

# Vision Transformer

생성형 AI Study  
Luddite



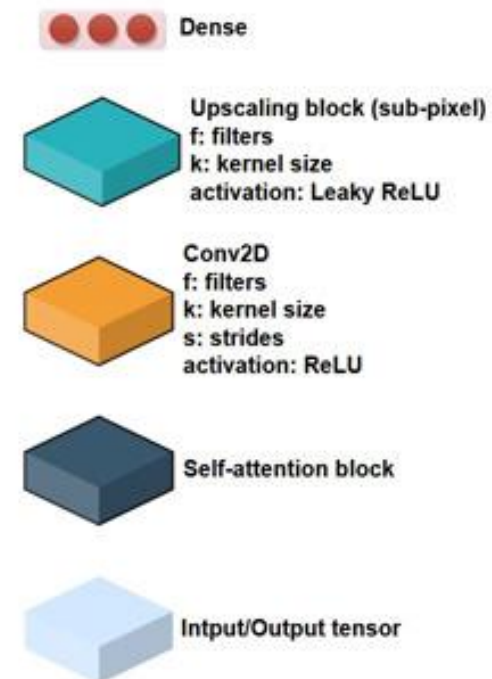
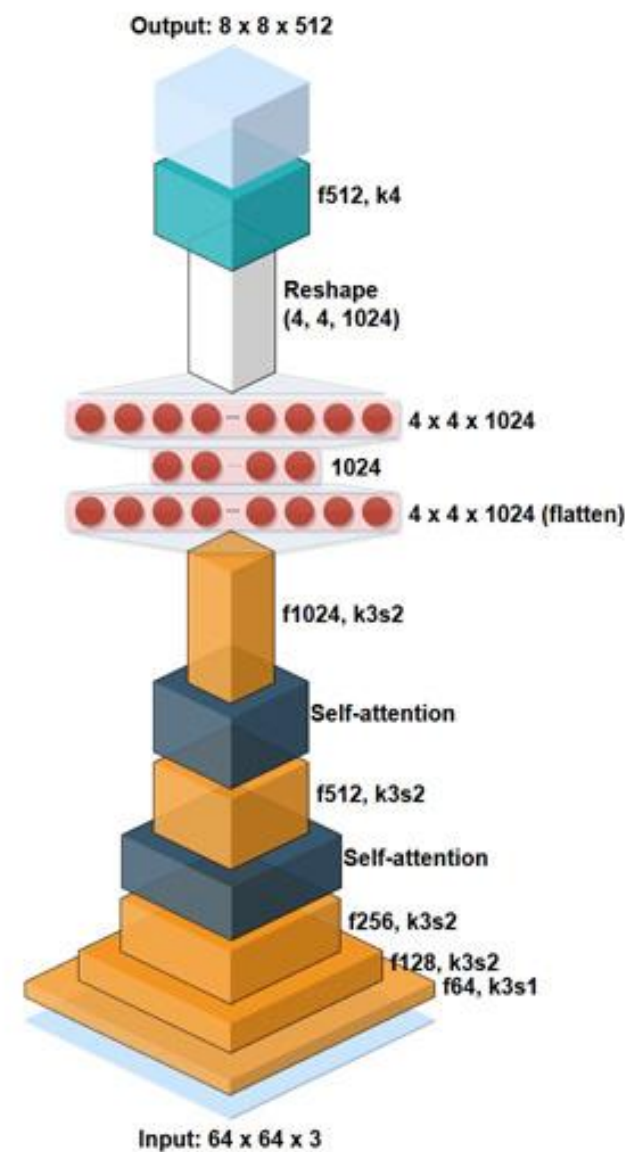
---

발표자 : 김채아

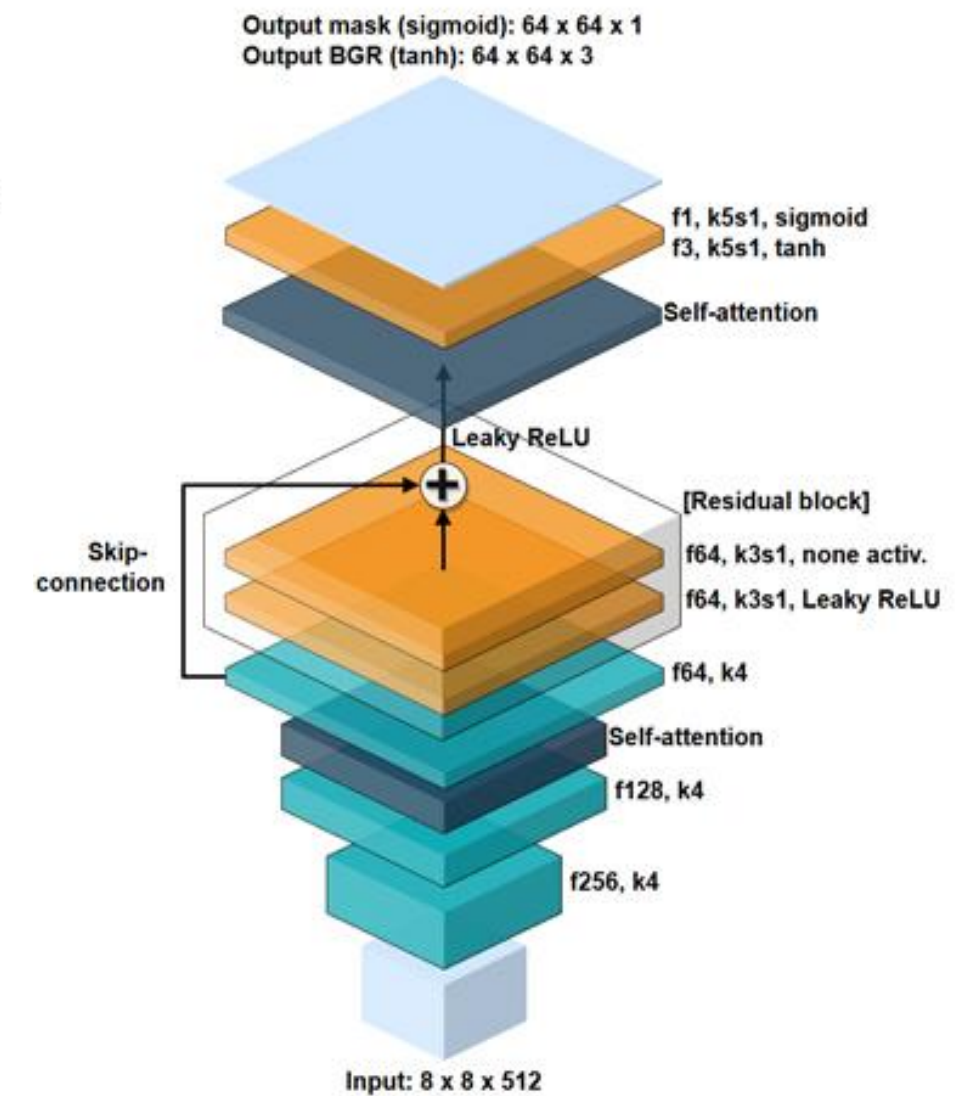
# Deepfake?

- Encoder-Decoder 구조
- Sequence data
- Self-attention, skip-connection

## Encoder



## Decoder



Faceswap-GAN

# Transformer

Sequential data를  
처리하고 인코딩 하는 모델

- Translation
- Image classification (ViT)
- Image detection (DETR)
- DALL-E, Deepfake 등 transformer 구조인  
self-attention 활용

# Transformer 구조

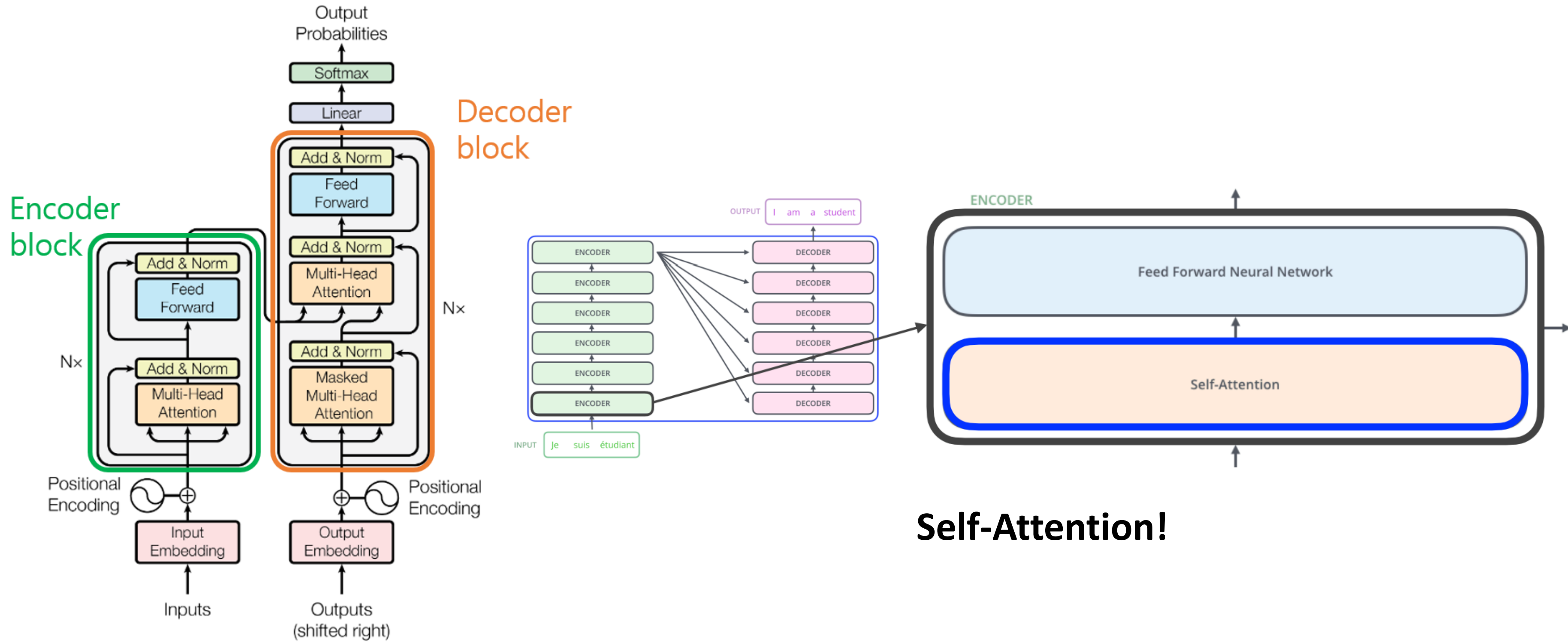
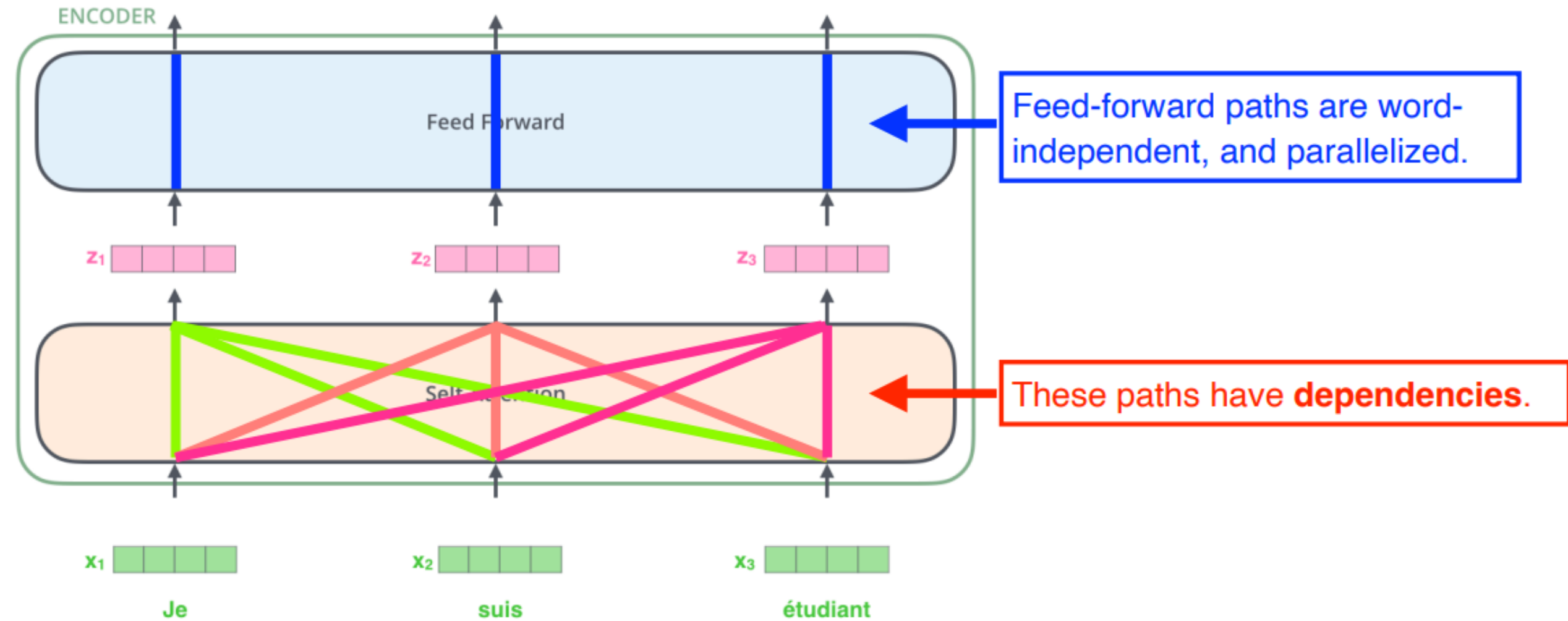


Figure 1: The Transformer - model architecture.

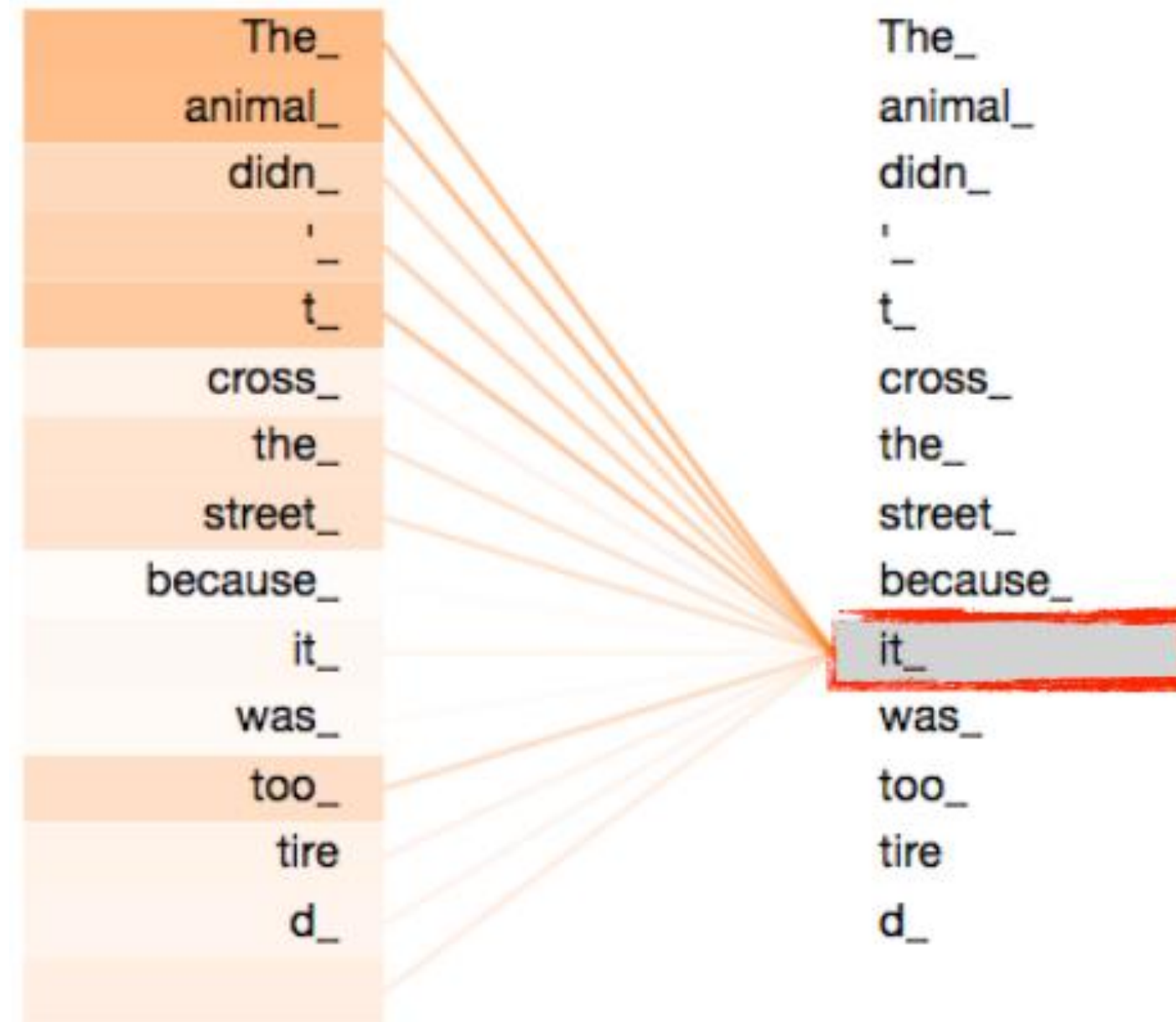
# Self-attention



n개의 입력  $x$ 가 주어지고  $z$ 벡터를 찾는데,  
 $i$ 번째  $x$ 벡터를  $z_i$ 로 바꿀 때 나머지  $n-1$ 개의  $x$ 벡터를 같이 고려

# Self-attention

The animal didn't cross the street because **it** was too tired.





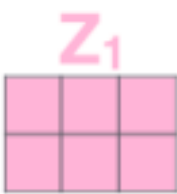
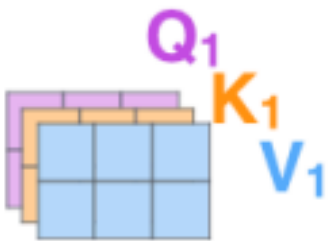
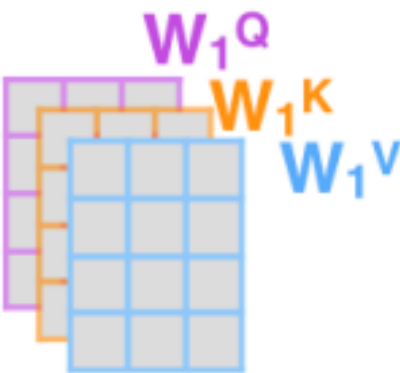
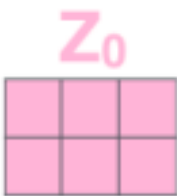
# Self-attention

- 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



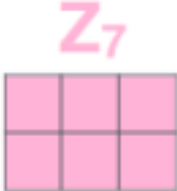
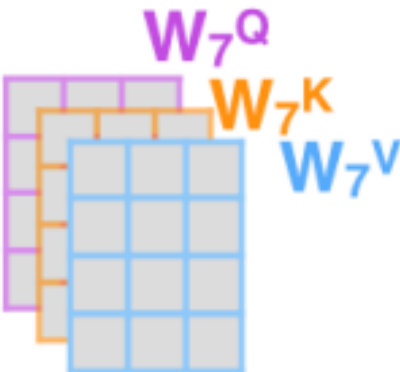
\* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one



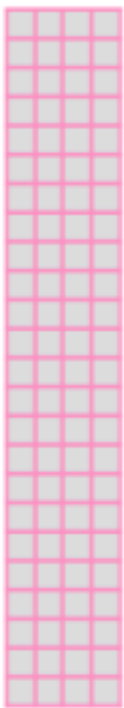
...

...

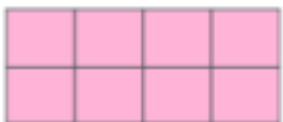
...



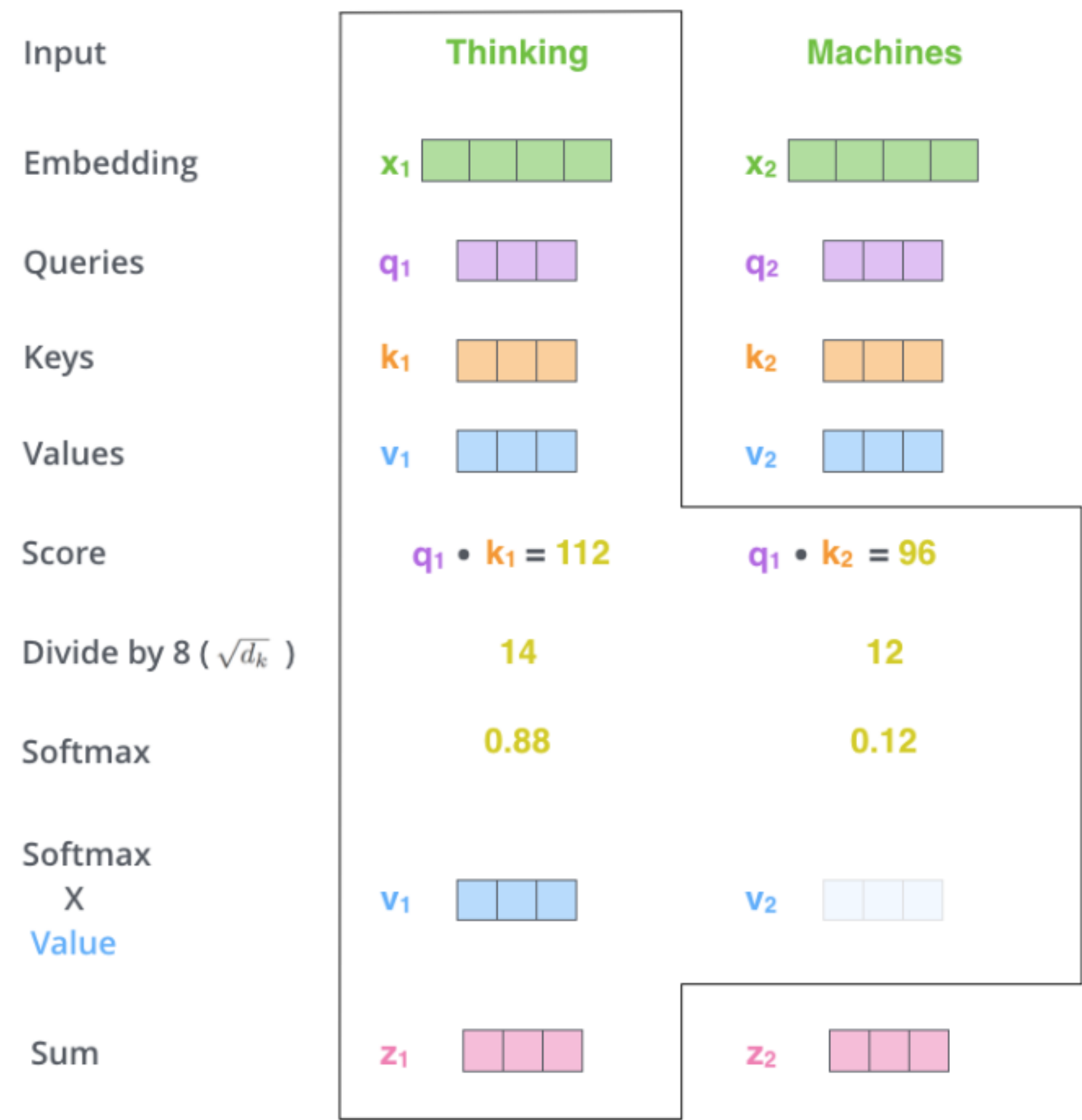
$W^O$



$Z$



# Self-attention



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

=  $Z$



# Vision Transformer

---

NLP task에서 SOTA를 찍었던  
모델인 Transformer를  
CV분야에 적용한 방법론

- CNN보다 inductive bias가 부족

데이터가 적으면 ResNet보다 정확도 낮음

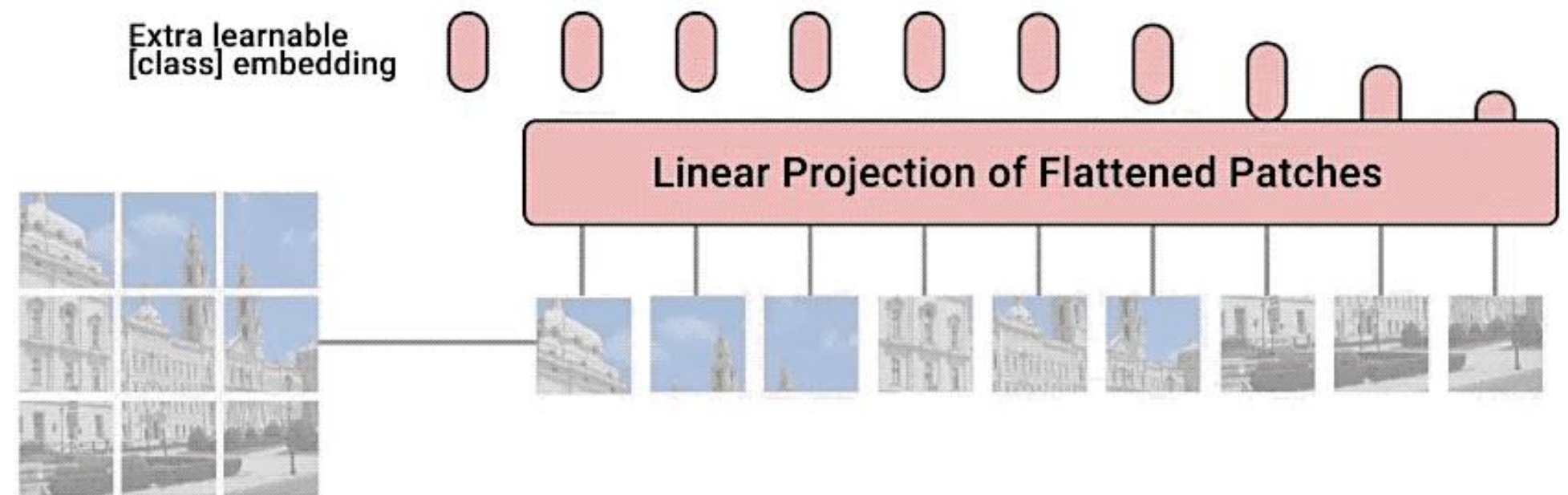
-> 더 많은 데이터를 요구

- Large scale dataset 활용

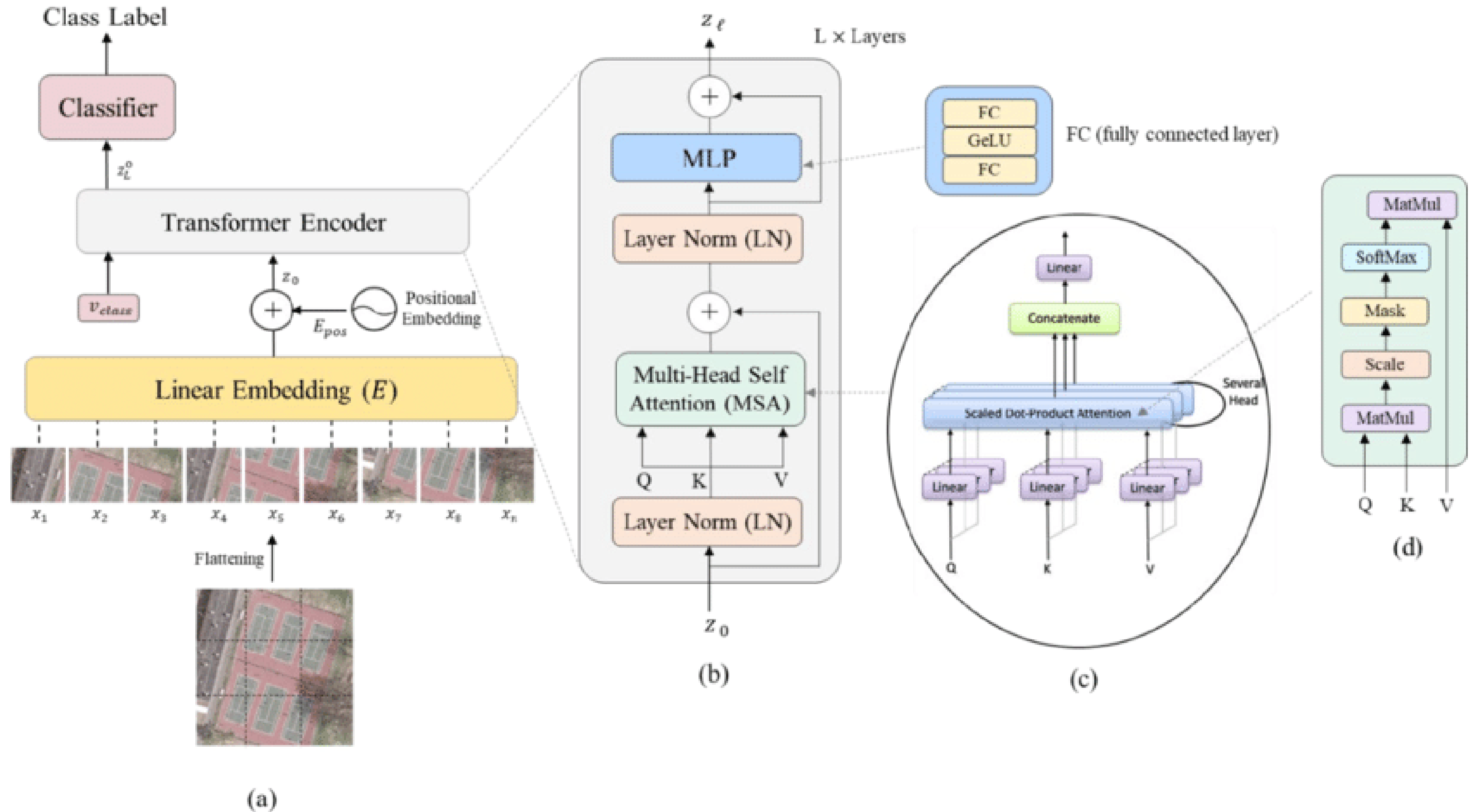
CNN보다 확연히 높은 성능 보임

# Method of ViT

- (1) 이미지 patch로 단위화
- (2) Linear Projection of Flattened Patches
- (3) Class token 추가, positional embedding
- (4) Transformer Encoder 통과
- (5) MLP Head 통과
- (6) Class 분류



# Method of ViT



감사합니다

