# Modern NLP Models

Prof. Kuan-Ting Lai

2022/5/16
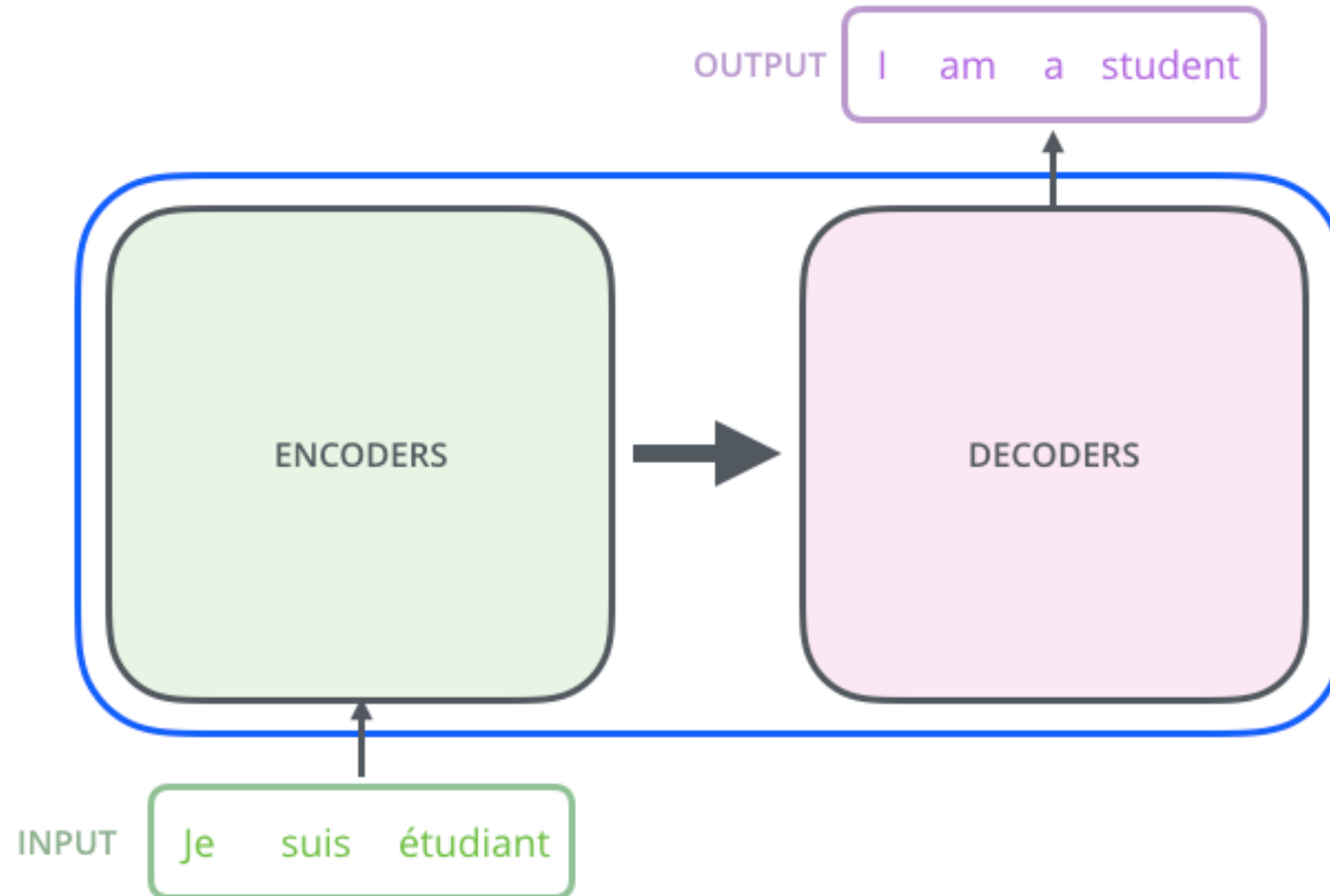
# Transformer Illustration
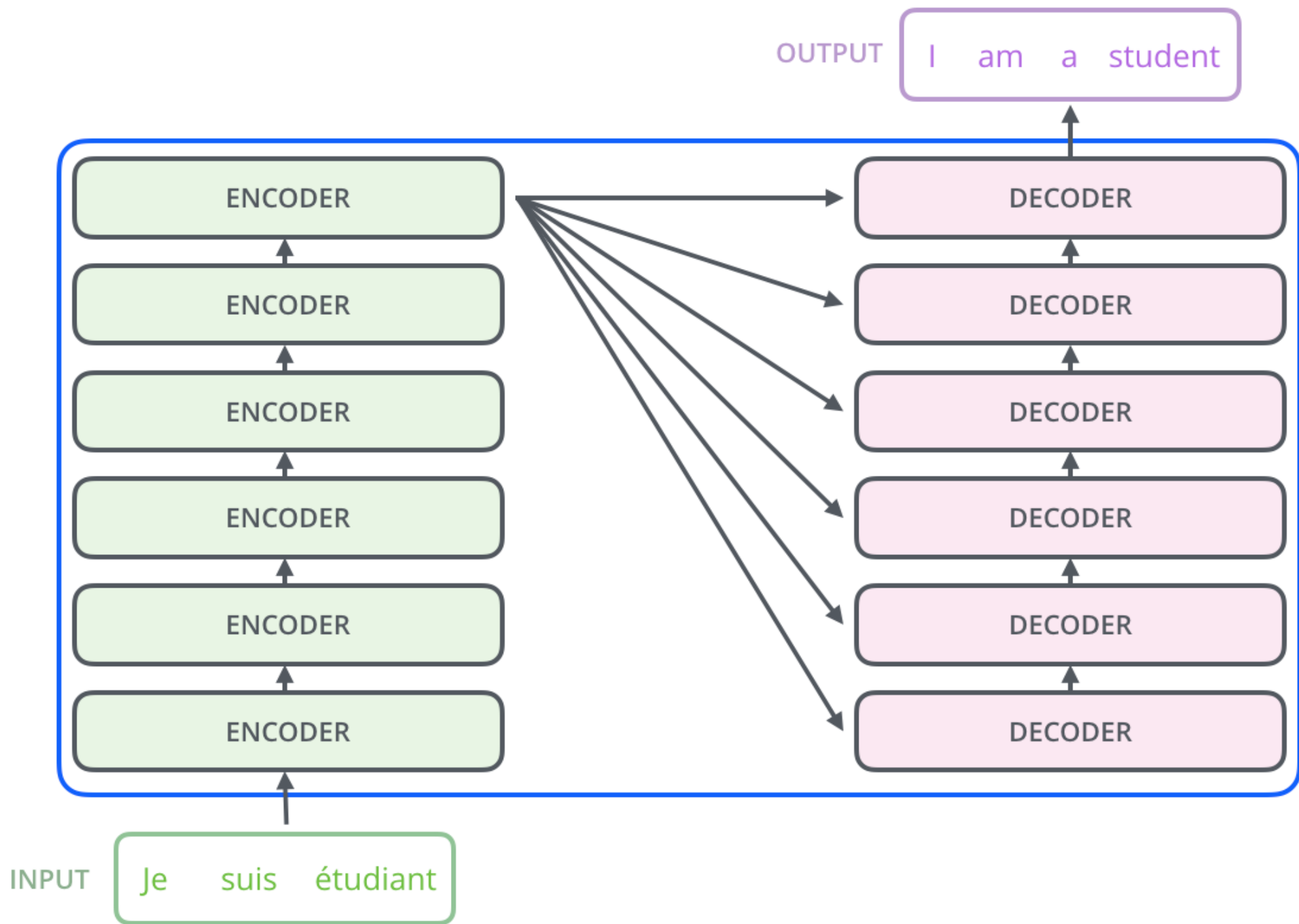
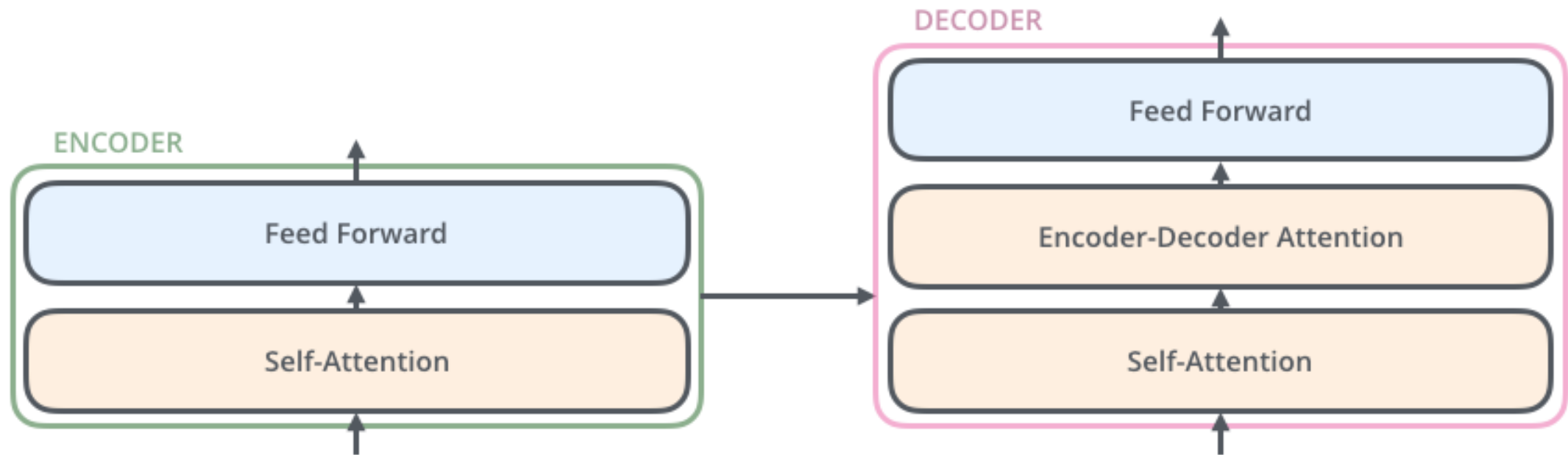- http://jalammar.github.io/illustrated-transformer/
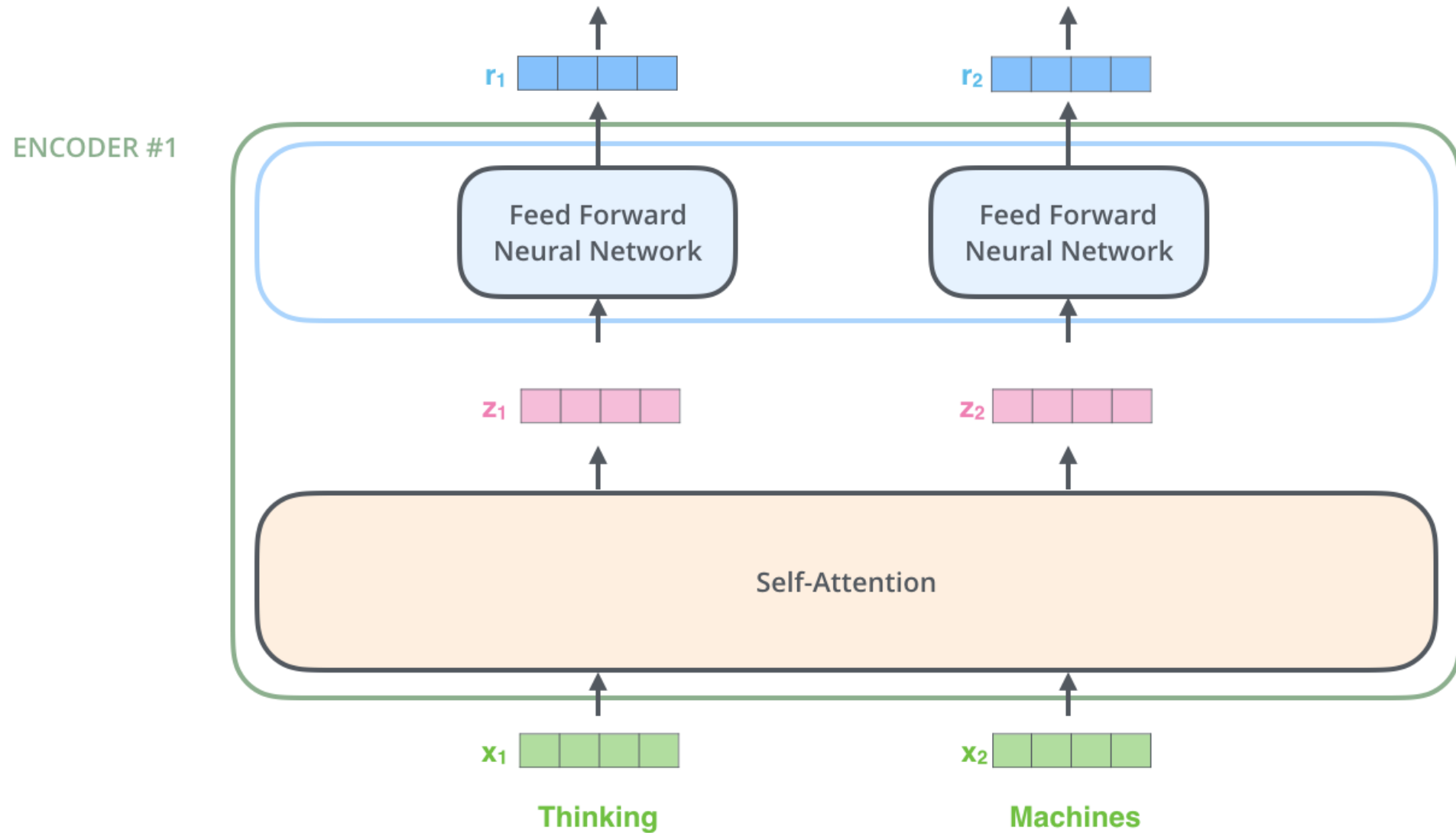
# Encoder and Decoder

# Structure of the Encoder and Decoder

- Self-attention
- Encoder-decoder attention

# Start Encoding

# Calculate Values

- Calculate dot product of key and value vector

- Multiply each value vector by the Softmax score

- Sum up the weighted value vectors $v_1$ and $v_2$

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

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

# Final Output of the Self-attention Module

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

$$= \underset{Z}{\blacksquare}$$

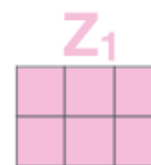# The Beast with Multiple Heads
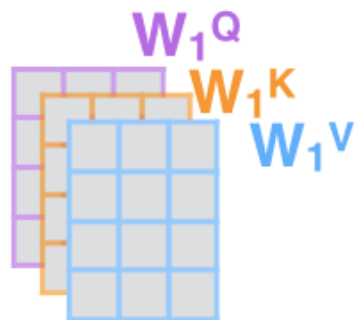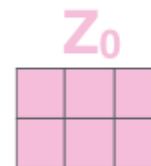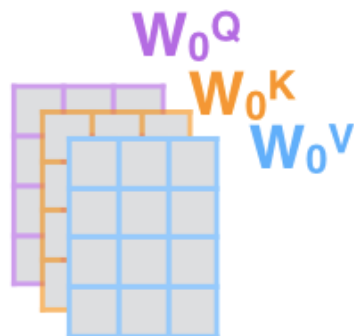
1) This is our input sentence*

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

$X$

$W_0^Q$
$W_0^K$
$W_0^V$

$Q_0$
$K_0$
$V_0$

$Z_0$
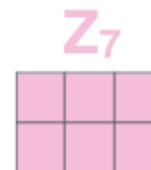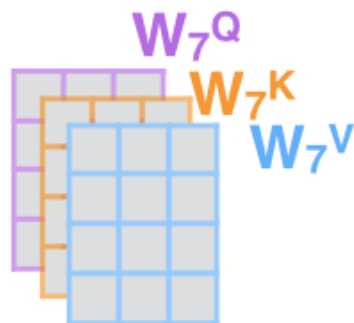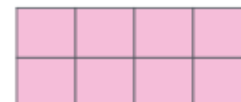
$W^O$

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

$W_1^Q$
$W_1^K$
$W_1^V$

$Q_1$
$K_1$
$V_1$

$Z_1$

$Z$

$R$

...

...

...

$W_7^Q$
$W_7^K$
$W_7^V$

$Q_7$
$K_7$
$V_7$

$Z_7$

# Positional Encoding

- Use sine and cosine functions of different frequencies
  - pos: word position
  - i: dimension index
  - $d_{model}$ = 512

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

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

| POSITIONAL ENCODING | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 1 | 1 | | 0.84 | 0.0001 | 0.54 | 1 | | 0.91 | 0.0002 | -0.42 | 1 |

+      +      +

EMBEDDINGS    $x_1$      $x_2$      $x_3$

INPUT     Je      suis      étudiant

# Encoding Variable-length Sentences

| Sequence | Index of token | Positional Encoding Matrix | | | |
|:---:|:---:|:---:|:---:|:---:|:---:|
| I | 0 | $P_{00}$ | $P_{01}$ | ... | $P_{0d}$ |
| am | 1 | $P_{10}$ | $P_{11}$ | ... | $P_{1d}$ |
| a | 2 | $P_{20}$ | $P_{21}$ | ... | $P_{2d}$ |
| Robot | 3 | $P_{30}$ | $P_{31}$ | ... | $P_{3d}$ |

Positional Encoding Matrix for the sequence 'I am a robot'

# Positional Encoding Layer in Transformers

- k: Position of an object in input sequence, $0 \le k < L/2$

- d: Dimension of the output embedding space

- P(k,j): Position function for mapping a position k in the input sequence to index (k, j) of the positional matrix

- n: User defined scalar. Set to 10,000 by the authors of [Attention is all You Need](#).

- i: Used for mapping to column indices $0 \le i < d/2$. A single value of i maps to both sine and cosine functions

$$P(k, 2i) = \sin\left(\frac{k}{n^{2i/d}}\right)$$
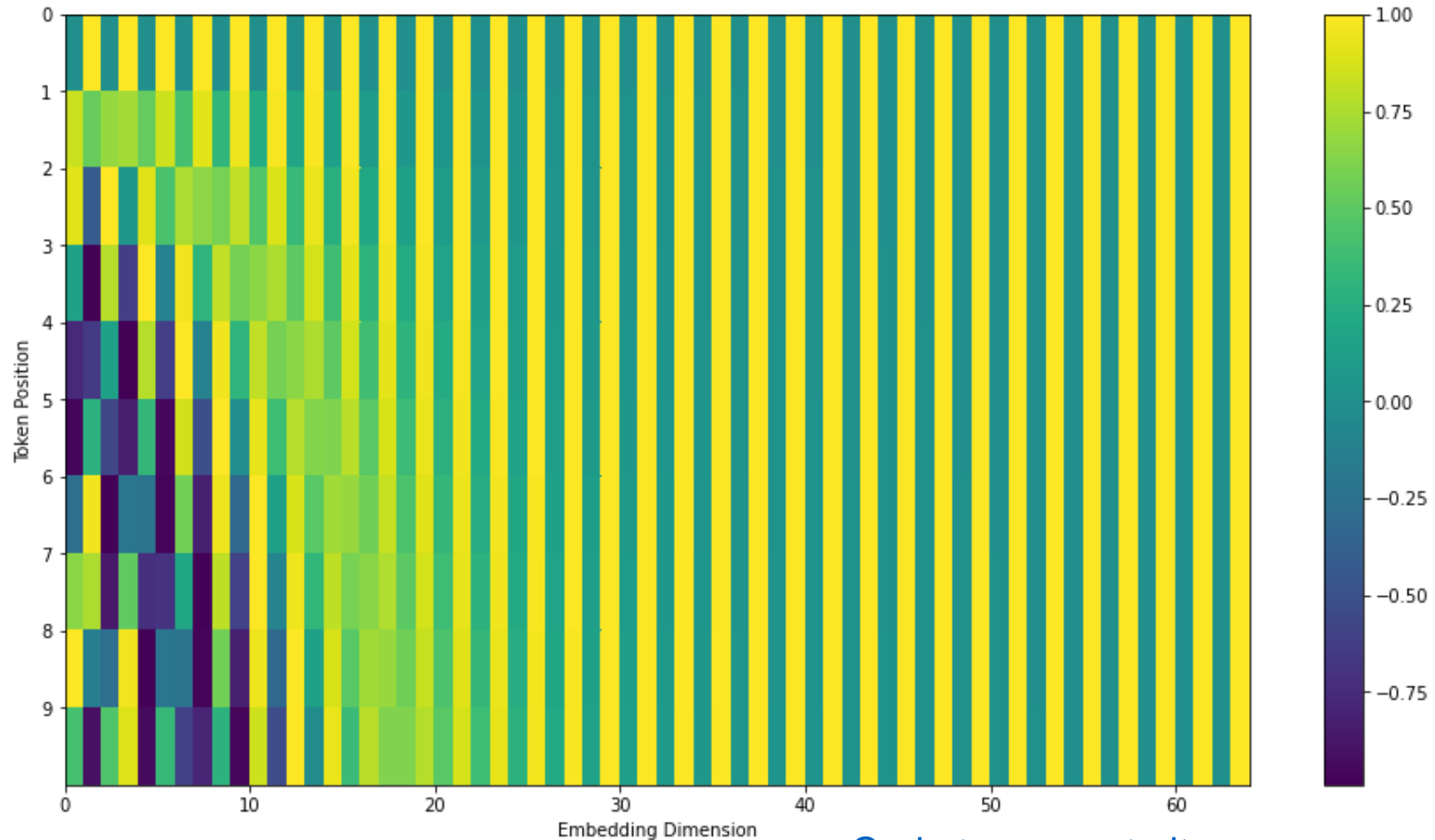
$$P(k, 2i + 1) = \cos\left(\frac{k}{n^{2i/d}}\right)$$

| Sequence | Index of token, k | Positional Encoding Matrix with d=4, n=100 | | | |
|---|---|---|---|---|---|
| | | i=0 | i=0 | i=1 | i=1 |
| I | 0 | $P_{00}=\sin(0)$ $= 0$ | $P_{01}=\cos(0)$ $= 1$ | $P_{02}=\sin(0)$ $= 0$ | $P_{03}=\cos(0)$ $= 1$ |
| am | 1 | $P_{10}=\sin(1/1)$ $= 0.84$ | $P_{11}=\cos(1/1)$ $= 0.54$ | $P_{12}=\sin(1/10)$ $= 0.10$ | $P_{13}=\cos(1/10)$ $= 1.0$ |
| a | 2 | $P_{20}=\sin(2/1)$ $= 0.91$ | $P_{21}=\cos(2/1)$ $= -0.42$ | $P_{22}=\sin(2/10)$ $= 0.20$ | $P_{23}=\cos(2/10)$ $= 0.98$ |
| Robot | 3 | $P_{30}=\sin(3/1)$ $= 0.14$ | $P_{31}=\cos(3/1)$ $= -0.99$ | $P_{32}=\sin(3/10)$ $= 0.30$ | $P_{33}=\cos(3/10)$ $= 0.96$ |

Positional Encoding Matrix for the sequence 'I am a robot'

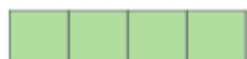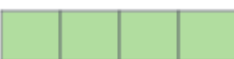# Visualizing Positional Encoding



Code to generate it

# Transformer Decoder

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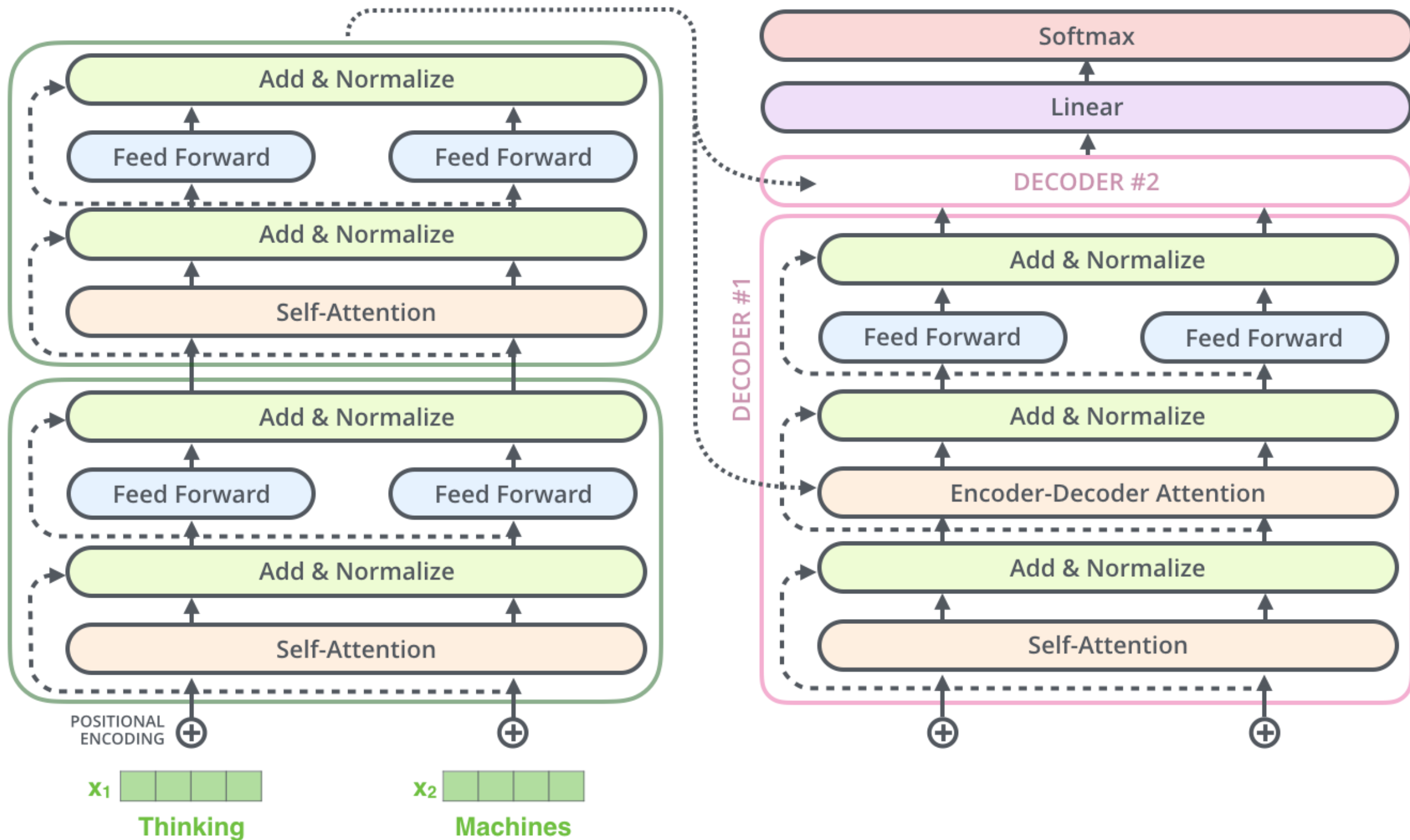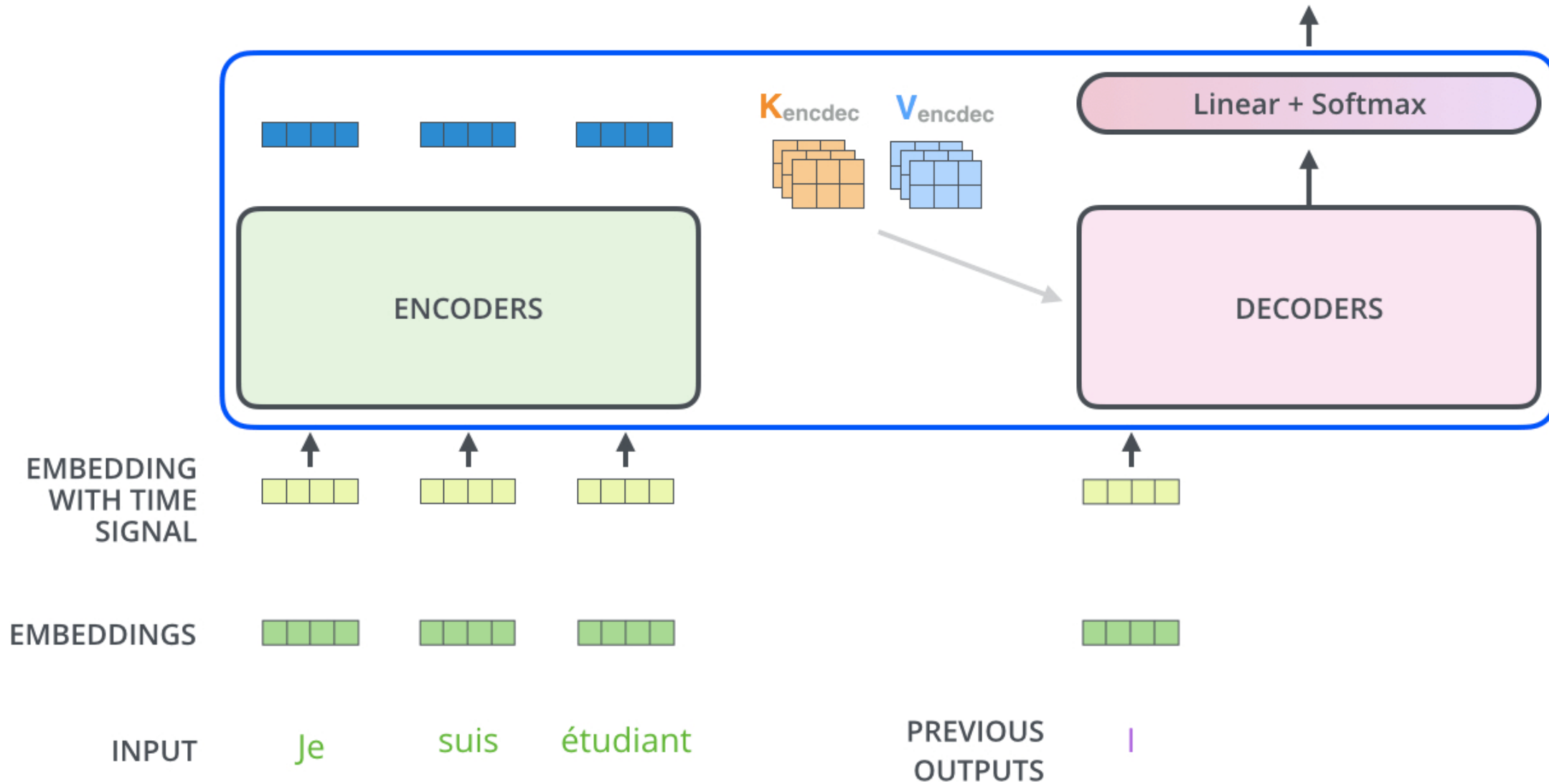
Decoding time step: 1 (2) 3 4 5 6    OUTPUT    I

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

# Output of Decoder

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 |
|-----|-----|-----|-----|-----|-----|

# Trained Model Outputs

# Latest NLP Models (2018 - )

Generative Pre-trained Transformer (GPT)

Embeddings from Language Models (ELMo)

Bidirectional Encoder Representations from Transformers (BERT)

# BERT: Bidirectional Encoder Representations from Transformers (2019)

- Use "Masked Language Model" to train the bidirectional transformer encoder
  - Randomly masked out some tokens and train models to predict them

- Fine-tuning on different tasks

- Achieved state-of-the-art results on multiple NLP tasks

**BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding**

**Jacob Devlin**    **Ming-Wei Chang**    **Kenton Lee**    **Kristina Toutanova**
Google AI Language
{jacobdevlin,mingweichang,kentonl,kristout}@google.com

https://arxiv.org/abs/1810.04805

# 1 - Semi-supervised training on large amounts of text (books, wikipedia..etc).

The model is trained on a certain task that enables it to grasp patterns in language. By the end of the training process, BERT has language-processing abilities capable of empowering many models we later need to build and train in a supervised way.

**Semi-supervised Learning Step**
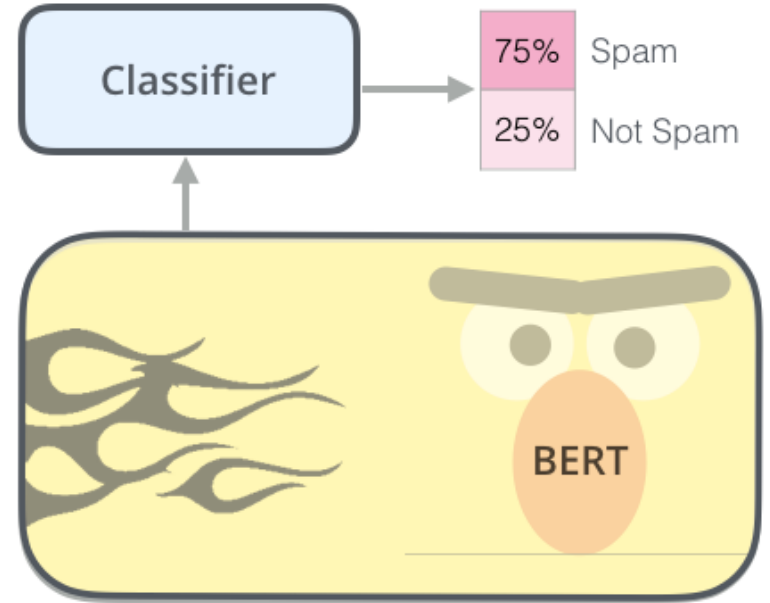
**Model:** BERT

**Dataset:** WIKIPEDIA *Die freie Enzyklopädie*

**Objective:** Predict the masked word (langauge modeling)

# 2 - Supervised training on a specific task with a labeled dataset.

**Supervised Learning Step**

Classifier → 75% Spam / 25% Not Spam

**Model:** (pre-trained in step #1) BERT

**Dataset:**

| Email message | Class |
|---|---|
| Buy these pills | Spam |
| Win cash prizes | Spam |
| Dear Mr. Atreides, please find attached… | Not Spam |

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

Spam Email Classification

85% Spam
15% Not Spam

Classifier
(Feed-forward neural network + softmax)

1 2 3 4 • • • 512

BERT

1 2 3 4 • • • 512

[CLS] Help Prince Mayuko
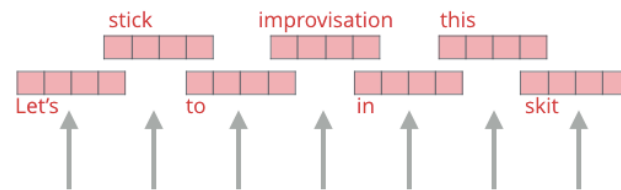
# Embeddings from Language Models (ELMo)

- Consider how words vary across contexts
- Use sentence as input and encoded it by bi-directional LSTM

ELMo Embeddings

stick    improvisation    this

Let's    to    in    skit

ELMo

Words to embed

Let's    stick    to improvisation in    this    skit

Deep contextualized word representations

Matthew E. Peters[†], Mark Neumann[†], Mohit Iyyer[†], Matt Gardner[†],
{matthewp,markn,mohiti,mattg}@allenai.org

Christopher Clark[*], Kenton Lee[*], Luke Zettlemoyer[†*]
{csquared,kentonl,lsz}@cs.washington.edu

[†]Allen Institute for Artificial Intelligence
[*]Paul G. Allen School of Computer Science & Engineering, University of Washington

# Use Bi-LSTM to create Word Embedding



1- Concatenate hidden layers

Forward Language Model

Backward Language Model

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

X   $s_2$

X   $s_1$

X   $s_0$

Let's   stick   to   Let's   stick   to

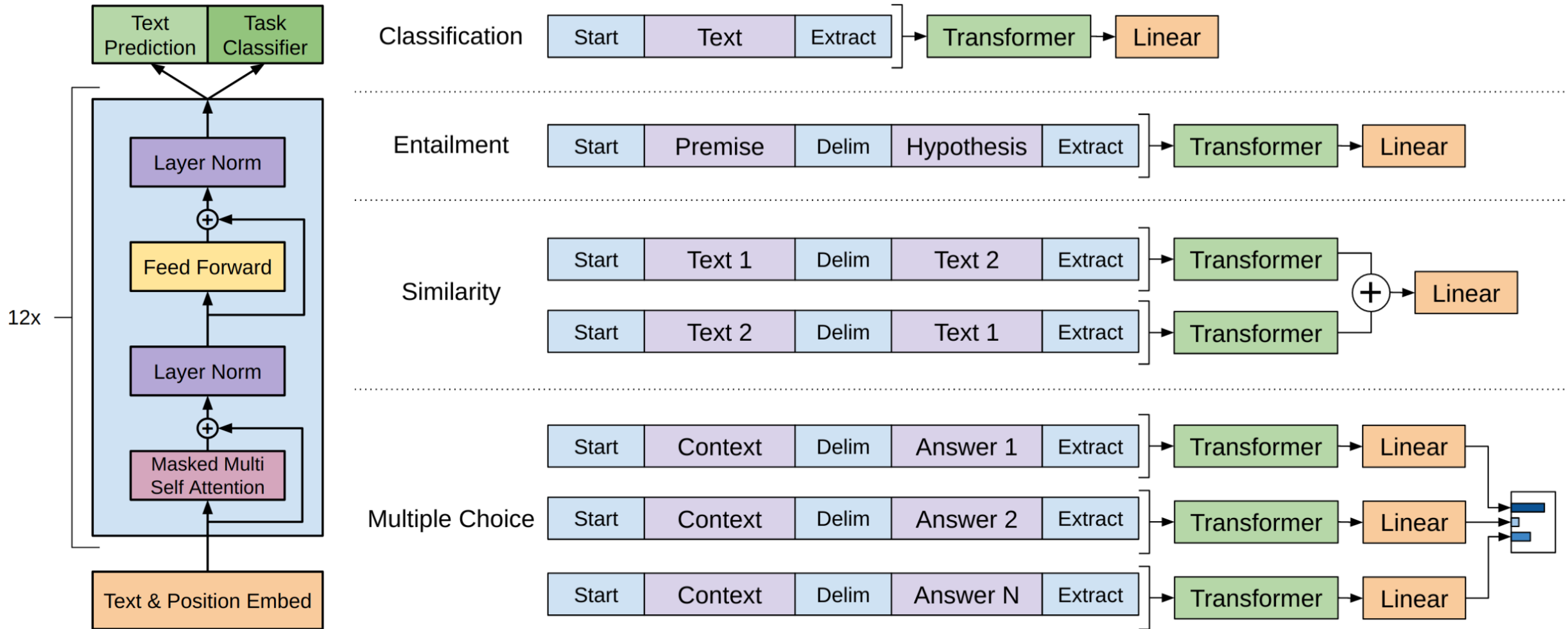3- Sum the (now weighted) vectors

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

27

# OpenAI GPT: Pre-training Transformer Decoders

- Unsupervised pre-train transform decoders for predicting the next word (GPT: Generative Pre-Training)

- Use 12 Transformer decoders in GPT-1
  - GPT-1: Improving Language Understanding with Unsupervised Learning (2018)
  - GPT-2: Better Language Models and Their Implications (2019)
  - GPT-3: Language Models are Few-Shot Learners (2020)

# OpenAI GPT for Different Tasks

# OpenAI GPT-2

- Pre-trained using 40GB of Internet text

- Scale-up of GPT with 10X parameters trained with 10X data

- Other tricks
  - Layer normalization was moved to the input of each sub-block
  - An additional layer normalization was added after the final self-attention block

| Parameters | Layers | $d_{model}$ |
|------------|--------|-------------|
| 117M       | 12     | 768         |
| 345M       | 24     | 1024        |
| 762M       | 36     | 1280        |
| 1542M      | 48     | 1600        |

https://openai.com/blog/better-language-models/

# Size does Matter! GPT-3

- 175 Billion Parameters!

- 175×4=700GB

- 55 years and $4,600,000 to train - even with the lowest priced GPU cloud on the market.



COMPARISON: NLP PRE-TRAINED MODELS

PARAMETER SIZE (MILION)

| MODEL | Parameter Size |
|-------|----------------|
| ALBERT-BASE | 12 |
| ALBERT-LARGE | 18 |
| DistilBERT | 66 |
| BERT-BASE | 110 |
| XLNeT-BASE | 110 |
| RoBERTa-BASE | 125 |
| BERT-LARGE | 340 |
| XLNET-LARGE | 340 |
| RoBERTa-LARGE | 355 |
| GPT-2 | 1,500 |
| T5 | 11,000 |
| Turing-NLG | 17,000 |
| GPT-3 | 175,000 |

# GPT3 Demo (gpt3demo.com )

← → ⟳   🔒 transformer.huggingface.co

# huggingface.co

# Write With **Transformer**

## Get a modern neural network to auto-complete your thoughts.

This web app, built by the Hugging Face team, is the official demo of the 🤗 /transformers repository's text generation capabilities.

Ⓞ Star     51,543

## Checkpoints

### 🐴 DistilGPT-2

The student of the now ubiquitous GPT-2 does not come short of its teacher's expectations. Obtained by distillation, DistilGPT-2 weighs 37% less, and is twice as fast as its OpenAI counterpart, while keeping the same generative power. Runs smoothly on an iPhone 7. The dawn of lightweight generative transformers?

**Start writing**     **More info**

### 🤗 Arxiv-NLP

# huggingface.co

# References

1. https://lilianweng.github.io/lil-log/2018/06/24/attention-attention.html

2. http://jalammar.github.io/illustrated-transformer/

3. http://jalammar.github.io/illustrated-bert/

4. Hong-Yi Lee, Transformer, 2019, https://www.youtube.com/watch?v=ugWDIIOHtPA