

BERT and The History Behind It

Daniil Osokin

Agenda

1. Some typical NLP tasks.
2. Sequence to sequence learning (Seq2Seq).
3. Transformer.
4. BERT.

Some typical NLP tasks

- Named entity recognition (NER), intent classification, slot filling (chatbots, virtual assistants):

Sentence	<i>show</i>	<i>flights</i>	<i>from</i>	<i>Boston</i>	<i>To</i>	<i>New</i>	<i>York</i>	<i>today</i>
Slots/Concepts	O	O	O	B-dept	O	B-arr	I-arr	B-date
Named Entity	O	O	O	B-city	O	B-city	I-city	O
Intent	<i>Find Flight</i>							
Domain	<i>Airline Travel</i>							

Identify known entities in a given text, classify the intent.

Some typical NLP tasks

- Question answering:

context: Beyoncé Giselle Knowles-Carter (/biːˈjɒnseɪ/ bee-YON-say) (born September 4, 1981) is an American singer, songwriter, record producer and actress. Born and raised in Houston, Texas, she performed in various singing and dancing competitions as a child, and rose to fame in the late 1990s as lead singer of R&B girl-group Destiny's Child. Managed by her father, Mathew Knowles, the group became one of the world's best-selling girl groups of all time. Their hiatus saw the release of Beyoncé's debut album, **Dangerously in Love** (2003), which established her as a solo artist worldwide, earned five Grammy Awards and featured the Billboard Hot 100 number-one singles "Crazy in Love" and "Baby Boy".

question: What was the name of Beyoncé's first solo album?

answer: Dangerously in Love

- Neural machine translation (NMT):

Je m'appelle Daniil → My name is Daniil

Tweet sentiment extraction

Given a tweet and sentiment, predict words subsequence, which responsible for the sentiment.

tweet: I`d have responded, if I were going

sentiment: neutral

answer: I`d have responded, if I were going

tweet: Sooo SAD I will miss you here in San Diego!!!

sentiment: negative

answer: Sooo SAD

tweet: I really really like the song Love Story by Taylor Swift

sentiment: positive

answer: like

tweet: I`m sorry.

sentiment: negative

answer: I`m sorry.

tweet: is back home now gonna miss every one

sentiment: negative

answer: onna

tweet: hm... Both of us I guess...

sentiment: neutral

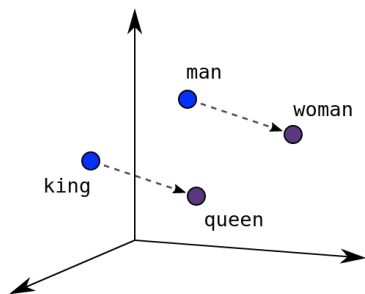
answer: hm... Both of us I guess...

From words to vectors

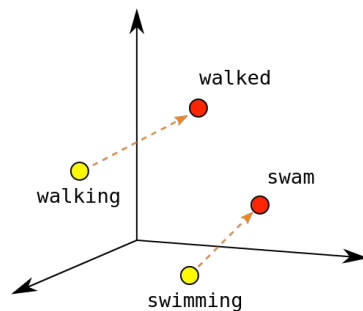
Q: How to map words to vectors?

From words to vectors

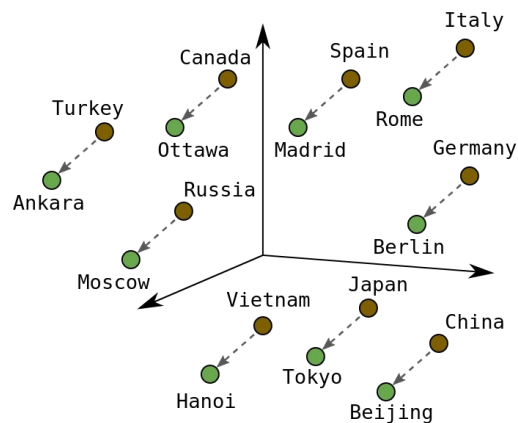
- One-hot encoding.
- Word2vec, GloVe.
- Embedding layer.



Male-Female



Verb Tense



Country-Capital

Slide name is hidden

1. $W^{\text{hx}}_X = ?$

Slide name is hidden

1. $W^{hx}x = ?$
2. $\text{sigm}(W^{hx}x) = ?$

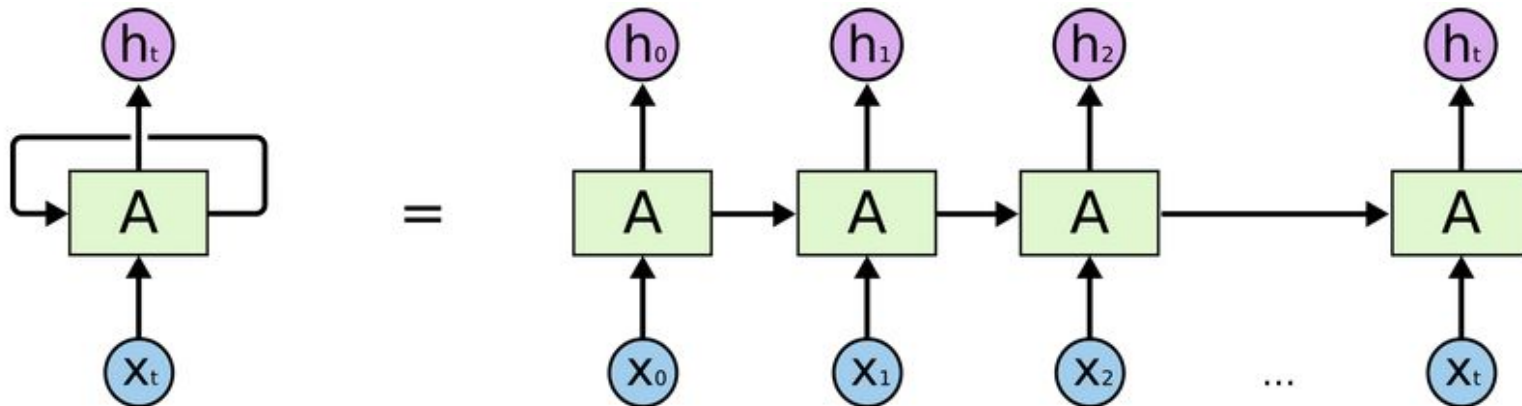
Slide name is hidden

1. $W^{hx}x = \text{FC}$.
2. $\text{sigm}(W^{hx}x) = \text{FC} + \text{activation}$.
3. $\text{sigm}(W^{hx}x_t + W^{hh}h_{t-1}) = ?$

Recurrent neural network

1. $W^{hx}x = \text{FC}$.
2. $\text{sigm}(W^{hx}x) = \text{FC} + \text{activation}$.
3. $\text{sigm}(W^{hx}x_t + W^{hh}h_{t-1}) = h_t$, recurrent unit.

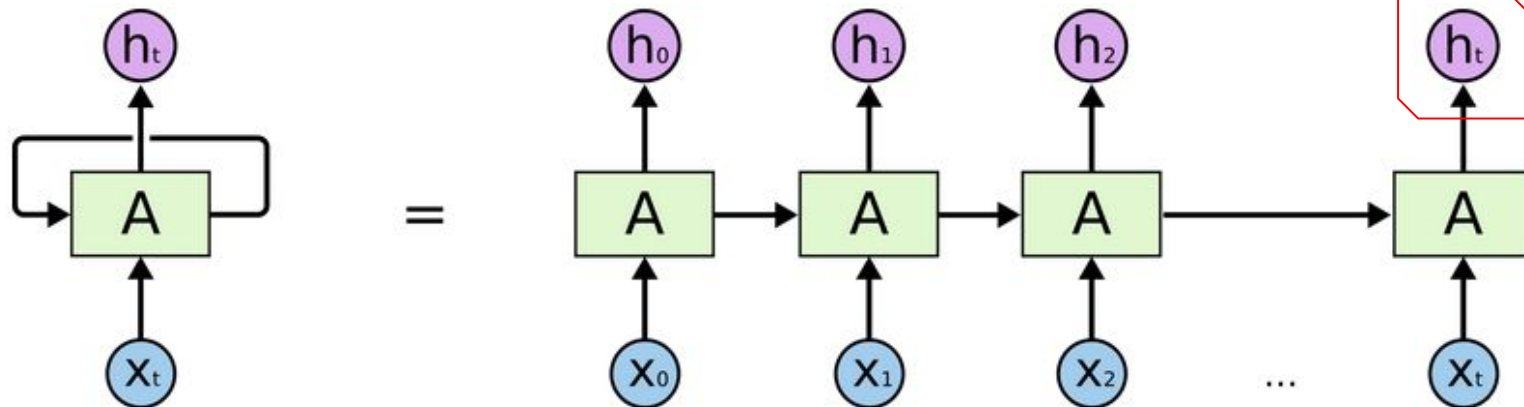
$y = W^{yh}h_t$ - output, e.g. logits for intent classification.



Recurrent neural network

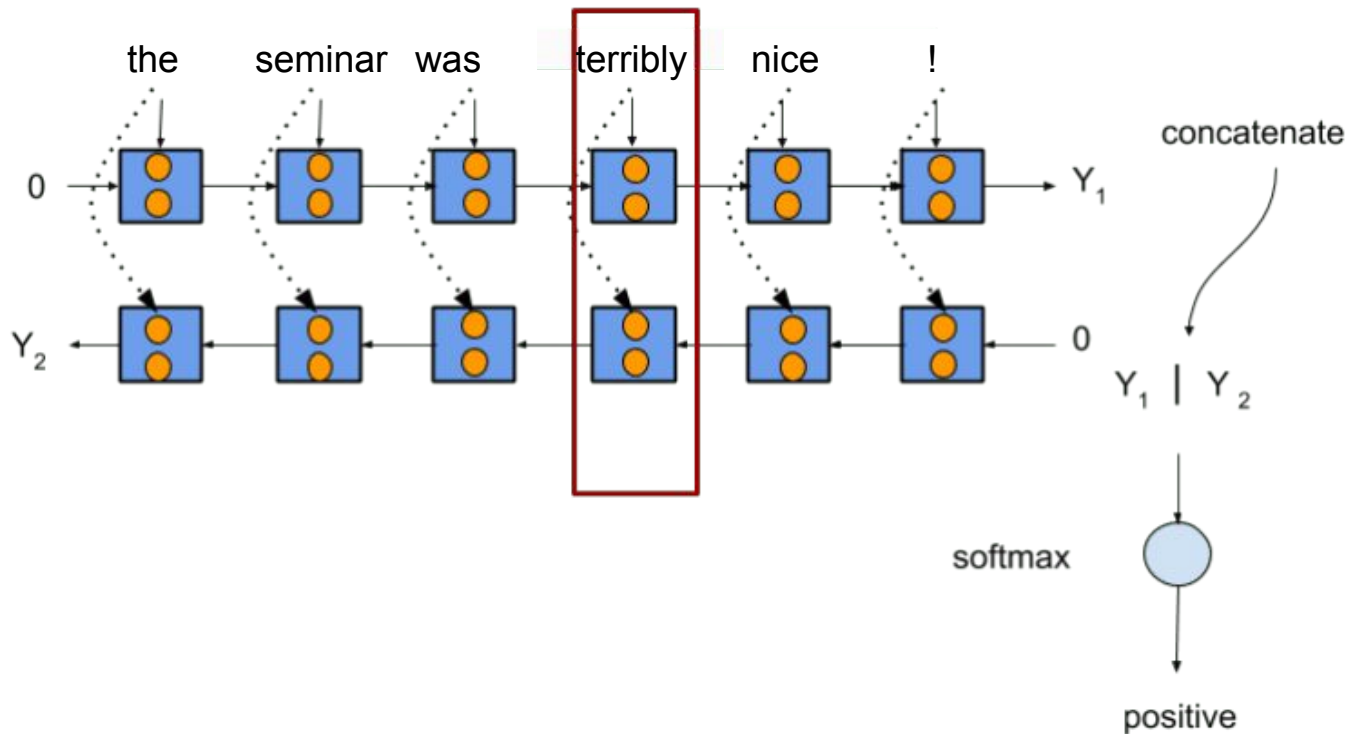
1. $W^{hx}x = \text{FC}$.
2. $\text{sigm}(W^{hx}x) = \text{FC} + \text{activation}$.
3. $\text{sigm}(W^{hx}x_t + W^{hh}h_{t-1}) = h_t$, recurrent unit.

$y = W^{yh}h_t$ - output, e.g. logits for intent classification.

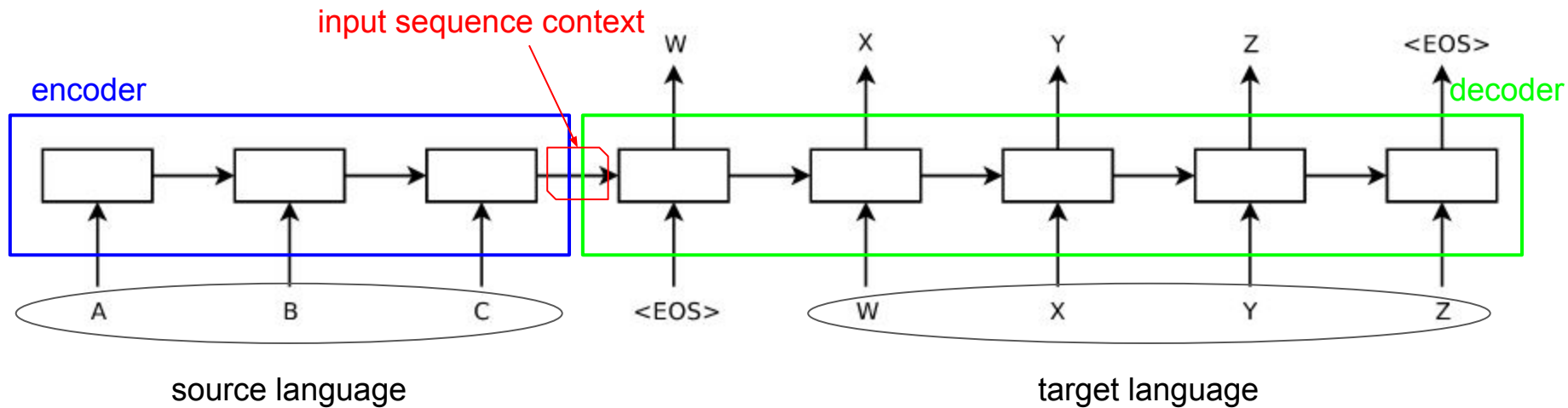


Bidirectional recurrent neural network

Consider review: “The seminar was terribly nice!”, the task is to predict sentiment.

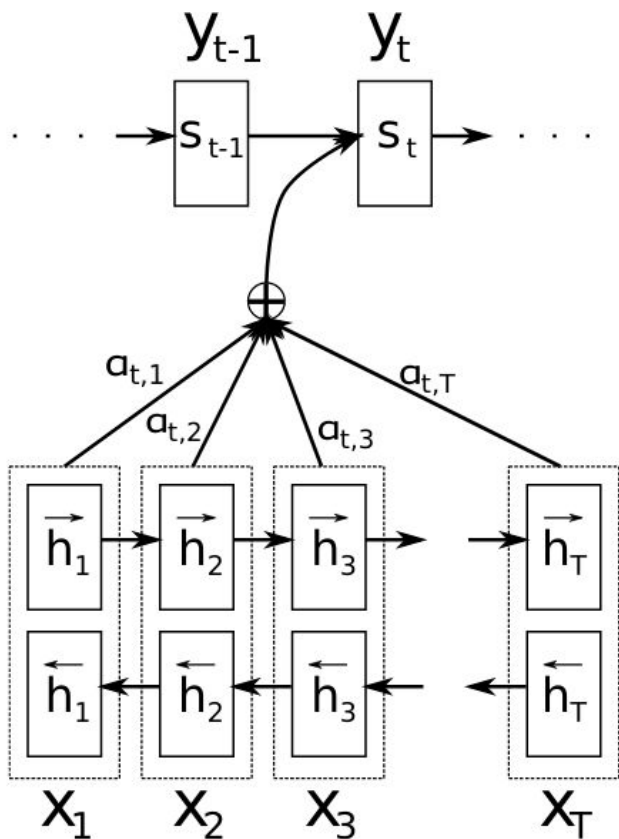


Seq2Seq: text translation



Encoder and decoder are RNNs (LSTM, GRU).

Seq2Seq: text translation with attention



Own context vector c_i for each target y_i .

$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j$$

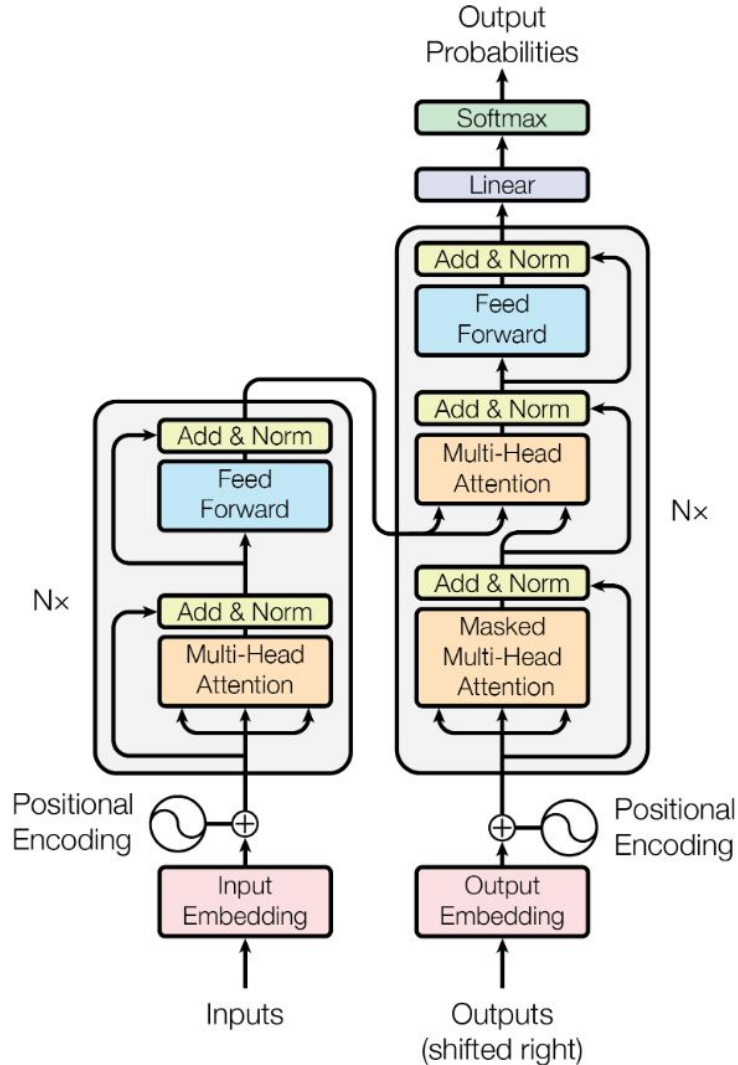
$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})}$$

$$e_{ij} = a(s_{i-1}, h_j)$$

$a(s, h)$ - model (fc) which scores how well the inputs around position j and the output at position i match.

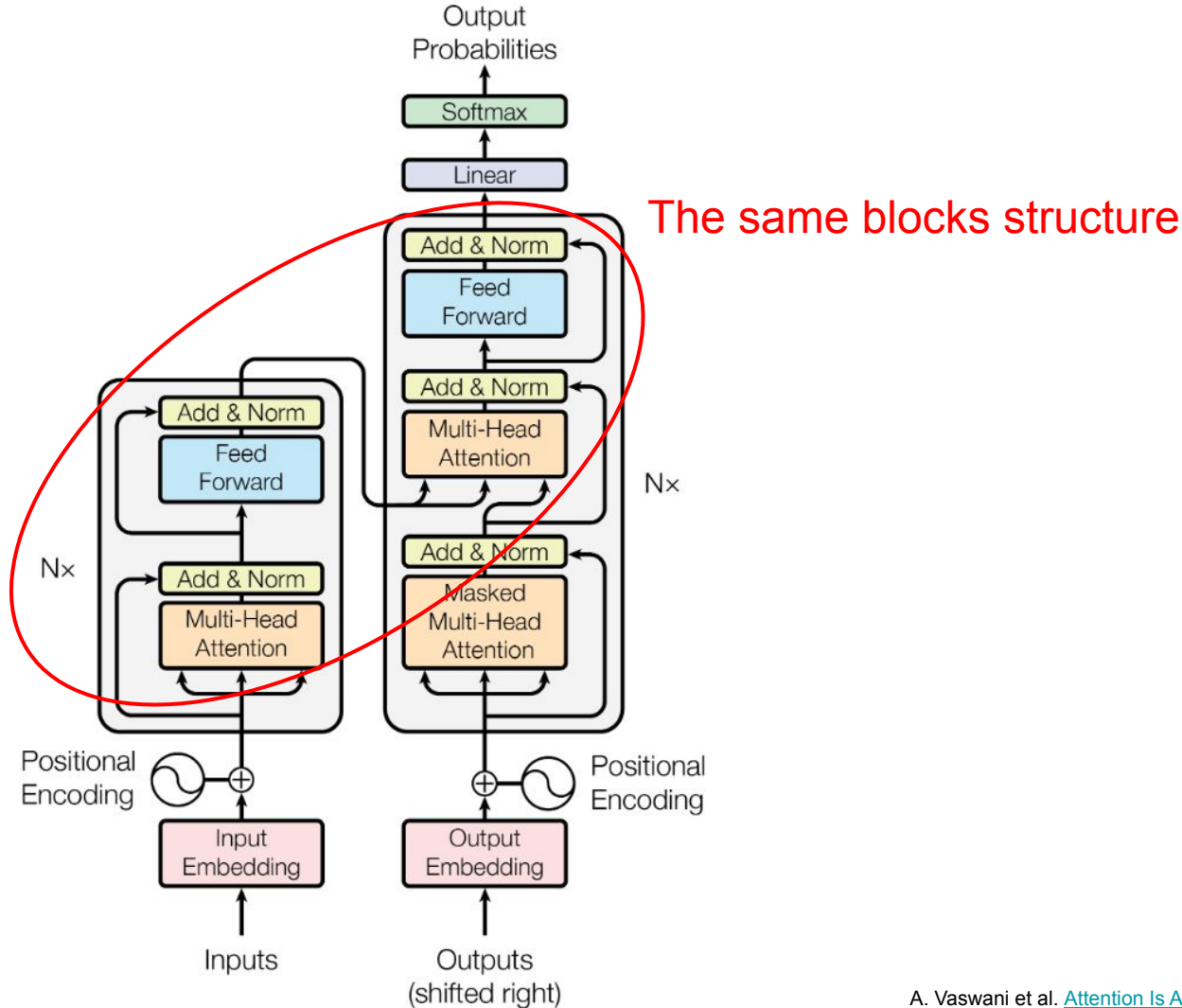
Transformer

Looks simple, huh?

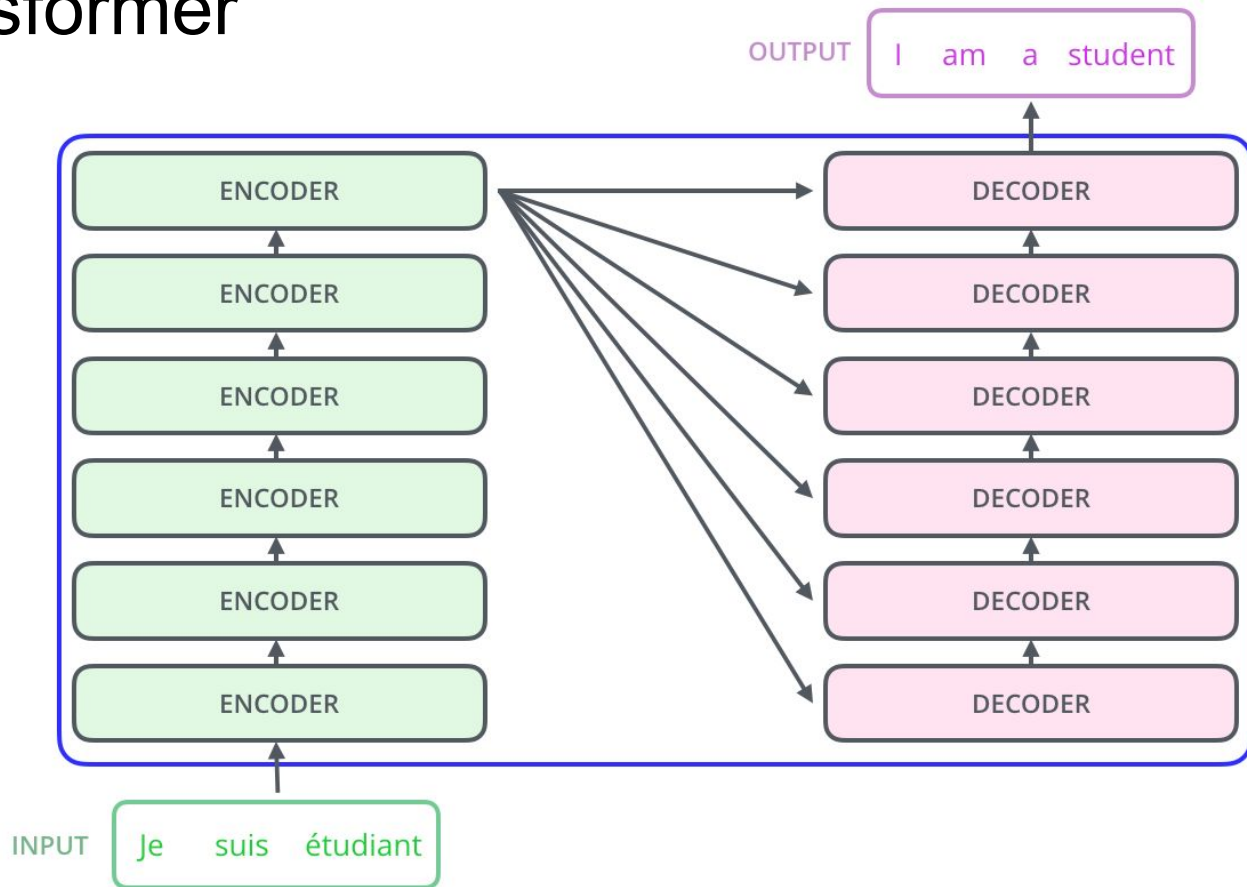


Transformer

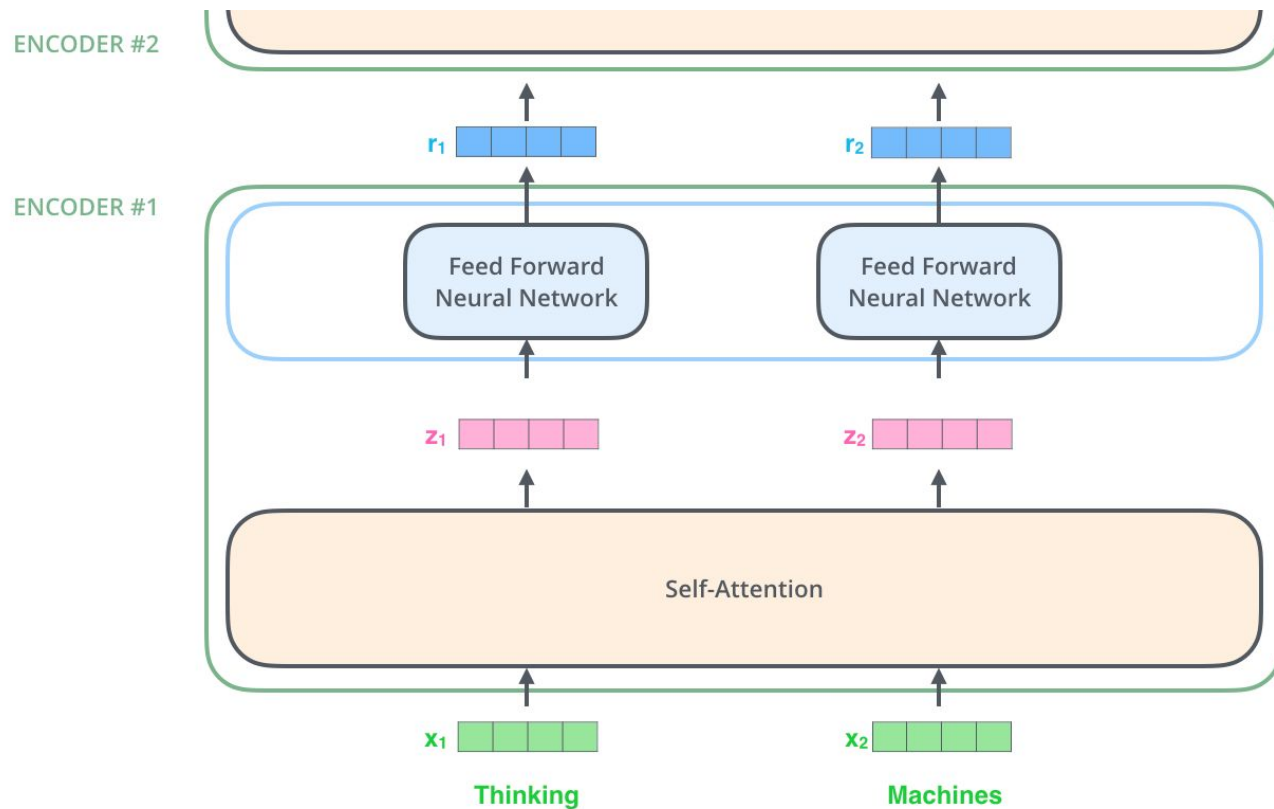
Looks simple, huh?



Transformer



Transformer: encoder



Transformer: self-attention

Input: $X = (x_1; x_2; \dots;$

$x_{\text{max_num_words}})$,

$X \in \mathbb{R}^{\text{max_num_words} \times \text{embedding_dim}}$

Input

Embedding

Queries

Keys

Values

Score

Divide by 8 ($\sqrt{d_k}$)

Softmax

Softmax

X

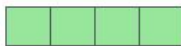
Value

Sum

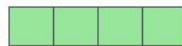
Thinking

Machines

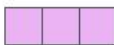
x_1



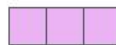
x_2



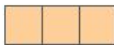
q_1



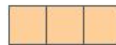
q_2



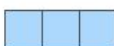
k_1



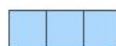
k_2



v_1



v_2



Embedding
projections
with learnable
weights (just
matmul)

$q_1 \cdot k_1 = 112$

$q_1 \cdot k_2 = 96$

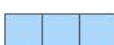
14

12

0.88

0.12

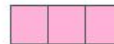
v_1



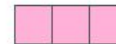
v_2



z_1



z_2



Transformer: self-attention

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

Input

Embedding

Queries

Keys

Values

Score

Divide by 8 ($\sqrt{d_k}$)

Softmax

Softmax

X

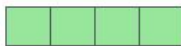
Value

Sum

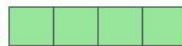
Thinking

Machines

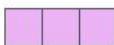
x_1



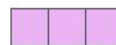
x_2



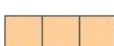
q_1



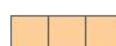
q_2



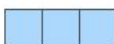
k_1



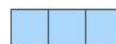
k_2



v_1



v_2



Embedding
projections
with learnable
weights (just
matmul)

$$q_1 \cdot k_1 = 112$$

$$q_1 \cdot k_2 = 96$$

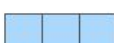
14

12

0.88

0.12

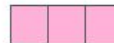
v_1



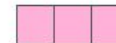
v_2



z_1



z_2



Transformer: self-attention

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

It is bidirectional.

Input

Embedding

Queries

Keys

Values

Score

Divide by 8 ($\sqrt{d_k}$)

Softmax

Softmax

X

Value

Sum

Thinking

Machines

x_1

x_2

q_1

q_2

k_1

k_2

v_1

v_2

Embedding
projections
with learnable
weights (just
matmul)

$q_1 \cdot k_1 = 112$

$q_1 \cdot k_2 = 96$

14

12

0.88

0.12

v_1

v_2

z_1

z_2

Transformer: self-attention

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

It is bidirectional.

It is multi-headed.

Input

Embedding

Queries

Keys

Values

Score

Divide by 8 ($\sqrt{d_k}$)

Softmax

Softmax

X

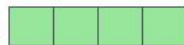
Value

Sum

Thinking

Machines

x_1 

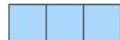
x_2 

q_1 

q_2 

k_1 

k_2 

v_1 

v_2 

Embedding
projections
with learnable
weights (just
matmul)

$q_1 \cdot k_1 = 112$

$q_1 \cdot k_2 = 96$

14

12

0.88

0.12

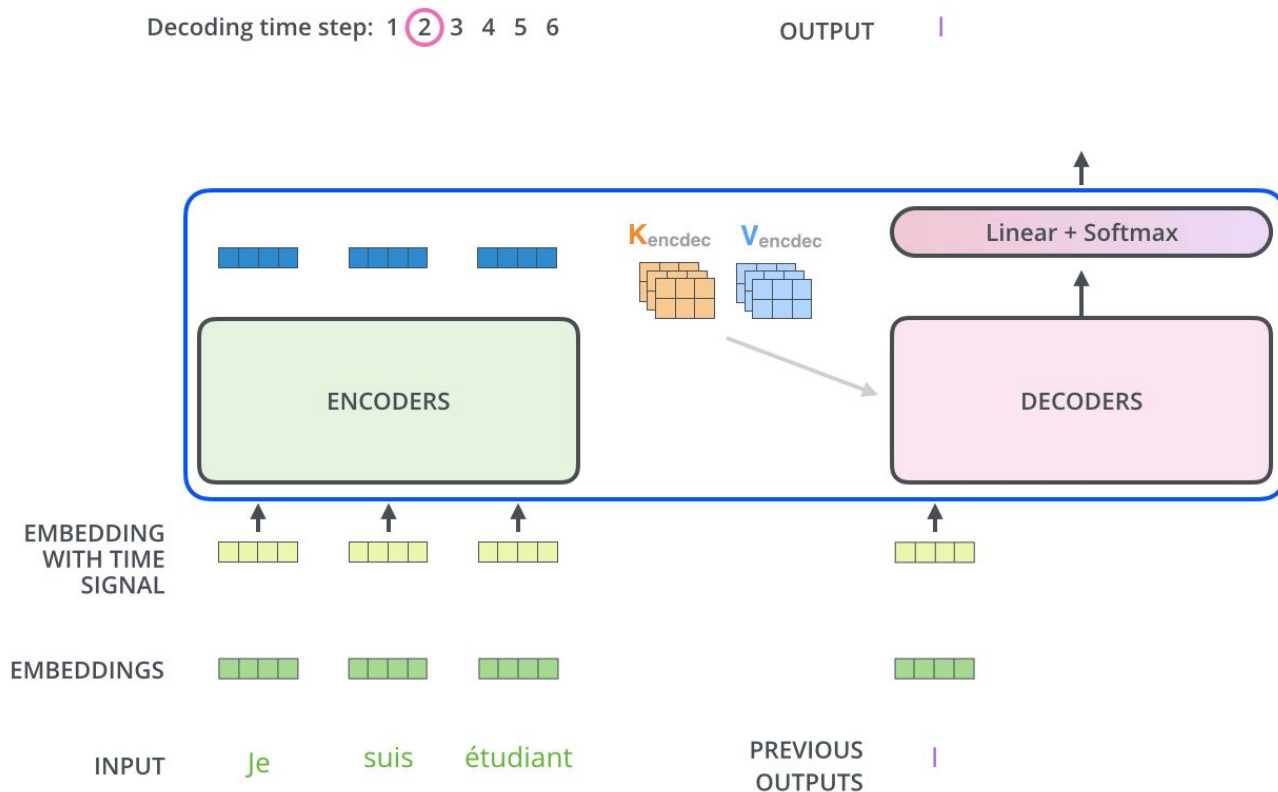
v_1 

v_2 

z_1 

z_2 

Transformer: decoder

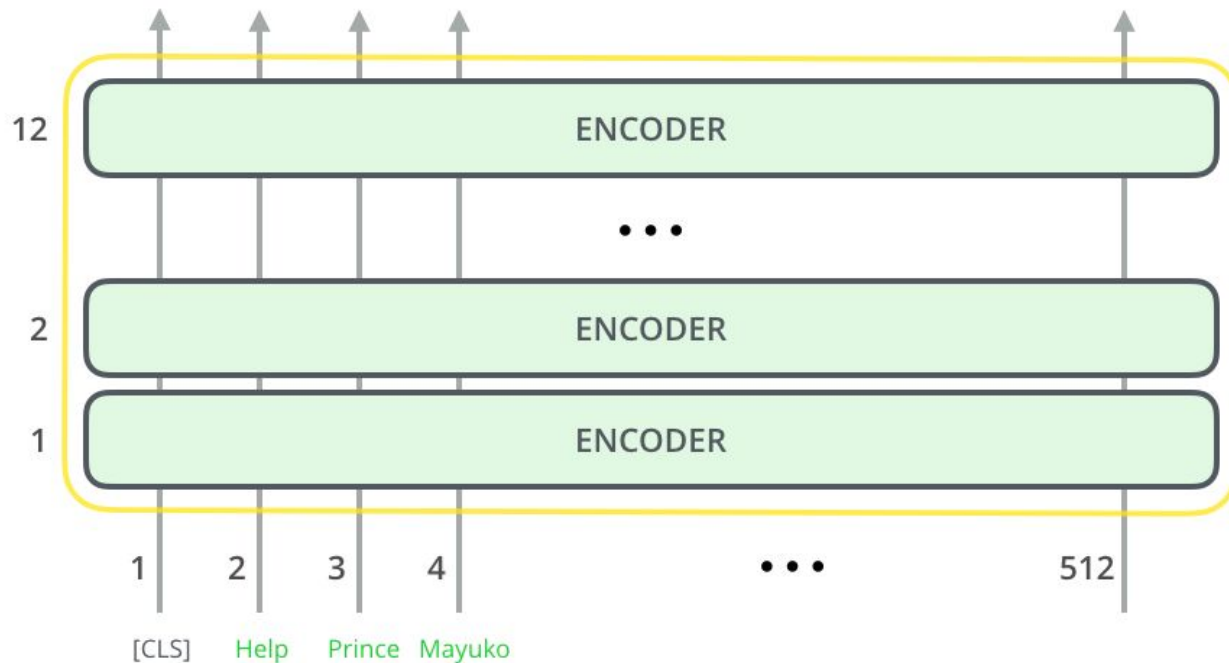


BERT: Bidirectional Encoder Representations from Transformers

How many transformers in BERT?

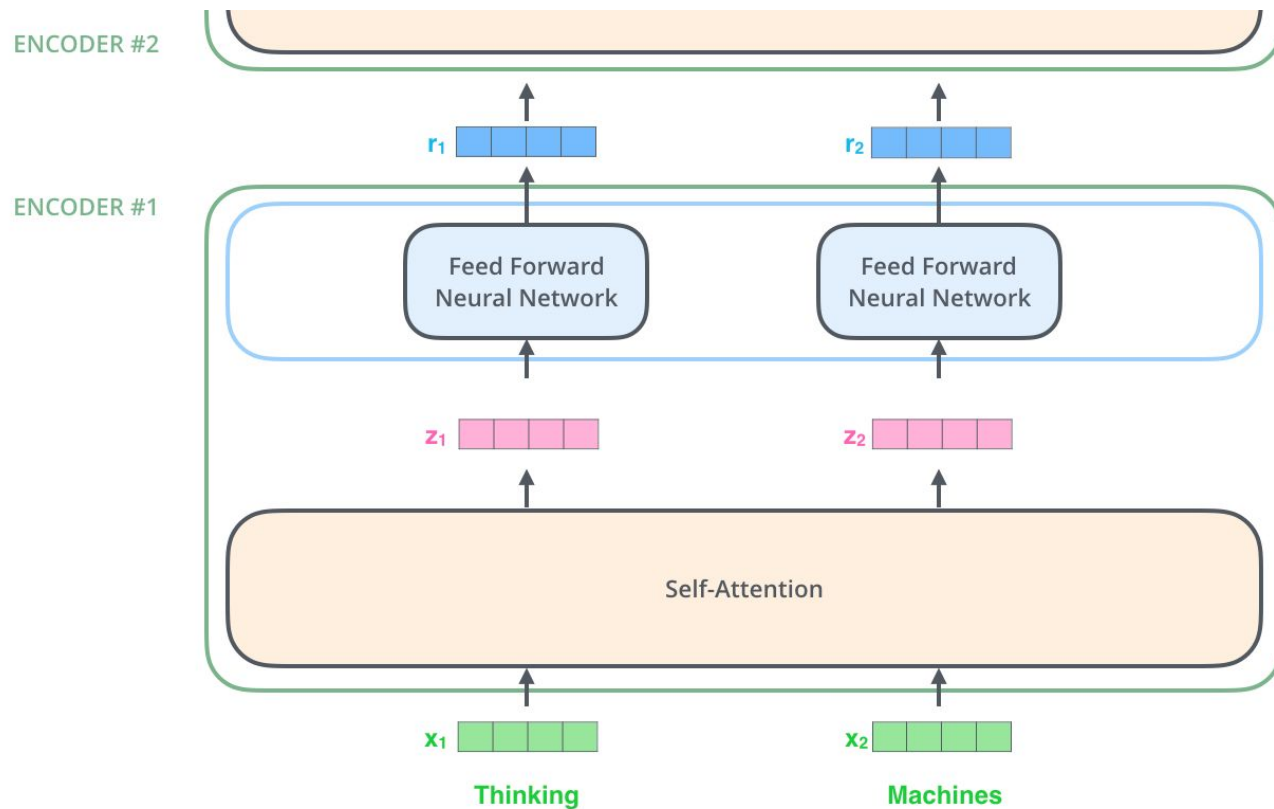


BERT: Bidirectional Encoder Representations from Transformers

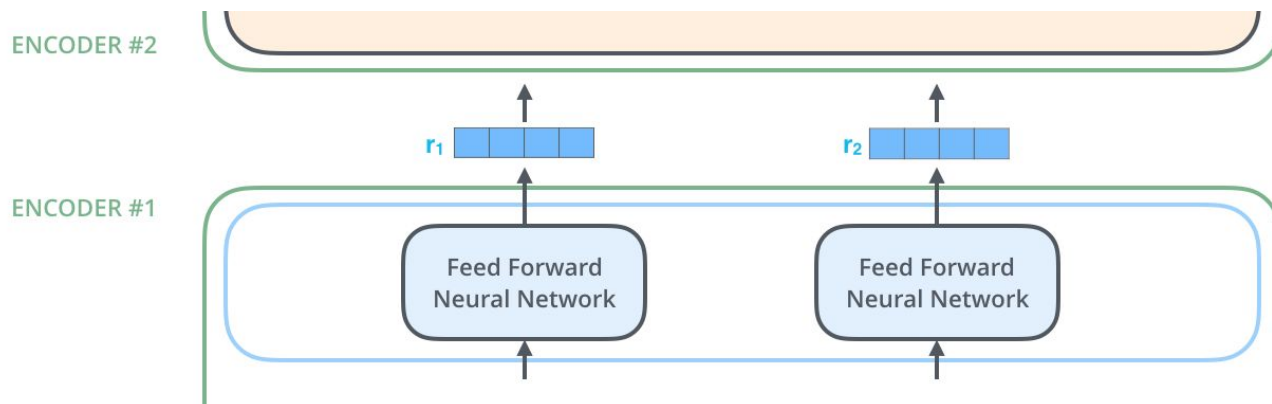


BERT

Transformer: encoder



Transformer: encoder



BERT Base: number of encoder layers = 12, Total Parameters = 110M

BERT Large: number of encoder layers = 24, Total Parameters = 340M

x1
Thinking

x2
Machines

BERT text mapping example

Tokenizer: split text into subwords.

```
from transformers import BertTokenizer, RobertaTokenizer
if __name__ == '__main__':
    sentence = 'On the way to Malaysia...no internet access to Twit'
    for tokenizer_class, tokenizer_name in zip(
        [BertTokenizer, RobertaTokenizer],
        ['bert-base-uncased', 'roberta-base']):
        tokenizer = tokenizer_class.from_pretrained(tokenizer_name)
        tokens = tokenizer.tokenize(sentence)
        print('{}:'.format(tokenizer_name))
        print(tokens)
bert-base-uncased:
['on', 'the', 'way', 'to', 'malaysia', '.', '.', '.', 'no', 'internet', 'access', 'to', 't', '##wi', '##t']
roberta-base:
['On', 'Ġthe', 'Ġway', 'Ġto', 'ĠMalaysia', 'Ġ...', 'no', 'Ġinternet', 'Ġaccess', 'Ġto', 'ĠTw', 'it']
```

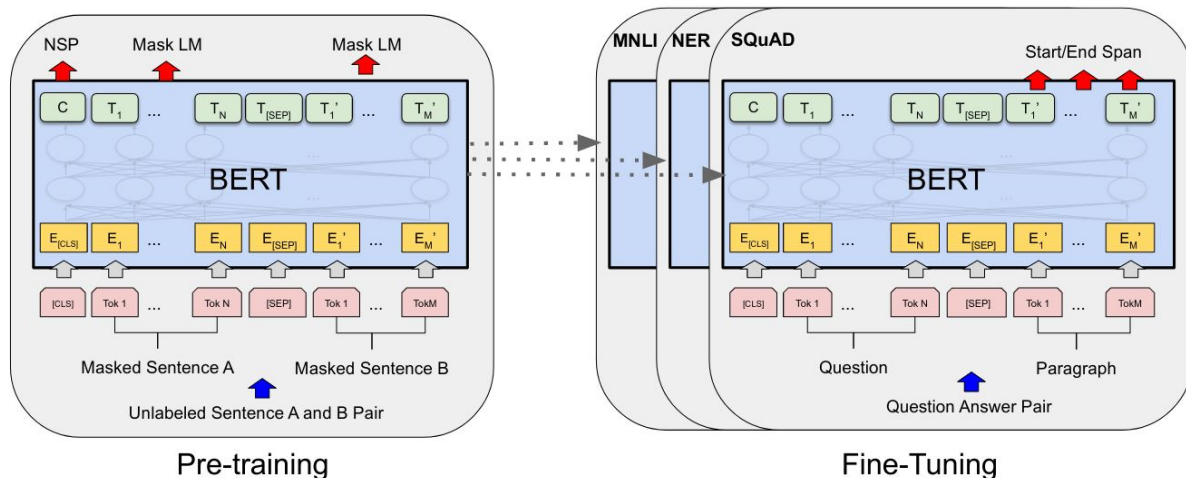
BERT: pre-training

Masked language modeling:

- Mask some percentage of the input tokens at random, and then predict those masked tokens.

Next sentence prediction:

- Understand the relationship between two sentences. Binary classification if in a given sentences pair the second sentence follows the first.



BERT: fine-tuning

```
class TweetModel(transformers.BertPreTrainedModel):
    def __init__(self, conf):
        super().__init__(conf)
        self.bert = transformers.BertModel.from_pretrained('bert-base-uncased', config=conf)
        self.drop_out = nn.Dropout(0.1)
        self.l0 = nn.Linear(768 * 2, 2)

    def forward(self, tokens_ids, mask, tokens_type_ids):
        last_hidden_state, pooler_output, hidden_states = self.bert(
            tokens_ids,
            attention_mask=mask,
            token_type_ids=tokens_type_ids
        )
```

BERT: fine-tuning

```
class TweetModel(transformers.BertPreTrainedModel):
    def __init__(self, conf):
        super().__init__(conf)
        self.bert = transformers.BertModel.from_pretrained('bert-base-uncased', config=conf)
        self.drop_out = nn.Dropout(0.1)
        self.l0 = nn.Linear(768 * 2, 2)

    def forward(self, tokens_ids, mask, tokens_type_ids):
        last_hidden_state, pooler_output, hidden_states = self.bert(
            tokens_ids,
            attention_mask=mask,
            token_type_ids=tokens_type_ids
        )
```

Input text tokens
(batch_size x max_tokens_num)

BERT: fine-tuning

```
class TweetModel(transformers.BertPreTrainedModel):
    def __init__(self, conf):
        super().__init__(conf)
        self.bert = transformers.BertModel.from_pretrained('bert-base-uncased', config=conf)
        self.drop_out = nn.Dropout(0.1)
        self.l0 = nn.Linear(768 * 2, 2)

    def forward(self, tokens_ids, mask, tokens_type_ids):
        last_hidden_state, pooler_output, hidden_states = self.bert(
            tokens_ids,
            attention_mask=mask,
            token_type_ids=tokens_type_ids
        )
```

Mask with 1 at positions of real tokens and 0 for padding
(batch_size x max_tokens_num)

BERT: fine-tuning

```
class TweetModel(transformers.BertPreTrainedModel):
    def __init__(self, conf):
        super().__init__(conf)
        self.bert = transformers.BertModel.from_pretrained('bert-base-uncased', config=conf)
        self.drop_out = nn.Dropout(0.1)
        self.l0 = nn.Linear(768 * 2, 2)

    def forward(self, tokens_ids, mask, tokens_type_ids):
        last_hidden_state, pooler_output, hidden_states = self.bert(
            tokens_ids,
            attention_mask=mask,
            token_type_ids=tokens_type_ids
        )
```

Mask to identify different input sequences with 0 for 1st and 1 for 2nd sequence (batch_size x max_tokens_num)

BERT: fine-tuning

```
class TweetModel(transformers.BertPreTrainedModel):
    def __init__(self, conf):
        super().__init__(conf)
        self.bert = transformers.BertModel.from_pretrained('bert-base-uncased', config=conf)
        self.drop_out = nn.Dropout(0.1)
        self.l0 = nn.Linear(768 * 2, 2)

    def forward(self, tokens_ids, mask, tokens_type_ids):
        last_hidden_state, pooler_output, hidden_states = self.bert(
            tokens_ids,
            attention_mask=mask,
            token_type_ids=tokens_type_ids
        )
```

← Last encoder output
(batch_size x max_tokens_num x embedding_dim)

BERT: fine-tuning

```
class TweetModel(transformers.BertPreTrainedModel):
    def __init__(self, conf):
        super().__init__(conf)
        self.bert = transformers.BertModel.from_pretrained('bert-base-uncased', config=conf)
        self.drop_out = nn.Dropout(0.1)
        self.l0 = nn.Linear(768 * 2, 2)

    def forward(self, tokens_ids, mask, tokens_type_ids):
        last_hidden_state, pooler_output, hidden_states = self.bert(
            tokens_ids,
            attention_mask=mask,
            token_type_ids=tokens_type_ids
        )
```

Processed last encoder output for the 1st token
(batch_size x embedding_dim)

BERT: fine-tuning

```
class TweetModel(transformers.BertPreTrainedModel):
    def __init__(self, conf):
        super().__init__(conf)
        self.bert = transformers.BertModel.from_pretrained('bert-base-uncased', config=conf)
        self.drop_out = nn.Dropout(0.1)
        self.l0 = nn.Linear(768 * 2, 2)

    def forward(self, tokens_ids, mask, tokens_type_ids):
        last_hidden_state, pooler_output, hidden_states = self.bert(
            tokens_ids,
            attention_mask=mask,
            token_type_ids=tokens_type_ids
        )
```

Output from embeddings and all encoders
Tuple of tensors (batch_size x max_tokens_num x embedding_dim)

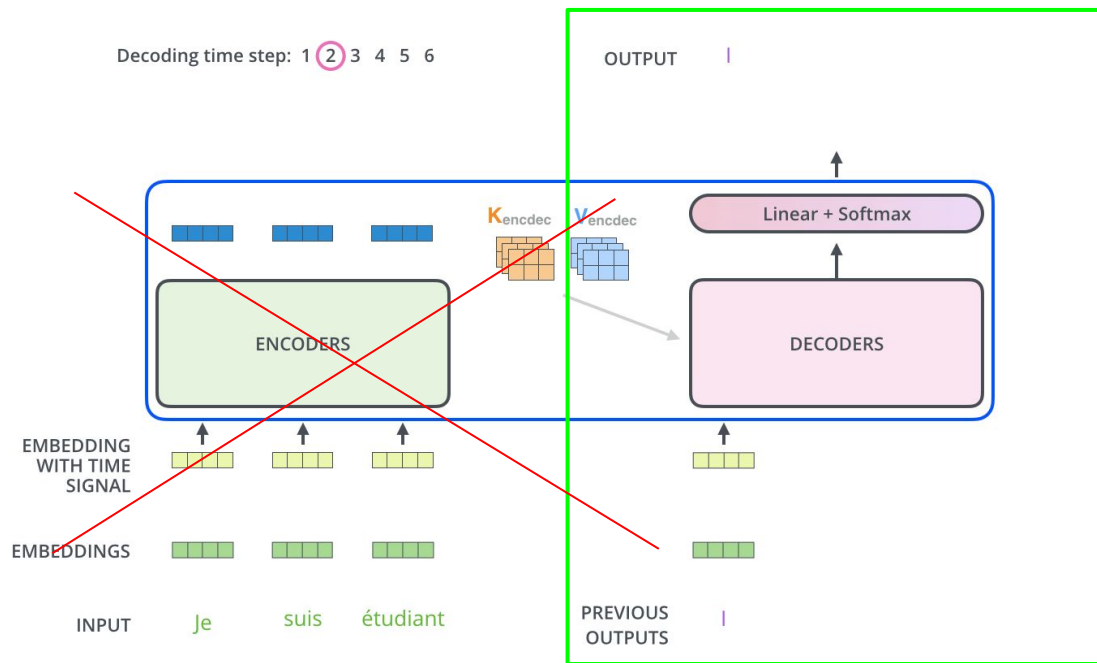
GPT: Improving Language Understanding by Generative Pre-Training

GPT: Improving Language Understanding by Generative Pre-Training

It is a decoder from transformer.

Try it out (GPT-2):

<https://demo.allennlp.org/n-ext-token-lm?text=Lobachevsky>



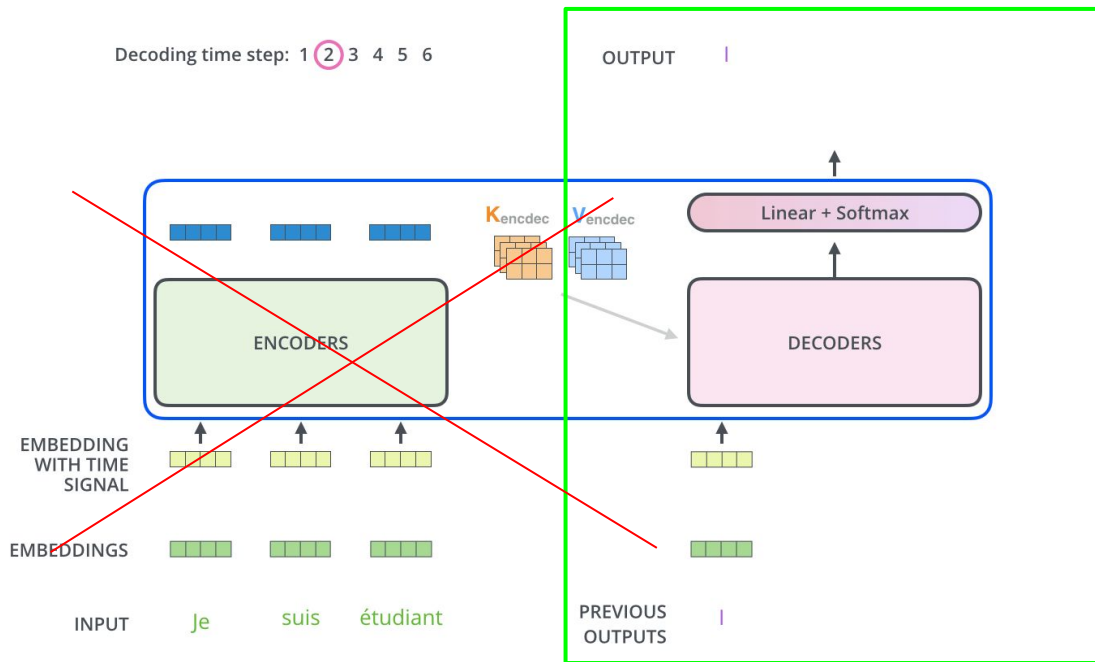
GPT: Improving Language Understanding by Generative Pre-Training

It is a decoder from transformer.

Try it out (GPT-2):

<https://demo.allennlp.org/n-ext-token-lm?text=Lobachevsky>

GPT-3 has 175B parameters.



Summary

- Seq2Seq: encode all input text into single context vector.
- Attention: make context vector adaptive to position of output word.
- Self-attention: shorten distance between words.
- BERT: large scale pre-training.

Thank you for the attention!



- <https://towardsdatascience.com/lost-in-translation-found-by-transformer-46a16bf6418f>
- <https://towardsdatascience.com/bert-for-dummies-step-by-step-tutorial-fb90890ffe03>
- <https://mccormickml.com/2019/07/22/BERT-fine-tuning/>
- <https://www.analyticsvidhya.com/blog/2019/06/understanding-transformers-nlp-state-of-the-art-models/>
- <http://jalamar.github.io/illustrated-gpt2/>
- https://pytorch.org/tutorials/intermediate/seq2seq_translation_tutorial.html
- <https://lilianweng.github.io/lil-log/2018/06/24/attention-attention.html>
- <https://ai.googleblog.com/2020/03/more-efficient-nlp-model-pre-training.html>
- <https://huggingface.co/transformers/>

