

CIS 6930 Topics in Computing for Data Science

Week 8: Transformers (1)+(2)

10/21/2021
Yoshihiko (Yoshi) Suhara

2pm-3:20pm

Week 8: Transformers

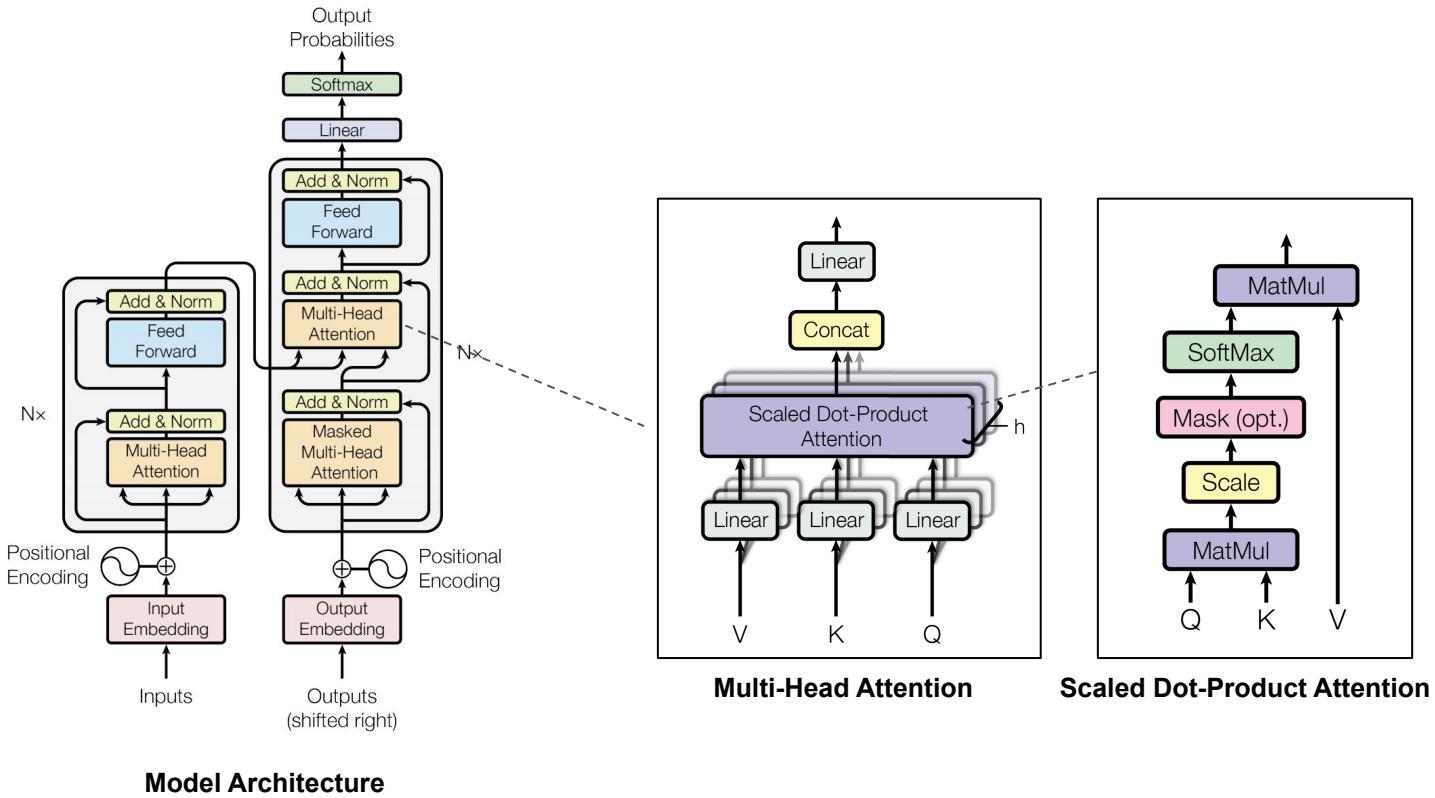
- **Week 8: Transformers**
- **Week 9: Pre-trained Language Models**
- Week 10: More Machine Learning Techniques
- Week 11: More Deep Learning Techniques for NLP (Text generation, Text summarization, Information Extraction etc.)
- Week 12: Advanced Techniques and Challenges
- Week 13: Final project presentations

Attention!

- Today's content is extremely important for understanding modern Deep Learning techniques
- The content might be technically difficult and complicated
 - Don't worry. The concept is relatively simple
- Please make sure that you follow the explanation of each step
 - Don't hesitate to interrupt and ask questions if you don't



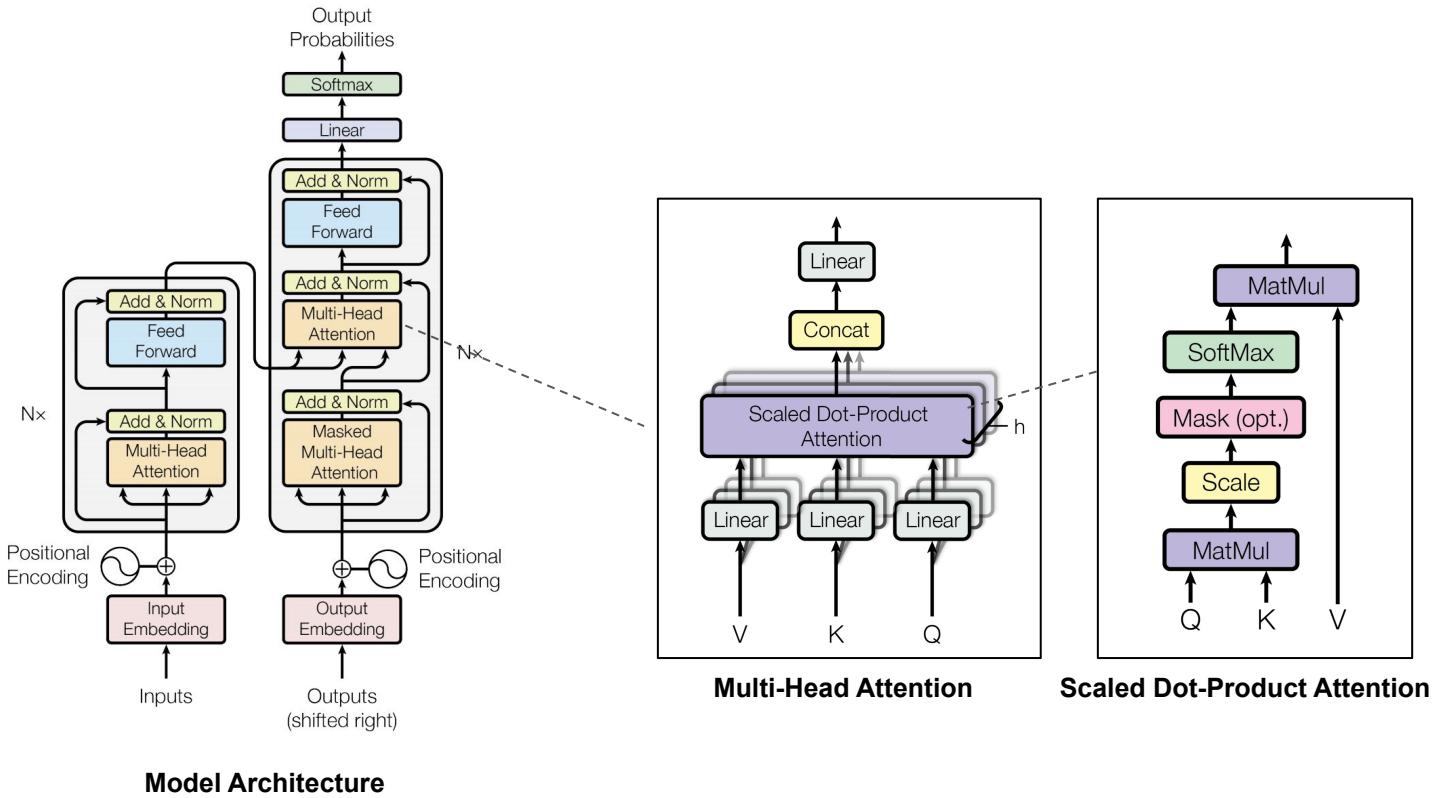
This Week's Goal: Understanding the Transformer Architecture



The Transformer: Basic Concepts

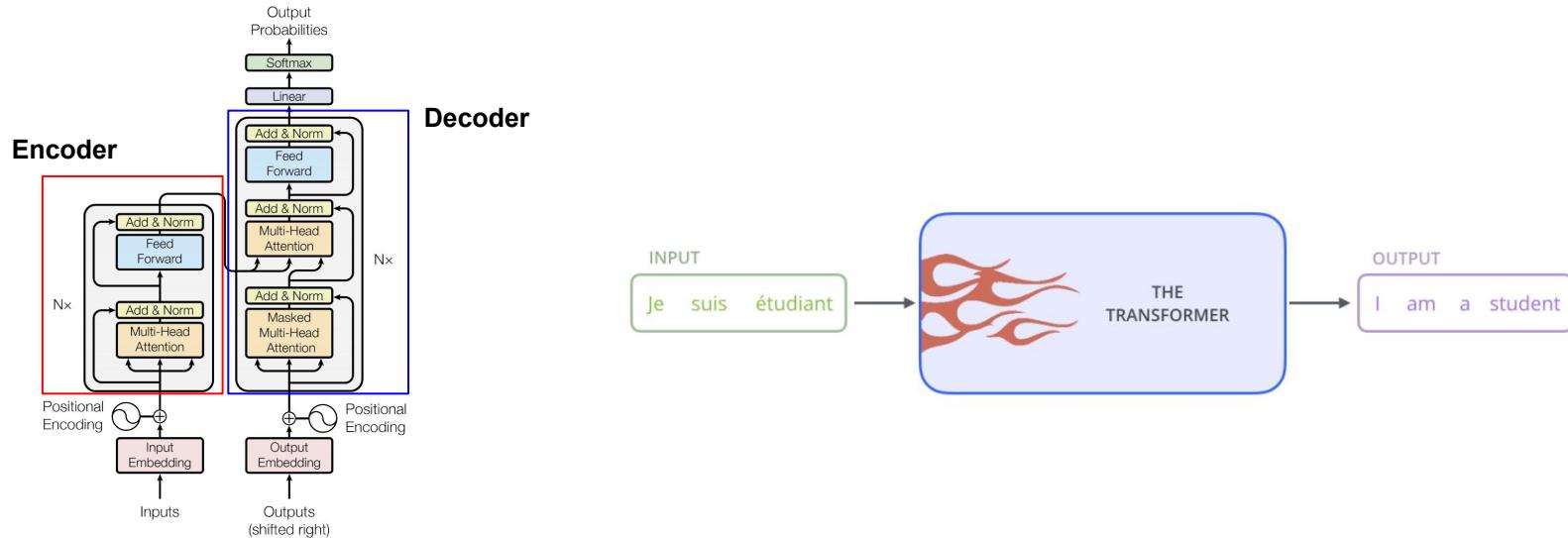
The figures and explanations are from the following excellent blog article (highly recommended)
[The Illustrated Transformer – Jay Alammar – Visualizing machine learning one concept at a time.](#)

The Transformer



What is the Transformer?

- Transformer is an **Encoder-Decoder model!**
 - Can be used for any sequence-to-sequence generation tasks
- Transformer = **Transformer Encoder + Transformer Decoder**

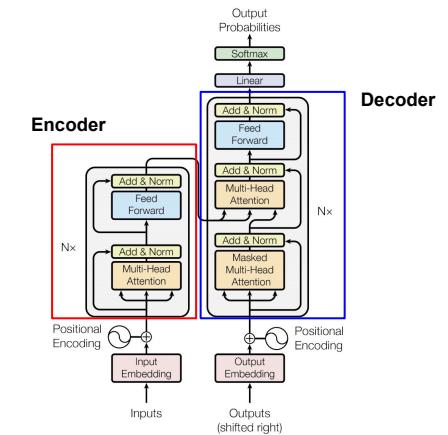
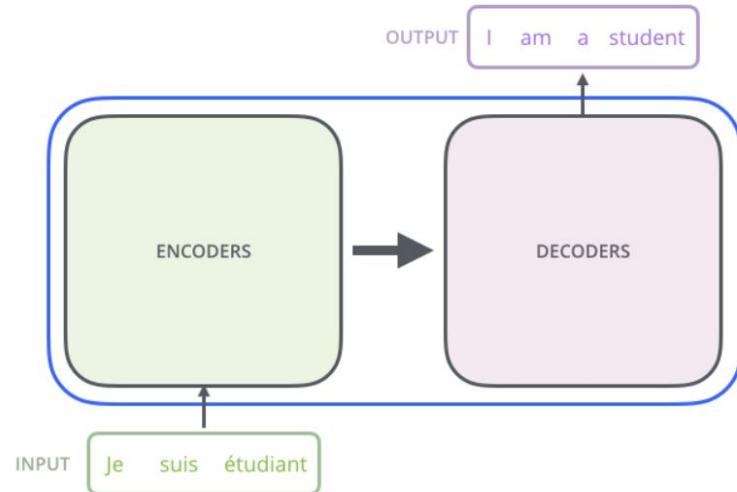


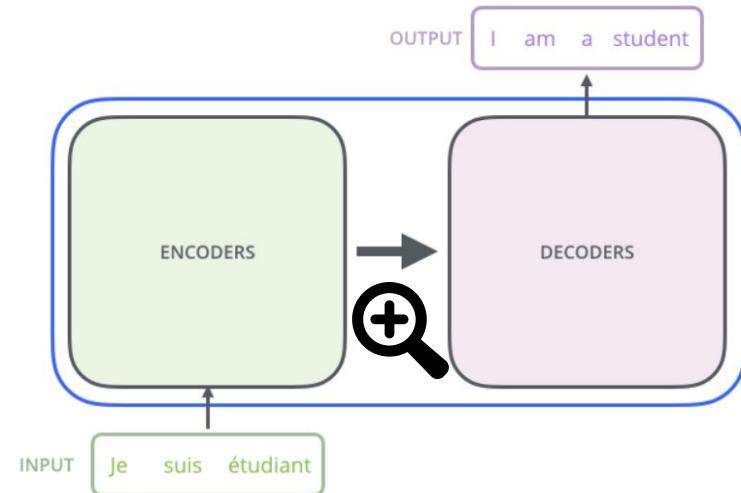
Transformer



Transformer

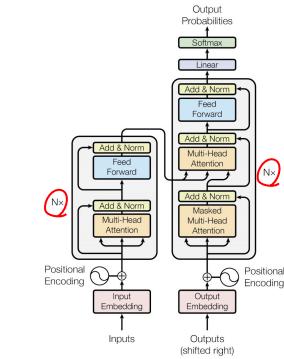
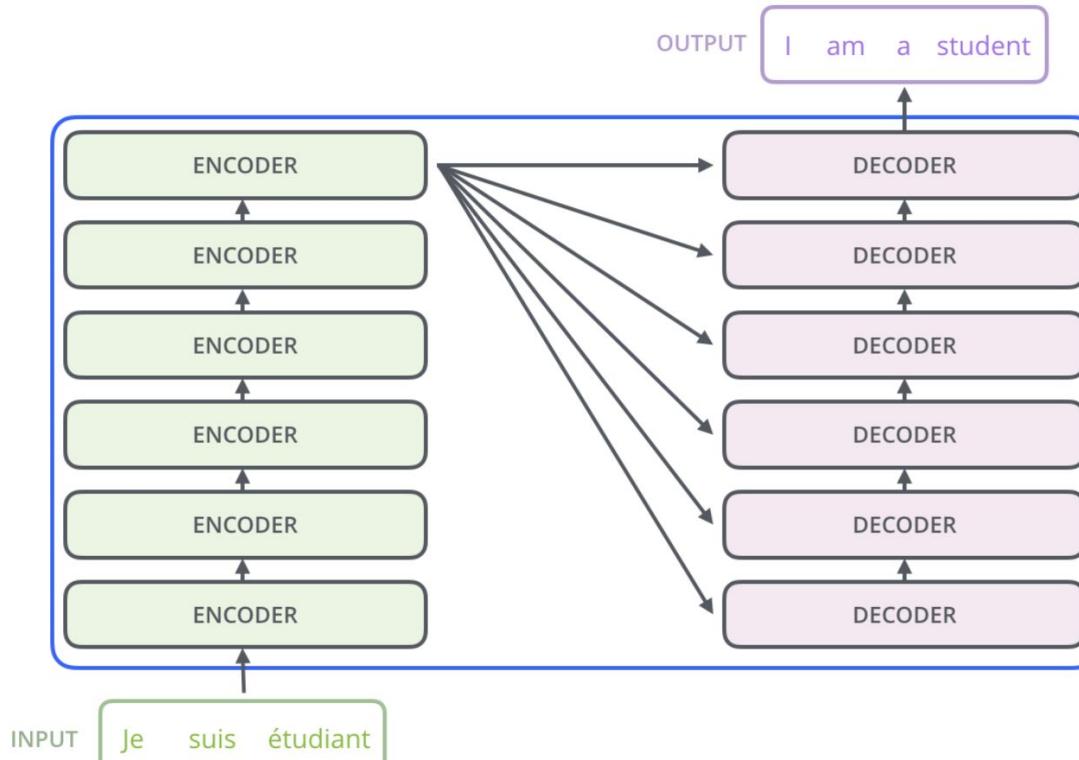






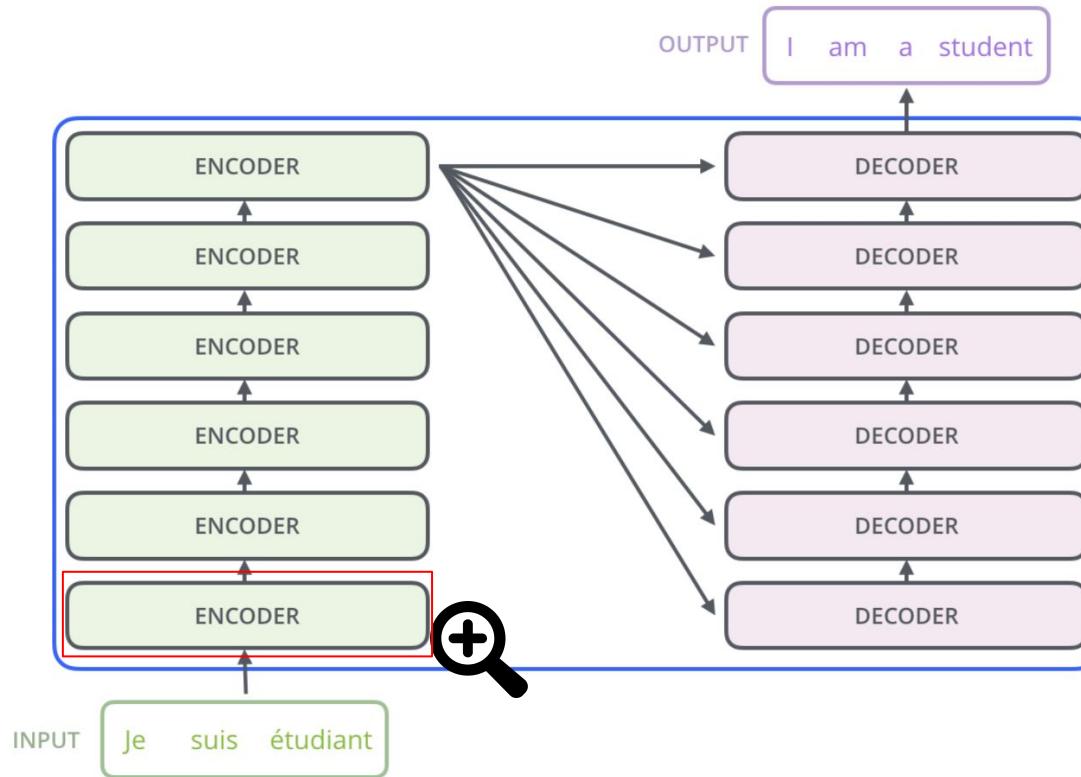
Encoder/Decoder

= Stack of Encoder/Decoder Blocks



Encoder/Decoder

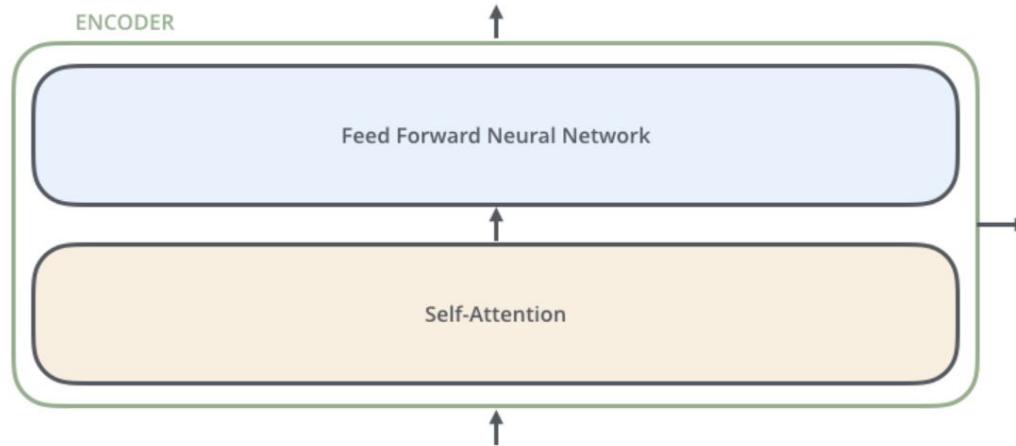
= Stack of Encoder/Decoder Blocks



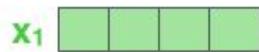


Encoder Block

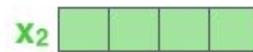
= Self-Attention + Feed Forward NN



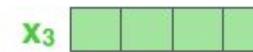
Step 1: Input Tokens → Token Embeddings



Je



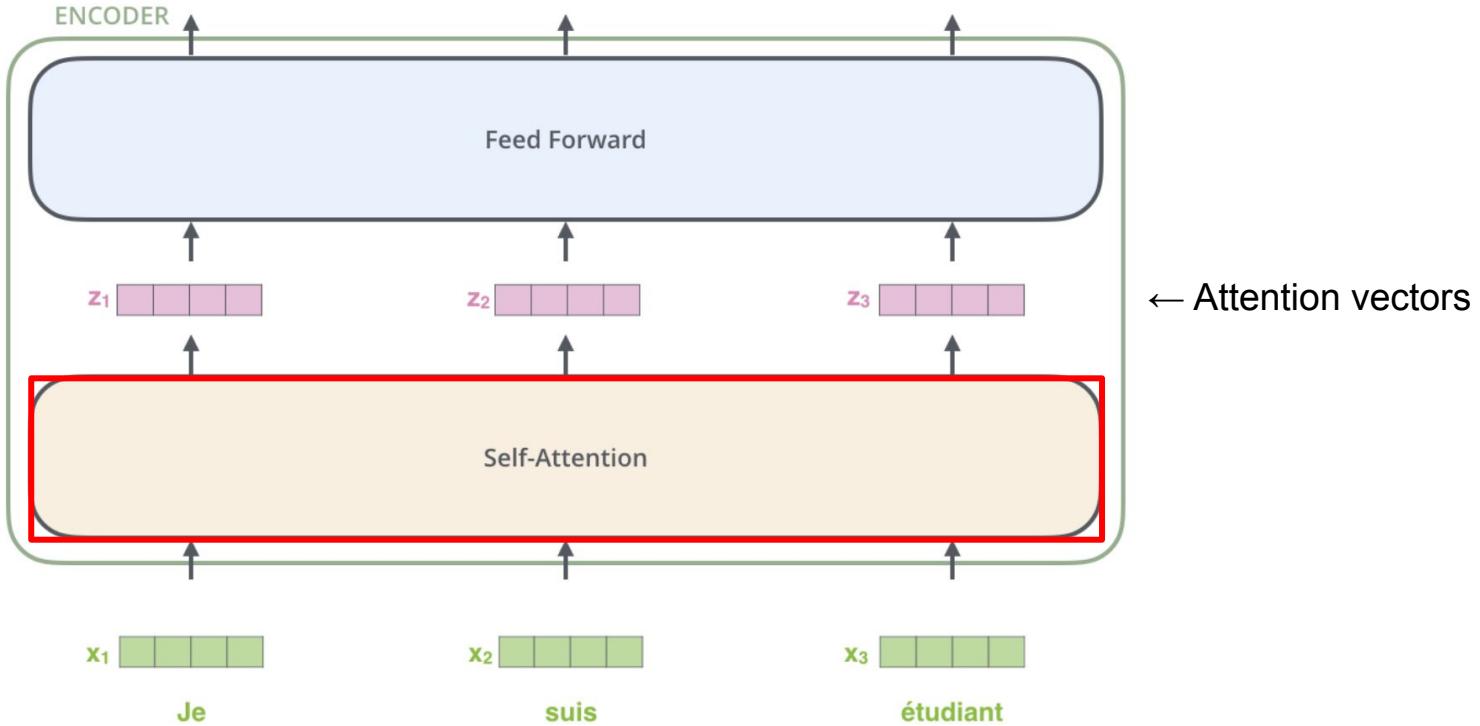
suis



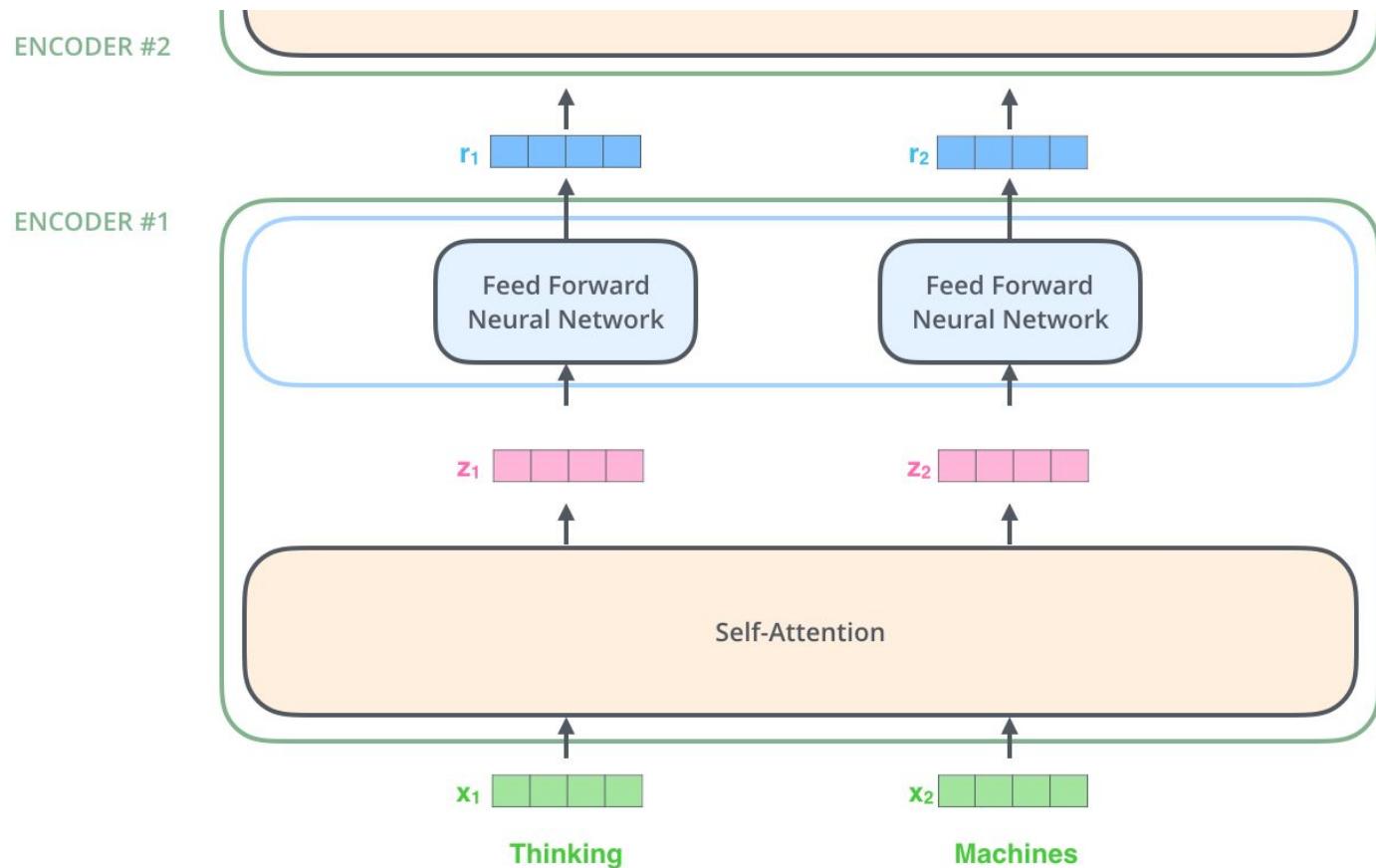
étudiant

Will come back to the positional encoding later

Step 2: Token Embeddings → Attention Vectors



Step 3: Attention Vectors → Output Vectors



Repeat Steps 1-3 N Times

Encoder #2 takes r_* as if token embeddings (e.g., x_*)

ENCODER #2

ENCODER #1

Same dimensions so that we can use the same encoder architecture as the next step

r_1

r_2

Feed Forward
Neural Network

Feed Forward
Neural Network

z_1

z_2

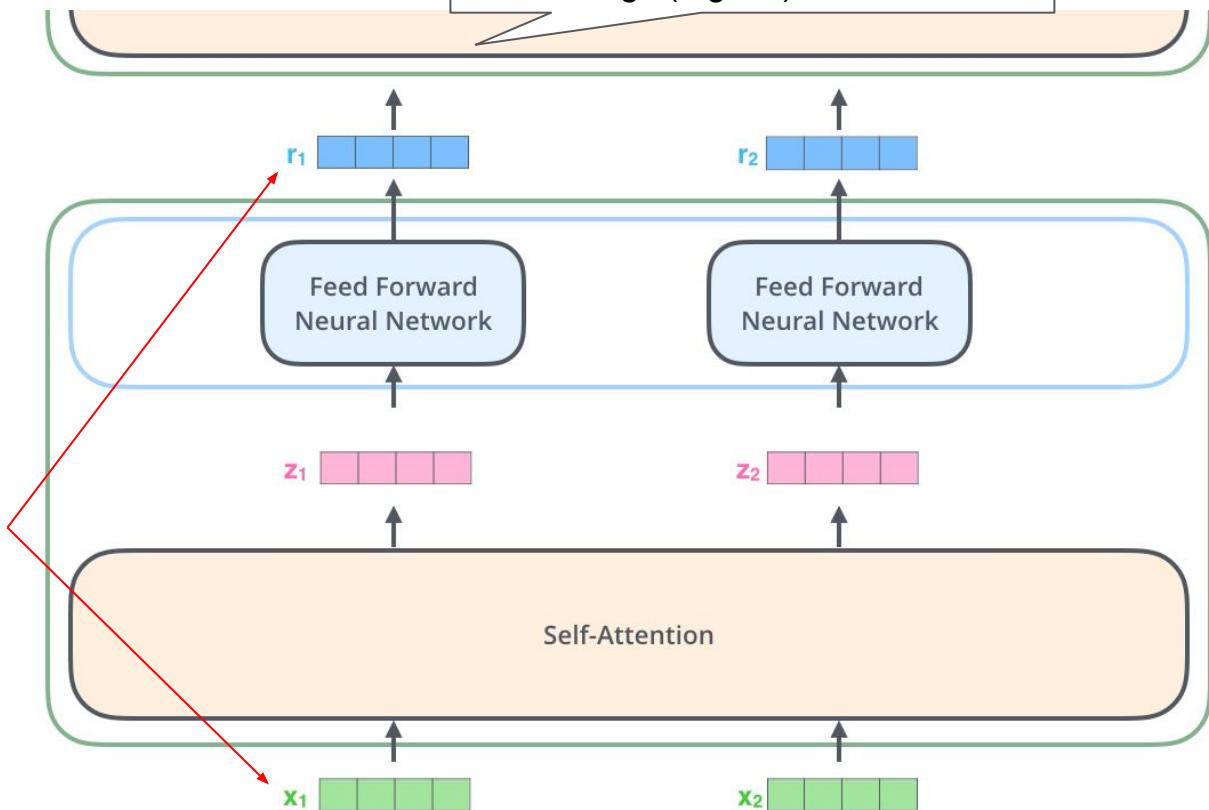
Self-Attention

x_1

x_2

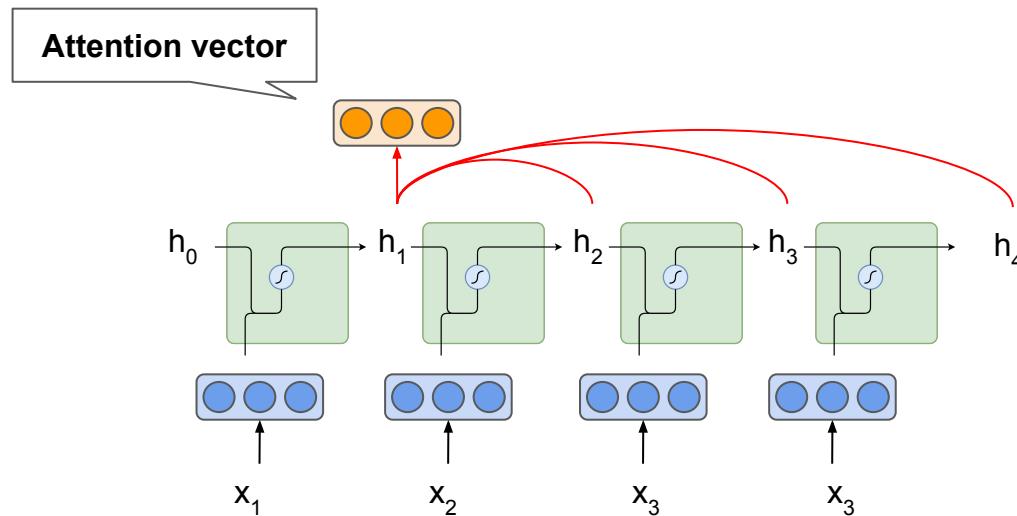
Thinking

Machines



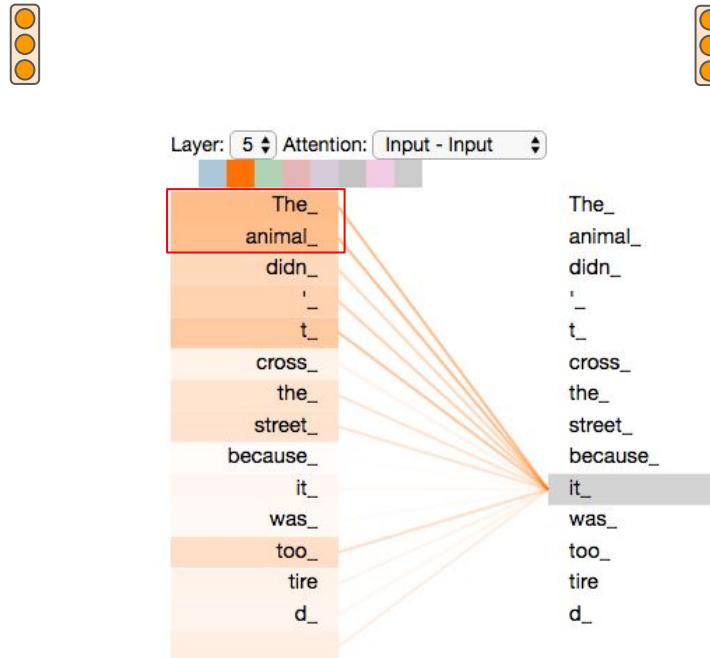
Self-Attention in Detail

Recap: Self-Attention for RNNs



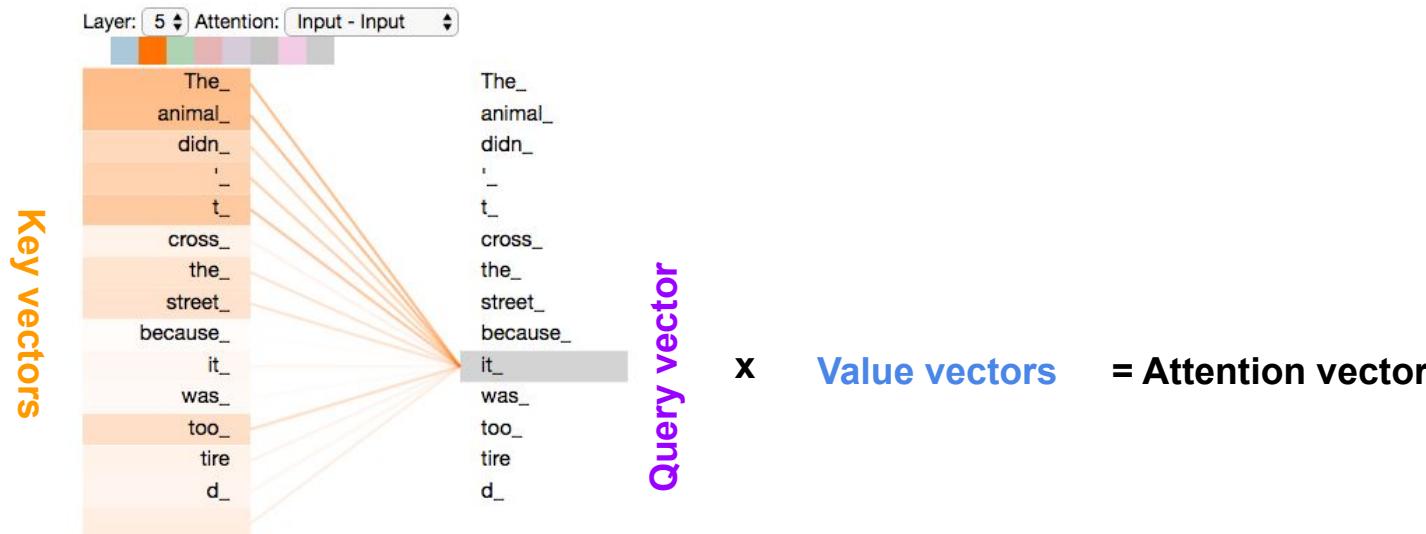
Self-Attention: Key Idea

- Self-Attention looks at **all tokens** in the input sequence to calculate **the target token embedding** based on the context
- Example: “The animal didn’t cross the street because **it** was too tired.”



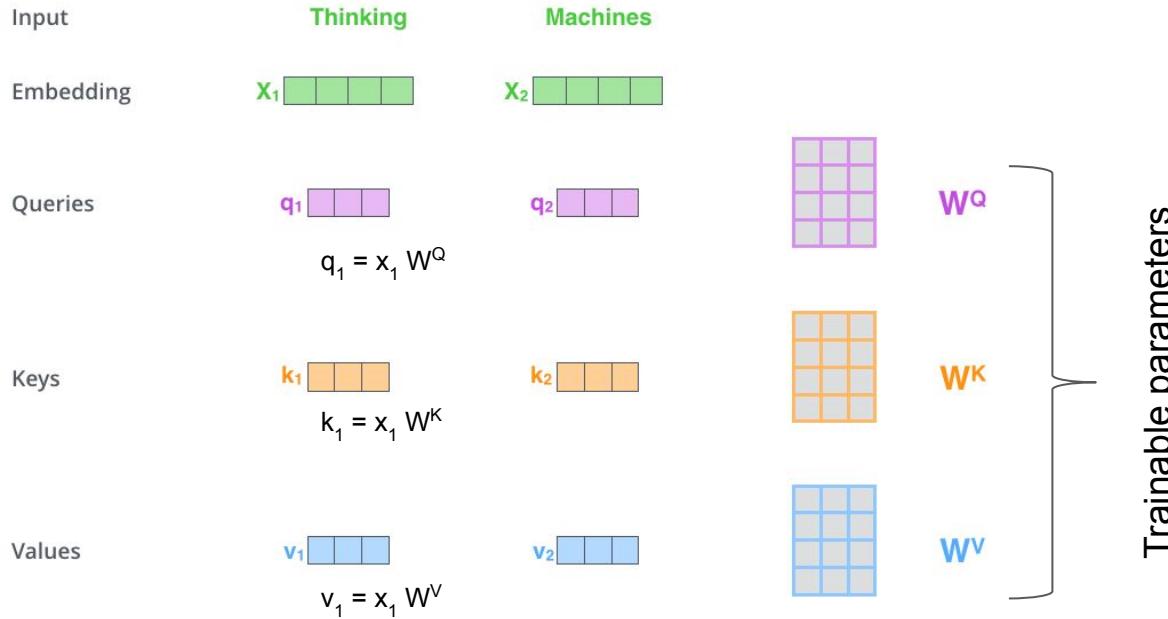
Self-Attention: Key, Query, Value Vectors

- The Self-Attention Layer encodes token embeddings by averaging **Value vectors** of input tokens based on the similarities between **Key vectors** of input tokens and **Query vector** of the target token

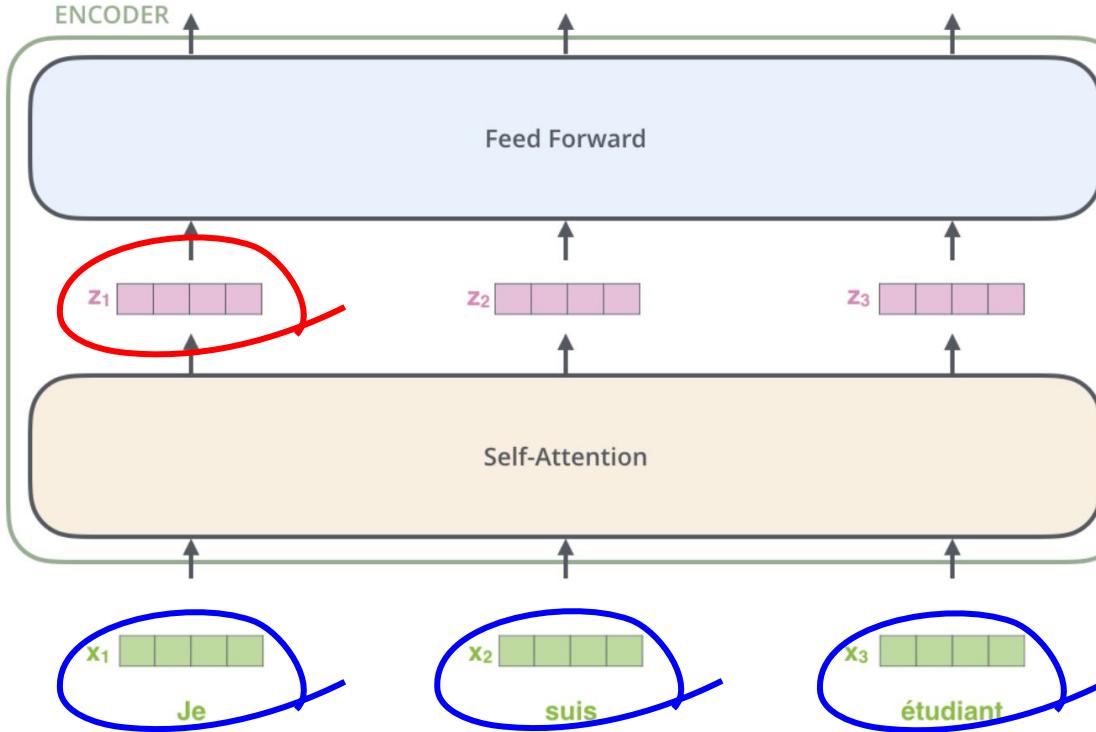


Token Embeddings → Query, Key, Value Vectors

- Self-attention layer has **transformation matrices** for Query, Key, and Value vectors



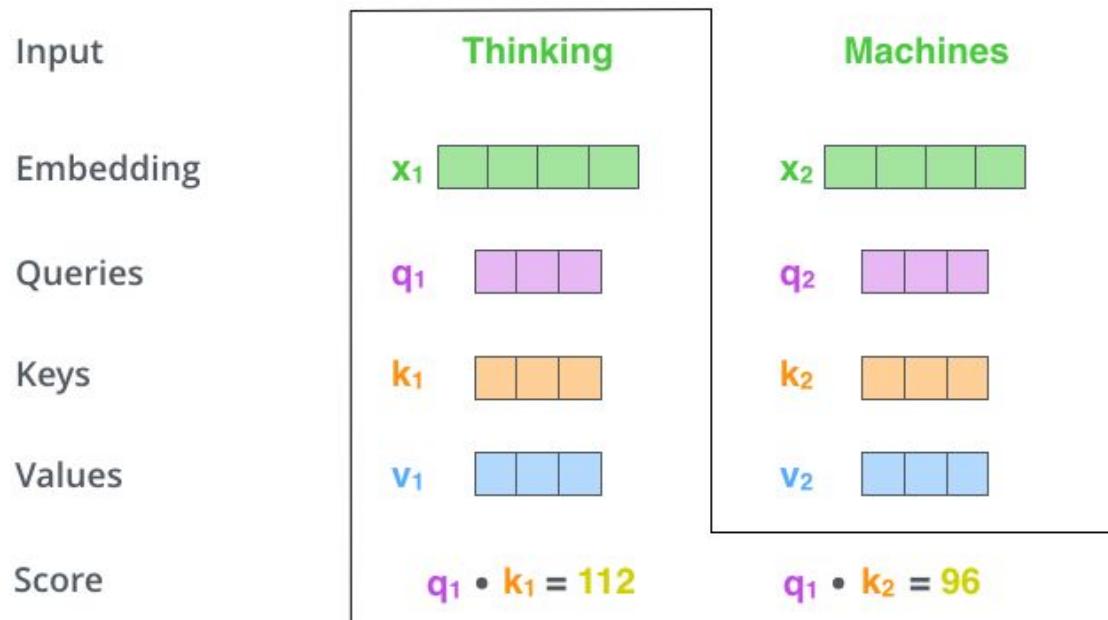
Calculating z_1 from Input Token Embeddings x_{1-3}



Example: Calculating Attention Vector for “Thinking”

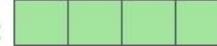
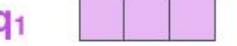
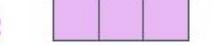
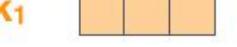
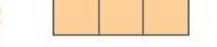
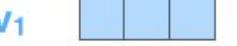
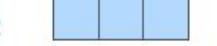
Step 1: Query-Key Similarity Calculation

- The inner product of q_1 and k_*



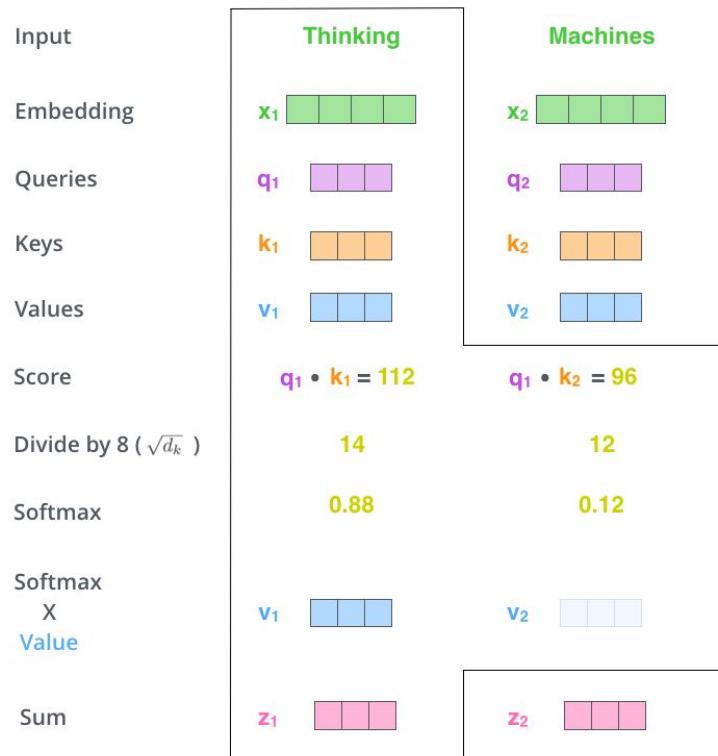
Example: Calculating Attention Vector for “Thinking”

Step 2: Scaling and normalization

Input		
Embedding	x_1 	x_2 
Queries	q_1 	q_2 
Keys	k_1 	k_2 
Values	v_1 	v_2 
Score	$q_1 \cdot k_1 = 112$	$q_1 \cdot k_2 = 96$
Scaling	Divide by 8 ($\sqrt{d_k}$) <small>d_k: Dimension of key vector</small>	14 12
Normalization	Softmax	0.88 0.12

Example: Calculating Attention Vector for “Thinking”

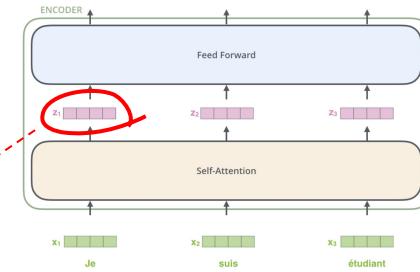
Step 3: Weighted Average of Value Vectors



Example: Calculating Attention Vector for “Thinking”

Step 3: Weighted Average of Value Vectors

Input	Thinking	Machines
Embedding	x_1	x_2
Queries	q_1	q_2
Keys	k_1	k_2
Values	v_1	v_2
Score	$q_1 \cdot k_1 = 112$	$q_1 \cdot k_2 = 96$
Divide by 8 ($\sqrt{d_k}$)	14	12
Softmax	0.88	0.12
Softmax X Value	v_1	v_2
Sum	z_1	z_2



Understanding Self-Attention in Matrix Calculation

$$X \times W^Q = Q$$

Diagram illustrating the calculation of Query (Q) matrix. An input matrix X (green, 3x4) is multiplied by a weight matrix W^Q (purple, 4x4). The result is a Query matrix Q (purple, 3x3).

$$X \times W^K = K$$

Diagram illustrating the calculation of Key (K) matrix. An input matrix X (green, 3x4) is multiplied by a weight matrix W^K (orange, 4x4). The result is a Key matrix K (orange, 3x3).

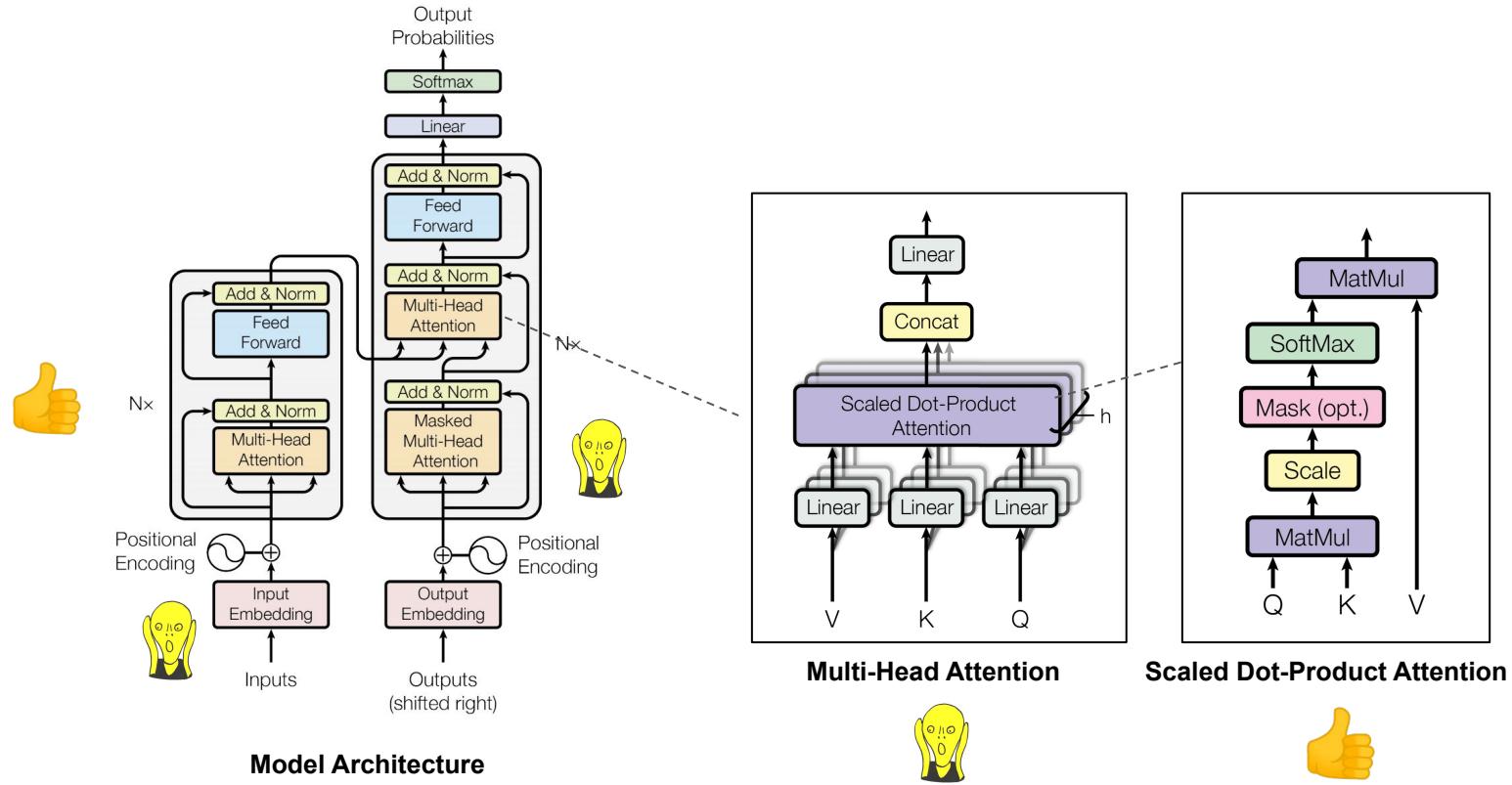
$$X \times W^V = V$$

Diagram illustrating the calculation of Value (V) matrix. An input matrix X (green, 3x4) is multiplied by a weight matrix W^V (blue, 4x4). The result is a Value matrix V (blue, 3x3).

$$\text{softmax}\left(\frac{Q \times K^T}{\sqrt{d_k}}\right) = Z$$

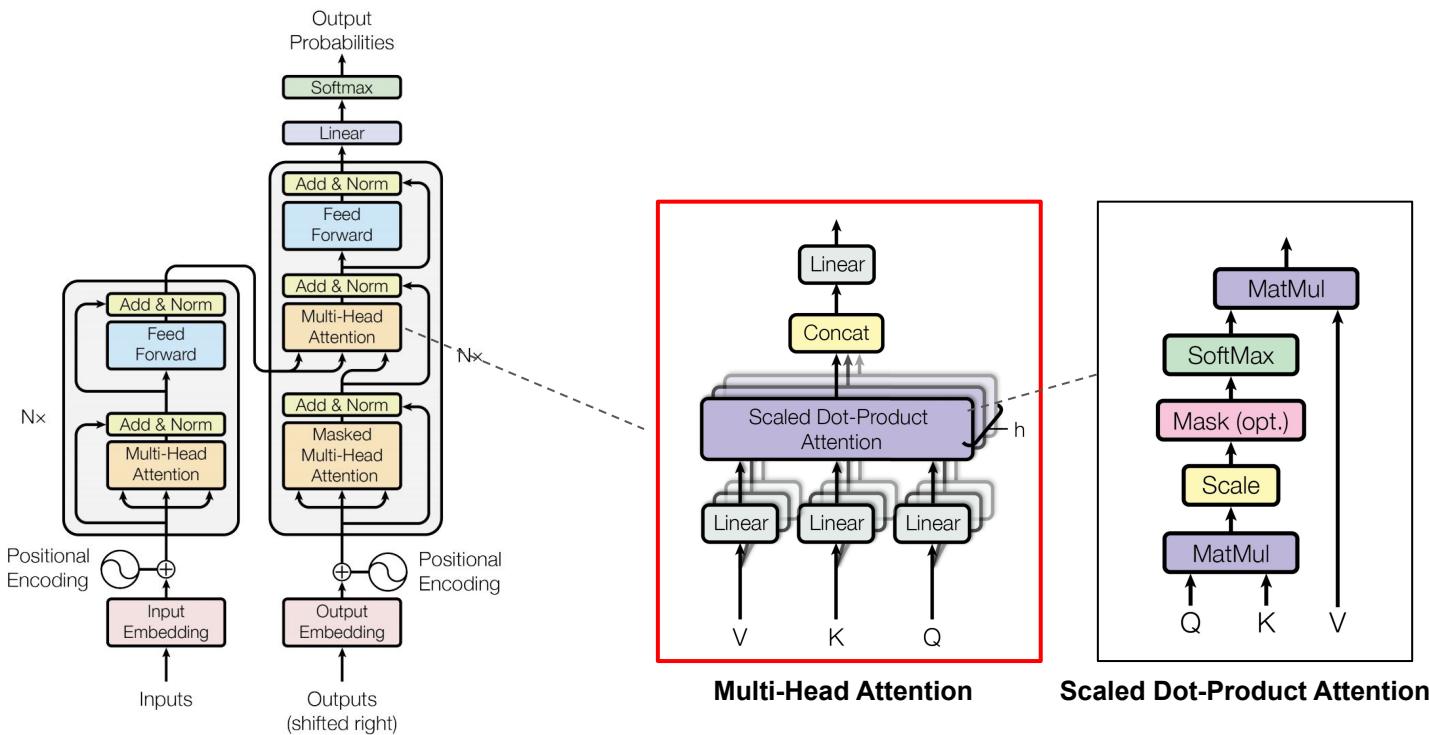
Diagram illustrating the calculation of the attention weights matrix Z. The Query matrix Q (purple, 3x3) and the transpose of the Key matrix K^T (orange, 3x3) are multiplied, and the result is divided by $\sqrt{d_k}$. This is passed through a softmax function to produce the attention weights matrix Z (pink, 3x3).

The Transformer



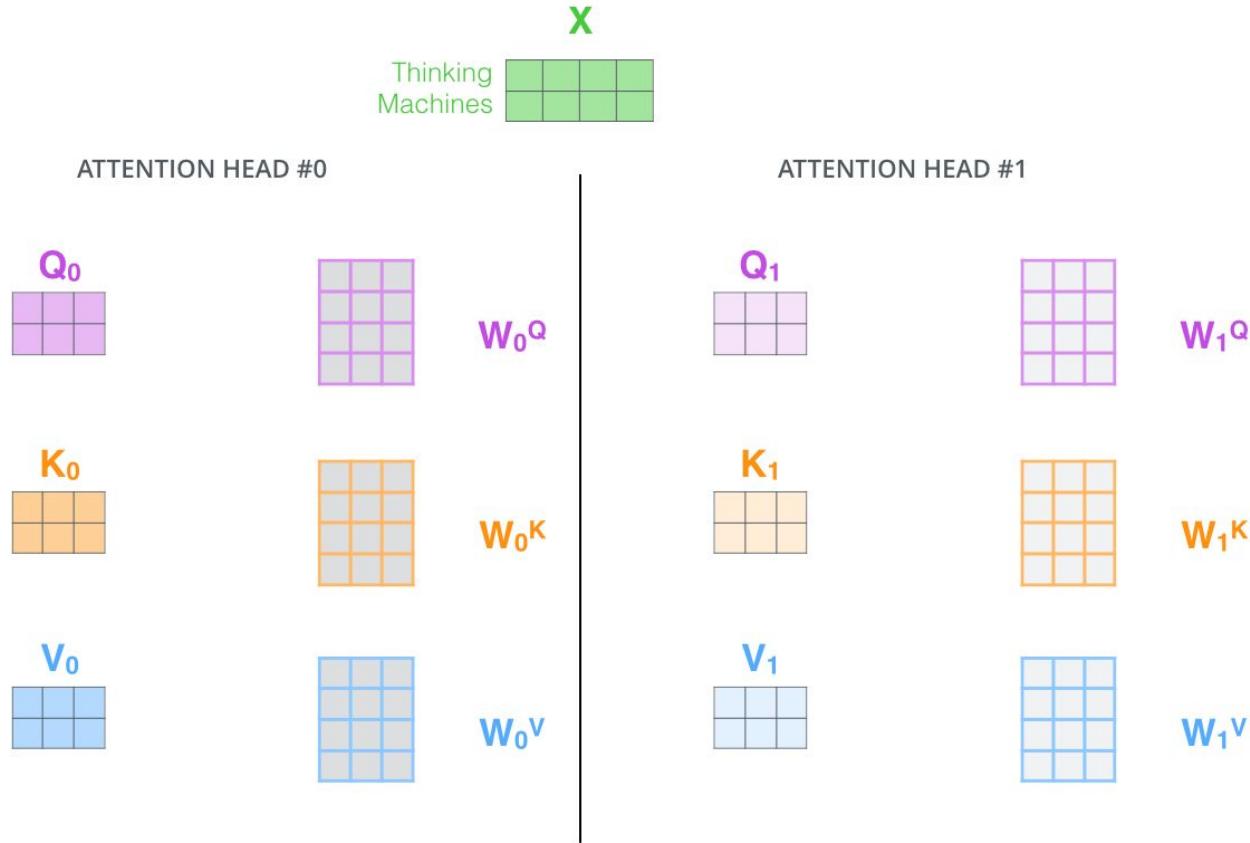
Multi-head Attention!

Multi-head Attention for Multiple K, Q, V Representations



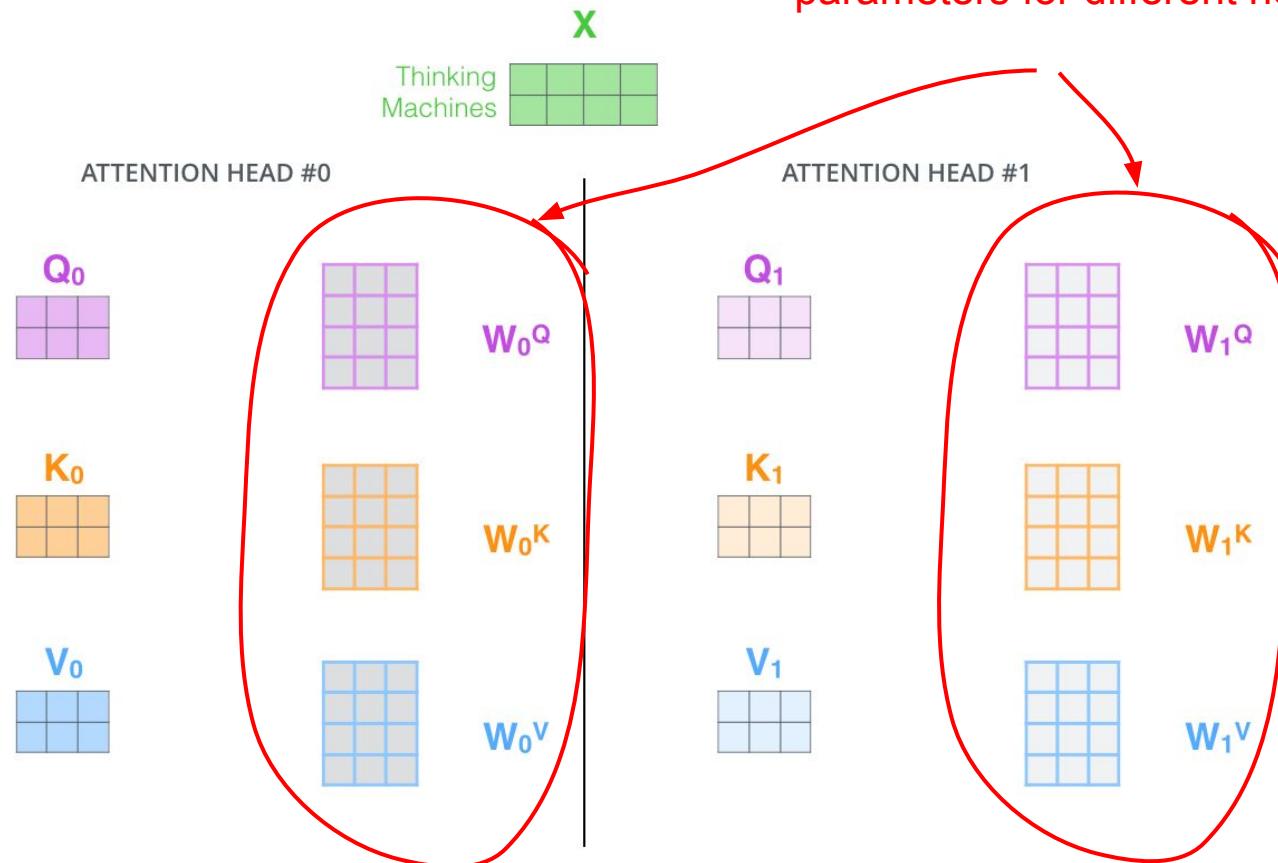
Model Architecture

2-head Self-Attention

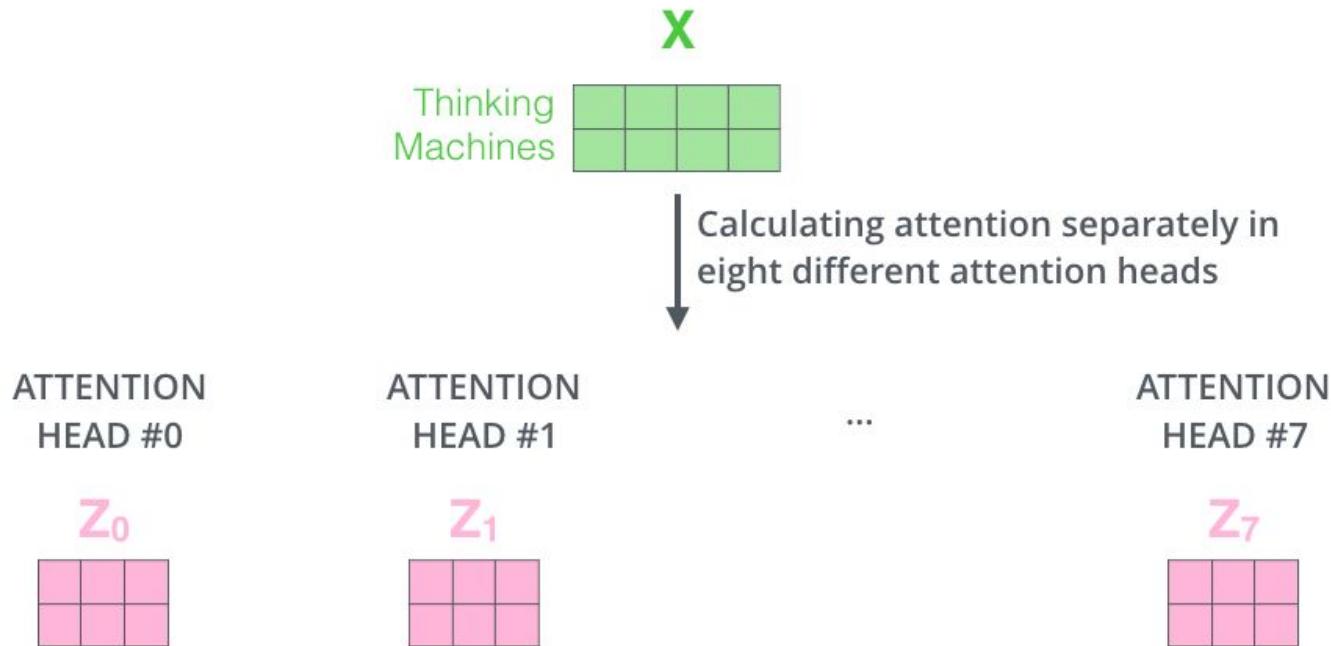


2-head Self-Attention

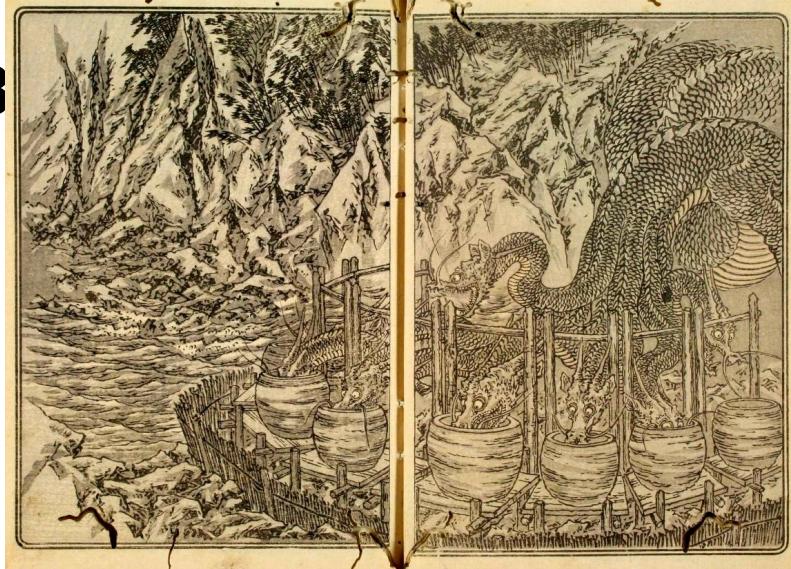
[Important!] Different trainable parameters for different heads



8-head Self-Attention = 8 Attention Vectors



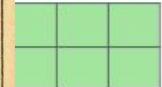
8



Yamata no Orochi (やまたの オロチ 八岐大蛇)
A legendary 8-head 8-tail Japanese dragon
https://en.wikipedia.org/wiki/Yamata_no_Orochi

Attention Vectors

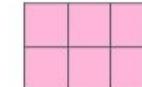
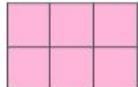
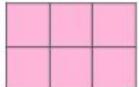
X



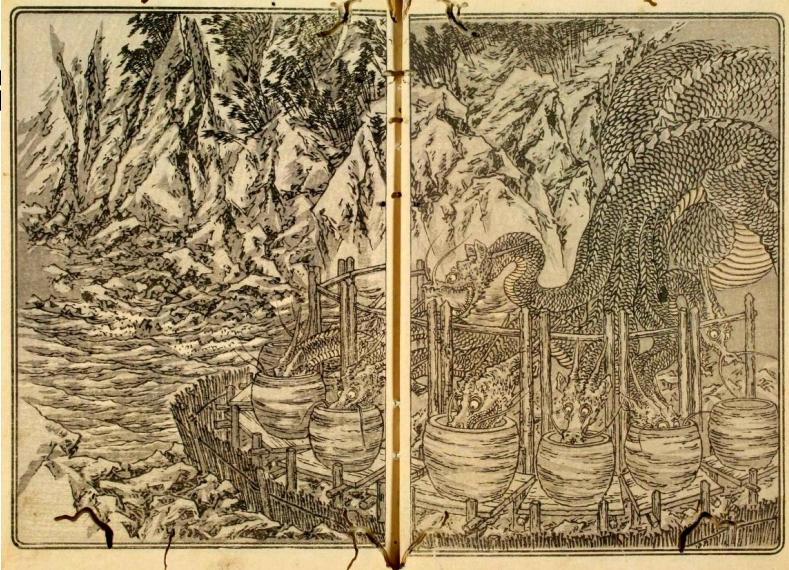
Calculating attention separately in
eight different attention heads

...

ATTENTION
HEAD #7

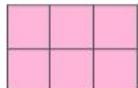
 Z_7  Z_0  Z_1 

8

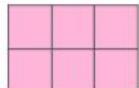


Yamata no Orochi (やまとのおろち 八岐大蛇)
A legendary 8-head 8-tail Japanese dragon
https://en.wikipedia.org/wiki/Yamata_no_Orochi

Z_0



Z_1



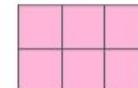
A



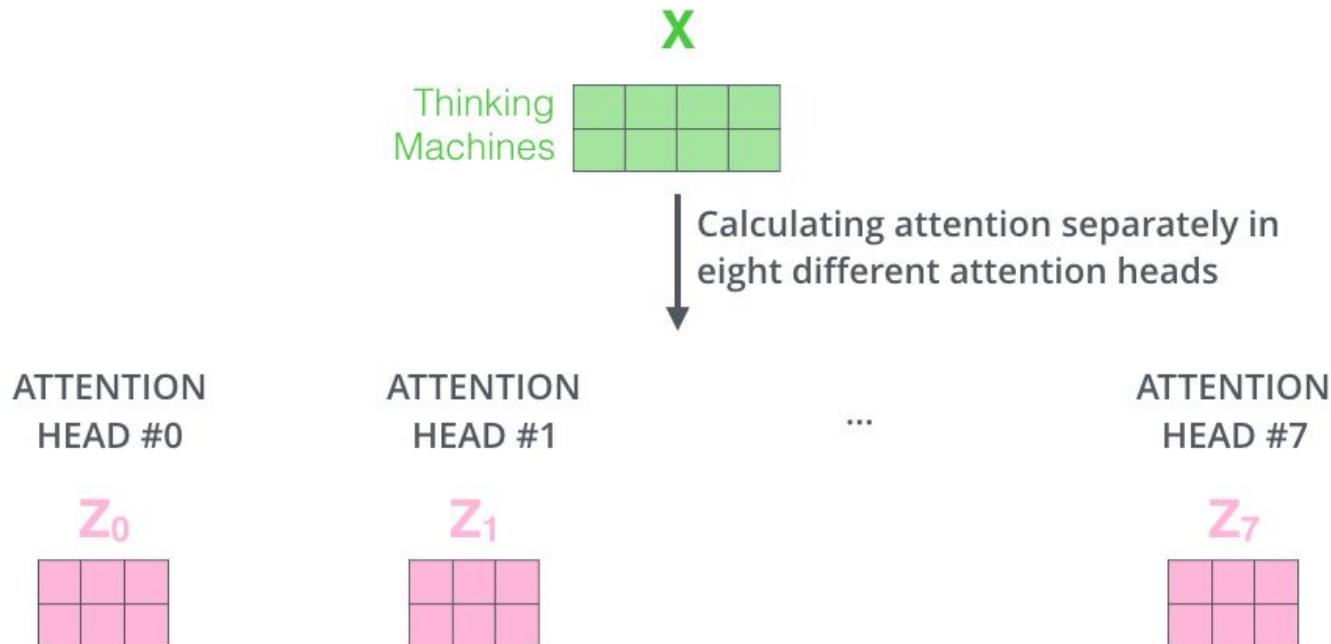
cf. Hydra (Greek myth)
Multi-head serpent

https://en.wikipedia.org/wiki/Lernaean_Hydra

Z_7



8-head Self-Attention = 8 Attention Vectors



Q. How do we aggregate them to get a single z_1 ?

Concatenation & Linear-Transformation

1) Concatenate all the attention heads



($\text{num_tokens}, \text{attention_dim} * \text{num_heads}$)

2) Multiply with a weight matrix W^o that was trained jointly with the model

\times



W^o

3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN

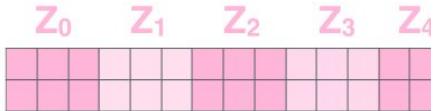
$$= \begin{matrix} Z \\ \hline \end{matrix} \quad \begin{matrix} \text{---} \\ \hline \end{matrix} \quad \begin{matrix} \text{---} \\ \hline \end{matrix}$$

($\text{num_tokens}, \text{attention_dim}$)

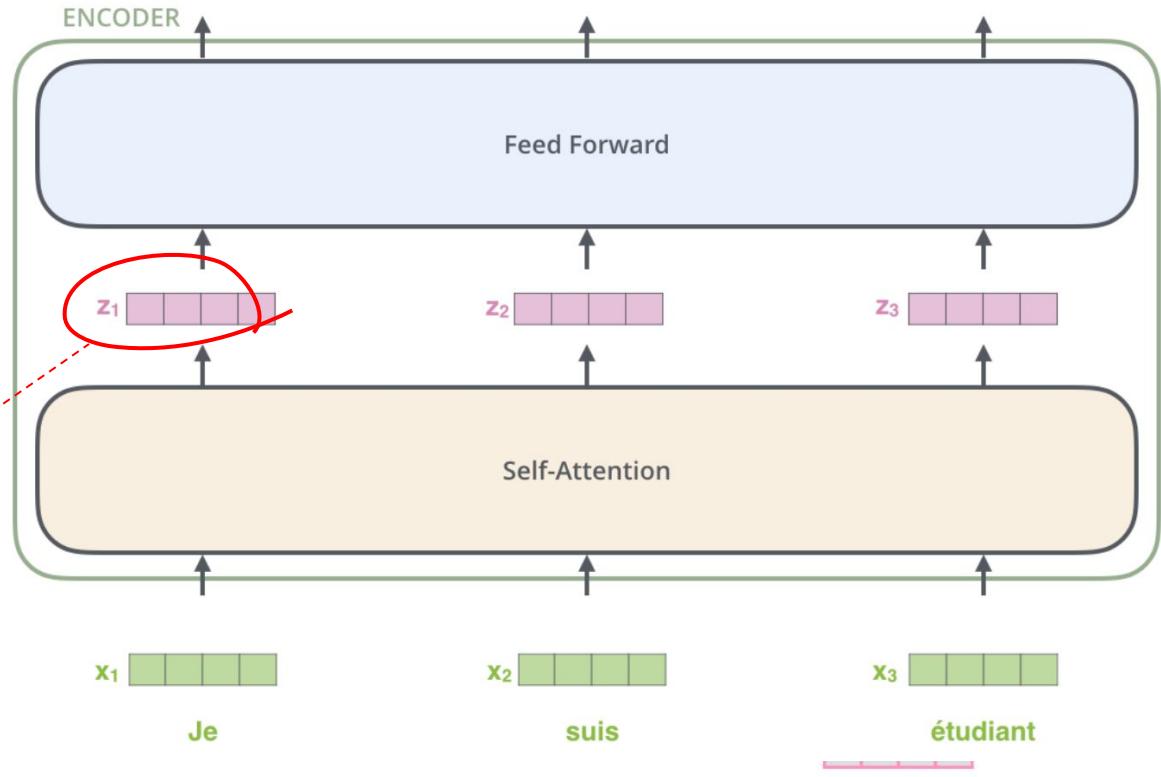
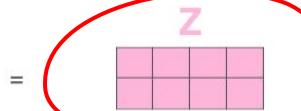
($\text{num_heads} * \text{attention_dim}, \text{attention_dim}$)

Concatenation & Linear-Transformation

1) Concatenate all the attention heads

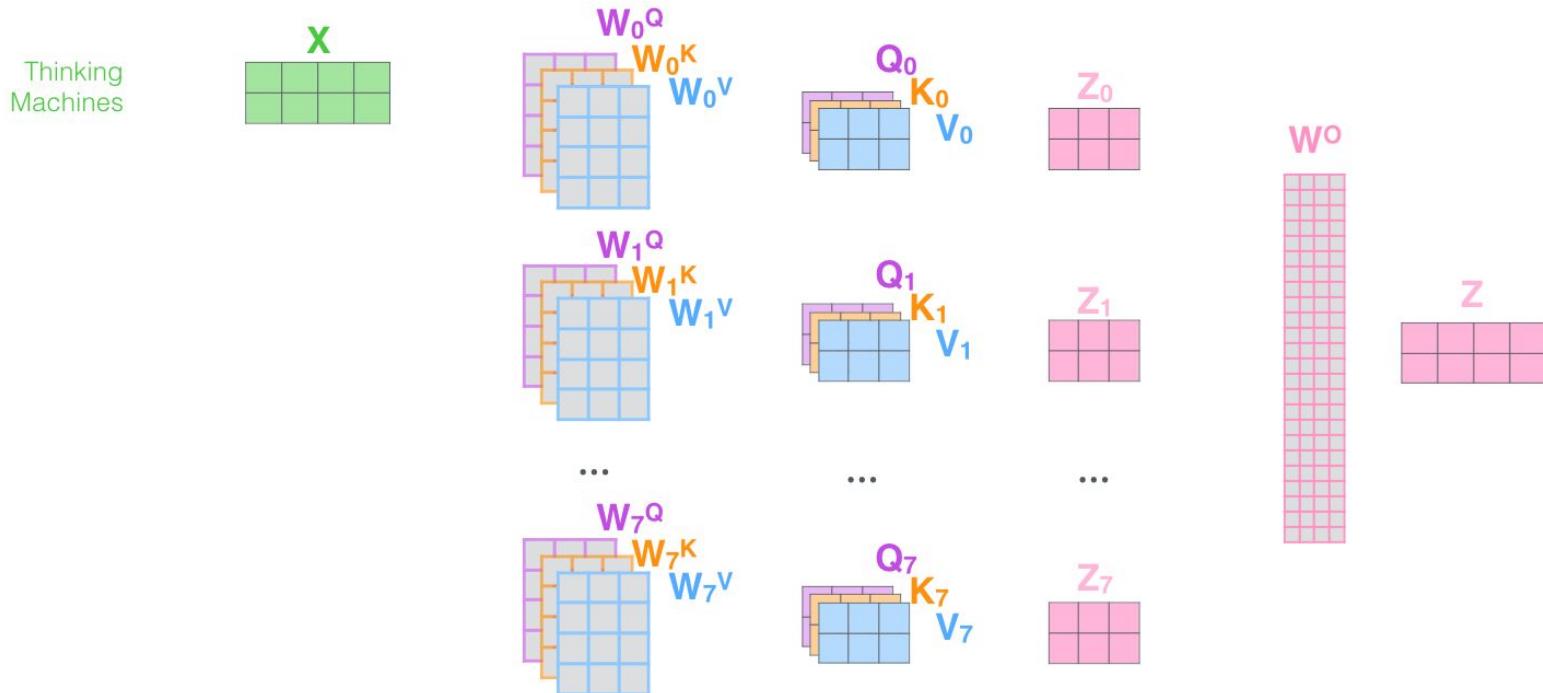


3) The result would be the Z matrix from all the attention heads. We can

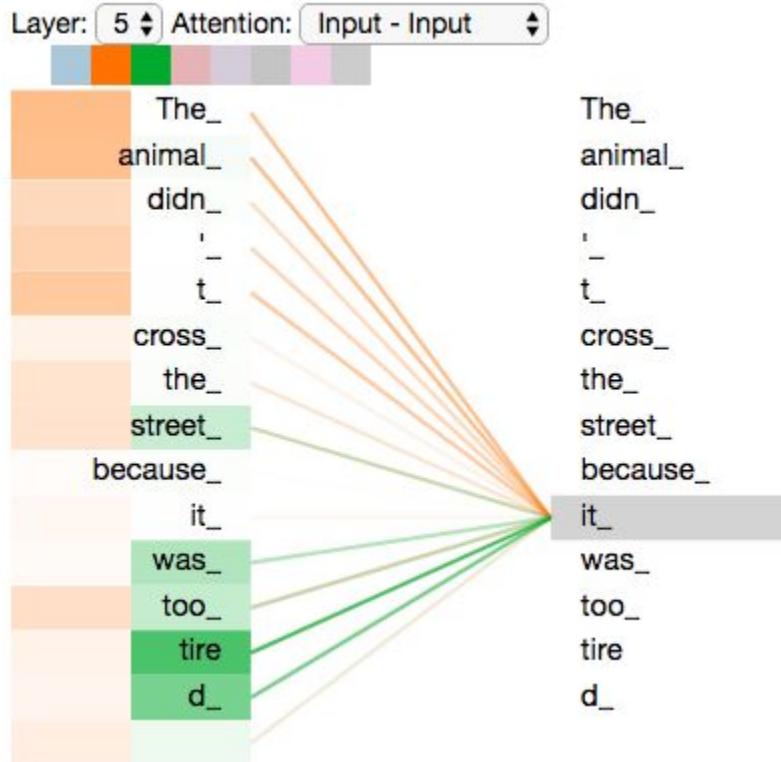


Recap: How Multi-head Self-Attention Works

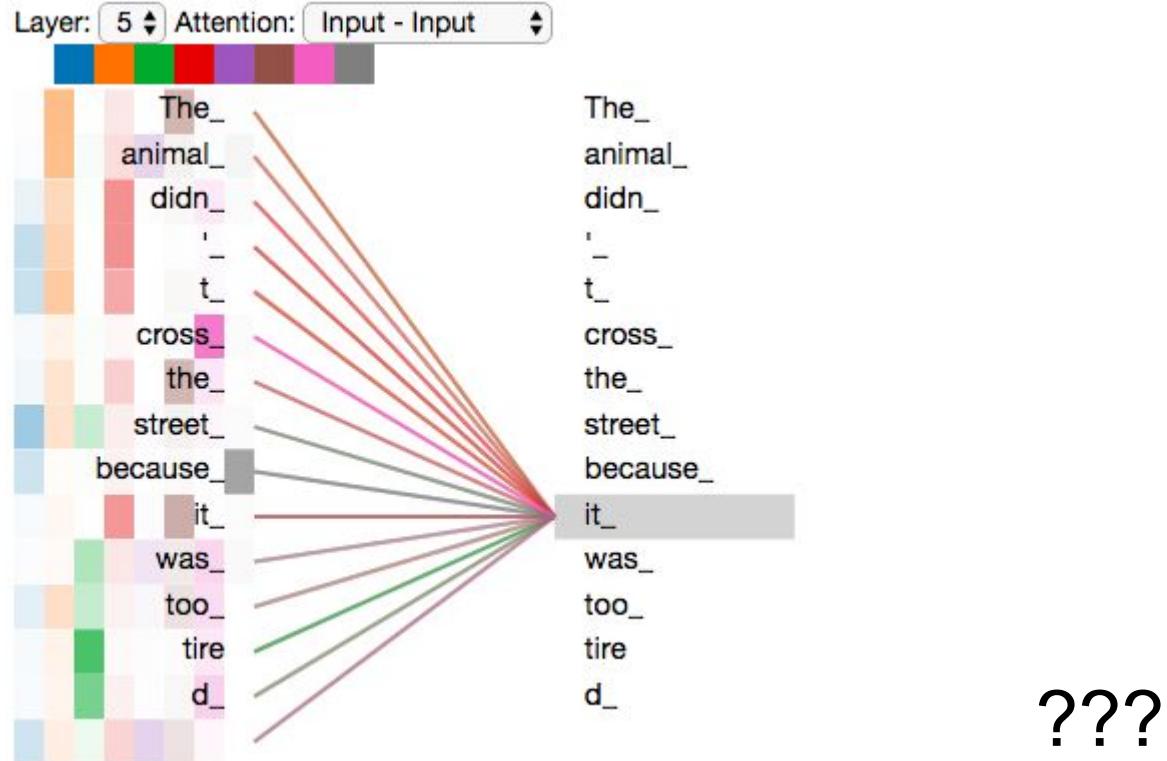
- 1) This is our input sentence* X
- 2) We embed each word* R
- 3) Split into 8 heads. We multiply X or R with weight matrices W_0^Q, W_0^K, W_0^V
- 4) Calculate attention using the resulting $Q/K/V$ matrices
- 5) Concatenate the resulting Z matrices, then multiply with weight matrix W^o to produce the output of the layer



In Theory: Different Heads Capture Different Types of Contextual Information



In Reality: Multiple Attention Heads are Hard to Interpret



Attention is not Explanation

Sarthak Jain

Northeastern University

jain.sar@husky.neu.edu

Byron C. Wallace

Northeastern University

b.wallace@northeastern.edu

Attention is not Explanation

Sarthak Jain

Northeastern University

jain.sar@husky.neu.edu

Attention is not not Explanation

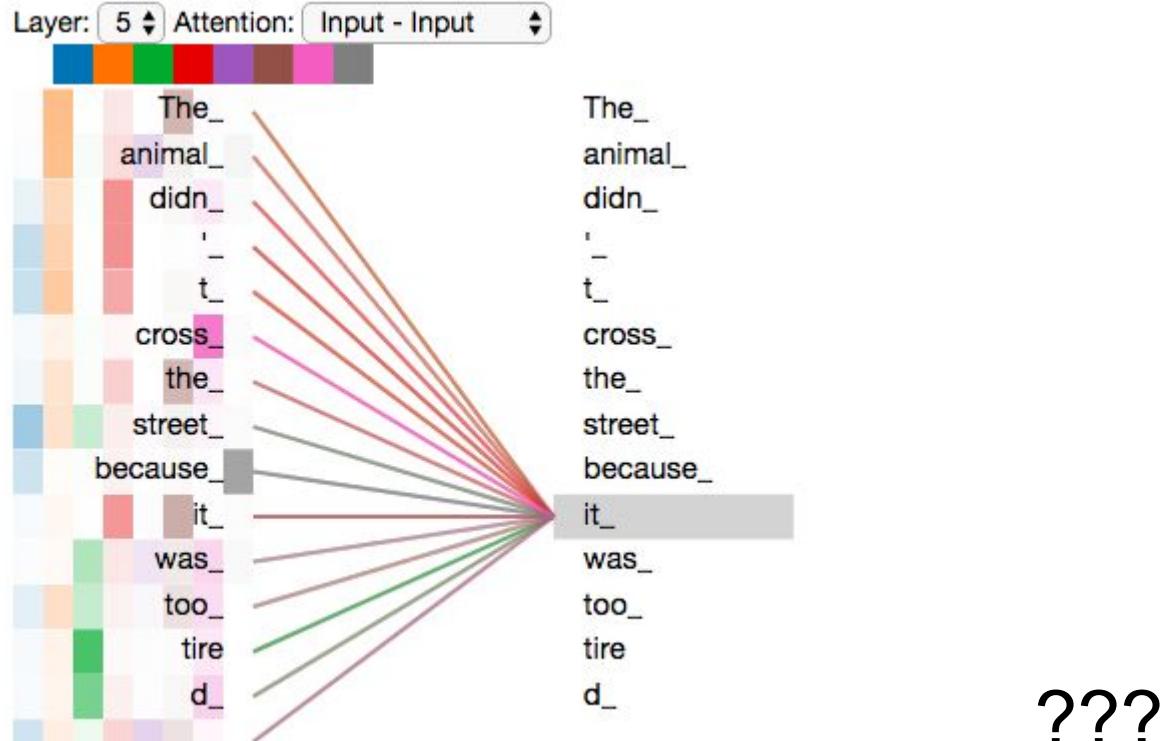
Sarah Wiegreffe*

School of Interactive Computing
Georgia Institute of Technology
saw@gatech.edu

Yuval Pinter*

School of Interactive Computing
Georgia Institute of Technology
uvp@gatech.edu

In Reality: Multiple Attention Heads are Hard to Interpret



Q. Do we really need so many attention heads?

Q. Do we really need so many attention heads?

- Short answer is **No**

ACL '19

Analyzing Multi-Head Self-Attention: Specialized Heads Do the Heavy Lifting, the Rest Can Be Pruned

Elena Voita^{1,2} David Talbot¹ Fedor Moiseev^{1,5} Rico Sennrich^{3,4} Ivan Titov^{3,2}

¹Yandex, Russia ²University of Amsterdam, Netherlands

³University of Edinburgh, Scotland ⁴University of Zurich, Switzerland

⁵Moscow Institute of Physics and Technology, Russia

{lena-voita, talbot, femoiseev}@yandex-team.ru
rico.sennrich@ed.ac.uk ititov@inf.ed.ac.uk

NeurIPS '19

Are Sixteen Heads Really Better than One?

Paul Michel
Language Technologies Institute
Carnegie Mellon University
Pittsburgh, PA
pmichel1@cs.cmu.edu

Omer Levy
Facebook Artificial Intelligence Research
Seattle, WA
omerlevy@fb.com

Graham Neubig
Language Technologies Institute
Carnegie Mellon University
Pittsburgh, PA
gneubig@cs.cmu.edu

BlackBoxNLP '19@ACL '19

What Does BERT Look At? An Analysis of BERT's Attention

Kevin Clark[†] Urvashi Khandelwal[†] Omer Levy[‡] Christopher D. Manning[†]

[†]Computer Science Department, Stanford University

[‡]Facebook AI Research

{kevclark, urvashik, manning}@cs.stanford.edu
omerlevy@fb.com

the simple weighted average. In this paper we make the surprising observation that even if models have been trained using multiple heads, in practice, a large percentage of attention heads can be removed at test time without significantly impacting performance. In fact, some layers can even be reduced to a single head. We further examine greedy algorithms for pruning down models, and

(My opinion)

We could tune the number of heads

But, overshooting may not be that harmful

Note: Conventional Attention Techniques for RNN

- To obtain **context vector c_t** as auxiliary information for the decoder

Similarity

$$c_t = \sum_{i=1} \alpha_{t,i} h_i$$

$$\alpha_{t,i} = \text{align}(y_t, x_i)$$

$$= \frac{\exp(\text{score}(s_{t-1}, h_i))}{\sum_{i'=1}^n \exp(\text{score}(s_{t-1}, h_{i'}))}$$

; Context vector for output y_t

; How well two words y_t and x_i are aligned.

; Softmax of some predefined alignment score..

Note: Conventional Attention Techniques for RNN

- To obtain **context vector c_t** as auxiliary information for the decoder

The diagram illustrates the computation of the context vector c_t and attention weights $\alpha_{t,i}$. It shows the formula for c_t and the definition of $\alpha_{t,i}$, with various components labeled:

$$c_t = \sum_{i=1} \alpha_{t,i} h_i$$

Annotations for the formula:

- "Similarity" points to $\alpha_{t,i}$.
- "Value" points to h_i .
- A speech bubble says "= hidden state of the encoder".
- "Query" points to s_{t-1} .
- "Key" points to h_i .

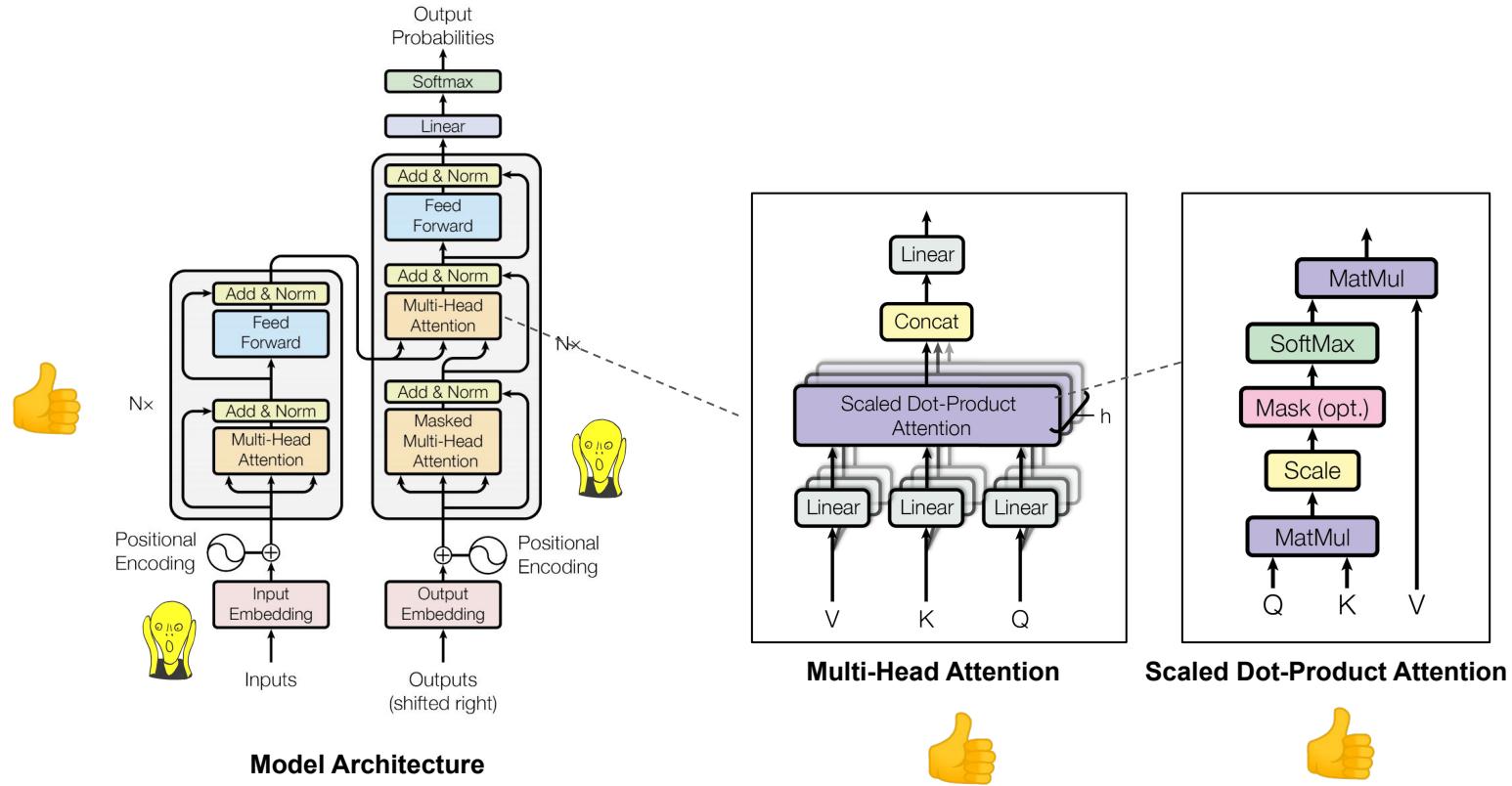
Annotations for $\alpha_{t,i}$:

- ; Context vector for output y_t
- ; How well two words y_t and x_i are aligned.
- ; Softmax of some predefined alignment score..

Note: Various Score Functions

Name	Query	Key	Alignment score function	Citation
Content-base attention	s_t	\mathbf{h}_i	$\text{score}(s_t, \mathbf{h}_i) = \text{cosine}[s_t, \mathbf{h}_i]$	Graves2014
Additive(*)			$\text{score}(s_t, \mathbf{h}_i) = \mathbf{v}_a^\top \tanh(\mathbf{W}_a[s_t; \mathbf{h}_i])$	Bahdanau2015
Location-Base			$\alpha_{t,i} = \text{softmax}(\mathbf{W}_a s_t)$ Note: This simplifies the softmax alignment to only depend on the target position.	Luong2015
General			$\text{score}(s_t, \mathbf{h}_i) = s_t^\top \mathbf{W}_a \mathbf{h}_i$ where \mathbf{W}_a is a trainable weight matrix in the attention layer.	Luong2015
Dot-Product			$\text{score}(s_t, \mathbf{h}_i) = s_t^\top \mathbf{h}_i$	Luong2015
Scaled Dot-Product(^)			$\text{score}(s_t, \mathbf{h}_i) = \frac{s_t^\top \mathbf{h}_i}{\sqrt{n}}$ Note: very similar to the dot-product attention except for a scaling factor; where n is the dimension of the source hidden state.	Vaswani2017

The Transformer



Remaining Topics

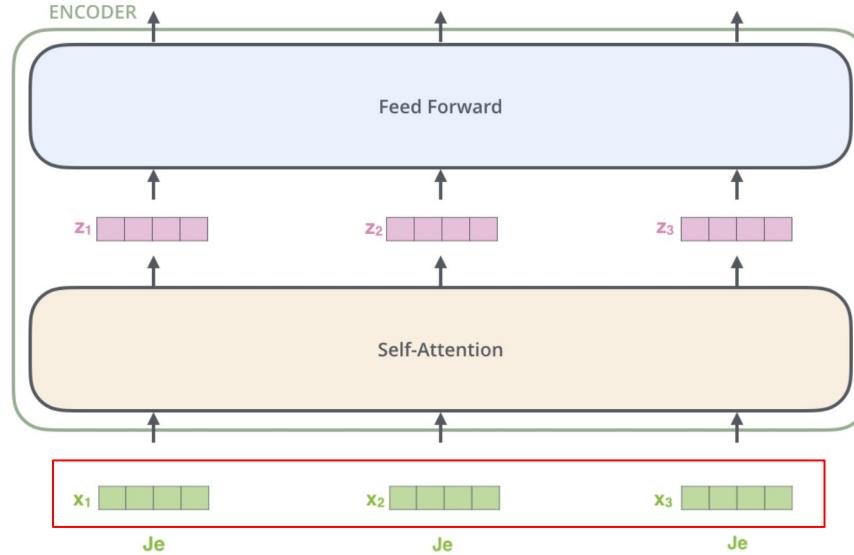
- Positional Encoding
- Residual Connection & Layer normalization
- Transformer Decoder

Remaining Topics

- **Positional encoding**
- **Residual Connection & Layer normalization**
- Transformer Decoder

Why Positional Encoding?

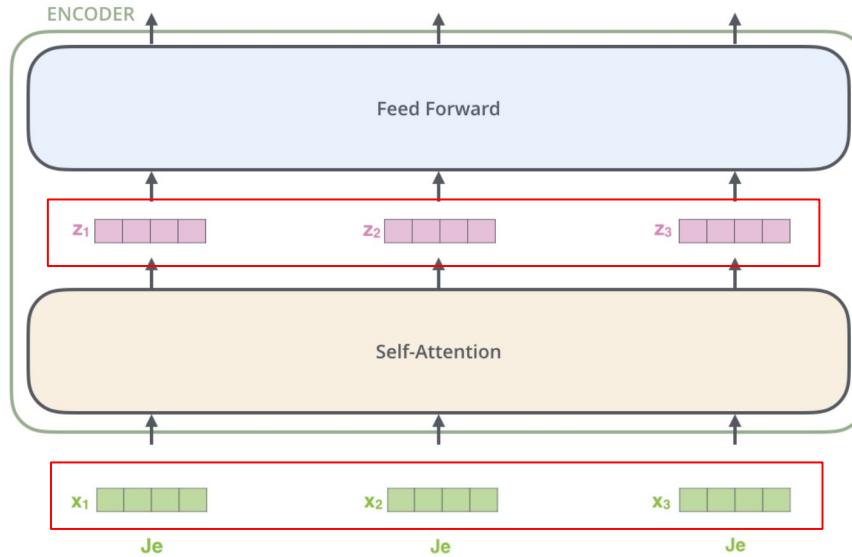
- Transformer Encoder (i.e., self-attention) does **NOT consider any positional information**



If the input token embeddings were the same

Why Positional Encoding?

- Transformer Encoder (i.e., self-attention) does **NOT consider any positional information**
 - In one sense, it can be considered as “**contextualized**” bag-of-words representations

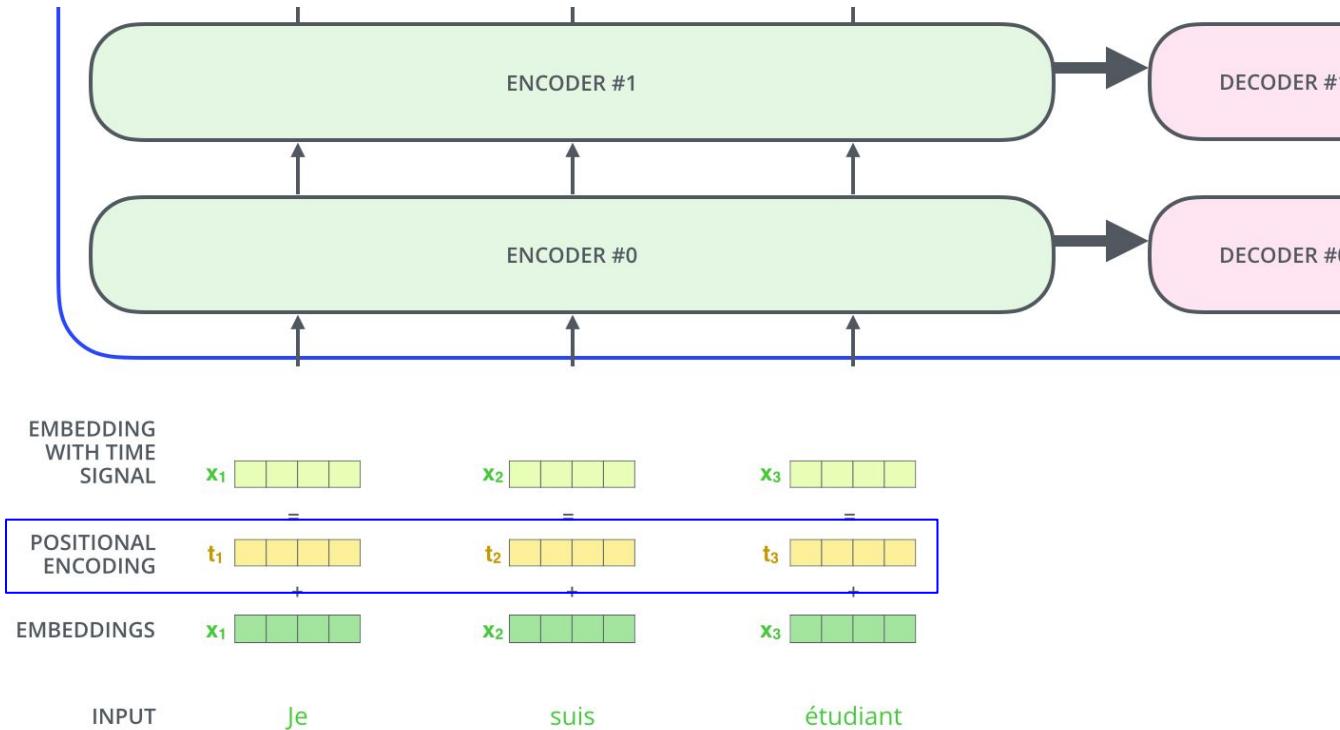


The corresponding attention vectors would become identical

If the input token embeddings were the same

Adding Positional Encoding to Make Position-aware

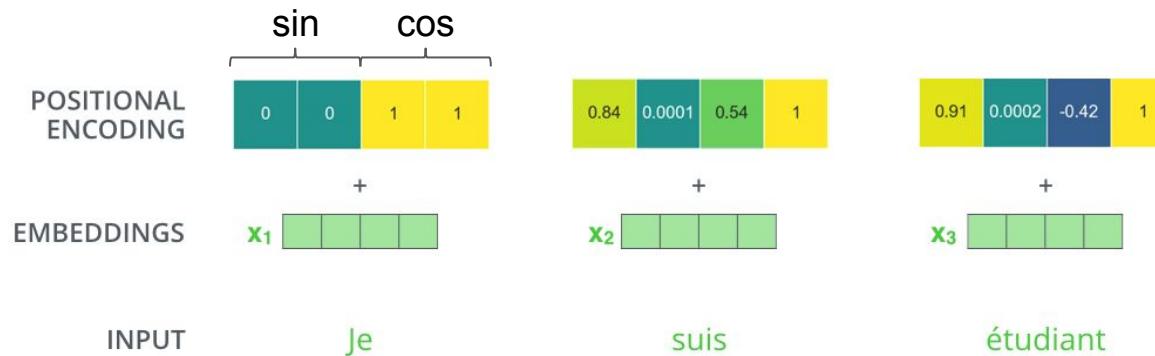
- Just adding **position information** to input word embeddings



Positional Encoding in Transformer

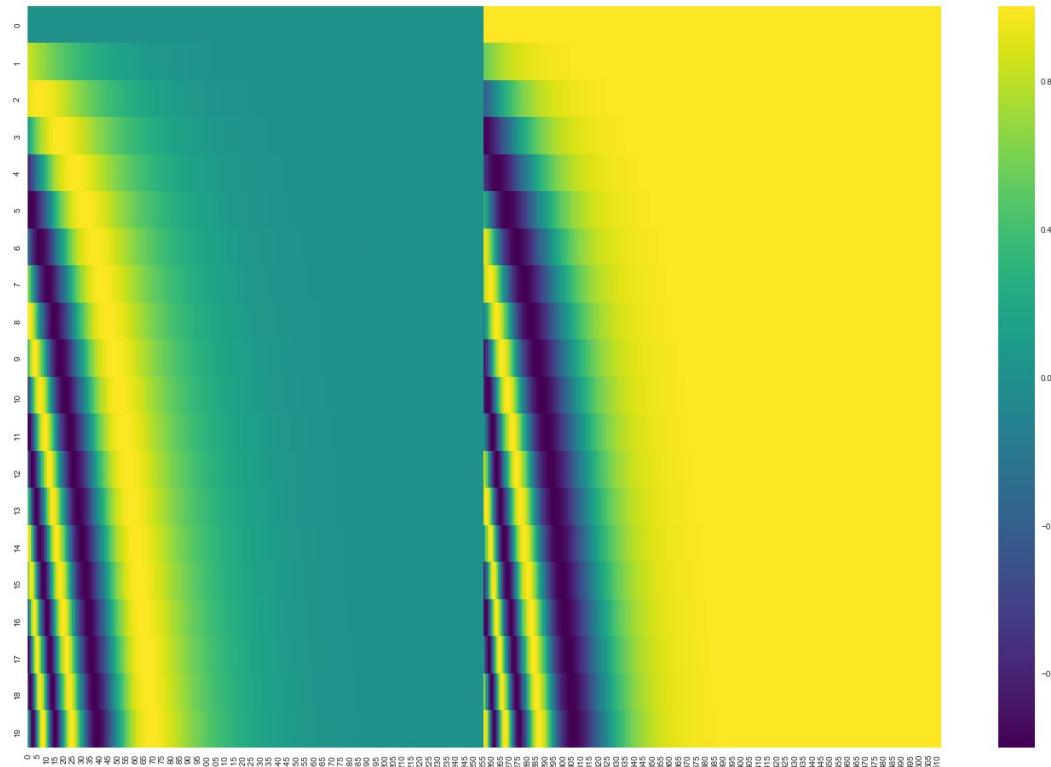
$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$



How Actual Positional Encoding Values Look Like

20 tokens & embedding size of 512



Input Embedding Layer: Pytorch Example

```
class TransformerModel(nn.Module):

    def __init__(self, ntoken: int, d_model: int, nhead: int, d_hid: int,
                 nlayers: int, dropout: float = 0.5):
        super().__init__()
        self.model_type = 'Transformer'
        self.pos_encoder = PositionalEncoding(d_model, dropout)
        encoder_layers = TransformerEncoderLayer(d_model, nhead, d_hid, dropout)
        self.transformer_encoder = TransformerEncoder(encoder_layers, nlayers)
        self.encoder = nn.Embedding(ntoken, d_model)
        self.d_model = d_model
        self.decoder = nn.Linear(d_model, ntoken)

        self.init_weights()

    def forward(self, src: Tensor, src_mask: Tensor) -> Tensor:
        """
        Args:
            src: Tensor, shape [seq_len, batch_size]
            src_mask: Tensor, shape [seq_len, seq_len]

        Returns:
            output Tensor of shape [seq_len, batch_size, ntoken]
        """

        src = self.encoder(src) * math.sqrt(self.d_model)
        src = self.pos_encoder(src)
        output = self.transformer_encoder(src, src_mask)

        return output
```

```
class PositionalEncoding(nn.Module):

    def __init__(self, d_model: int, dropout: float = 0.1, max_len: int = 5000):
        super().__init__()
        self.dropout = nn.Dropout(p=dropout)

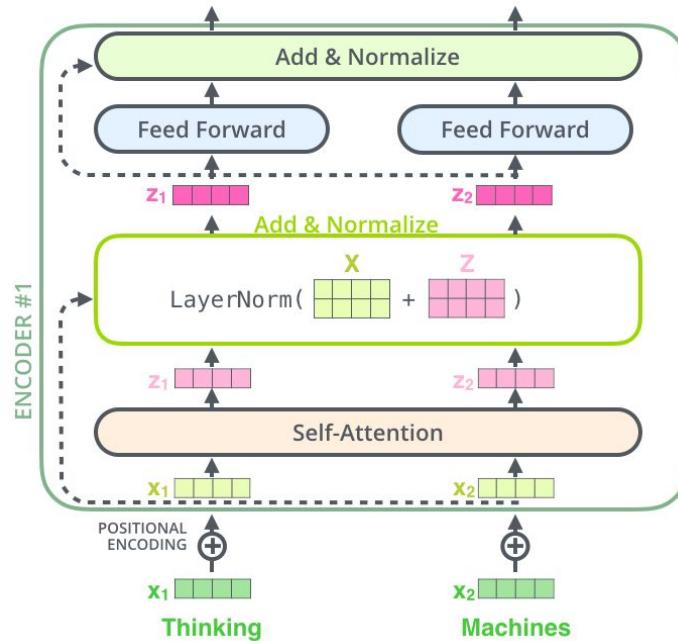
        position = torch.arange(max_len).unsqueeze(1)
        div_term = torch.exp(torch.arange(0, d_model, 2) * (-math.log(10000.0) / d_model))
        pe = torch.zeros(max_len, 1, d_model)
        pe[:, 0, 0::2] = torch.sin(position * div_term)
        pe[:, 0, 1::2] = torch.cos(position * div_term)
        self.register_buffer('pe', pe)

    def forward(self, x: Tensor) -> Tensor:
        """
        Args:
            x: Tensor, shape [seq_len, batch_size, embedding_dim]
        """

        x = x + self.pe[:x.size(0)]
        return self.dropout(x)
```

Residual Connection & Layer Normalization

- Apply **Residual Connection** (i.e., adding x_*) before **Layer Normalization** (*1)



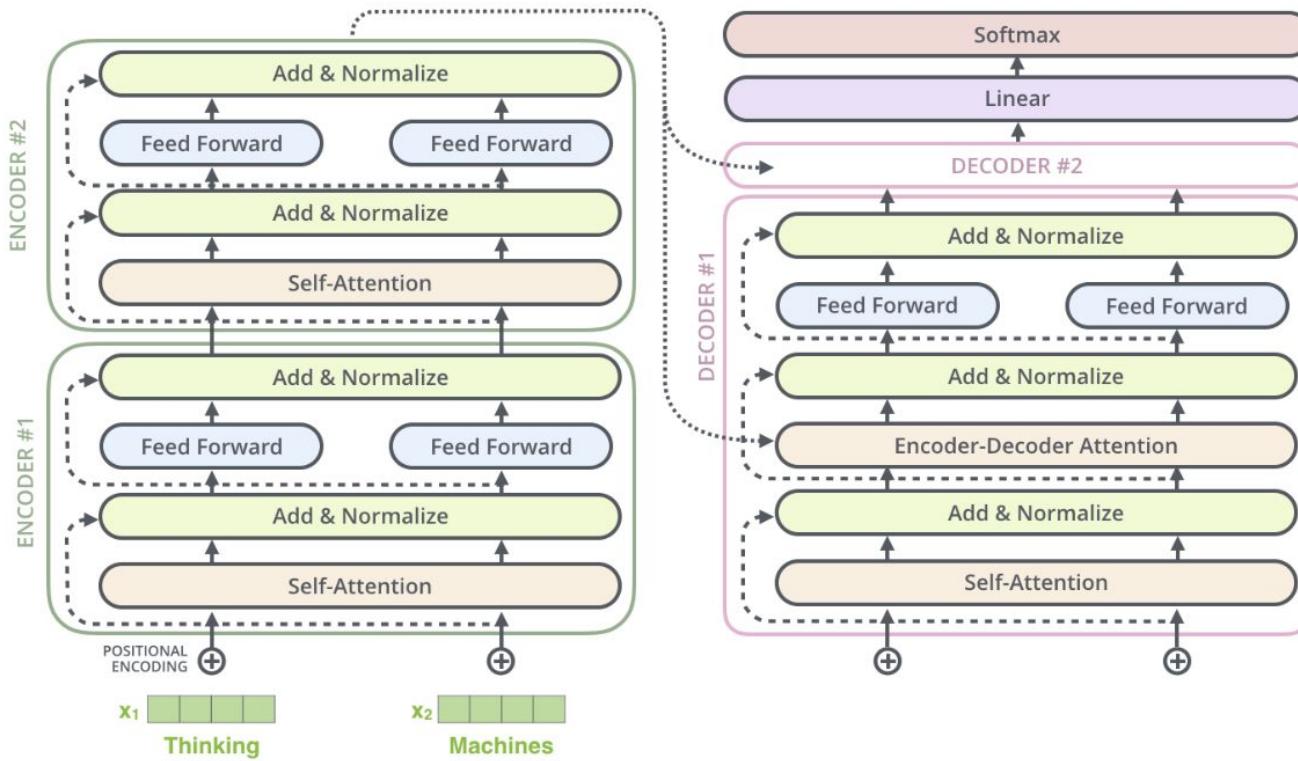
(*1) Layer normalization: Normalization over input tokens (cf. Batch Normalization)

Questions?

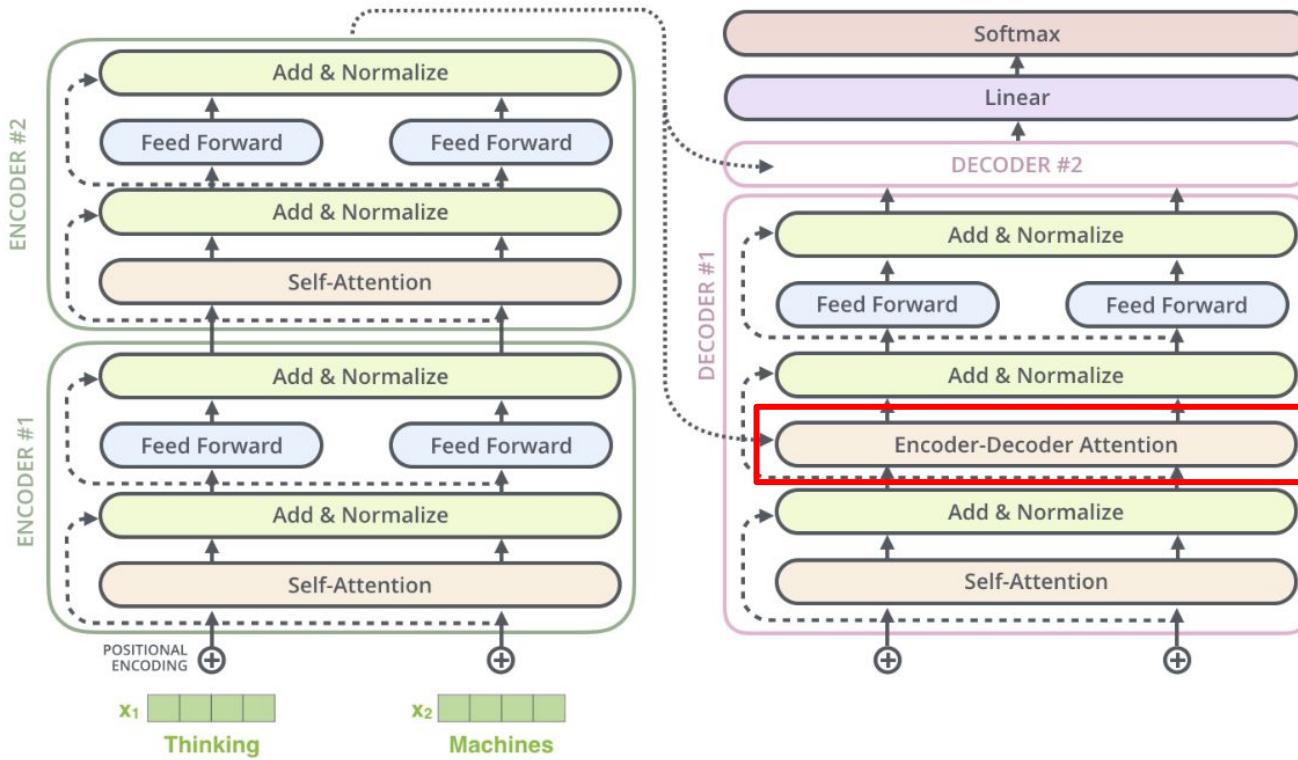
Remaining Topics

- Positional encoding
- Residual Connection & Layer normalization
- **Transformer Decoder**

Transformer (2 Encoders & 2 Decoders)

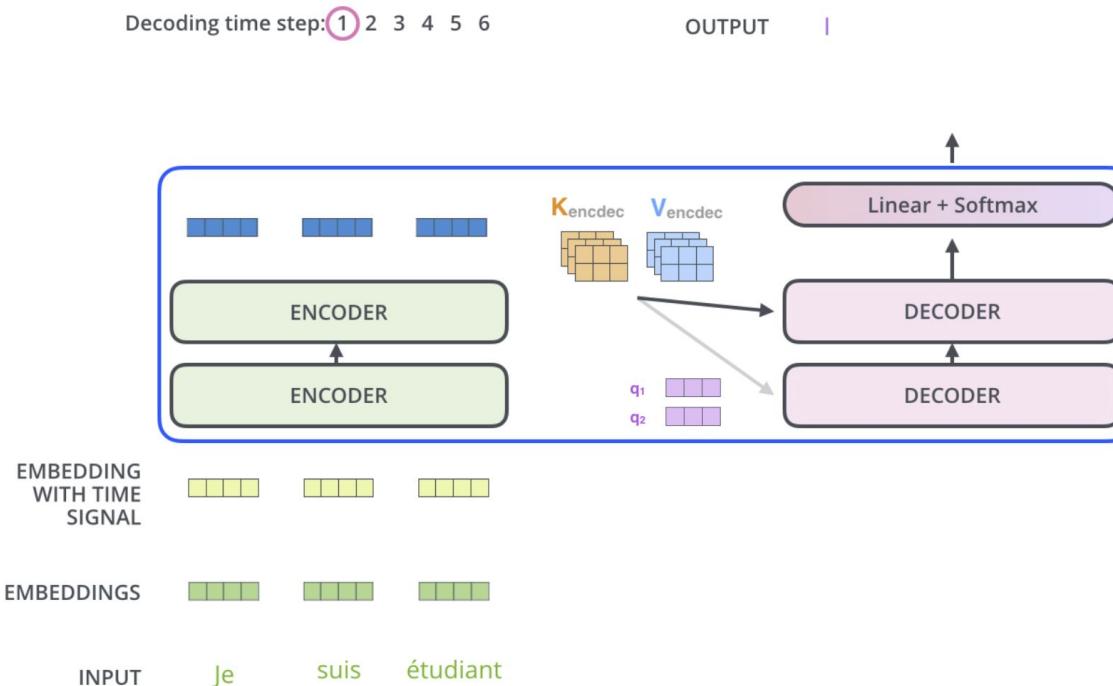


Transformer (2 Encoders & 2 Decoders)



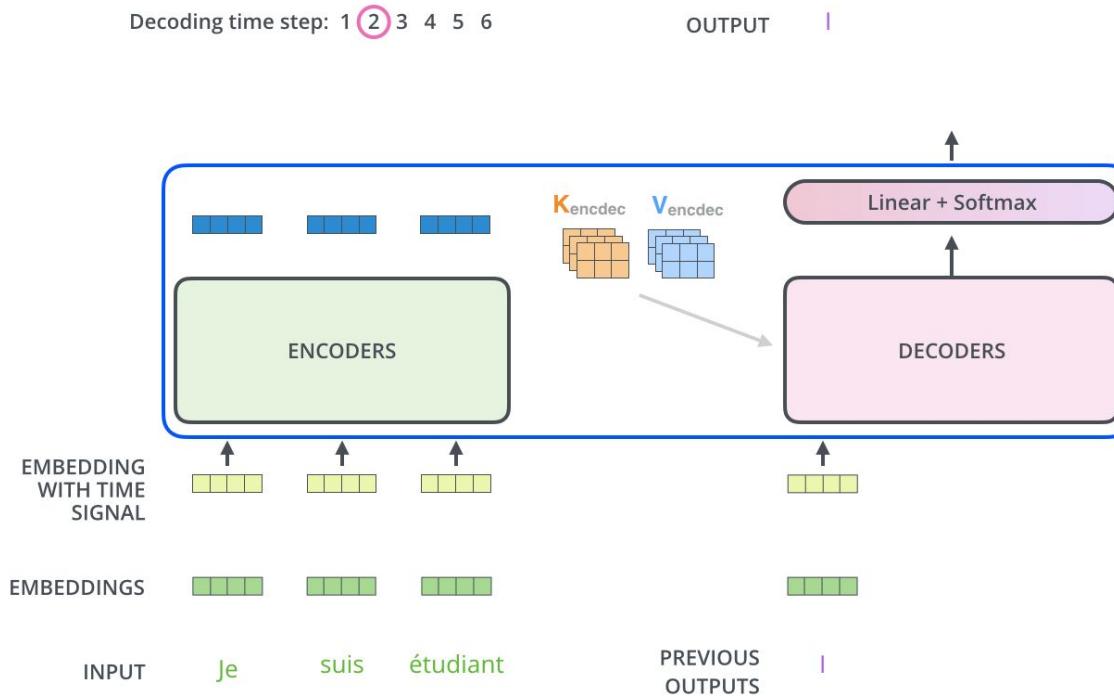
Decoder: Encoder-Decoder Attention

- The same “Query-Key-Value” style attention where **Key** and **Value vectors** are from Encoder and **Query vectors** are from Decoder

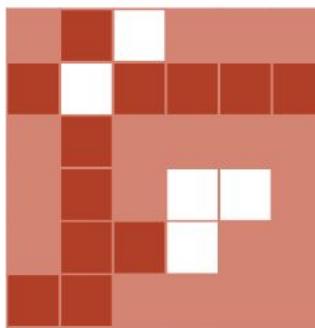


Decoder: “Masked” Self-Attention

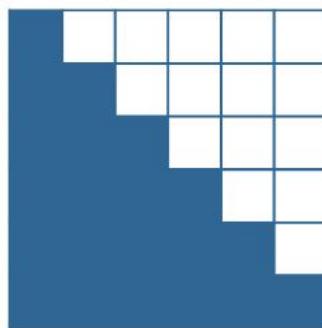
- The self-attention layer (in the Decoder) is **only allowed to attend to earlier positions** in the output sequence



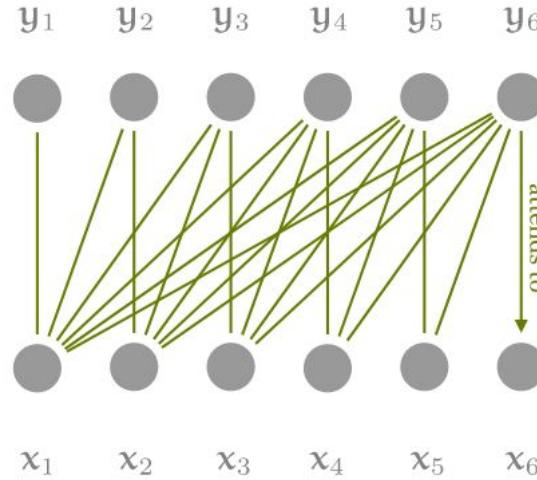
Masked Attention to Prevent “Previous” Step From Attending “Future” Steps



raw attention weights

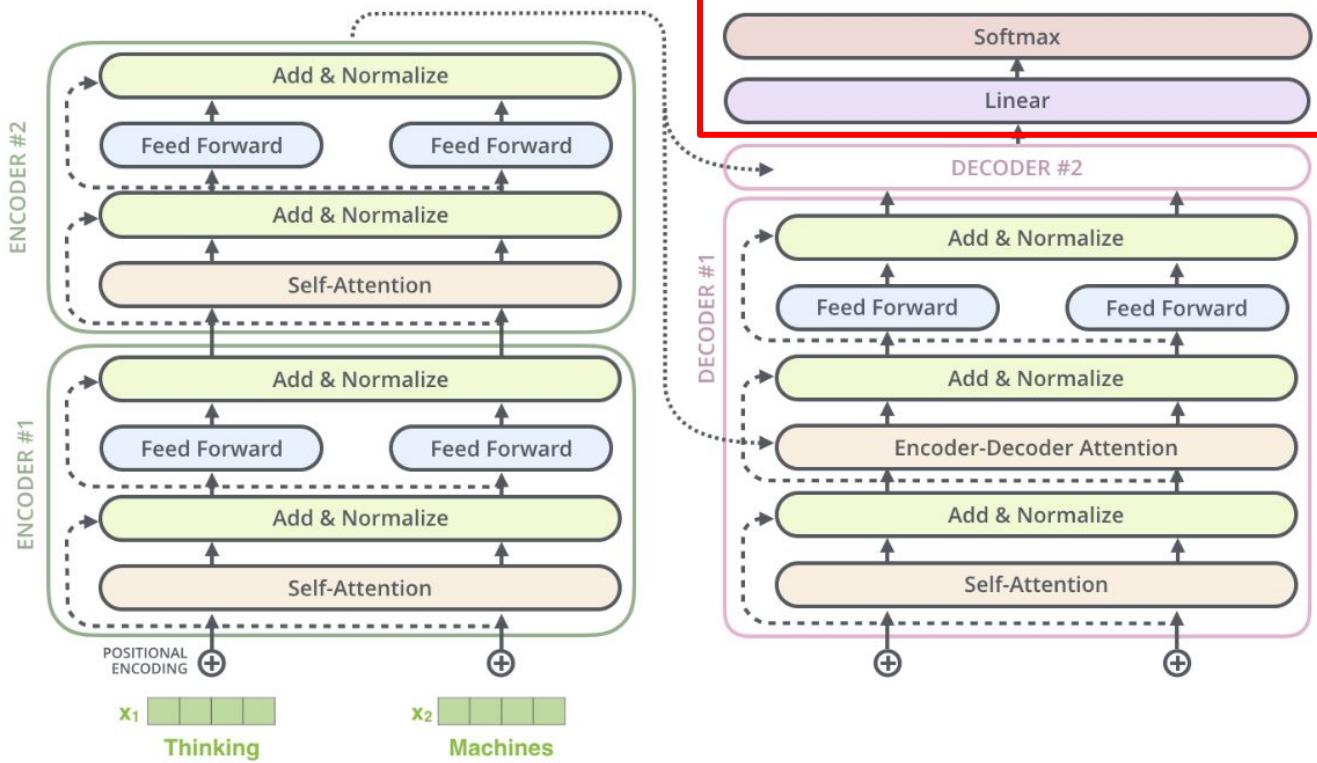


mask



Decoder

Transformer (2 Encoders & 2 Decoders)



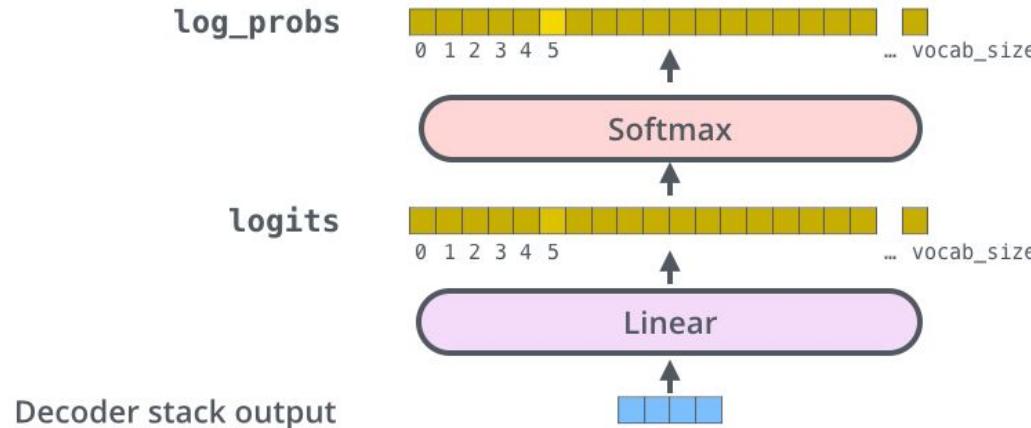
Output Layer

Which word in our vocabulary
is associated with this index?

am

Get the index of the cell
with the highest value
(`argmax`)

5



Output Layer: Training with Ground-truth Labels aka Teacher-forcing Training

Target Model Outputs

Output Vocabulary: a am I thanks student <eos>

position #1	0.0	0.0	1.0	0.0	0.0	0.0
position #2	0.0	1.0	0.0	0.0	0.0	0.0
position #3	1.0	0.0	0.0	0.0	0.0	0.0
position #4	0.0	0.0	0.0	0.0	1.0	0.0
position #5	0.0	0.0	0.0	0.0	0.0	1.0



Trained Model Outputs

Output Vocabulary: a am I thanks student <eos>

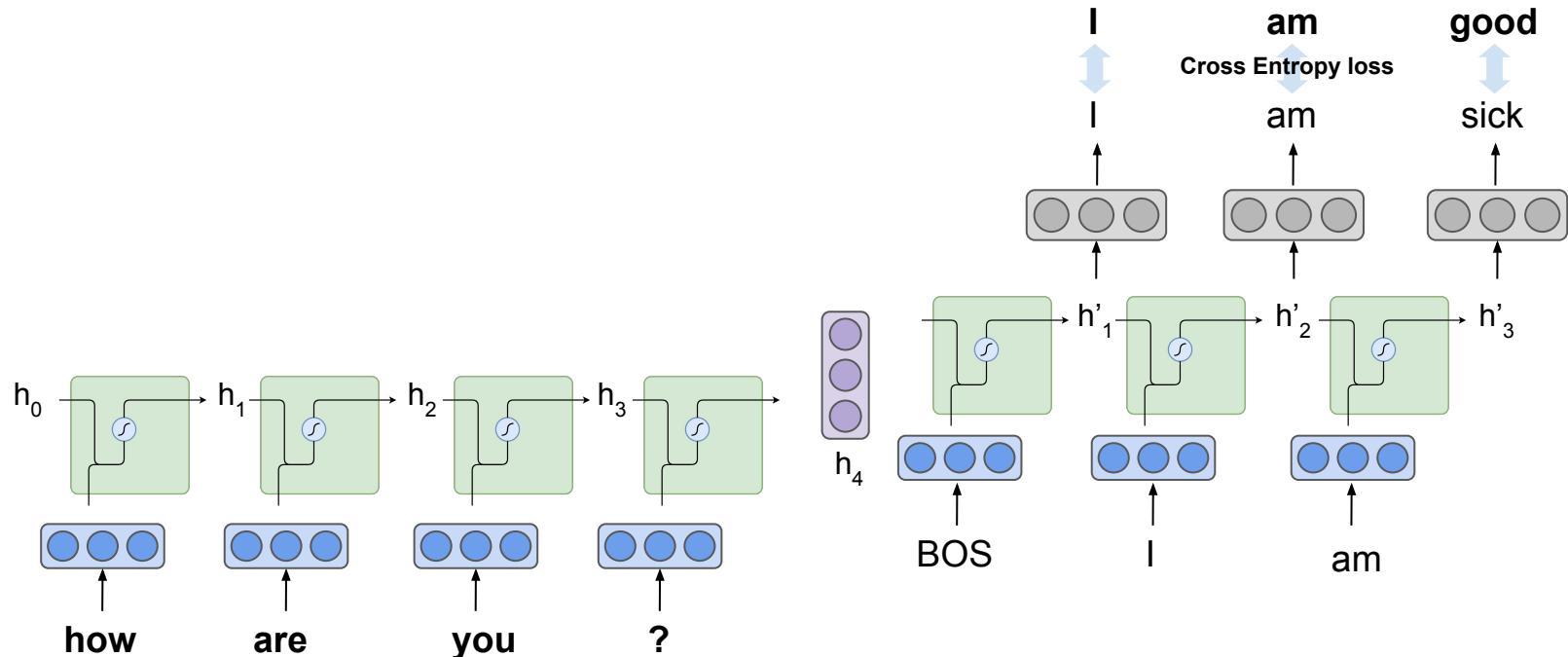
position #1	0.01	0.02	0.93	0.01	0.03	0.01
position #2	0.01	0.8	0.1	0.05	0.01	0.03
position #3	0.99	0.001	0.001	0.001	0.002	0.001
position #4	0.001	0.002	0.001	0.02	0.94	0.01
position #5	0.01	0.01	0.001	0.001	0.001	0.98



a am I thanks student <eos>

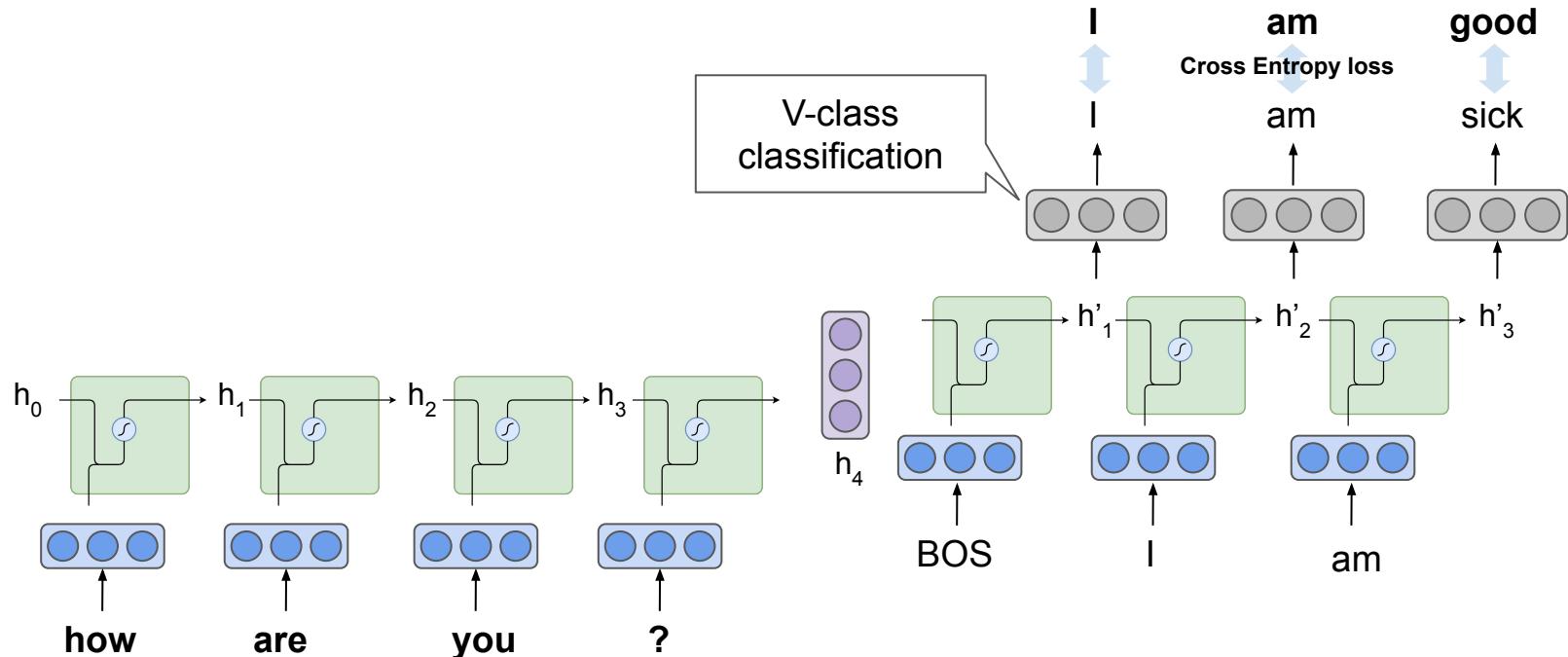
Training Encoder-Decoder Models: Training Data

- Training data: **Pairs of text sequences** (called a **parallel corpus**)
 - e.g., (“How are you?”, “I am good”)



Training Encoder-Decoder Models: “Teacher Forcing”

- Train the model to output **ground-truth outputs**



Note: Transformer Needs Decoding Algorithm

- Any decoding algorithm (e.g., Greedy/Beam Search) can be used
- Again, the Transformer is just an **encoder-decoder model**
- Although the Transformer does not have any recurrent architecture, the decoding procedure is still conducted in an iterative manner

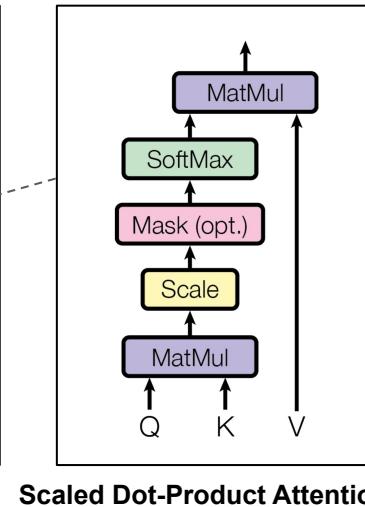
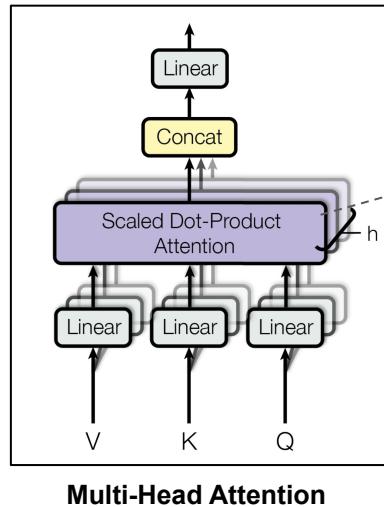
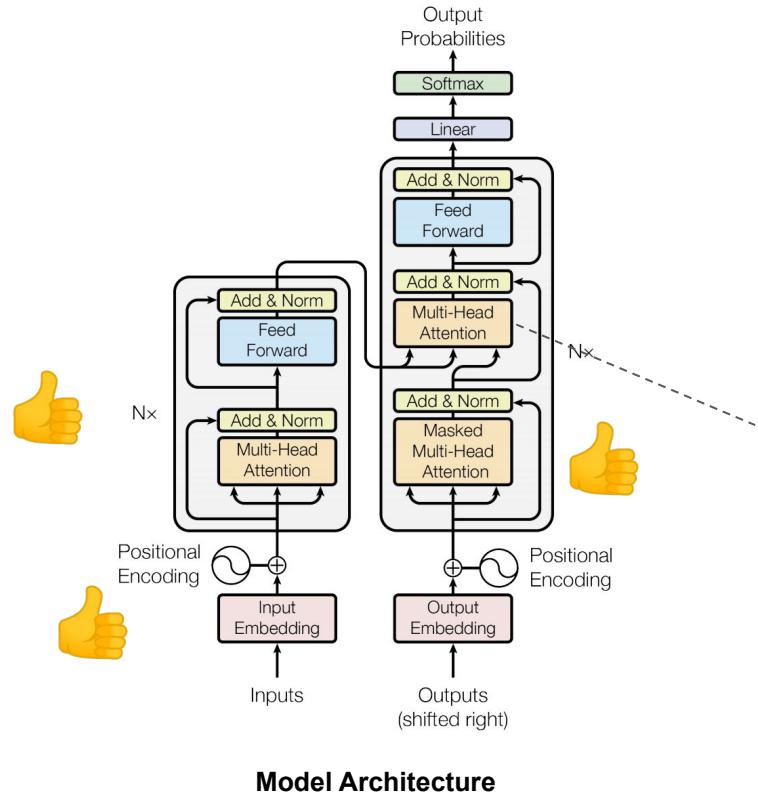
Parameter Initialization: Xavier Initialization

- Sampling initial parameter values from the Normal Distribution, where the standard deviation is **the number of parameters in the previous layer**
 - i.e., More parameters in the prev layer → Lower variance → Smaller absolute values

$$W_{i,j}^{[l]} = \mathcal{N} \left(0, \frac{1}{n^{[l-1]}} \right)$$

Detailed derivation: <https://cs230.stanford.edu/section/4/>

The Transformer



Key Takeaway

- Model Architecture
 - The Transformer is an **Encoder-Decoder model**
 - Transformer Encoder + Transformer Decoder (not explained yet)
 - Transformer Encoder = Stack of **Transformer Encoder Blocks**
 - Transformer Encoder Block = **Self-Attention Layer** + Feed Forward NN
- **Self-Attention Layer**
 - Token Embeddings → Query, Key, Value Vectors → Attention Vectors
 - Multi-head Attention
- Additional techniques
 - Positional Encoding
 - Residual Connection + Layer Normalization

Let's Take A Look At The Original Paper

Attention Is All You Need

Ashish Vaswani*

Google Brain

avaswani@google.com

Noam Shazeer*

Google Brain

noam@google.com

Niki Parmar*

Google Research

nikip@google.com

Jakob Uszkoreit*

Google Research

usz@google.com

Llion Jones*

Google Research

llion@google.com

Aidan N. Gomez* †

University of Toronto

aidan@cs.toronto.edu

Łukasz Kaiser*

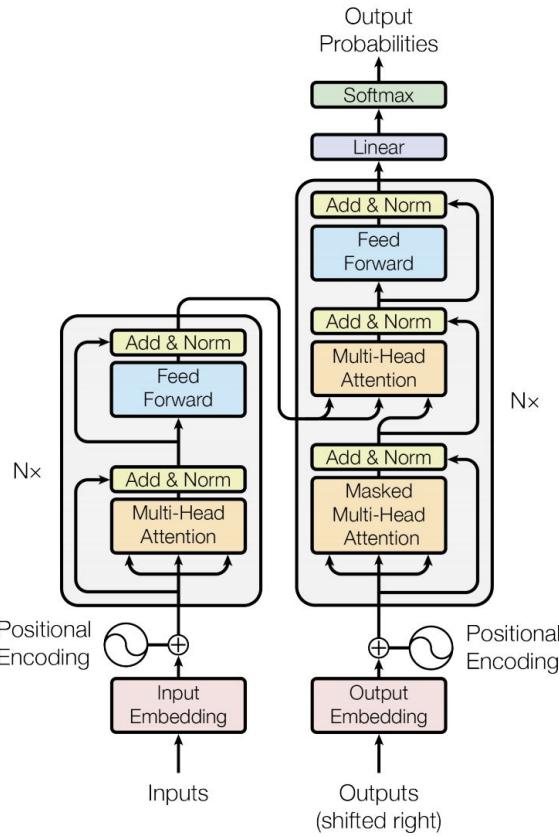
Google Brain

lukaszkaiser@google.com

Illia Polosukhin* ‡

illia.polosukhin@gmail.com

Model Descriptions in the Original Paper

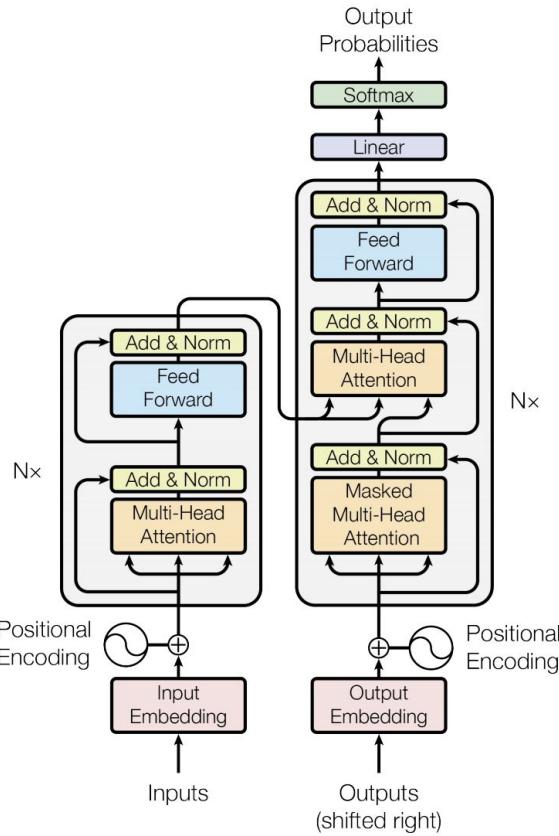


3.1 Encoder and Decoder Stacks

Encoder: The encoder is composed of a stack of $N = 6$ identical layers. Each layer has two sub-layers. The first is a multi-head self-attention mechanism, and the second is a simple, position-wise fully connected feed-forward network. We employ a residual connection [11] around each of the two sub-layers, followed by layer normalization [1]. That is, the output of each sub-layer is $\text{LayerNorm}(x + \text{Sublayer}(x))$, where $\text{Sublayer}(x)$ is the function implemented by the sub-layer itself. To facilitate these residual connections, all sub-layers in the model, as well as the embedding layers, produce outputs of dimension $d_{\text{model}} = 512$.

Decoder: The decoder is also composed of a stack of $N = 6$ identical layers. In addition to the two sub-layers in each encoder layer, the decoder inserts a third sub-layer, which performs multi-head attention over the output of the encoder stack. Similar to the encoder, we employ residual connections around each of the sub-layers, followed by layer normalization. We also modify the self-attention sub-layer in the decoder stack to prevent positions from attending to subsequent positions. This masking, combined with fact that the output embeddings are offset by one position, ensures that the predictions for position i can depend only on the known outputs at positions less than i .

Model Descriptions in the Original Paper

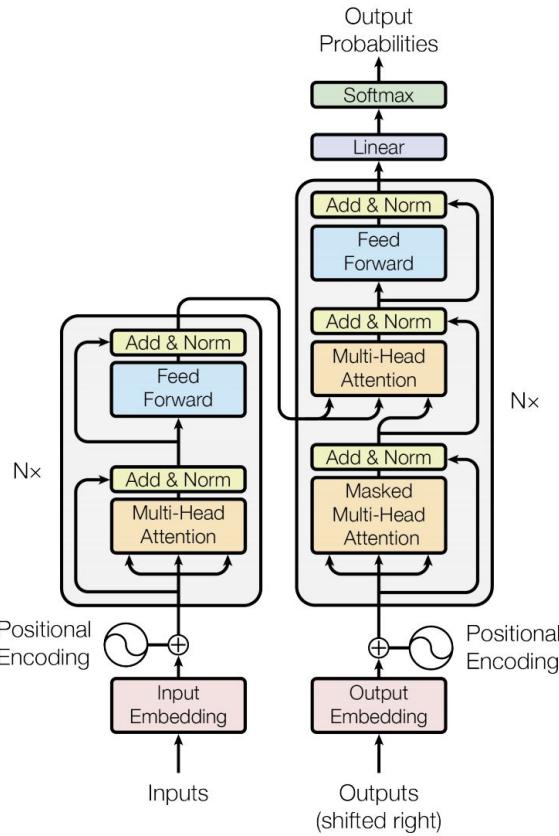


3.1 Encoder and Decoder Stacks

Encoder: The encoder is composed of a stack of $N = 6$ identical layers. Each layer has two sub-layers. The first is a multi-head self-attention mechanism, and the second is a simple, position-wise fully connected feed-forward network. We employ a residual connection [11] around each of the two sub-layers, followed by layer normalization [1]. That is, the output of each sub-layer is $\text{LayerNorm}(x + \text{Sublayer}(x))$, where $\text{Sublayer}(x)$ is the function implemented by the sub-layer itself. To facilitate these residual connections, all sub-layers in the model, as well as the embedding layers, produce outputs of dimension $d_{\text{model}} = 512$.

Decoder: The decoder is also composed of a stack of $N = 6$ identical layers. In addition to the two sub-layers in each encoder layer, the decoder inserts a third sub-layer, which performs multi-head attention over the output of the encoder stack. Similar to the encoder, we employ residual connections around each of the sub-layers, followed by layer normalization. We also modify the self-attention sub-layer in the decoder stack to prevent positions from attending to subsequent positions. This masking, combined with fact that the output embeddings are offset by one position, ensures that the predictions for position i can depend only on the known outputs at positions less than i .

Model Descriptions in the Original Paper

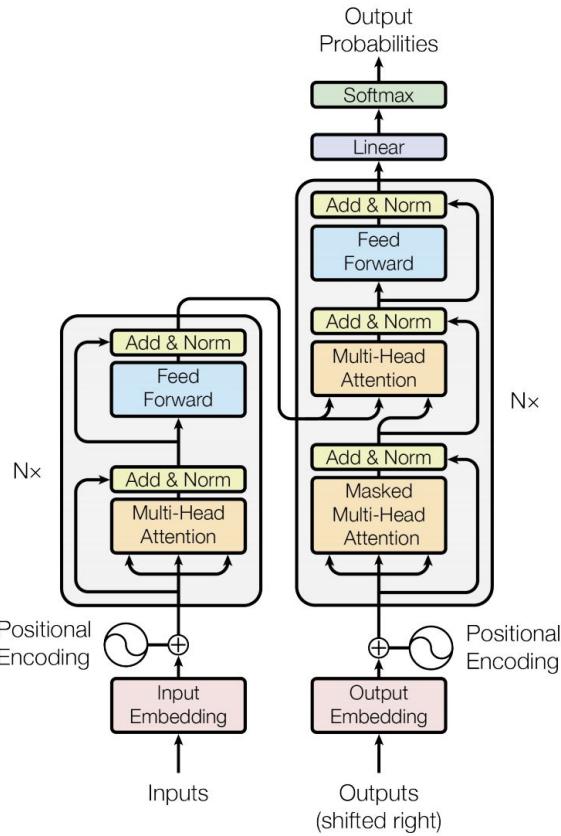


3.1 Encoder and Decoder Stacks

Encoder: The encoder is composed of a stack of $N = 6$ identical layers. Each layer has two sub-layers. The first is a multi-head self-attention mechanism, and the second is a simple, position-wise fully connected feed-forward network. We employ a residual connection [11] around each of the two sub-layers, followed by layer normalization [1]. That is, the output of each sub-layer is $\text{LayerNorm}(x + \text{Sublayer}(x))$, where $\text{Sublayer}(x)$ is the function implemented by the sub-layer itself. To facilitate these residual connections, all sub-layers in the model, as well as the embedding layers, produce outputs of dimension $d_{\text{model}} = 512$.

Decoder: The decoder is also composed of a stack of $N = 6$ identical layers. In addition to the two sub-layers in each encoder layer, the decoder inserts a third sub-layer, which performs multi-head attention over the output of the encoder stack. Similar to the encoder, we employ residual connections around each of the sub-layers, followed by layer normalization. We also modify the self-attention sub-layer in the decoder stack to prevent positions from attending to subsequent positions. This masking, combined with fact that the output embeddings are offset by one position, ensures that the predictions for position i can depend only on the known outputs at positions less than i .

Model Descriptions in the Original Paper

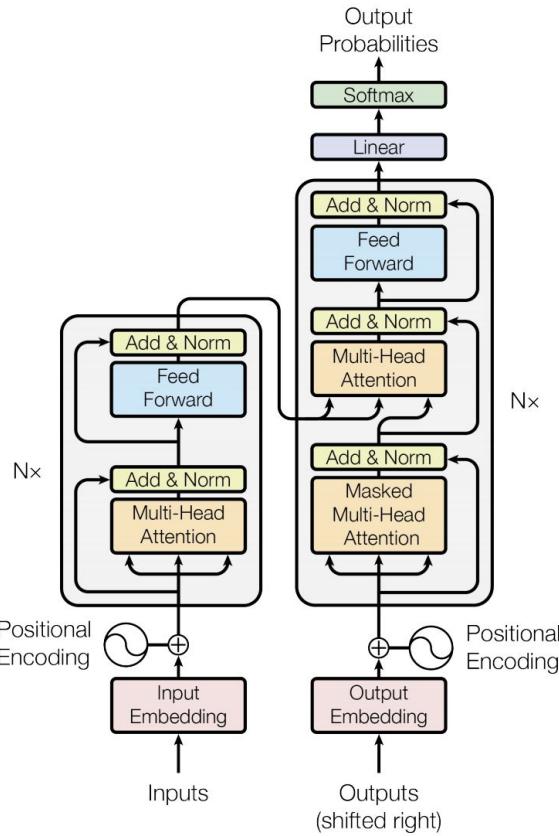


3.1 Encoder and Decoder Stacks

Encoder: The encoder is composed of a stack of $N = 6$ identical layers. Each layer has two sub-layers. The first is a multi-head self-attention mechanism, and the second is a simple, position-wise fully connected feed-forward network. We employ a residual connection [11] around each of the two sub-layers, followed by layer normalization [1]. That is, the output of each sub-layer is $\text{LayerNorm}(x + \text{Sublayer}(x))$, where $\text{Sublayer}(x)$ is the function implemented by the sub-layer itself. To facilitate these residual connections, all sub-layers in the model, as well as the embedding layers, produce outputs of dimension $d_{\text{model}} = 512$.

Decoder: The decoder is also composed of a stack of $N = 6$ identical layers. In addition to the two sub-layers in each encoder layer, the decoder inserts a third sub-layer, which performs multi-head attention over the output of the encoder stack. Similar to the encoder, we employ residual connections around each of the sub-layers, followed by layer normalization. We also modify the self-attention sub-layer in the decoder stack to prevent positions from attending to subsequent positions. This masking, combined with fact that the output embeddings are offset by one position, ensures that the predictions for position i can depend only on the known outputs at positions less than i .

Model Descriptions in the Original Paper

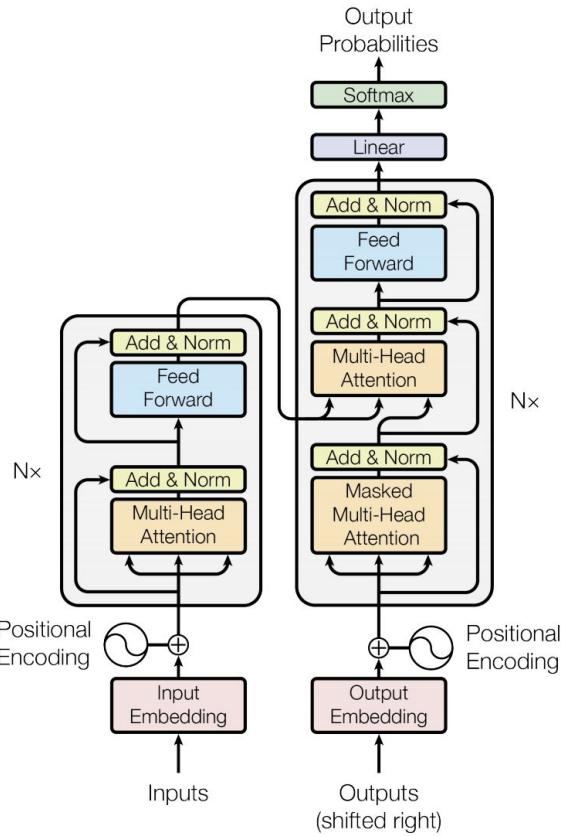


3.1 Encoder and Decoder Stacks

Encoder: The encoder is composed of a stack of $N = 6$ identical layers. Each layer has two sub-layers. The first is a multi-head self-attention mechanism, and the second is a simple, position-wise fully connected feed-forward network. We employ a residual connection [11] around each of the two sub-layers, followed by layer normalization [1]. That is, the output of each sub-layer is $\text{LayerNorm}(x + \text{Sublayer}(x))$, where $\text{Sublayer}(x)$ is the function implemented by the sub-layer itself. To facilitate these residual connections, all sub-layers in the model, as well as the embedding layers, produce outputs of dimension $d_{\text{model}} = 512$.

Decoder: The decoder is also composed of a stack of $N = 6$ identical layers. In addition to the two sub-layers in each encoder layer, the decoder inserts a third sub-layer, which performs multi-head attention over the output of the encoder stack. Similar to the encoder, we employ residual connections around each of the sub-layers, followed by layer normalization. We also modify the self-attention sub-layer in the decoder stack to prevent positions from attending to subsequent positions. This masking, combined with fact that the output embeddings are offset by one position, ensures that the predictions for position i can depend only on the known outputs at positions less than i .

Model Descriptions in the Original Paper

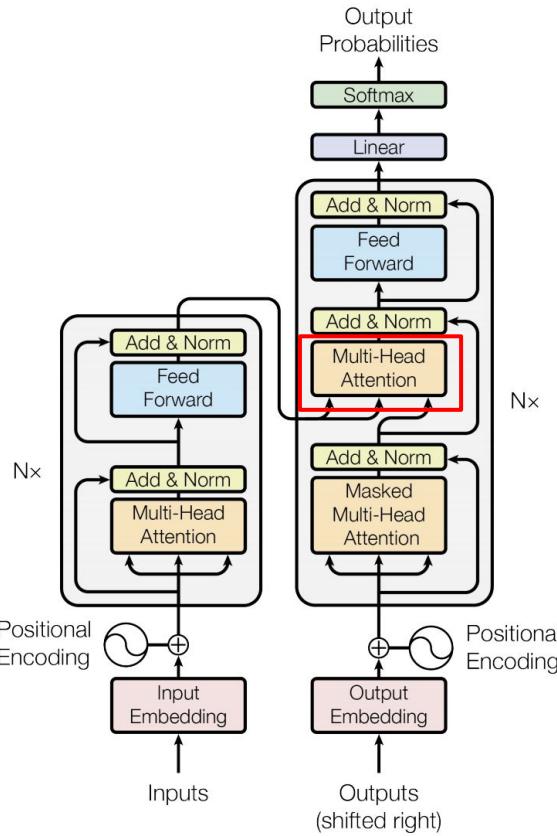


3.1 Encoder and Decoder Stacks

Encoder: The encoder is composed of a stack of $N = 6$ identical layers. Each layer has two sub-layers. The first is a multi-head self-attention mechanism, and the second is a simple, position-wise fully connected feed-forward network. We employ a residual connection [11] around each of the two sub-layers, followed by layer normalization [1]. That is, the output of each sub-layer is $\text{LayerNorm}(x + \text{Sublayer}(x))$, where $\text{Sublayer}(x)$ is the function implemented by the sub-layer itself. To facilitate these residual connections, all sub-layers in the model, as well as the embedding layers, produce outputs of dimension $d_{\text{model}} = 512$.

Decoder: The decoder is also composed of a stack of $N = 6$ identical layers. In addition to the two sub-layers in each encoder layer, the decoder inserts a third sub-layer, which performs multi-head attention over the output of the encoder stack. Similar to the encoder, we employ residual connections around each of the sub-layers, followed by layer normalization. We also modify the self-attention sub-layer in the decoder stack to prevent positions from attending to subsequent positions. This masking, combined with fact that the output embeddings are offset by one position, ensures that the predictions for position i can depend only on the known outputs at positions less than i .

Model Descriptions in the Original Paper

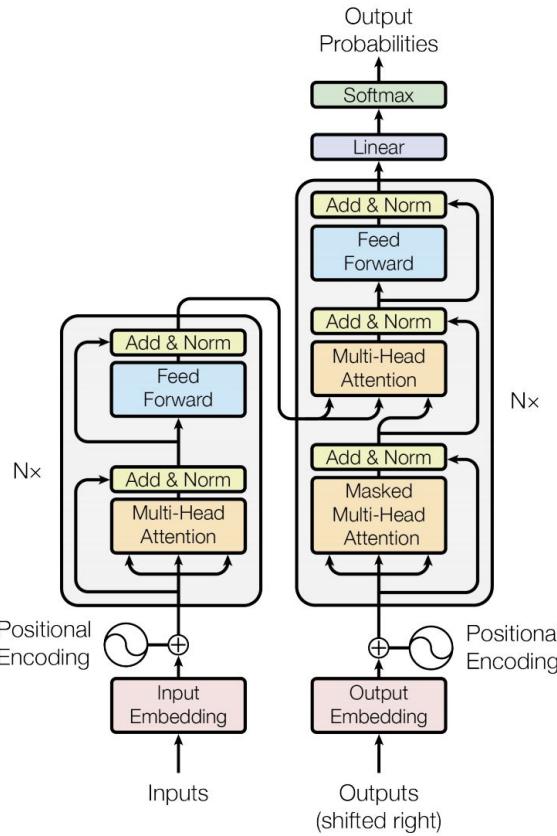


3.1 Encoder and Decoder Stacks

Encoder: The encoder is composed of a stack of $N = 6$ identical layers. Each layer has two sub-layers. The first is a multi-head self-attention mechanism, and the second is a simple, position-wise fully connected feed-forward network. We employ a residual connection [11] around each of the two sub-layers, followed by layer normalization [1]. That is, the output of each sub-layer is $\text{LayerNorm}(x + \text{Sublayer}(x))$, where $\text{Sublayer}(x)$ is the function implemented by the sub-layer itself. To facilitate these residual connections, all sub-layers in the model, as well as the embedding layers, produce outputs of dimension $d_{\text{model}} = 512$.

Decoder: The decoder is also composed of a stack of $N = 6$ identical layers. In addition to the two sub-layers in each encoder layer, the decoder inserts a third sub-layer, which performs multi-head attention over the output of the encoder stack. Similar to the encoder, we employ residual connections around each of the sub-layers, followed by layer normalization. We also modify the self-attention sub-layer in the decoder stack to prevent positions from attending to subsequent positions. This masking, combined with fact that the output embeddings are offset by one position, ensures that the predictions for position i can depend only on the known outputs at positions less than i .

Model Descriptions in the Original Paper

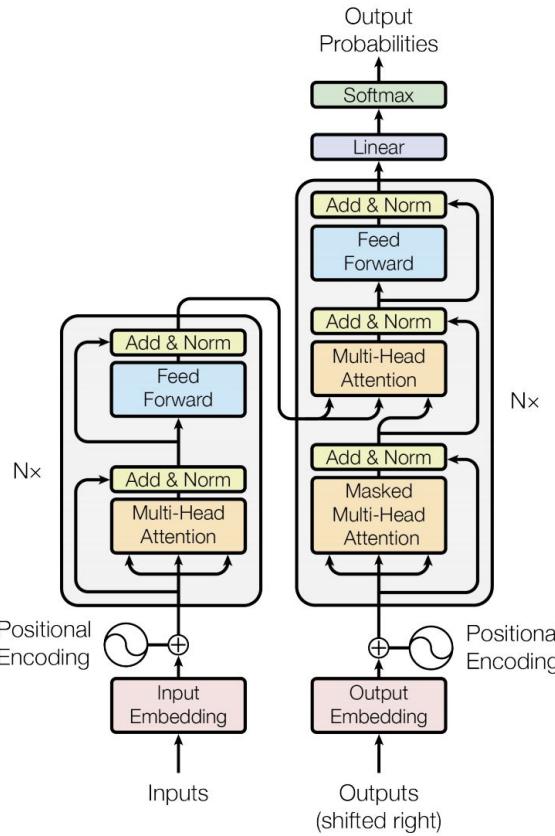


3.1 Encoder and Decoder Stacks

Encoder: The encoder is composed of a stack of $N = 6$ identical layers. Each layer has two sub-layers. The first is a multi-head self-attention mechanism, and the second is a simple, position-wise fully connected feed-forward network. We employ a residual connection [11] around each of the two sub-layers, followed by layer normalization [1]. That is, the output of each sub-layer is $\text{LayerNorm}(x + \text{Sublayer}(x))$, where $\text{Sublayer}(x)$ is the function implemented by the sub-layer itself. To facilitate these residual connections, all sub-layers in the model, as well as the embedding layers, produce outputs of dimension $d_{\text{model}} = 512$.

Decoder: The decoder is also composed of a stack of $N = 6$ identical layers. In addition to the two sub-layers in each encoder layer, the decoder inserts a third sub-layer, which performs multi-head attention over the output of the encoder stack. Similar to the encoder, we employ residual connections around each of the sub-layers, followed by layer normalization. We also modify the self-attention sub-layer in the decoder stack to prevent positions from attending to subsequent positions. This masking, combined with fact that the output embeddings are offset by one position, ensures that the predictions for position i can depend only on the known outputs at positions less than i .

Model Descriptions in the Original Paper



3.1 Encoder and Decoder Stacks

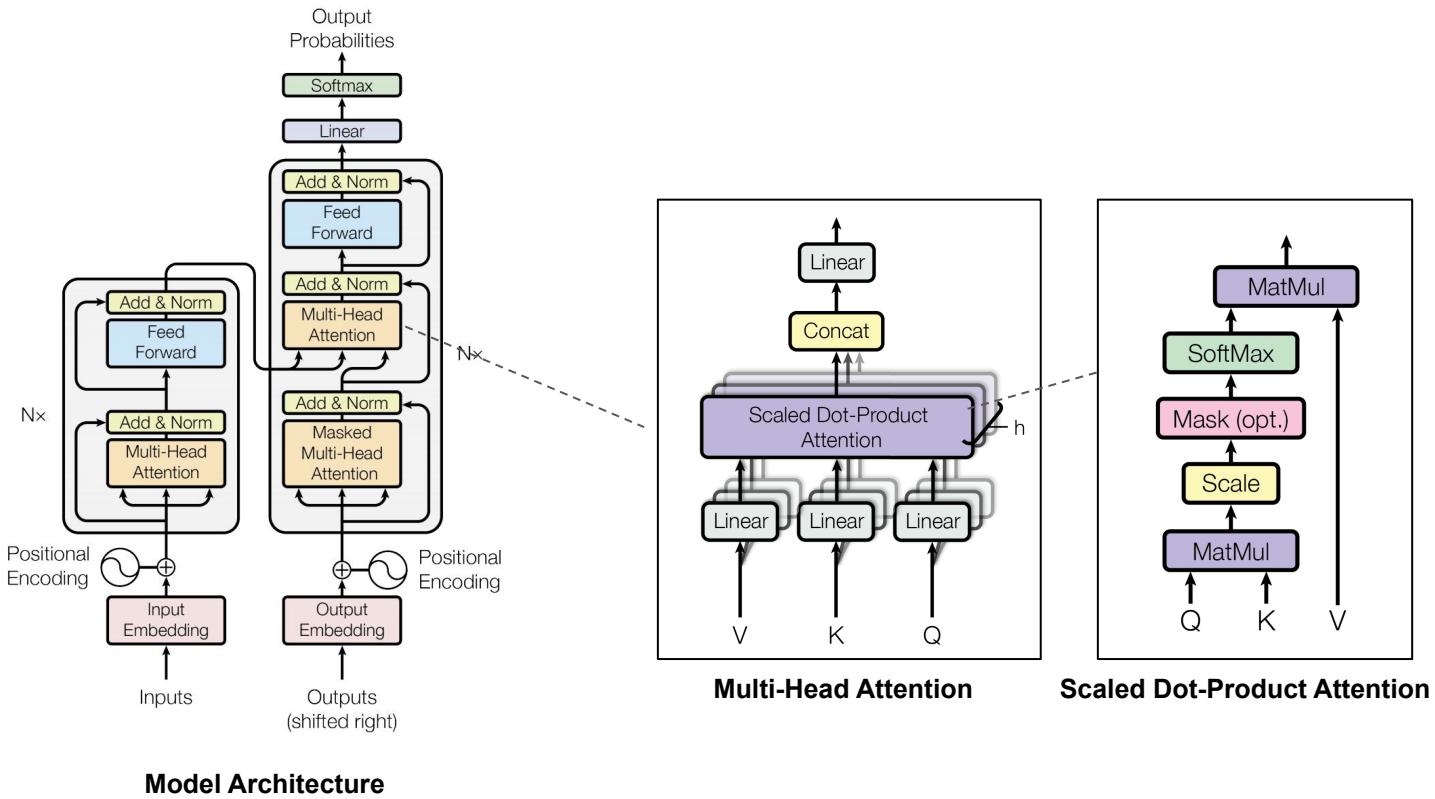
Encoder: The encoder is composed of a stack of $N = 6$ identical layers. Each layer has two sub-layers. The first is a multi-head self-attention mechanism, and the second is a simple, position-wise fully connected feed-forward network. We employ a residual connection [11] around each of the two sub-layers, followed by layer normalization [1]. That is, the output of each sub-layer is $\text{LayerNorm}(x + \text{Sublayer}(x))$, where $\text{Sublayer}(x)$ is the function implemented by the sub-layer itself. To facilitate these residual connections, all sub-layers in the model, as well as the embedding layers, produce outputs of dimension $d_{\text{model}} = 512$.

Decoder: The decoder is also composed of a stack of $N = 6$ identical layers. In addition to the two sub-layers in each encoder layer, the decoder inserts a third sub-layer, which performs multi-head attention over the output of the encoder stack. Similar to the encoder, we employ residual connections around each of the sub-layers, followed by layer normalization. We also modify the self-attention sub-layer in the decoder stack to prevent positions from attending to subsequent positions. This masking, combined with fact that the output embeddings are offset by one position, ensures that the predictions for position i can depend only on the known outputs at positions less than i .

PyTorch Code Reading

- [Transformer](#)
- [TransformerEncoder](#)
- [TransformerEncoderLayer](#)
- [MultiheadAttention](#)
- [TransformerDecoder](#)
- [TransformerDecoderLayer](#)

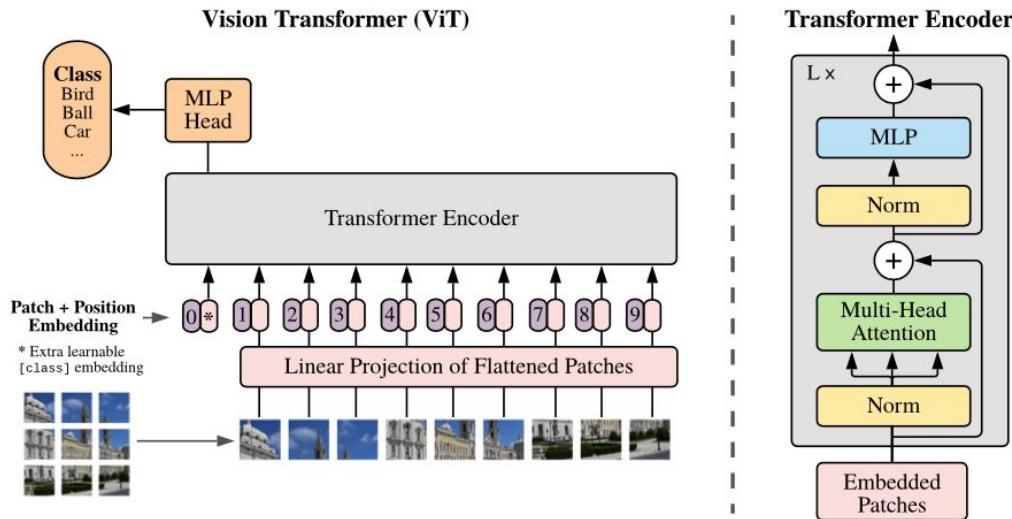
This Week's Goal: Understanding the Transformer Architecture 😊



Transformers for Computer Vision

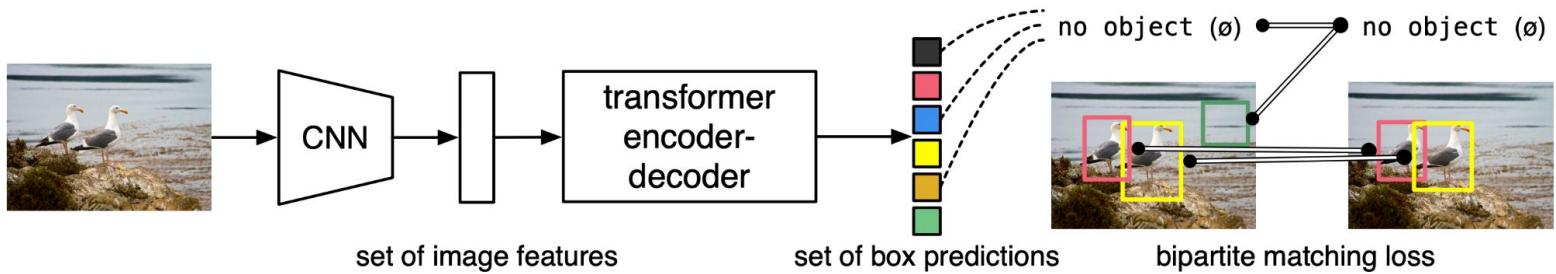
Vision Transformer (ViT) [Dosovitskiy et al. 2021]

- RGB Image → N Patches → Flattening patches
 - Each patch = $(3 * P^2)$ -dimensional vector



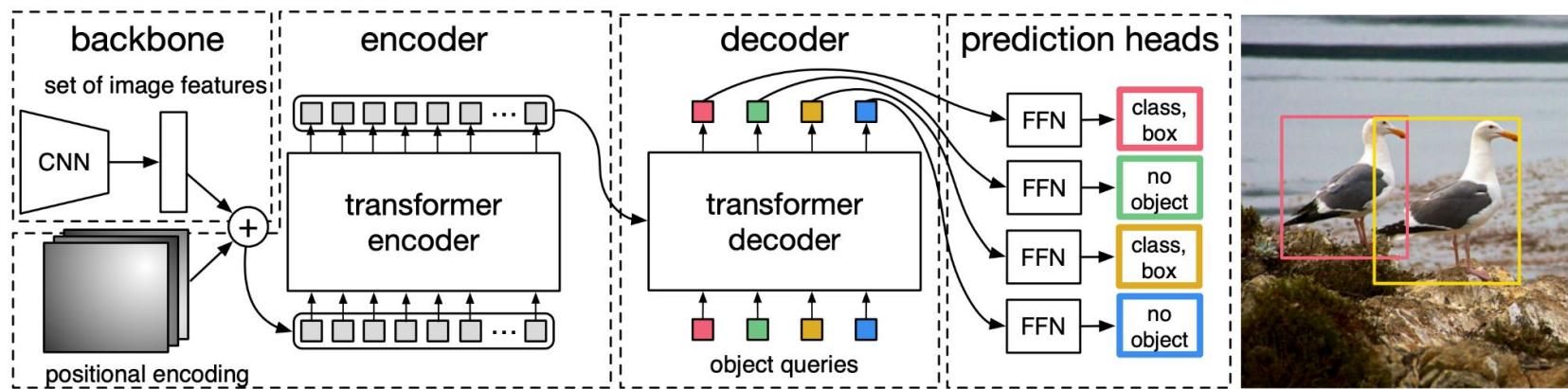
Detection Transformer (DETR) [Carion et al. 2020]

- CNN + Transformer for object detection

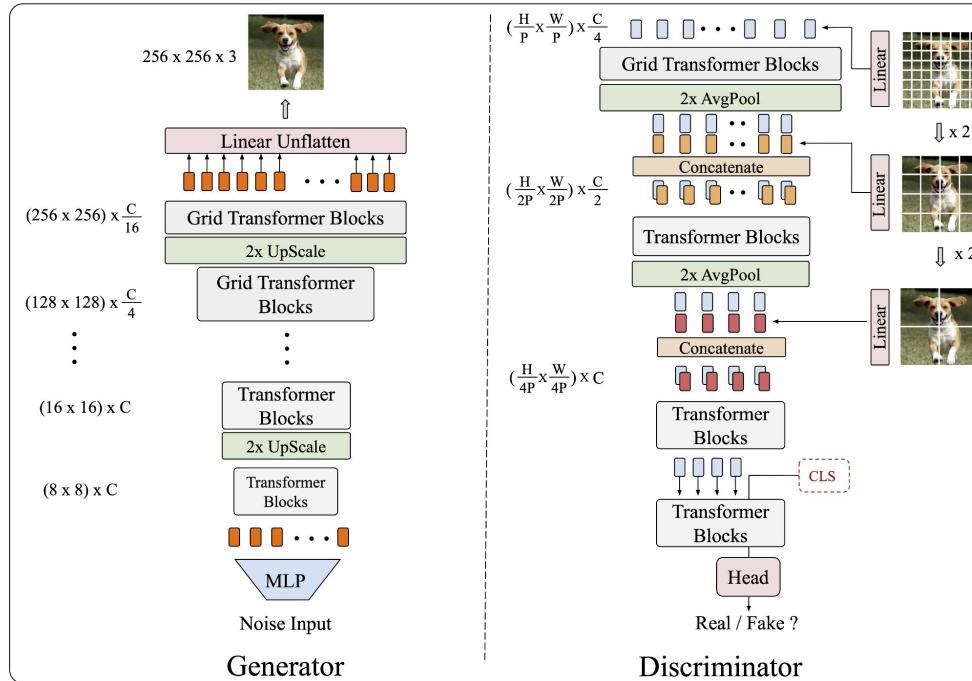


Detection Transformer (DETR) [Carion et al. 2020]

- CNN + Transformer for object detection



TransGAN [Jian et al. 2021]

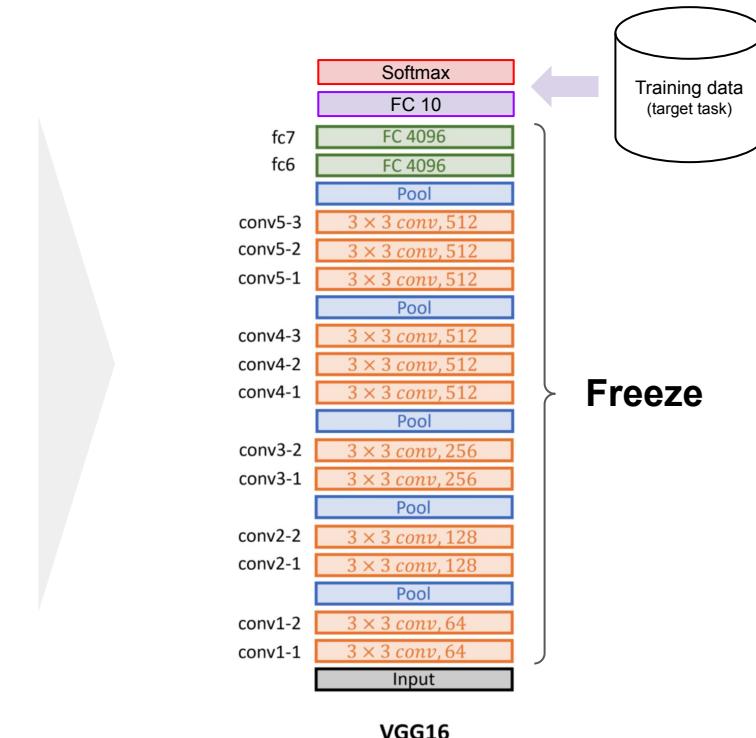
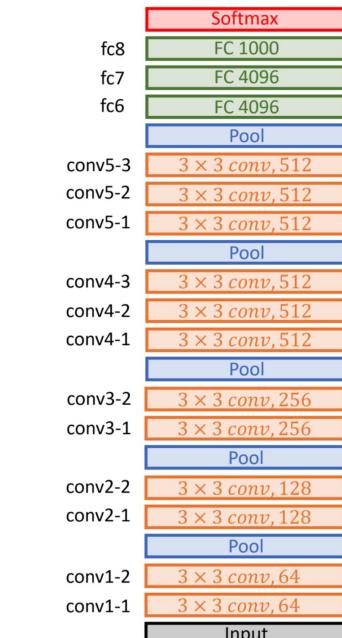
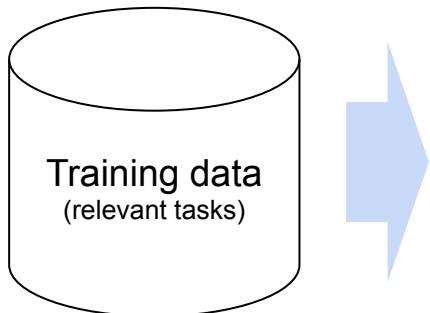


Yifan Jiang, Shiyu Chang, Zhangyang Wang, “TransGAN: Two Pure Transformers Can Make One Strong GAN, and That Can Scale Up”, arXiv 2021.

Pre-trained Language Models

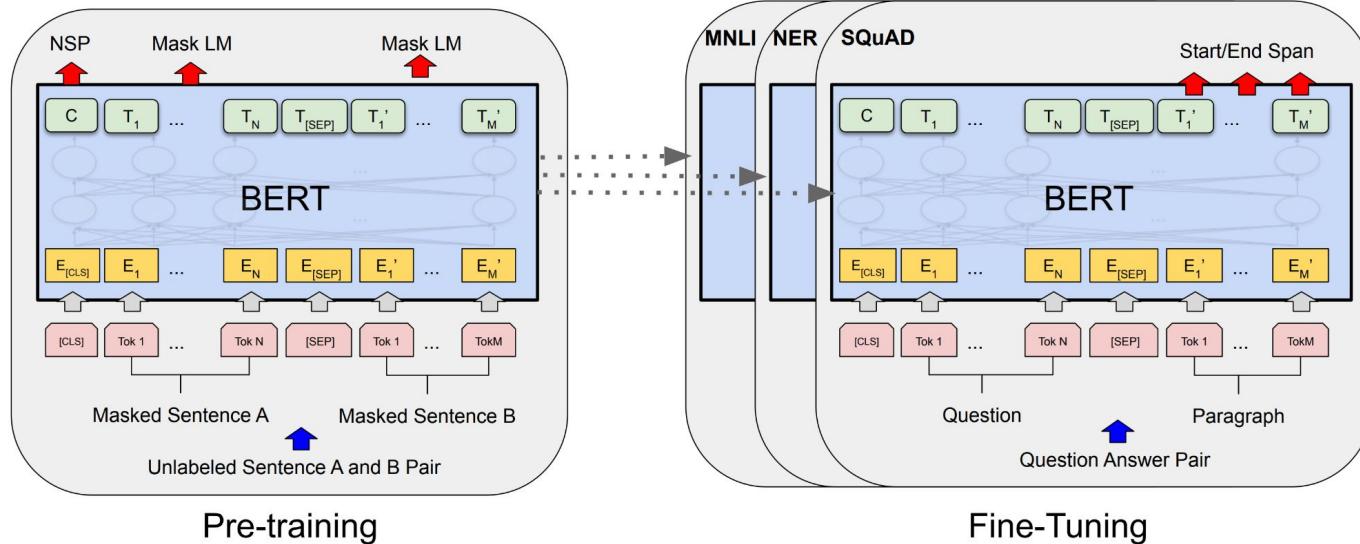
What are Pre-training and Fine-tuning?

- Example for Computer Vision

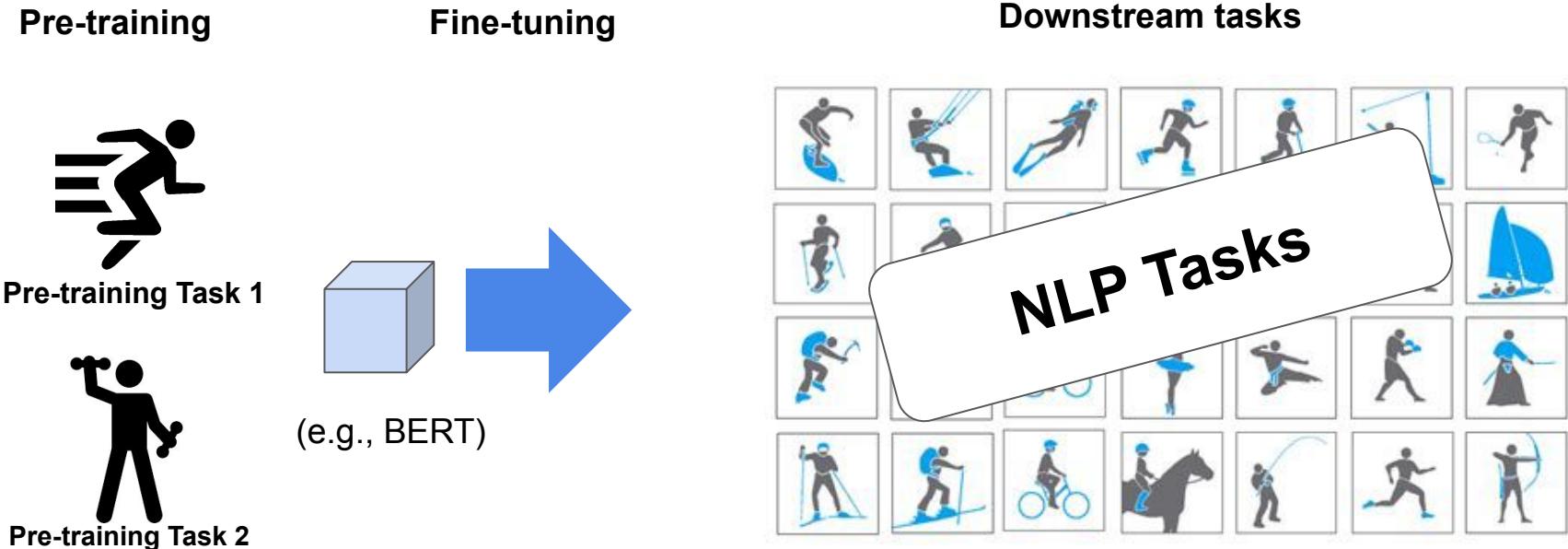


What are Pre-training and Fine-tuning?

- Example for NLP

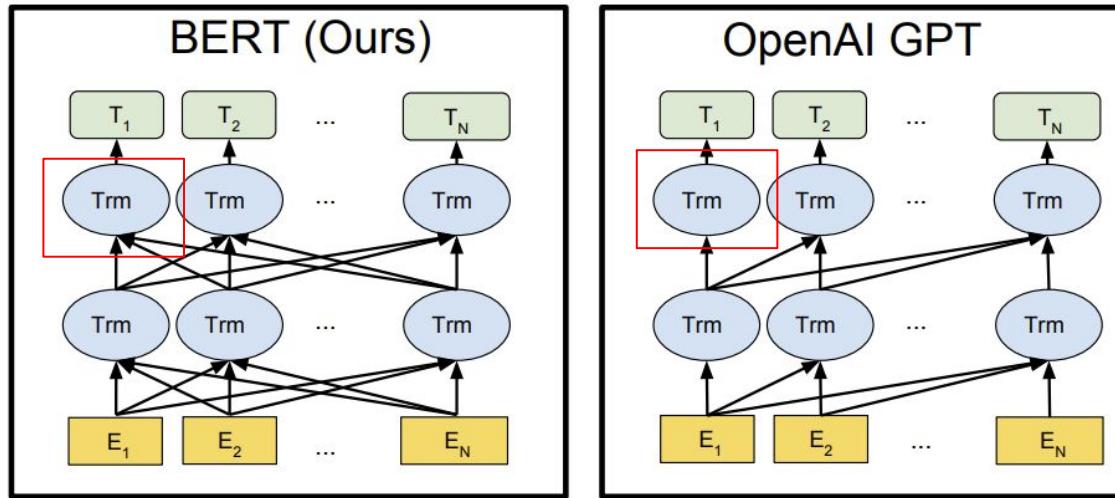


Pre-training & Fine-tuning: Intuition



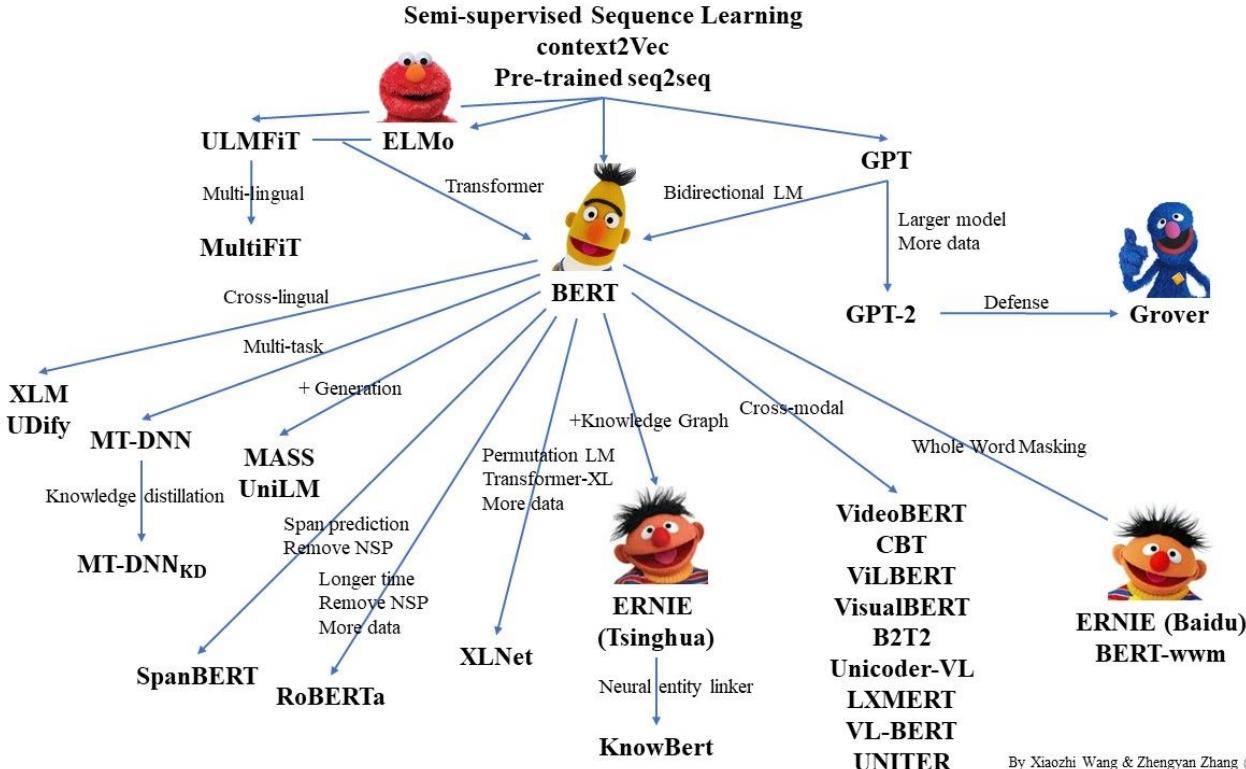
BERT and GPT

- BERT = Bi-directional Transformers (Transformer **Encoder**)
- GPT = left-to-right Transformers (Transformer **Decoder**)

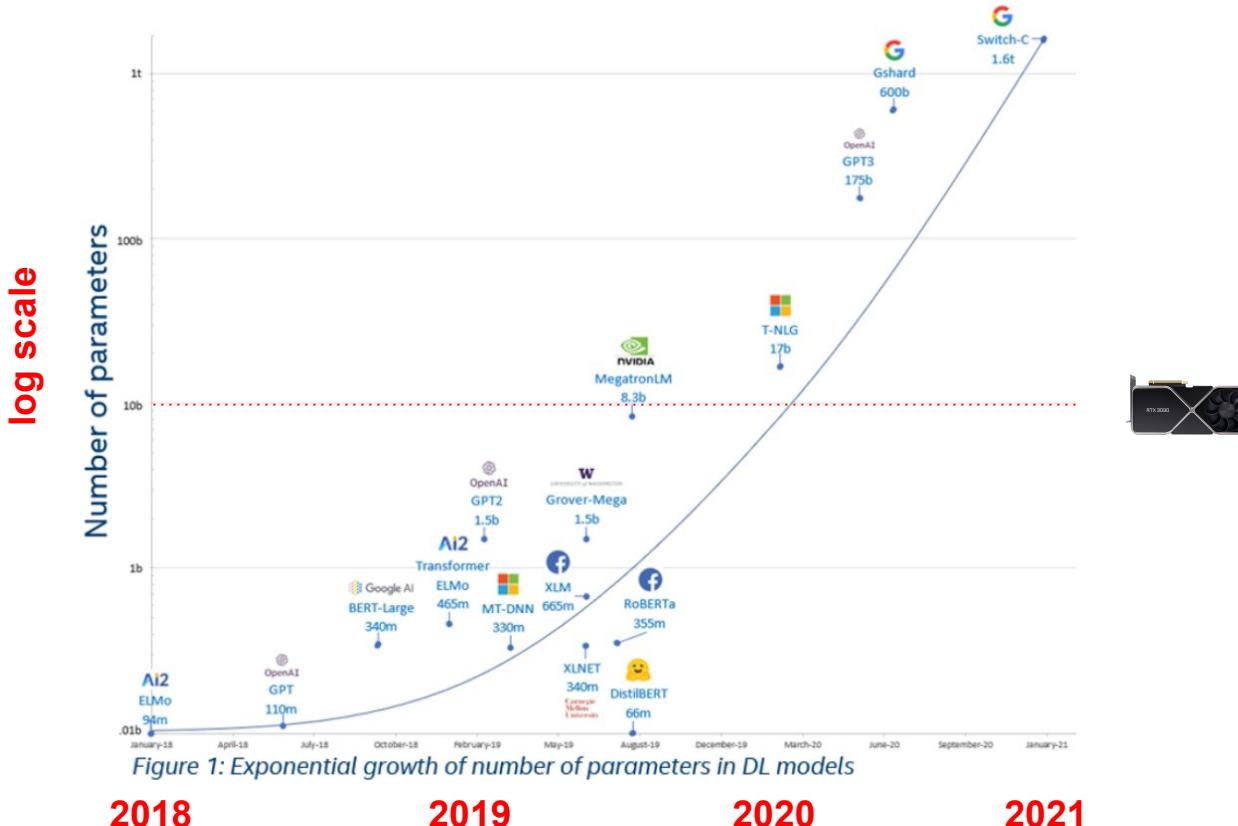


BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding
GPT: Generative Pre-Training

Family Tree of Pre-trained LMs



Pre-trained Model War: Rapid Increase in Model Size



Hugging Face's Transformers Library

<https://github.com/huggingface/transformers>

- You can use a variety of pre-trained language models just with a few lines of Python code

```
>>> from transformers import AutoTokenizer, AutoModel  
  
>>> tokenizer = AutoTokenizer.from_pretrained("bert-base-uncased")  
>>> model = AutoModel.from_pretrained("bert-base-uncased")  
  
>>> inputs = tokenizer("Hello world!", return_tensors="pt")  
>>> outputs = model(**inputs)
```

Summary

- The Transformer architecture
 - Encoder & Decoder
 - Encoder-decoder attention
- Brief introduction to Transformers for Computer Vision
 - Vision Transformer
 - Detection Transformer
 - TransGAN
- Brief overview of pre-trained Language Models
 - BERT and GPT are Transformer Encoder and Decoder
 - (Almost) all pre-trained models today use the Transformer architecture

Questions?

Further Reading

- [The Illustrated Transformer – Jay Alammar – Visualizing machine learning one concept at a time.](#)
- [Multi-head attention mechanism: "queries", "keys", and "values," over and over again](#)