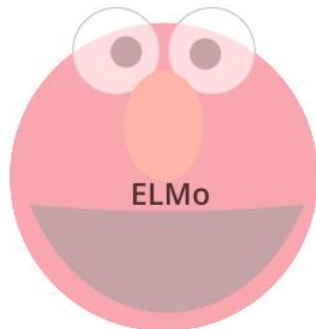# Lecture 6:
# How NLP Cracked Transfer Learning

**Anastasia Ianina**

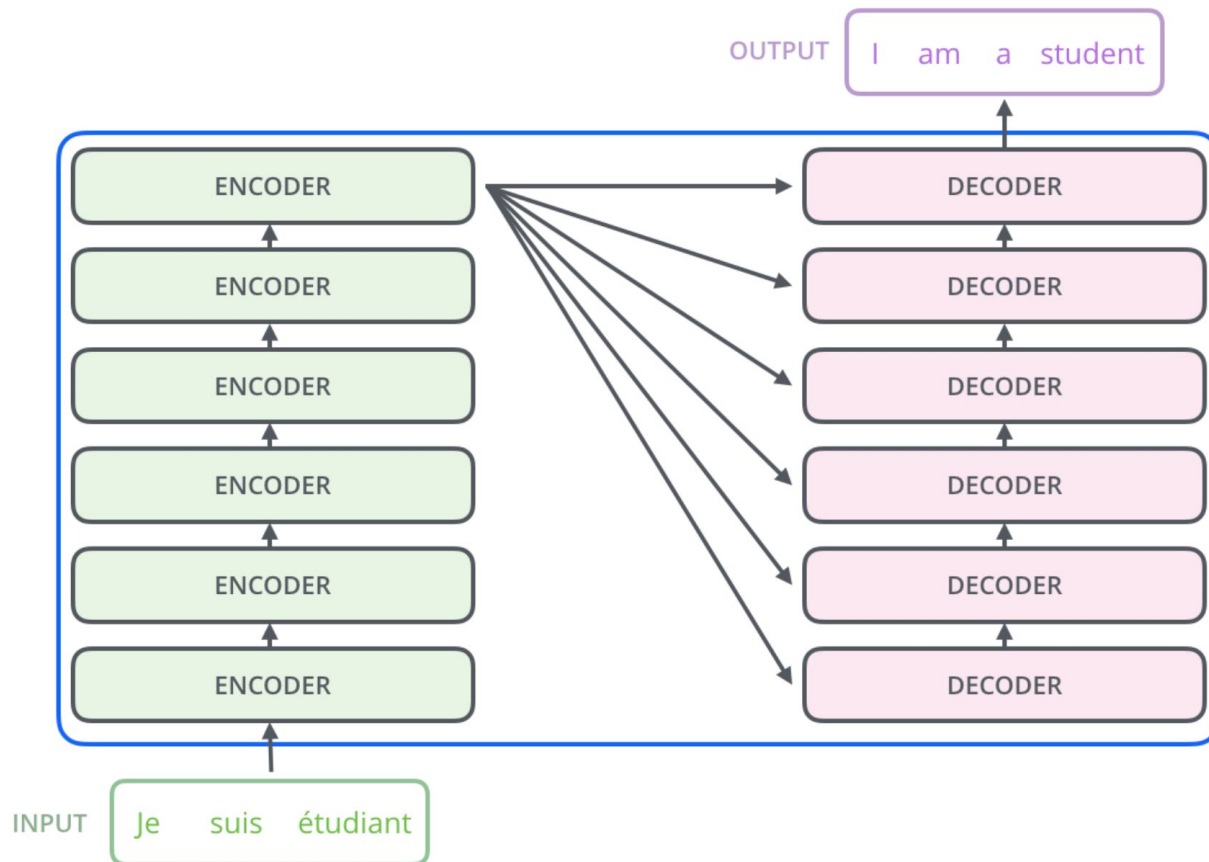Harbour.Space University
15.07.2019, Barcelona, Spain

# Outline

1. Transformer: recap
2. OpenAI Transformer
3. ELMO
4. BERT
5. Q & A

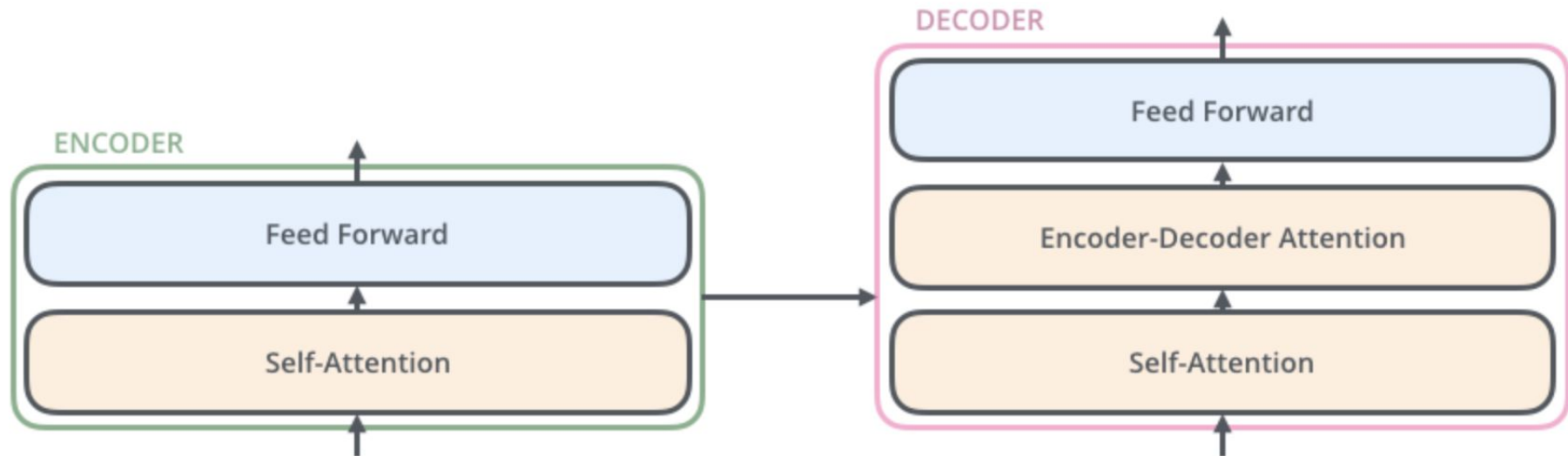Based on: http://web.stanford.edu/class/cs224n/slides/cs224n-2019-lecture13-contextual-representations.pdf

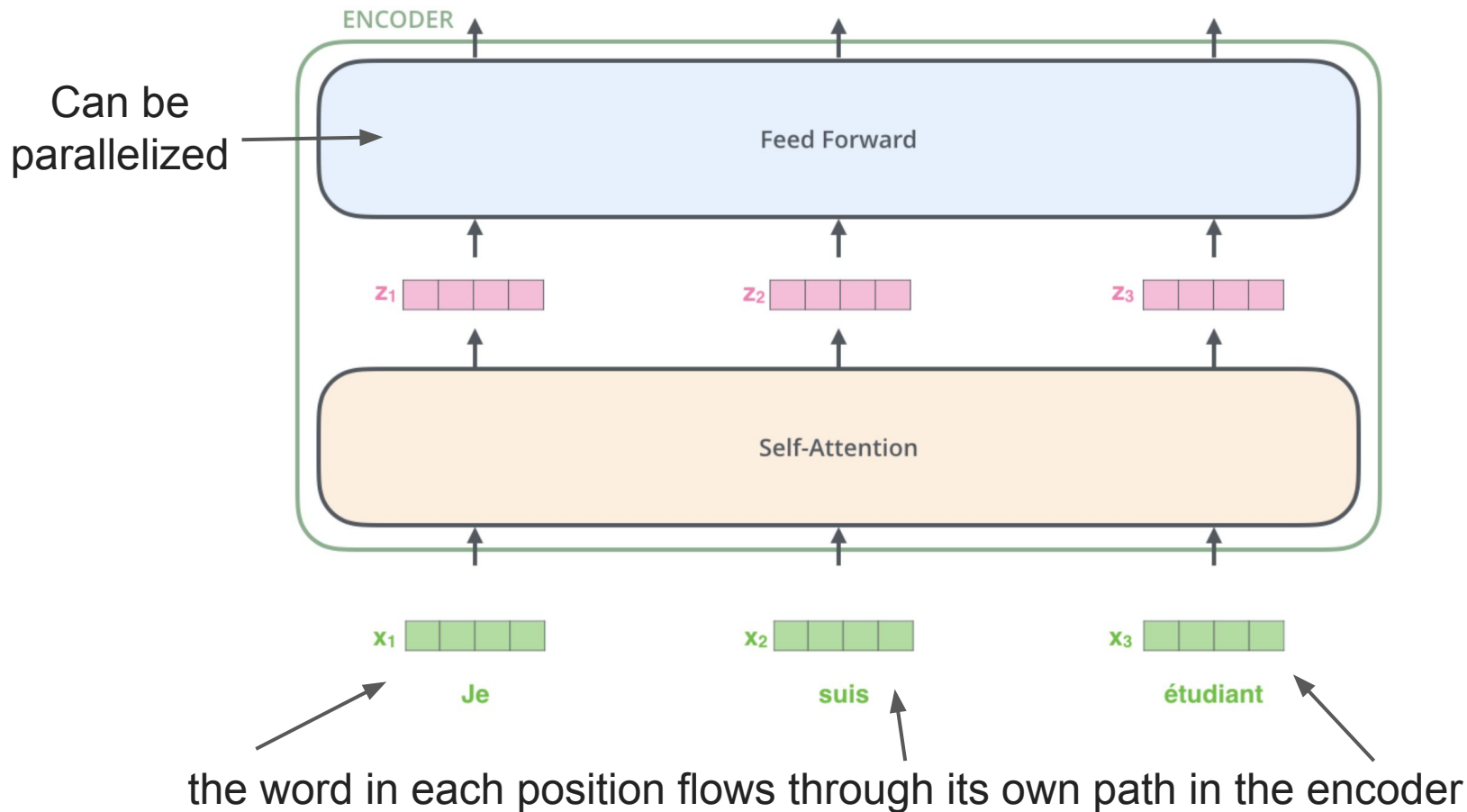https://jalammar.github.io/illustrated-transformer/
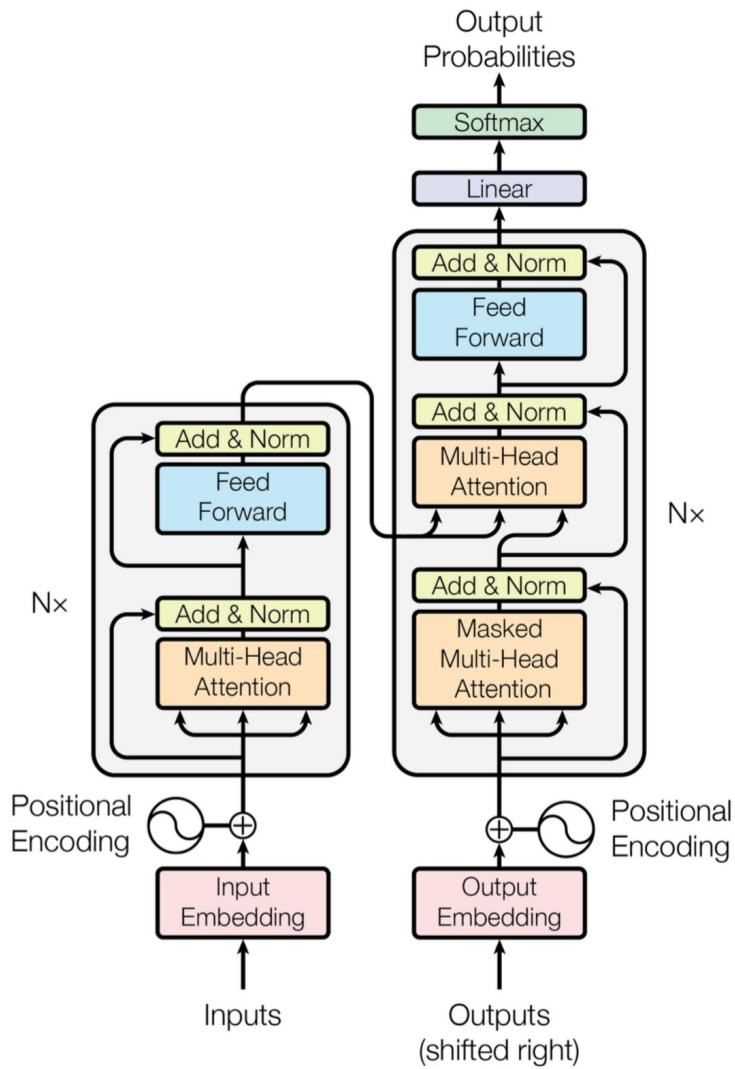
http://jalammar.github.io/illustrated-bert/

# The Transformer: recap

# The Transformer



OUTPUT: I am a student

ENCODER → DECODER (×6 layers each)

INPUT: Je suis étudiant

# The Transformer

ENCODER

Can be parallelized →

Feed Forward

$z_1$ | | | | 　$z_2$ | | | | 　$z_3$ | | | |

Self-Attention

$x_1$ | | | | 　$x_2$ | | | | 　$x_3$ | | | |

Je 　 suis 　 étudiant

the word in each position flows through its own path in the encoder

# The Transformer: recap



Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

Nx

Nx

Positional Encoding

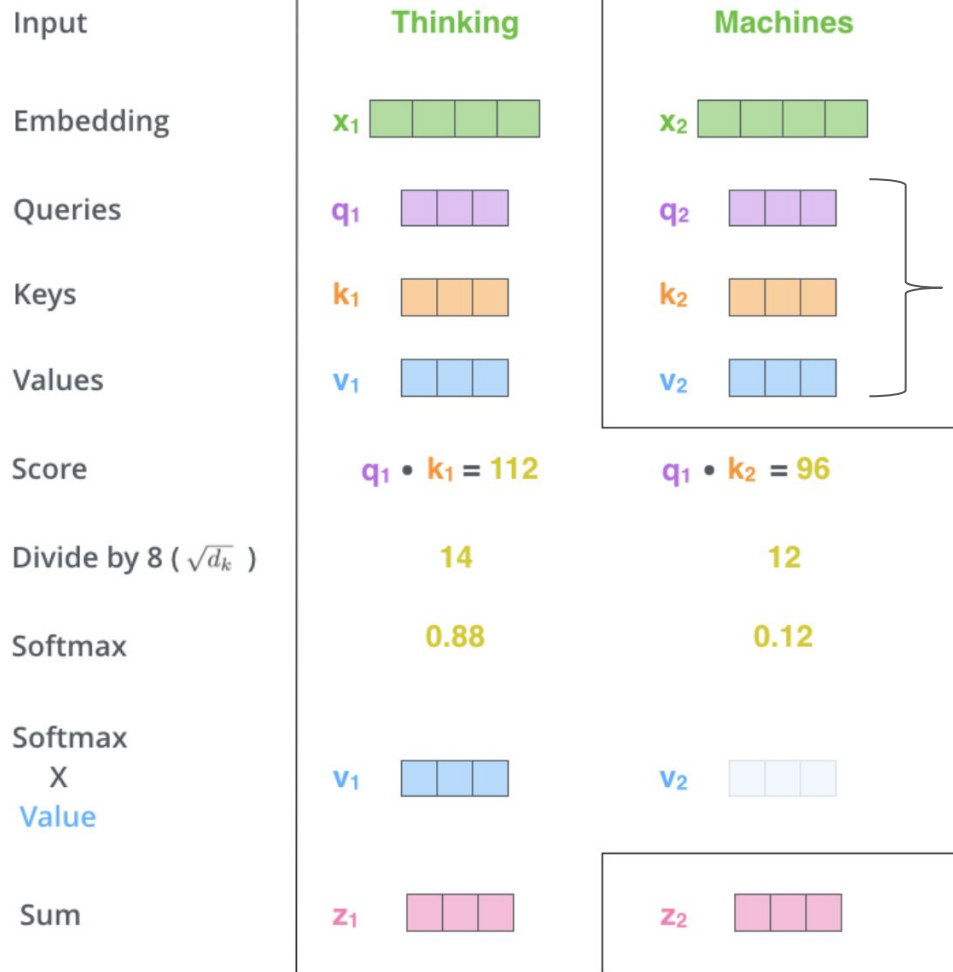Positional Encoding

Input Embedding

Output Embedding

Inputs

Outputs (shifted right)

- Proposed in the paper "Attention is All You Need" (Ashish Vaswani et al.)
- No recurrent or convolutional neural networks -> just attention
- Uses Multi-Head **self-attention** concept

8

# Self-Attention: recap

# Self-Attention in Detail

| | **Thinking** | **Machines** | |
|---|---|---|---|
| Input | | | |
| Embedding | $x_1$ | $x_2$ | |
| Queries | $q_1$ | $q_2$ | |
| Keys | $k_1$ | $k_2$ | **STEP 1:** create Query, Key, Value |
| Values | $v_1$ | $v_2$ | |
| Score | $q_1 \cdot k_1 = 112$ | $q_1 \cdot k_2 = 96$ | **STEP 2:** calculate scores |
| Divide by 8 ( $\sqrt{d_k}$ ) | 14 | 12 | **STEP 3:** divide by $\sqrt{d_k}$ |
| Softmax | 0.88 | 0.12 | **STEP 4:** softmax |
| Softmax X Value | $v_1$ | $v_2$ | **STEP 5:** multiply each value vector by the softmax score |
| Sum | $z_1$ | $z_2$ | **STEP 6:** sum up the weighted value vectors |

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

$$= Z$$

# Multi-Head Attention

**X**

Thinking
Machines

ATTENTION HEAD #0

$Q_0$
$W_0^Q$

$K_0$
$W_0^K$

$V_0$
$W_0^V$

ATTENTION HEAD #1

$Q_1$
$W_1^Q$

$K_1$
$W_1^K$

$V_1$
$W_1^V$

# Multi-Head Attention

# Multi-Head Attention

1) Concatenate all the attention heads

$Z_0$  $Z_1$  $Z_2$  $Z_3$  $Z_4$  $Z_5$  $Z_6$  $Z_7$

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

X

$W^O$

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

Z

=

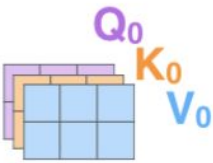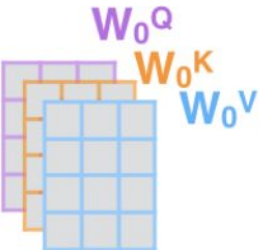**1) This is our input sentence***

**2) We embed each word***

**3) Split into 8 heads. We multiply $X$ or $R$ with weight matrices**

**4) Calculate attention using the resulting $Q$/$K$/$V$ matrices**

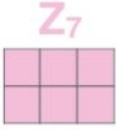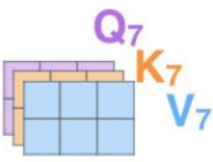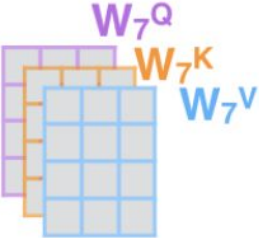**5) Concatenate the resulting $Z$ matrices, then multiply with weight matrix $W^O$ to produce the output of the layer**
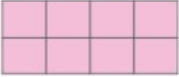
Thinking Machines

$X$
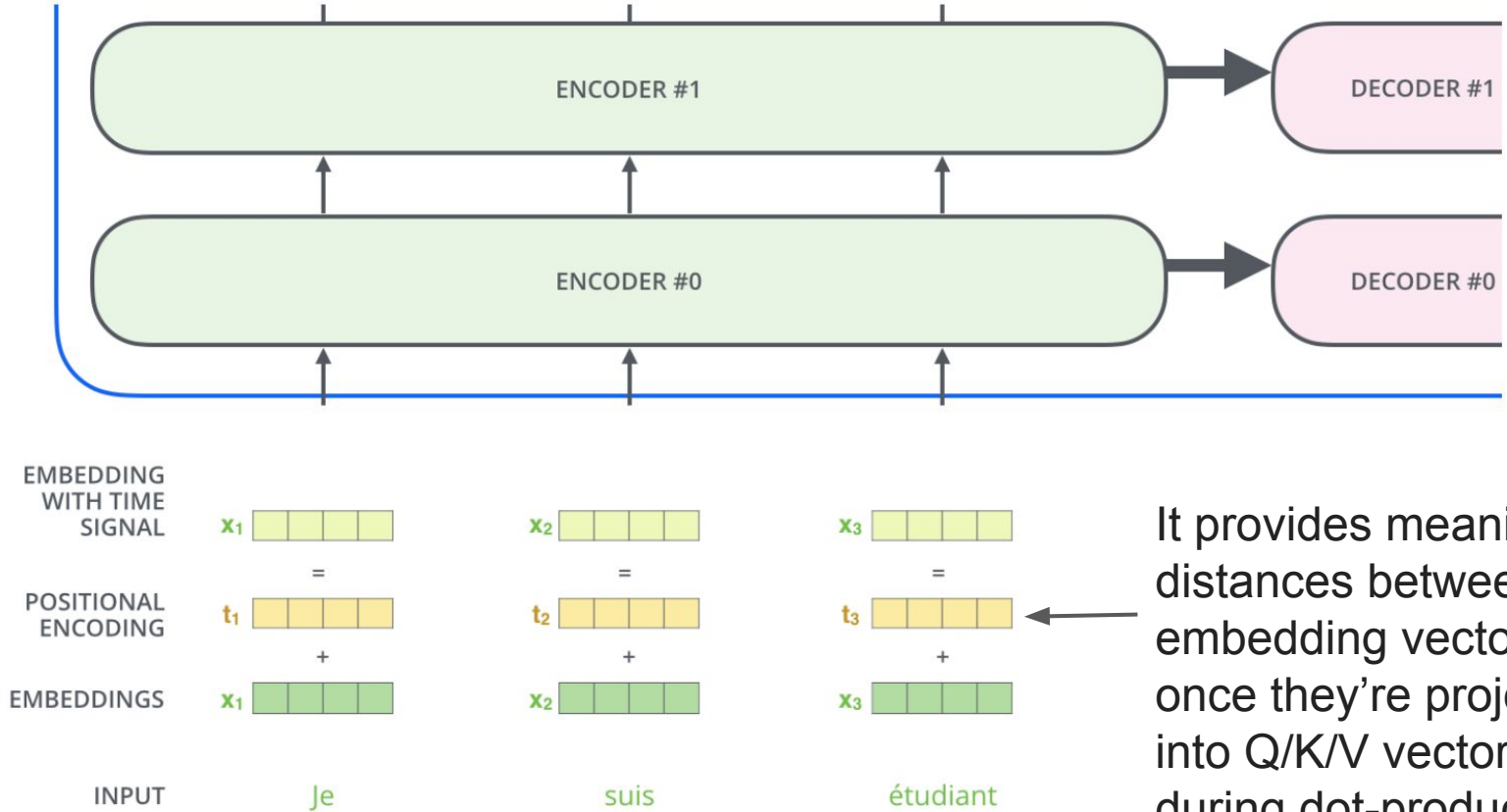
* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one

$R$

$W_0{}^Q$
$W_0{}^K$
$W_0{}^V$

$W_1{}^Q$
$W_1{}^K$
$W_1{}^V$

...

$W_7{}^Q$
$W_7{}^K$
$W_7{}^V$

$Q_0$
$K_0$
$V_0$

$Q_1$
$K_1$
$V_1$

...

$Q_7$
$K_7$
$V_7$

$Z_0$
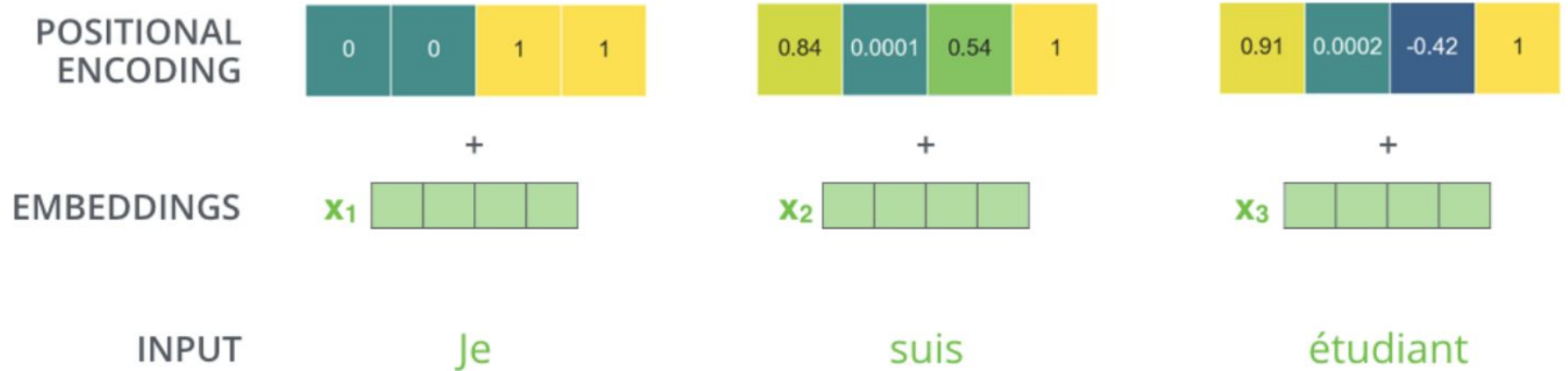
$Z_1$

...

$Z_7$

$W^O$

$Z$
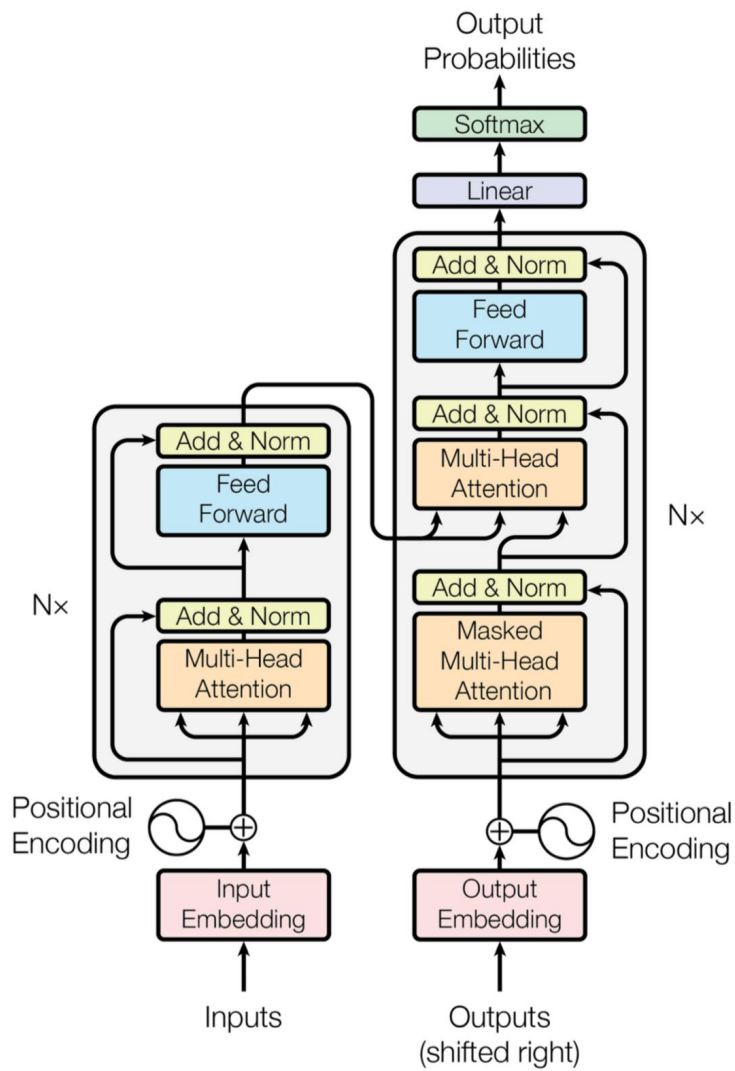
# Positional Encoding

# Positional Encoding



It provides meaningful distances between the embedding vectors once they're projected into Q/K/V vectors and during dot-product attention
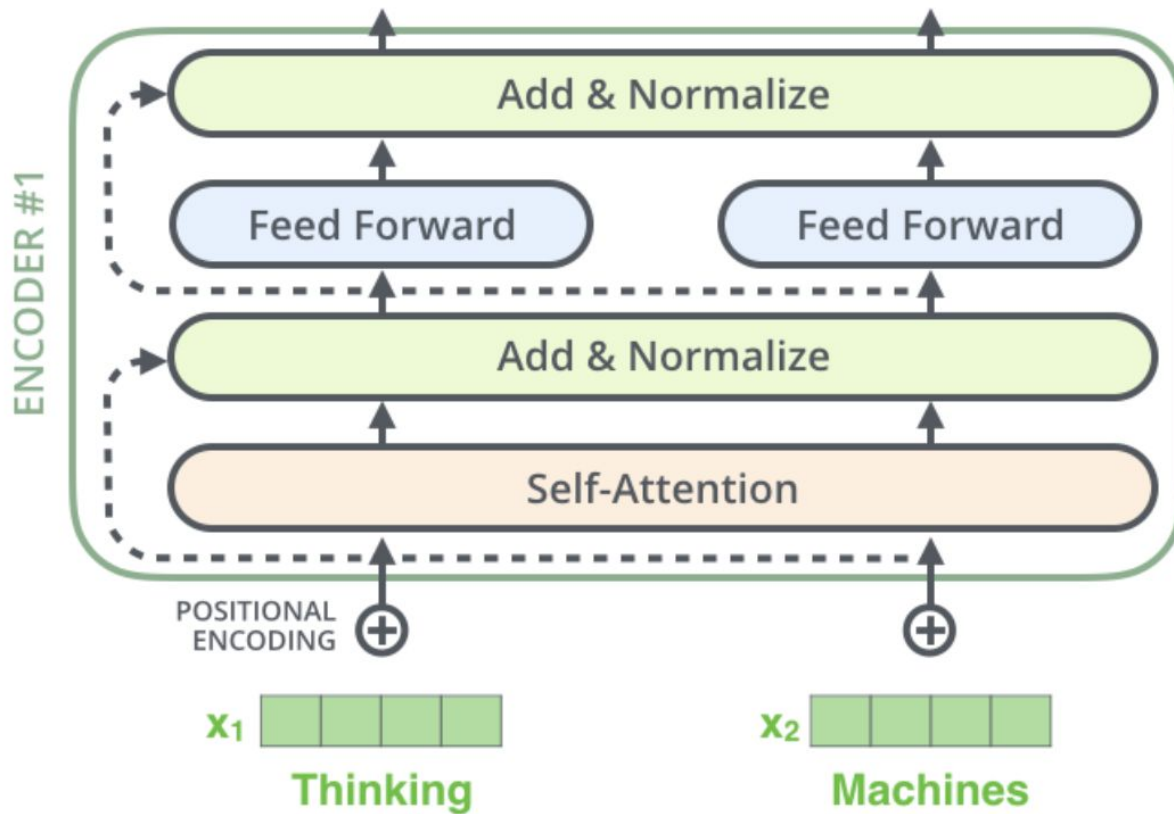
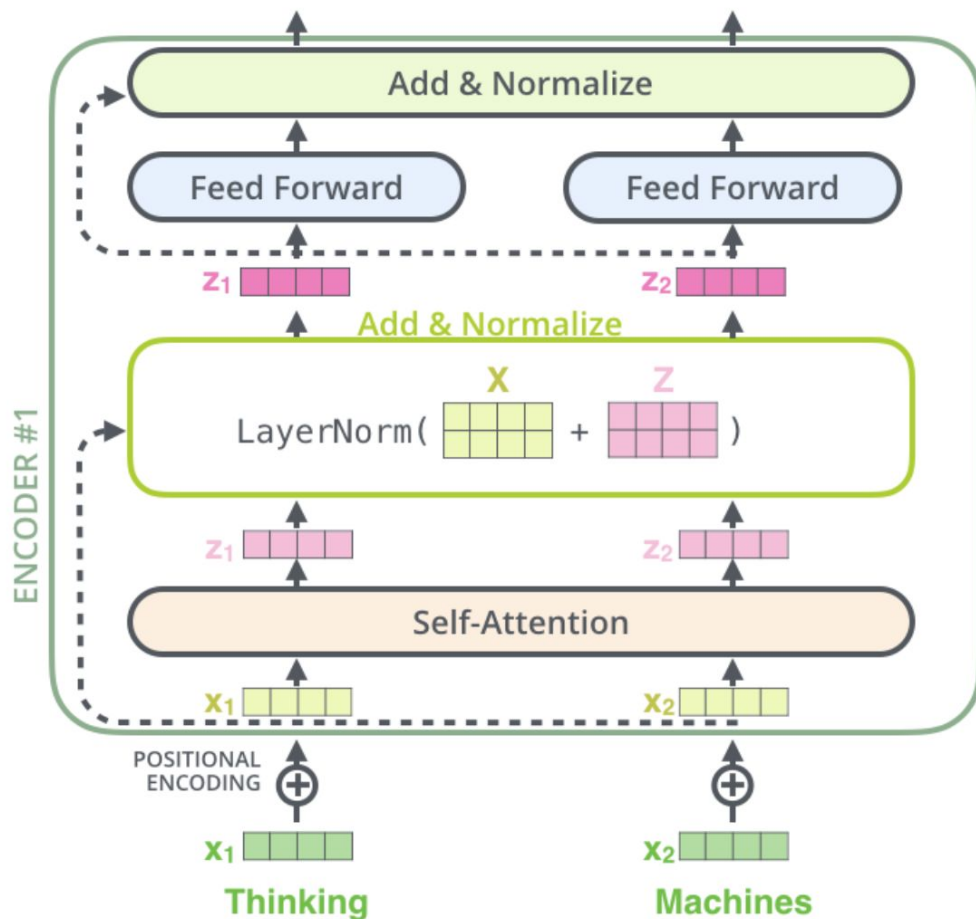# Positional Encoding

# Layer Normalization

Output Probabilities
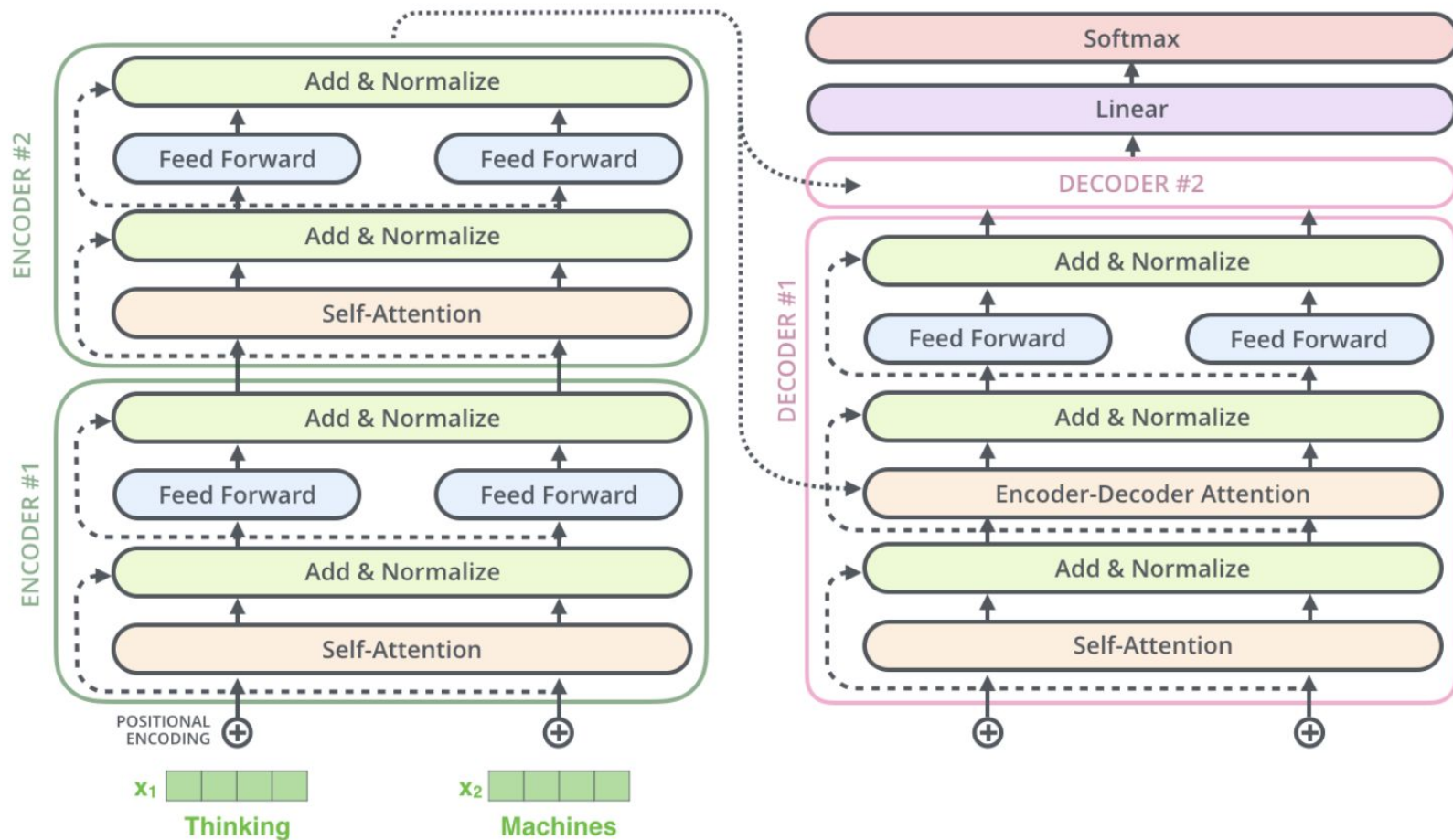
Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

N×

N×

Positional Encoding

Positional Encoding

Input Embedding

Output Embedding

Inputs

Outputs (shifted right)

20

# Layer Normalization

# Layer Normalization

# Layer Normalization

# The Decoder

# The Decoder Side

Decoding time step: (1) 2  3  4  5  6          OUTPUT
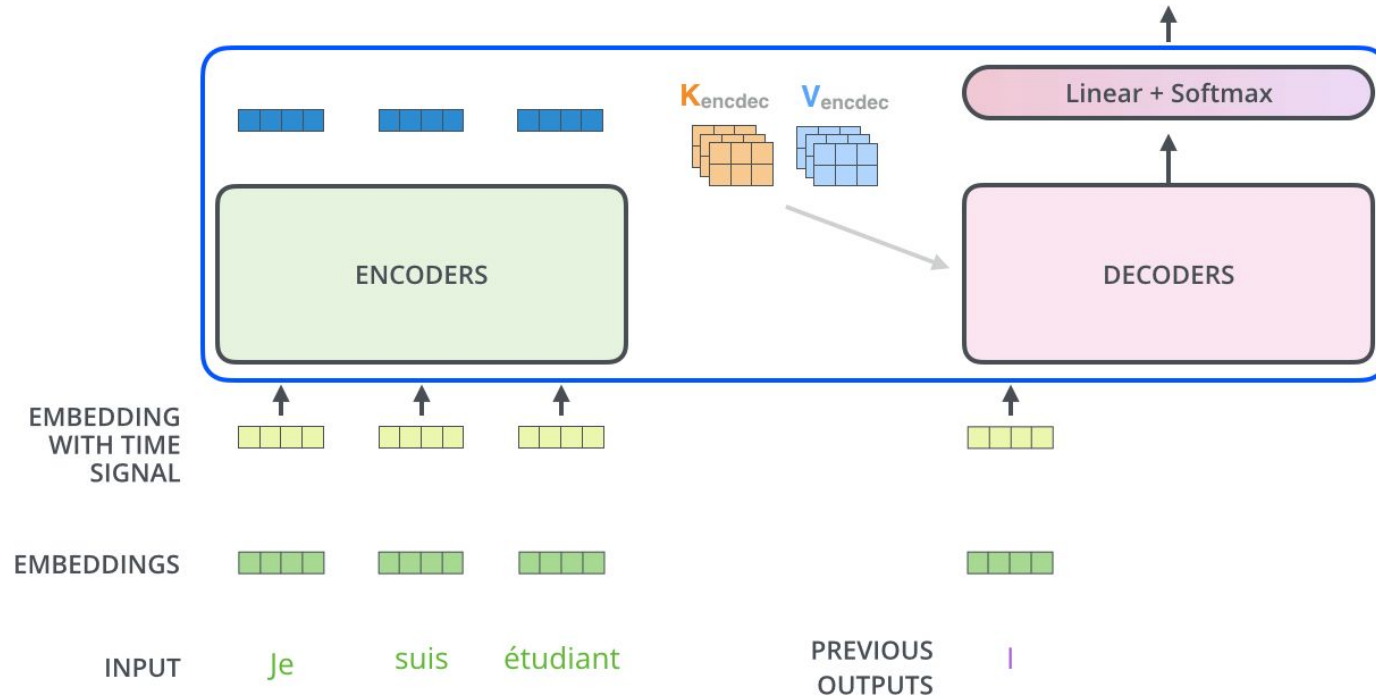


EMBEDDING WITH TIME SIGNAL

EMBEDDINGS

INPUT     Je        suis      étudiant

# The Decoder Side

Decoding time step: 1 ② 3  4  5  6          OUTPUT          I



ENCODERS          $K_{encdec}$  $V_{encdec}$          Linear + Softmax

DECODERS

EMBEDDING WITH TIME SIGNAL

EMBEDDINGS

INPUT          Je          suis          étudiant          PREVIOUS OUTPUTS          I

Which word in our vocabulary is associated with this index?

am

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

5

log_probs

0 1 2 3 4 5                                    ... vocab_size

Softmax

logits

0 1 2 3 4 5                                    ... vocab_size

Linear

Decoder stack output

# Recap of Training

### Output Vocabulary

| WORD | a | am | I | thanks | student | <eos> |
|------|---|----|----|--------|---------|-------|
| INDEX | 0 | 1 | 2 | 3 | 4 | 5 |

One-hot encoding of the word "am"

| | | | | | |
|---|---|---|---|---|---|
| 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 |

Example: one-hot encoding of our output vocabulary

# Recap of Training

**Target Model Outputs**

Output Vocabulary:  a   am   I   thanks   student   <eos>

| | 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 |

a   am   I   thanks   student   <eos>

**Trained Model Outputs**

Output Vocabulary:  a   am   I   thanks   student   <eos>

| | 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>

# Loss Function

Kullback-Leibler Divergence:

$$D_{KL}(P||Q) = \sum_x p(x) log \frac{p(x)}{q(x)}$$

Cross-Entropy:

$$H(p, q) = -\sum_x p(x) \log q(x)$$

# **OpenAI Transformer:** Pre-training Decoder for Language Modeling

- The Encoder-Decoder structure of the transformer made it perfect for machine translation
- But what about sentence classification?
- **Main goal: pre-train a language model that can be fine-tuned for other tasks**

# OpenAI Transformer



Possible classes: All English words

| | |
|---|---|
| 0.1% | Aardvark |
| ... | ... |
| 10% | Improvisation |
| ... | ... |
| 0% | Zyzzyva |

FFNN + Softmax

12 DECODER

...

2 DECODER

1 DECODER

Let's    stick    to

Differences from vanilla Transformer:

- no encoder
- decoder layers would not have the encoder-decoder attention sublayer
- Pre-train the model on predicting the next word using massive (unlabeled) datasets

# OpenAI Transformer



- During pre-training phase layers have been tuned to reasonably handle language
- Now let's use it for downstream tasks (e.g. sentence classification)

# Input transformations for different tasks

# ELMo: context that matters

# ELMo: contextualized word embeddings

*"Why not give it an embedding based on the context it's used in – to both capture the word meaning in that context as well as other contextual information?"*

Peters et. al., 2017, McCann et. al., 2017, and yet again Peters et. al., 2018 in the ELMo paper

# ELMo - deep contextualized word representations

ELMo

**What does it stand for?**

1. **E**xpedited **L**abour **M**arket **O**pinion
2. **E**lectric **L**ight **M**achine **O**rganization
3. **E**nough **L**et's **M**ove **O**n

**What does it stand for?**

1. **E**xpedited **L**abour **M**arket **O**pinion
2. **E**lectric **L**ight **M**achine **O**rganization
3. **E**nough **L**et's **M**ove **O**n

4. **Embeddings from Language Models**

# ELMo: contextualized word embeddings

# ELMo: Contextualized word embeddings

ELMo
Embeddings

stick    improvisation    this

Let's    to    in    skit

ELMo

Words to embed

Let's    stick    to improvisation in    this    skit

- uses a bi-directional LSTM trained on Language Modeling task
- a model can learn without labels

# Bidirectional Language Models (biLMs)



biLMs consist of forward and backward LMs:

- forward:

$$p(t_1, t_2, ..., t_N) = \prod_{k=1}^{N} p(t_k | t_1, t_2, ..., t_{k-1})$$

- Backward:

$$p(t_1, t_2, ..., t_N) = \prod_{k=1}^{N} p(t_k | t_{k+1}, t_{k+2}, ..., t_N)$$

LSTM predicts next word in both directions to build biLMs

# ELMo: main pipeline

Embedding of "stick" in "Let's stick to" - Step #1

Forward Language Model

LSTM Layer #2

LSTM Layer #1

Embedding

Let's    stick    to

Backward Language Model

Let's    stick    to

# ELMo: main pipeline

ELMo represents a word as a linear combination of corresponding hidden layers:

Embedding of "stick" in "Let's stick to" - Step #2

| 1- Concatenate hidden layers | Forward Language Model | Backward Language Model |

2- Multiply each vector by a weight based on the task

$\times \quad s_2$

$\times \quad s_1$

$\times \quad s_0$

Let's        stick        to              Let's        stick        to

3- Sum the (now weighted) vectors

ELMo embedding of "stick" for this task in this context

ELMo can be integrated to almost all neural NLP tasks with simple concatenation to the embedding layer

- Pretrained ELMo models: [http://allennlp.org/elmo](http://allennlp.org/elmo)
- AllenNLP is a library on the top of PyTorch

- Higher levels seems to catch semantics while lower layer probably capture syntactic features

# BERT

Bidirectional Encoder Representations from Transformers

# BERT

Input
**Features**

Output
**Prediction**

Help Prince Mayuko Transfer Huge Inheritance

BERT

**Classifier**
(Feed-forward neural network + softmax)

| 85% | Spam |
| 15% | Not Spam |

# BERT: base and large



BERT_BASE

BERT_LARGE

# BERT vs. Transformer



| | THE TRANSFORMER | Base BERT | Large BERT |
|---|---|---|---|
| Encoders | 6 | 12 | 24 |
| Units in FFN | 512 | 768 | 1024 |
| Attention Heads | 8 | 12 | 16 |

Identical to the Transformer up until this point

Identical to the Transformer up until this point

Why is BERT so special?

# Model outputs

Each position outputs a vector



BERT

For sentence classification we focus on the first position (that we passed [CLS] token to)

85% Spam

15% Not Spam

Classifier
(Feed-forward neural network + softmax)

1    2    3    4    • • •    512

This vector can now
be used as the input
for a classifier

BERT

1    2    3    4    • • •    512

[CLS]  Help  Prince  Mayuko

# Similar to CNN concept!

BERT: pre-training

Use the output of the masked word's position to predict the masked word

Possible classes:
All English words

| 0.1% | Aardvark |
| ... | ... |
| 10% | Improvisation |
| ... | ... |
| 0% | Zyzzyva |

FFNN + Softmax

1 2 3 4 5 6 7 8 ••• 512

BERT

Randomly mask 15% of tokens

1 2 3 4 5 6 7 8 ••• 512

[CLS] Let's stick to [MASK] in this skit

Input

[CLS] Let's stick to improvisation in this skit

- "Masked Language Model" approach

- To make BERT better at handling relationships between multiple sentences, the pre-training process includes an additional task:

  *"Given two sentences (A and B),is B likely to be the sentence that follows A, or not?"*

BERT: pre-training

Predict likelihood that sentence B belongs after sentence A

1%  IsNext

99%  NotNext

FFNN + Softmax

1  2  3  4  5  6  7  8  •••  512

BERT

Tokenized Input

1  2  3  4  5  6  7  8  •••  512

[CLS]  the  man  [MASK]  to  the  store  [SEP]

Input

[CLS] the man [MASK] to the store [SEP] penguin [MASK] are flightless birds [SEP]

Sentence A          Sentence B

# BERT: fine-tuning for different tasks



(a) Sentence Pair Classification Tasks:
MNLI, QQP, QNLI, STS-B, MRPC,
RTE, SWAG

(b) Single Sentence Classification Tasks:
SST-2, CoLA

# BERT: fine-tuning for different tasks



(c) Question Answering Tasks:
SQuAD v1.1

(d) Single Sentence Tagging Tasks:
CoNLL-2003 NER

# BERT for feature extraction

**Generate Contextualized Embeddings**

The output of each encoder layer along each token's path can be used as a feature representing that token.
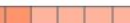
But which one should we use?

# BERT for feature extraction

What is the best contextualized embedding for "Help" in that context?
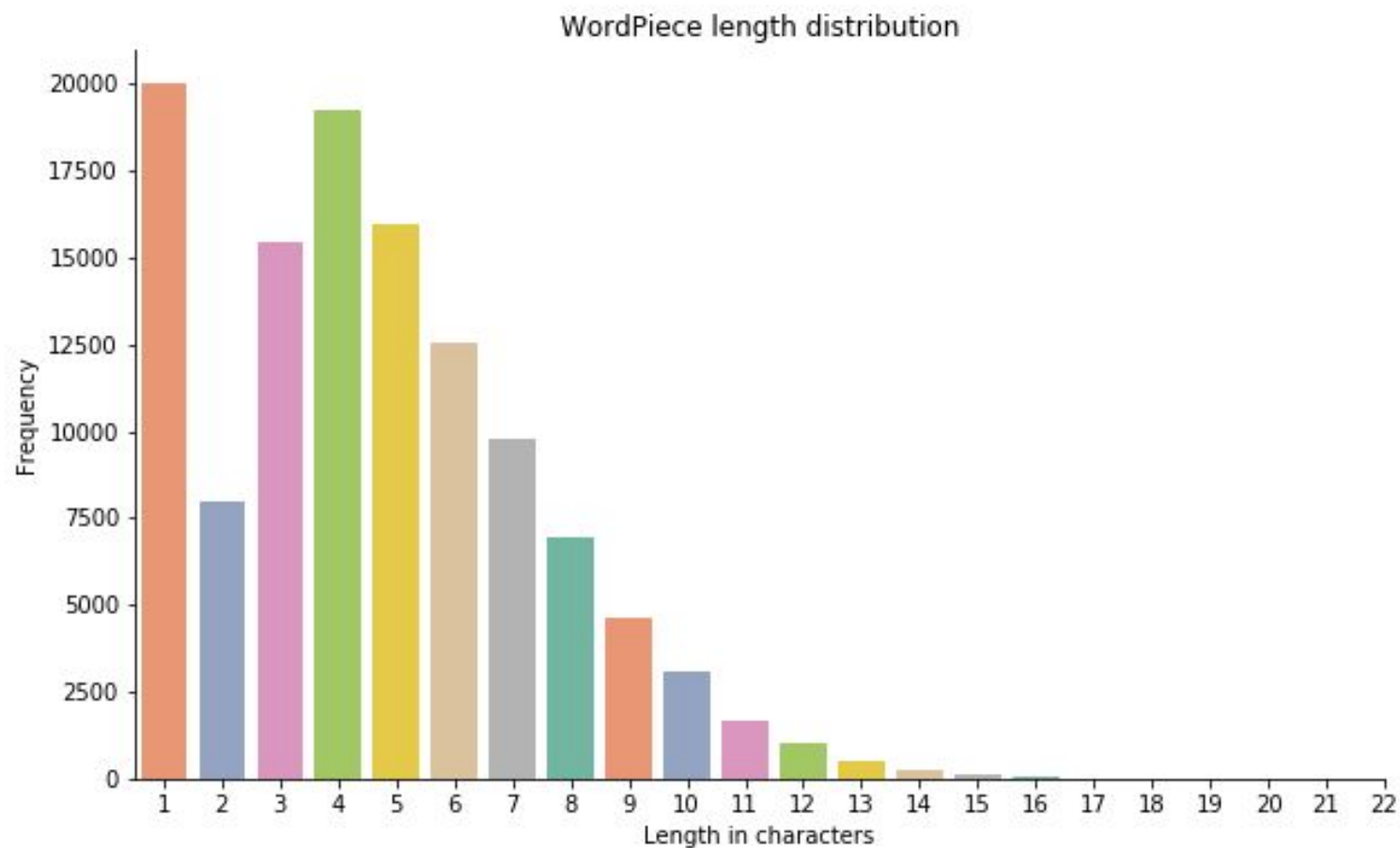For named-entity recognition task CoNLL-2003 NER



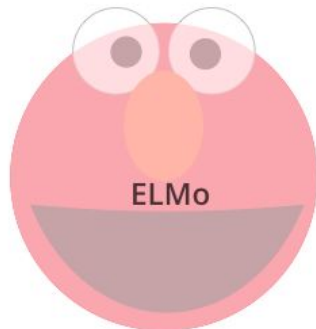| | Dev F1 Score |
|---|---|
| First Layer | 91.0 |
| Last Hidden Layer | 94.9 |
| Sum All 12 Layers | 95.5 |
| Second-to-Last Hidden Layer | 95.6 |
| Sum Last Four Hidden | 95.9 |
| Concat Last Four Hidden | 96.1 |

## Example: Unaffable -> un, ##aff, ##able

- Single model for 104 languages with a large shared vocabulary (119,547 WordPiece model)
- Non-word-initial units are prefixed with ##
- The first 106 symbols: constants like PAD and UNK
- 36.5% of the vocabulary are non-initial word pieces
- The alphabet consists of 9,997 unique characters that are defined as word-initial (C) and continuation symbols (##C), which together make up 19,994 word pieces
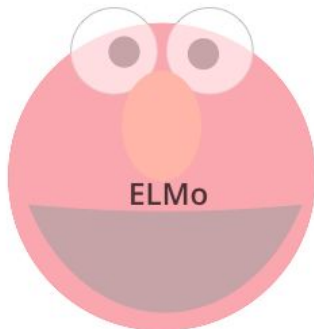- The rest are multicharacter word pieces of various length.

WordPiece length distribution

# BERT: overview

- [BERT repo](#)
- [Try out BERT on TPU](#)
- [WordPieces Tokenizer](#)
- [PyTorch Implementation of BERT](#)

THE TRANSFORMER

ULM-FiT

Idea for the individual project :)

OpenAI Transformer

ELMo

BERT