# Distinguished Engineering

## Transformer 5/5
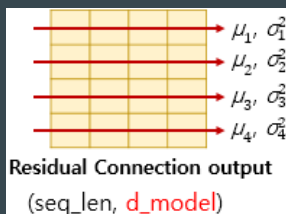## - Transformer, the dark knight

● ● ●

BW

# Plan

- Prologue, seq2seq

- **Attention, please**

- Transformer, a new hope

- Transformer, revenge of the fallen

- Transformer, vision

# Transformer

- Residual connection, Layer Normalization
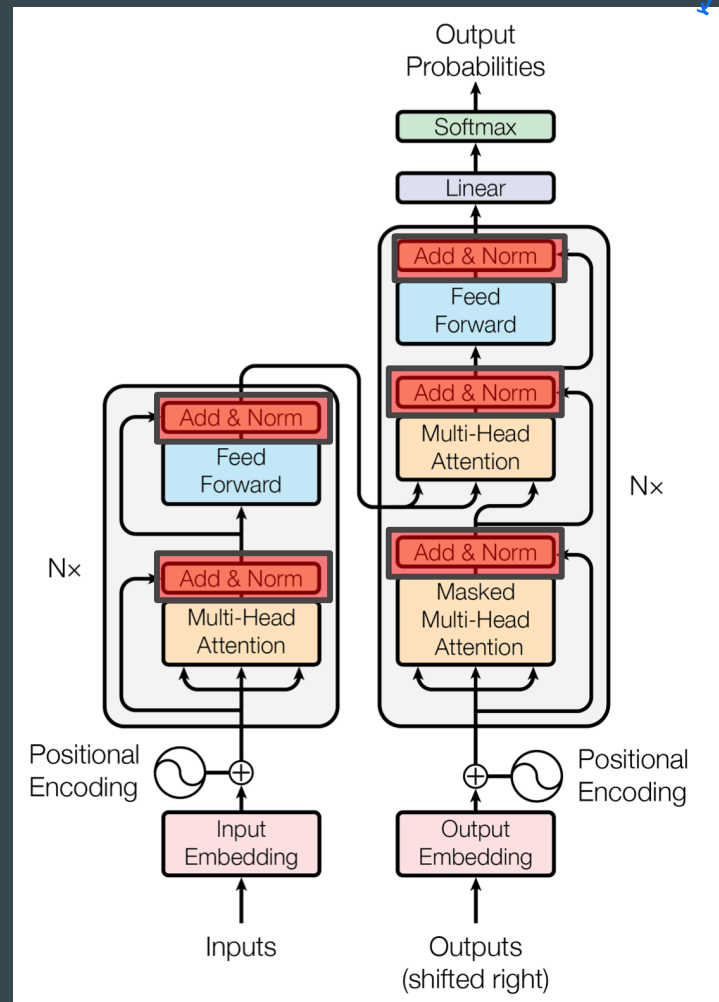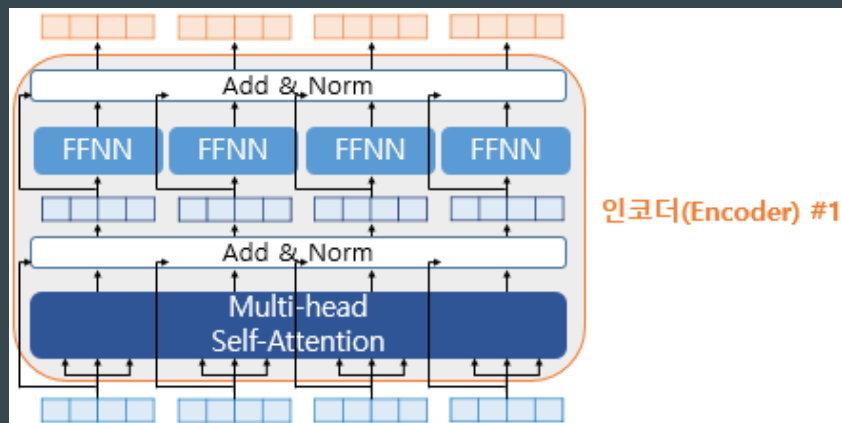


Residual Connection output
(seq_len, **d_model**)

$$ln_i = LayerNorm(x_i)$$

$$\hat{x}_{i,k} = \frac{x_{i,k} - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}}$$
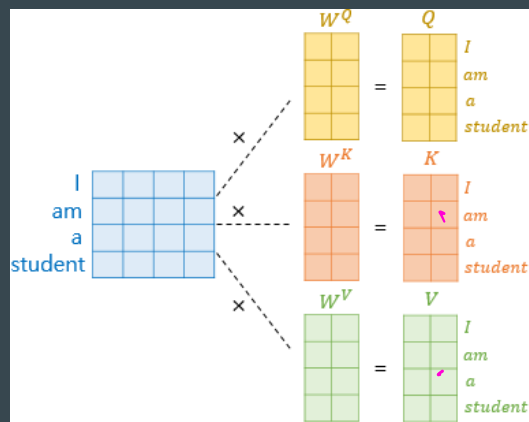
$$ln_i = \gamma\hat{x}_i + \beta = LayerNorm(x_i)$$

# Transformer

Scaled Dot-Product Attention

Multi-Head Attention

- ## Multi-head Attention

$W^Q$  $Q$

I am a student

$W^K$  $K$

I am a student

$W^V$  $V$

I am a student

$Q$  ×  $K^T$  =  I am a student

I am a student

Attention Value Matrix $a$

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

$$Attention(Q, K, V) = softmax(\frac{QK^T}{\sqrt{d_k}})V$$

$d_{model} = d_v \times \text{num\_heads}$

$a_0$ $a_1$ $a_2$ $a_4$ $a_5$ $a_6$ $a_7$ $a_8$

concatenate

$d_{model} = d_v \times \text{num\_heads}$

$W^O$

seq_len

concatenated matrix

$d_v \times \text{num\_heads}$

$d_{model}$

= $d_{model}$ seq_len

**Multi-head attention matrix**

- ## FFNN

$x$

$F_1 = xW_1 + b_1$

활성화 함수 : $ReLU$
$F_2 = max(0, F_1)$

$F_3 = F_2 W_2 + b_2$

인코더(Encoder) #2

FFNN  FFNN  FFNN  FFNN

Multi-head
Self-Attention

인코더(Encoder) #1

인코더 #1의 출력
인코더 #2의 입력

인코더 #1의 입력

# An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale, ICLR2021

**Vision Transformer (ViT)**

**Transformer Encoder**

The MLP contains two layers with a GELU non-linearity.

$$\mathbf{z}_0 = [\mathbf{x}_{\text{class}}; \mathbf{x}_p^1\mathbf{E}; \mathbf{x}_p^2\mathbf{E}; \cdots ; \mathbf{x}_p^N\mathbf{E}] + \mathbf{E}_{pos}, \qquad \mathbf{E} \in \mathbb{R}^{(P^2 \cdot C) \times D}, \ \mathbf{E}_{pos} \in \mathbb{R}^{(N+1) \times D} \quad (1)$$

$$\mathbf{z}'_\ell = \text{MSA}(\text{LN}(\mathbf{z}_{\ell-1})) + \mathbf{z}_{\ell-1}, \qquad \ell = 1 \ldots L \quad (2)$$

$$\mathbf{z}_\ell = \text{MLP}(\text{LN}(\mathbf{z}'_\ell)) + \mathbf{z}'_\ell, \qquad \ell = 1 \ldots L \quad (3)$$

$$\mathbf{y} = \text{LN}(\mathbf{z}_L^0) \quad (4)$$

| Model | Layers | Hidden size $D$ | MLP size | Heads | Params |
|-------|--------|-----------------|----------|-------|--------|
| ViT-Base | 12 | 768 | 3072 | 12 | 86M |
| ViT-Large | 24 | 1024 | 4096 | 16 | 307M |
| ViT-Huge | 32 | 1280 | 5120 | 16 | 632M |

Table 1: Details of Vision Transformer model variants.

# ViT

- +
  - Simple (almost) architecture & well-proved performance
  - Scalability
  - Less Training Time
  - Excellent performance

- -

  - Less inductive bias —> requires more data!!
    - Less data, less performance

# ?

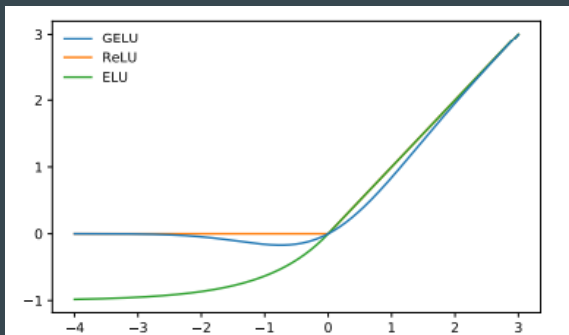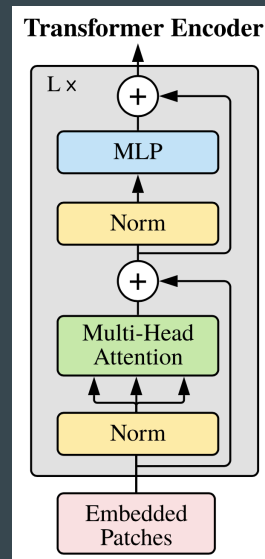- Order

- Embedding
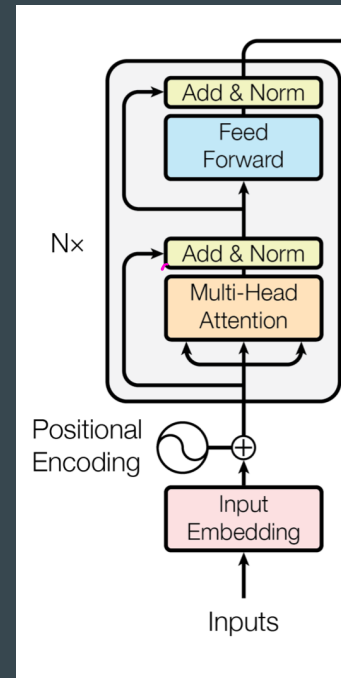
- Activation
  - GELU vs ReLU



Figure 1: The GELU ($\mu = 0, \sigma = 1$), ReLU, and ELU ($\alpha = 1$).

$$\mathrm{GELU}(x) = xP(X \leq x) = x\Phi(x) = x \cdot \frac{1}{2}\left[1 + \mathrm{erf}(x/\sqrt{2})\right]$$
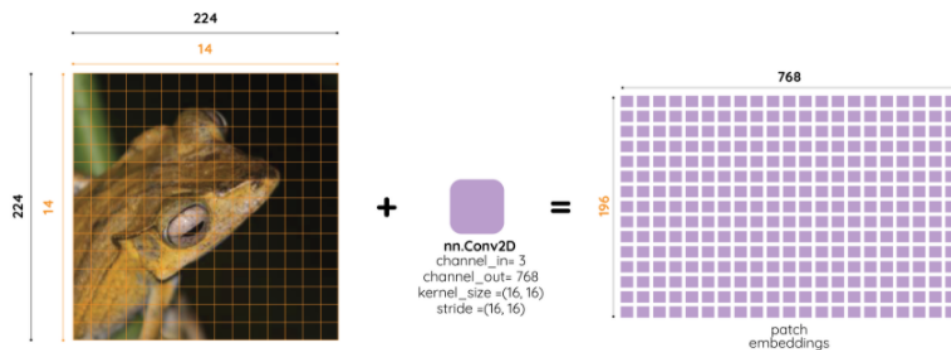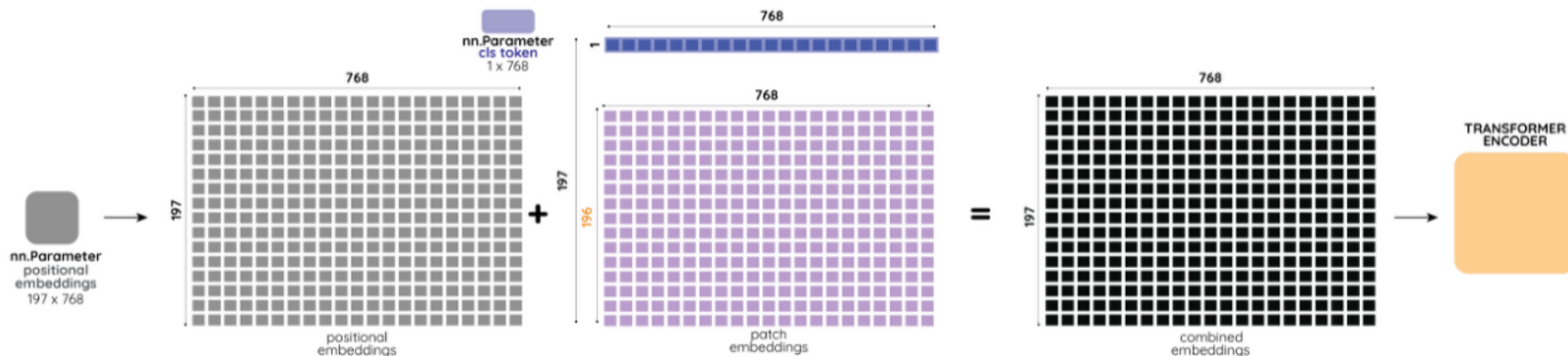


**ViT**



**Transformer**

# ViT

- Positional Embedding



```
# Input image [B, C, H, W]
x = torch.randn(1, 3, 224, 224)

# 2D conv
conv = nn.Conv2d(3, 768, 16, 16)
x = conv(x) # [B, 768, 14, 14]
x = x.reshape(B, -1, 196).transpose(1, 2) # [B, 196, 768]
```

Everything can be found in
https://wikidocs.net/book/2155
https://github.com/ukairia777/tensorflow-nlp-tutorial
https://hongl.tistory.com/232
https://dev-woong.tistory.com/38