# Transformer Models

김윤식
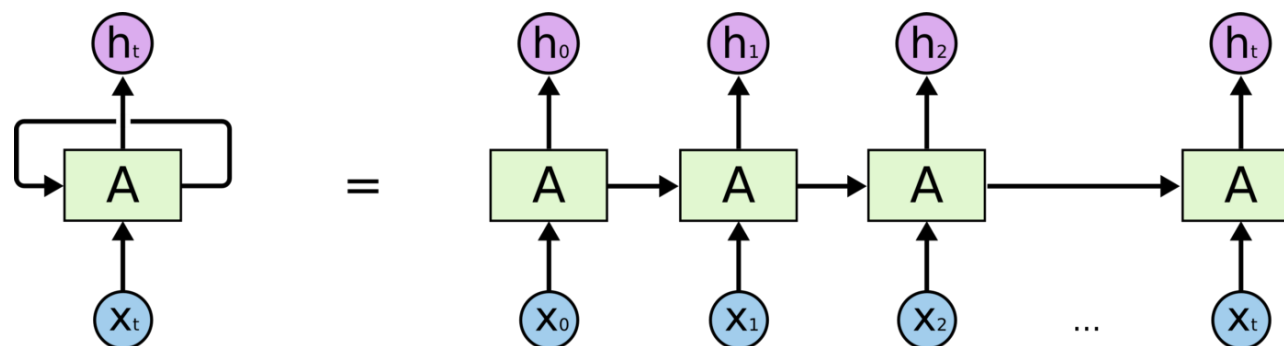
# Introduction

- Machine translation, Machine text generation, etc.: Seq2Seq

- Previously used RNN, LSTM, ⋯



- New approach: "Attention"

# Overview
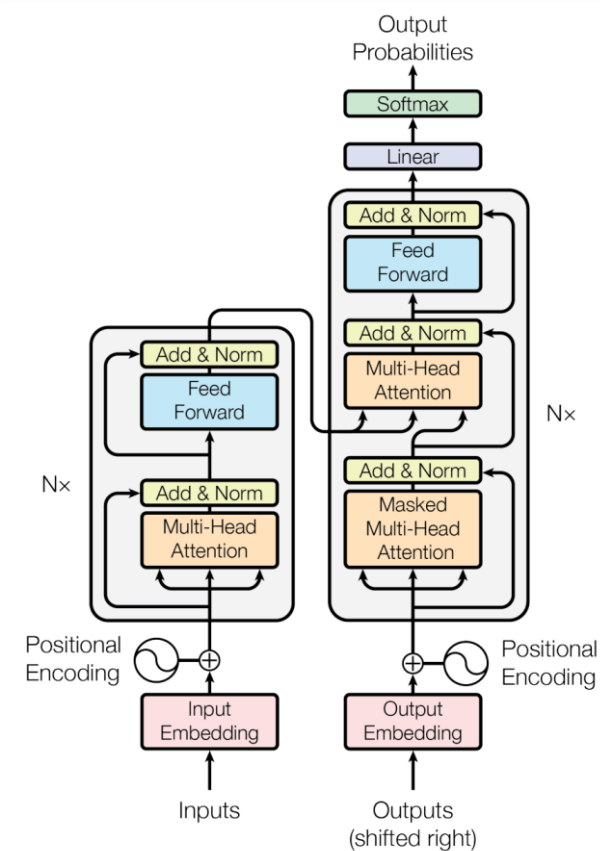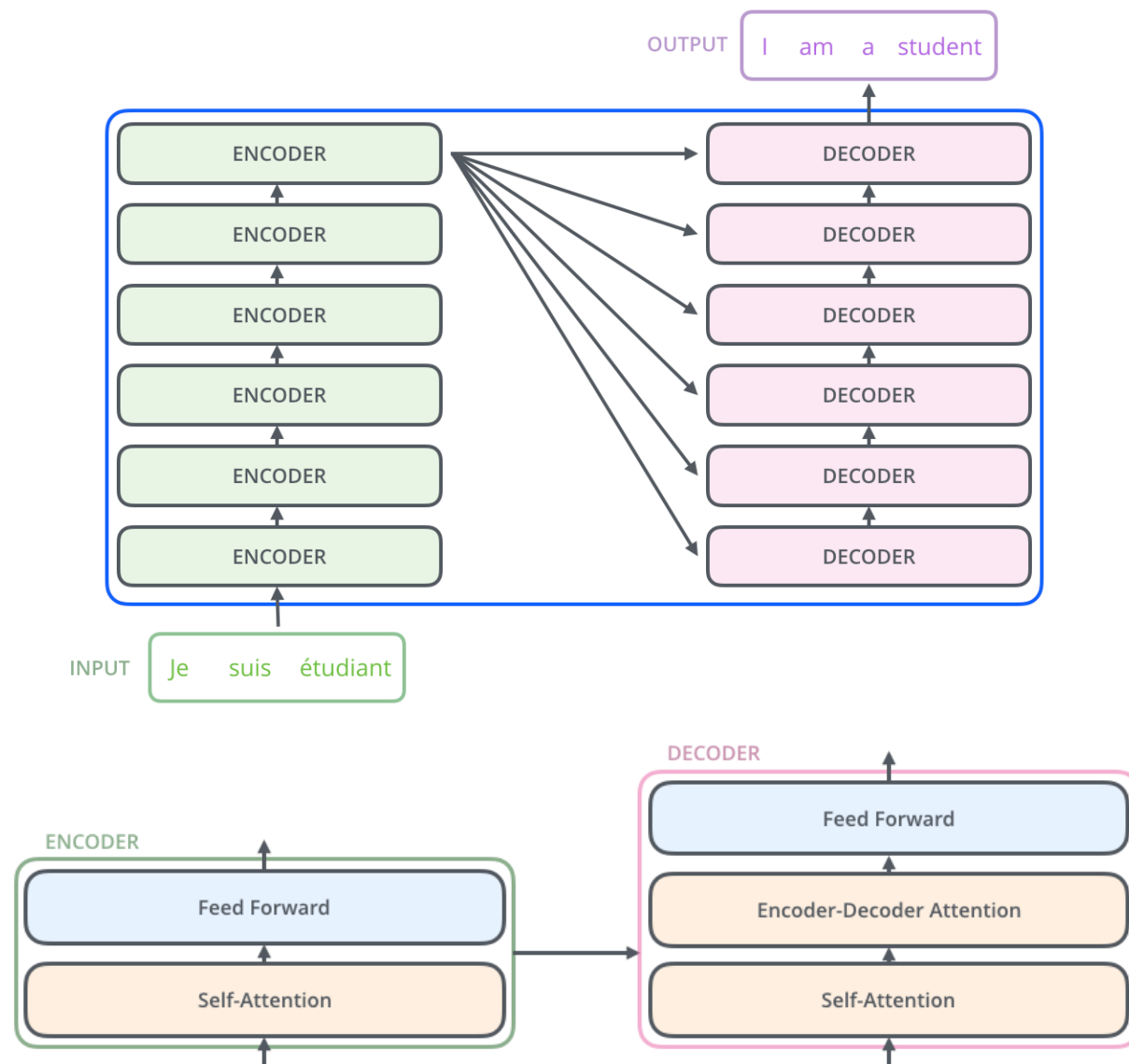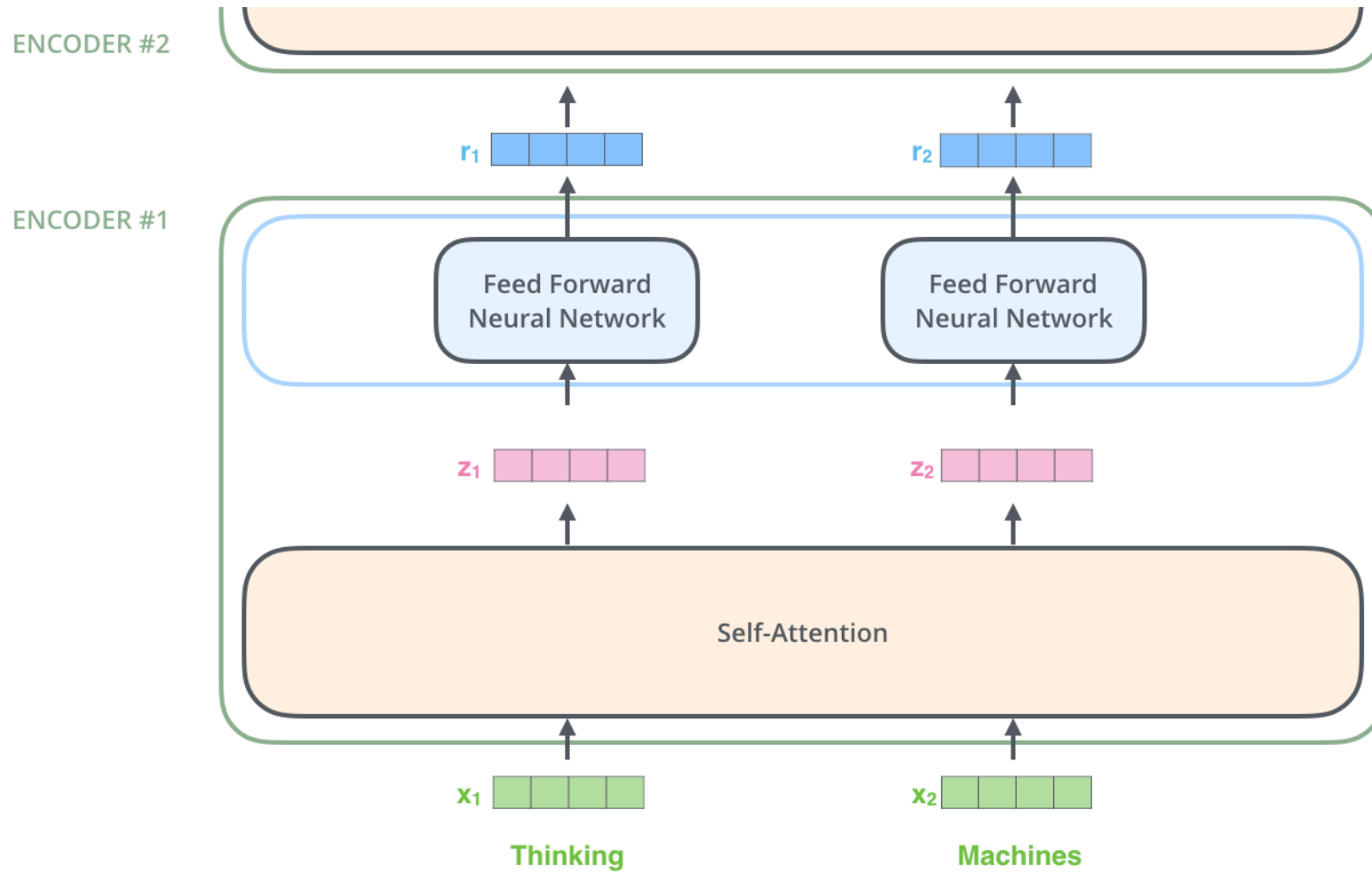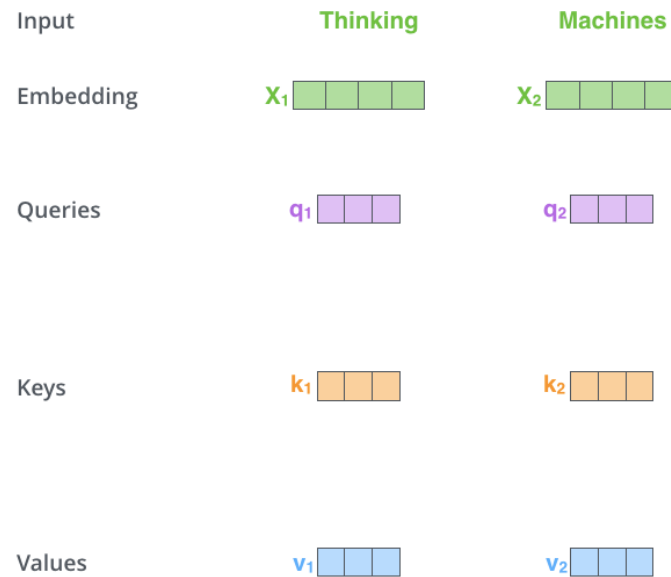
OUTPUT | I am a student

ENCODER
ENCODER
ENCODER
ENCODER
ENCODER
ENCODER

DECODER
DECODER
DECODER
DECODER
DECODER
DECODER

INPUT | Je suis étudiant

**ENCODER**

Feed Forward

Self-Attention

**DECODER**

Feed Forward

Encoder-Decoder Attention

Self-Attention

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

Nx

Nx

Positional Encoding

Positional Encoding

Input Embedding

Output Embedding

Inputs

Outputs (shifted right)

Figure 1: The Transformer - model architecture.

# Encoder

ENCODER #2

ENCODER #1

$r_1$

$r_2$

Feed Forward
Neural Network

Feed Forward
Neural Network

$z_1$

$z_2$

Self-Attention

$x_1$

$x_2$

**Thinking**

**Machines**

# Self-attention

# Multi-headed Self-attention

1) This is our input sentence*

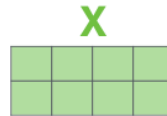2) We embed each word*

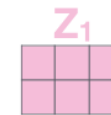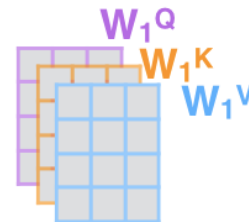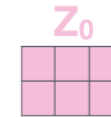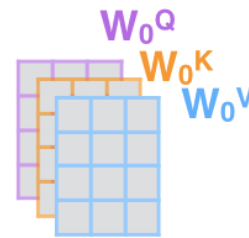3) Split into 8 heads. We multiply $X$ or $R$ with weight matrices

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
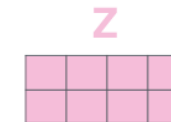
Thinking Machines

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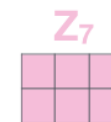
* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one

$R$

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

$Q_0$
$K_0$
$V_0$

$Z_0$

$W^O$

$Z$

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

$Q_1$
$K_1$
$V_1$

$Z_1$

...

...

...

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

$Q_7$
$K_7$
$V_7$

$Z_7$

# Positional Encoding

- Injecting information about the positions of tokens

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

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

# Overall structure of an Encoder
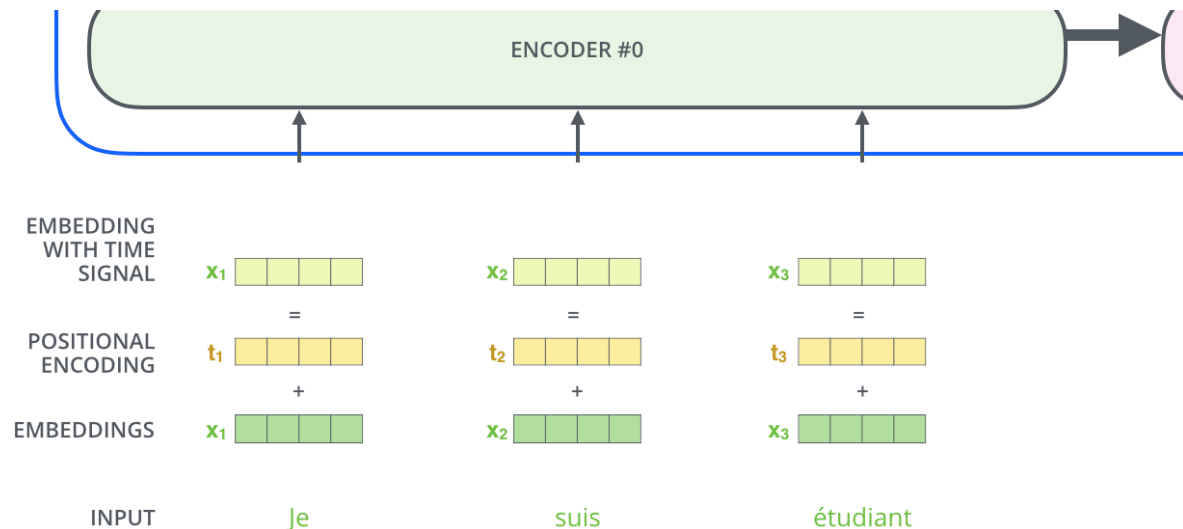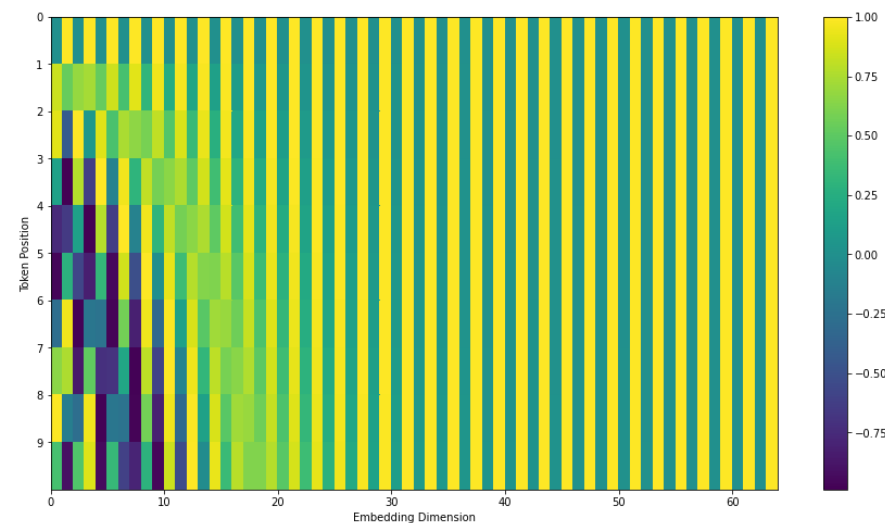
- Residual connections

- Layer normalization

# Decoder



**Encoder-Decoder Attention:**
Same as Self-attention, but K and V vectors are from the encoder outputs

**"Masked" Attention:**
Since decoders should only consider earlier positions in the sequence, future positions are "masked" by setting them to $-\infty$ before the softmax step

9

# Linear and Softmax layers

Which word in our vocabulary is associated with this index?

am

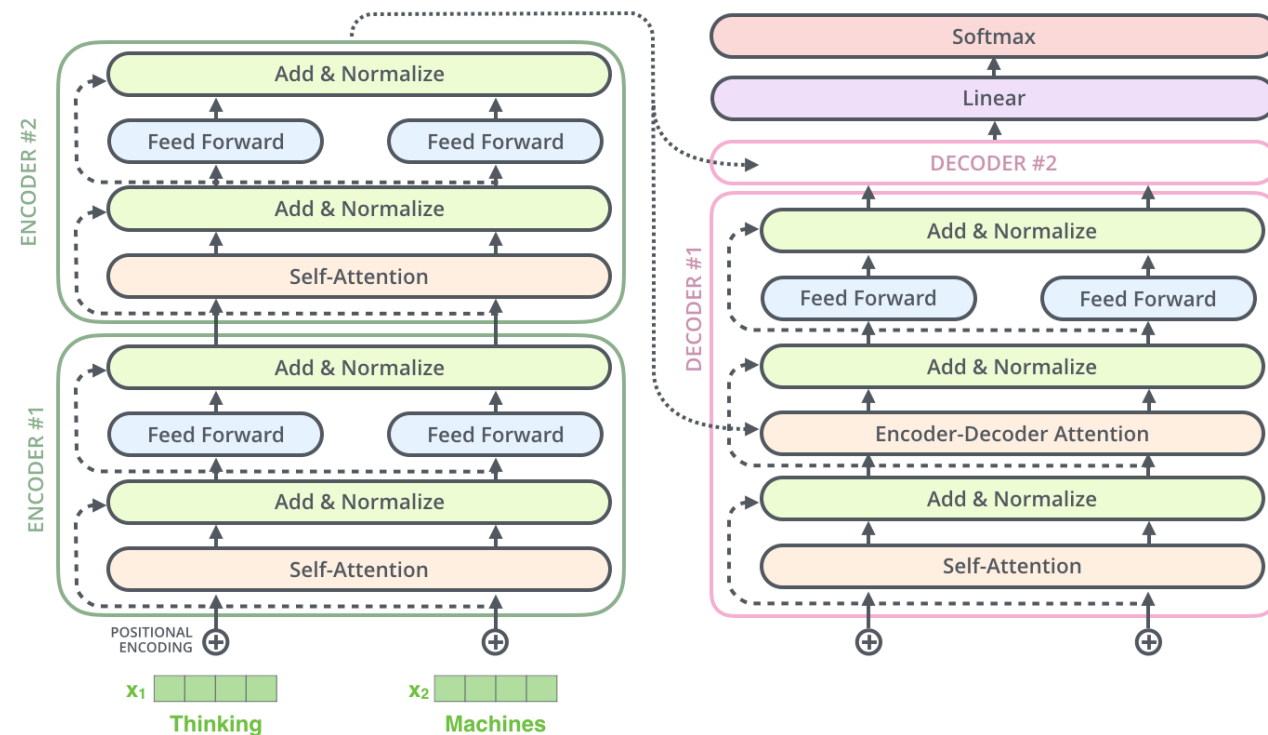Get the index of the cell with the highest value (`argmax`)

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

# Training the model

**Target Model Outputs**

Output Vocabulary:  a   am   I   thanks   student   <eos>

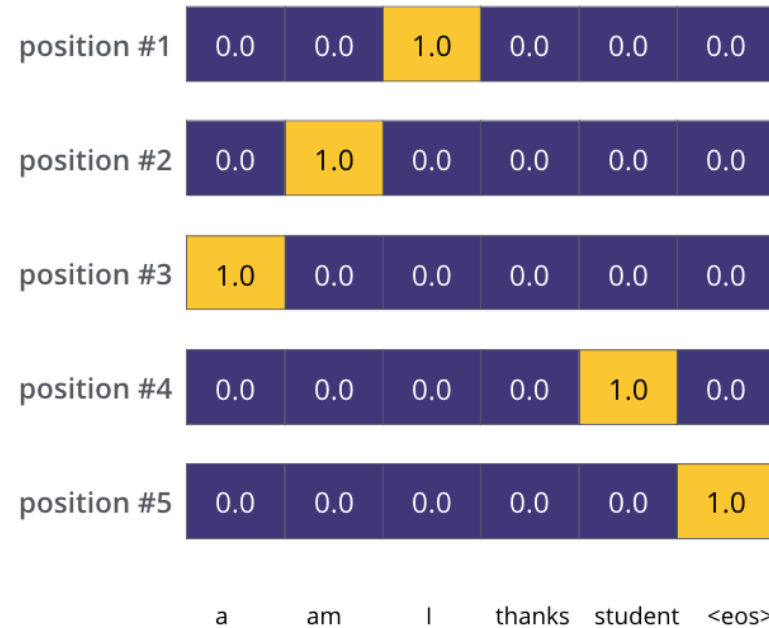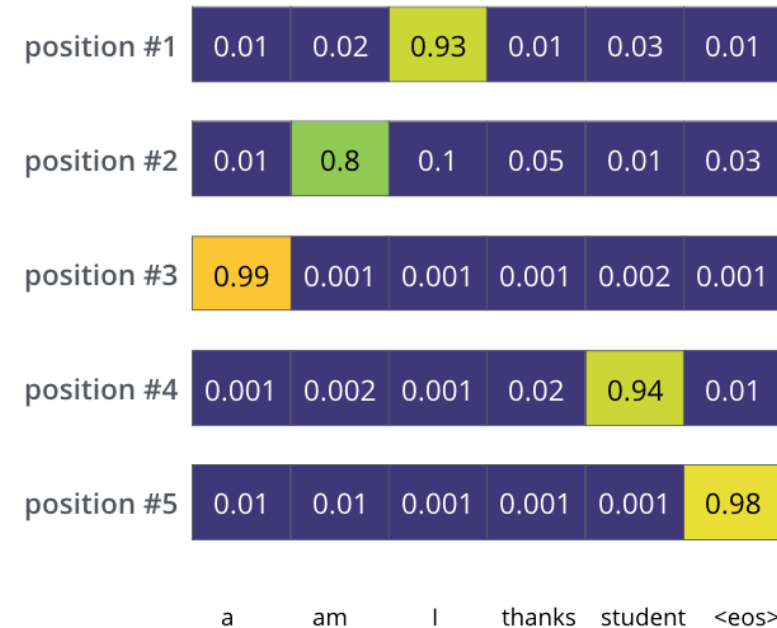| | a | am | I | thanks | student | <eos> |
|---|---|---|---|---|---|---|
| position #1 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 |
| position #2 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| position #3 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| position #4 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 |
| position #5 | 0.0 | 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>

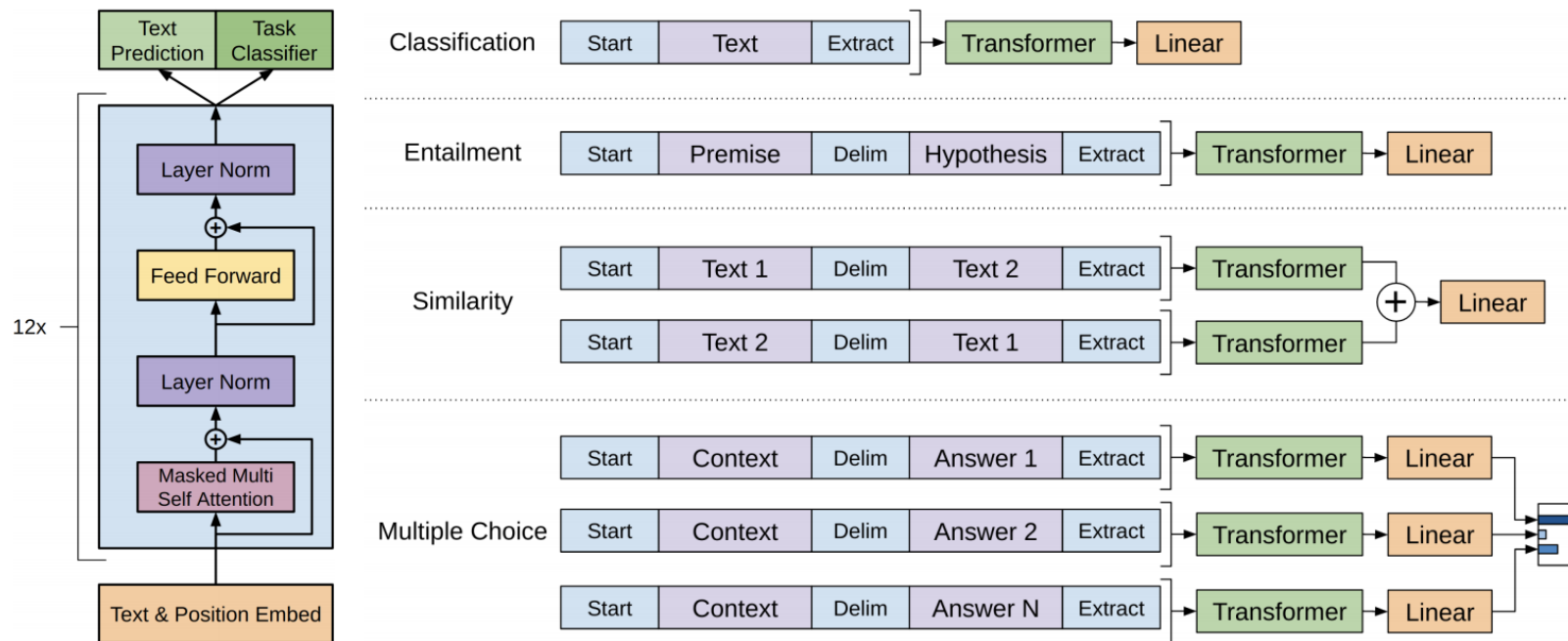| | a | am | I | thanks | student | <eos> |
|---|---|---|---|---|---|---|
| position #1 | 0.01 | 0.02 | 0.93 | 0.01 | 0.03 | 0.01 |
| position #2 | 0.01 | 0.8 | 0.1 | 0.05 | 0.01 | 0.03 |
| position #3 | 0.99 | 0.001 | 0.001 | 0.001 | 0.002 | 0.001 |
| position #4 | 0.001 | 0.002 | 0.001 | 0.02 | 0.94 | 0.01 |
| position #5 | 0.01 | 0.01 | 0.001 | 0.001 | 0.001 | 0.98 |

a   am   I   thanks   student   <eos>

ex. Cross-entropy loss

# ※ GPT (Generative Pre-Training)

- General model for many NLP tasks

- Only uses Decoder part of the Transformer architecture

# References

- Attention Is All You Need (paper): https://arxiv.org/abs/1706.03762

- The Illustrated Transformer (blog post): https://jalammar.github.io/illustrated-transformer/

- Improving Language Understanding by Generative Pre-Training (paper): https://s3-us-west-2.amazonaws.com/openai-assets/research-covers/language-unsupervised/language_understanding_paper.pdf