



# ViT

## 1. Concepts

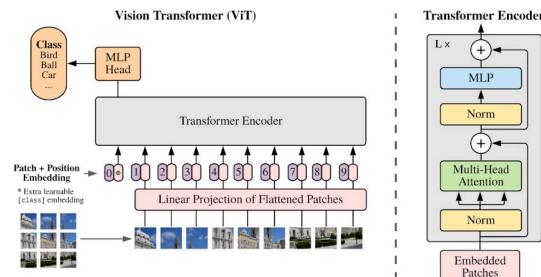
CNNs inductive bias

translation equivariance  
locality

inductive bias: The inherent assumption or conditions that a model possesses by default.

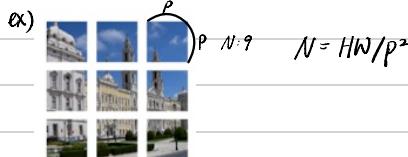
Due to Transformer doesn't have these inductive biases, it has to learn from data.  
Therefore, ViT needs more data.

## 2. Methods

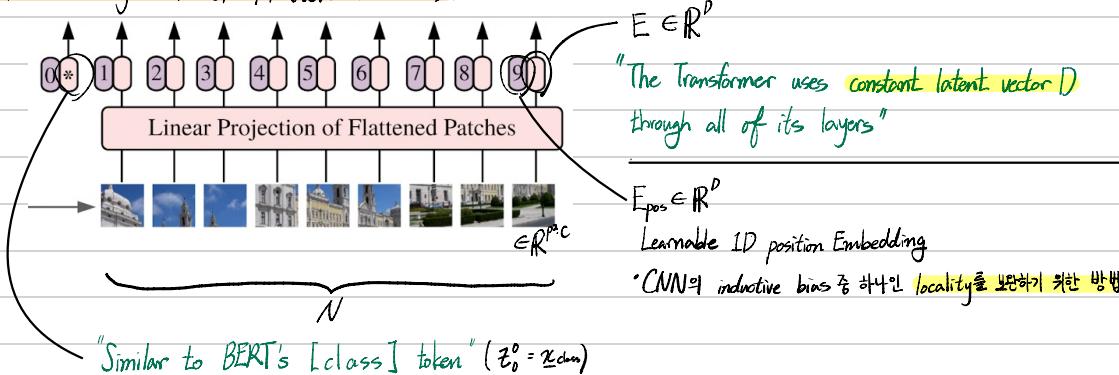


### ① Patch Partitioning

$X \in \mathbb{R}^{H \times W \times C}$  reshape  $\rightarrow X \in \mathbb{R}^{N \times (P^2 \cdot C)}$    
 raw image flattened 2D Patches.  $N$ : number of patches / length of input sequence.  
 $P$ : patches' H, W.

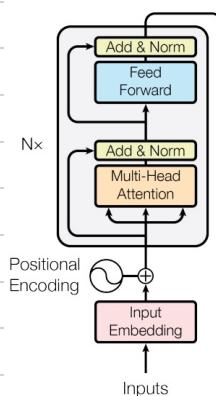


## ② Linear Projection of Flatten Patches.

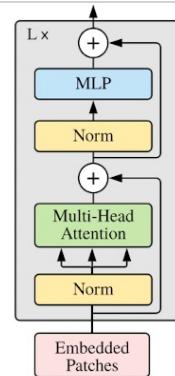


## ③ Transformer Encoder Architecture.

Encoder of Trm



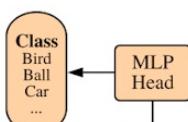
Encoder of ViT.



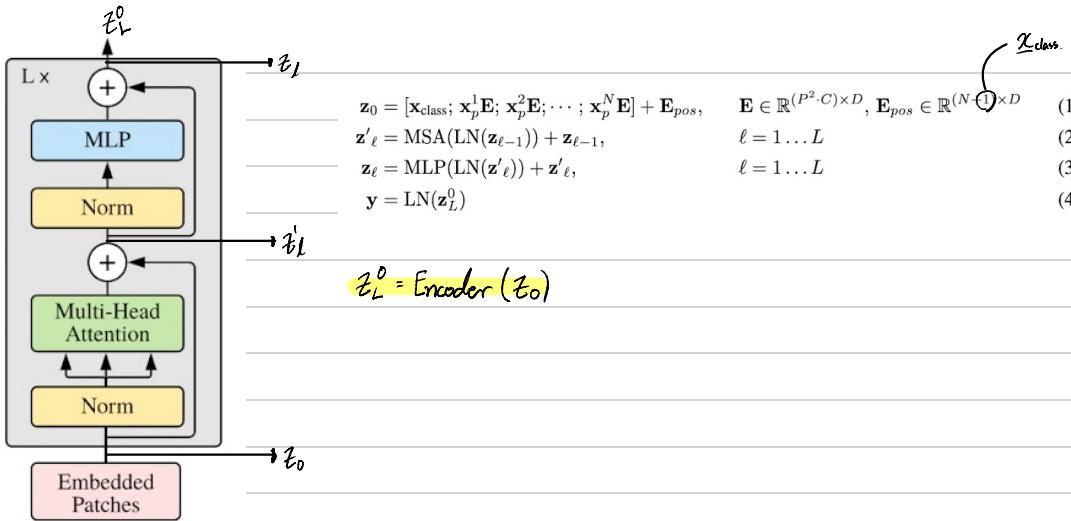
## Pre-Layer Normalization's Strength

- Improvement training stability.
- Residual Connection은 꼭 필요!

## ④ Pre-training vs. Fine-tuning.



"MLP with one hidden layer at pre-training time,  
single linear layer at fine-tuning time."



# Swin Transformer

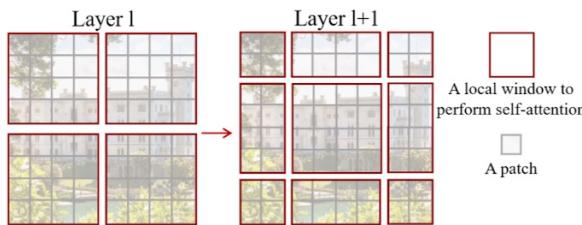
## 1. Concepts

### ① hierarchical feature maps.

Swin transformer model can conveniently leverage advanced techniques for dense prediction, such as feature pyramid networks (FPN) or U-Net.  
Tasks whose predictions are made for every location in the input data, such as images or videos.

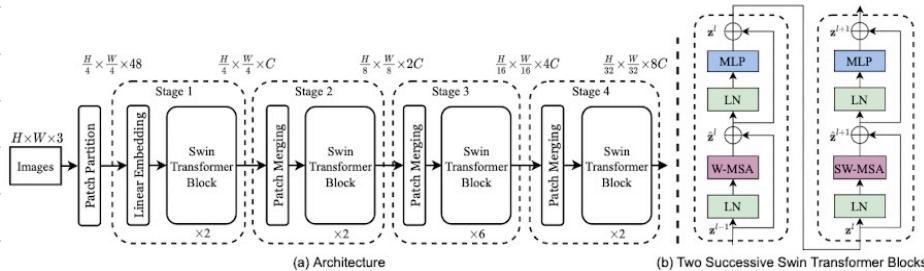
### ② linear computational complexity.

Linear computational complexity is achieved by computing self-attention locally within non-overlapping windows that partition an image.



This is an illustration of the shifted windows. The self-attention computation in the new windows crosses the boundaries of the previous windows in layer  $l$ , providing connections among them.

## 2. Methods



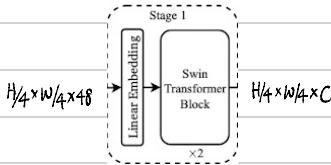
### ① Patch Partition.

- The input RGB image is divided into non-overlapping  $4 \times 4$  patches. (a.k.a. tokens).

$$\text{raw data} \in \mathbb{R}^{H \times W \times 3} \Rightarrow \underset{\substack{\text{patch} \\ \text{partition}}}{\mathbb{R}^{H/4 \times W/4 \times 48}} \Rightarrow \underset{\substack{\text{linear} \\ \text{embedding}}}{\mathbb{R}^{H/4 \times W/4 \times C}}$$

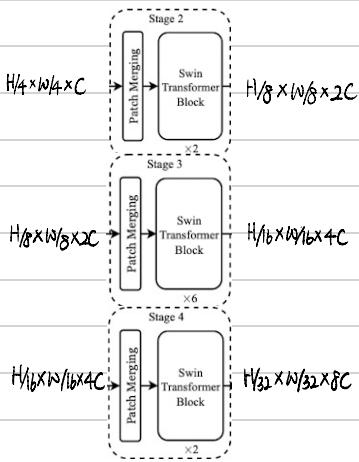
## ② Hierarchical Representation.

### Stage 1.



- Starts with the  $H/4 \times W/4$  token map created from the input image.
- Applies a linear embedding layer to generate the initial feature map.

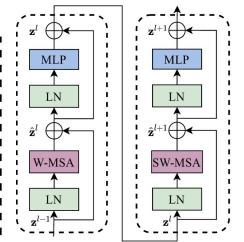
### Stage 2, 3, 4.



In stage 2, 3, and 4, the resolution of the feature map is progressively halved, while the feature dimension is doubled at each stage using patch merging layers.

This hierarchical process allows the model to capture increasingly abstract and high-level representations, with Swin Transformer blocks applied at each stage to refine the features.

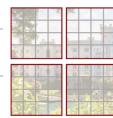
## ③ Swin Transformer Blocks.



Two Successive Swin Transformer Blocks

### 1. Window-based Multi-head Self-Attention (W-MSA)

Computes self-attention over windows to reduce computational complexity.



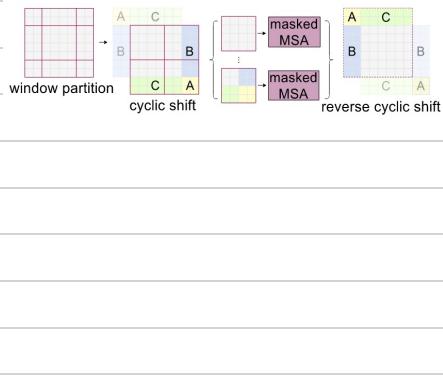
### 2. Shift Window Multi-head Self-Attention (SW-MSA)

Shifts the window partitions between layers to introduce cross-window connections.



## ④ Cyclic Shift

The shift window partitioning increases the number of windows, leading to inefficient computation and memory usage.



### 1. Cyclic Shifting

Windows that are moved out of one side are wrapped around to the opposite side.

### 2. Masking

ensure that self-attention calculations are restricted to the original shifted window regions, maintaining the correct connections

# Stable Diffusion

## 1 Concepts

### The problems of GAN

- mode collapse: When a generative model is trained, it tends to focus on a specific distribution of the data, resulting in a lack of diversity. It repeatedly generates similar images.
- training instabilities.

### Solution by using DMs.

The paper emphasize that Diffusion Models can effectively learn various modes of the data distribution through their mode covering behavior.

{ likelihood-based learning method. (ex. VAEs, AR model, DMs)  
gradual denoising process.  
training stability: DMs doesn't use adversarial loss.

### The problems of DMs.

- Requiring massive computational resources.
- Evaluating a trained model is also

### Solution proposed in the paper: (LDMs)

As with any likelihood-based model, learning can be roughly divided into two stages:

#### 1. Perceptual compression.

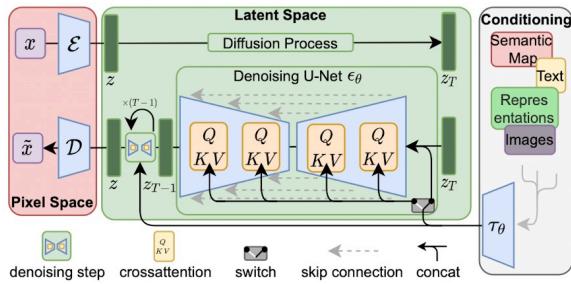
Training an autoencoder which provides a lower-dimensional representational space which is perceptually equivalent to the data space. (aka latent space)

#### 2. Semantic compression.

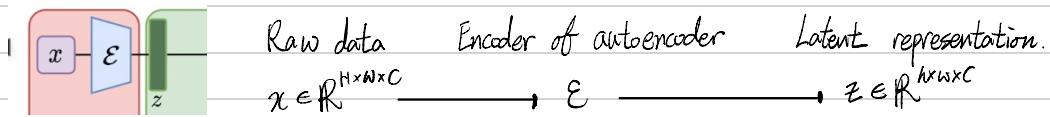
Training DMs in the learned latent space, which exhibits better scaling properties with respect to the spatial dimensionality.

"Moreover, we design a general-purpose conditioning mechanism based on cross-attention, enabling multi-modal training"

## 2. Methods

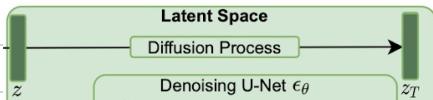


### ① Perceptual Compression.



- The original resolution  $H, W$  is reduced, and the number of channels  $C$  is increased.

### ② Forward Diffusion Process.

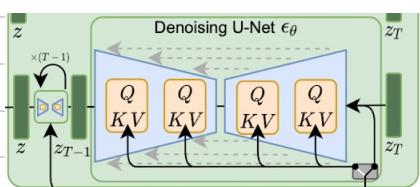


$$z_t = \sqrt{1 - \beta_t} z_{t-1} + \sqrt{\beta_t} \epsilon \quad \text{if } t \text{ is a scheduling parameter controlling the noise level.}$$

- The Forward Diffusion Process progressively adds noise  $\epsilon \sim N(0, 1)$  to the latent representation  $z_0$  producing  $z_T$ .

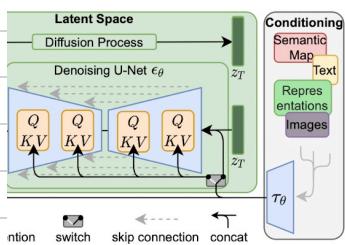
### ③ Reverse Diffusion Process.

#### 1. Using a UNet Architecture.



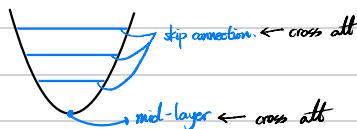
input:  $z_t$  target:  $\epsilon(z_t, t)$

## 2. Cross Attention.



- Conditioning (e.g., Text, Images, ...) is incorporated into UNet via Cross Attention mechanism.

- Cross Attention is applied in the mid-layers and skip connections.

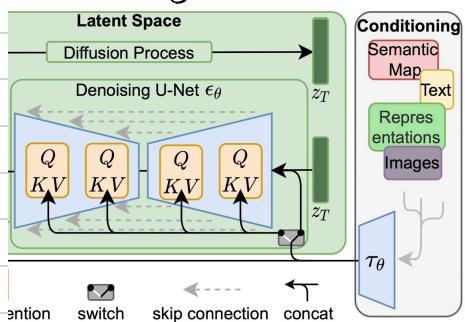


## 3. Loss Function.

- The UNet is trained by minimizing the L2 loss between the predicted noise  $\epsilon_\theta(z_t, t)$  and the ground truth noise  $\epsilon$ :

$$L = E_{\epsilon(x), \epsilon \sim N(0, 1), t} [\|\epsilon - \epsilon_\theta(z_t, t)\|_2^2]$$

## ④ Conditioning



A domain specific encoder  $T_\theta$  that projects  $y$  to an intermediate representation  $T_\theta(y) \in \mathbb{R}^{M \times d_r}$   $M$ : num of input tokens.

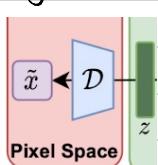
$$Q = W_Q^{(i)} \varphi_i(z_t) \quad K = W_K^{(i)} \cdot T_\theta(y) \quad V = W_V^{(i)} \cdot T_\theta(y)$$

$$W_Q^{(i)} \in \mathbb{R}^{d_{\text{latent}} \times d_r}, \quad W_K^{(i)} \in \mathbb{R}^{d_{\text{latent}} \times d_r}, \quad W_V^{(i)} \in \mathbb{R}^{d_{\text{latent}} \times d_r}$$

$$\Rightarrow \text{softmax}\left(\frac{\partial K^T}{\partial t}\right) \cdot V \in \mathbb{R}^{d_{\text{latent}} \times d_{\text{latent}}}$$

$$\text{Maybe } W_O \in \mathbb{R}^{N \times d_r}$$

## ⑤ Image Generation.



Latent representation. Decoder of autoencoder

$$z \in \mathbb{R}^{h \times w \times C}$$

$$\longrightarrow D$$

Raw data

$$\tilde{x} \in \mathbb{R}^{H \times W \times C}$$

- The Autoencoder's Decoder reconstructs the latent representation  $z_0$  back into the original image space.