

Visiontransformer: AN IMAGE is WORTH 16×16 words:

① Input Image $X \in \mathbb{R}^{h \times w \times c}$ = $\mathbb{R}^{480 \times 480 \times 3}$
height, width, channels

② Patch extraction

The image is divided into patches

$X = \text{Input Image}$
 $\in \mathbb{R}^{480 \times 480 \times 3}$

$(\frac{H}{P})(\frac{W}{P})$ P is the size of the patches
 $P = 16 \times 16$ - Patch Size

$P=16 \times 16$ -Patch Size

In our Example we get

$$\left(\frac{480}{16}\right) * \left(\frac{480}{16}\right) = 30 \times 30 = 900 \text{ patches.}$$

Now we put these 900 patches in a Row &
and each patch has $(\text{Patch size}-x)(\text{Patch size}-y)(\text{channels})$
 $[16 \times 16 \times 3 = 768]$

This gives us patch embeddings.

$\text{Patchembd} \in \mathbb{R}^{900 \times 768}$

Image $X \xrightarrow[\text{patches}]{} \text{patch embedding, flatten}$

IR 480x480x3 → IR 900x768

$R^{H \times W \times 3} \rightarrow IR^{\# \text{of patches} \times (P \times P \times C)}$ $P = \text{patch size}$
 $C = \# \text{of channels}$

N = # of patches when image is divided.

③ Linear projection of flatten spaces.

③ Linear projection of flatten spaces.
Linear projection is applied to the flatten patch embeddings.

A learnable linear projection is applied to patch embeddings $\mathbb{R}^{N \times (P^2 \times C)}$ to produce features $\mathbb{R}^{(P^2 \times C) \times D}$. "D" can be 512 or any dim.

Shape After Operation N

NXD

where $N = \# \text{of patches}$

D = 512 or any dimension -

④ Positional Encoding

We have divided the image into patches and now we are adding positional encoding to these flatten patches that are projected to D dimensions.

Why bother with positional encoding?

We have flatten our image and stacked all the patches in row. Now there is no information preserved about the position of the patch no spatial information is retained when the patches are flatten. We just multiply the $(h' \times w' \times 3)$ here $h' = w'$ = patch size and stack them.

Adding positional encoding to the flatten patches to maintain the position information among the patches.

ternable

The dimension of the position encoding is equal to the flatten projected patch size

$$\in \mathbb{R}^{N \times D}$$

flatten projected patch

$$\mathbb{R}^{N \times D}$$

⊕ position encoding → position encoded embed

$$\mathbb{R}^{N \times D}$$

⑤ Adding CLS-token

Since we are trying to predict the "image class" with transformer we want to add an classification class to the flatten outputs since this will help us make the prediction about the class we use this for our loss and optimizer will update the parameters for this task.

pos-encoded-embedding $\stackrel{\text{Stack}}{\oplus}$ CLS-token → Input to transformer

$$\in \mathbb{R}^{N \times D}$$

$$\in \mathbb{R}^{1 \times D}$$

$$\in \mathbb{R}^{(N+1) \times D} \quad +1 = \text{CLS-token}$$

⑥ Transformer Encoder Layer:

Input to the transformer Encoder layer is [flatten patches · projected to $\mathbb{R}^{N \times D}$ · D dimension + positional encoding] stack Gls-token for class prediction.

$$\mathbb{R}^{N \times (P^2 \times C)}$$

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

→ Input to transformer $\mathbb{R}^{(N+1) \times D}$

→ Linear Projection (Query, key, values) (Q, K, V) (W, W, W)

→ multi head attention Concatenate

→ first (QK^T) → Query into keys heads

→ Second $\left(\frac{QK}{\sqrt{dK}}\right)^T$ → \sqrt{dK} for numerical stability

→ Third Softmax $\left(\frac{QK}{\sqrt{dK}}\right)$ Attention Scores

→ fourth (Softmax $\times \sqrt{dK}$) Projections.

→ fifth (Softmax $\times \sqrt{dK}$) Concatenate

$$\mathbb{R}^{(N+1) \times d}, \mathbb{R}^{(N+1) \times d}, \mathbb{R}^{D \times D}, \mathbb{R}^{D \times D}, \mathbb{R}^{D \times D}$$

$$\text{Softmax} \left(\frac{QK^T}{\sqrt{dK}} \right) * V$$

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

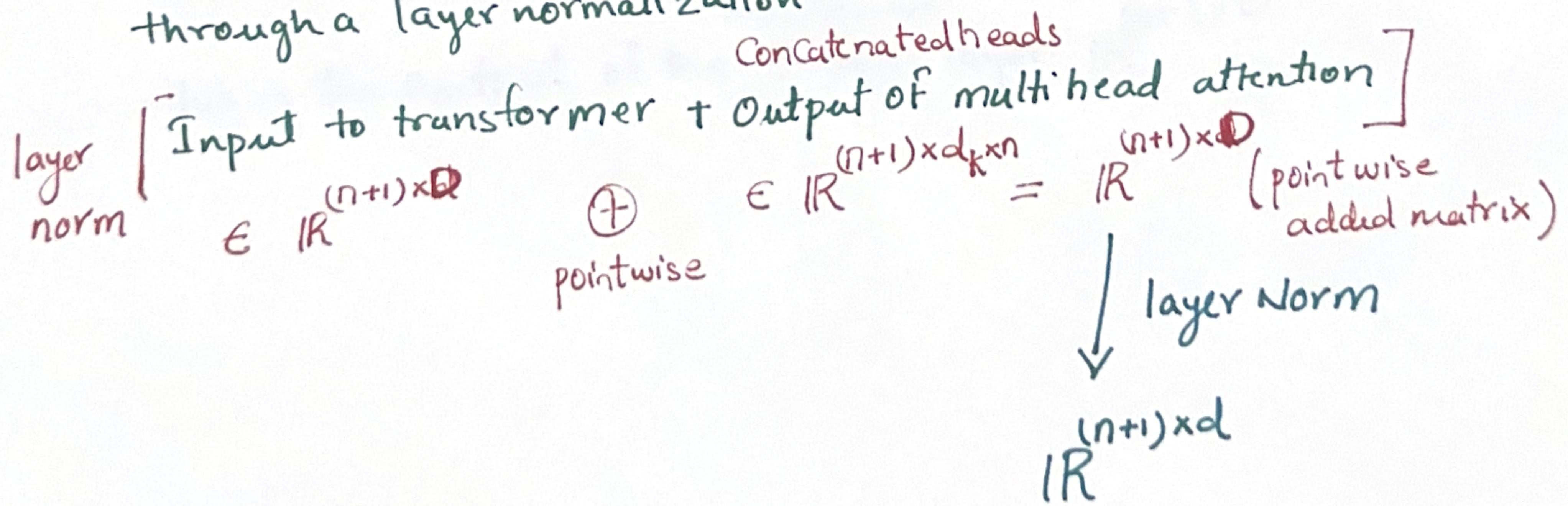
$$\mathbb{R}^{(N+1) \times (N+1)} \times \mathbb{R}^{(N+1) \times d}$$

$$dK = \frac{D}{n}$$

$n = \text{no of heads in mHA}$

⑦ Residual Connections & LayerNorm

The output of the attention/mHA is added back to the input and this residual connection is then passed through a layer normalization



⑧ Feed forward network

1) first linear layer (Expand) $\mathbb{R}^{(n+1) \times D}$ $\mathbb{R}^{D \times \text{Expand}}$
 $[\text{layer Norm}] \times [W_{\text{first linear projection}}] = \mathbb{R}^{(n+1) \times \text{Exp}}$
 Relu Activation (Activation function) $\text{Relu}[\text{first}] \text{fmax}(0, \text{value}) = \mathbb{R}^{(n+1) \times \text{exp}}$

Second linear layer (Contract Back) $[\text{Relu}] \times [w_{\text{second linear projection}}]$
 \downarrow
 $\mathbb{R}^{(n+1) \times \text{Expand}} \times \mathbb{R}^{\text{Expand} \times D}$
 $\mathbb{R}^{(n+1) \times D}$

Residual Connection & layer Normalization

Apply layer Norm $(\text{FNN_second linear layer } \in \mathbb{R}^{(n+1) \times D}) + (\text{Normalized outputs from layer Norm Block Before fNN}) \in \mathbb{R}^{(n+1) \times D}$

The output of fNN is added to the normalized output from the previous layer pointwise

After the point wise addition we apply layernorm to the ~~the~~ pointwise sum.

The dimension after the layernorm to the pointwise sum remains the same $\mathbb{R}^{(n+1) \times D}$

This is the output of the Transformer Encoder layer. $\mathbb{R}^{(n+1) \times D}$

⑨ Classification:

Now from the output of transformer_encoder_output
 $\in \mathbb{R}^{(n+1) \times D}$

we extract the cls-token vector

Since this cls-token vector was a trainable parameter during training we take this vector

$$\mathbb{R}^{1 \times D} = \text{cls-token (vector)}$$

Then we pass this through a linear-layer-mlp that is of dimension of $\mathbb{R}^{D \times \text{number of classes}}$.

$$\begin{array}{ccc} \text{cls-token} & \times & \text{weight_for_last_mlp} \\ \mathbb{R}^{1 \times D} & \times & \mathbb{R}^{D \times \text{no of classes}} \end{array} = \mathbb{R}^{1 \times \#\text{of classes}}$$