



16/2/2022

## BERT Research Series

### BERT

- Released in October 2018
  - Surpassed human performance in language understanding
- \* Transfer learning

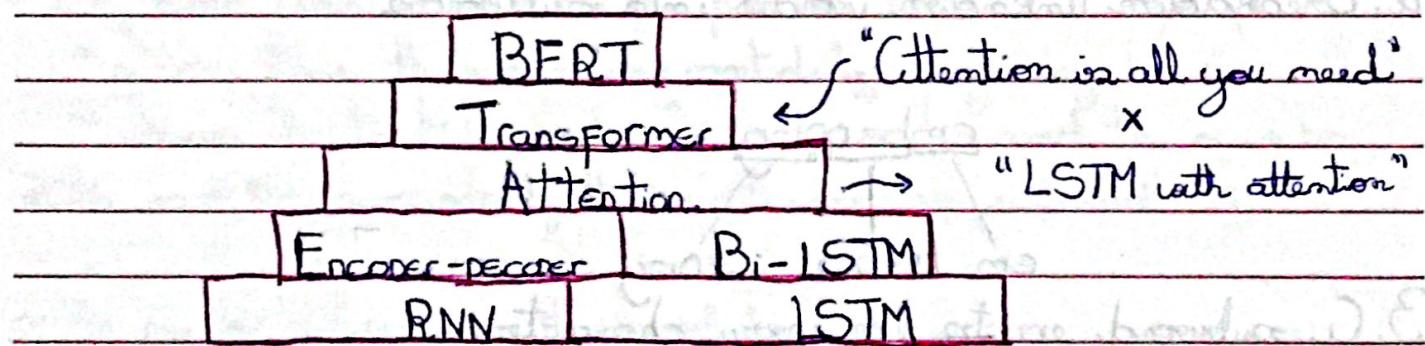
→ "Take the model, modify it by changing the output layer, doing a finetune training"

Leverage pretraining

"BERT is huge and expensive to train from scratch."

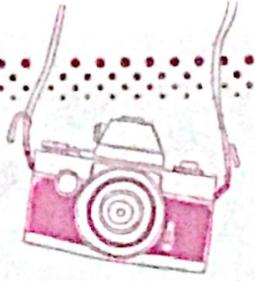
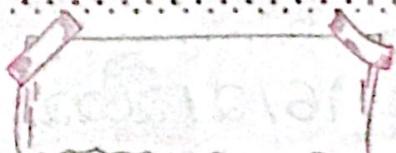
"BERT maintain"

→ "self-attention"



My playlist today: #





Word embedding

Feature representation of a word

$\langle 0.4, 0.6, \dots, 0.5 \rangle$

BERT = 30K tokens  $\times$  768 features

dictionary = word string  $\rightarrow$  word id

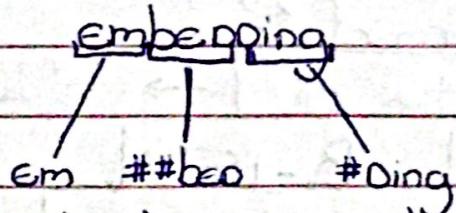
distance between word embeddings  
↓

distance between word meanings  
=

Word similarity

BERT's vocabulary

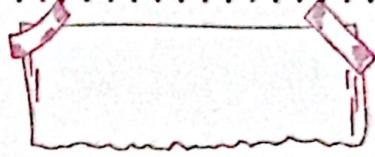
1. Bert is pretrained.  $\rightarrow$  fixed vocabulary
2. Break down unknown words into subwords



3. A subword exists for every character

~ all subwords start with  $\#\#$  except for the first subword

AS



From perusing the vocabulary, we saw that:

- The first 999 tokens (1-indexed) appear to be reserved and are most of the form [unseen957].

1 - [PAD]

101 - [UNK]

102 - [CLS]

103 - [SEP]

104 - [MASK]

- Rows 1000 - 1996 appear to be a dump of individual characters.

"They don't appear to be sorted by frequency (e.g. the letters of the alphabet are all sequence).

- The first word is "the" at the position 1997

\* From there, the words are sorted by frequency

\* From The top ~18 words are whole words, and the number 2016 is ##\$, presumably the most common subword

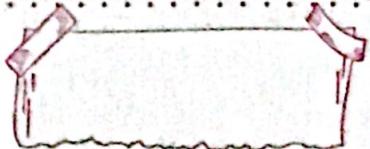
~ 3K person names in the vocabulary

~ many numbers

Meow  
Meow



My playlist today #



BERT already encodes a lot about our language

only needs finetuning

### Advantages of Fine-tuning

1. Quicker development
2. Less data required
3. Better results

✗ Not train specific

"Plug and finetune"

### BERT shortcomings

1. BERT is very large
  - a. after finetuning
  - b. slow inference

Layer	Weights
Embedding	~24M
Transformers × 12	→ $7 \times 12 = 85$
TOTAL	~109M

### Distillation

Remove weights while maintaining accuracy

### 2. Jargon (domain-specific language)

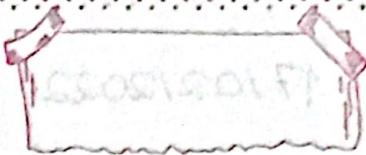
### 3. Not all NLP applications

Yes

- Classification
- NER
- POS
- Extractive QA

No

- Language model
- Text generation
- Translation



## Required formatting

## > Special tokens

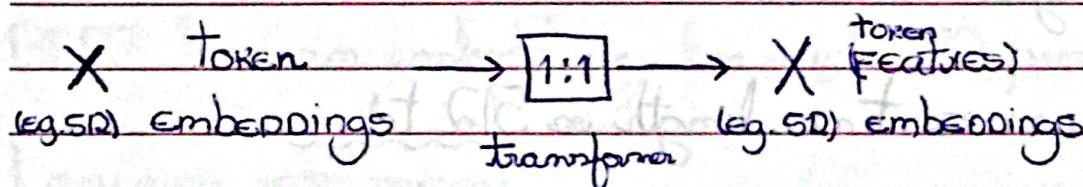
[SEP] separation - separates two sentences

[CLS] classification - must be the 1<sup>st</sup> token

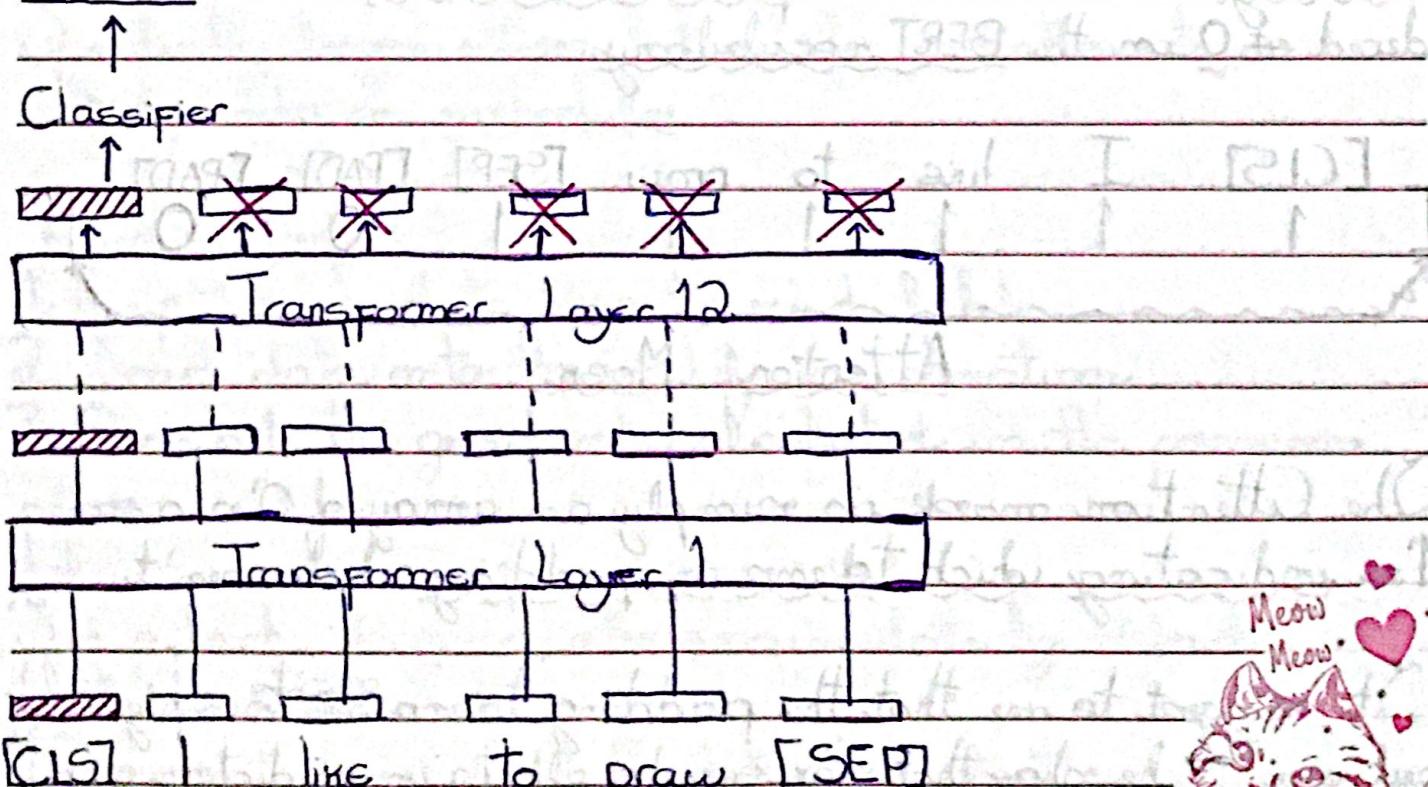
## Architecture

BERT = stack of 12 transformers

each transformer



## Prediction



My playlist today: # each token has its embeddings



17/02/2022



"The first token of every sequence is always a special classification token ([CLS]). The hidden state corresponding to this token is used as the aggregate representation for classification tasks." (from the BERT paper)

Sequence length & Attention masks

→ Does BERT handle with sentences with varying lengths?

BERT has two constraints:

1. All sentences must be padded or truncated to a single, fixed length.

2. The maximum sentence length is 512 tokens

Padding is done with a special [PAD] token, which is indexed at 0 in the BERT vocabulary.

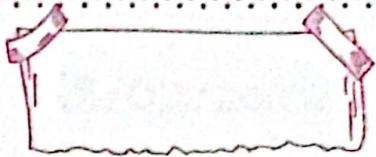
[CLS]	I	like	to	draw	[SEP]	[PAD]	[PAD]
1	1	1	1	1	1	0	0

Attention Mask

The Attention mask is simply an array of 0's and 1's indicating which tokens are padding which aren't.

"It suggest to me that the padding tokens ~~are~~ simply skipped over → he show that max seq len affects in prediction acc."

We can do it!



## Preprocessing steps

1. Convert sentence into list of vocabulary ID's
2. Pad or truncate
3. Calculate attention masks

## BertForSequenceClassification

Fine-tuning

Load the pretrained BERT model with a single linear classification on top

BERT has an embedding for each 512 positions

## Learning rate decay

Initially larger steps, later smaller steps  
↳ depends on scheduler

## Training loop

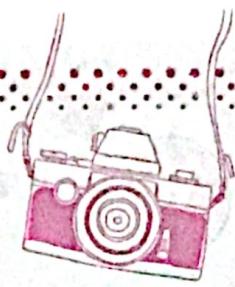
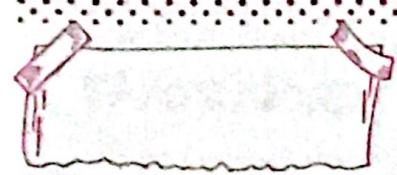
1. Unpack our data inputs and labels
2. Load data onto the GPU for acceleration
3. Clear out the gradient calculated in the previous pass

In pytorch, the gradients accumulate by default  
(better for RNN)

4. Forward pass. I feed the data through the

My playlist today: #





network.)

5. Backward pass (backpropagation)

6. Tell the network to update parameters

7. Track variables for monitoring progress.

### Evaluation loop

1. Unpack our data inputs and labels

2. Load data onto GPU for acceleration

3. Forward pass

4. Compute loss on our validation data and track various  
losses for monitoring progress

Edits: "drop out" and "batch normalization" behave differently  
during training and test.\*

Batch contains three pytorch tensors:

[0]: input ids

[1]: attention masks

[2]

gradient clipping avoids gradients to explode.

Ass

18/02/2022

BERT = Encoder <sup>x12</sup> + Decoder

Transformer = Encoder <sup>x6</sup> + Decoder <sup>x6</sup>

Transformers Illustrated - Jay Alammar

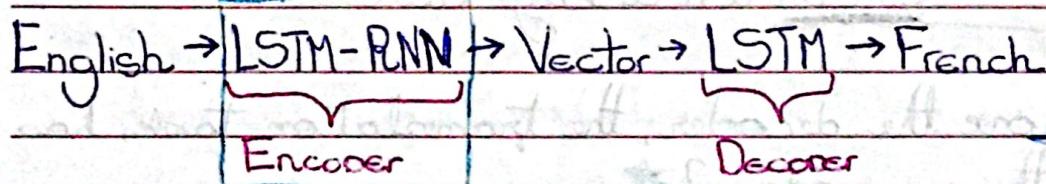
Transformers use attention in a different way  
self-attention

Neural Machine Translation (NMT)

→ Don't use hand engineered features

Transformers outperformed Google NMT

Before Google Translate



Probability 2<sup>nd</sup> word =  $\mathcal{P}(1^{\text{st}} \text{ word in English}, 2^{\text{nd}} \text{ word vector})$   
in French

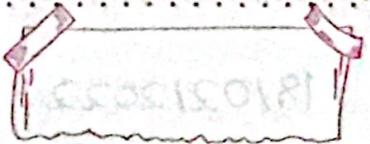
Stopping technique → [EOS] End of Sentence in vocabulary  
when [EOS] is more likely to be the next word

Google added attention

Attention mechanism tells which word should be focused on

My playlist today, #





"The biggest benefit from Transformers comes from how the Transformer lends itself to parallelization."

Transformer = Encoder <sup>x6</sup> + Decoder <sup>x6</sup>

- > stacks of encoder encoders and decoders
- > same architecture but different weights

→ This first half is interpreting well the input data, could we use it for other tasks instead of translation?

↳ BERT's idea

- Use Transformer's encoder
- 12 encoders instead of 6
- Train in a huge dataset

→ If we remove the decoder, the translation task, how do we train those encoders? \*

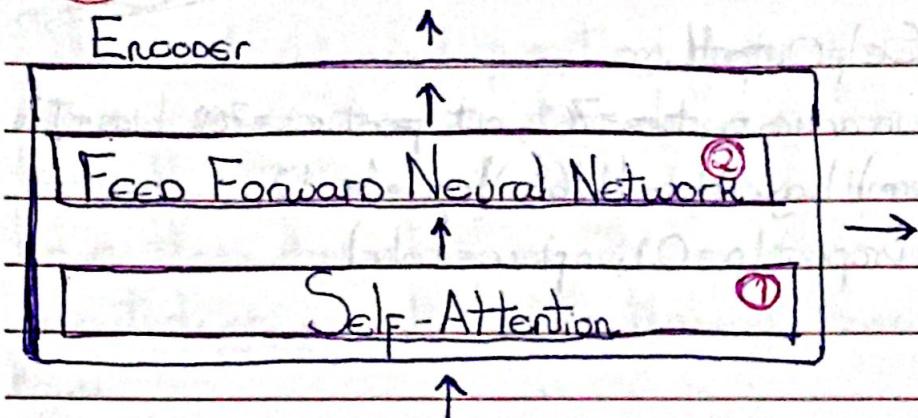
"fake-tasks"

1. Masked Language Model (MLM) - predict the masked word
2. Next Sentence Prediction (NSP)

↳ Was sentence B said immediately after sentence A?

You can apply BERT for different tasks by changing the output layer.

We can do it!



## Feed Forward Neural Network

↳ standard MLP dense neural network layer

~ From BERT HF parameters

(embeddings): BertEmbeddings | vocab

{ word embeddings: Embedding(30522, 768, padding\_idx=0)

position embeddings: Embedding(512, 768)

(token-type embeddings) Embedding(2, 768)

\* (LayerNorm): LayerNorm(768, eps=1e-12)

(dropout): Dropout(p=0.1, inplace=False)

)

dimension = 768 =  $6 \times 128$  (or  $512 + 256$ ), 50% longer than original

\* Layer Normalization (2016) = reduce training time

(attention): BertAttention |

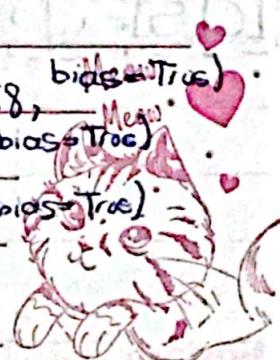
(self): BertSelfAttention |

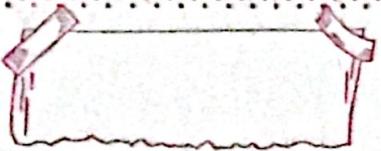
① (query): Linear(in\_features=768, out\_features=768, bias=True)

(key): Linear(in\_features=768, out\_features=768, bias=True)

(value): Linear(in\_features=768, out\_features=768, bias=True)

My playlist today: #



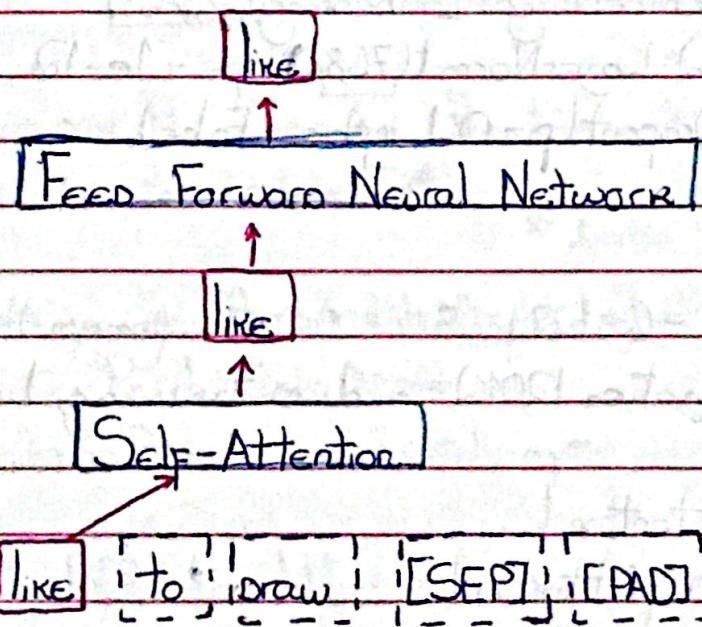


(output): BertSelfOutput

(dense): Linear(in\_features=768, out\_features=768, bias=True)  
② (LayerNorm): LayerNorm(768, ) {eps=1e-12},  
(dropout): Dropout(p=0.1, inplace=False)  
)  
}

## Self-Attention

A layer that helps the encoder look at other words in the input sentences as it encodes a specific word. The outputs of the self-attention are feed-forward neural network.

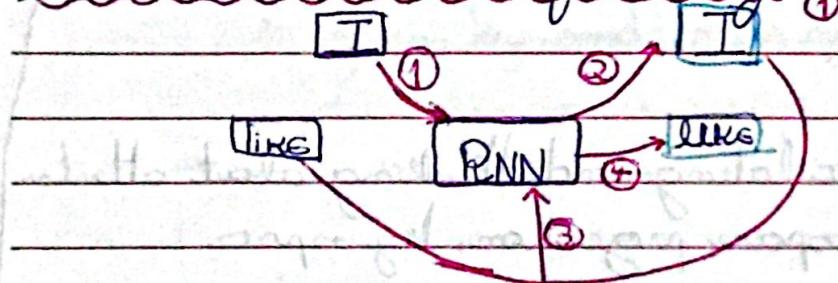


A/B

20/02/2022

The word in each position flows through its own path in the encoder. There are dependencies between these paths in the self attention layer. The feed forward layer does not have these dependencies. However, the various paths can be executed in parallel while flowing through the feed forward layer.

RNN has to run sequentially.  $I \rightarrow \text{like} \rightarrow \text{to}$



The animal didn't cross the street because it was too tired.

Self-attention allows to associate it with animal

### Self-Attention

Given an input word, score every other word by relevance

### Self-Attention in detail

vectors:

Word embedding  $x_i$   $\rightarrow$  query  $q_i$   
 $\rightarrow$  key  $k_i$   
 $\rightarrow$  value  $v_i$

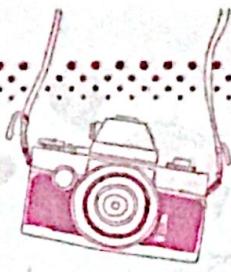
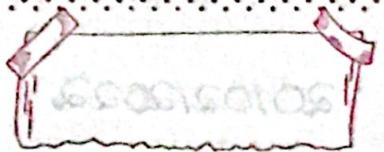
Thinking machines

$x_1$   $x_2$  = Meow

Meow

My playlist today #





1<sup>st</sup> step - Calculate query, key, and values

$$q_1 \quad k_1 \quad v_1$$

Eg. "Thinking machines" and encoding ( $\overset{w^Q}{\text{thinking}}$ )

$$x_1 \quad x_2$$

Thinking ( $x_1$ )

$$\begin{aligned} \rightarrow x_1 \cdot W^Q &= q_1 \\ \rightarrow x_1 \cdot W^K &= k_1 \\ \rightarrow x_1 \cdot W^V &= v_1 \end{aligned}$$

Machines ( $x_2$ )

$$\begin{aligned} \rightarrow x_2 \cdot W^Q &= q_2 \\ \rightarrow x_2 \cdot W^K &= k_2 \\ \rightarrow x_2 \cdot W^V &= v_2 \end{aligned}$$

$W^Q$  (query weights),  $W^K$ ,  $W^V$

- abstractions that help calculating and thinking about attention
- projection features: query space projection, key space ...

2<sup>nd</sup> step - Calculate scores

Score of word  $w_i$  while encoding word  $w_j$  is given by:

$$\text{score}(i, j) = q_j \cdot k_i$$

$$\text{Score}(\text{Thinking}) = q_1 \cdot k_1$$

$$\text{Score}(\text{Machines}) = q_2 \cdot k_2$$

Query space is about the word being encoded.

Key space is about the word being calculated

3<sup>rd</sup> and 4<sup>th</sup> steps - Score Normalizing

(softmax)

$$= \text{softmax} \left( \frac{\text{score}(i, j)}{\text{len}(k_i)} \right)$$

Divide scores by  $\sqrt{\text{len}(k_i)}$  (square root of key dimension) and normalize through softmax leading to more stable gradients.

FORONI

We can do it!

AS

Value space is about word relevance

5<sup>th</sup> - multiply by  $\alpha_i$ , the value vector by the softmax scores. The intuition here is to drown out irrelevant words.

6<sup>th</sup> - sum up weighted value vector

~ All those calculations are made in matrix for faster processing:

1<sup>st</sup> calculate query, key, and value matrices

$$X \cdot W^Q = Q$$

$$X \cdot W^K = K$$

$$X \cdot W^V = V$$

2<sup>nd</sup> to 6<sup>th</sup> steps

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

### Multi-headed attention

Improves the performance of attention layer in two ways:

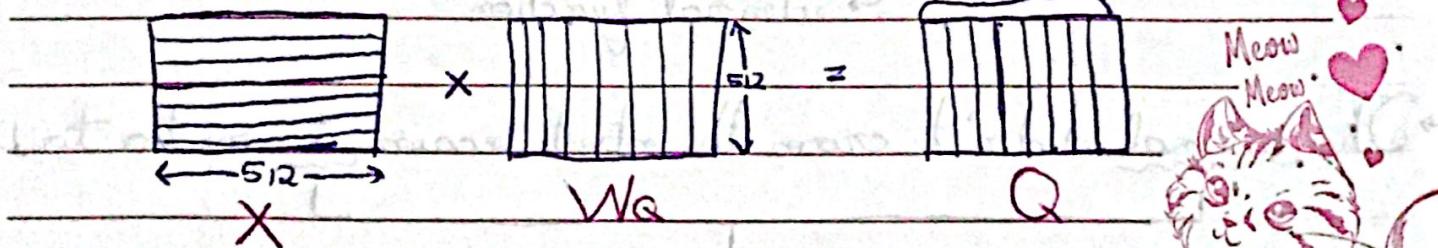
1. Expands the model's ability to focus on different positions

2. Multiple "representation subspaces"  $\rightarrow$  find different local minimums

- Multiple sets of query/key/value: 8 attentions - 8 sets

- Each randomly initialized

$$8 \text{ heads} \rightarrow 512 = 8 \times 64$$



My playlist today: #





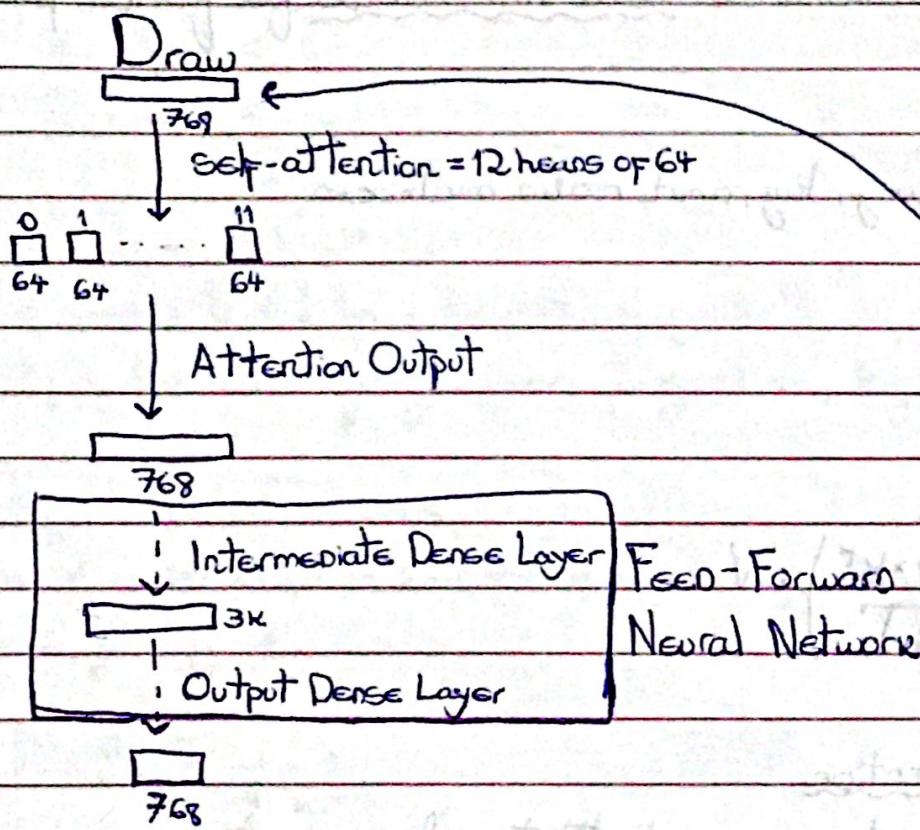
## Transformer

- 8 heads len 64
- len 512 embeddings

## BERT

- 12 heads len 64
- len 768 embeddings

## Feed Forward Neural Network



## Positional Encoding

- RNN does it naturally.

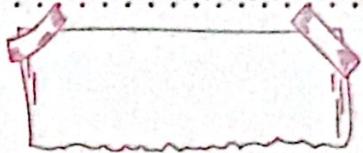
what is fed to the model

embeddings + positional encoding = encoding with signal  
↳ external function

→ The animal didn't cross the street because it was too tired.

Positional encoding considers this distance

We can do it!

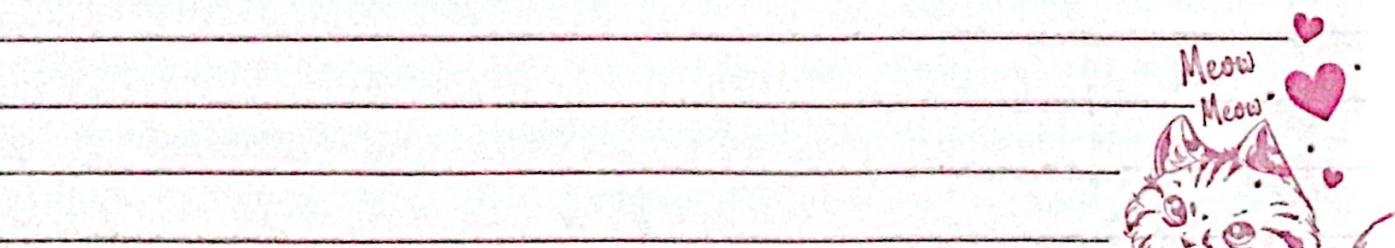


## BERT "Pre-training" tasks

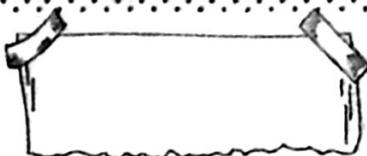
- Masked Language Model (MLM) → cloze task - very known in NLP
- Next Sentence Prediction (NSP)

The training loss is the sum of the mean MLM likelihood and the NSP likelihood.

They are example such that combined length is  $\leq 512$  tokens.



My playlist today: #



## BERT "Pre-training" tasks

- Masked Language Model (MLM) → close task - very known in NLP
- Next Sentence Prediction (NSP)

The training loss is the sum of the mean MLM likelihood and the NSP likelihood.

They are samples such that combined length is  $\leq 512$  tokens.

My playlist today: #

Meow  
Meow

