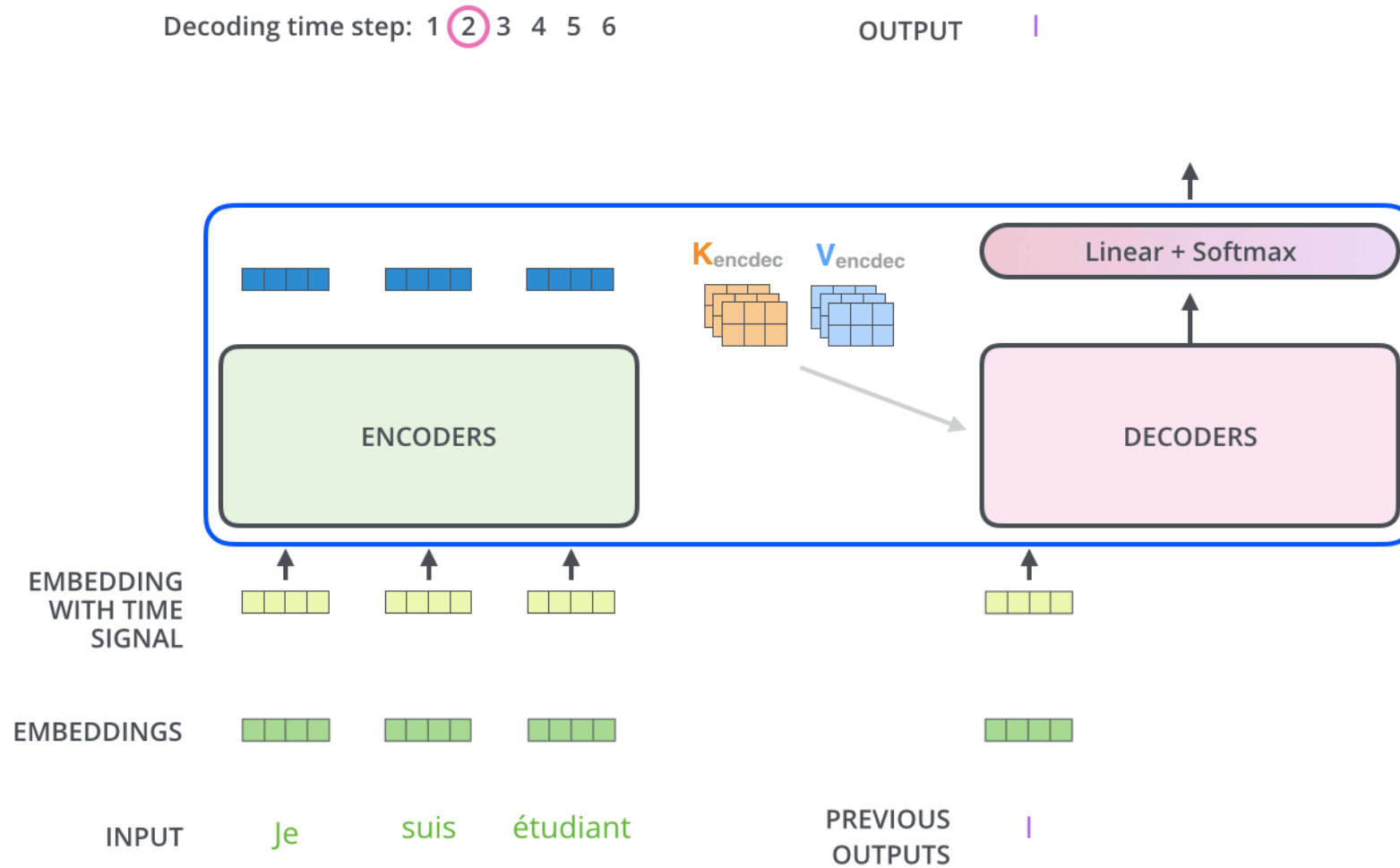
$$\text{where } \text{head}_i = \text{Attention}(Q_i^Q, K_i^K, V_i^V)$$

Where the projections are parameter matrices  $W_i^Q \in \mathbb{R}^{d_{\text{model}} \times d_k}$ ,  $W_i^K \in \mathbb{R}^{d_{\text{model}} \times d_k}$ ,  $W_i^V \in \mathbb{R}^{d_{\text{model}} \times d_v}$  and  $W^O \in \mathbb{R}^{hd_v \times d_{\text{model}}}$ .

# Transformer



# Transformer



# Transformer

## Attention Pad Mask

	I	am	a	boy	[pad]	[pad]
I						
am						
a						
boy						
[pad]						
[pad]						

# Transformer

## Attention Decoder Mask

	I	am	a	boy	[pad]	[pad]
I						
am						
a						
boy						
[pad]						
[pad]						



# Text Classification + Contrastive Loss

$$\mathcal{L}_{\text{s\_cl}} = -\frac{1}{T} \sum_{i=1}^N \sum_{j=1}^N \mathbb{1}_{y_i=y_j} \log \frac{\exp(\text{sim}(\mathbf{h}_i, \mathbf{h}_j)/\tau)}{\sum_{n=1}^N \exp(\text{sim}(\mathbf{h}_i, \mathbf{h}_n)/\tau)}$$

$$\mathcal{L}_{\text{intent}} = -\frac{1}{N} \sum_{j=1}^C \sum_{i=1}^N \log P(C_j | u_i)$$

$$\mathcal{L}_{\text{stage2}} = \mathcal{L}_{\text{s\_cl}} + \lambda' \mathcal{L}_{\text{intent}}$$