# Natural Language Processing: Transformers

HSE Faculty of Computer Science

Machine Learning and Data-Intensive Systems
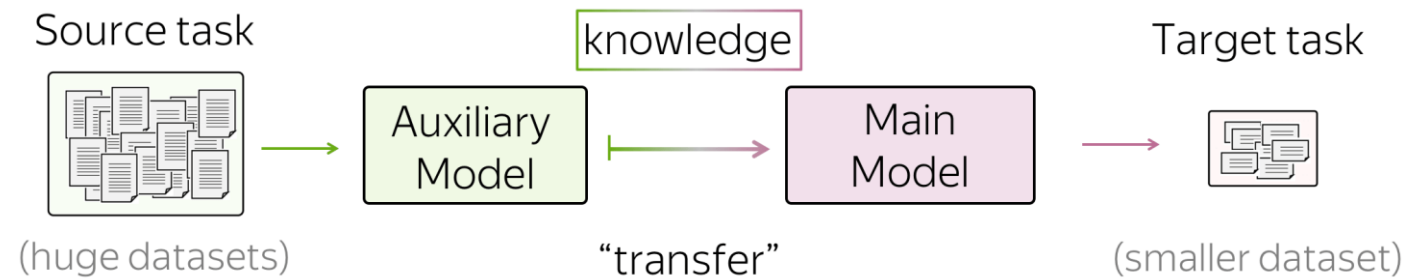
Murat Khazhgeriev
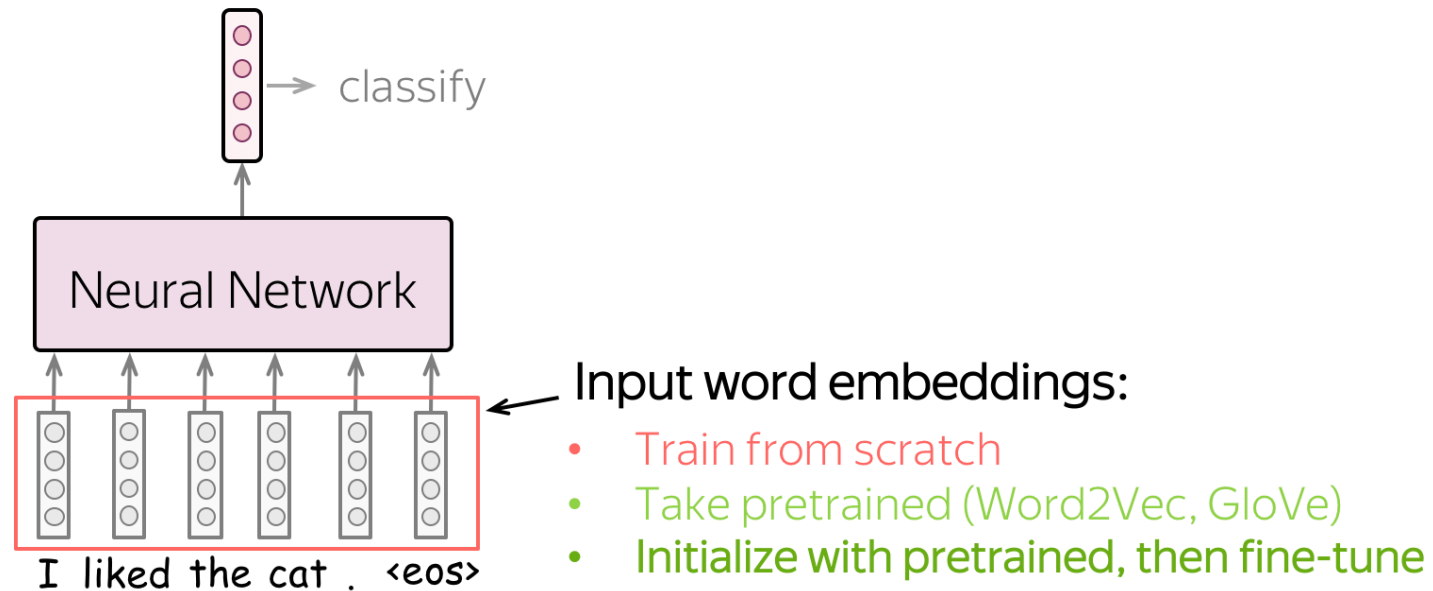
# Table of Content

- **The power of transfer learning**
- From word-specific to contextual embeddings
- Transformer architecture overview
- BERT
- GPT

# Training a model on a simple task can benefit a downstream one



Source: https://lena-voita.github.io/nlp_course/transfer_learning.html

# Training a model on a downstream task can be useful for another



Input word embeddings:
- Train from scratch
- Take pretrained (Word2Vec, GloVe)
- Initialize with pretrained, then fine-tune

Source: https://lena-voita.github.io/nlp_course/transfer_learning.html

# Training a model on a downstream task can be useful for another



- Train from scratch

What they will know:

May be not enough to learn relationships between words

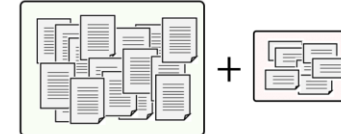- Take pretrained (Word2Vec, GloVe)

What they will know:

Know relationships between words, but are **not** specific to the task

- Initialize with pretrained, then fine-tune

What they will know:

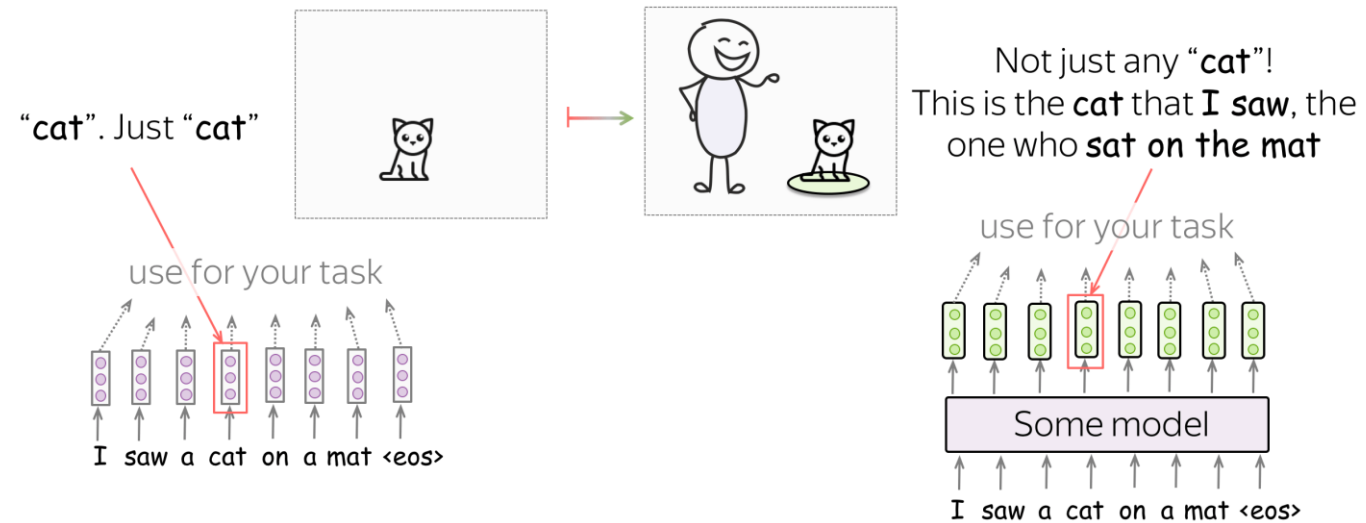Know relationships between words and adapted for the task

"**Transfer**" knowledge from a huge unlabeled corpus to your task-specific model

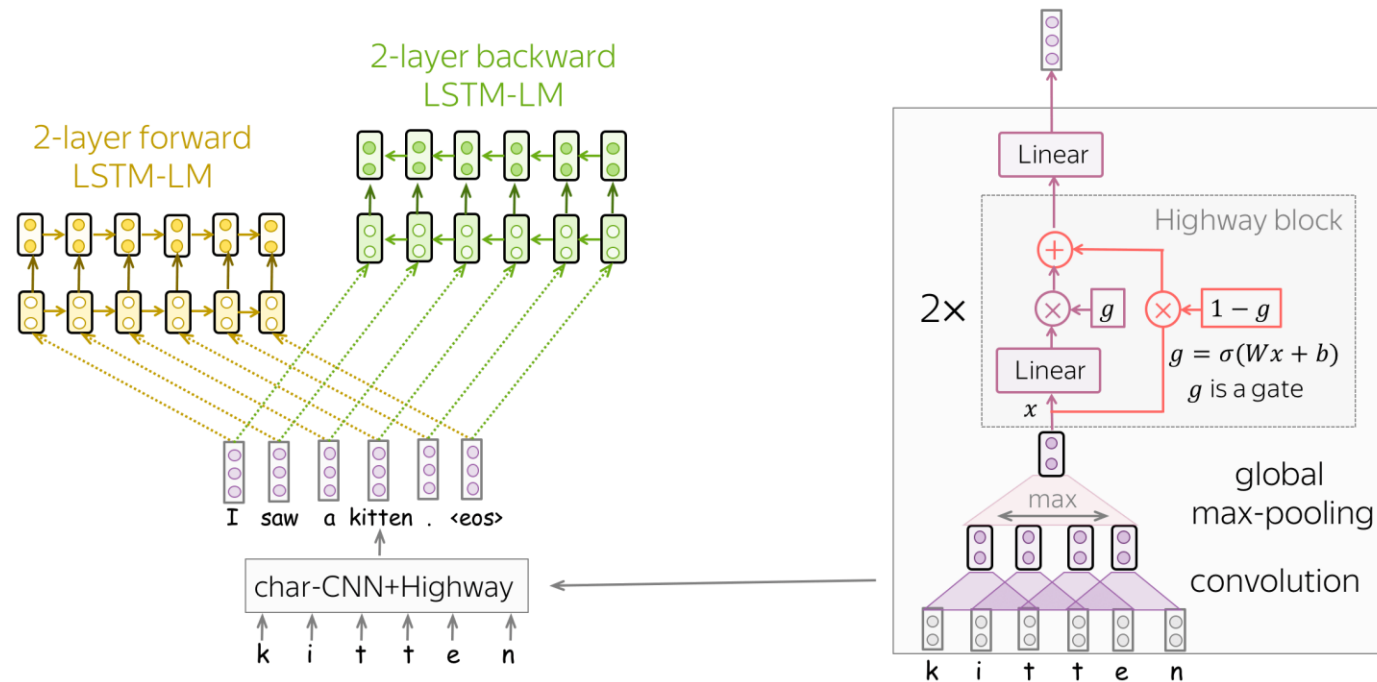Source: https://lena-voita.github.io/nlp_course/transfer_learning.html

# Table of Content

- The power of transfer learning
- **From word-specific to contextual embeddings**
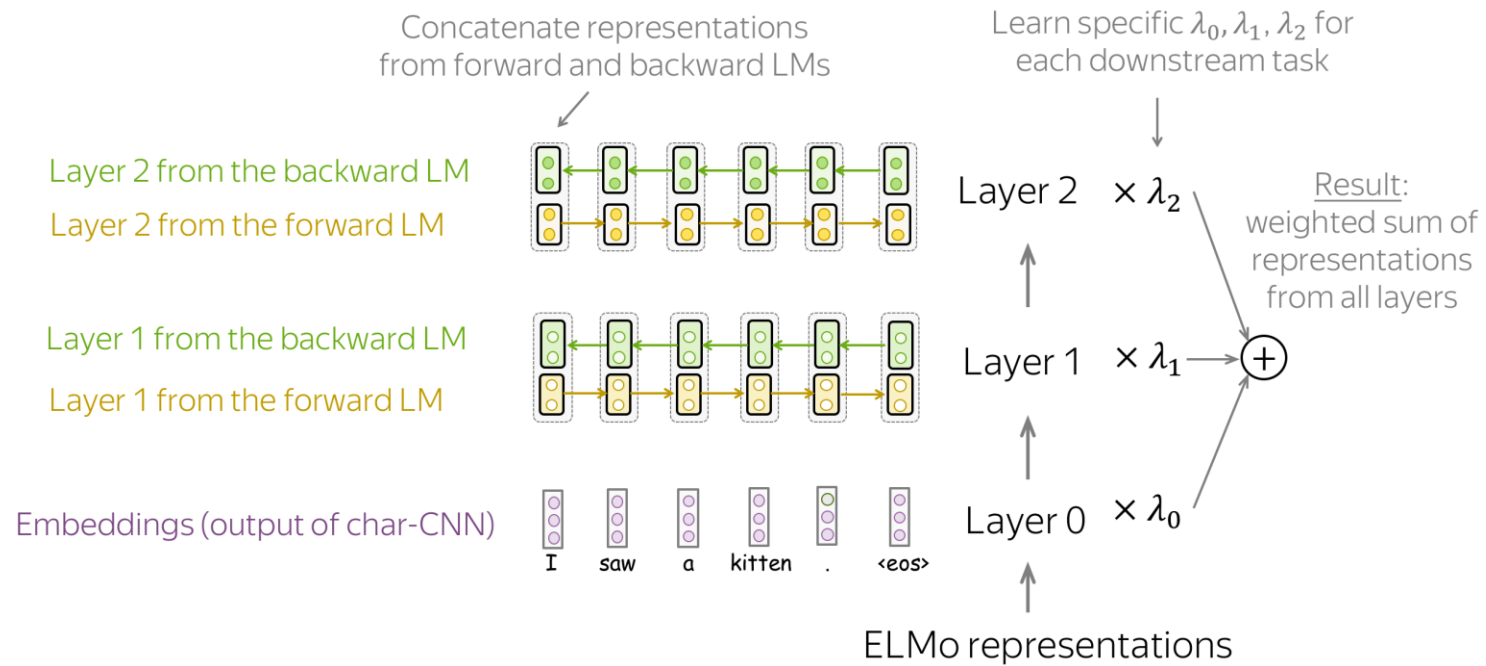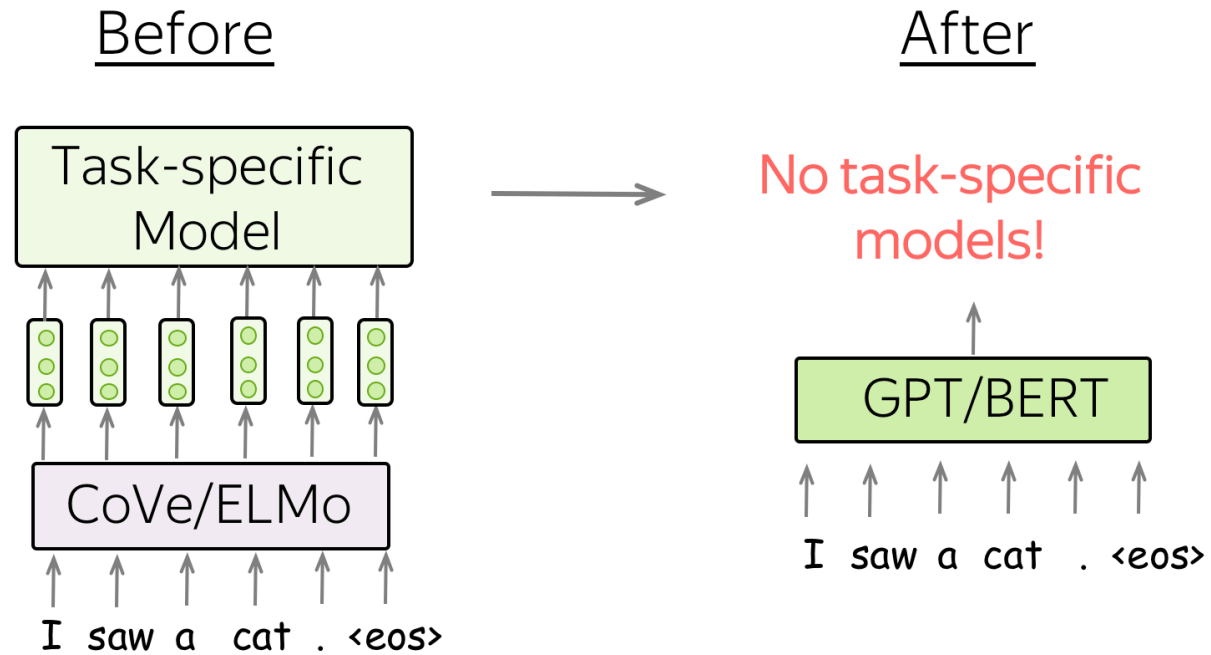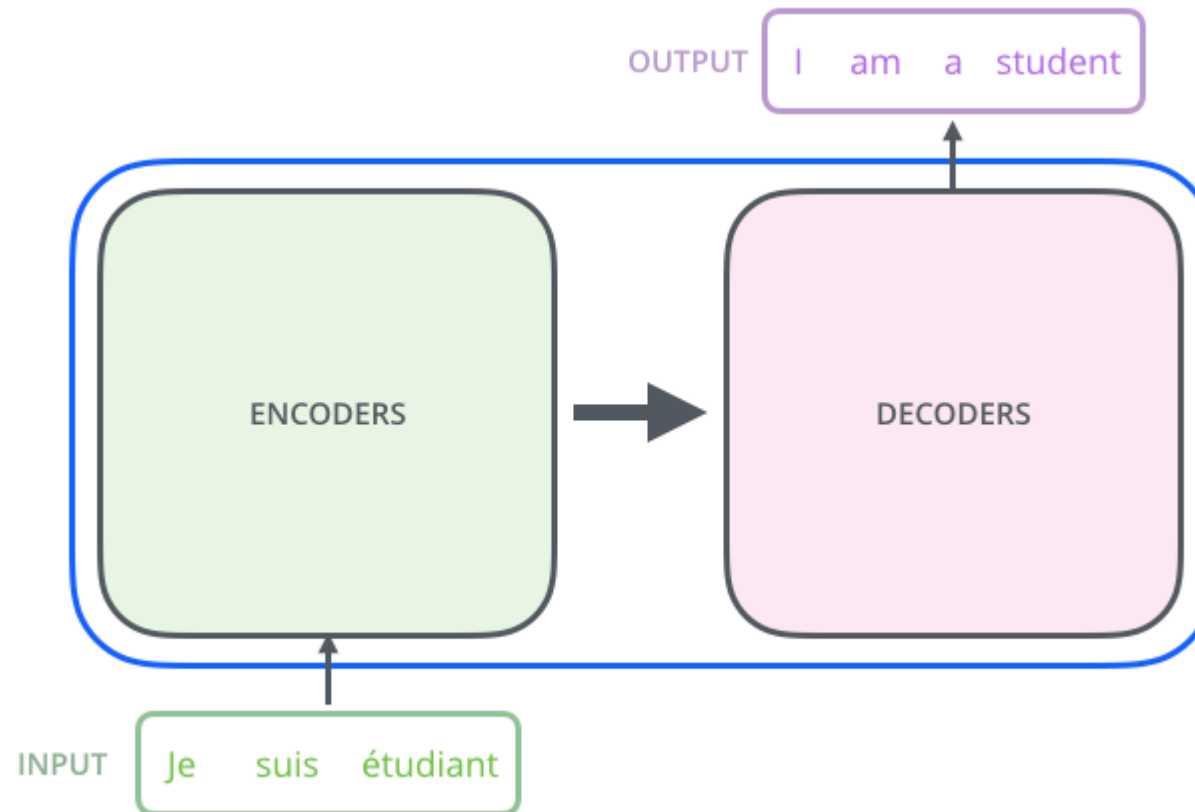- Transformer architecture overview
- BERT
- GPT

# Not just a cat, but the cat!



Source: https://lena-voita.github.io/nlp_course/transfer_learning.html

# Train a "translator" from word-specific to "contextual" space



Source: https://lena-voita.github.io/nlp_course/transfer_learning.html

# Multiple layers to capture low-level and high-level context



Source: https://lena-voita.github.io/nlp_course/transfer_learning.html

# From embedding generator to a universal model

Before

After

Task-specific
Model

No task-specific
models!

GPT/BERT

CoVe/ELMo

I  saw  a  cat  .  <eos>

I  saw  a  cat  .  <eos>

# Table of Content

- The power of transfer learning
- From word-specific to contextual embeddings
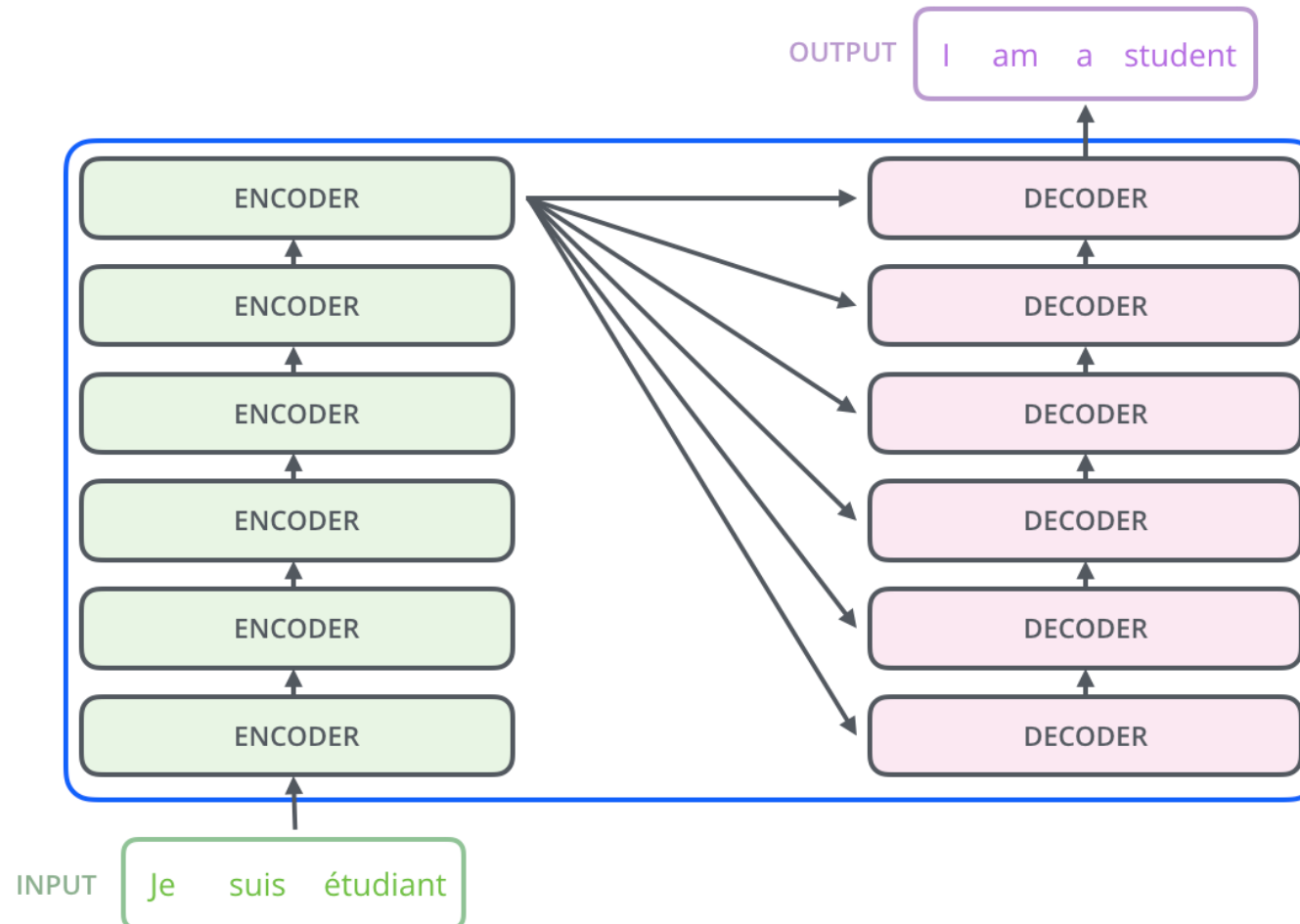- **Transformer architecture overview**
- BERT
- GPT

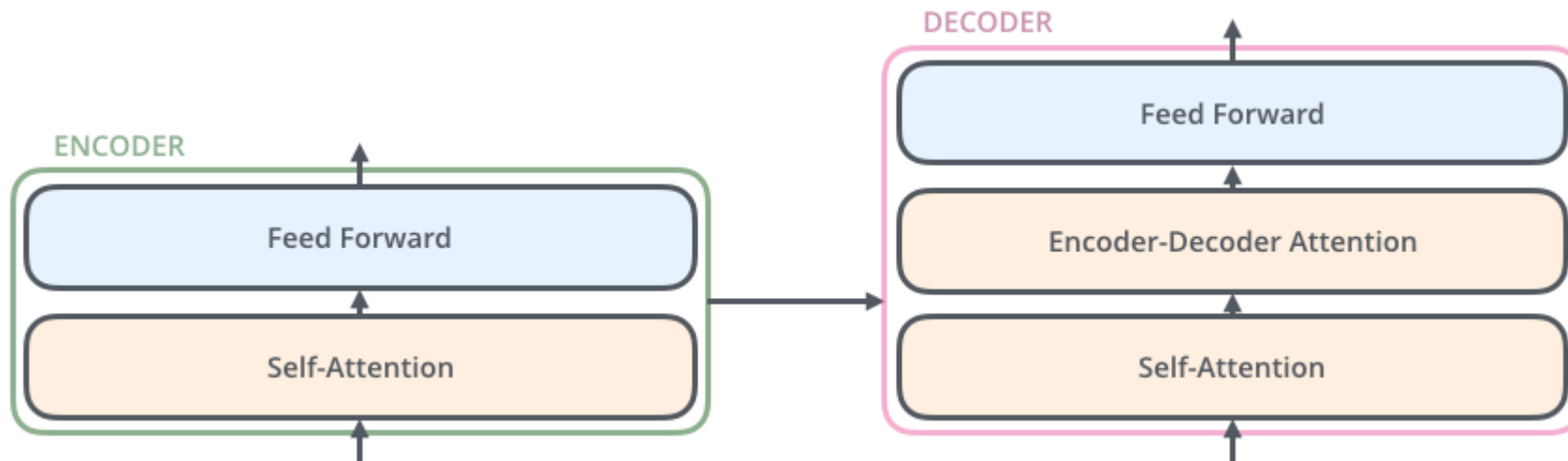# Transformer is an example of Encoder-Decoder architecture

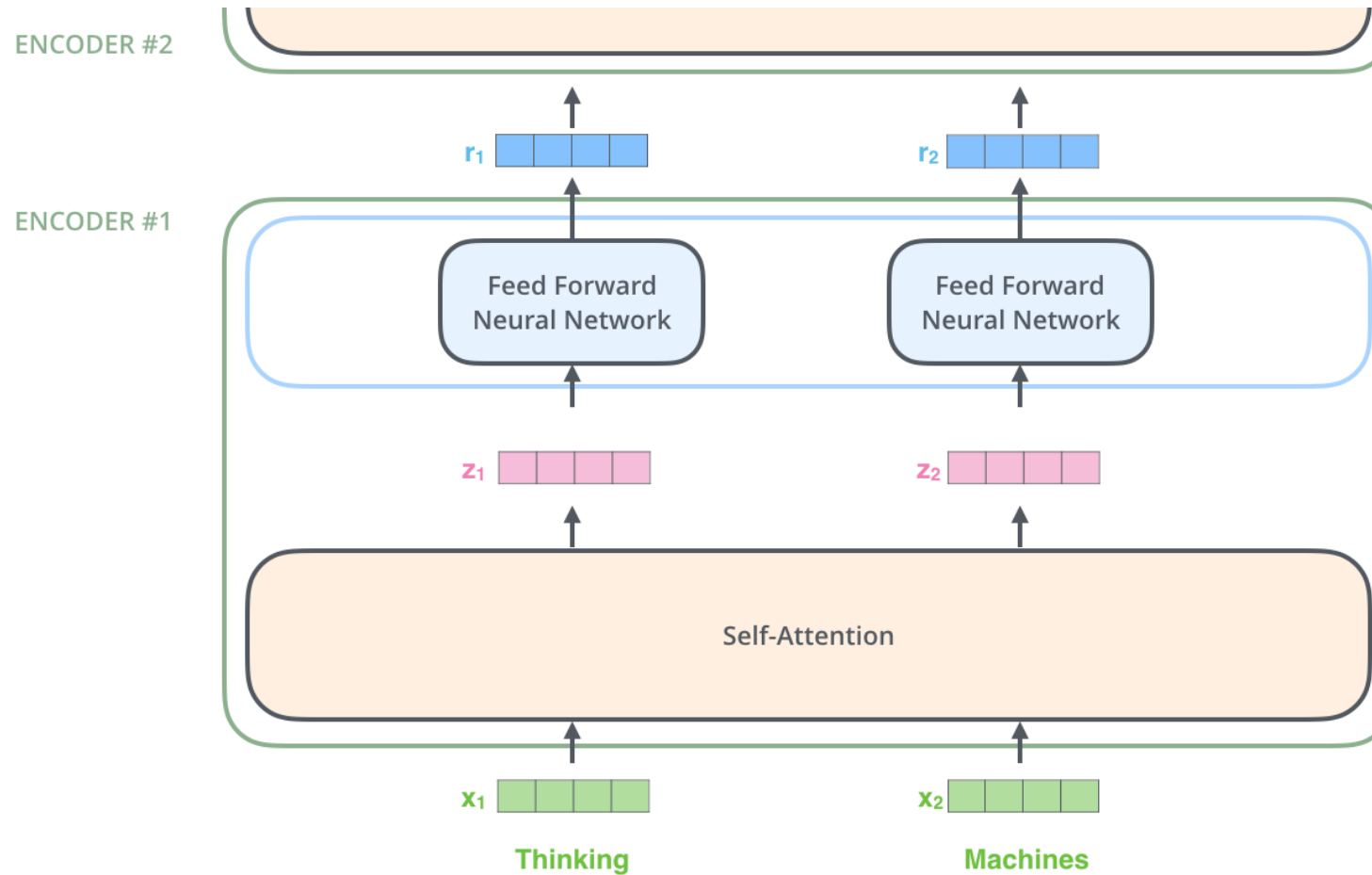# Transformer is an example of Encoder-Decoder architecture



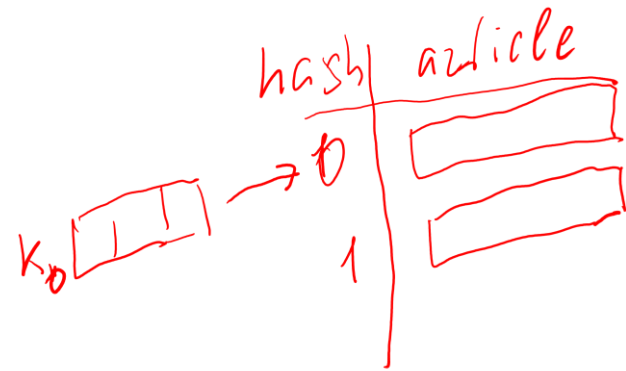Source: https://jalammar.github.io/illustrated-transformer/

# Transformer is an example of Encoder-Decoder architecture



Source: https://jalammar.github.io/illustrated-transformer/

# There is always two of them: the Attention and the FFN



Source: https://jalammar.github.io/illustrated-transformer/

# There is always two of them: the Attention and the FFN



Source: https://jalammar.github.io/illustrated-transformer/

# Inside the self-Attention



| | 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 | $z_1$ | $z_2$ |

Source: https://jalammar.github.io/illustrated-transformer/

# Inside the self-Attention: Matrix View

# Inside the self-Attention: Matrix View

$$Q = \begin{bmatrix} - q_1 - \\ - q_2 - \end{bmatrix}$$

$$k^T = \begin{bmatrix} | & | \\ k_1 & k_2 \\ | & | \end{bmatrix}$$

$$QK^T = \begin{bmatrix} \langle q_1, k_1 \rangle & \langle q_1, k_2 \rangle \\ \langle q_2, k_1 \rangle & \langle q_2, k_2 \rangle \end{bmatrix} =$$
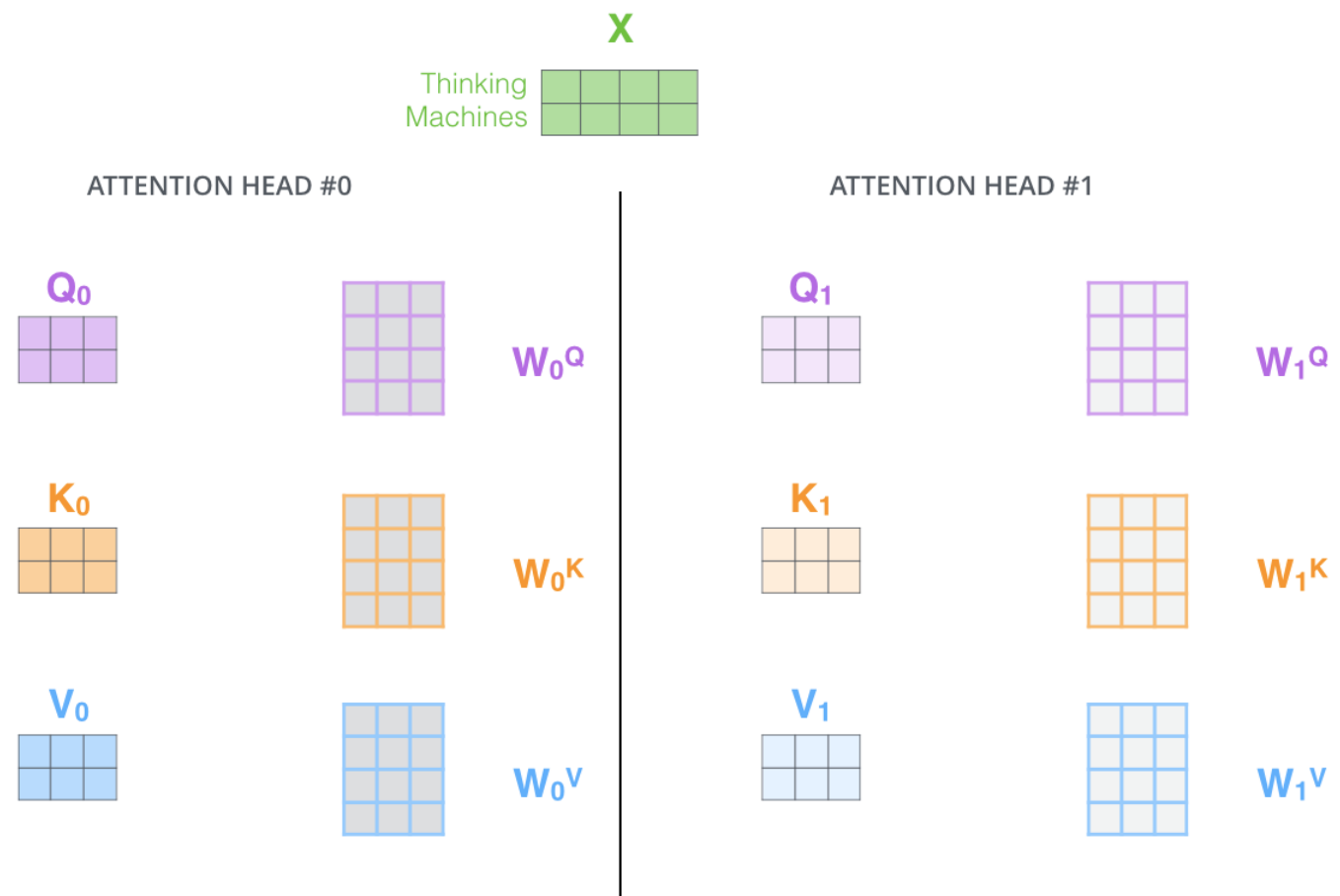
$$\text{softmax}\left( \frac{Q \times K^T}{\sqrt{d_k}} \right) V$$

**Q** **K**ᵀ **V**

$q_1$   $k_1$   $v_1$   $v_2$

**Z**

Source: https://jalammar.github.io/illustrated-transformer/

# A beast with many heads

# A beast with many heads

# A beast with many heads

1) Concatenate all the attention heads

$Z_0$  $Z_1$  $Z_2$  $Z_3$  $Z_4$  $Z_5$  $Z_6$  $Z_7$
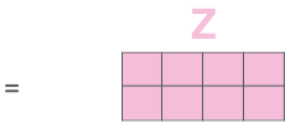
2) Multiply with a weight
matrix $W^O$ that was trained
jointly with the model

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

=

Source: https://jalammar.github.io/illustrated-transformer/

# A multi-head attention overview
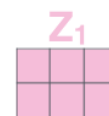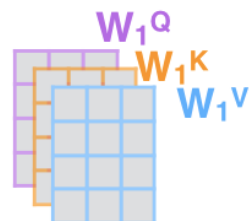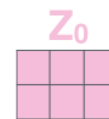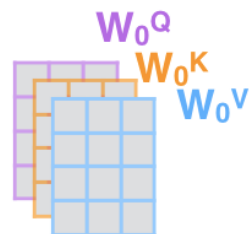
1) This is our
input sentence*

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
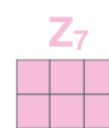produce the output of the layer

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

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

$Q_1$
$K_1$
$V_1$

$Z_1$

$Z$

...

...

...

$R$

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

$Q_7$
$K_7$
$V_7$

$Z_7$

Source: https://jalammar.github.io/illustrated-transformer/

# Each head focuses on a specific representation

# As opposed to RNNs, Transformers do not track the position implicitly

# Absolute Positional Encoding

$$\mathrm{PE}(pos, 2i) = \sin\left(pos/10000^{2i/d_{model}}\right)$$

$$\mathrm{PE}(pos, 2i + 1) = \cos\left(pos/10000^{2i/d_{model}}\right)$$

# What it looks like



Source: https://kazemnejad.com/blog/transformer_architecture_positional_encoding/

# Skip Connection and Layer Normalization for a robust training



Source: https://kazemnejad.com/blog/transformer_architecture_positional_encoding/

# Bringing it all together



Source: https://kazemnejad.com/blog/transformer_architecture_positional_encoding/

# There is always two of them: the Attention and the FFN



Source: https://jalammar.github.io/illustrated-transformer/

# Attention is all you need

# Attention is all you need (but not really)



Source: https://kazemnejad.com/blog/transformer_architecture_positional_encoding/

# Table of Content

- The power of transfer learning
- From word-specific to contextual embeddings
- Transformer architecture overview
- **BERT**
- GPT

# BERT is just an Encoder part



Image source: https://jalammar.github.io/illustrated-transformer/

# Using encoder as an embedding generator



I  saw  a  grey  cat  on  a  mat  .  <eos>

Model architecture:

- Transformer's encoder

What is special about it:

- Training objectives
  - MLM: Masked language modeling
  - NSP: Next sentence prediction

- The way it is used
  - No task-specific models

Source: https://lena-voita.github.io/nlp_course/transfer_learning.html

# Objective one: tell whether the two sequences are consecutive

**[CLS]**: Special token
- Training time: predict if sentences are consecutive or not (Next Sentence Prediction /NSP objective)
- Test time: downstream tasks (e.g., classification)

[SEP]: Special token-separator

[CLS] My dog is very cute [SEP] He likes playing in the garden with me [SEP]

Segment A          Segment B

Training on pairs of sentences: either consecutive or random (50%/50%)

Source: https://lena-voita.github.io/nlp_course/transfer_learning.html

# Using encoder as an embedding generator



Training time: predict if sentences are consecutive (NSP objective)
Test time: classification

Training time: MLM objective

Model (Transformer encoder)

several layers

Input

positions
0 1 2 3 4 …

segments
A A A A A B B B B B

tokens
[CLS] My dog is …

Training on pairs of sentences: either consecutive or random (50%/50%)

[CLS] My dog is cute [SEP] He is very fluffy [SEP]

Segment A          Segment B

Source: https://lena-voita.github.io/nlp_course/transfer_learning.html

# Objective two: Masked Language Modeling



Source: https://lena-voita.github.io/nlp_course/transfer_learning.html
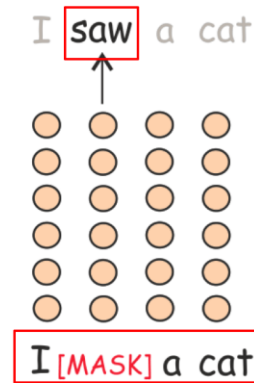
# LM vs. MLM

### Language Modeling

- Target: next token

- Prediction: $P(* | \text{I saw})$
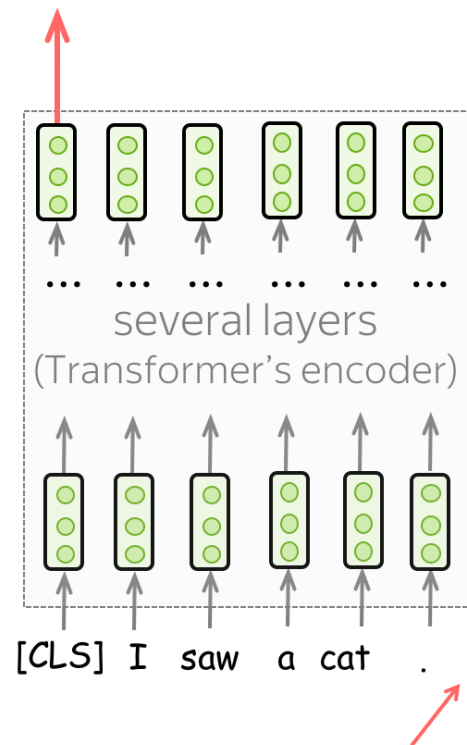


left-to-right, does
not see future

### Masked Language Modeling

- Target: current token (the true one)

- Prediction: $P(* | \text{I [MASK] a cat})$



sees the whole text, but
something is corrupted

Source: https://lena-voita.github.io/nlp_course/transfer_learning.html

# Single Sentence Classification



Source: https://lena-voita.github.io/nlp_course/transfer_learning.html

# Sentence Pair Classification



Source: https://lena-voita.github.io/nlp_course/transfer_learning.html

# Question Answering

# Input tagging

# Table of Content

- The power of transfer learning

- From word-specific to contextual embeddings

- Transformer architecture overview

- BERT

- **GPT**

# GPT is just a Decoder part



Image source: https://jalammar.github.io/illustrated-transformer/

# Decoder as a universal model

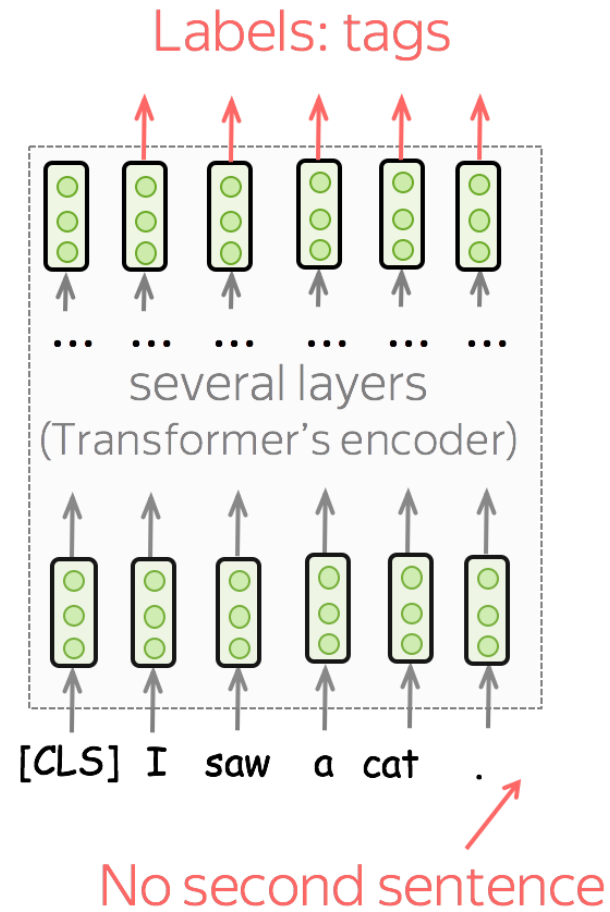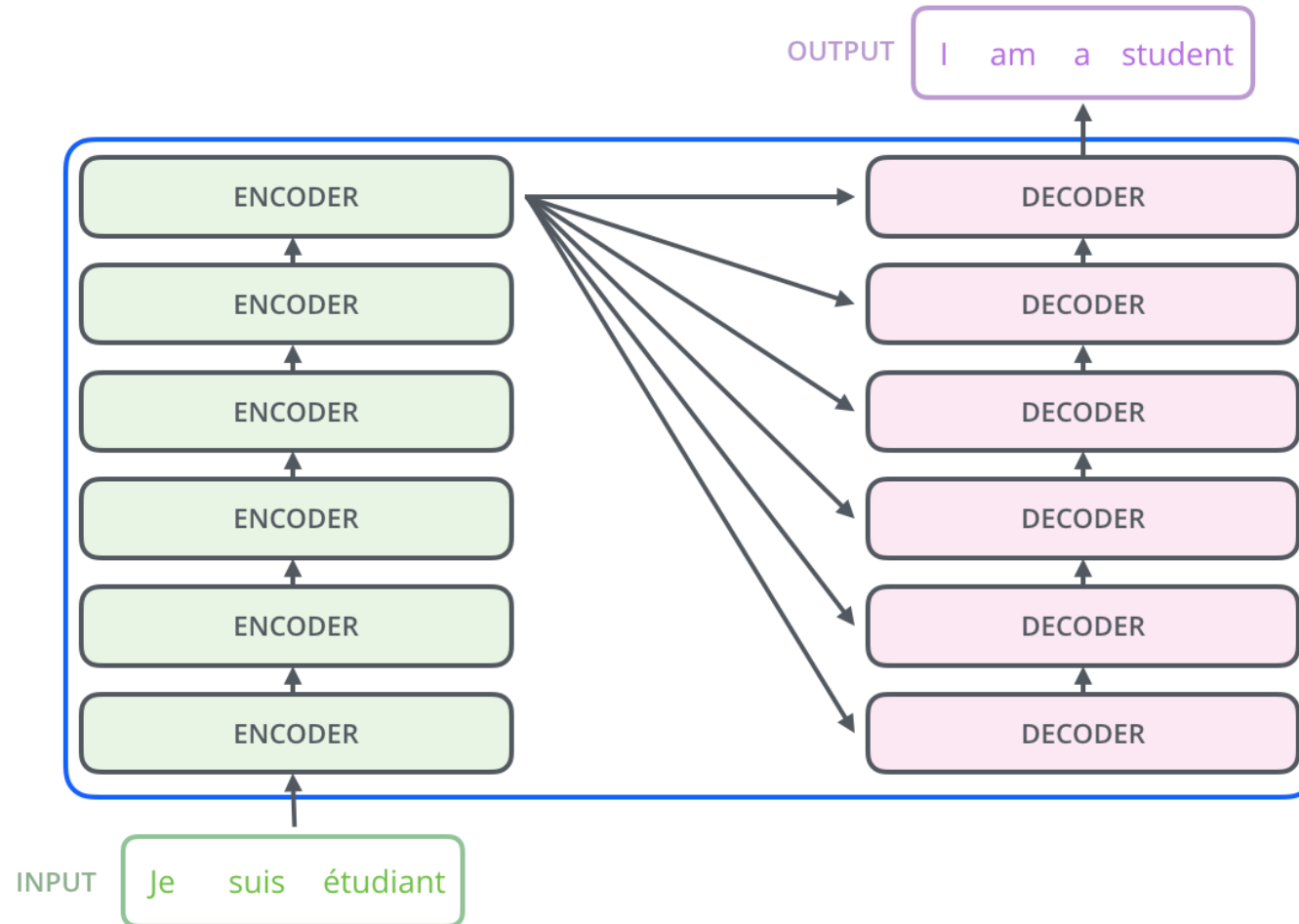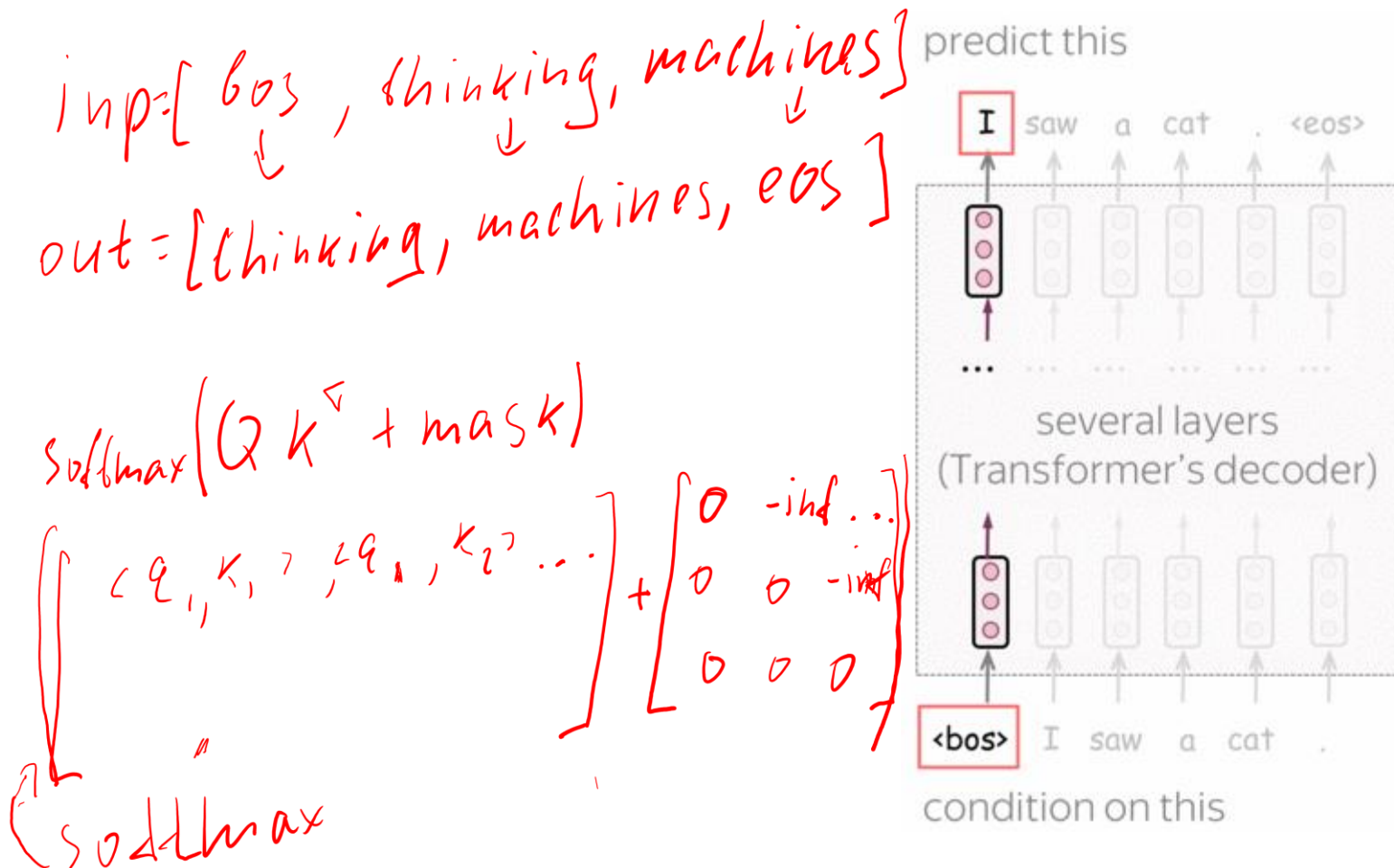inp=[ bos , thinking, machines]

out=[thinking, machines, eos]

$$\text{softmax}\left( Q k^{\nabla} + \text{mask} \right)$$

$$\left[ \begin{array}{c} \langle q_1, k_1 \rangle , \langle q_1, k_2 \rangle \cdots \\ \end{array} \right] + \left[ \begin{array}{ccc} 0 & -\inf & \cdots \\ 0 & 0 & -\inf \\ 0 & 0 & 0 \end{array} \right]$$

softmax

predict this

I saw a cat . \<eos\>

several layers
(Transformer's decoder)

\<bos\> I saw a cat .

condition on this

Source: https://cdn.openai.com/research-covers/language-unsupervised/language_understanding_paper.pdf

# Decoder as a universal model



Source: https://cdn.openai.com/research-covers/language-unsupervised/language_understanding_paper.pdf