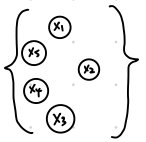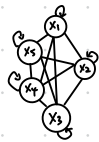We can think about a **set** in 2 different ways.

(1)



Graph without any edges

(2)



Message Passing
: Allow every node talk
  to every other node

: This is the core Principle behind transformer networks

Graph that is fully connected

**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., $\Omega = V$, 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 $N \times k$.  $X = (X_1, ..., X_N)^T$

· $i^{th}$ 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.

$$f\left( \begin{array}{c} X_1 \\ X_5 \\ X_4 \\ X_3 \end{array} \right) = y = f\left( \begin{array}{c} X_5 \\ X_1 \\ X_4 \\ X_3 \end{array} \right)$$

Symmetry group  $G$: n-element Permutation group $\Sigma_n$

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 : $F(PX) = P F(X)$

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

# Deep Sets

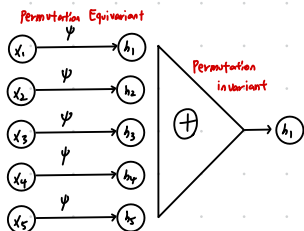$$f(X) = \phi\left(\bigoplus_{i \in V} \psi(x_i)\right)$$

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

**Permutation Equivariant**



**Permutation invariant**

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

---

Now, (2) · Sets as fully connected graphs



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

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

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

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

---

# Basic Self-Attention

· Input: Sequence of tensors $x_1, x_2, ..., x_t$

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

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

· $W$ 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 $j$

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

# Pictorial View of basic self attention



$$y_2 = \sum_j W_{2j} X_j = W_{21}X_1 + W_{22}X_2 + W_{23}X_3 + W_{24}X_4$$

$$= \frac{exp(x_2^T x_1)}{\sum_j exp(x_2^T x_j)} X_1 + \frac{exp(x_2^T x_2)}{\sum_j exp(x_2^T x_j)} X_2 + \frac{exp(x_2^T x_3)}{\sum_j exp(x_2^T x_j)} X_3 + \frac{exp(x_2^T x_4)}{\sum_j exp(x_2^T x_j)} X_4 = 0.1 X_1 + 0.8 X_2 + 0 + 0.1 X_4$$

⭐ $X_i^T X_j$ : 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 <mark>to compute attention weights</mark> for its own output $y_i$
- Key: Compared to every other vector to <mark>compute attention weight $w_{ij}$</mark> for output $y_j$
- Value: Summed with other vectors <mark>to form the result of the attention weighted sum</mark>



Key
$W_k x_j$

$x_i$

Query
$W_q x_i$

$w_{ij}$

$x_j$  Value
$W_v x_j$

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

Learning the matrices → learning attention

$$q_i = W_q x_i \qquad W'_{ij} = q_i^T k_j$$
$$K_i = W_k x_i \qquad w_{ij} = Softmax(W'_{ij})$$
$$v_i = W_v x_i \qquad y_i = \sum_j w_{ij} v_j$$
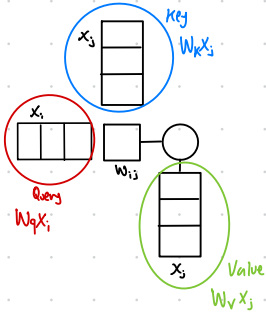
$W_q, W_k, W_v$ : Trainable matrices, such as MLP

$x_i$



$x_j$



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

$w_{ij} = Softmax(W'_{ij})$ : Scalar value.

$y_i = \sum_j w_{ij} v_j$

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

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

$x_j$



$q_i = W_q x_i$

$$= \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix} \begin{bmatrix} x_1^i \\ x_2^i \\ x_3^i \end{bmatrix}$$

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

$$= \begin{bmatrix} h_1 \\ h_2 \end{bmatrix}_q$$

$\quad 2\times 1$

$K_i = W_k x_i$

$$= \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix} \begin{bmatrix} x_1^i \\ x_2^i \\ x_3^i \end{bmatrix}$$

$$= \begin{bmatrix} h_1 \\ h_2 \end{bmatrix}_k$$

$\quad 2\times 1$

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

$v_j = W_k x_j$

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

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

$y_i = \sum_j w_{ij} v_j$

$X : (B, T, L)$

self.key $(L, h) \rightarrow B, T, H$

self.query $(L, h) \rightarrow B, T, H$

$W = Q \otimes h^T \quad [[[ \overbrace{\phantom{xxx}} ], \; [ \overbrace{\phantom{xxx}} ], \; [ \vdots ]]]^T$

$B, T\times H \quad \alpha, H\times T$

$: B \times T\times T$

$y_2$

Pictorial View of $W_q, W_k, W_v$

$W_K X_3$

$W_V X_3$

$W_q X_2$

$X$

$W_{23}$

$X_2$

$X_3$

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.

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

Encoder

Add & Norm

Multi-Head Attention

Add & Norm

Multi-Head Attention

Size (B, Sequence, emb_dim)

Positional Embedding

Input Embedding

Inputs

Pictorial View of MultiHead

$d_{model}$

$Q$ $(seq, d_{model})$  x  $W^{iQ}$ $(d_{model}, d_{model})$  =  $Q'$ $(seq, d_{model})$   $Q_1 | Q_2 | Q_3 | Q_4$  seq

Input $(seq, d_{model})$

After Pos + input Embedding

$K$ $(seq, d_{model})$  x  $W^k$ $(d_{model}, d_{model})$  =  $K'$ $(seq, d_{model})$   $k_1 | k_2 | k_3 | k_4$

$V$ $(seq, d_{model})$  x  $W^V$ $(d_{model}, d_{model})$  =  $V'$ $(seq, d_{model})$   $v_1 | v_2 | v_3 | v_4$

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

head$_i = $ Attention$(QW_i^Q, KW_i^k, VW_i^v)$

$\begin{pmatrix} H \\ E \\ A \\ D \\ 1 \end{pmatrix} \begin{pmatrix} H \\ E \\ A \\ D \\ 2 \end{pmatrix} \begin{pmatrix} H \\ E \\ A \\ D \\ 3 \end{pmatrix} \begin{pmatrix} H \\ E \\ A \\ D \\ 4 \end{pmatrix}$   $H$ $(seq, d_{model})$

seq

$d_v$

MultiHead$(Q, K, V) = $ Concat$(h_1, ..., h_h)$

# Transformer

· Self-attention layer → Layer normalization + Dense Layer



$$
\begin{matrix}
x_1 & y_1 & z_1 \\
x_2 & y_2 & z_2 \\
x_3 & y_3 & z_3 \\
x_4 & y_4 & z_4
\end{matrix}
\quad
\begin{matrix}
& & \\
+ & & = \\
& & \\
& &
\end{matrix}
\quad
\overset{\text{layer}}{\underset{\text{norm}}{}}(\vec{z}) \rightarrow
\begin{matrix}
\text{normalized } z_1 \\
\text{'' } z_2 \\
\text{'' } z_3 \\
\text{'' } z_4
\end{matrix}
\quad +
\quad
\begin{matrix}
\text{MLP (normalized } z_1) \\
\vdots \\
\vdots \\
\text{MLP (normalized } z_4)
\end{matrix}
$$

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! ≠ a Jedi master Yoda is!

## Position embedding



input sequence   Word embedding   Position embedding        Output sequence
this
movie
is
great

In "Attention is All you need" Paper,

$$
\vec{P_t}^{(i)} = f(t)^{(i)} := \begin{cases} \sin(w_k \cdot t) & \text{if } i=2k \\ \cos(w_k \cdot t) & \text{if } i=2k+1 \end{cases}
$$

$$
\vec{P_t} = \begin{bmatrix} \sin(w_1 \cdot t) \\ \cos(w_1 \cdot t) \\ \vdots \\ \sin(w_{d/2} \cdot t) \\ \cos(w_{d/2} \cdot t) \end{bmatrix}_{d \times 1}
$$

where $d$ be the encoding dimension $(d \overset{=}{>} 0)$
(embedding)

where $w_k = \dfrac{1}{10000^{2k/d}}$

Example)

Sentence 1:

$$PE(Pos, 2i) = \sin \frac{Pos}{10000^{\frac{2i}{d}}}$$

$$PE(Pos, 2i+1) = \cos \frac{Pos}{10000^{\frac{2i}{d}}}$$



| Your | Cat | Is |
|---|---|---|
| PE(0,0) | PE(1,0) | PE(2,0) |
| PE(0,1) | PE(1,1) | PE(2,1) |
| PE(0,2) | PE(1,2) | ⋮ |
| PE(0,3) | ⋮ | |
| ⋮ | | |
| PE(0,510) | PE(1,510) | PE(2,510) |
| PE(0,511) | PE(1,511) | PE(2,511) |

$d \in \mathbb{R}^{512 \times 1}$

$\vec{P_1}$    $\vec{P_2}$    $\vec{P_2}$

# Why trigonometric functions?

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

# Layer Normalization



H, W

Channel

N: Batch

Layer Normalization

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

Example of layer norm:

Batches of 3 items:

Item ①

| |
|---|
| 50.147 |
| 3314.8 |
| ... |
| |
| 8941.2 |
| 1994.7 |

$\mu_1$
$\sigma_1^2$

Item ②

| |
|---|
| 1990.2 |
| 688.3 |
| ... |
| ... |
| 27.4 |
| 94.11 |

$\mu_2$
$\sigma_2^2$

Item ③

| |
|---|
| 183.7 |
| 174.878 |
| ... |
| ... |
| 10923.7 |
| 1004.88 |

$\mu^3$
$\sigma_3^2$

$\Rightarrow \quad \hat{x}_j = \dfrac{x_j - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$

Pictorially building Transformer w/ depth 2

Output Probabilities ——→ CrossEntropy (output, target)

Softmax (x)

Linear (x) : (Batch, seq_len, d_model)
——→ (Batch, seq_len, tgt_vocab_size)

X = LayerNorm (X+x')

Feed
Forward
(for non-linearity)

Add & Norm

Feed Forward

X = LayerNorm (X+x')

Q':=key(x), Q':=key(x)
K':=key(x), K':=key(x)
V':=Value(x), V':=Value(x)
1. W=(Q·K^T), 1. W=(Q·K^T)
2. Masking W, 2. Masking W
3. X:= W·V, 3. X:= W·V

Add & Norm

Multi-Head Attention

Value  Key     Query

LayerNorm (X+x')
X': Concat all x'

Q':=key(x), Q':=key(x)
K':=key(x), K':=key(x)
V':=Value(x), V':=Value(x)
1. W=(Q·K^T), 1. W=(Q·K^T)
2. Masking W, 2. Masking W
3. X:= W·V, 3. X:= W·V

Causal
Masking(x)
e.g.
1 0 0 ...0 0
1 1 0 ...0 0
1 1 1 ...0 0
1 1 1 ...0 0

Add & Norm

Masked Multi-Head Attention

encoder_output (B, seq_len, d_model)

X = LayerNorm (X+x')
x'

Feed
Forward
(for non-linearity)

x' X (B,S,d_model)

X = LayerNorm (X+x')
X': Concat all x'

Add & Norm

Feed Forward

Q':=key(x), Q':=key(x)
K':=key(x), K':=key(x)
V':=Value(x), V':=Value(x)
1. W=(Q·K^T), 1. W=(Q·K^T)
2. Masking W, 2. Masking W
3. X:= W·V, 3. X:= W·V

Q: (B, seq_len, d_model // h)
K: (B, seq_len, d_model//h)
V: (B, seq_len, d_model//h)
W: (B, seq_len, seq_len //h)
X': (B, seq_len, d_model//h)
X': Concat all X'

Add & Norm

Multi-Head Attention

X = LayerNorm (X+x')   X (B, Seq_len, d_model)
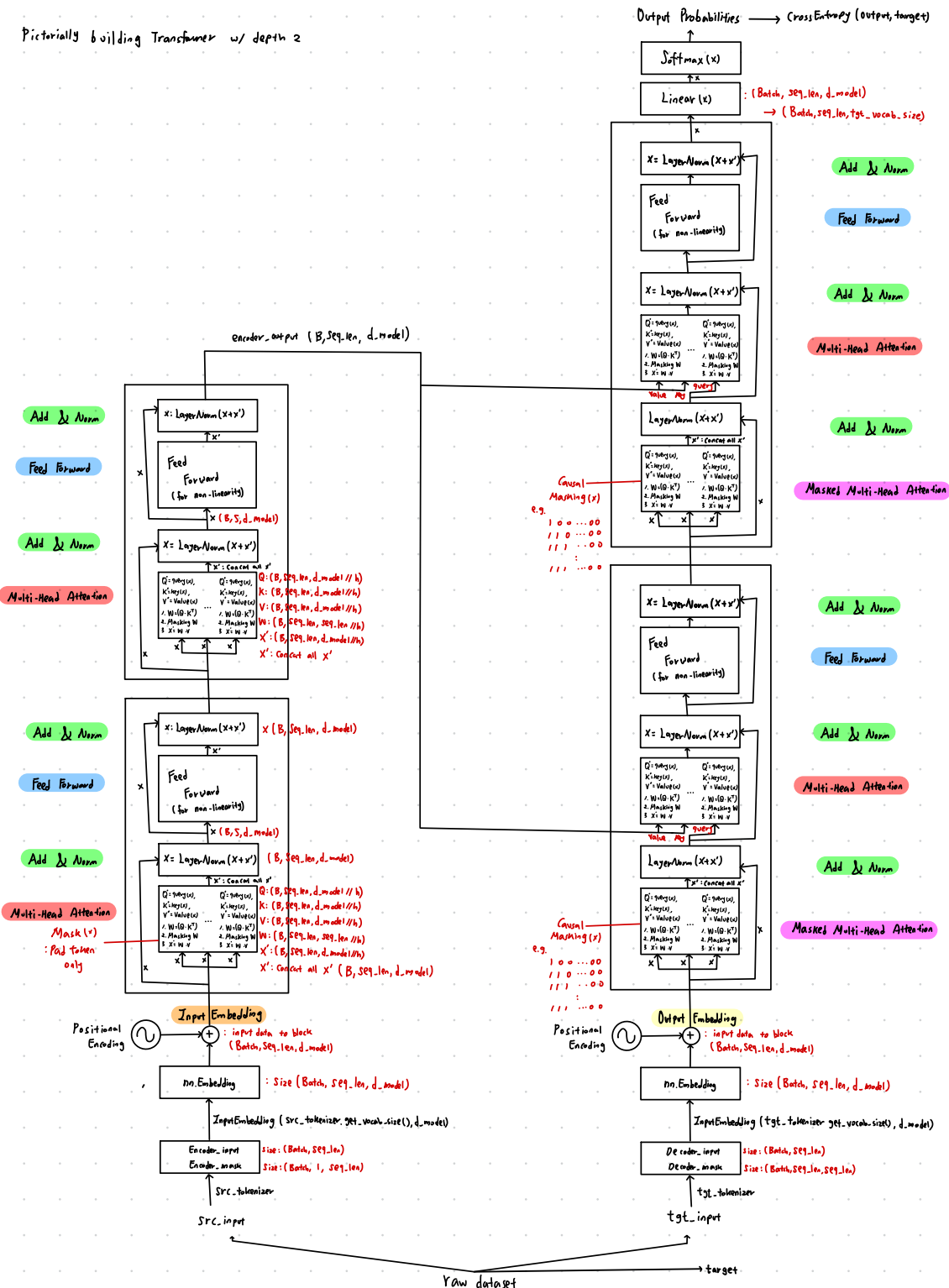x'

Feed
Forward
(for non-linearity)

x' X (B,S,d_model)

X = LayerNorm (X+x')   (B, Seq_len, d_model)
X': Concat all x'

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention
Mask(v):
: Pad token
only

Q':=key(x), Q':=key(x)
K':=key(x), K':=key(x)
V':=Value(x), V':=Value(x)
1. W=(Q·K^T), 1. W=(Q·K^T)
2. Masking W, 2. Masking W
3. X:= W·V, 3. X:= W·V

X': Concat all X' (B, seq_len, d_model)

X = LayerNorm (X+x')

Feed
Forward
(for non-linearity)

X = LayerNorm (X+x')

Add & Norm

Feed Forward

Add & Norm

Q':=key(x), Q':=key(x)
K':=key(x), K':=key(x)
V':=Value(x), V':=Value(x)
1. W=(Q·K^T), 1. W=(Q·K^T)
2. Masking W, 2. Masking W
3. X:= W·V, 3. X:= W·V

Multi-Head Attention

Value  Key     Query

LayerNorm (X+x')
X': Concat all x'

Q':=key(x), Q':=key(x)
K':=key(x), K':=key(x)
V':=Value(x), V':=Value(x)
1. W=(Q·K^T), 1. W=(Q·K^T)
2. Masking W, 2. Masking W
3. X:= W·V, 3. X:= W·V

Causal
Masking(x)
e.g.
1 0 0 ...0 0
1 1 0 ...0 0
1 1 1 ...0 0
1 1 1 ...0 0

Add & Norm

Masked Multi-Head Attention

Positional
Encoding  ∿  ⊕  Input Embedding
: input data to block
(Batch, seq_len, d_model)

Positional
Encoding  ∿  ⊕  Output Embedding
: input data to block
(Batch, seq_len, d_model)

nn.Embedding  : Size (Batch, seq_len, d_model)

InputEmbedding ( src_tokenizer. get_vocab_size(), d_model)

nn.Embedding  : Size (Batch, seq_len, d_model)

InputEmbedding (tgt_tokenizer. get_vocab_size(), d_model)

Encoder_input   size:(Batch, seq_len)
Encoder_mask   size:(Batch, 1, seq_len)

src_tokenizer

src_input

Decoder_input   size:(Batch, seq_len)
Decoder_mask   size:(Batch,seq_len,seq_len)

tgt_tokenizer

tgt_input

Raw dataset  ——→ target

# 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)

Linear projection to C-dim vector of predicted class scores

Output vectors

Transformer

Exact same as NLP Transformer!

Add positional embedding: learned D-dim vector per position

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

Linear projection to D-dimensional vector

N input patches, each of shape 3x16x16

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**

Class
Bird
Ball
Car
...

MLP
Head

Transformer Encoder

L x

MLP

Norm

Multi-Head
Attention

Norm

**Patch + Position Embedding**

* Extra learnable [class] embedding

0*   1   2   3   4   5   6   7   8   9

Linear Projection of Flattened Patches

Embedded Patches