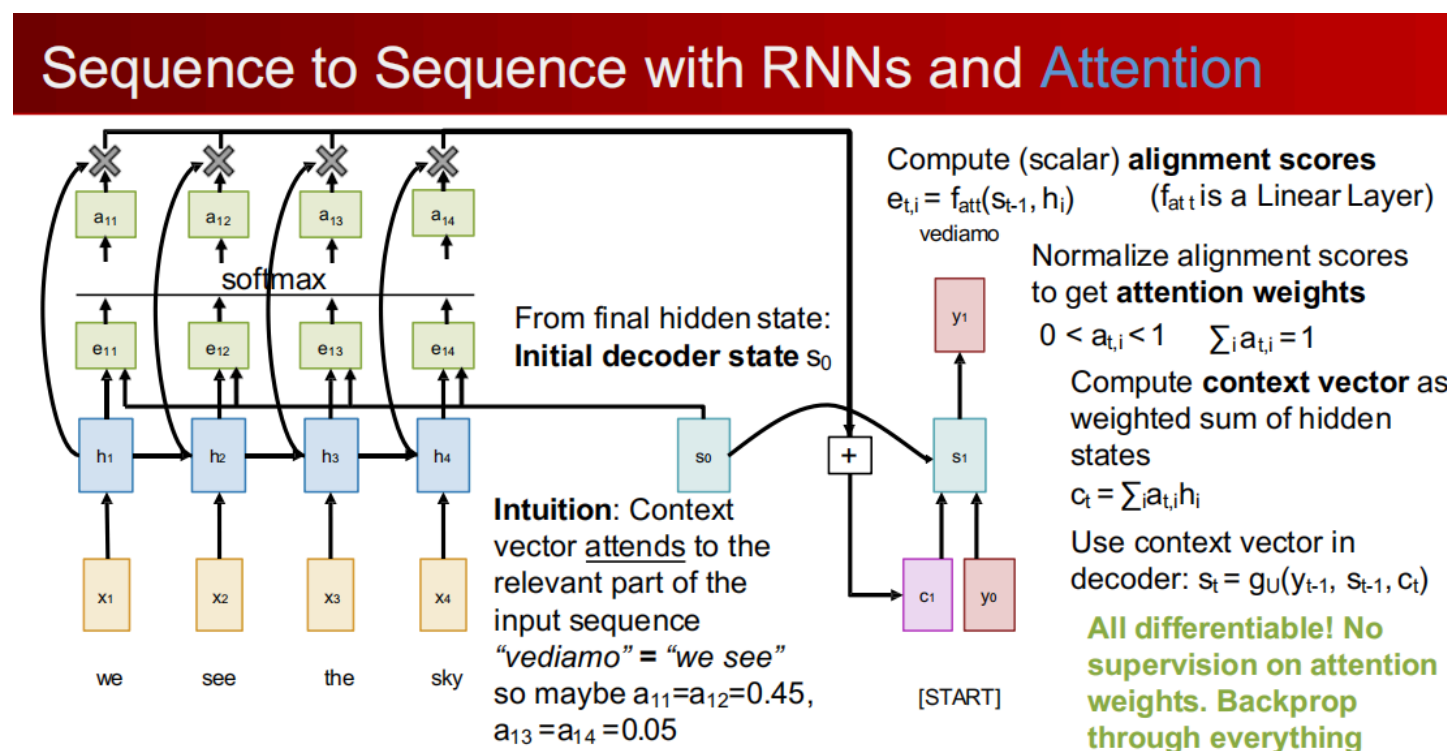


# Attention and Transformers

## Attention

attention is introduced to solve the limited representation ability of context token in sequence-to-sequence RNNs. It allows the decoder to see all the words in the input sequence at once, and focus on the relevant parts of the input sequence. That saves time(decoder predict all the words in the output sequence at once. But at eval time, it still needs to generate word by word) and improves ability(better context), but it also increase the storage.

## algorithm



Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR2015

**NOTE:** it's the hidden state in the decoder and the hidden state of the encoder that are used to calculate similarity between them.

# attention layer

## Attention Layer

### Inputs:

Query vector:  $\mathbf{Q} [N_Q \times D_Q]$

Data vectors:  $\mathbf{X} [N_X \times D_X]$

Key matrix:  $\mathbf{W}_K [D_X \times D_Q]$

Value matrix:  $\mathbf{W}_V [D_X \times D_V]$

### Computation:

Keys:  $\mathbf{K} = \mathbf{XW}_K [N_X \times D_Q]$

Values:  $\mathbf{V} = \mathbf{XW}_V [N_X \times D_V]$

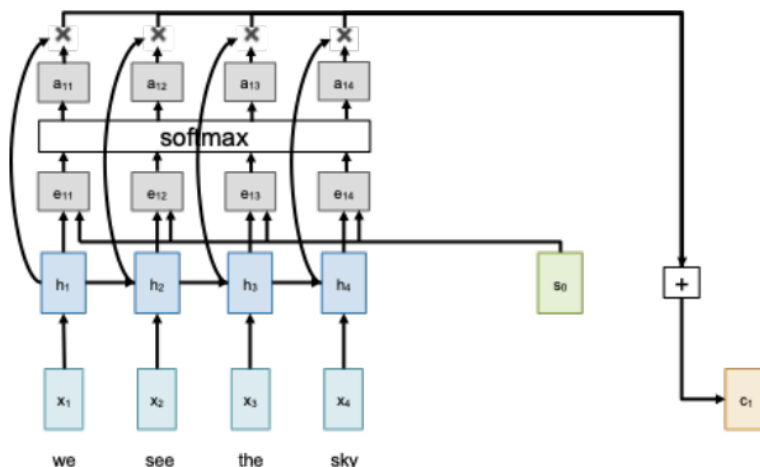
Similarities:  $\mathbf{E} = \mathbf{QK}^T / \sqrt{D_Q} [N_Q \times N_X]$

$$E_{ij} = \mathbf{Q}_i \cdot \mathbf{K}_j / \sqrt{D_Q}$$

Attention weights:  $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1) [N_Q \times N_X]$

Output vector:  $\mathbf{Y} = \mathbf{AV} [N_Q \times D_V]$

$$\mathbf{Y}_i = \sum_j \mathbf{A}_{ij} \mathbf{V}_j$$



### Changes

- Use scaled dot product for similarity
- Multiple **query** vectors
- Separate **key** and **value**

**NOTE:** why do we divide the attention weight by the square root of the hidden state dimension? Because it makes training more stable. If we make assumption that  $q(N_q \times D)$  and  $k(N_k \times D)$  are of normal distribution with variance equals to 1, then the production of  $q$  and  $k$  will have a variance of  $D$ . We divide the weight by the square root of  $D$  to make the product also has a variance of 1.

# self-attention layer

## Self-Attention Layer

### Inputs:

Input vectors:  $\mathbf{X} [N \times D_{in}]$

Key matrix:  $\mathbf{W}_K [D_{in} \times D_{out}]$

Value matrix:  $\mathbf{W}_V [D_{in} \times D_{out}]$

Query matrix:  $\mathbf{W}_Q [D_{in} \times D_{out}]$

Each **input** produces one **output**, which is a mix of information from all **inputs**

### Computation:

Queries:  $\mathbf{Q} = \mathbf{XW}_Q [N \times D_{out}]$

Keys:  $\mathbf{K} = \mathbf{XW}_K [N \times D_{out}]$

Values:  $\mathbf{V} = \mathbf{XW}_V [N \times D_{out}]$

Similarities:  $\mathbf{E} = \mathbf{QK}^T / \sqrt{D_Q} [N \times N]$

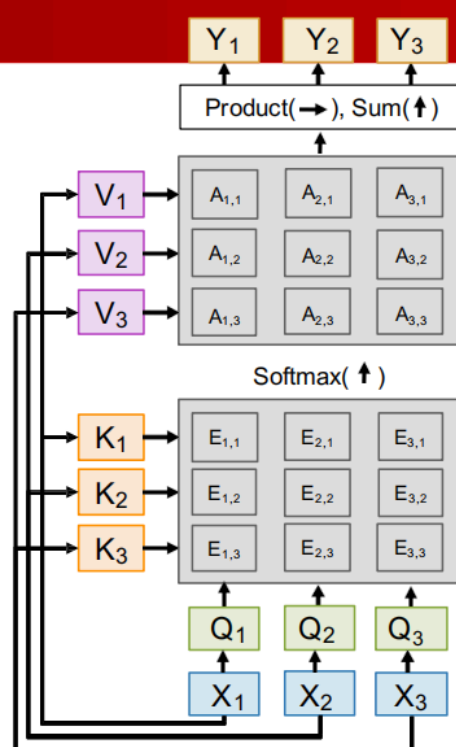
$$E_{ij} = \mathbf{Q}_i \cdot \mathbf{K}_j / \sqrt{D_Q}$$

Attention weights:  $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1) [N \times N]$

Output vector:  $\mathbf{Y} = \mathbf{AV} [N \times D_{out}]$

$$\mathbf{Y}_i = \sum_j \mathbf{A}_{ij} \mathbf{V}_j$$

Shapes get a little simpler:  
-  $N$  input vectors, each  $D_{in}$   
- Almost always  $D_Q = D_V = D_{out}$



**NOTE:** in previous algorithm, the key is actually not the word, but the hidden state, which means it's comparing the query to the sentence before this word. But in self-attention layer, the query is comparing to the word representation itself. This may seem good, but it actually loses some syntax information. In fact, naive self-attention layer is permutation equivariant, which means it doesn't care about the order of the words(or syntax). To solve this we introduce a positional encoding to add the order information.

**NOTE:** when training self-attention layer, we use a 'mask' (set the similarity to  $-\infty$  so that after softmax, the weight will be zero) to prevent the model from looking at the future words. That's because if we provide the whole sentence to the model, then the answer is just the next word, we are not sure if the model has understood the sentence or learned the trivial way to generate the next word to perform well at training time. This manner occurs also because we want to process all the inputs at once.

**NOTE:** what's the relationship between position encoding and masking? Though masking is not introduced to solve permutation equivariance problem, it still works in that way, so will it take the work of position encoding, and we don't need position encoding if we use masking? The answer is no. Masking is used in decoder to run predictions in parallel, but it's still poor at solving permutation equivariance problem. So example, in sentence 'abcde', if we are trying to decode 'e', we will calculate its similarity with 'abcd', but it's the same as doing so with 'bdac', so masking doesn't solve the problem completely.

**NOTE:** personal opinion: position encoding is used in encoder and mask is not used in encoder to better learn the word representation; position encoding is used in decoder to provide syntax information of output sequence, and maybe self-attention is used on already-generated outputs to better grasp the context before making cross-attention with the input sequence. Mask is used in decoder to run predictions in parallel.

# multiheaded attention

## Multiheaded Self-Attention Layer

Run H copies of Self-Attention in parallel

### Inputs:

Input vectors:  $X$  [ $N \times D_{in}$ ]

Key matrix:  $W_K$  [ $D_{in} \times D_{out}$ ]

Value matrix:  $W_V$  [ $D_{in} \times D_{out}$ ]

Query matrix:  $W_Q$  [ $D_{in} \times D_{out}$ ]

### Computation:

Queries:  $Q = XW_Q$  [ $N \times D_{out}$ ]

Keys:  $K = XW_K$  [ $N \times D_{out}$ ]

Values:  $V = XW_V$  [ $N \times D_{out}$ ]

Similarities:  $E = QK^T / \sqrt{D_Q}$  [ $N \times N$ ]  
 $E_{ij} = Q_i \cdot K_j / \sqrt{D_Q}$

Attention weights:  $A = \text{softmax}(E, \text{dim}=1)$  [ $N \times N$ ]

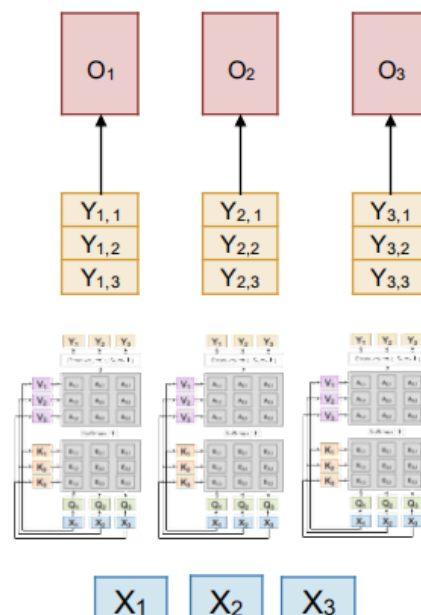
Output vector:  $Y = AX$  [ $N \times D_{out}$ ]

$$Y_i = \sum_j A_{ij} V_j$$

Output projection fuses data from each head

Stack up the H independent outputs for each input X

H = 3 independent self-attention layers (called heads), each with their own weights



## Multiheaded Self-Attention Layer

Run H copies of Self-Attention in parallel

### Inputs:

Input vectors:  $X$  [ $N \times D$ ]

Key matrix:  $W_K$  [ $D \times HD_H$ ]

Value matrix:  $W_V$  [ $D \times HD_H$ ]

Query matrix:  $W_Q$  [ $D \times HD_H$ ]

Output matrix:  $W_O$  [ $HD_H \times D$ ]

### Computation:

Queries:  $Q = XW_Q$  [ $H \times N \times D_H$ ]

Keys:  $K = XW_K$  [ $H \times N \times D_H$ ]

Values:  $V = XW_V$  [ $H \times N \times D_H$ ]

Similarities:  $E = QK^T / \sqrt{D_Q}$  [ $H \times N \times N$ ]

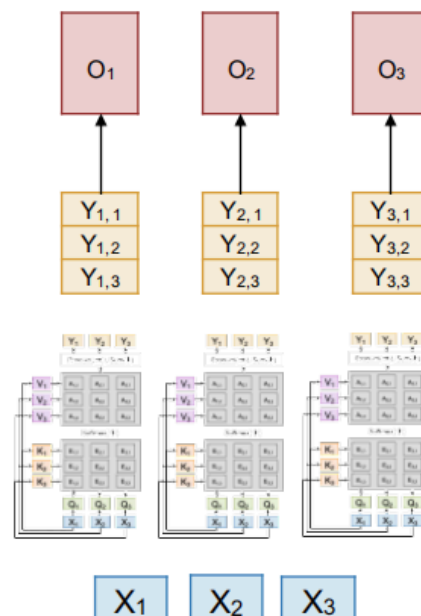
Attention weights:  $A = \text{softmax}(E, \text{dim}=1)$  [ $H \times N \times N$ ]

Head outputs:  $Y = AV$  [ $H \times N \times D_H$ ]  $\Rightarrow$  [ $N \times HD_H$ ]

Outputs:  $O = YW_O$  [ $N \times D$ ]

Each of the H parallel layers use a qkv dim of  $D_H$  = "head dim"

Usually  $D_H = D / H$ , so inputs and outputs have the same dimension



**NOTE:** multiheaded attention doesn't increase the calculation cost, but it greatly increase the expressiveness for the model. It allows the model to focus on different parts of the input sequence at different times.

## complexity

### Self-Attention is Four Matrix Multiplies!

#### Inputs:

Input vectors:  $\mathbf{X}$   $[N \times D]$

Key matrix:  $\mathbf{W}_K$   $[D \times HD_H]$

Value matrix:  $\mathbf{W}_V$   $[D \times HD_H]$

Query matrix:  $\mathbf{W}_Q$   $[D \times HD_H]$

Output matrix:  $\mathbf{W}_O$   $[HD_H \times D]$

#### Computation:

Queries:  $\mathbf{Q} = \mathbf{XW}_Q$   $[H \times N \times D_H]$

Keys:  $\mathbf{K} = \mathbf{XW}_K$   $[H \times N \times D_H]$

Values:  $\mathbf{V} = \mathbf{XW}_V$   $[H \times N \times D_H]$

Similarities:  $\mathbf{E} = \mathbf{QK}^T / \sqrt{D_Q}$   $[H \times N \times N]$

Attention weights:  $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$   $[H \times N \times N]$

Head outputs:  $\mathbf{Y} = \mathbf{AV}$   $[H \times N \times D_H] \Rightarrow [N \times HD_H]$

Outputs:  $\mathbf{O} = \mathbf{YW}_O$   $[N \times D]$

#### 1. QKV Projection

$[N \times D] [D \times 3HD_H] \Rightarrow [N \times 3HD_H]$

Split and reshape to get  $\mathbf{Q}$ ,  $\mathbf{K}$ ,  $\mathbf{V}$  each of shape  $[H \times N \times D_H]$

#### 2. QK Similarity

$[H \times N \times D_H] [H \times D_H \times N] \Rightarrow [H \times N \times N]$

#### 3. V-Weighting

$[H \times N \times N] [H \times N \times D_H] \Rightarrow [H \times N \times D_H]$

Reshape to  $[N \times HD_H]$

#### 4. Output Projection

$[N \times HD_H] [HD_H \times D] \Rightarrow [N \times D]$

Q1: how much compute does it take as the number of tokens(N) increases?

A1:  $O(N^2)$ , as a result of the second and third matrix multiplications.

Q2: how much memory does it take as the number of tokens(N) increases?

A2:  $O(N^2)$ , as a result of the attention weights.

S2: use Flash Attention can calculate the second and third step together without storing attention weights. This helps to reduce memory cost to  $O(N)$

## Three Ways of Processing Sequences

### 1. RNN

(+) Theoretically good at long sequences:  $O(N)$  compute and memory for a sequence of length N

(-) Not parallelizable. Need to compute hidden states sequentially

### 2. CNN

(+) Theoretically good at long sequences:  $O(N)$  compute and memory for a sequence of length N

(-) Not parallelizable. Need to compute hidden states sequentially

### 3. Self-Attention

(+) Great for long sequences; each output depends directly on all inputs

(+) Highly parallel, it's just 4 matmuls

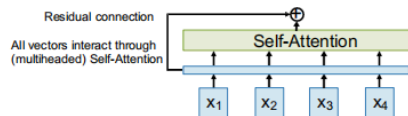
(-) Expensive:  $O(N^2)$  compute,  $O(N)$  memory for sequence of length N

# Transformer

## The Transformer

### Transformer Block

Input: Set of vectors  $x$



Viewari et al, "Attention is all you need," NeurIPS 2017

## The Transformer

### Transformer Block

Input: Set of vectors  $x$

Recall Layer Normalization (Baet al, 2016):

Given  $h_1, \dots, h_N$  (Shape:  $D$ )

scale:  $\gamma$  (Shape:  $D$ )

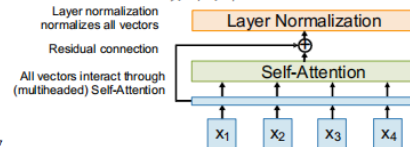
shift:  $\beta$  (Shape:  $D$ )

$\mu_i = (\sum_j h_{i,j})/D$  (scalar)

$\delta_i = (\sum_j (h_{i,j} - \mu_i)^2/D)^{1/2}$  (scalar)

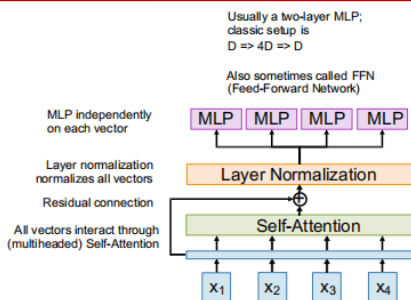
$z_i = (h_i - \mu_i)/\delta_i$

$y_i = \gamma \cdot z_i + \beta$



Viewari et al, "Attention is all you need," NeurIPS 2017

## The Transformer

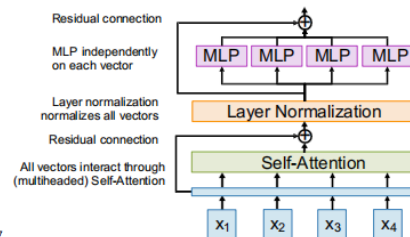


Viewari et al, "Attention is all you need," NeurIPS 2017

## The Transformer

### Transformer Block

Input: Set of vectors  $x$

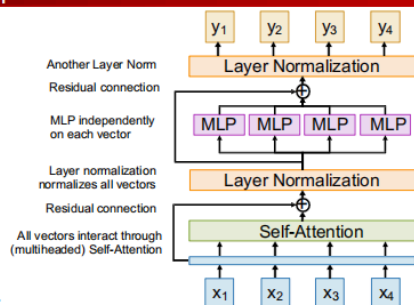


Viewari et al, "Attention is all you need," NeurIPS 2017

## The Transformer

### Transformer Block

Input: Set of vectors  $x$



Viewari et al, "Attention is all you need," NeurIPS 2017

## The Transformer

### Transformer Block

Input: Set of vectors  $x$

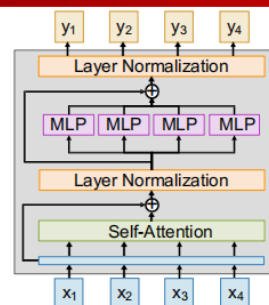
Output: Set of vectors  $y$

Self-Attention is the only interaction between vectors

LayerNorm and MLP work on each vector independently

Highly scalable and parallelizable, most of the compute is just 6 matmuls:

4 from Self-Attention  
2 from MLP



Viewari et al, "Attention is all you need," NeurIPS 2017

## NOTE:

1. layer normalization is applied to every word representation.
2. self-attention layer also has a FC layer at the end

# The Transformer

## Transformer Block

**Input:** Set of vectors  $x$

**Output:** Set of vectors  $y$

Self-Attention is the only interaction between vectors

LayerNorm and MLP work on each vector independently

Highly scalable and parallelizable, most of the compute is just 6 matmuls:

4 from Self-Attention

2 from MLP

A **Transformer** is just a stack of identical Transformer blocks!

They have not changed much since 2017... but have gotten a lot bigger

Original: [Vaswani et al, 2017]

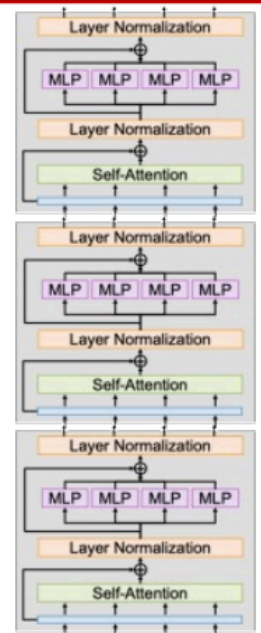
12 blocks,  $D=1024$ ,  $H=16$ ,  $N=512$   
213M params

GPT-2: [Radford et al, 2019]

48 blocks,  $D=1600$ ,  $H=25$ ,  $N=1024$   
1.5B params

GPT-3: [Brown et al, 2020]

96 blocks,  $D=12288$ ,  $H=96$ ,  $N=2048$   
175B params



Vaswani et al, "Attention is all you need," NeurIPS 2017

## NOTE:

1. parameters:

- self-attention:  $Q(D \times D)$ ,  $K(D \times D)$ ,  $V(D \times D)$ , output( $D \times D$ )
- feed-forward:  $FC1(D \times 4D)$ ,  $FC2(4D \times D)$

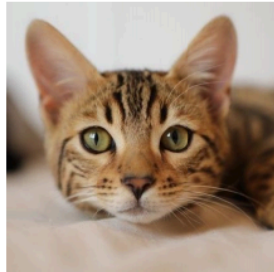
## Transformers for Language Modeling

- Learn an embedding matrix at the start of the model to convert words into vectors.
- Given vocab size  $V$  and model dimension  $D$ , it's a lookup table of shape  $[V \times D]$ .
- Use masked attention inside each transformer block so each token can only see the ones before it
- At the end, learn a projection matrix of shape  $[D \times V]$  to project each  $D$ -dim vector to a  $V$ -dim vector of scores for each element of the vocabulary.
- Train to predict next token using softmax + cross-entropy loss
- good for translation

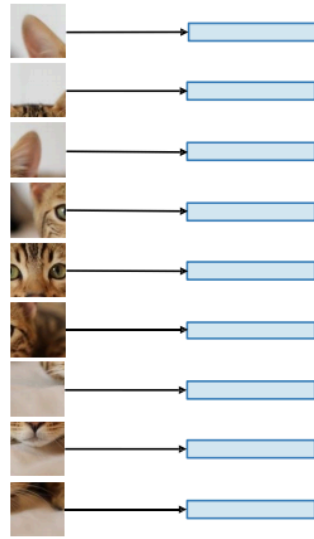


# Vision Transformers (ViT)

## Vision Transformers (ViT)



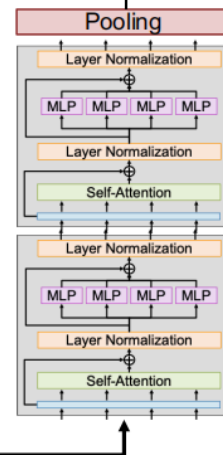
Input image:  
e.g. 224x224x3



Break into patches  
e.g. 16x16x3

Flatten and apply a linear  
transform  $768 \Rightarrow D$

Average pool  $N \times D$  vectors to  
 $1 \times D$ , apply a linear layer  
 $D \Rightarrow C$  to predict class scores



Transformer  
gives an output  
vector per patch

Don't use any  
masking; each  
image patch can  
look at all other  
image patches

Use positional  
encoding to tell  
the transformer  
the 2D position  
of each patch

$D$ -dim vector per patch  
are the input vectors to  
the Transformer

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

### NOTE:

1. another way to process image patch instead of flatten it is to apply conv layers with output channels equals to  $D$ .

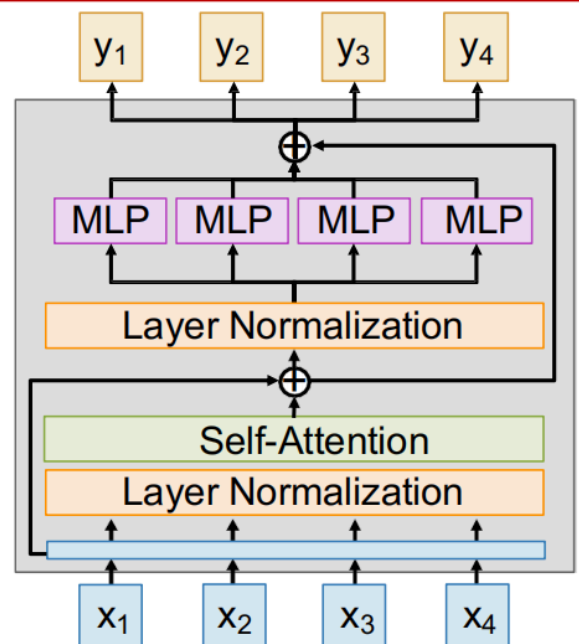
## Tweaking Transformers

### Pre-Norm Transformer

**Layer normalization** is outside  
the residual connections

Kind of weird, the model can't  
actually learn the identity function

**Solution:** Move layer  
normalization before the Self-  
Attention and MLP, inside the  
residual connections. Training is  
more stable.



Baevski & Auli, "Adaptive Input Representations for Neural Language Modeling", arXiv 2018



# RMSNorm

Replace Layer Normalization with Root-Mean-Square Normalization (RMSNorm)

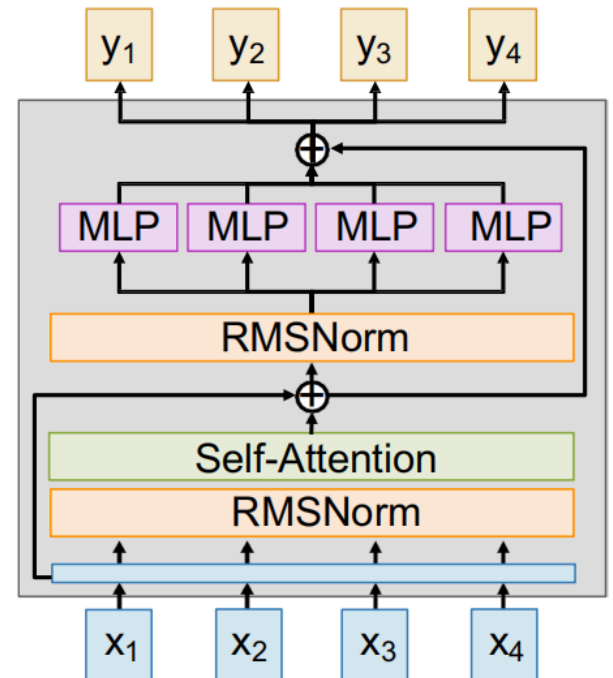
**Input:**  $x$  [shape  $D$ ]  
**Output:**  $y$  [shape  $D$ ]  
**Weight:**  $\gamma$  [shape  $D$ ]

$$y_i = \frac{x_i}{RMS(x)} * \gamma_i$$

$$RMS(x) = \sqrt{\epsilon + \frac{1}{N} \sum_{i=1}^N x_i^2}$$

Training is a bit more stable

Zhang and Sennrich, "Root Mean Square Layer Normalization", NeurIPS 2019



## NOTE:

1. RMS gives all word representations the same norm.

# SwiGLU MLP

## Classic MLP:

**Input:**  $X$  [ $N \times D$ ]  
**Weights:**  $W_1$  [ $D \times 4D$ ]  
 $W_2$  [ $4D \times D$ ]  
**Output:**  $Y = \delta(XW_1) W_2$  [ $N \times D$ ]

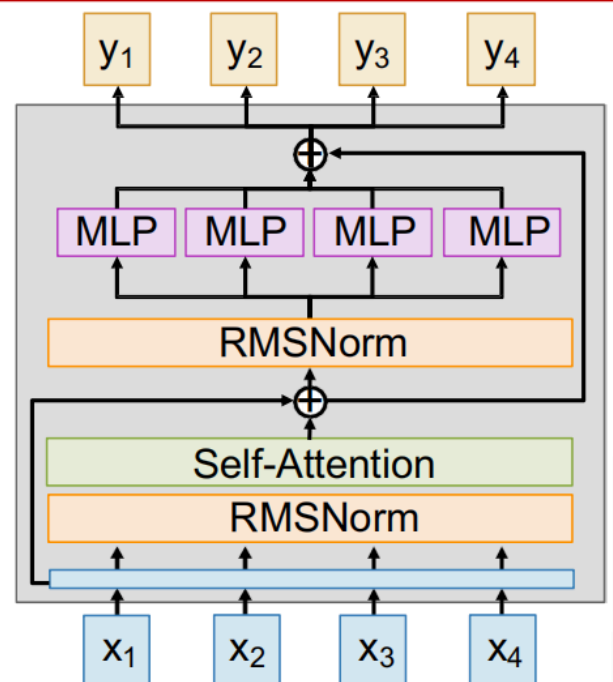
## SwiGLU MLP:

**Input:**  $X$  [ $N \times D$ ]  
**Weights:**  $W_1, W_2$  [ $D \times H$ ]  
 $W_3$  [ $D \times H$ ]

**Output:**

$$Y = (\delta(XW_1) \odot XW_2) W_3$$

Setting  $H = 8D/3$  keeps same total params



Shazeer, "GLU Variants Improve Transformers", 2020

## NOTE:

1. SwiGLU is a bit like LSTM. It learns some forget gate and output gate to control the information flow.

# Mixture of Experts (MoE)

Learn  $E$  separate sets of MLP weights in each block; each MLP is an *expert*

$W_1: [D \times 4D] \Rightarrow [E \times D \times 4D]$

$W_2: [4D \times D] \Rightarrow [E \times 4D \times D]$

Each token gets *routed* to  $A < E$  of the experts. These are the *active experts*.

Increases params by  $E$ ,  
But only increases compute by  $A$

All of the biggest LLMs today (e.g. GPT4o, GPT4.5, Claude 3.7, Gemini 2.5 Pro, etc) almost certainly use MoE and have  $> 1T$  params; but they don't publish details anymore

