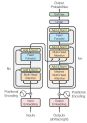


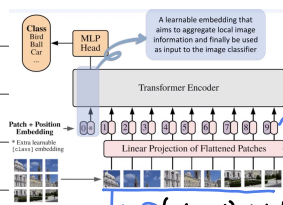
Vision Transformer ViT

- 기존의 Transformer은 아래 그림처럼 Input에는 문장벡터값이 들어갔다.



- 하지만 ViT로 넘어오면서 여러 부분이 바뀐다.

① Input



① (256, 256) 이미지를 (16, 16)으로 나눠서 256개의 패치 생성 후에

1개의 패치당 RGB를 고려해 768크기의 벡터 생성한다.

[256, 768]

③ 클래스 토큰 추가 → $(N+1, D)$

Position Encoding은 $(N+1, D)$ 사이즈를 만들어 그대로 더해준다. → $(N+1, D)$

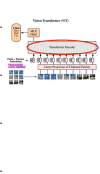
② [768, 모델파라미터 D] 와 연산을 통해 [256, D]로 만든다.

→ (N, D)

② Transformer Encoder

- 1번 반복하기 위해 입력과 출력의 크기가 같도록 유지한다.

- 기존의 Transformer Encoder에서는 Multi-Head 하고 Norm을 하는데 순서가 바뀌어 있다.



$z_L \in \mathbb{R}^{(N+1) \times D}$

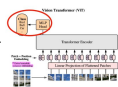
$z_0 \in \mathbb{R}^{(N+1) \times D}$

→ FC층 2개, 4 ELU 사용

- 입력과 입력의 누적정규분포의 곱을 사용한 형태이다.

- 모든 곳에 ReLU, tanh, sigmoid 함수가 아니라 Activation Function에 함수로 사용 가능

- 입력값이 다른 입력에 비해 얼마나 큰지 비율로 값이 조정되어 확률적 해석 가능



MLP Head: LN을 적용하고 FC를 거쳐 \hat{y} 를 생성

$\hat{y} \in \mathbb{R}^C$

C는 클래스 수

$z_L^0 \in \mathbb{R}^D$

클래스 토큰 부분만 분류에 사용

$z_L \in \mathbb{R}^{(N+1) \times D}$

→ 이렇게 나오면 목적에 맞게 활성화 함수를

정하고 모델을 학습한다.