

# ResNet & Transformer & Diffusion Policy

## 1. ResNet

1. ⑤ CNN

↳ gradient vanishing/exploding

degradation: 깊이 ↑ 성능 ↓

2.  $H(x)$  대신  $H(x) = F(x) + \gamma C \frac{z}{2}$   $\frac{z}{2}$   $\frac{z}{2}$   $\frac{z}{2}$

Short Cut

(skip-connection)

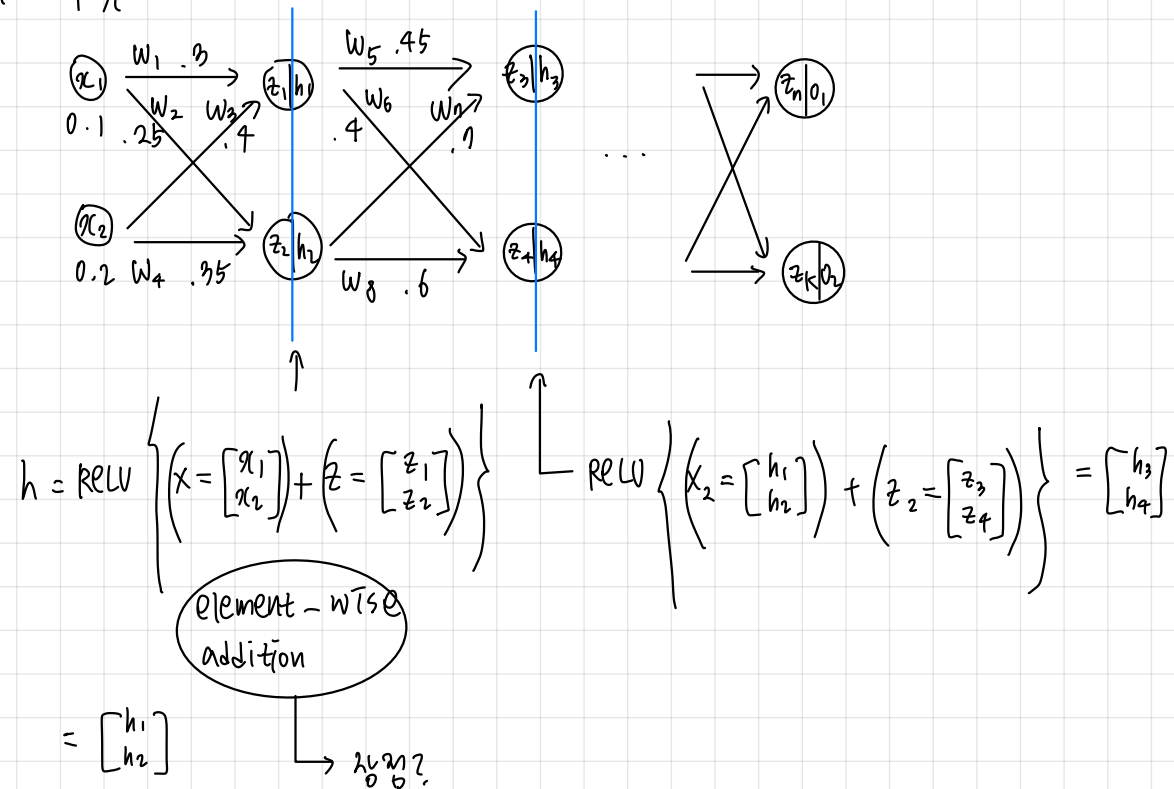
: 가중치  $X$ , bias  $X$

( $F(x), \gamma C$  채널 크기  $X$  :  $H(x) = F(x) + \underline{\underline{W \gamma C}}$ )

down-sampling

mapping  
방식 [ identity: zero-padding  
projection:  $1 \times 1$  conv ⑤

3. "+ 1"



1. 해당 레이어의 필요  $x \rightarrow$  입력으로 무시

↳

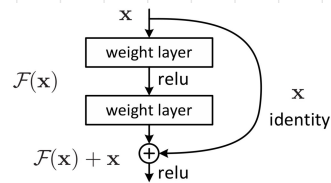


Figure 2. Residual learning: a building block.

보통 등대:  $H(x) = x + F(x) \approx x$

2. 역전파 gradient 때는 항상 인장 (gradient vanishing, exploding x)

↳  $\frac{\partial L}{\partial x} = \frac{\partial L}{\partial y} \cdot \left( \frac{\partial y}{\partial x} \right)$

Plain:  $y = F(x)$   
 $\rightarrow \frac{\partial F(x)}{\partial x}$

Residual:  $y = F(x) + x$   
 $\rightarrow \frac{\partial F(x) + x}{\partial x} = \frac{\partial F(x)}{\partial x} + \underline{\underline{1}}$

## 2. Transformer (Attention is all you need)

### 1. ㉓ RNN ㉓

recurrent

↳ 순차 의존성 :  $h_t = f(h_{t-1}, x_t)$

병렬화 x

### 2. Encoder & Decoder (with formula)

#### 2.1 multi-head attention

입력 seq

Encoder: Encoder self-attention  
FFN:  $w_2 \cdot \text{ReLU}(w_1 x + b_1) + b_2$

$z = (z_1, \dots, z_n)$

Decoder: decoder self-attention (Masked Multi-head attention)  
Encoder-Decoder attention  
FFN

$y = (y_1, \dots, y_n)$

↳ Residual Add & Norm  
" + x "

#### 2.1 attention (Q, K, V)

1) X: token embedding ( $T \times d_{\text{model}}$ ) + Positional Encoding

$$\begin{aligned} Q &= X \cdot W^Q \\ K &= X \cdot W^K \\ V &= X \cdot W^V \end{aligned} \rightarrow \frac{d_{\text{model}}}{h} \text{ (head)}$$

$$\begin{aligned} PE_{(pos, 2i)} &= \sin(pos/10000^{2i/d_{\text{model}}}) \\ PE_{(pos, 2i+1)} &= \cos(pos/10000^{2i/d_{\text{model}}}) \end{aligned}$$

3)

$$\sum \left\{ \text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \right\} = \text{weighted-sum (head 개수만큼)} \rightarrow \text{Concatenate} \rightarrow \text{Linear}$$

weight  
attention score S

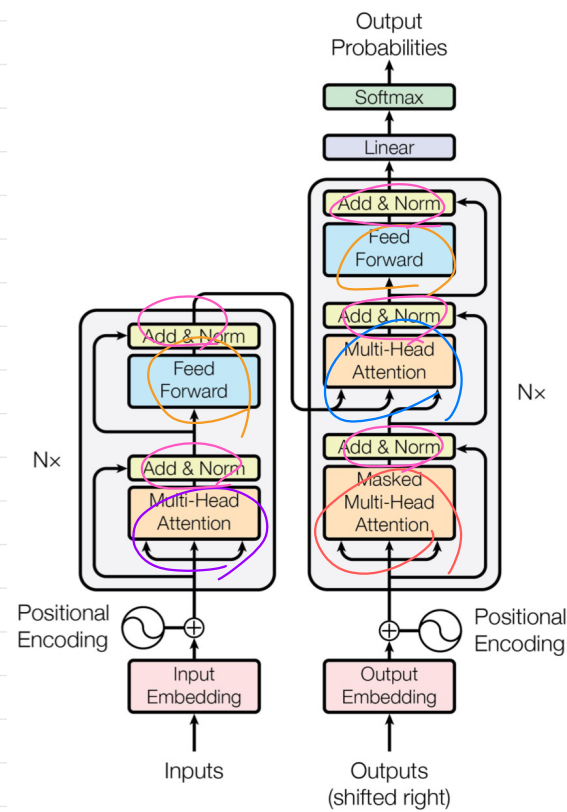


Figure 1: The Transformer - model architecture.

pos 0 1 2 3

ex) 나는 오늘 밥을 먹어.

→  $d_{\text{model}} = 10$  (토큰의 임베딩 차원)

→ "나는": token embedding =  $[e_1, \dots, e_{10}]$

+  
position embedding =  $[p_1, \dots, p_{10}] = [\sin(1/10000^{\frac{0}{10}}), \cos(1/10000^{\frac{0}{10}}), \sin(1/10000^{\frac{1}{10}}), \cos(1/10000^{\frac{1}{10}}), \dots]$

### 3. Diffusion policy

1. training: DDPM + policy  
testing: DDIM

↳ <training>

샘플 행동  $x_0 \rightarrow$  노이즈  $x_k = \sqrt{\alpha_k} x_0 + \sqrt{1 - \alpha_k} \epsilon$   $\rightarrow L = \|\epsilon - \epsilon_\theta(x_k, k, c)\|^2$  손실 (gradient)

$\alpha_k$ : hyperparameter  $\downarrow$   $\nabla_k$

$-\frac{1}{\nabla_k} \epsilon_\theta(x_k, k, c) = \nabla_{x_k} \log p(x_k | c)$  손실 (행동 시 가짜 손실)

<testing>: Langevin 동적학 (denoising)

$x_{k-1} = x_k + \alpha_k \cdot S_\theta(x_k, k, c) + \sigma_k \xi_k$

$\alpha_k = k_{inter}$   $\downarrow$  gradient 적용 정도

$\sigma_k$  or  $\xi_k$  또는 gaussian noise  
: denoising 정도를 평준화  
상태의 행동을 유지

### 2. Explicit / Implicit / Diffusion policy

1) Explicit

$x_0 \xrightarrow{F_\theta} \hat{a}$  : 바르 행동 예측

2) Implicit

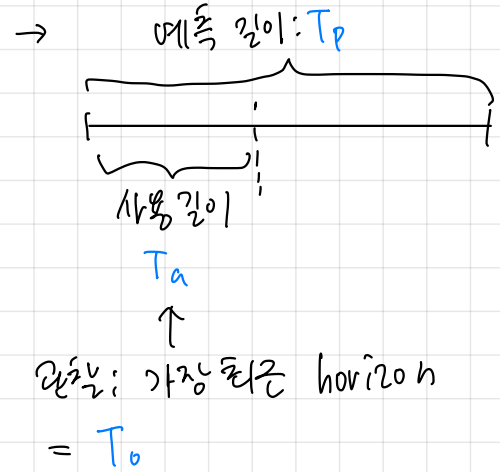
$o, a \xrightarrow{E_\theta} \arg \min E \rightarrow \hat{a}$

3) Diffusion policy

$o \rightarrow \epsilon_\theta(o, a) \rightarrow \nabla E \& k_{inter} \rightarrow \hat{a}$

3. 다음과 같은 행동 분포의 표현  
연속된 클러스터링  
인정적 행동

4. ① closed-loop action seq (receding)



② visual conditioning

→  $D_t$ 는 C 2D인 하인 클러스터링 X (visual Encoder 사용: resnet18)

③ CNN vs transformer

↓  
over-smoothing "토큰화"