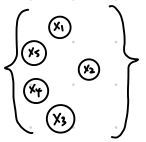


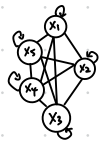
We can think about a **Set** in 2 different ways.

(1)



Graph without any edges

(2)



Graph that is fully connected

Message Passing  
: Allow every node talk to every other node  
: This is the core principle behind transformer networks

### Quick Recap

- Within linear algebra, each permutation defines a  $|V| \times |V|$  matrix.
- Such matrices are called permutation matrices.
- They have exactly one 1 in every row and column, zeroes elsewhere

Ex)

$$P_{(2,4,1,3)} X = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} -x_1 \\ -x_2 \\ -x_3 \\ -x_4 \end{bmatrix} = \begin{bmatrix} -x_2 \\ -x_4 \\ -x_1 \\ -x_3 \end{bmatrix}$$

### Learning on Sets

#### ① Setup

For now, assume our graph has no edges (i.e.,  $E = \emptyset$ , the set of nodes)

Let  $x_i \in \mathbb{R}^k$  be the features of node  $i$ . (i.e., our feature space is  $C = \mathbb{R}^k$ )

We can stack these features into a node feature matrix of shape  $N \times k$ .  $X = (x_1, \dots, x_N)^T$

$i^{th}$  row of  $X$  corresponds to  $x_i$

**Permutation invariant operator** is an operator ( $F$ ), that if we apply  $F$  to our set or a permuted version of a set, we're getting the same output.

$$f\left(\begin{pmatrix} x_1 \\ x_5 \\ x_4 \\ x_3 \\ x_2 \end{pmatrix}\right) = y = f\left(\begin{pmatrix} x_5 \\ x_3 \\ x_1 \\ x_2 \\ x_4 \end{pmatrix}\right)$$

Symmetry group  $G$ :  $n$ -element permutation group  $S_n$

Group element  $g \in G$ : permutation

### Permutation invariance

Want: function  $f(X)$  over sets that will not depend on the order

Equivalently: Applying a permutation matrix shouldn't modify result

$f(X)$  is permutation invariant if, for all permutation matrices  $P$ :  $f(PX) = f(X)$

### Permutation Equivariance:

• Permutation invariant models are good for **set-level outputs**. If we would like to **answer at the node level**, we need **permutation equivariant models**.

• Permutation equivariant functions:  $F(PX) = PF(X)$

↳ It doesn't matter if we do permute before  $F$  or later  $F$ .

## Deep Sets

$$f(X) = \phi\left(\bigoplus_{i \in V} \psi(x_i)\right)$$

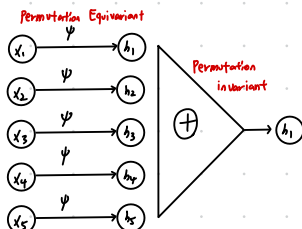
where  $\phi, \psi$  are (learnable) functions, e.g. MLPs

$\bigoplus$  denotes any permutation-invariant operator.

## General blueprint for learning on sets

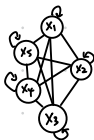
One way we can enforce locality in equivariant set functions is through a shared function  $\psi$  applied to every node in isolation:  $h_i = \psi(x_i)$ , and stack  $h_i$  into a matrix  $H = F(X)$

### Pictorial View



Note:  $\psi$ : a shared function to every node in isolation  
And, there's no message passing

## Now, (2): Sets as fully connected graphs



$$\text{Message Passing: } f(x_i) = \phi\left(x_i, \bigoplus_{j \in N_i} \psi(x_i, x_j)\right)$$

$\psi(x_i, x_j)$ : Message Passing between  $x_i$  and  $x_j$

$\bigoplus_{j \in N_i}$ : Message Aggregation

$\phi(x_i, z)$ : Node feature update

## Basic Self-Attention

Input: Sequence of tensors  $x_1, x_2, \dots, x_t$

Output: Sequence of tensors, each one a weighted sum of the input sequence:  $y_1, y_2, \dots, y_t$

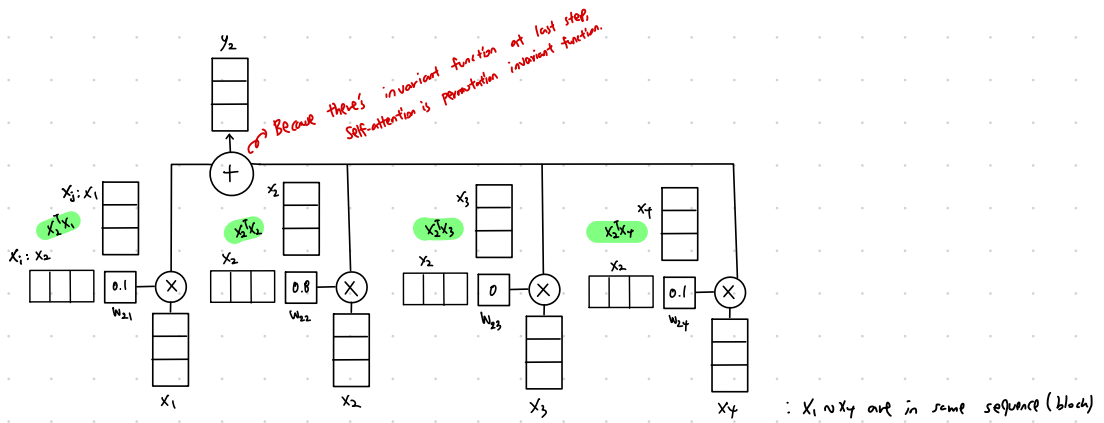
$$y_i = \sum_j w_{ij} x_j$$

$w$  in self-attention is not a learned weight, but a function of  $x_i$  and  $x_j$ :  $w'_{ij} = x_i^T x_j$

$w$  must sum to 1 over  $j$

$$w_{ij} = \frac{\exp w'_{ij}}{\sum_j \exp w'_{ij}} : \text{Softmax (we're applying softmax to the nodes in some sequence)}$$

# Pictorial View of basic self attention



$$y_2 = \sum_j w_{2j} x_j = w_{21}x_1 + w_{22}x_2 + w_{23}x_3 + w_{24}x_4$$

$$= \frac{\exp(x_2^T x_1)}{\sum_j \exp(x_2^T x_j)} x_1 + \frac{\exp(x_2^T x_2)}{\sum_j \exp(x_2^T x_j)} x_2 + \frac{\exp(x_2^T x_3)}{\sum_j \exp(x_2^T x_j)} x_3 + \frac{\exp(x_2^T x_4)}{\sum_j \exp(x_2^T x_j)} x_4 = 0.1x_1 + 0.8x_2 + 0 + 0.1x_4$$

★  $x_i^T x_j$ : Notion of similarity between  $x_i$  and  $x_j$

Therefore,  $w_{2j}$  tells how important  $x_j$  is to  $x_2$ . And  $y_2$  is computed based on its relationship with other nodes. This is "Attention"!

## Basic Self Attention

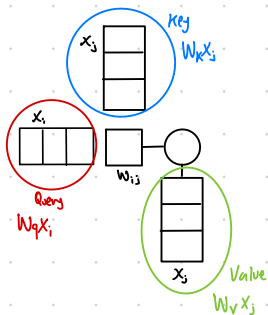
No learned weights

Order of sequence doesn't affect result of computations

Query, Key, Value

Every input vector  $x_i$  is used in 3 ways:

- Query: Compared to every other vector to compute attention weights for its own output  $y_i$ ;
- Key: Compared to every other vector to compute attention weight  $w_{ij}$  for output  $y_j$ ;
- Value: Summed with other vectors to form the result of the attention weighted sum

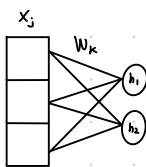
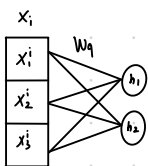


We can process each input vector to fulfill the three roles with matrix multiplication

Learning the matrices  $\rightarrow$  learning attention

$$\begin{aligned} q_i &= W_q x_i & w'_{ij} &= q_i^T k_j \\ k_i &= W_k x_i & w_{ij} &= \text{Softmax}(w'_{ij}) \\ v_i &= W_v x_i & y_i &= \sum_j w_{ij} v_j \end{aligned}$$

$W_q, W_k, W_v$  : Trainable matrices, such as MLP



And,  $w'_{ij} = q_i^T k_j$

$w_{ij} = \text{Softmax}(w'_{ij})$  : Scalar value.

$y_i = \sum_j w_{ij} v_j$

$$W_q = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix}$$

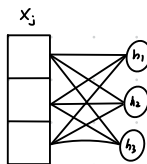
$$W_k = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix}$$

$$\begin{aligned} q_i &= W_q x_i \\ &= \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix} \begin{bmatrix} x_1^i \\ x_2^i \\ x_3^i \end{bmatrix} \\ &= \begin{bmatrix} h_1 \\ h_2 \end{bmatrix}_q \end{aligned}$$

$2 \times 3 \quad 3 \times 1$

$$\begin{aligned} k_j &= W_k x_j \\ &= \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix} \begin{bmatrix} x_1^j \\ x_2^j \\ x_3^j \end{bmatrix} \\ &= \begin{bmatrix} h_1 \\ h_2 \end{bmatrix}_k \end{aligned}$$

$2 \times 1$



$$W_v = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \end{bmatrix}$$

$v_j = W_v x_j$

$$= \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \end{bmatrix} \begin{bmatrix} x_1^j \\ x_2^j \\ x_3^j \end{bmatrix} : 3 \times 1$$

$y_i = \sum_j w_{ij} v_j$

$\star W_q$  and  $W_k$  must share same size, but not  $W_v$ .

$X : (B, T, L)$

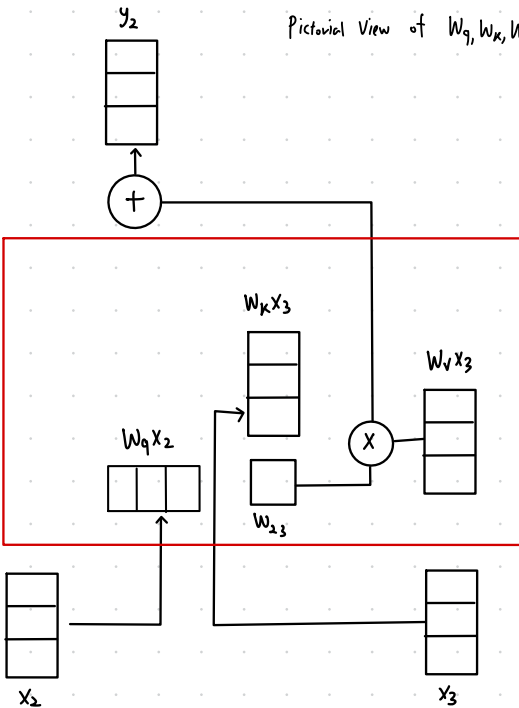
$\text{SelfAttn}(c, h) \rightarrow B, T, H$

$\text{SelfAttn}(c, h) \rightarrow B, T, H$

$W = Q \otimes K^T$   
 $\rightarrow \text{Tensor Mat}$   
 $: B \times L^2$

$\left( \begin{bmatrix} \text{---} \\ \text{---} \\ \vdots \end{bmatrix}, \begin{bmatrix} \text{---} \end{bmatrix} \right)^T$

Pictorial View of  $W_q, W_k, W_v$



### Multi-head attention

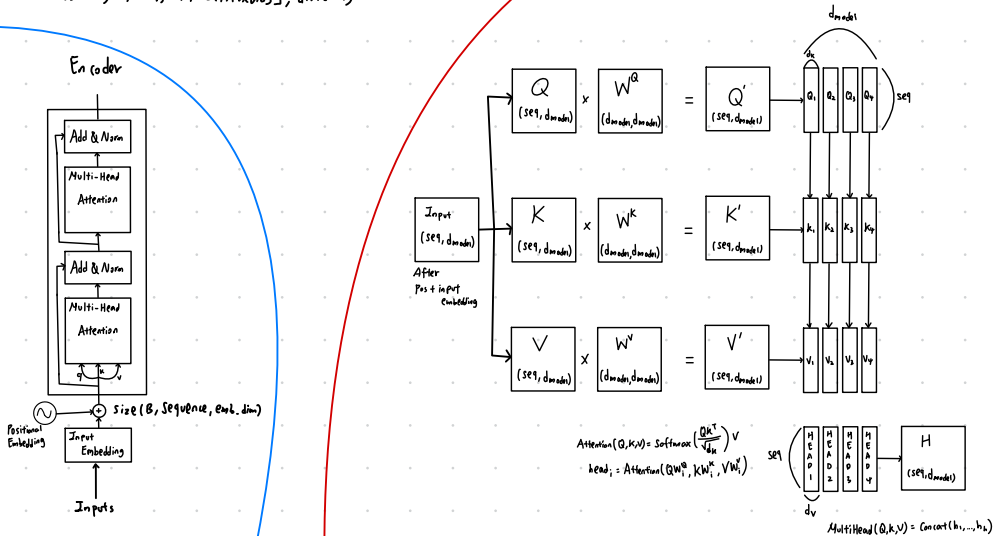
- Multiple "heads" of attention just means learning different sets of  $W_q, W_k$ , and  $W_v$  matrices simultaneously.
- Implemented as just a single matrix

Code: `nn.ModuleList([Head(...) for _ in range(num_heads)])`

And when you return you need concatenate them.

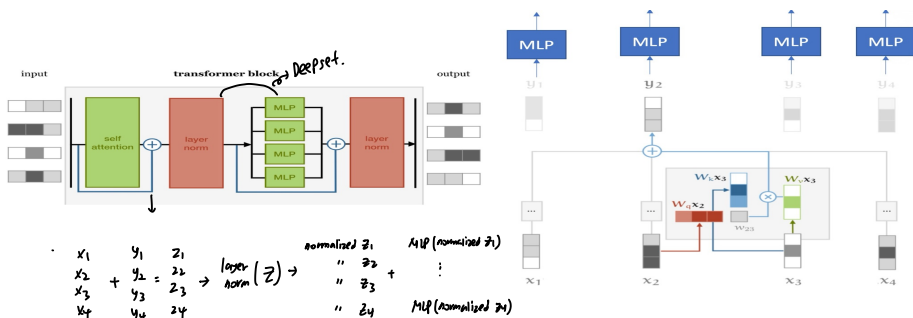
But = `torch.cat([h(x) for h in self.heads], dim=-1)`

Pictorial View of MultiHead



# Transformer

Self-attention layer  $\rightarrow$  Layer normalization  $\rightarrow$  Dense Layer



However, there are problems that order of sequence affect the result of computation, such as MLP.

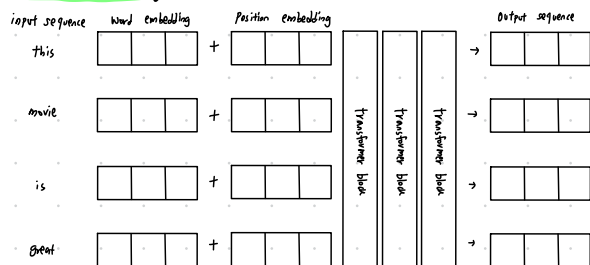
: Let's encode each vector with position

Text, Signals, and Images are not sets.

: There is an inherent ordering on most domains we deal with

Yoda is a Jedi master!  $\neq$  a Jedi master Yoda is!

## Position embedding



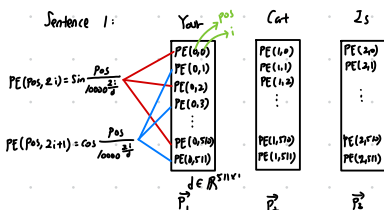
In "Attention is All you need" paper,

$$\vec{P}_k^{(i)} = H(k)^{(i)} := \begin{cases} \sin(W_k \cdot t) & \text{if } i = 2k \\ \cos(W_k \cdot t) & \text{if } i = 2k+1 \end{cases}$$

$$\text{where } W_k = \frac{1}{10000^{\frac{2k}{d}}}$$

$$\vec{P}_k = \begin{bmatrix} \sin(W_k \cdot t) \\ \cos(W_k \cdot t) \\ \vdots \\ \sin(W_k \cdot t) \\ \cos(W_k \cdot t) \end{bmatrix} \text{ dxl where } d \text{ be the encoding dimension } (d \approx 0) \text{ (embedding)}$$

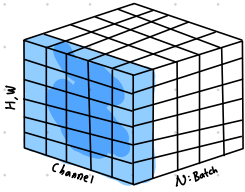
Example)



Why trigonometric functions?

: Trigonometric functions like sin and cos naturally represent a pattern that the model can recognize as continuous.

## Layer Normalization



### Layer Normalization

Layer Norm: Calculate  $\mu, \sigma^2$  of Channel per each batch.

Example of layer norm:

Batches of 3 items:

Item ①

50.149
3314.8
...
...
8941.2
1444.7

$\mu_1$   
 $\sigma_1^2$

Item ②

1990.2
688.3
...
...
27.4
94.18

$\mu_2$   
 $\sigma_2^2$

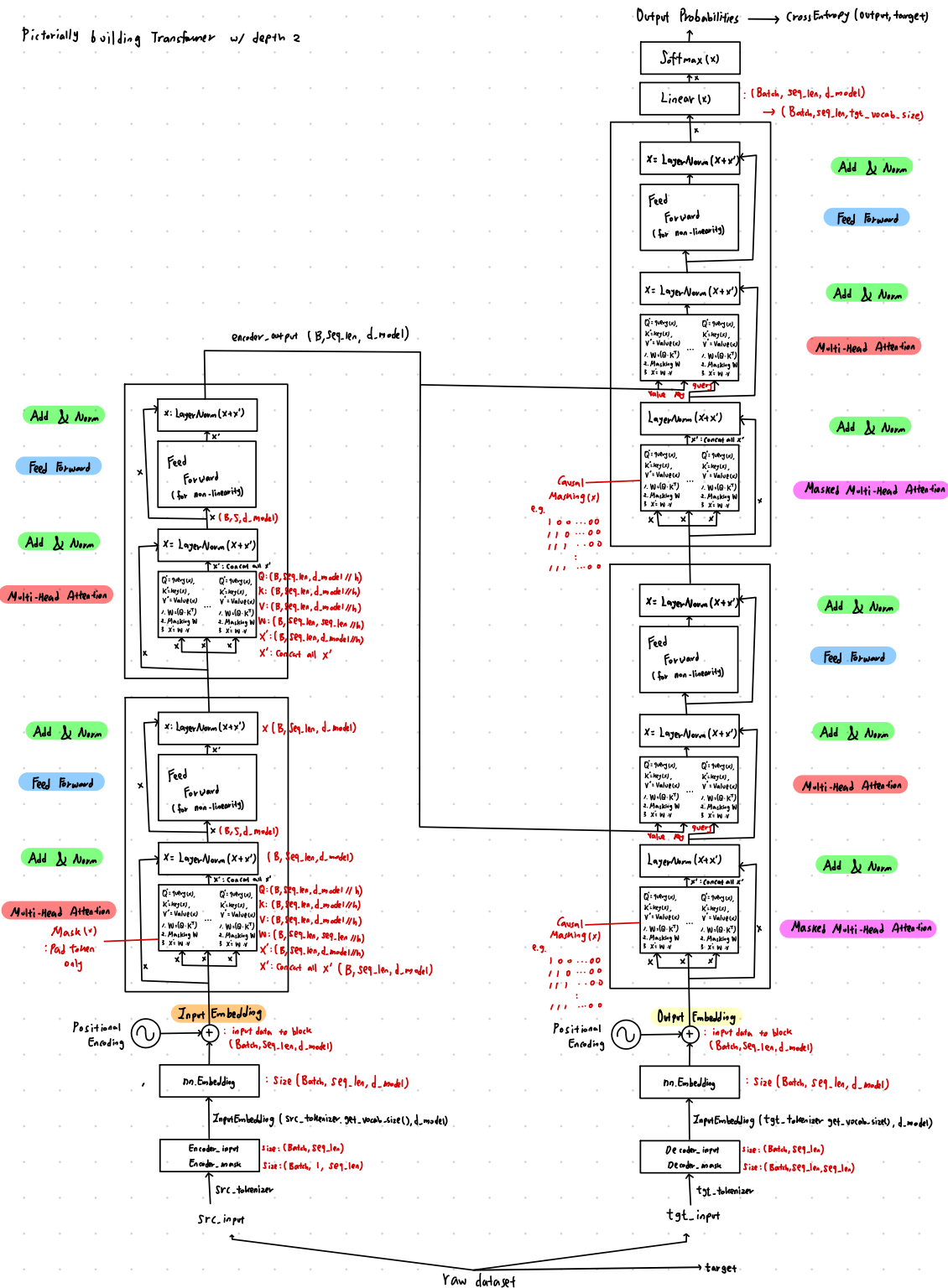
Item ③

182.7
174.878
...
...
10923.7
1004.88

$\mu_3$   
 $\sigma_3^2$

$$\Rightarrow \hat{x}_j = \frac{x_j - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

Pictorially building Transformer w/ depth 2





## Vision Transformer (ViT)

In practice: take 224x224 input image,  
divide into 14x14 grid of 16x16 pixel  
patches (or 16x16 grid of 14x14 patches)

Each attention matrix has  $14^4 = 38,416$   
entries, takes 150 KB  
(or 65,536 entries, takes 256 KB)

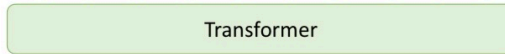
Output vectors

Exact same as  
NLP Transformer!

Add positional  
embedding: learned D-  
dim vector per position

Linear projection to  
D-dimensional vector

N input patches, each  
of shape 3x16x16



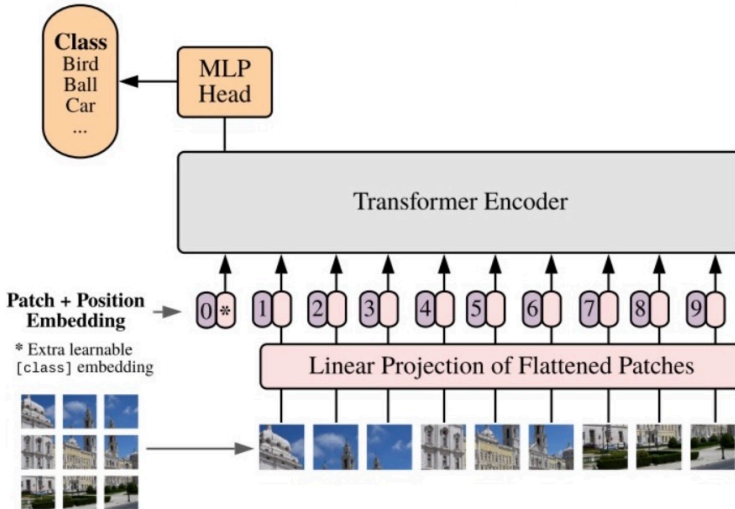
Linear projection  
to C-dim vector  
of predicted  
class scores

Special extra input:  
**classification token**  
(D dims, learned)

Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

[Cat image](#) is free for commercial  
use under a [Pixabay license](#)

## Vision Transformer (ViT)



## Transformer Encoder

