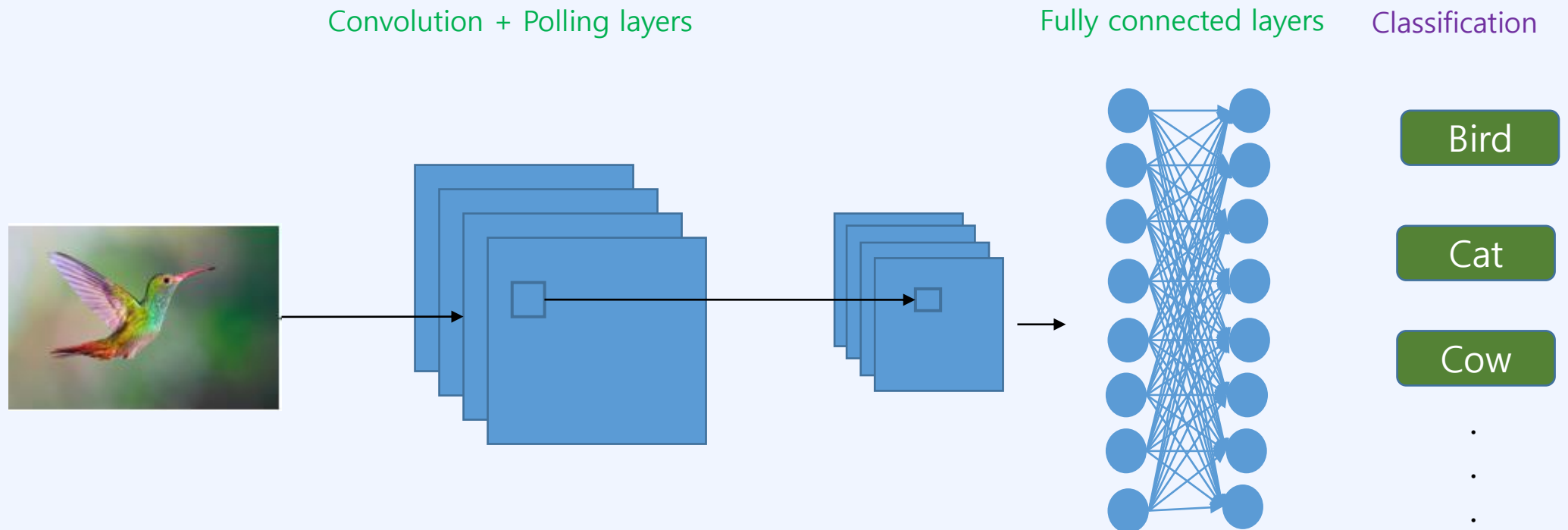


Ch1. ViT

An image is worth 16x16 words - Transformers for image recognition at scale

CNN to ViT

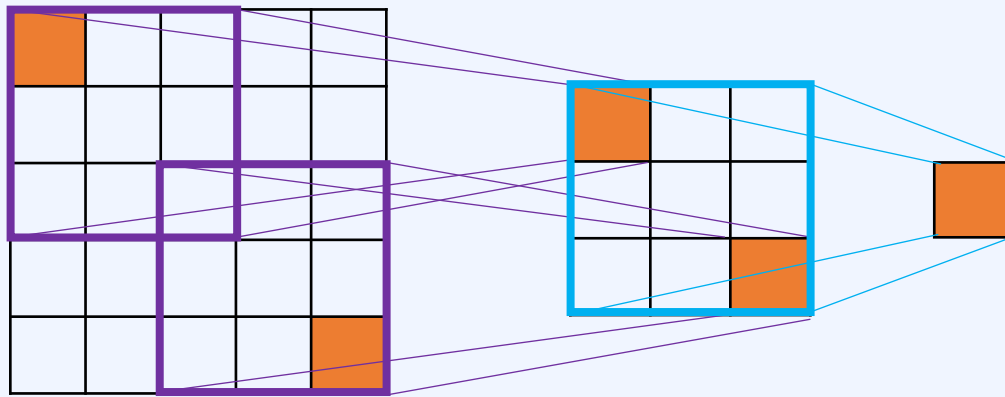
- CNN(Convolutional Neural Network)
 - Computer vision 분야에서 많이 사용
 - Input image의 공간정보를 유지한 채 학습



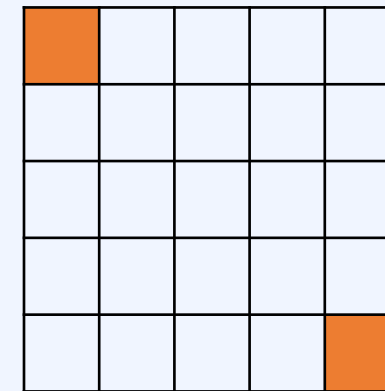
CNN to ViT

- Transformer & CNN

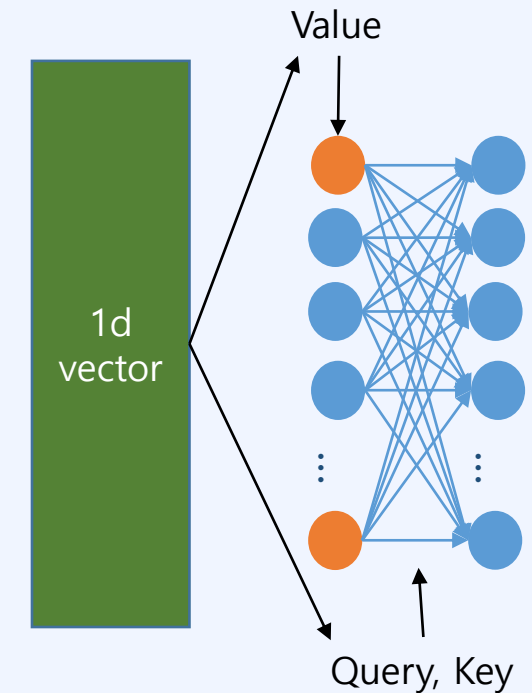
- CNN – Image 전체의 정보를 압축하기 위해 여러 개의 layer를 통과
- Transformer – 하나의 layer로 전체 image 정보를 압축



CNN

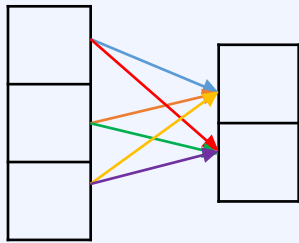


Transformer

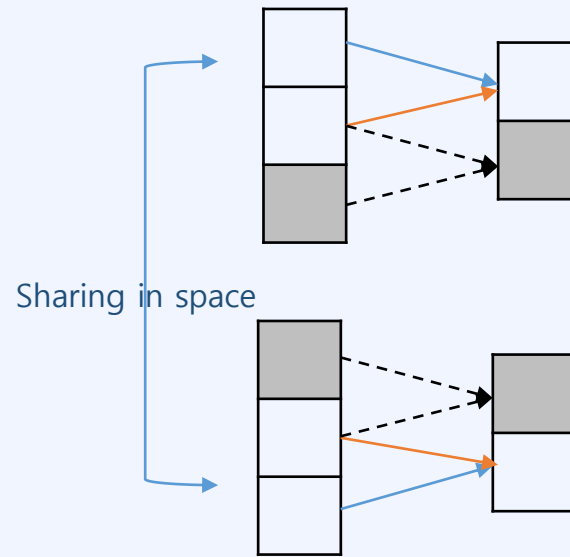


Transformer

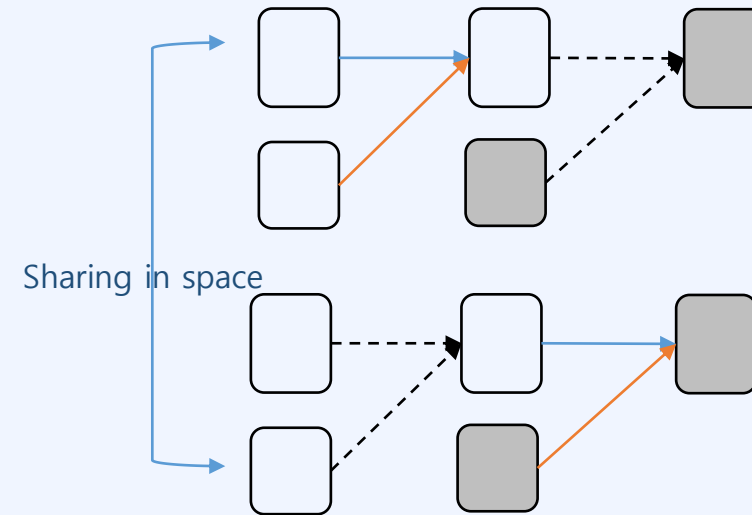
- Inductive bias
 - 주어지지 않은 입력의 출력을 예측



Fully connected



Convolution



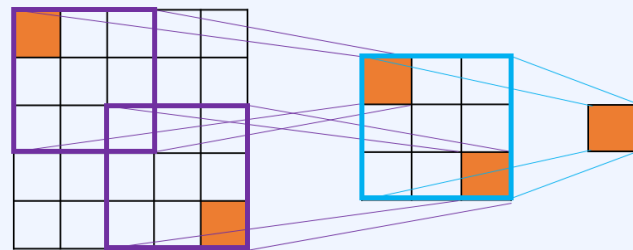
Recurrent

Transformer

- Inductive bias

CNN

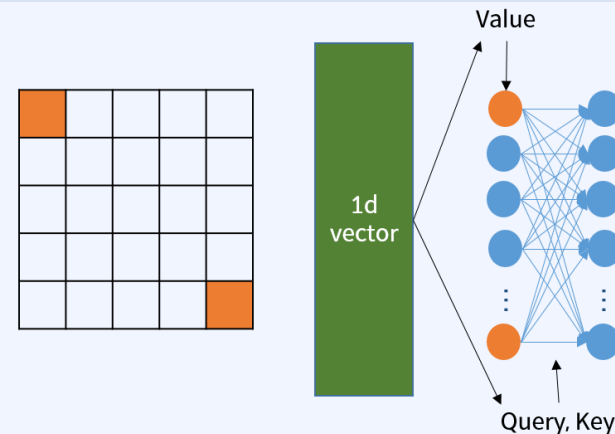
- Convolution filter 사용
- 지역적인 정보 유지o
- 학습 후, 고정된 Weight을 사용



CNN

Transformer

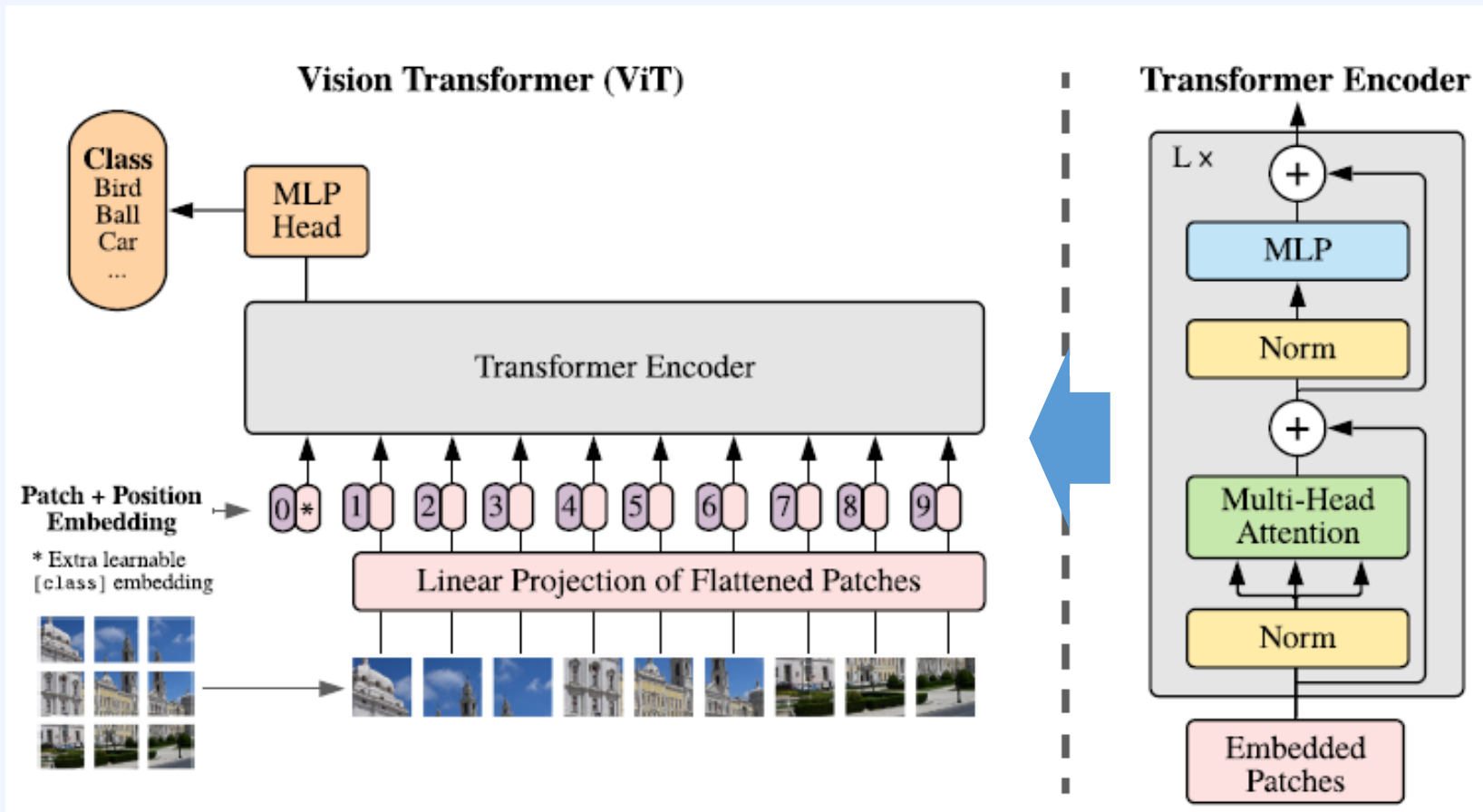
- 임베딩에 의한 벡터 변환 후, Self attention
- 지역적인 정보 유지x
- 학습 후에도 input vector에 따라 Weight이 달라짐
- Inductive bias ↓



Transformer

ViT(Vision in Transformer)

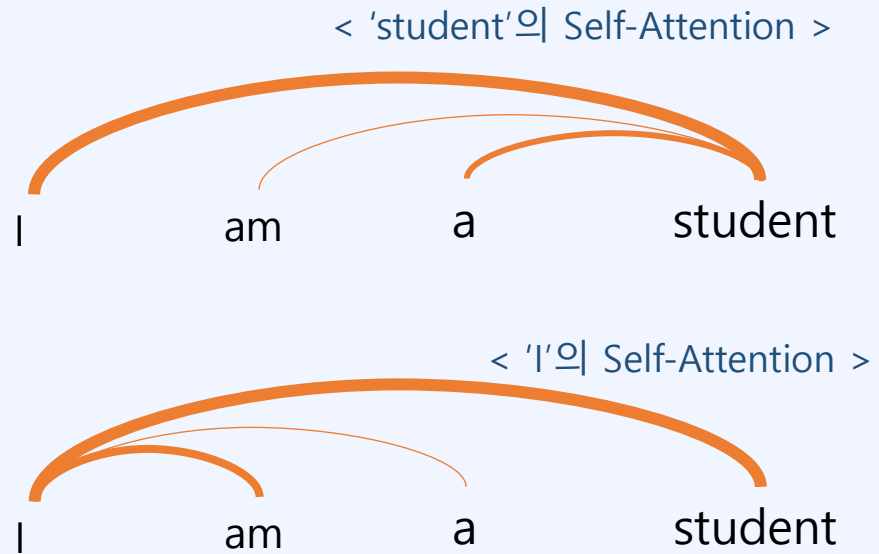
- Architecture



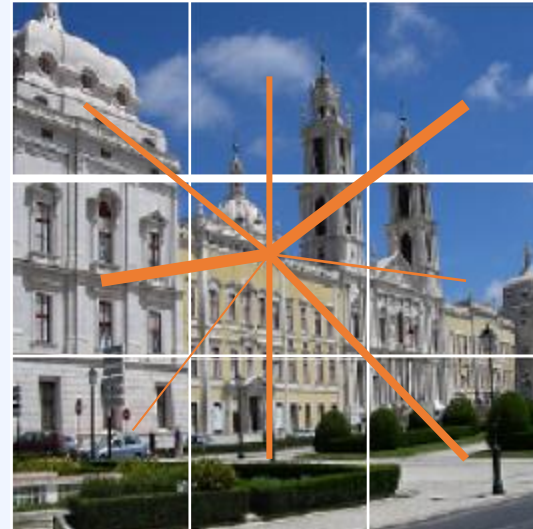
ViT 동작

Self attention

- NLP

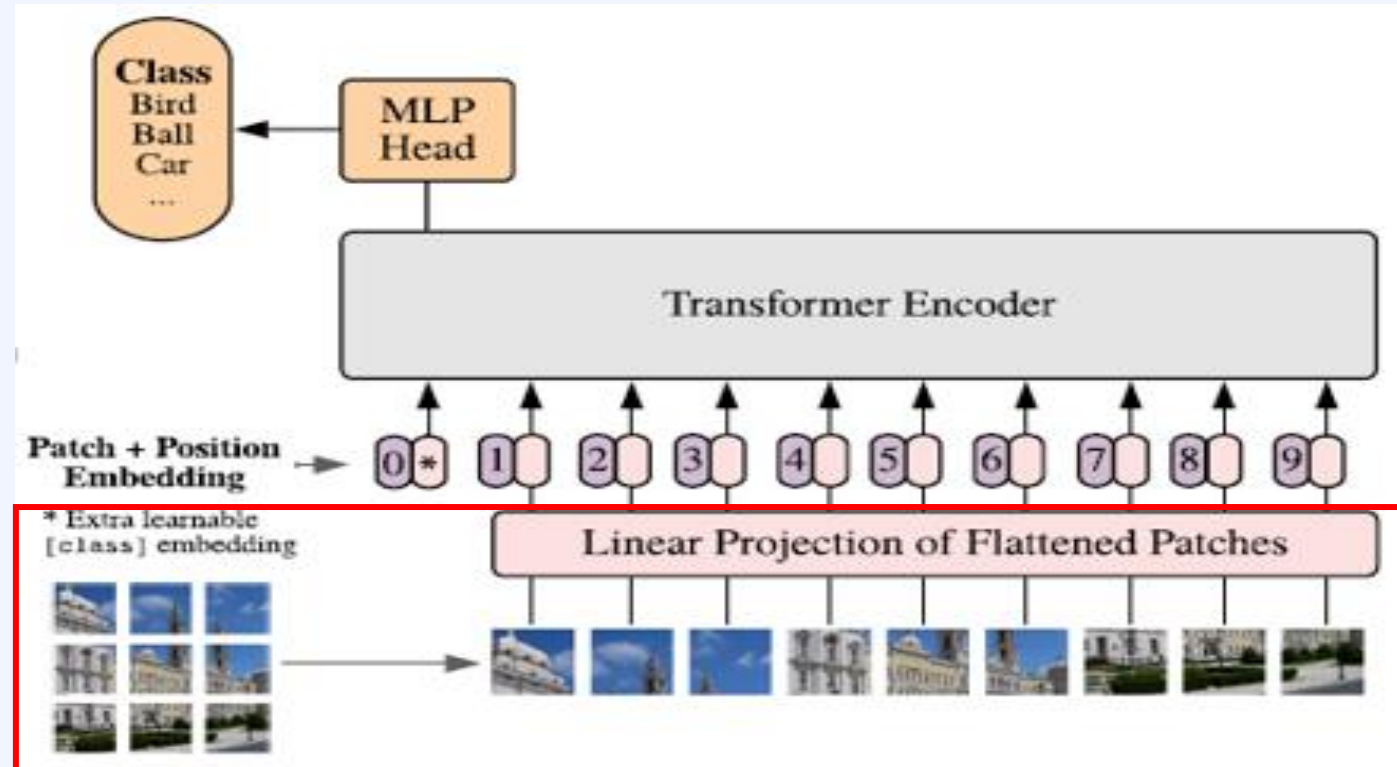


- Vision



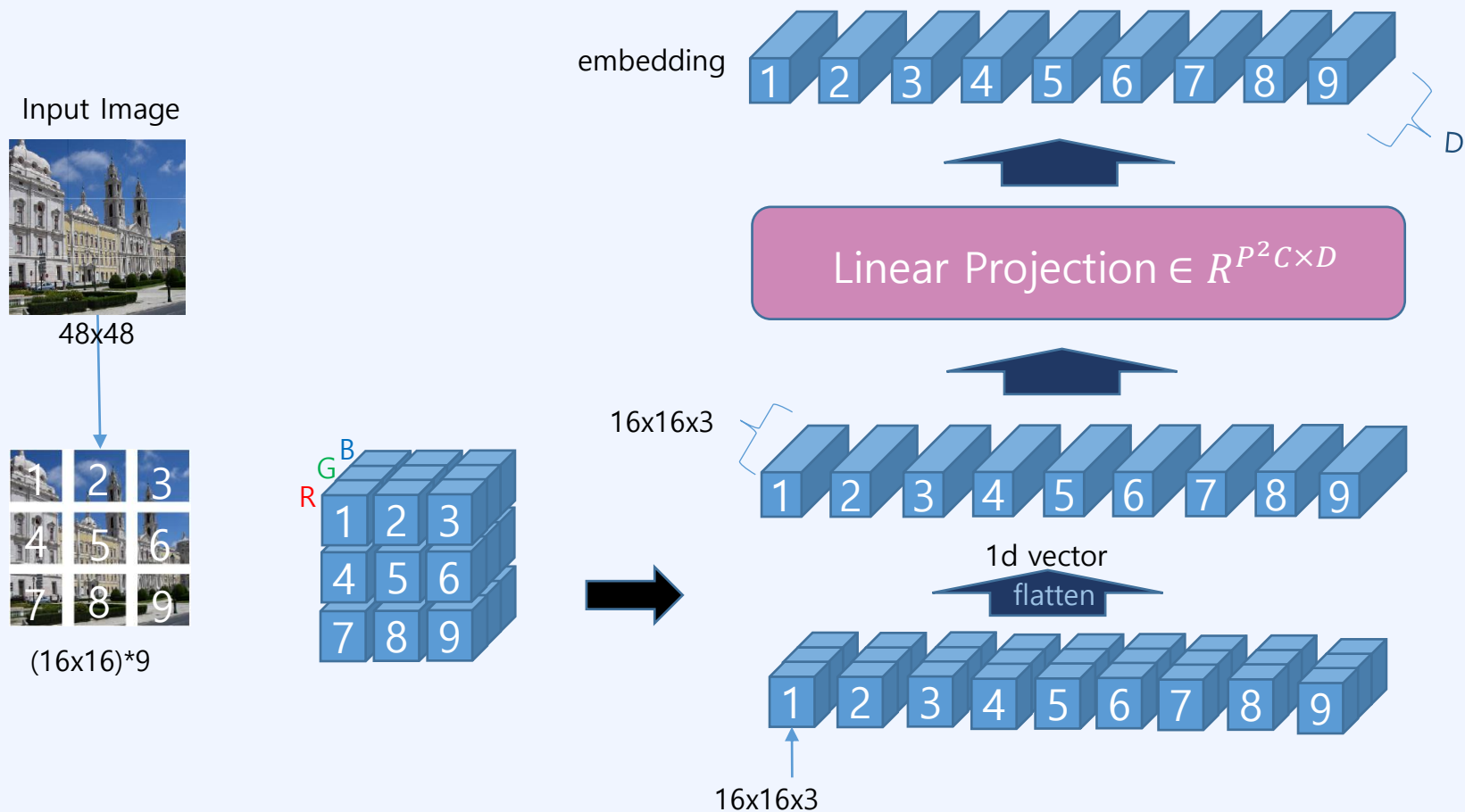
ViT 동작

- Step 1 - Patch embedding



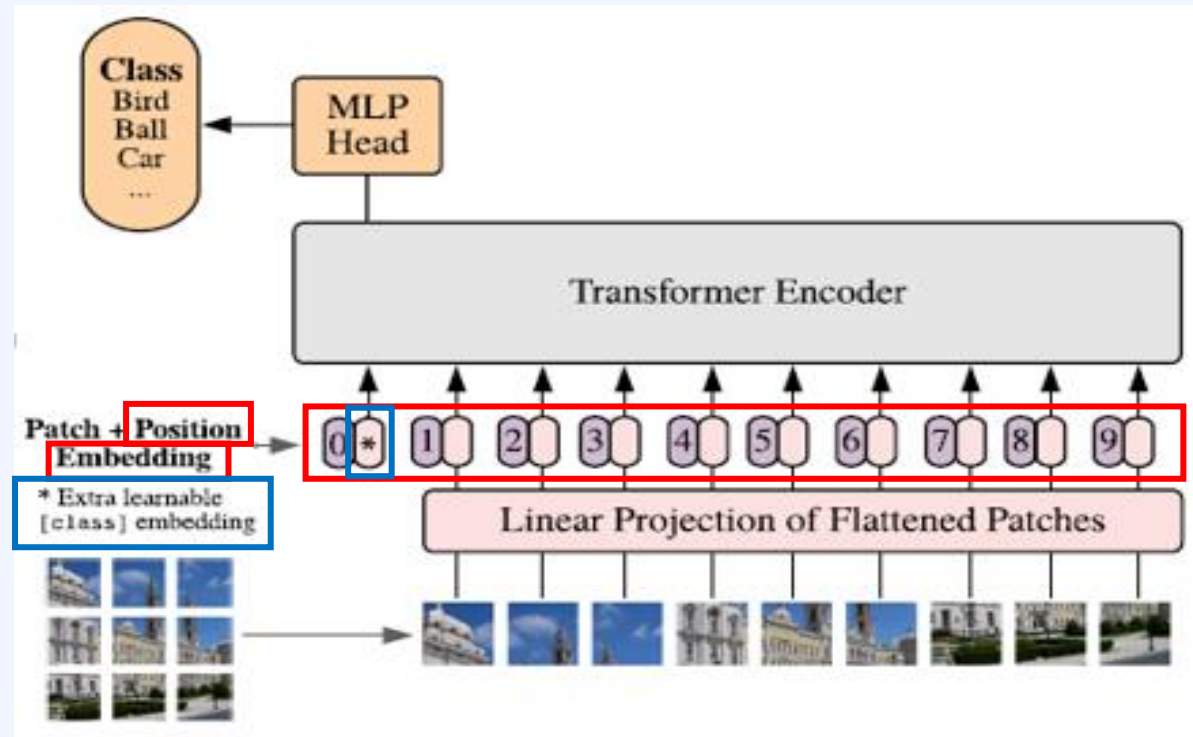
ViT 동작

- Patch embedding(example-ViT/16)



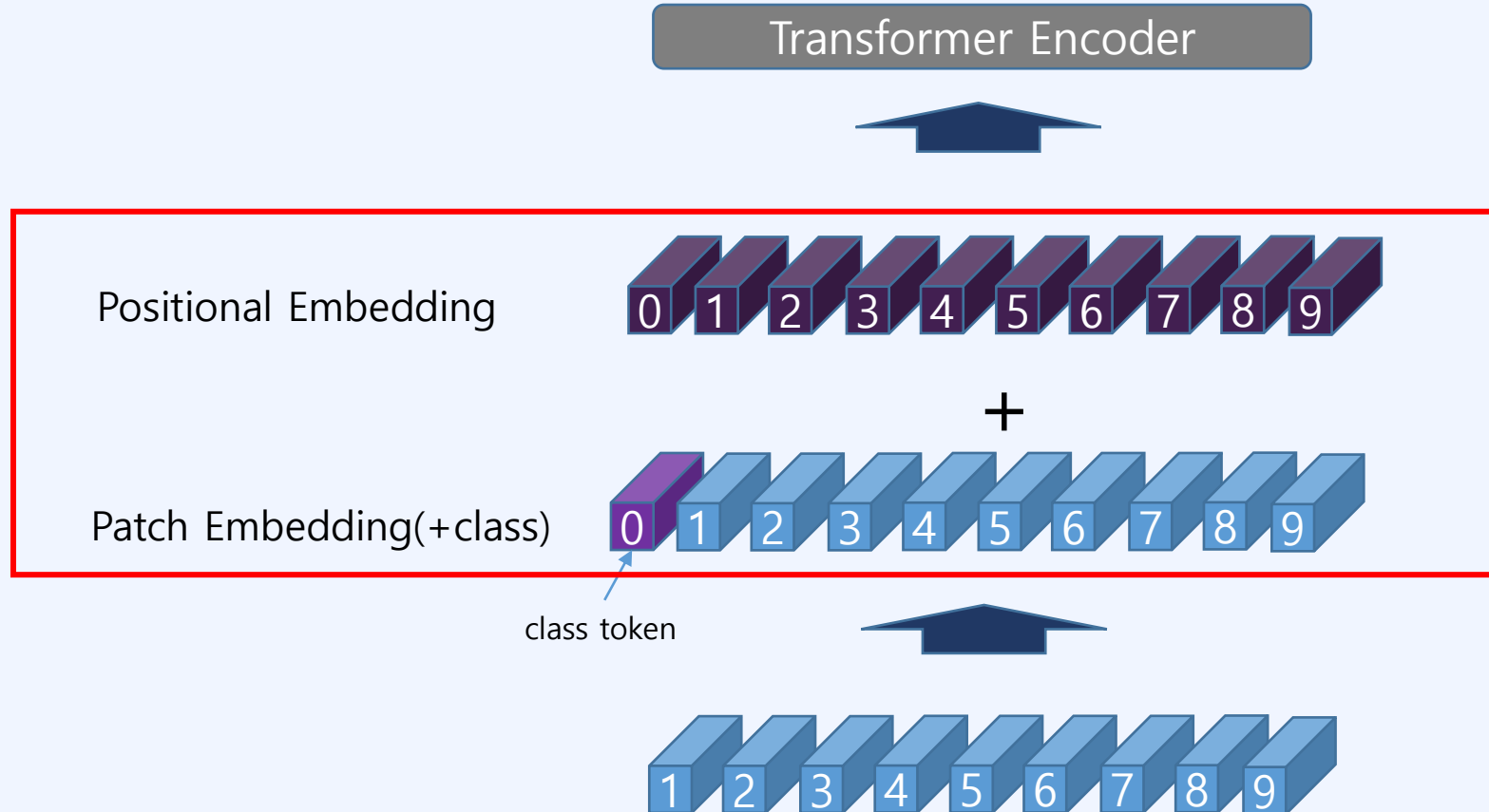
ViT 동작

- Step2 - Embedding patch + Positional Embedding



ViT 동작

- Patch embedding + Positional Embedding (example-ViT/16)



Positional embedding

- 주기함수

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

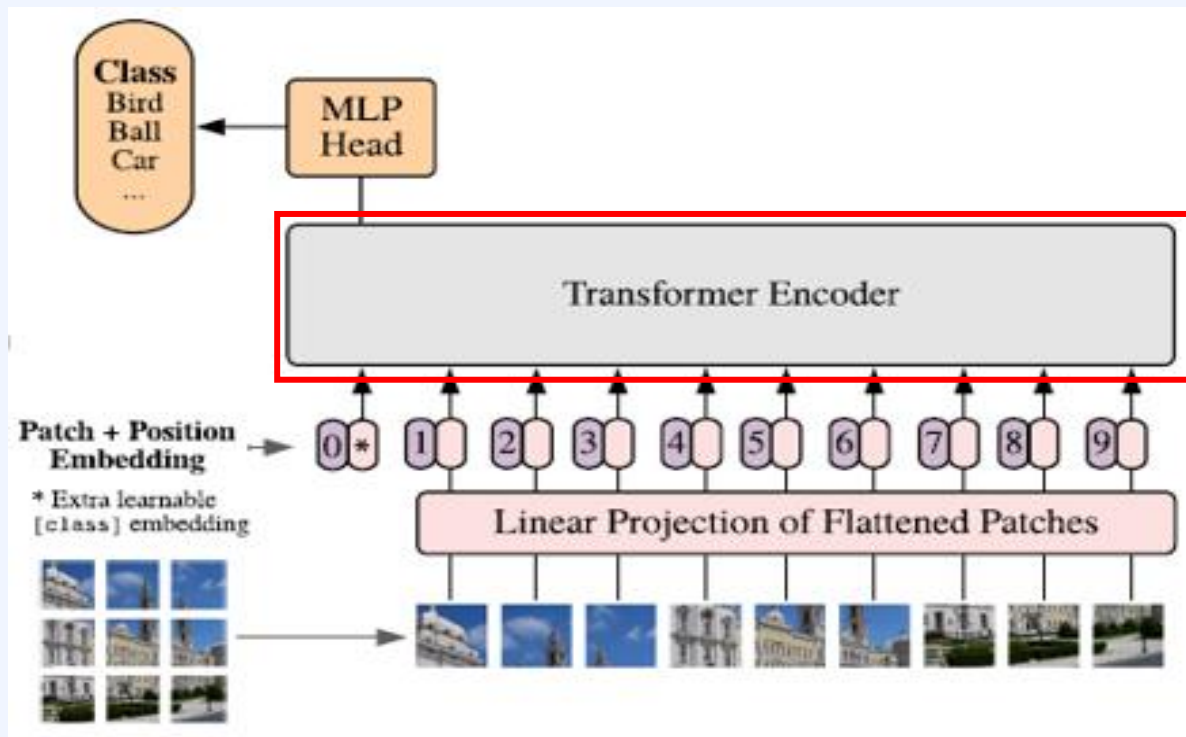
$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{model}})$$

- 예시 (Sequence : 나는 학생 입니다) ($d_{model}=4, pos=1, i=0,1,2,3$)

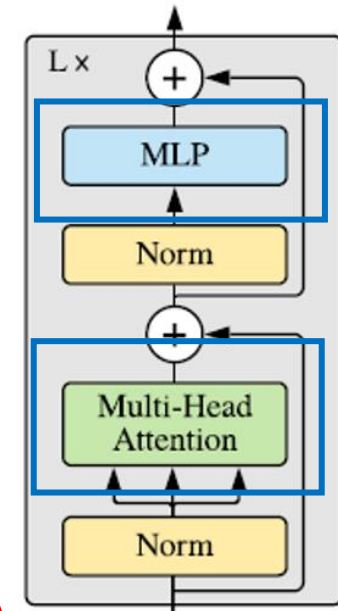
	d_{model}					
나는						
학생	0.3	3.1	3.1	3.1		
입니다						
student	embedding				+	Positional embedding
		0.25	0.99	0.01	0.99	

ViT 동작

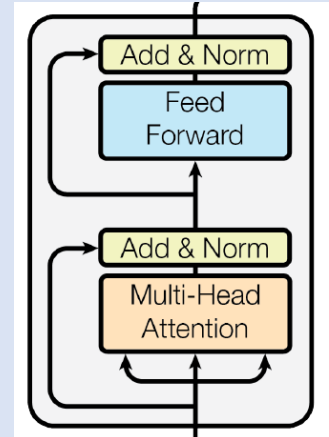
- Step 3 – Transformer encoder



Transformer Encoder



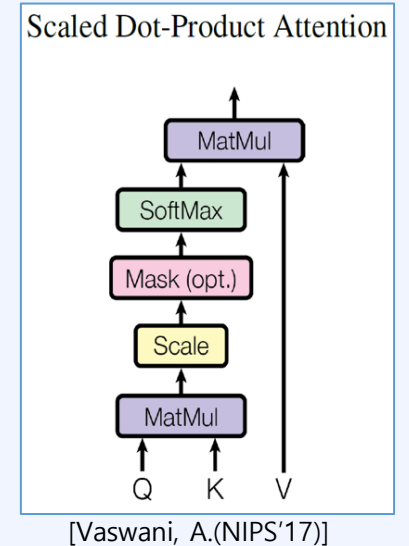
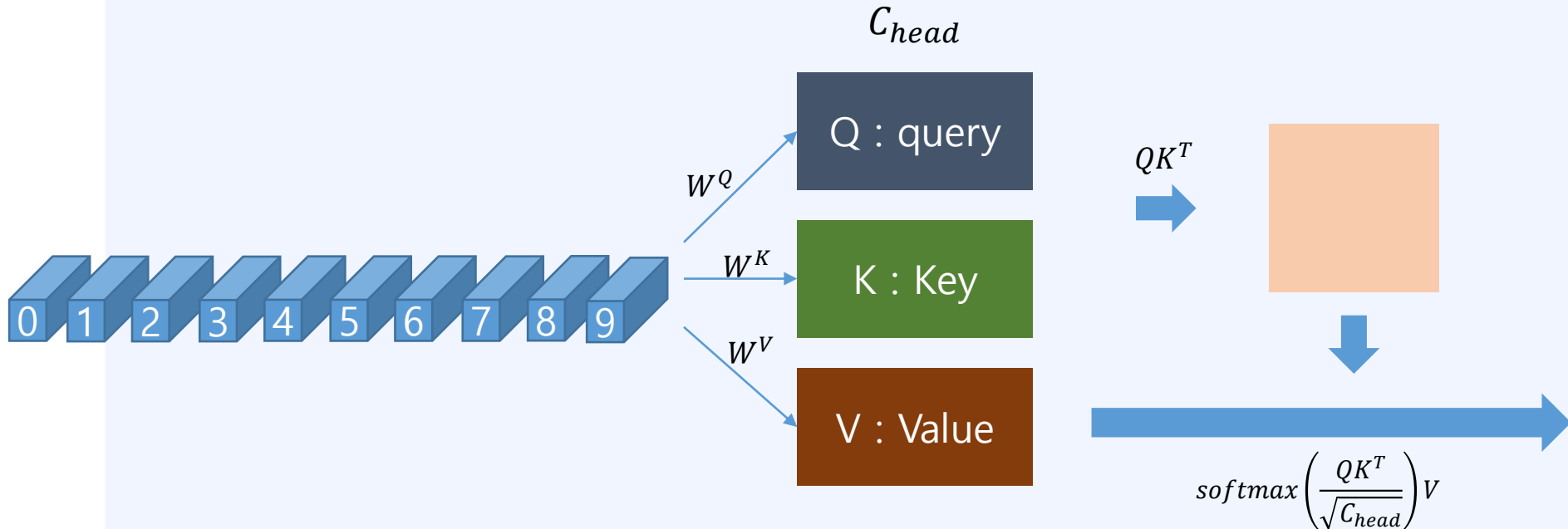
Transformer encoder (standard)



[Vaswani, A.(NIPS'17)]

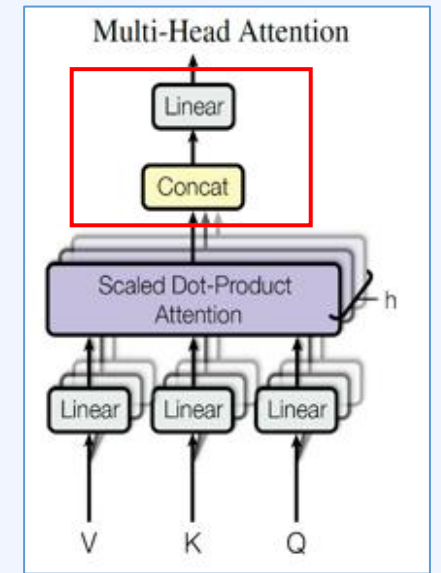
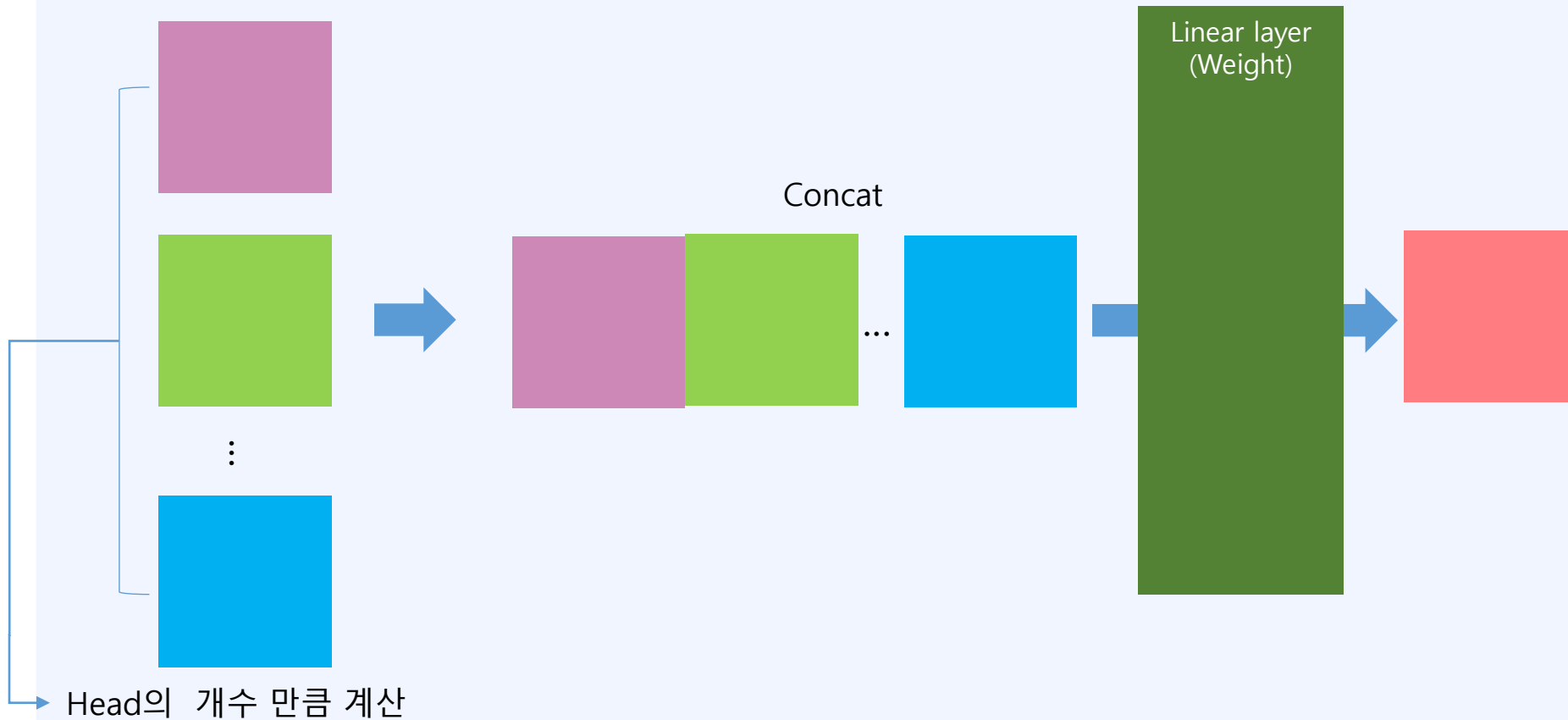
ViT 동작

- Multi-Head Attention
 - Scaled Dot-Product Attention



ViT 동작

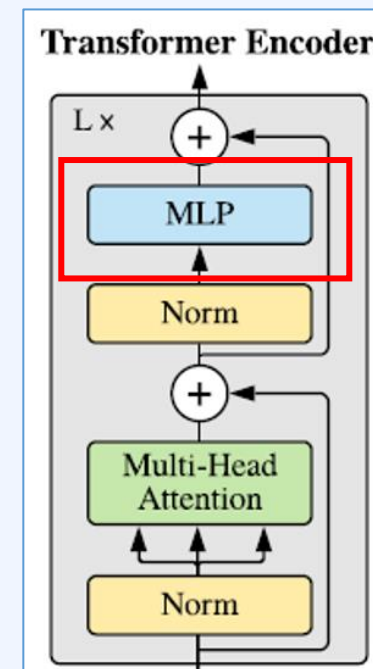
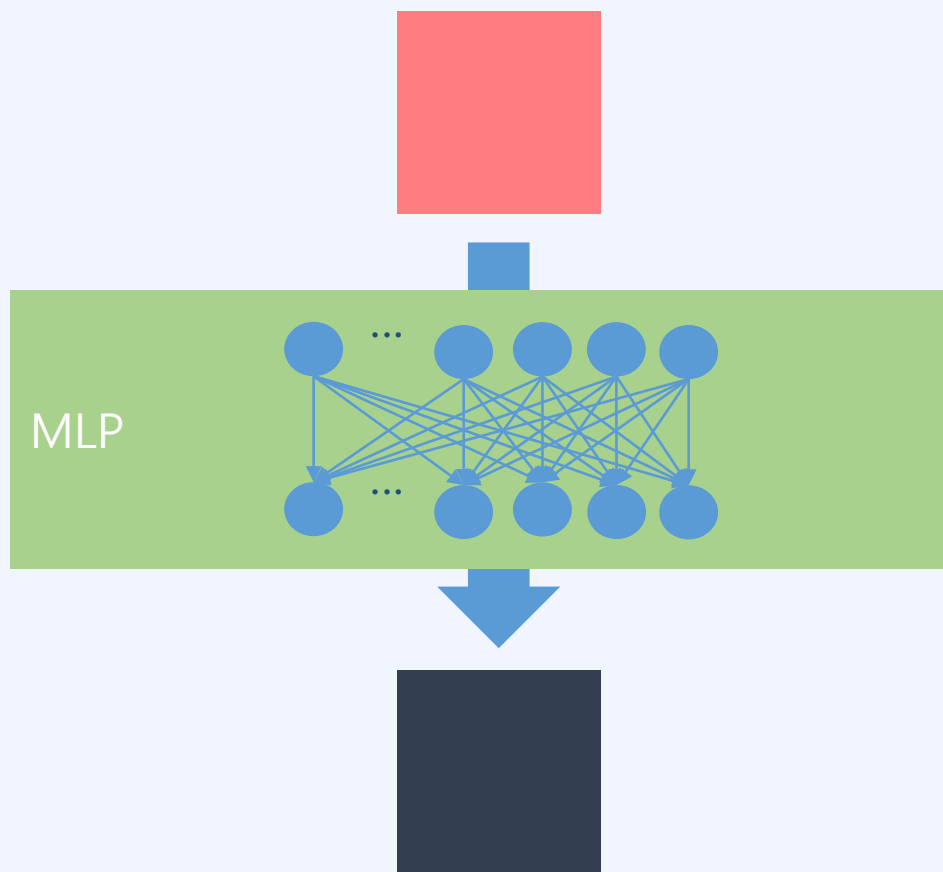
• Multi-Head Attention



[Vaswani, A.(NIPS'17)]

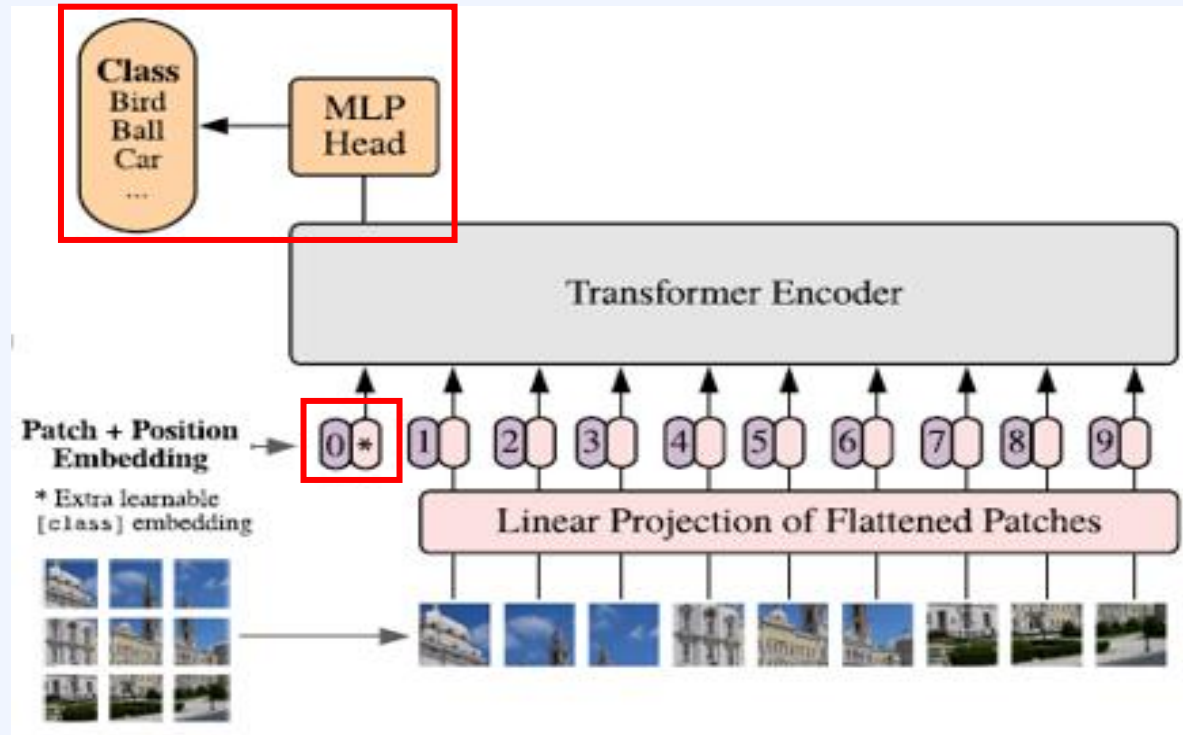
ViT 동작

- MLP



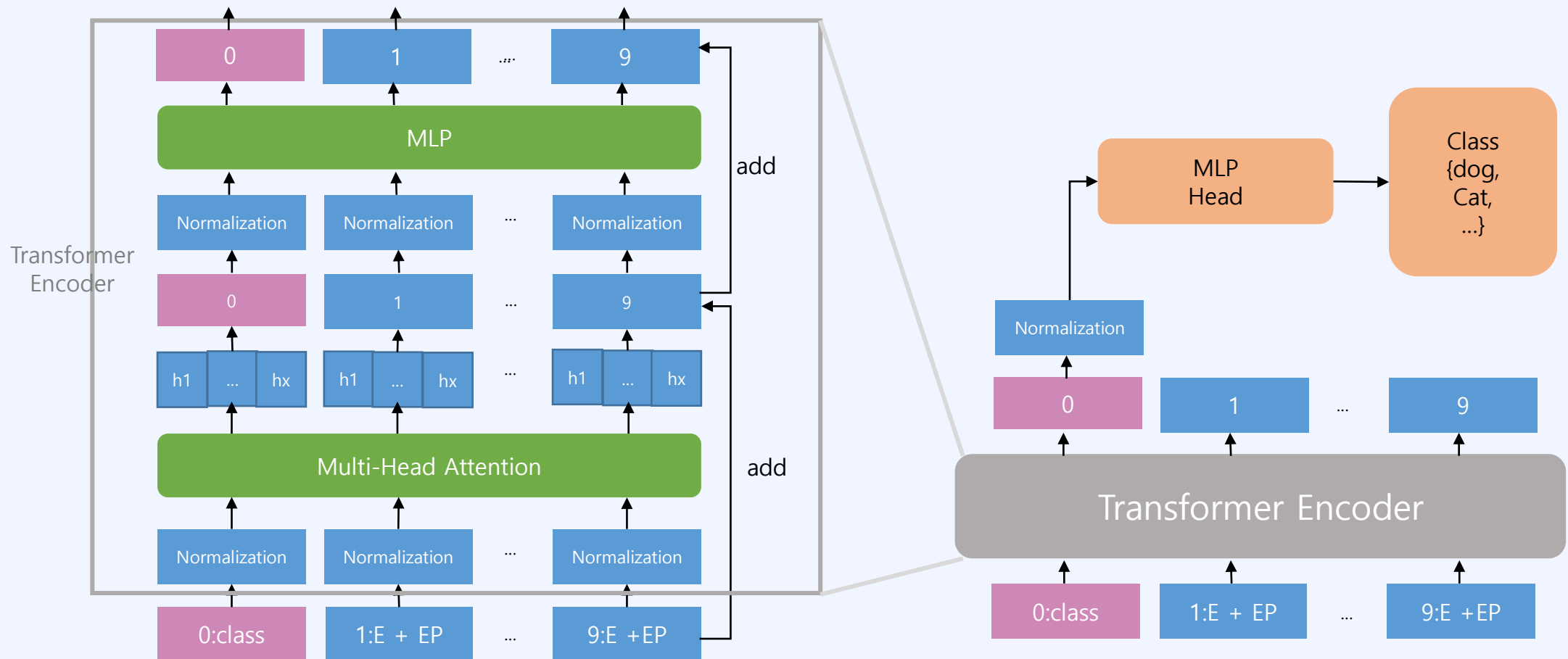
ViT 동작

- Step 4 – MLP Head



ViT 동작

- MLP Head



Summary

