

SAISON
23/24



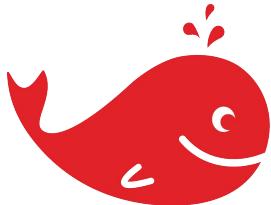
Formation
Introduction au
Deep Learning
Séquence n°9

“Attention is All You Need : Transformers”



FIDLE





FIDLE

Formation

Introduction au Deep Learning

Questions and answers :

<https://fidle.cnrs.fr/q2a>

Accompanied by :

AI Support (dream) Team of IDRIS

Directed by :

Agathe, Baptiste et Yanis - UGA/DAPI

Léo, Thibaut, Kamel - IDRIS



Formation

Introduction au Deep Learning



<https://fidle.cnrs.fr/listeinfo>

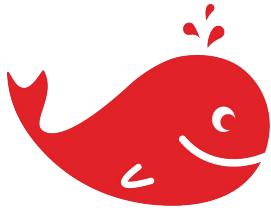
Fidle information list

Agoria

<http://fidle.cnrs.fr/agoria>

AI exchange list

New !



FIDLE

Formation

Introduction au Deep Learning

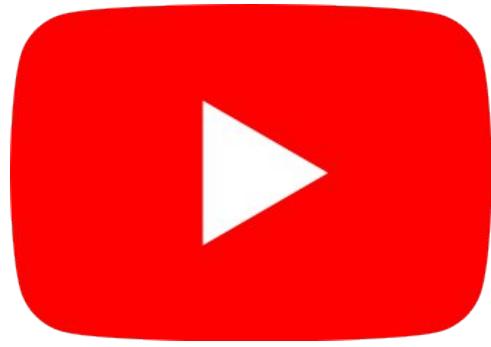


<https://listes.services.cnrs.fr/wws/info/devlog>
List of ESR* « Software developers » group



<https://listes.math.cnrs.fr/wws/info/calcul>
List of ESR* « Calcul » group

(*) ESR is Enseignement Supérieur et Recherche, french universities and public academic research organizations



YouTube

Abonnez-vous !



<https://fidle.cnrs.fr/youtube>

<https://www.youtube.com/@CNRS-FIDLE>



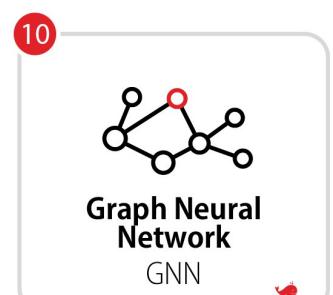
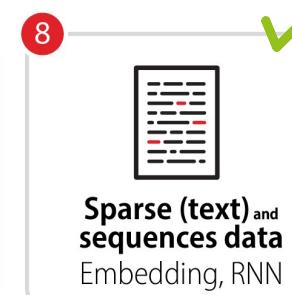
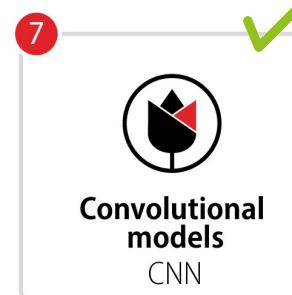
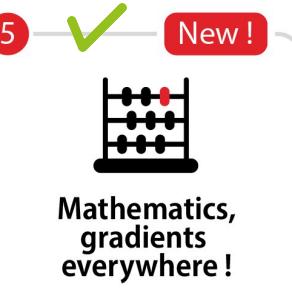
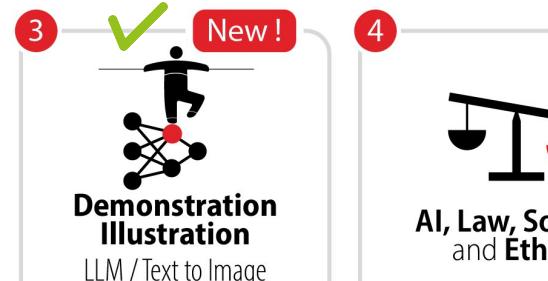
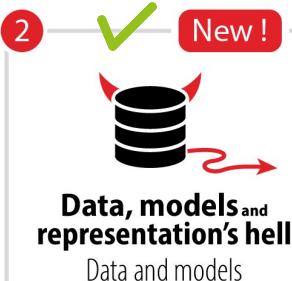
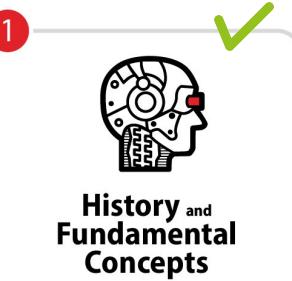
Le magazine de actualité en IA

Vendredi 16 février, 10h00

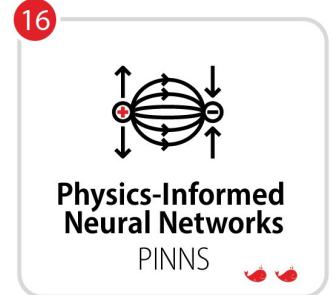
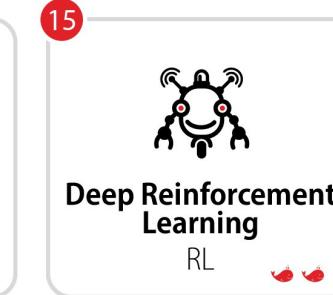
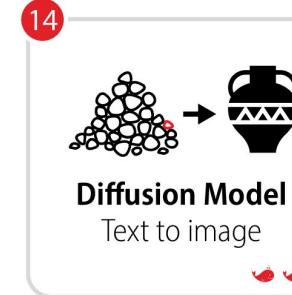
<http://www.idris.fr/panoramia.html>

<https://www.youtube.com/@IDRISCNRS>

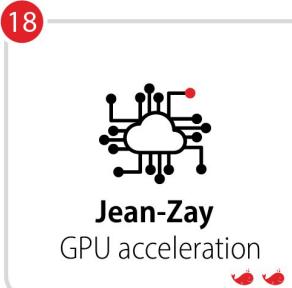
Bases, Concepts et Enjeux



L'IA comme un outil,



Acteur de l'IA



9



«Attention is
All You Need»

Transformers



1

What is a **Transformer**? The magic
of the **Attention Mechanism**

2

The different **Transformers**
architectures

3

Pre-training and Foundation
Models

4

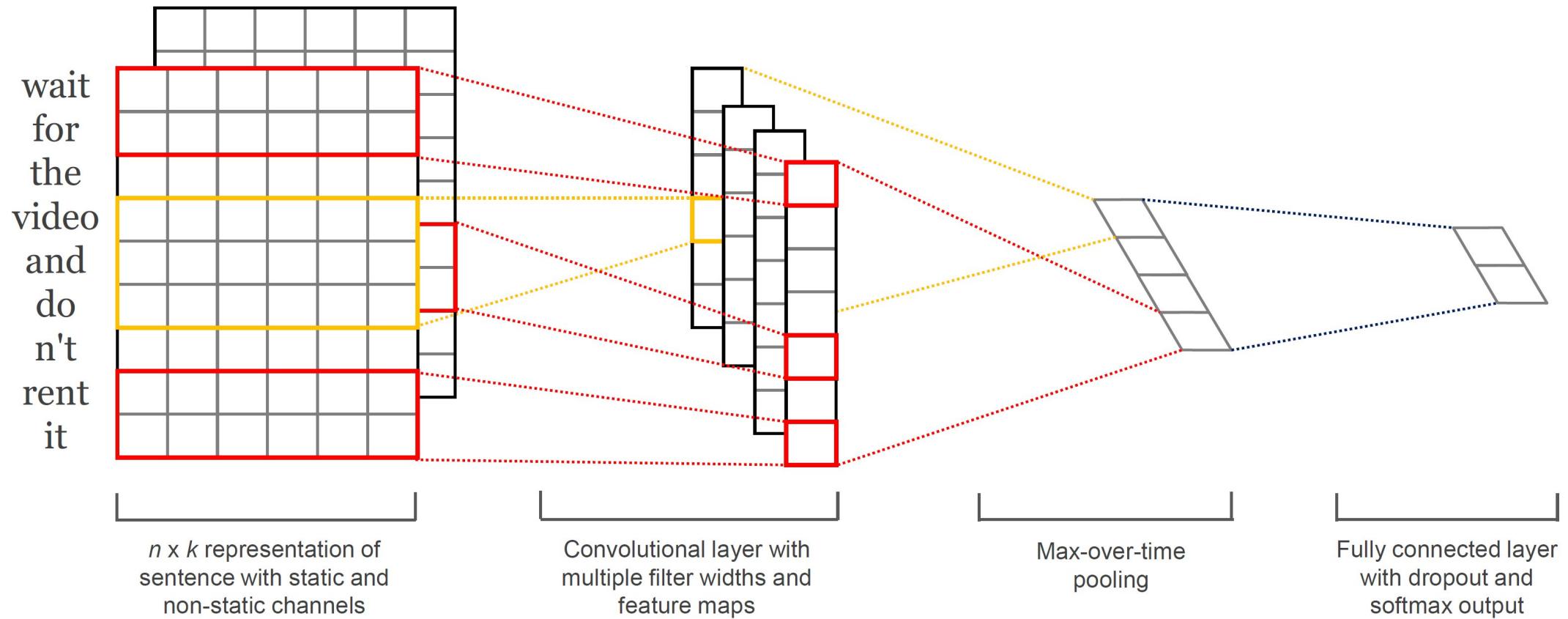
Specialization of Foundation
Models (especially LLM)

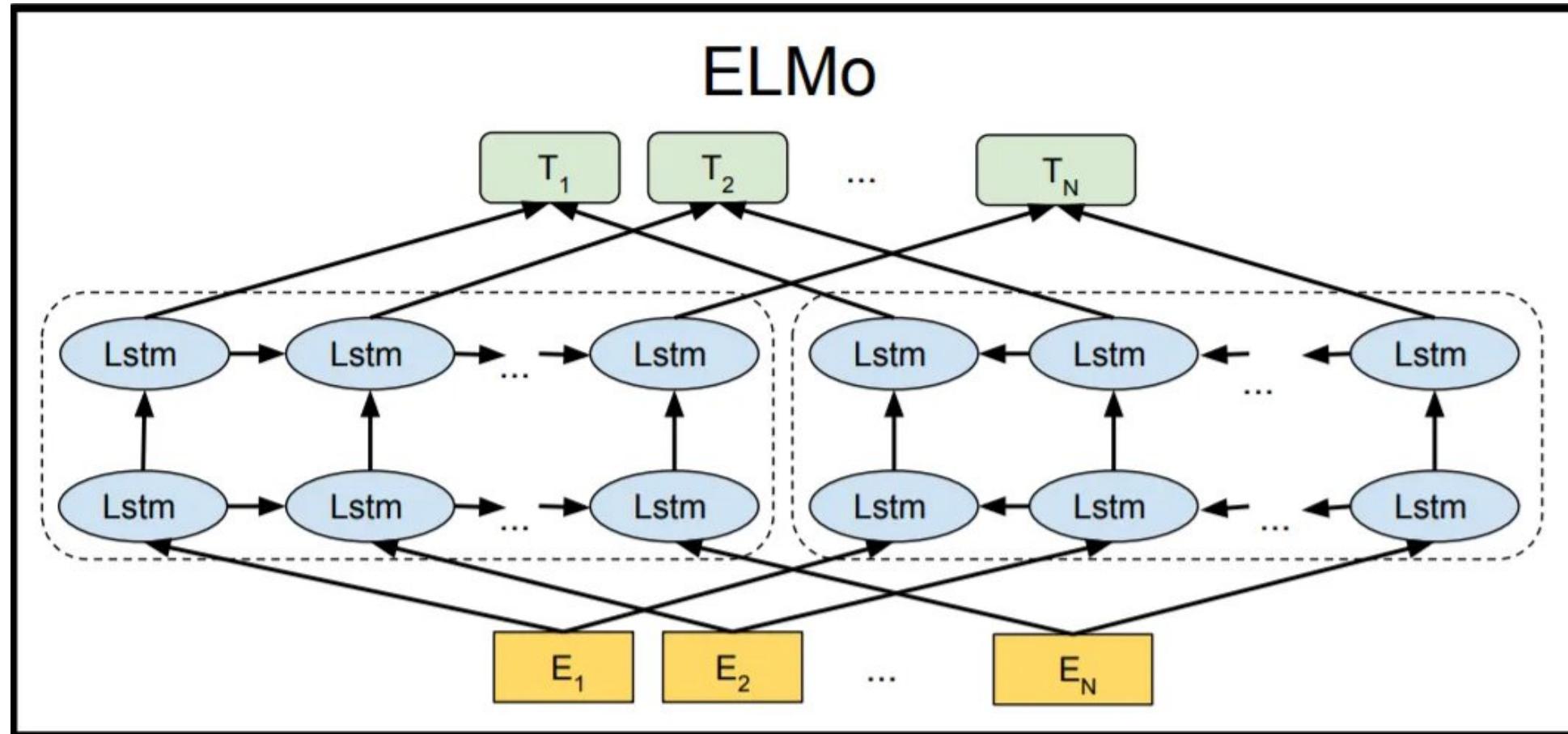
5

Example: IMDB Reviews

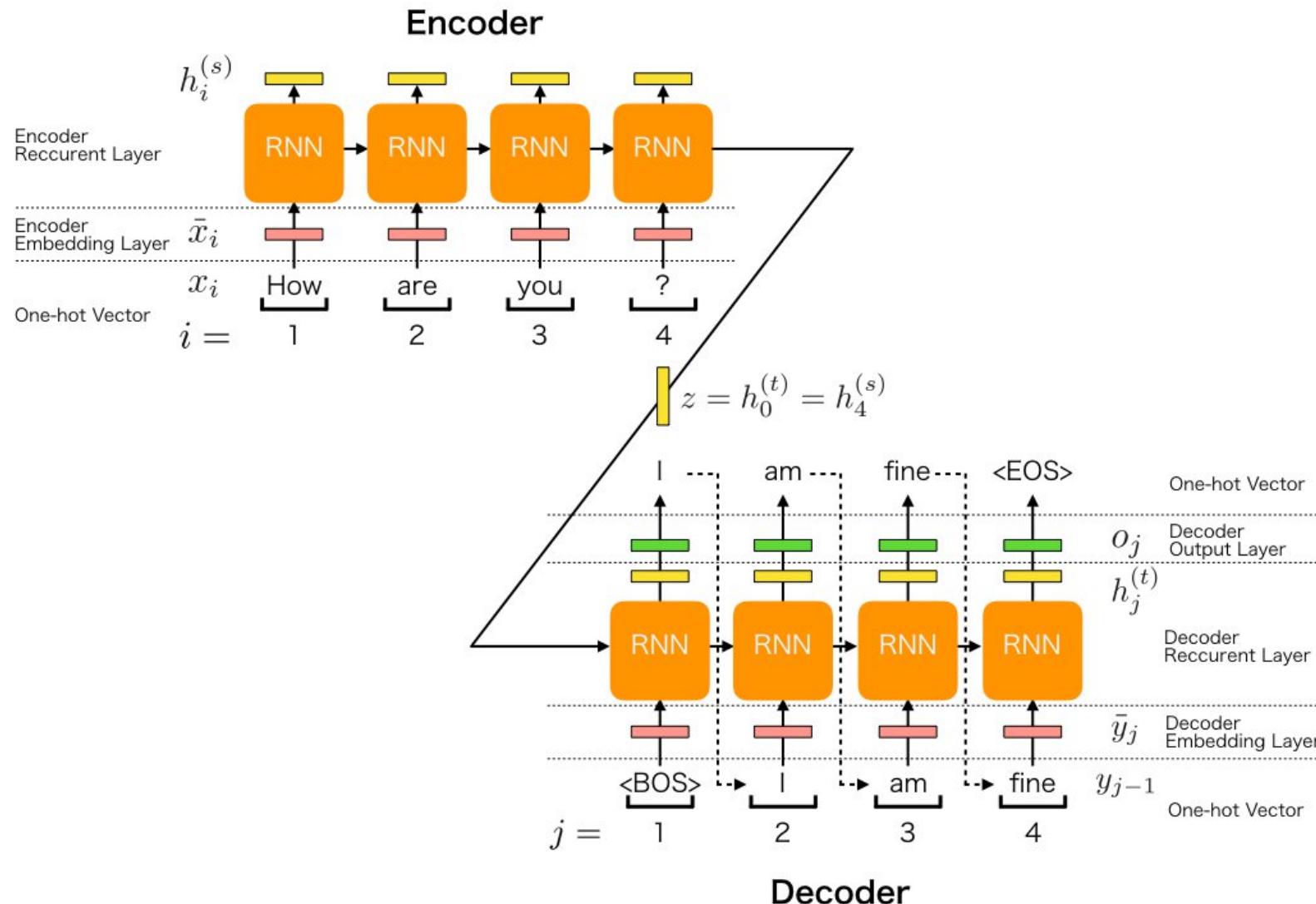


NLP back in the days





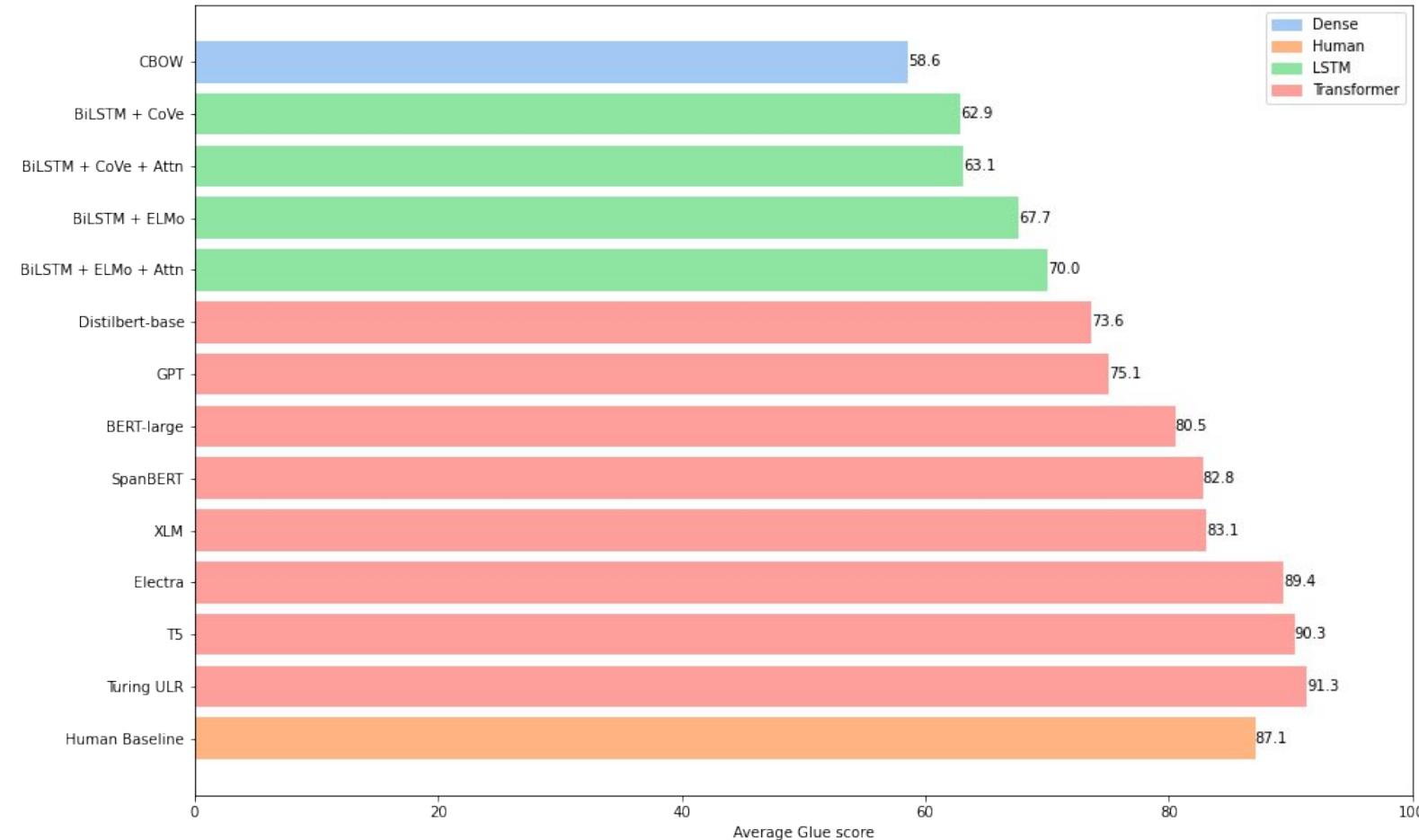
NLP back in the days



What do we want ?

- Process sequences (ideally the entire sentence)
- Easy to distribute on multiple GPUs
- Faster training than with RNN
- Initially for NLP tasks
- Allows to train huge models on gigantic datasets
- Allows for a pretraining session to pool trainings (at least partially) for multiple tasks

The King is dead. Long live the King!



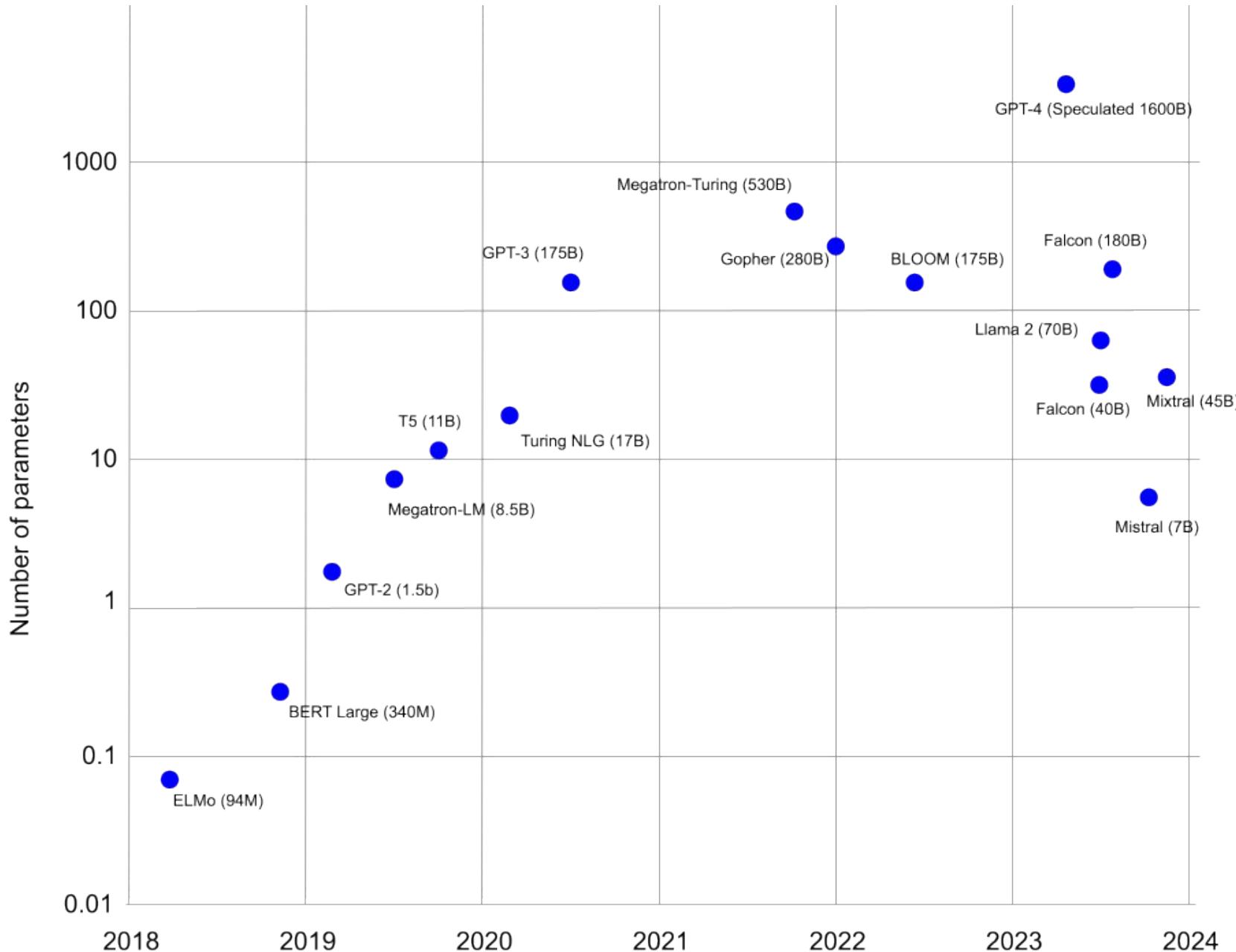
• This article is more than 5 months old

ChatGPT better than undergraduates at solving SAT problems, study suggests

Researchers at UCLA found GPT-3 solved 80% of reasoning problems correctly compared with 60% of humans



Size evolution of transformers



9



«Attention is
All You Need»

Transformers



1

What is a **Transformer**? The magic
of the **Attention Mechanism**

2

The different **Transformers**
architectures

3

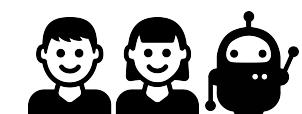
**Pre-training and Foundation
Models**

4

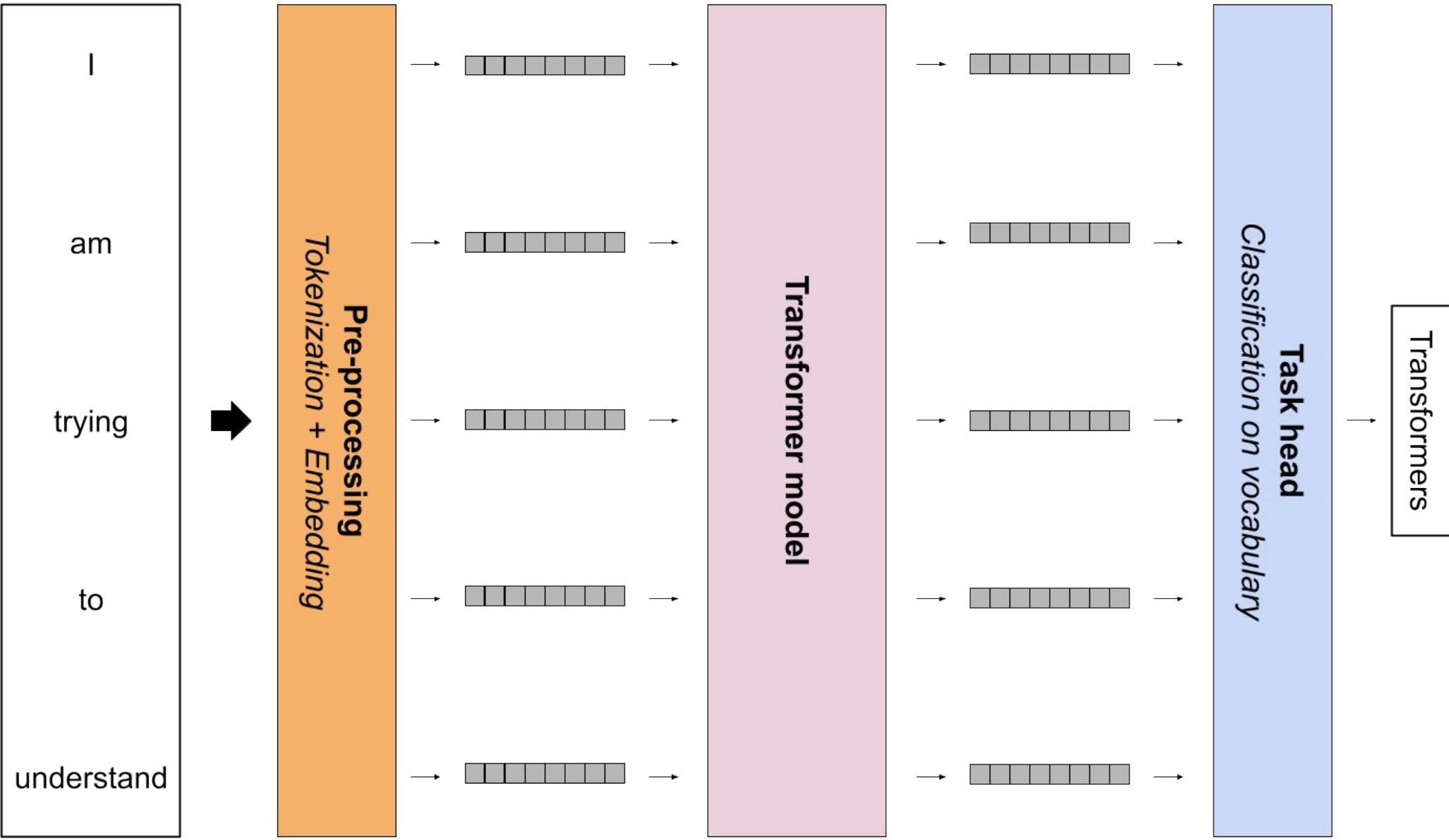
**Specialization of Foundation
Models** (especially LLM)

5

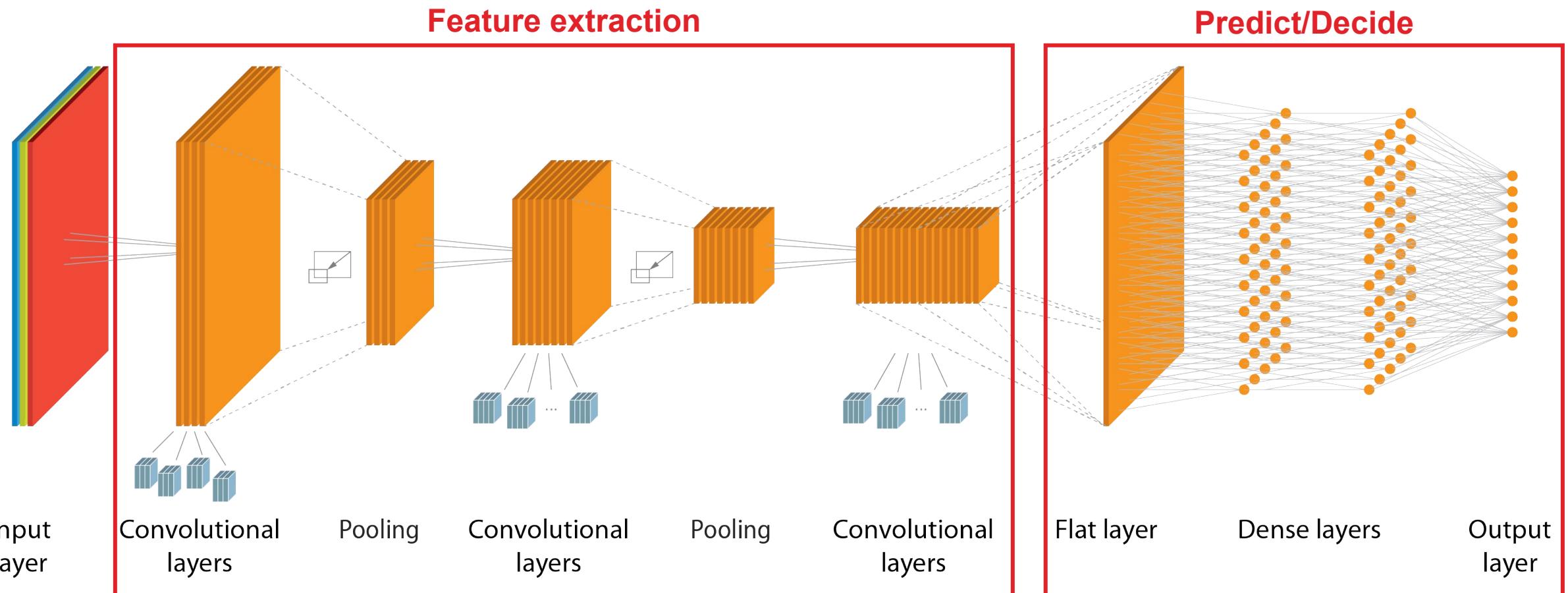
Example: IMDB Reviews



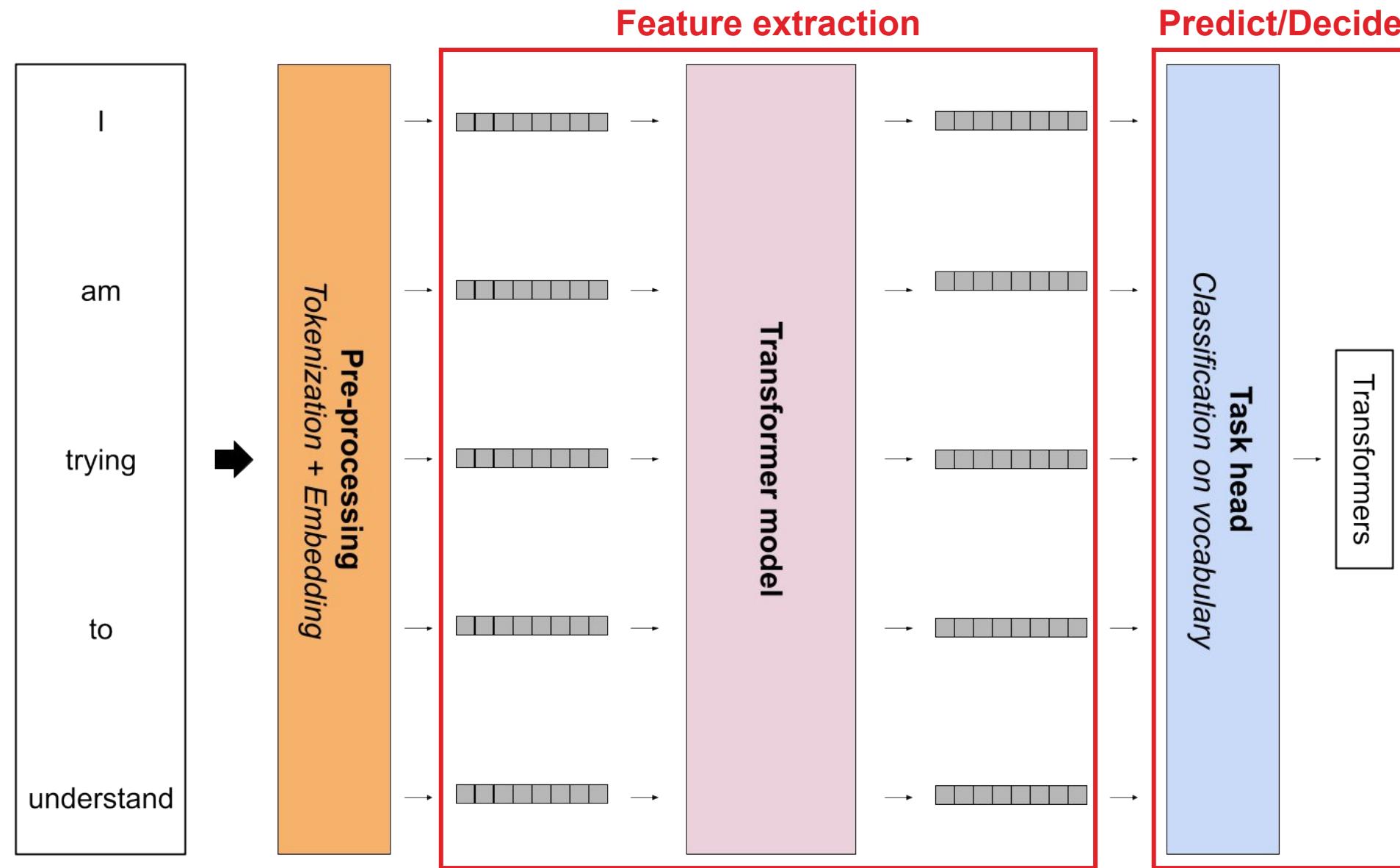
Example of NLP system



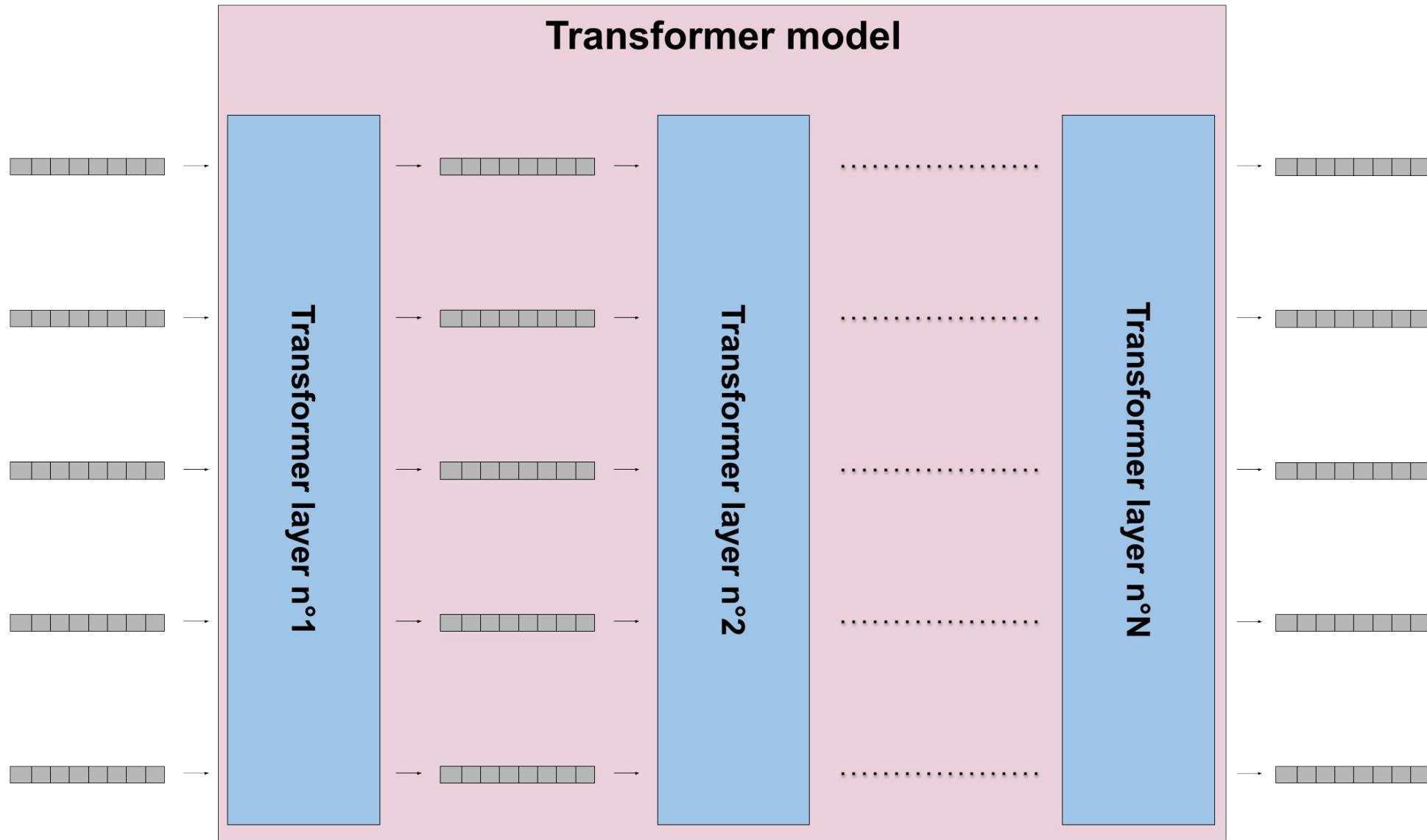
Feature extraction reminder



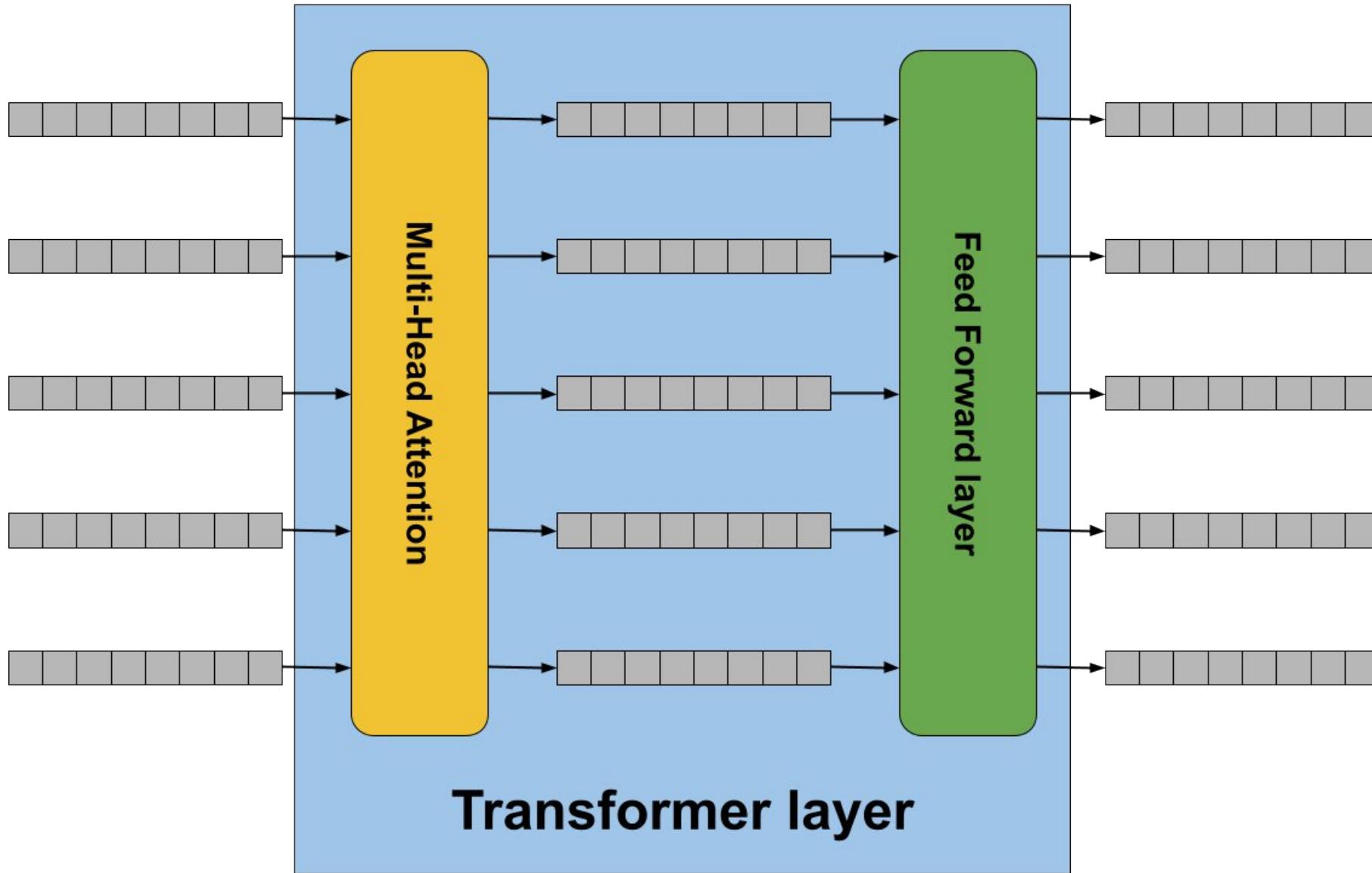
Feature extraction with Transformers



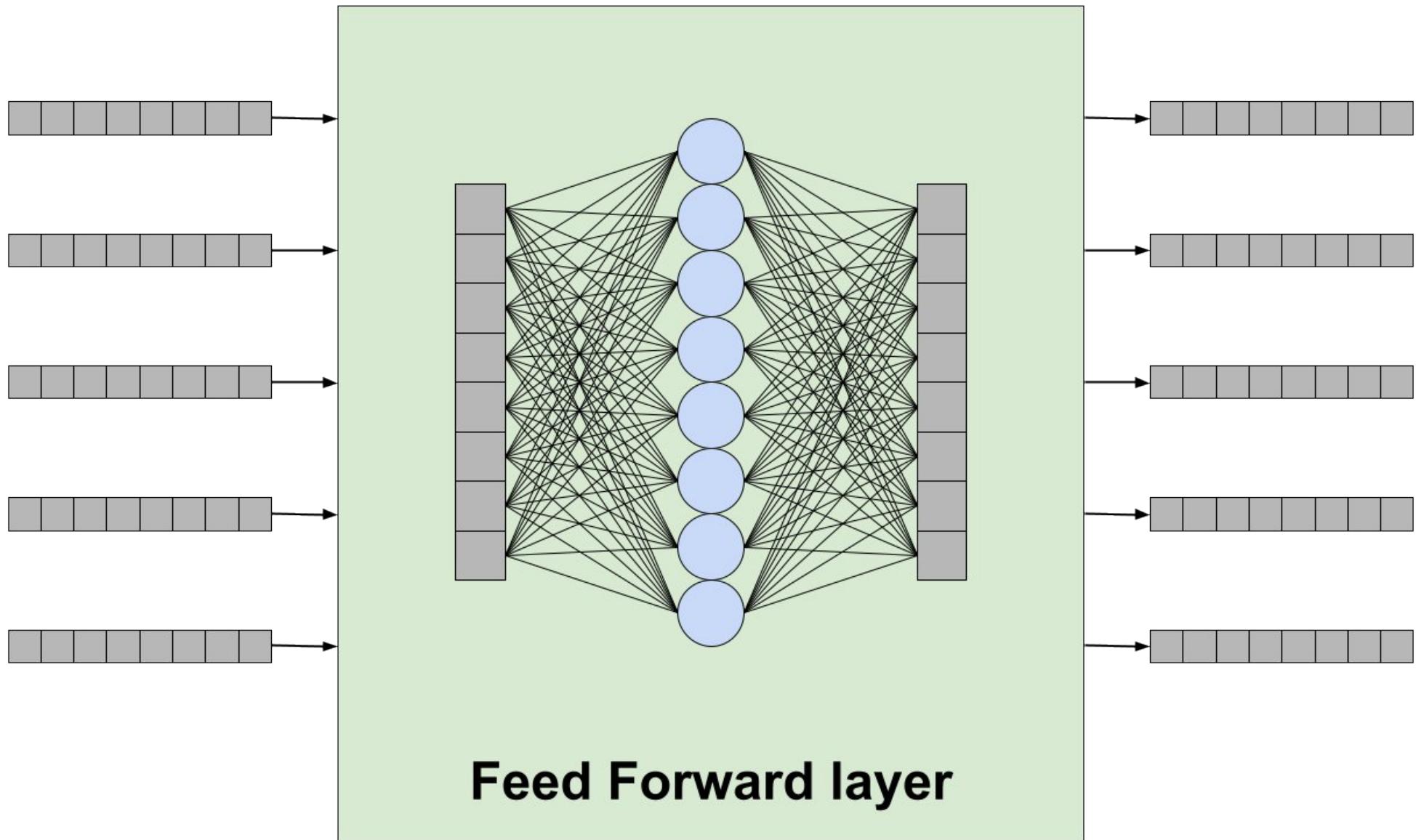
Vanilla Transformer architecture



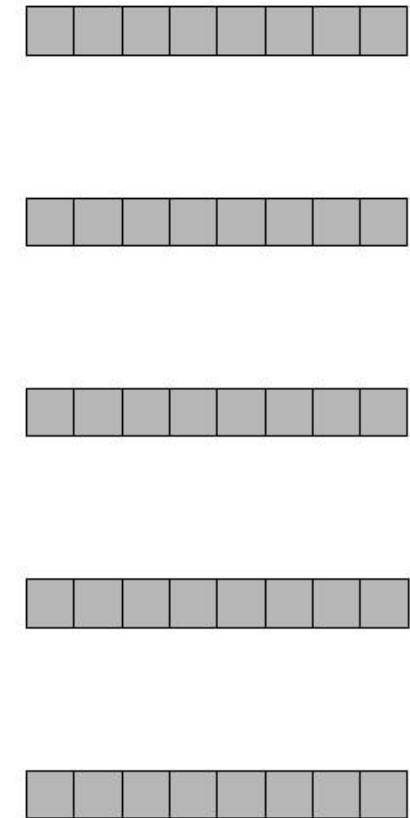
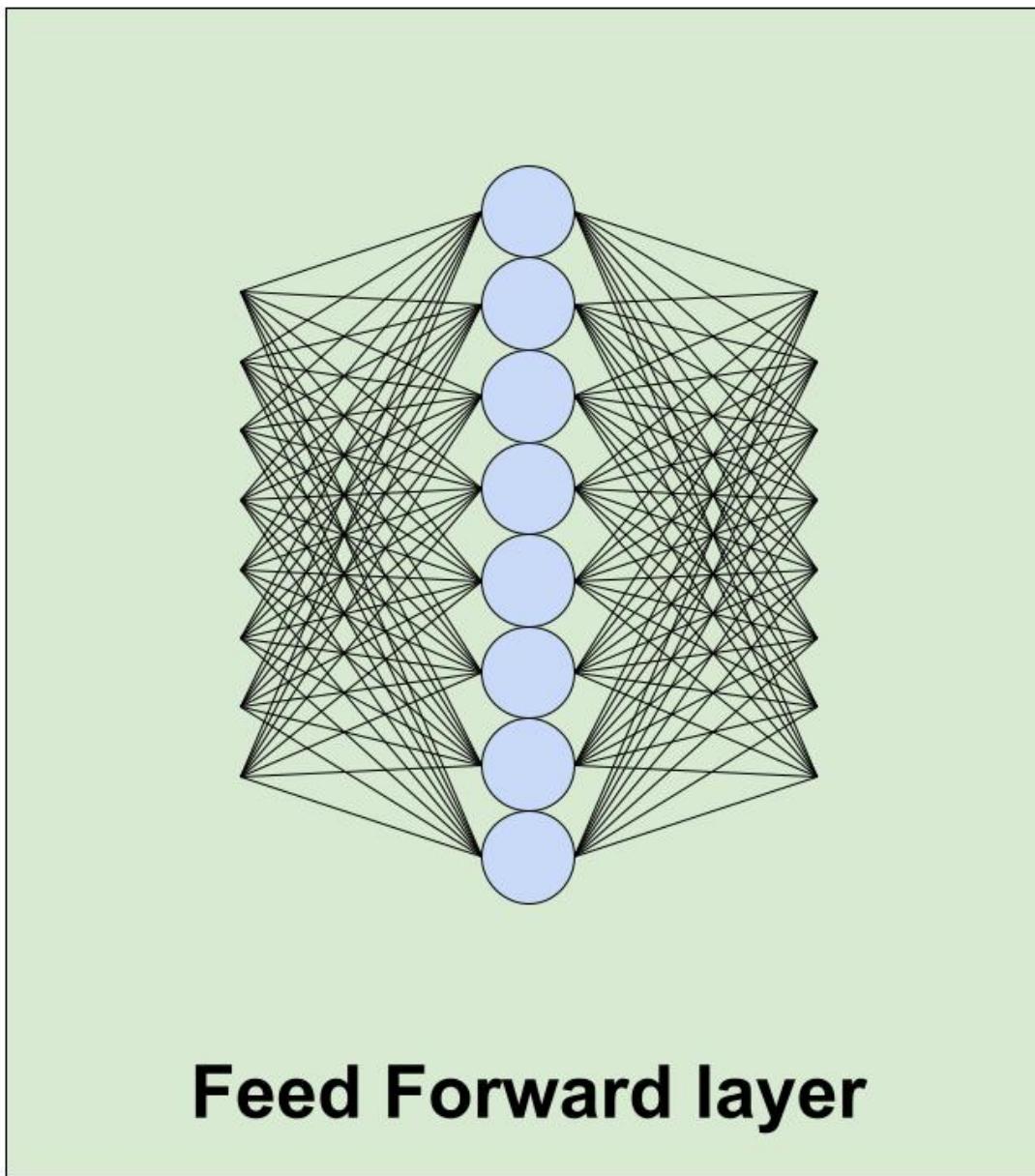
Transformer layer simplified



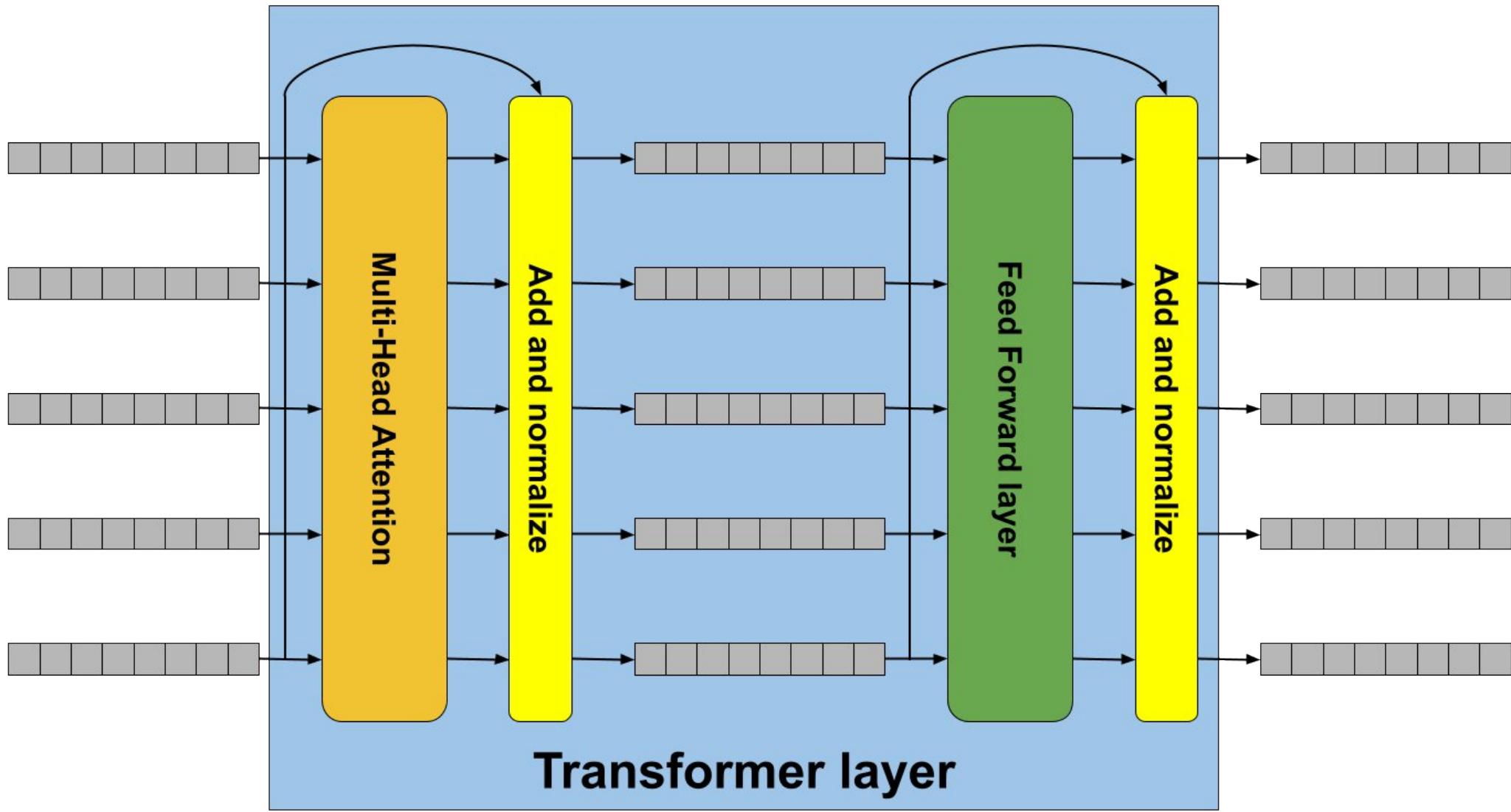
Feed Forward layer



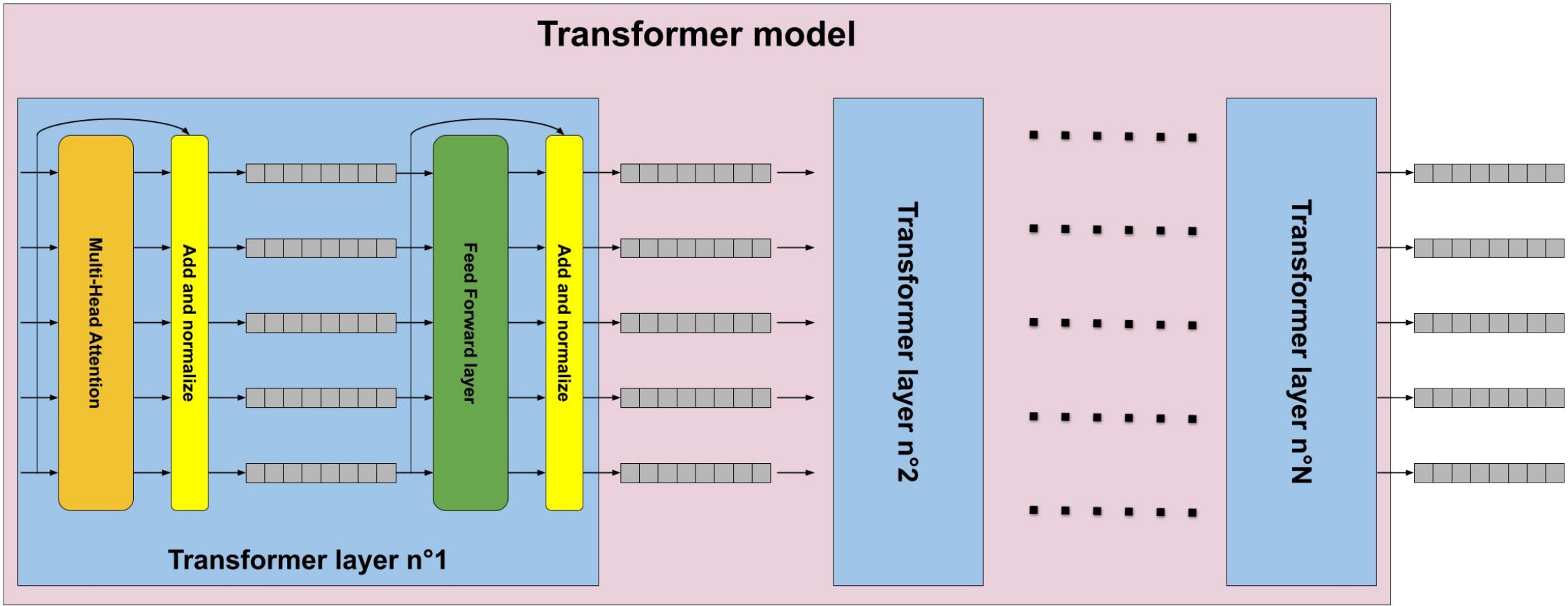
Feed Forward layer (animated)



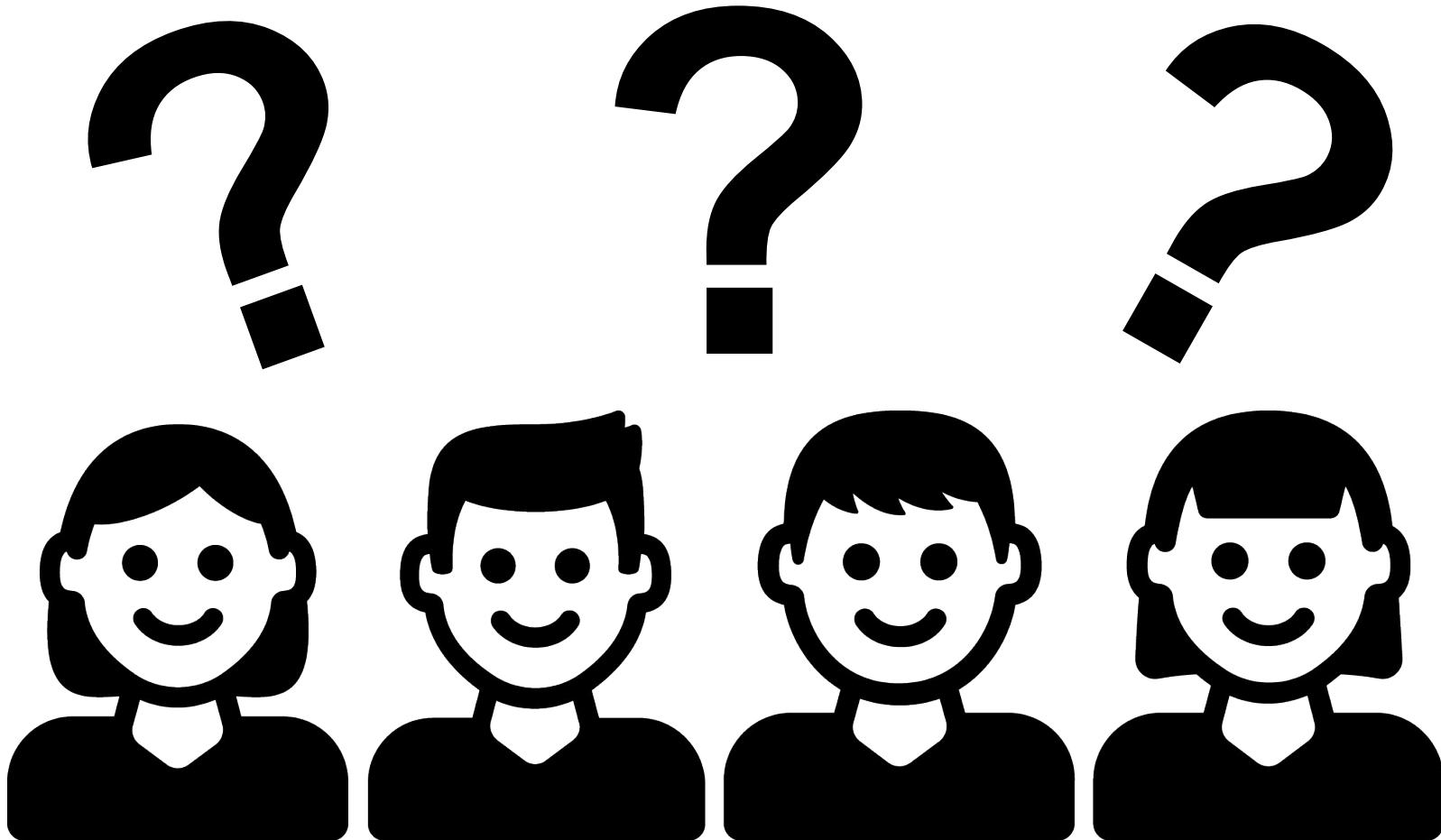
Transformer layer



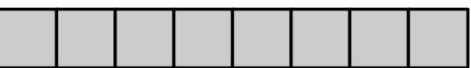
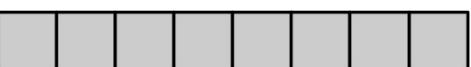
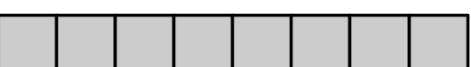
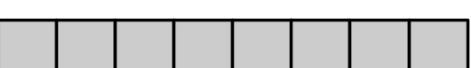
Vanilla Transformer architecture summary

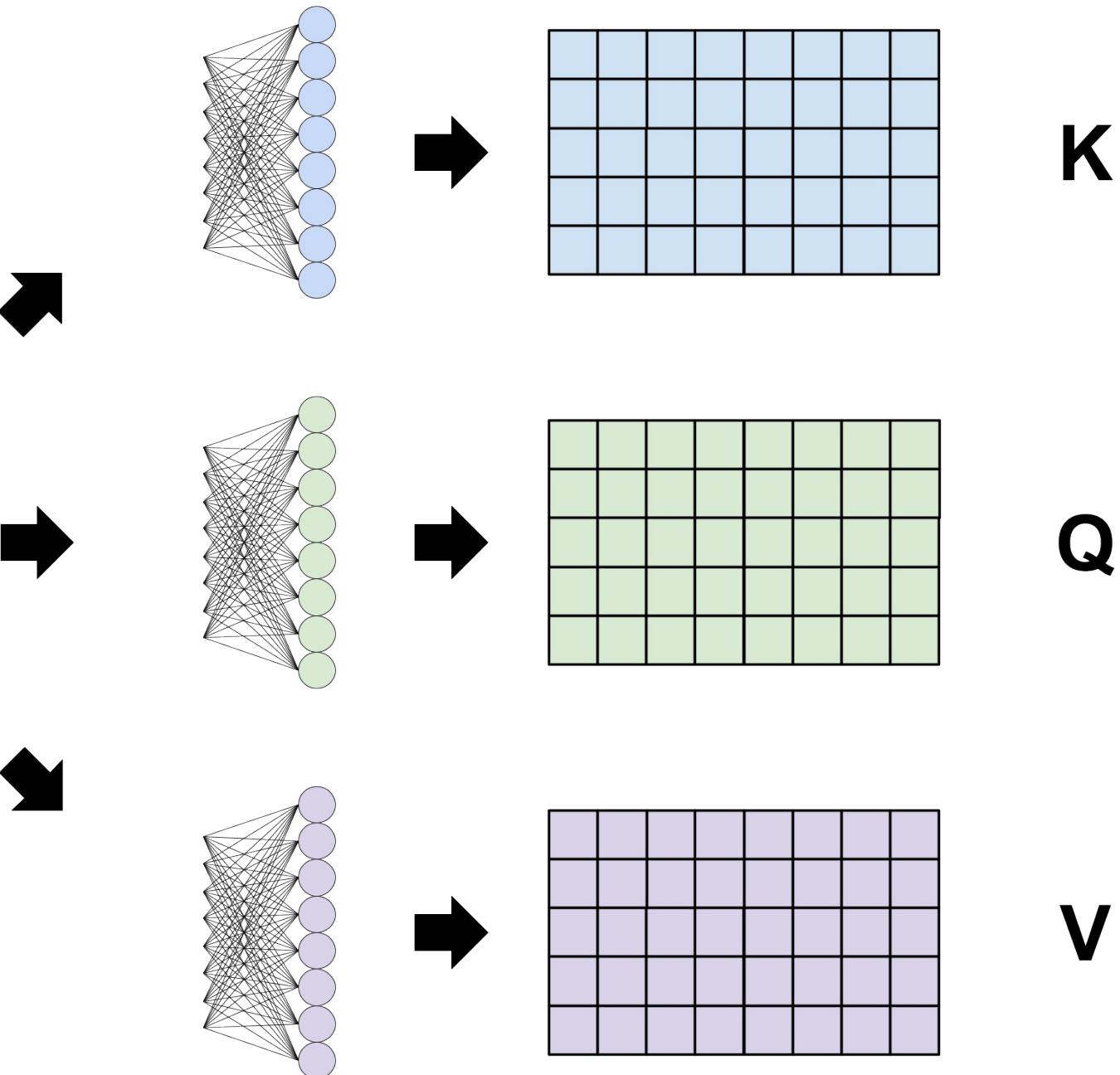


Question break #1

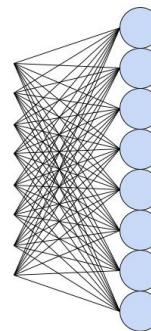


Attention explained - 1

I
A horizontal vector consisting of 8 gray squares.
am
A horizontal vector consisting of 8 gray squares.
trying
A horizontal vector consisting of 8 gray squares.
to
A horizontal vector consisting of 8 gray squares.
understand
A horizontal vector consisting of 8 gray squares.

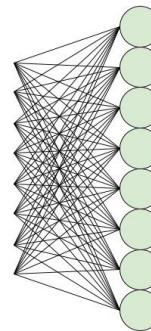


Attention explained - 1 (animated)



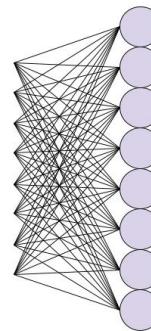
understand

K



understand

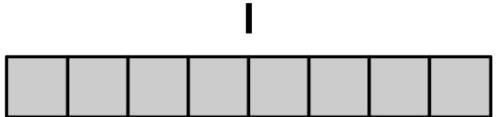
Q



understand

V

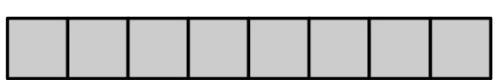
Attention explained - 1 (full)



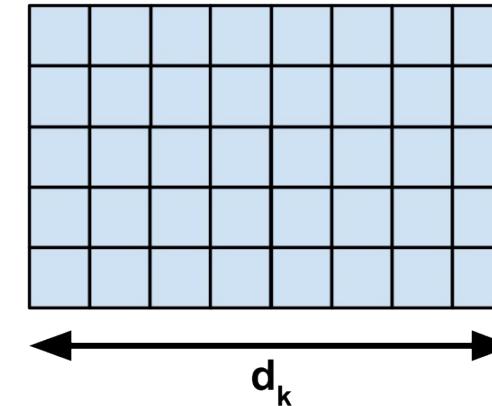
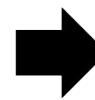
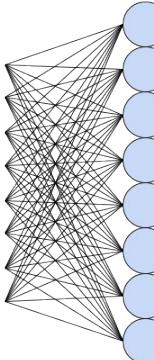
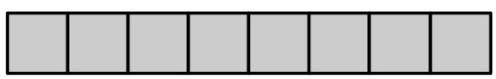
trying



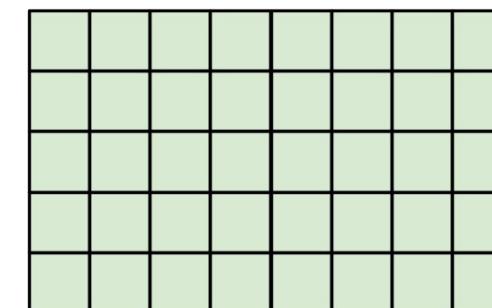
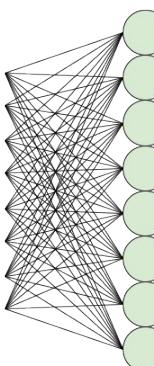
to



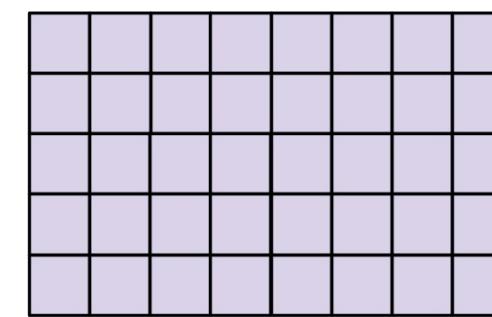
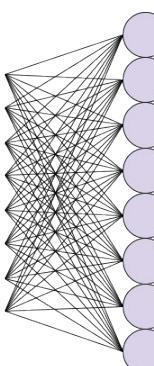
understand



$$K/\sqrt{d_k}$$



$$Q$$



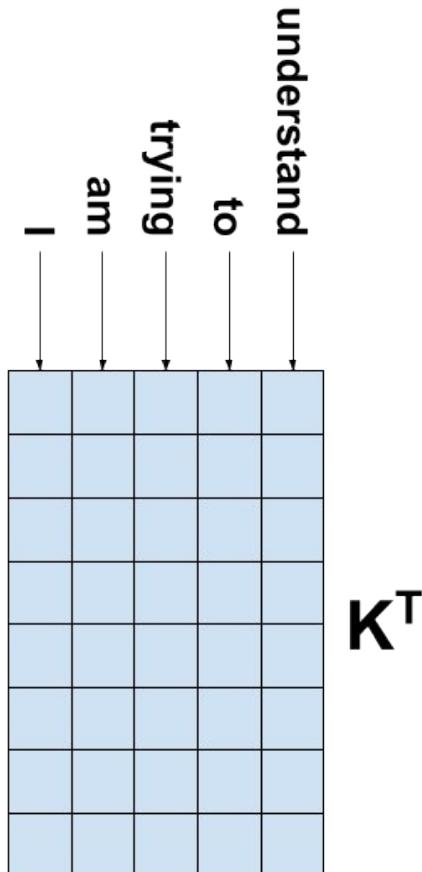
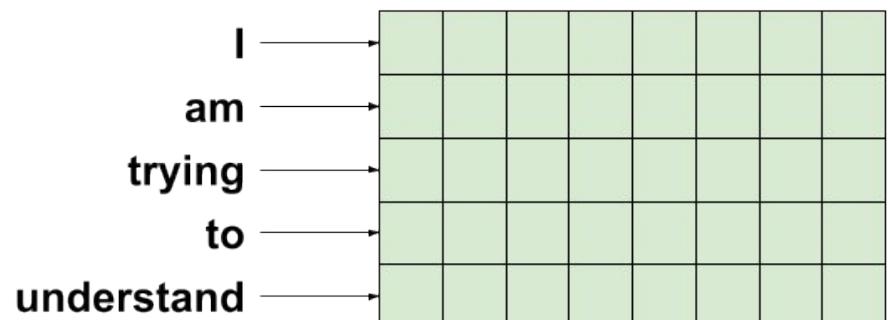
$$V$$

Attention explained - 2

Attention matrix scores

element $\in [-\infty, +\infty]$

Q

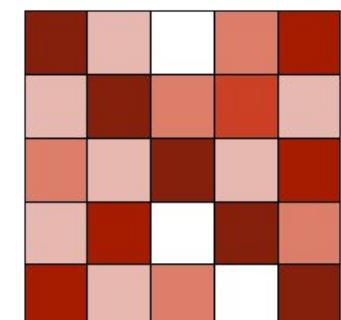


K^T

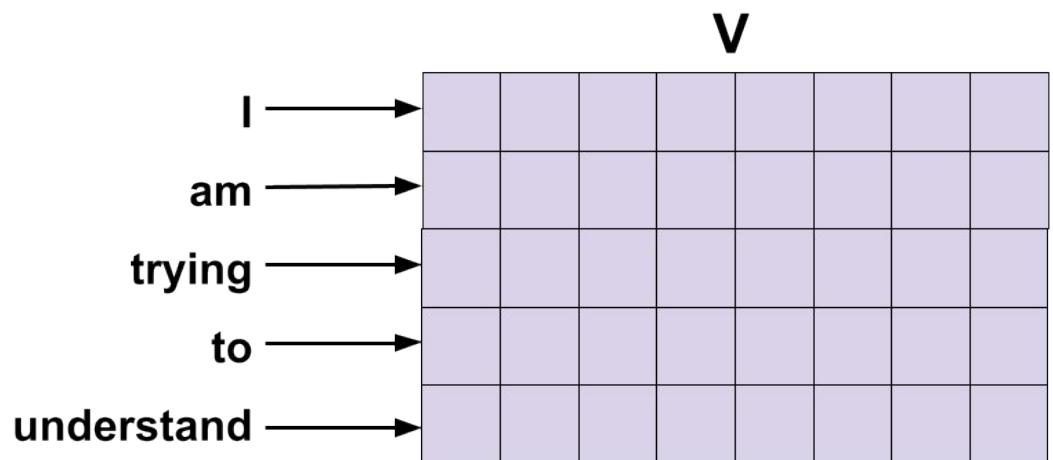
Attention matrix

element $\in [0,1]$
sum(line) = 1

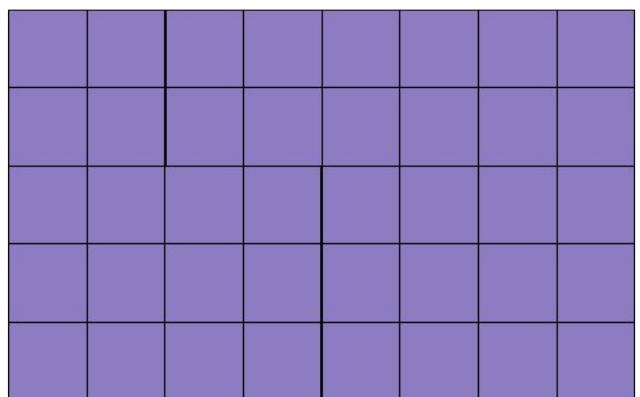
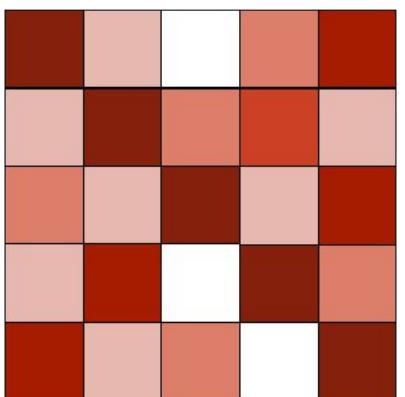
Softmax



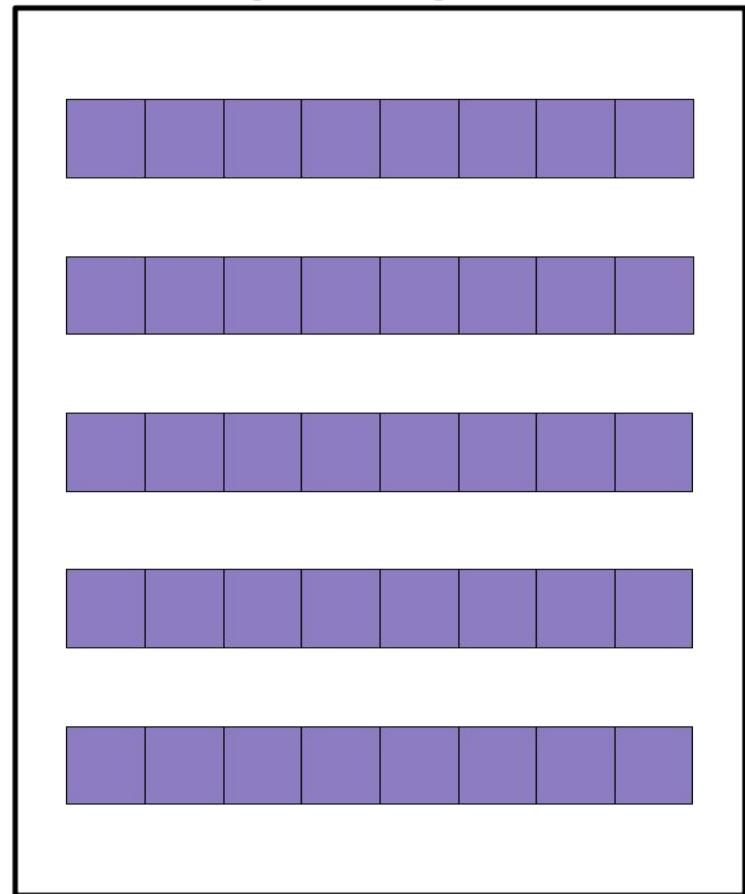
Attention explained - 3



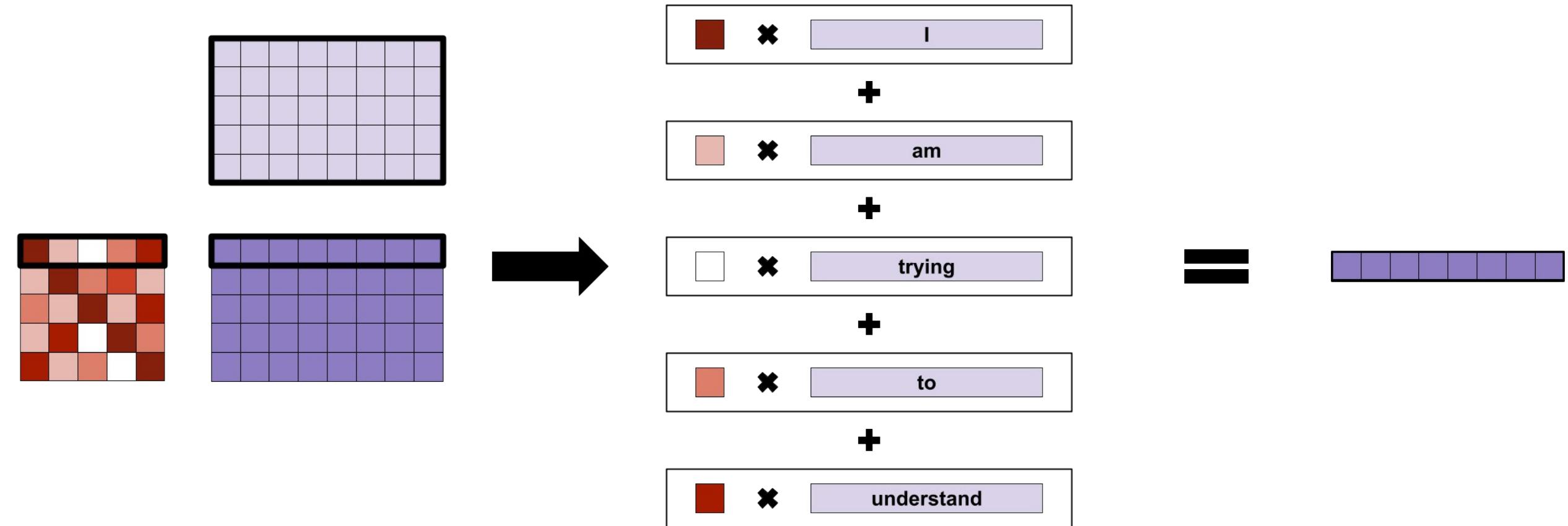
Attention matrix



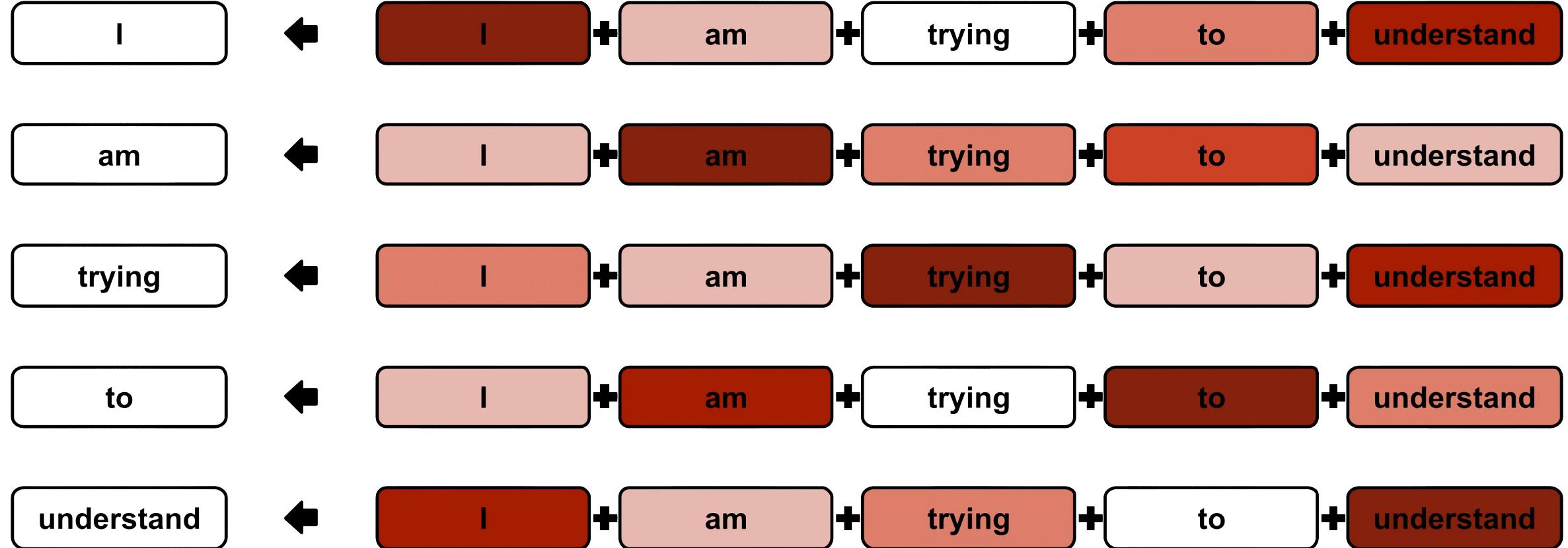
Output sequence



Attention explained - 3 (detailed)



Intuition behind the attention mechanism

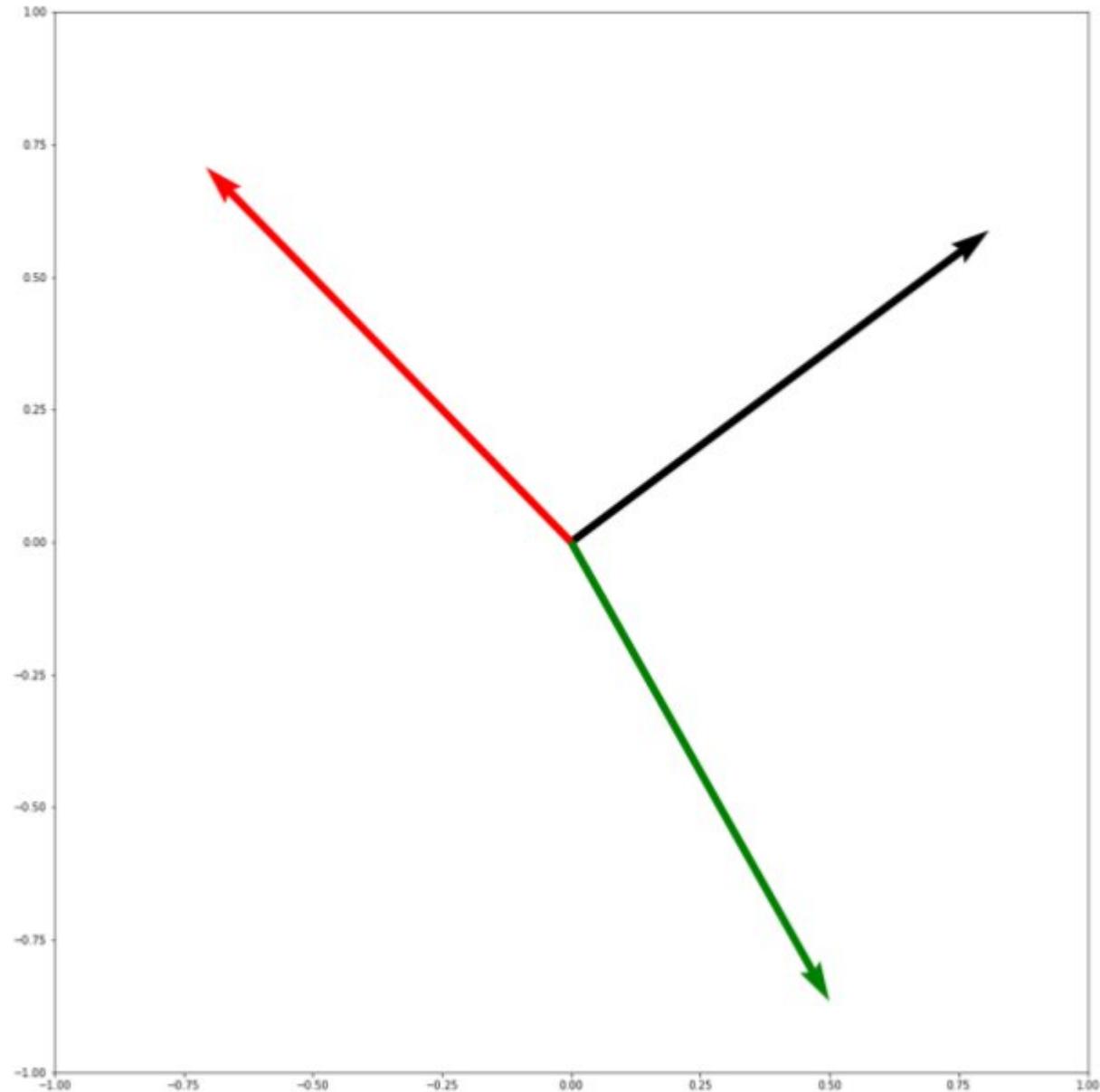


Intuition behind the attention mechanism - 0

The big dog



The : (0.50, -0.87)
big : (-0.70, 0.70)
dog : (0.81, 0.59)



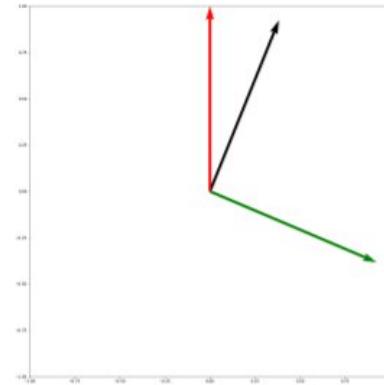
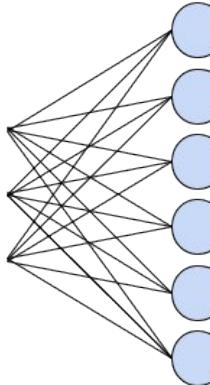
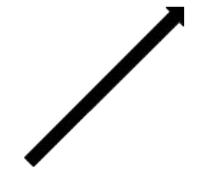
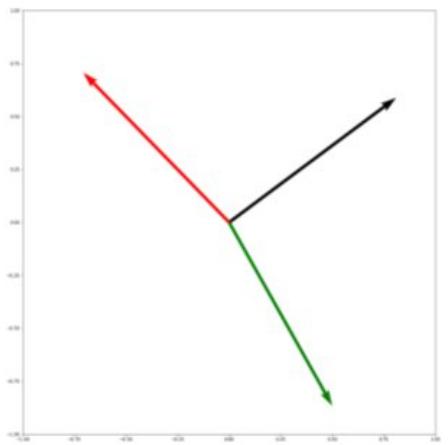
Intuition behind the attention mechanism - 1

0.50	-0.87
------	-------

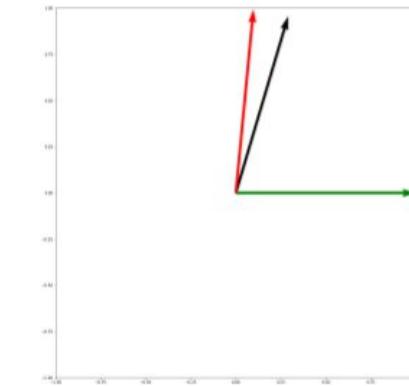
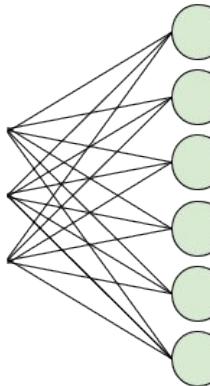
-0.70	0.70
-------	------

0.81	0.59
------	------

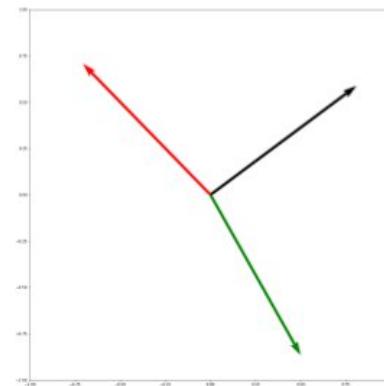
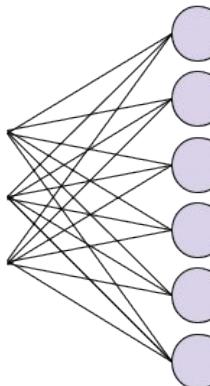
=



K

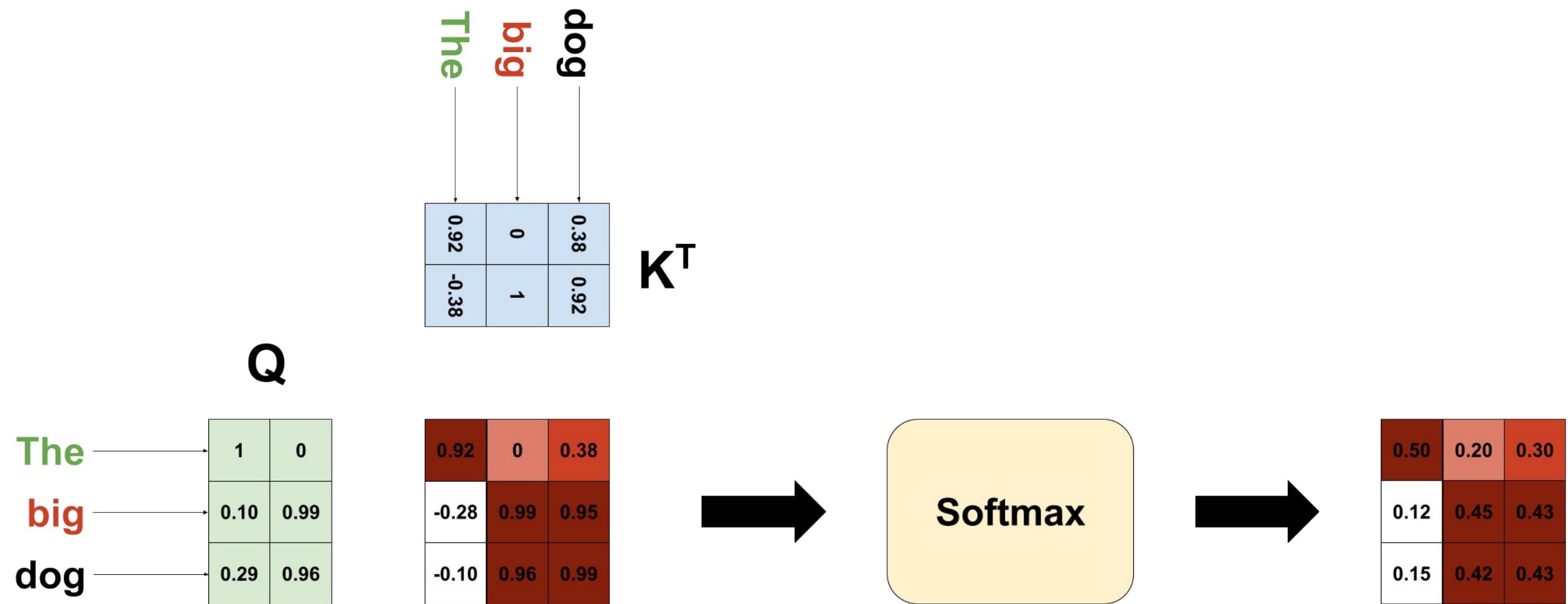


Q



V

Intuition behind the attention mechanism - 2



Intuition behind the attention mechanism - 3

0.50	0.20	0.30
0.12	0.45	0.43
0.15	0.42	0.43

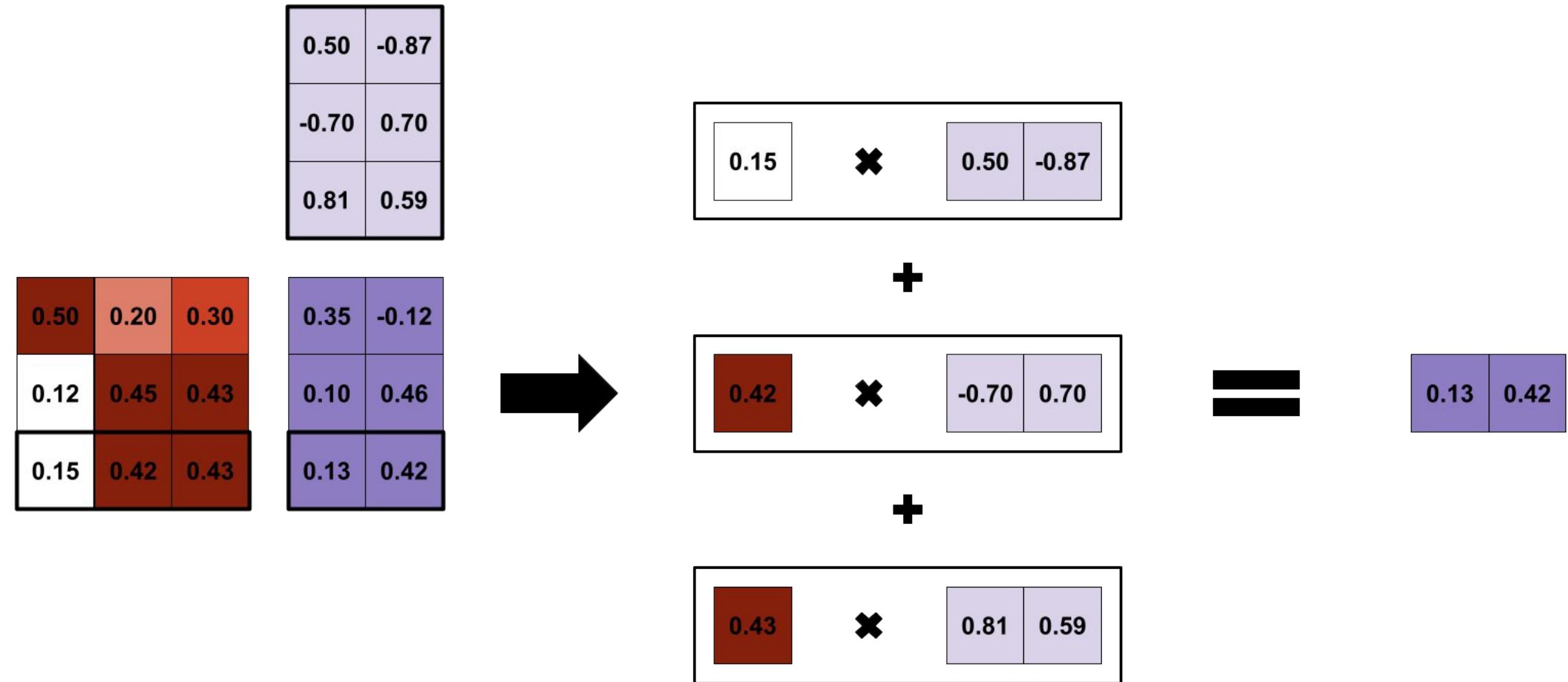


0.50	-0.87
-0.70	0.70
0.81	0.59

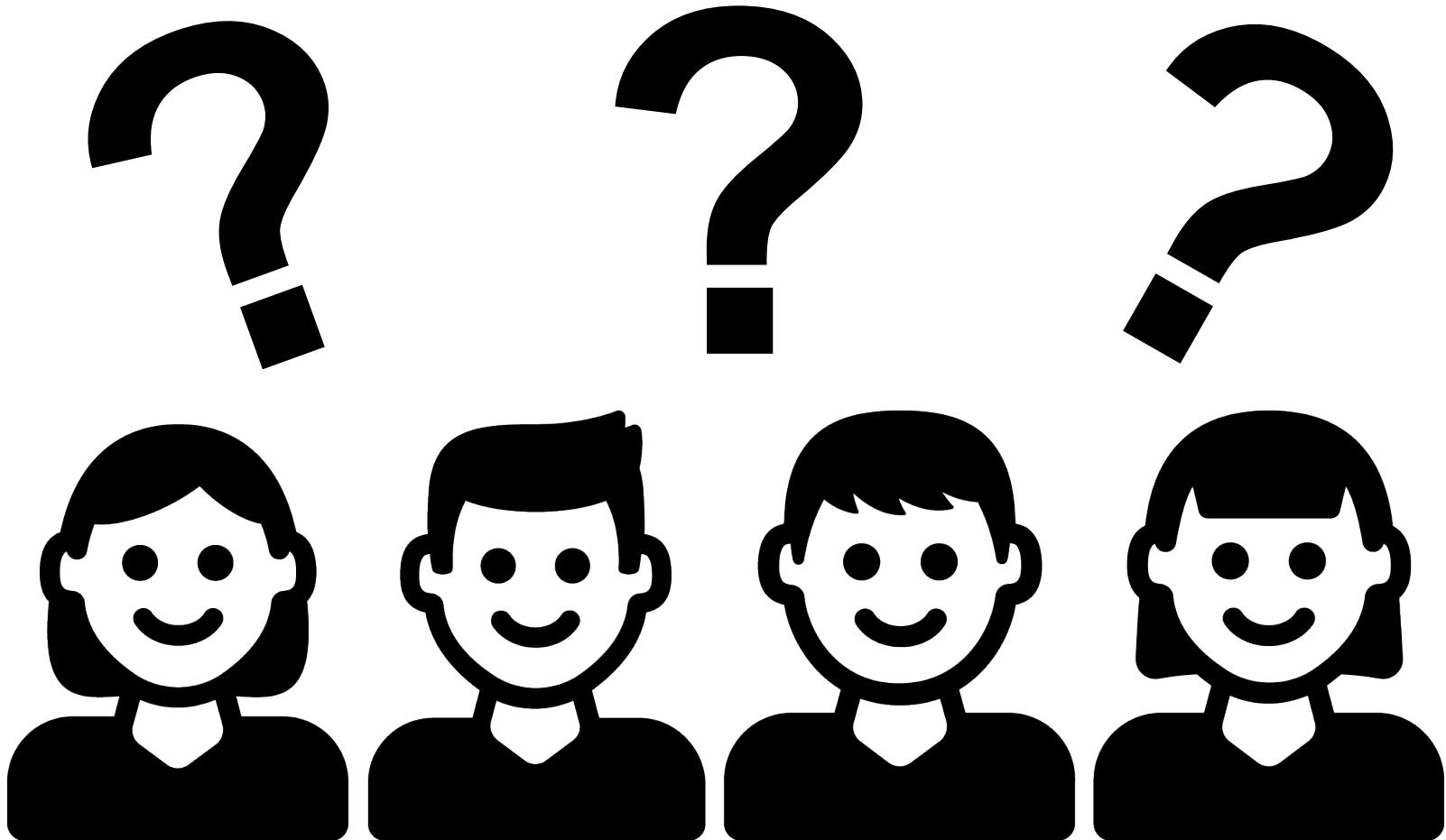


0.35	-0.12
0.10	0.46
0.13	0.42

Intuition behind the attention mechanism - 3 (detailed)



Question break #2



Multi-head Attention explained - 1

I



am



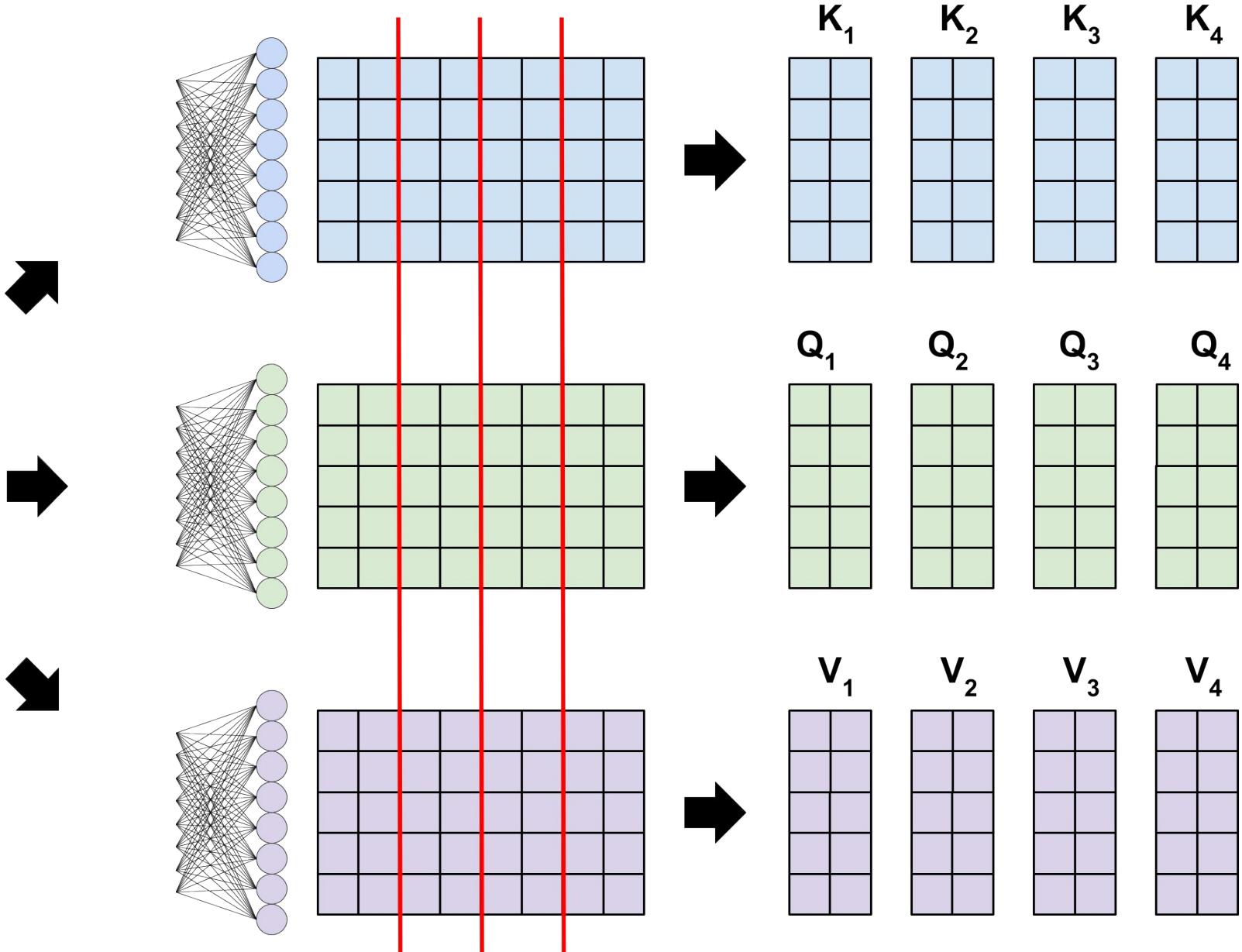
trying



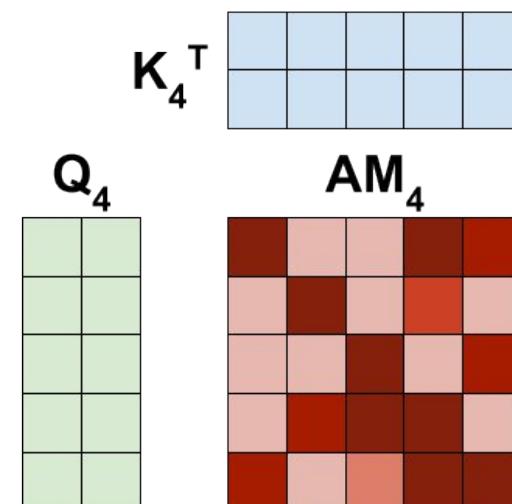
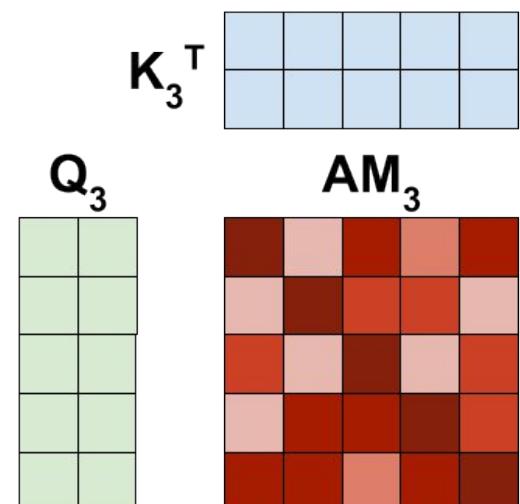
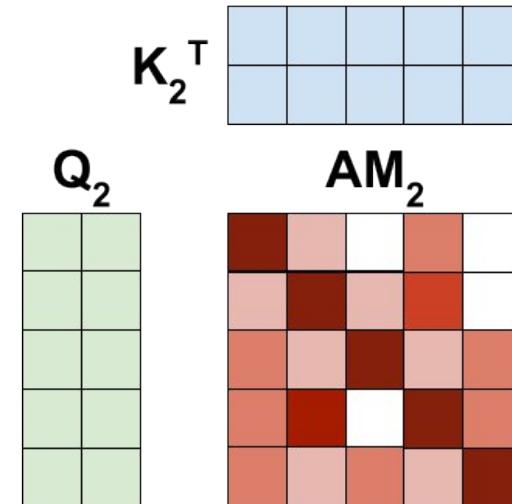
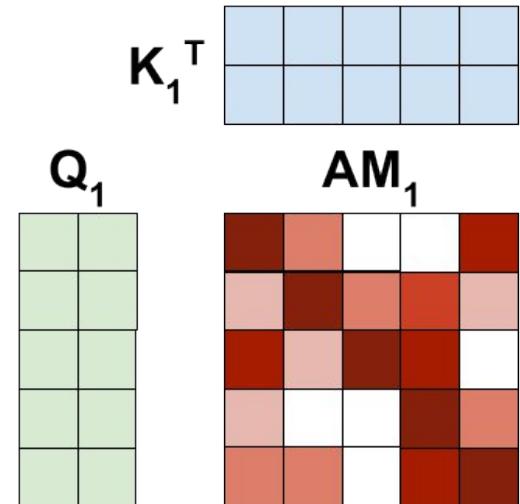
to



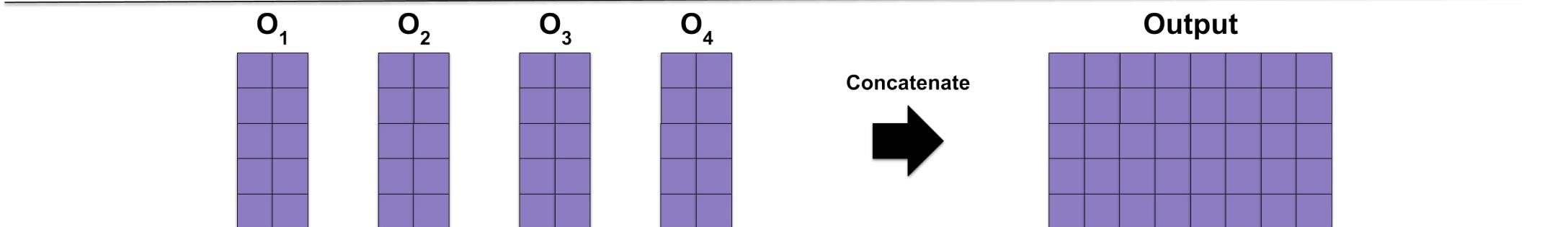
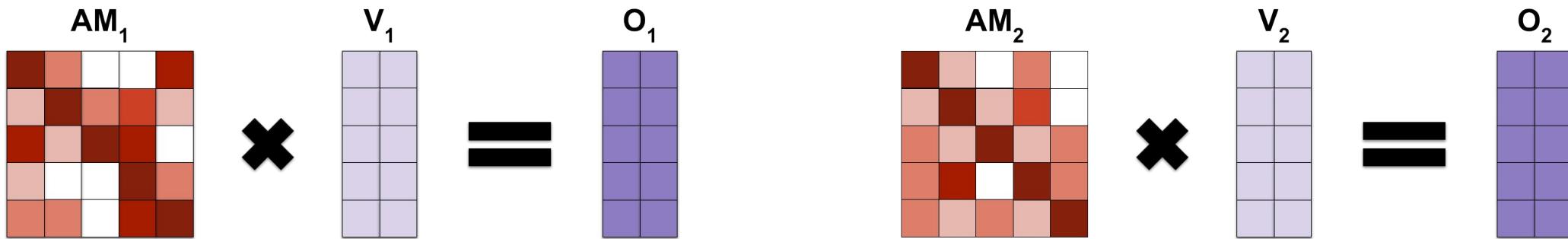
understand



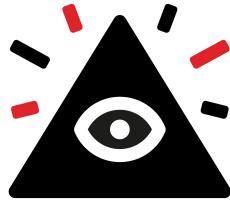
Multi-head Attention explained - 2



Multi-head Attention explained - 3



9



«Attention is
All You Need»

Transformers



1

What is a **Transformer**? The magic
of the **Attention Mechanism**

2

The different **Transformers**
architectures

3

**Pre-training and Foundation
Models**

4

**Specialization of Foundation
Models** (especially LLM)

5

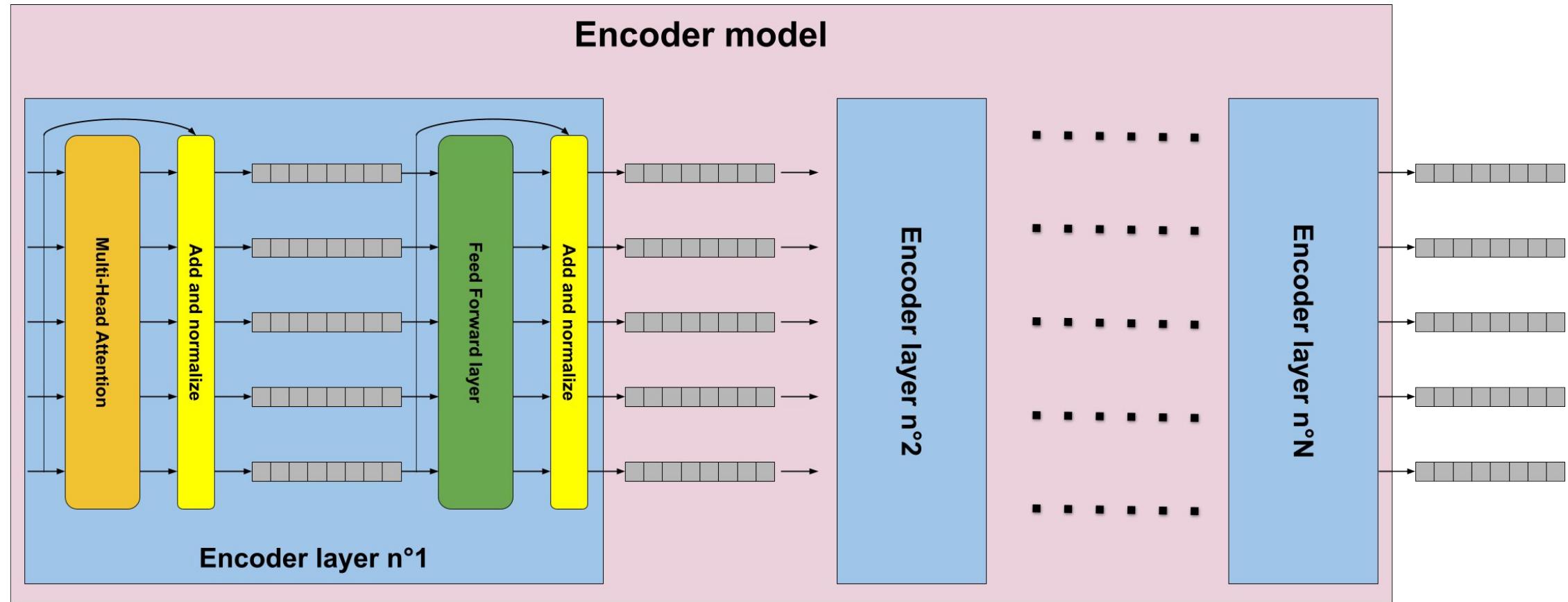
Example: IMDB Reviews



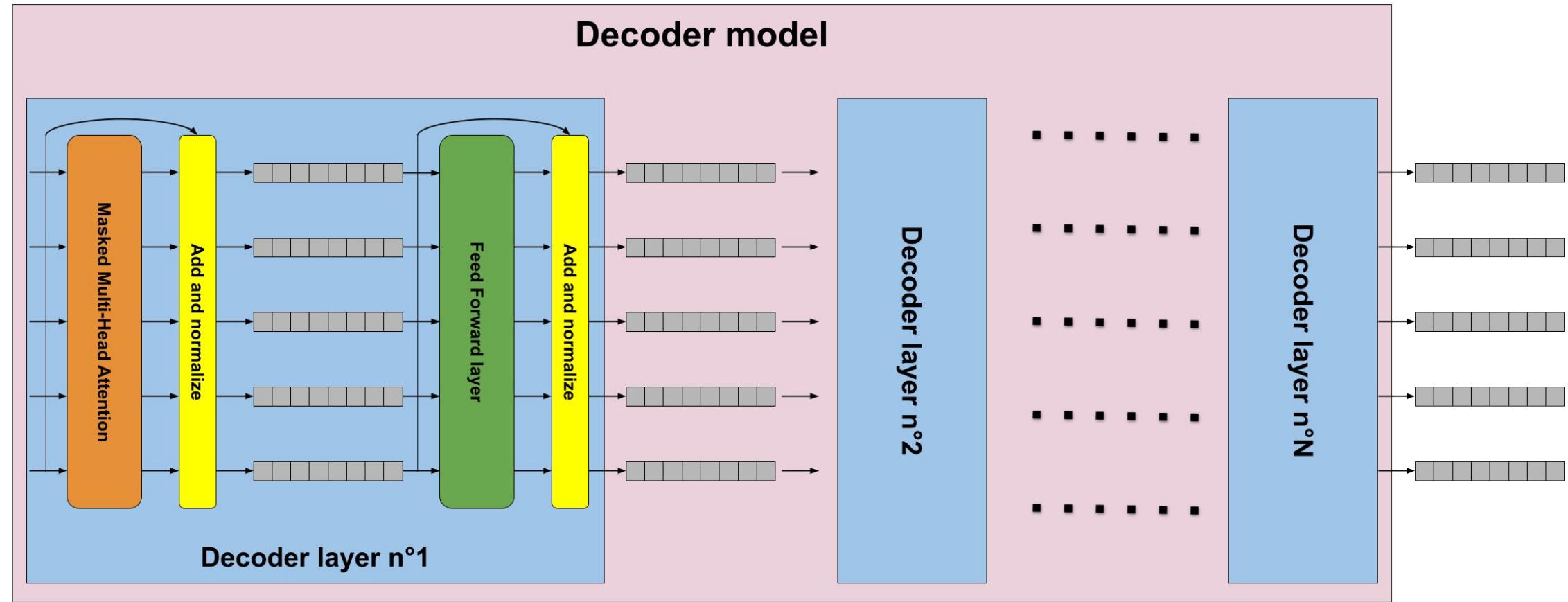
Encoder-Decoder



BERT / Encoder / Auto-encoding



GPT / Decoder / Auto-regressive



Bidirectional vs Unidirectional attention

Bidirectional attention for Encoder

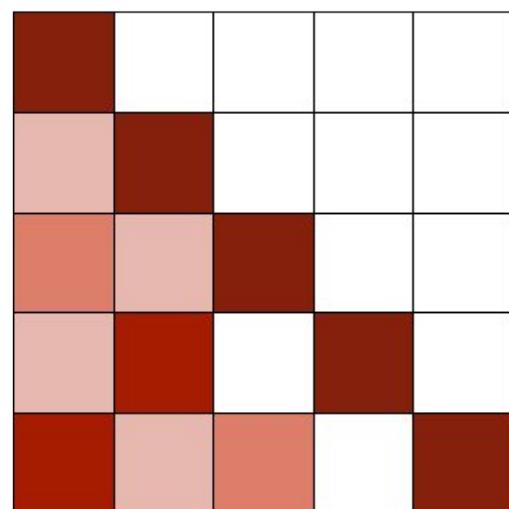
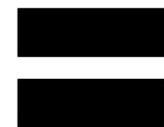
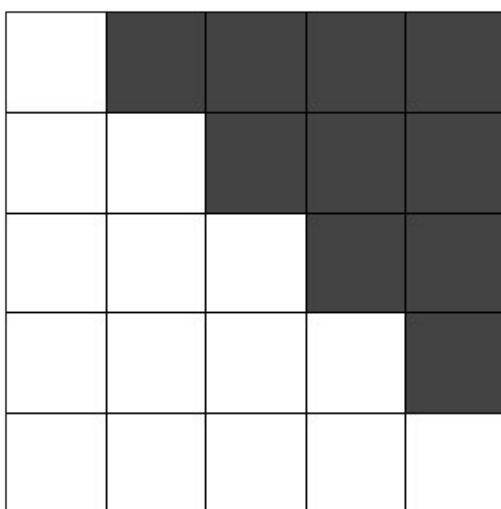
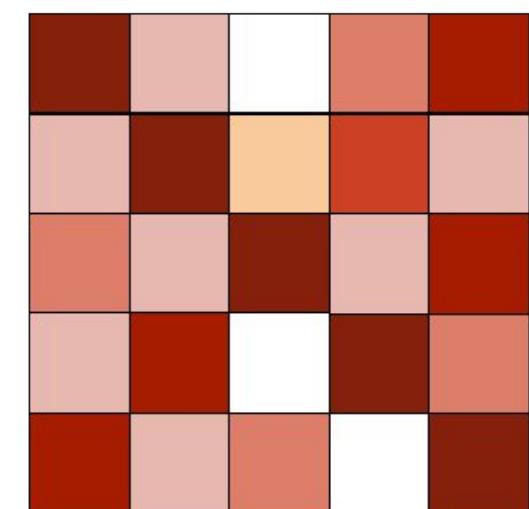


Unidirectional attention for Decoder

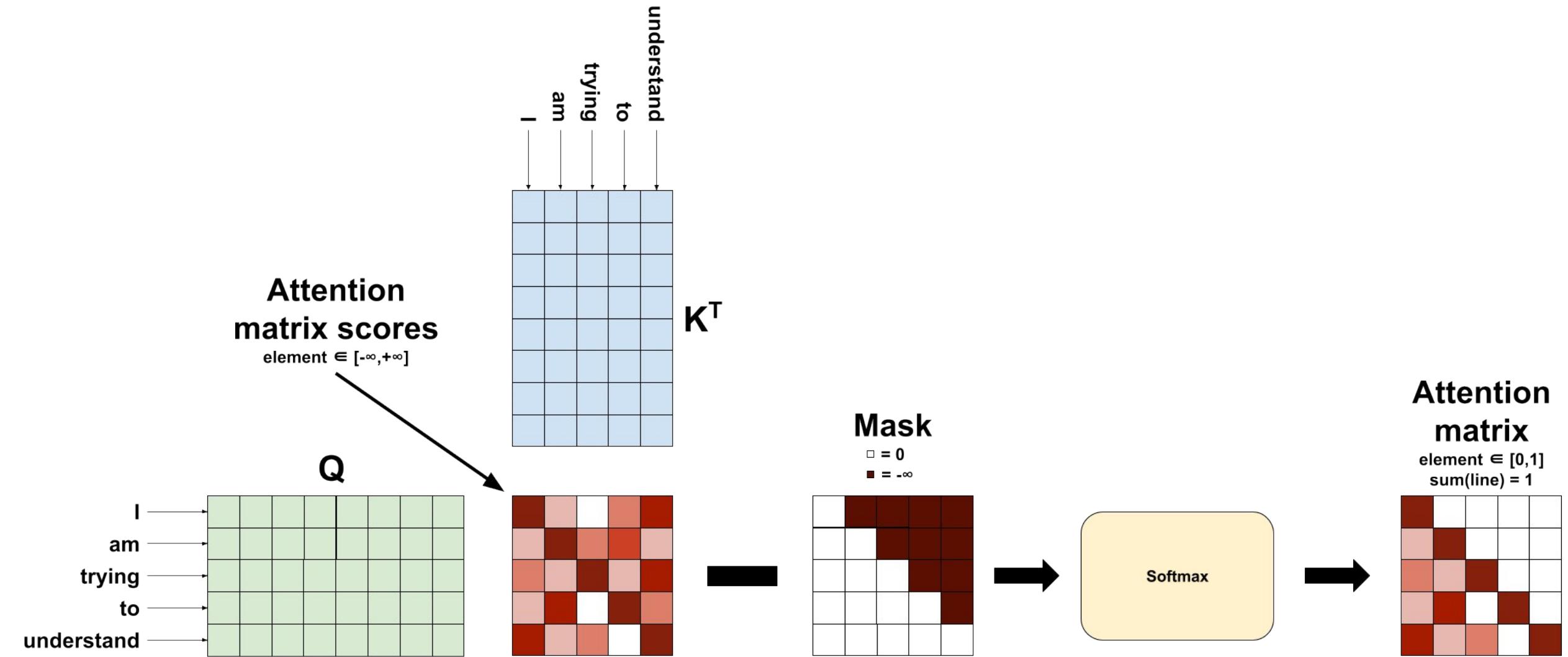


VS

Mask attention

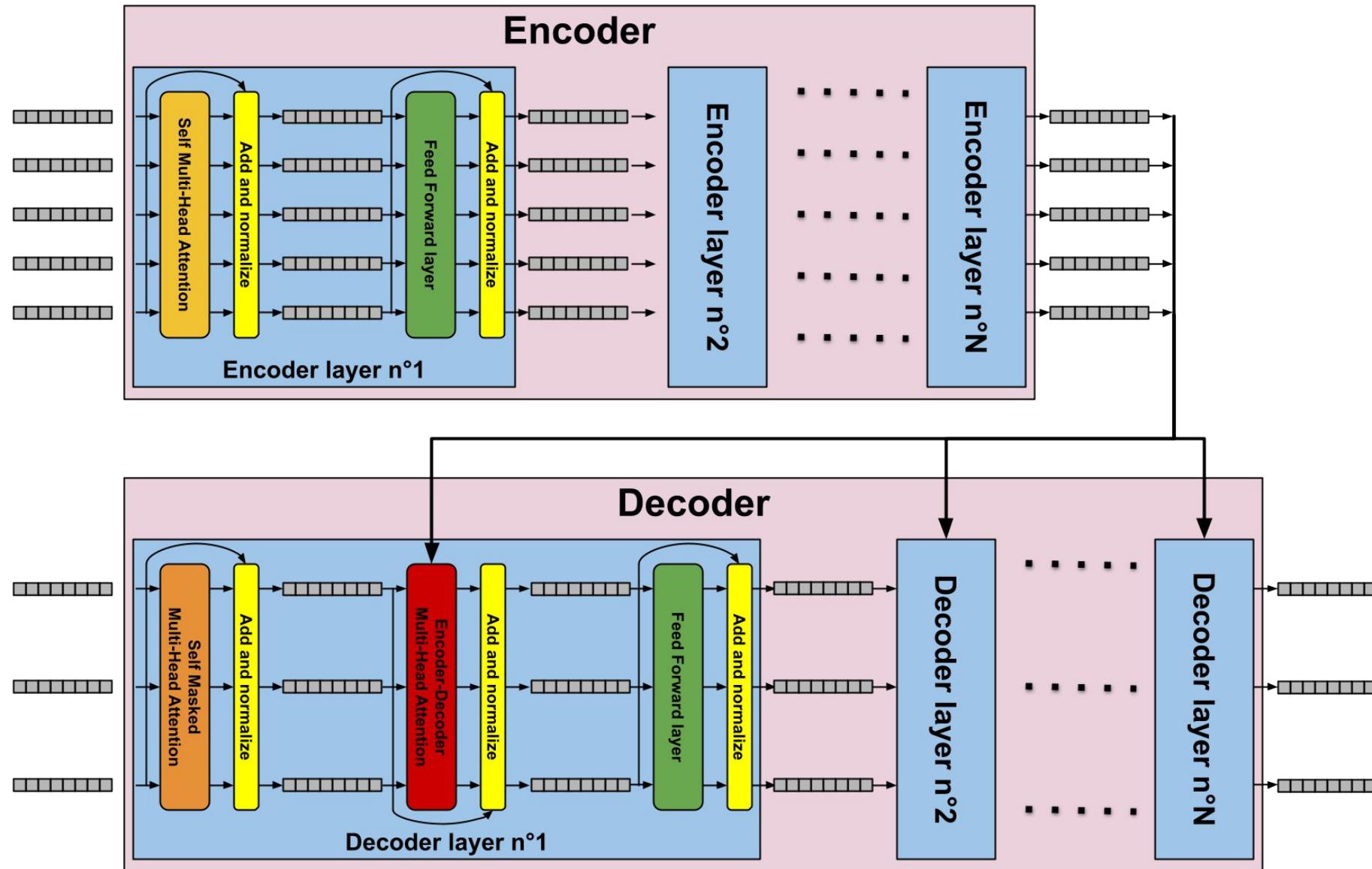


Unidirectional attention detailed

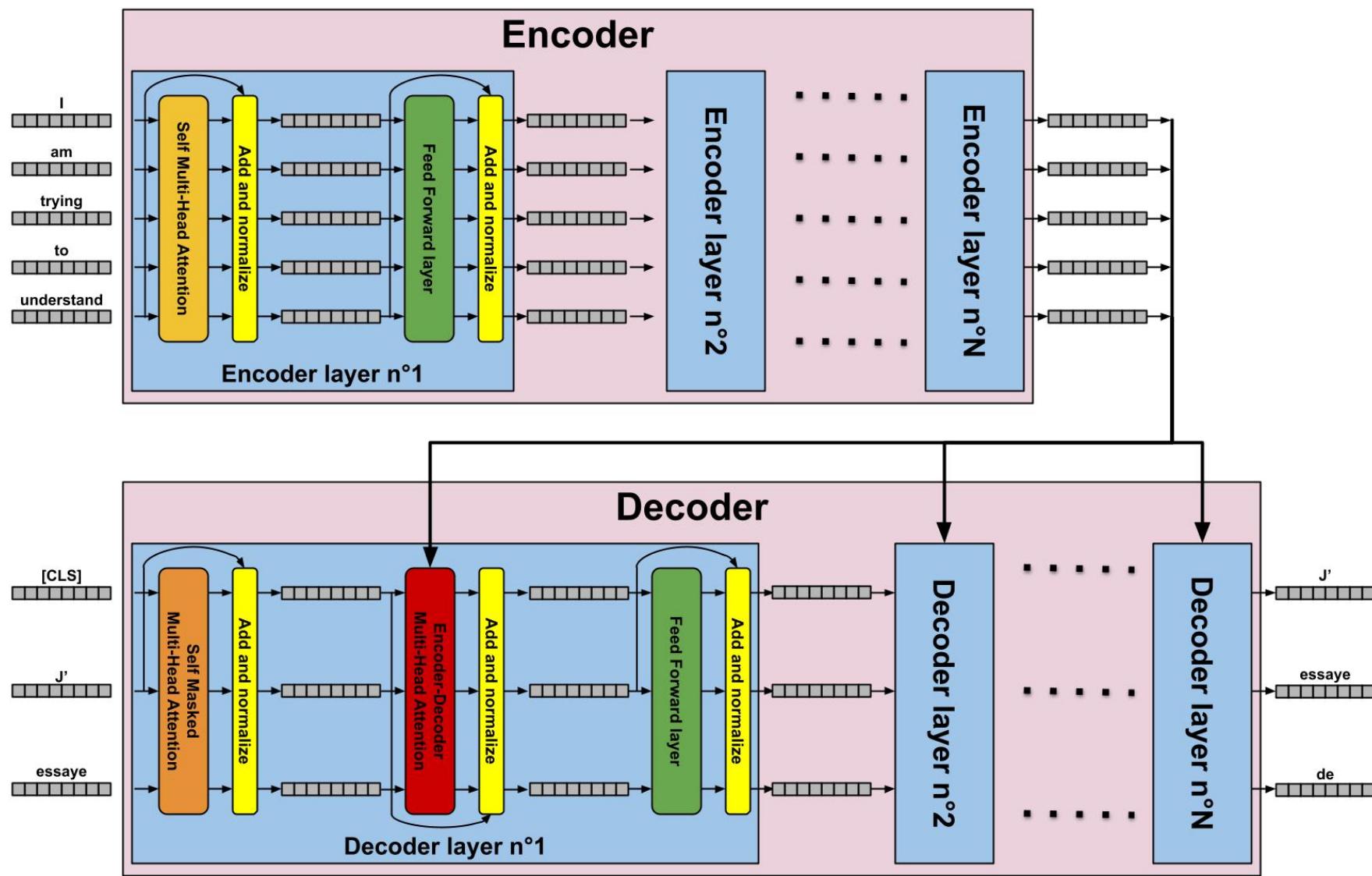


Encoder-Decoder architecture

T5 / Encoder-Decoder / Sequences-to-Sequences

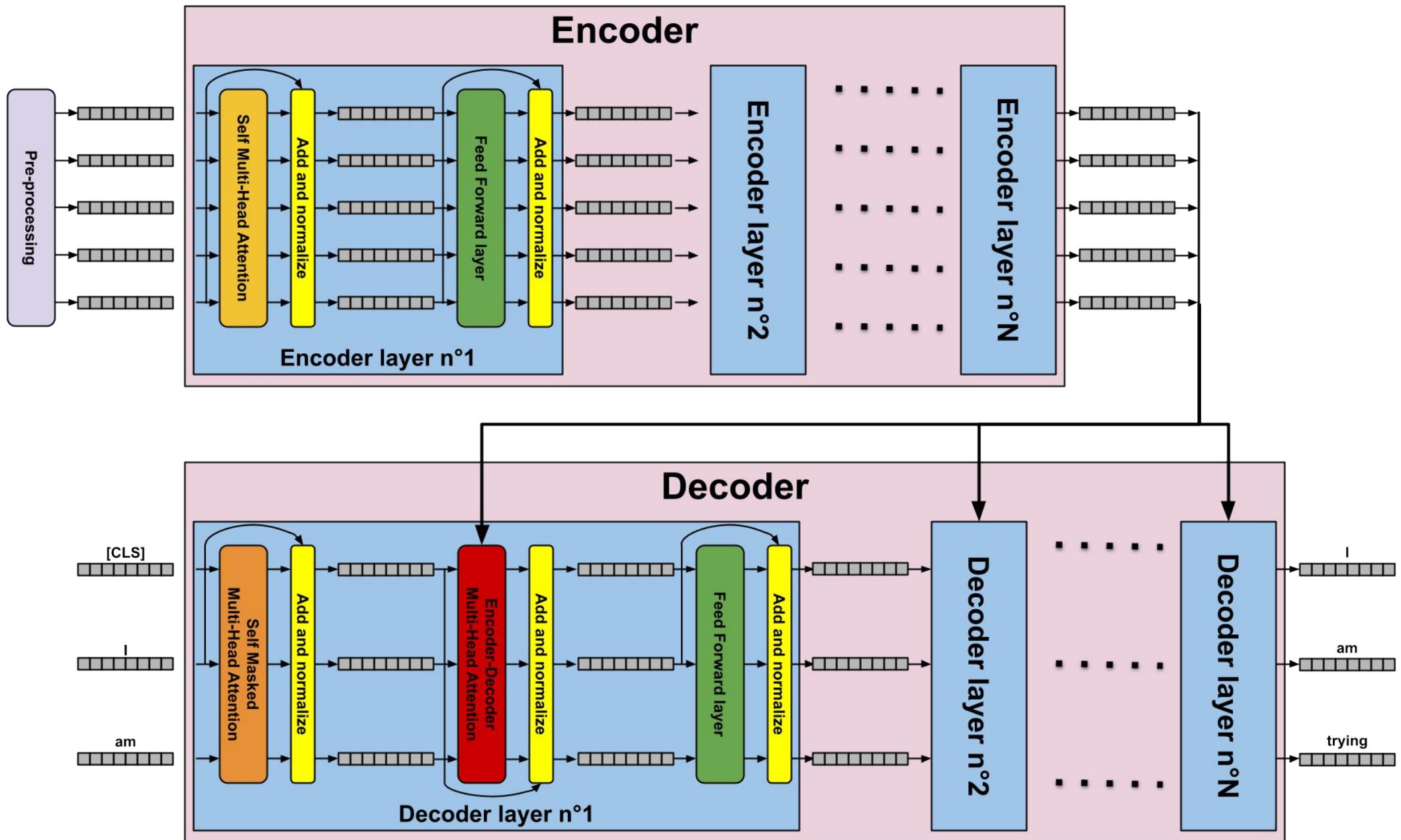


Encoder-Decoder translation example

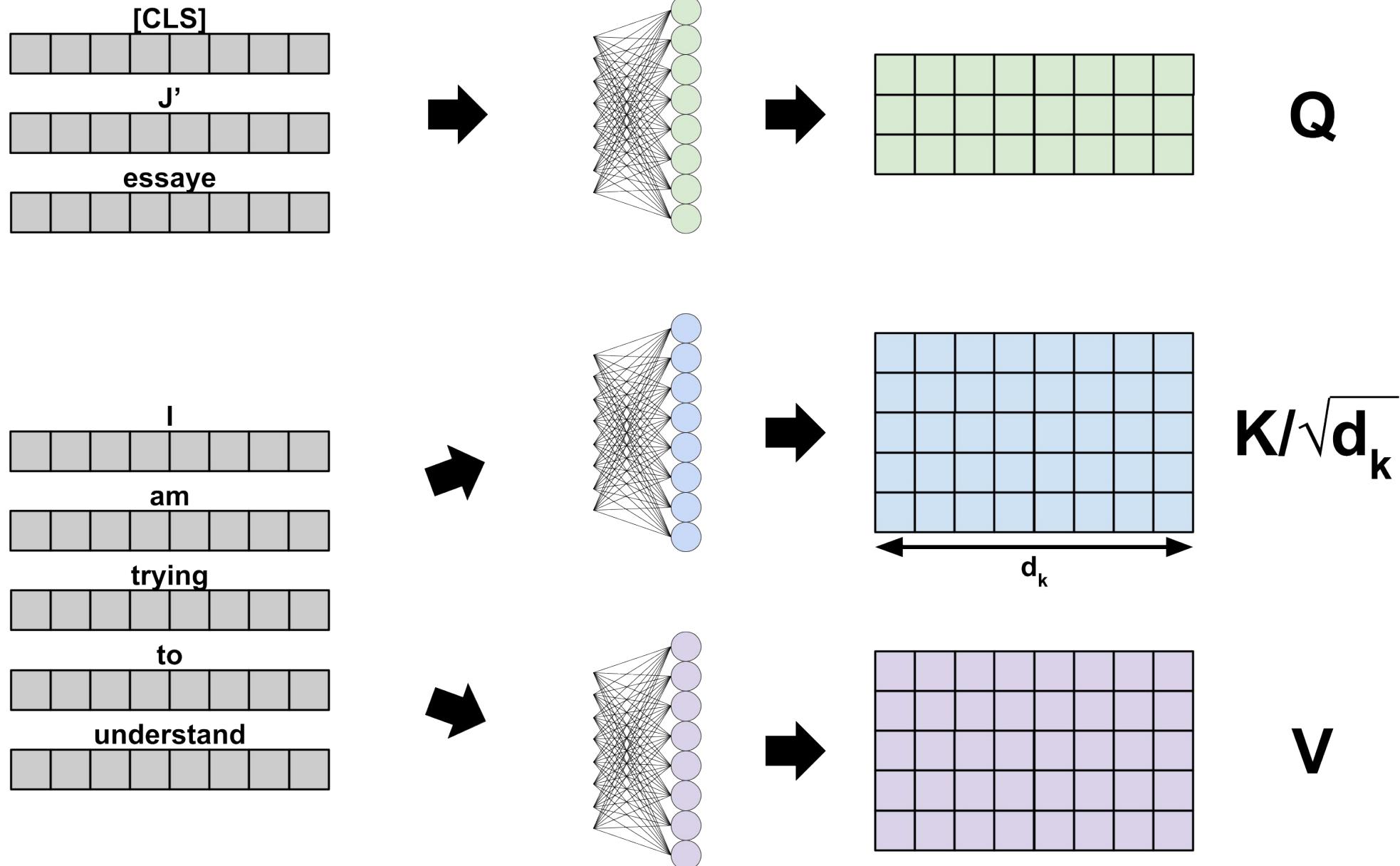


Encoder-Decoder whisper example

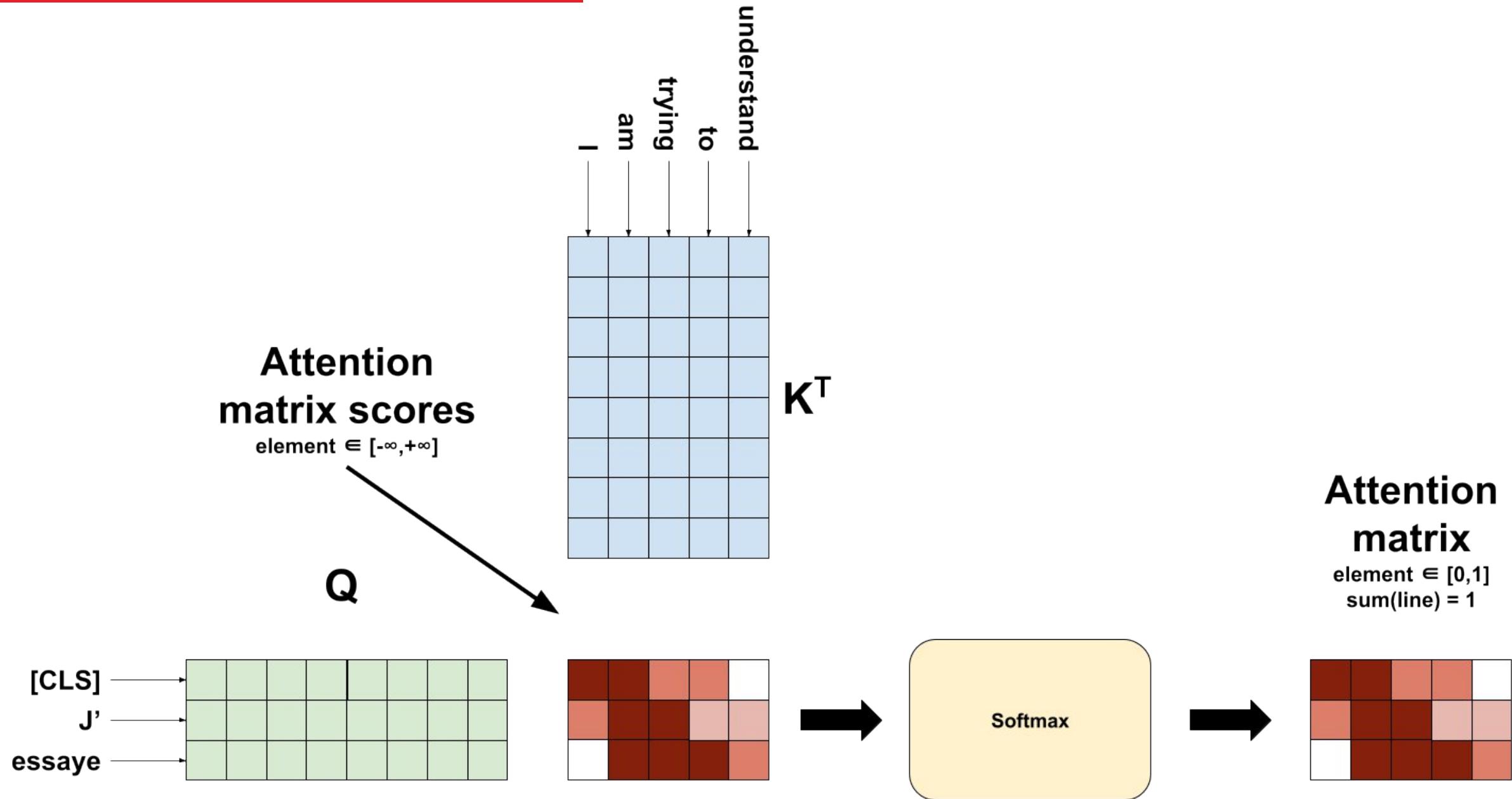
听得见



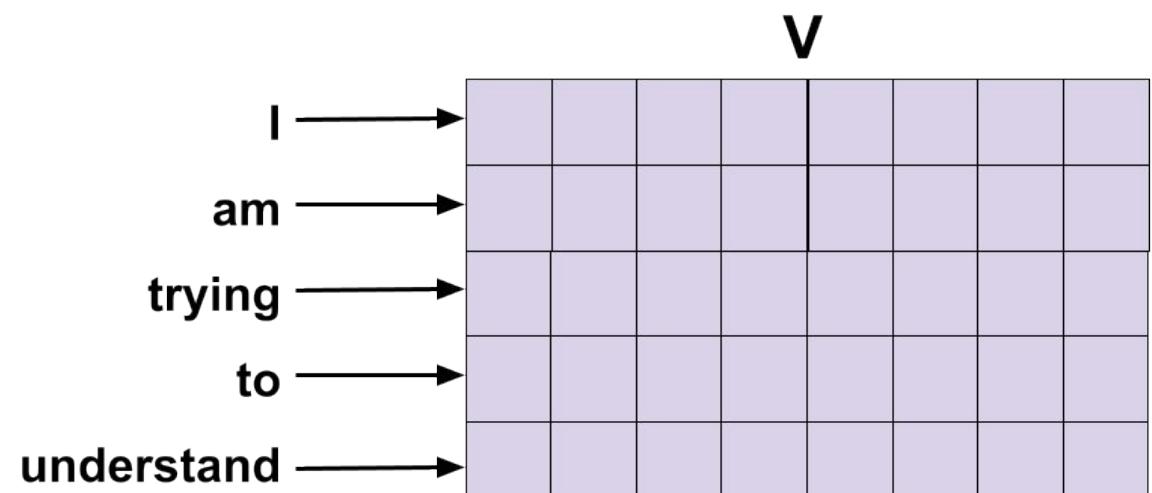
Encoder-Decoder attention - 1



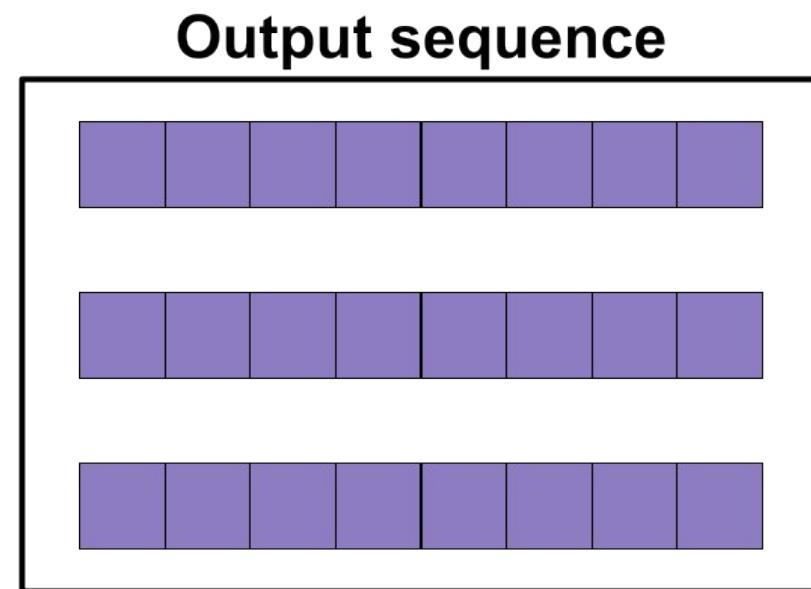
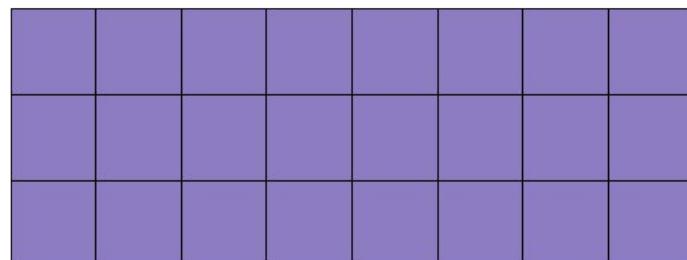
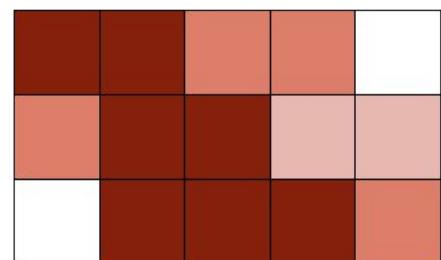
Encoder-Decoder attention - 2



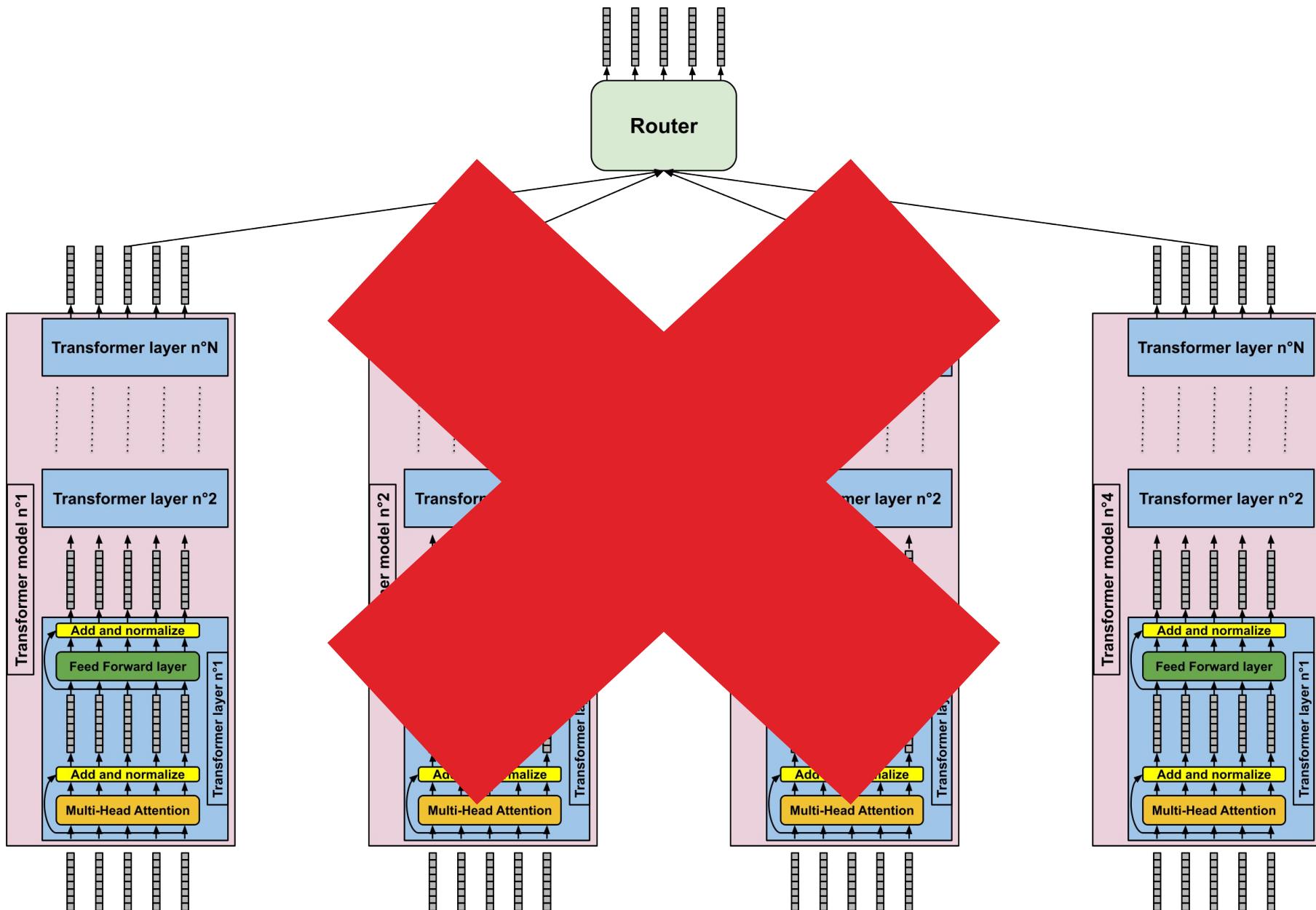
Encoder-Decoder attention - 3



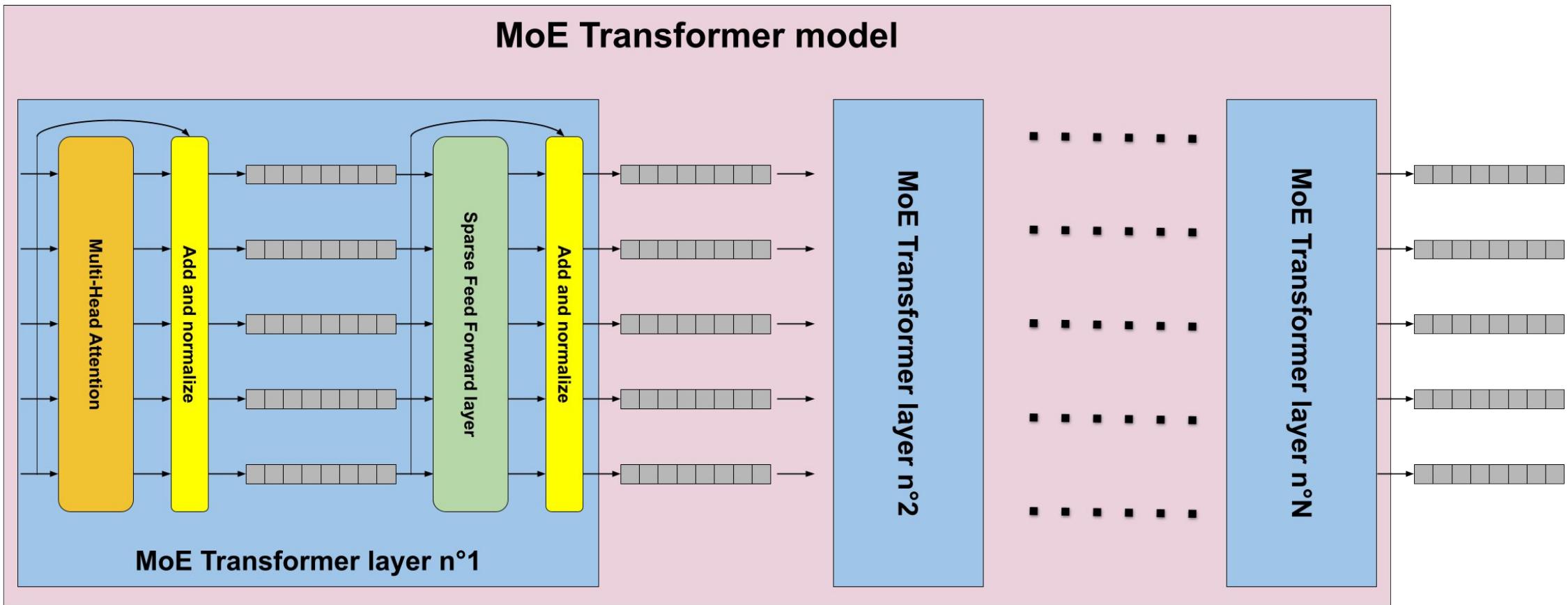
Attention matrix



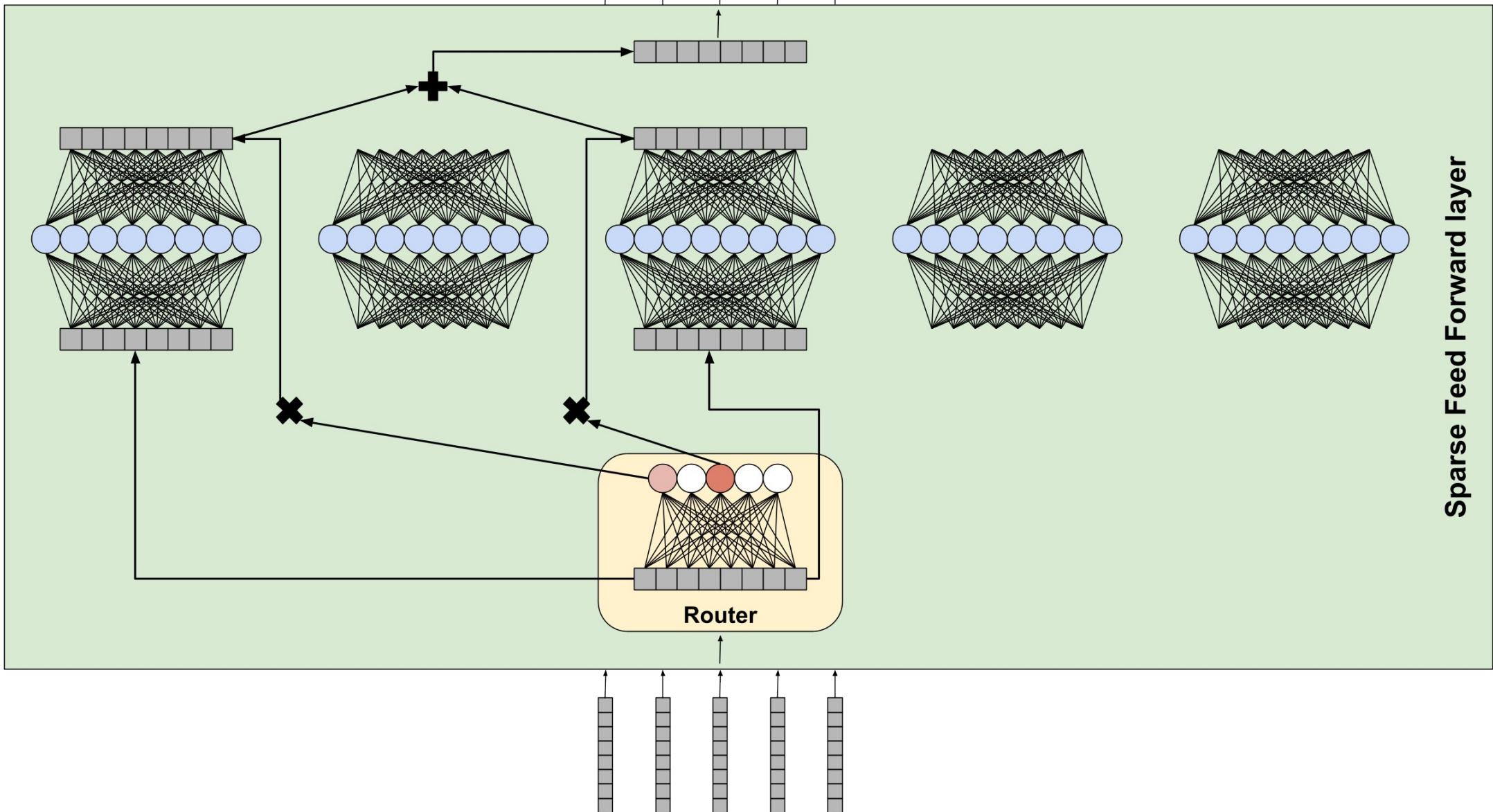
Misunderstanding about Mixture of Experts



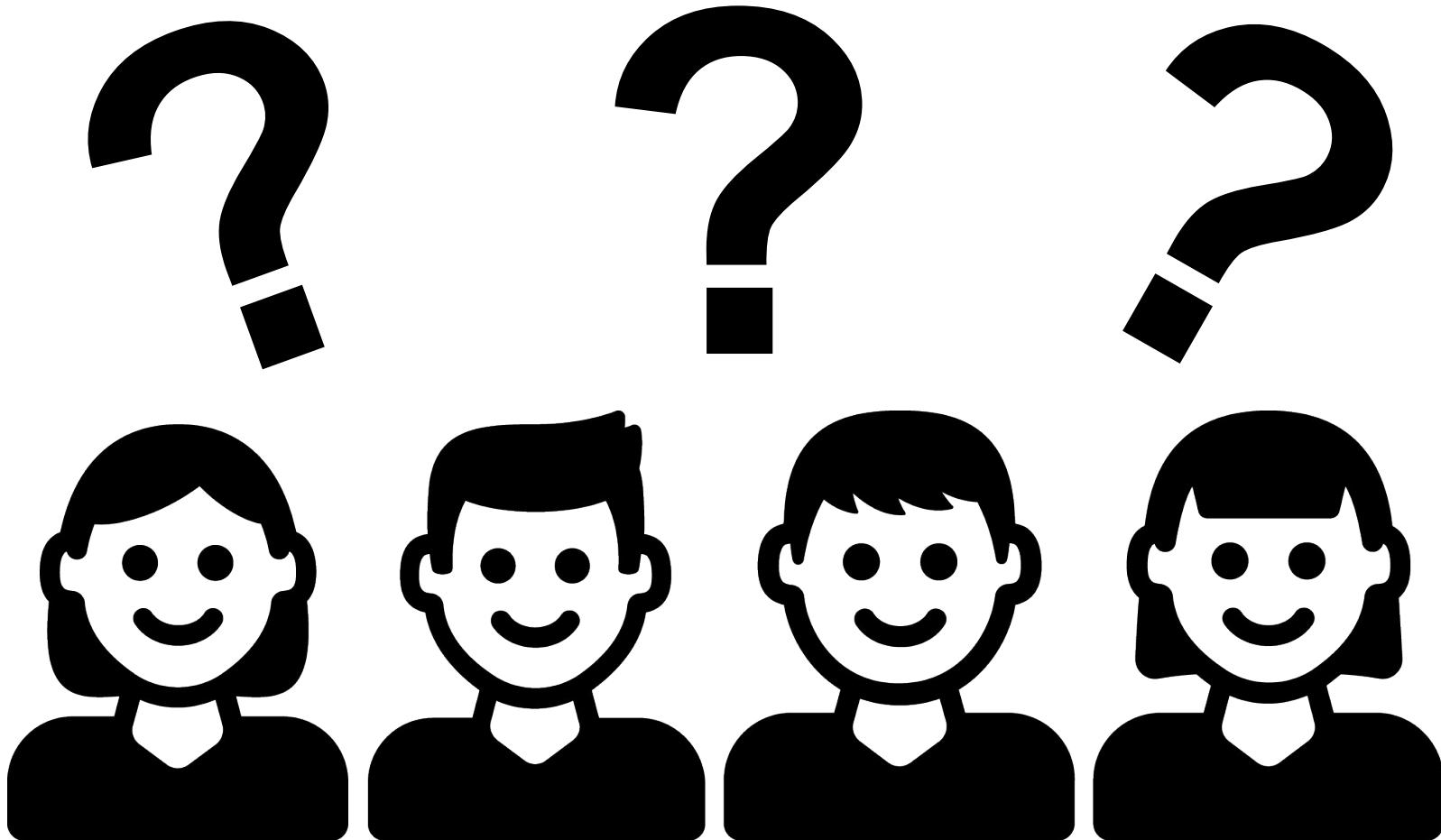
Mixture of Expert Transformer Architecture



Sparse Feed Forward Layer



Question break #3



9



«Attention is
All You Need»

Transformers



1

What is a **Transformer**? The magic
of the **Attention Mechanism**

2

The different **Transformers**
architectures

3

**Pre-training and Foundation
Models**

4

**Specialization of Foundation
Models** (especially LLM)

5

Example: IMDB Reviews



Training a language model

I want to make a chatbot to help people better understand the civil code.



I want to make a bot which can filter out respectful comments for a real reddit experience.



Training a language model

What do you need to train a large language model ?

A truckload of data

Transformers are ravenous. You need to feed them with a substantial and hard-to-come-by dataset.



A mighty compute infrastructure

GPUs with high throughput and large VRAM can execute this training in a reasonable amount of time.



A copious amount of electricity

Storage, memory and compute power consume a lot of electricity.

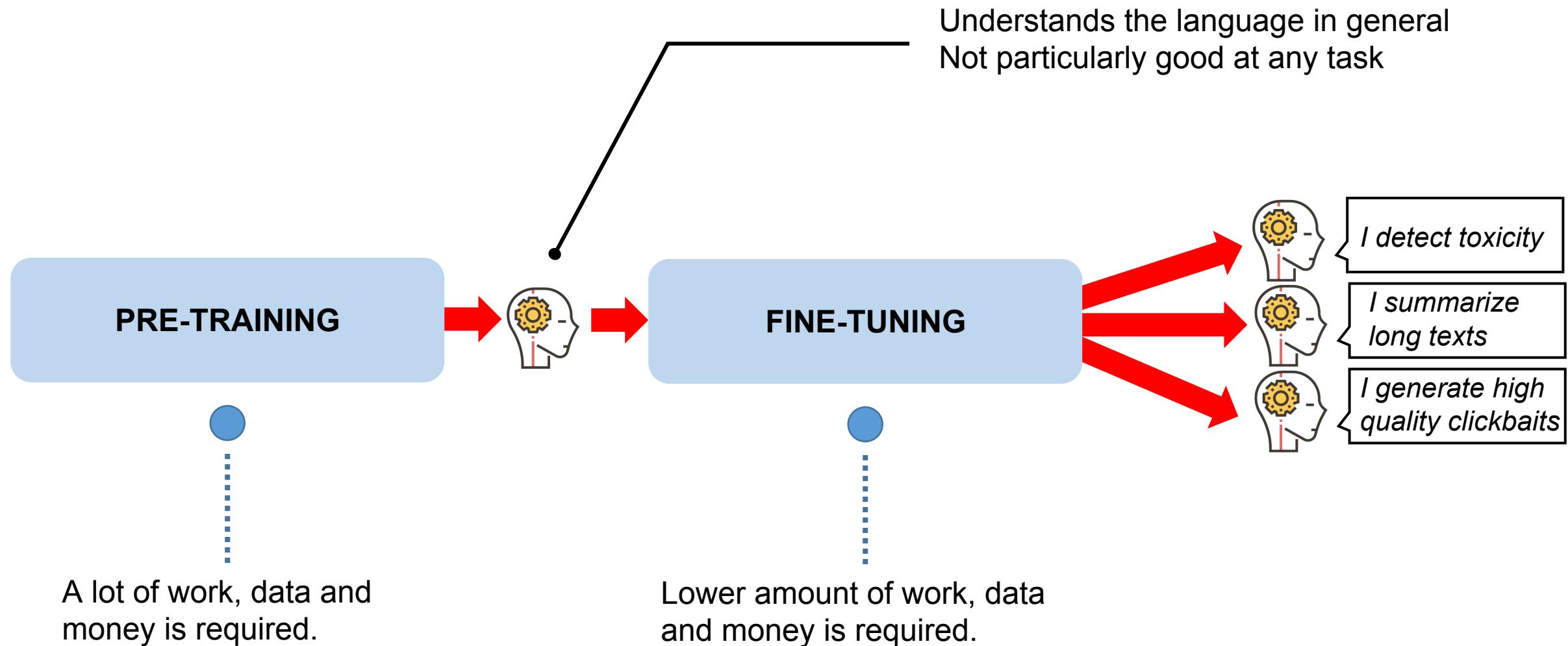


An abundance of manpower

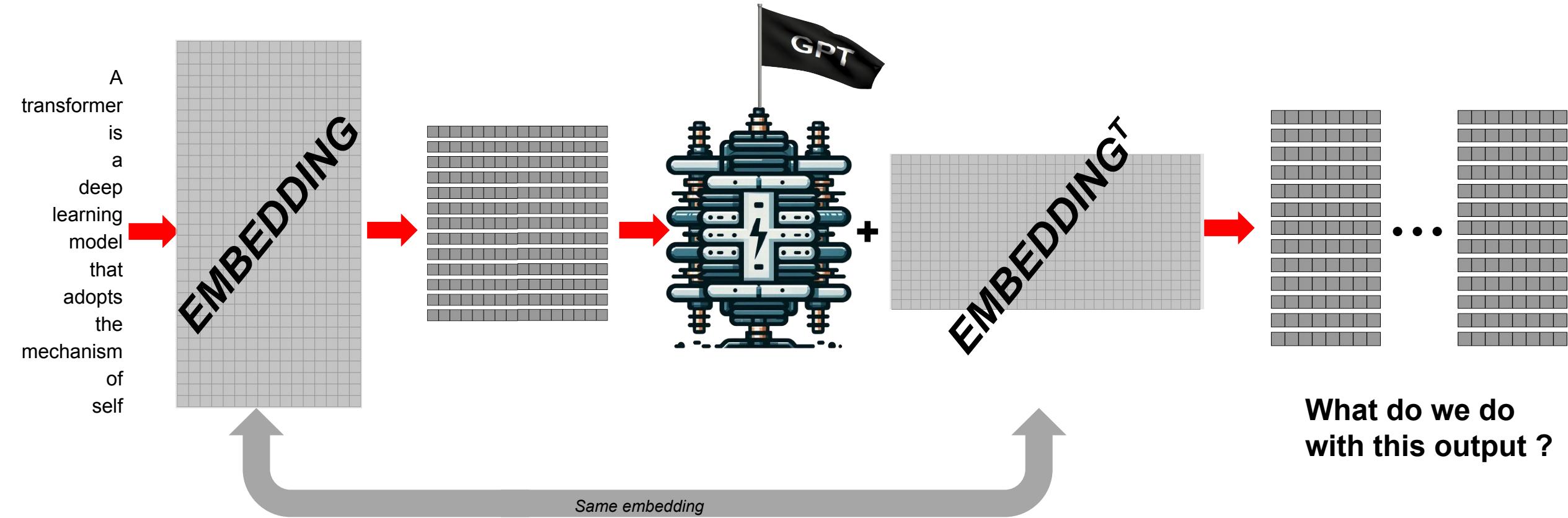
Cleaning the dataset, making experiments, monitoring SOTA advancements is a lot of work.



Training a language model

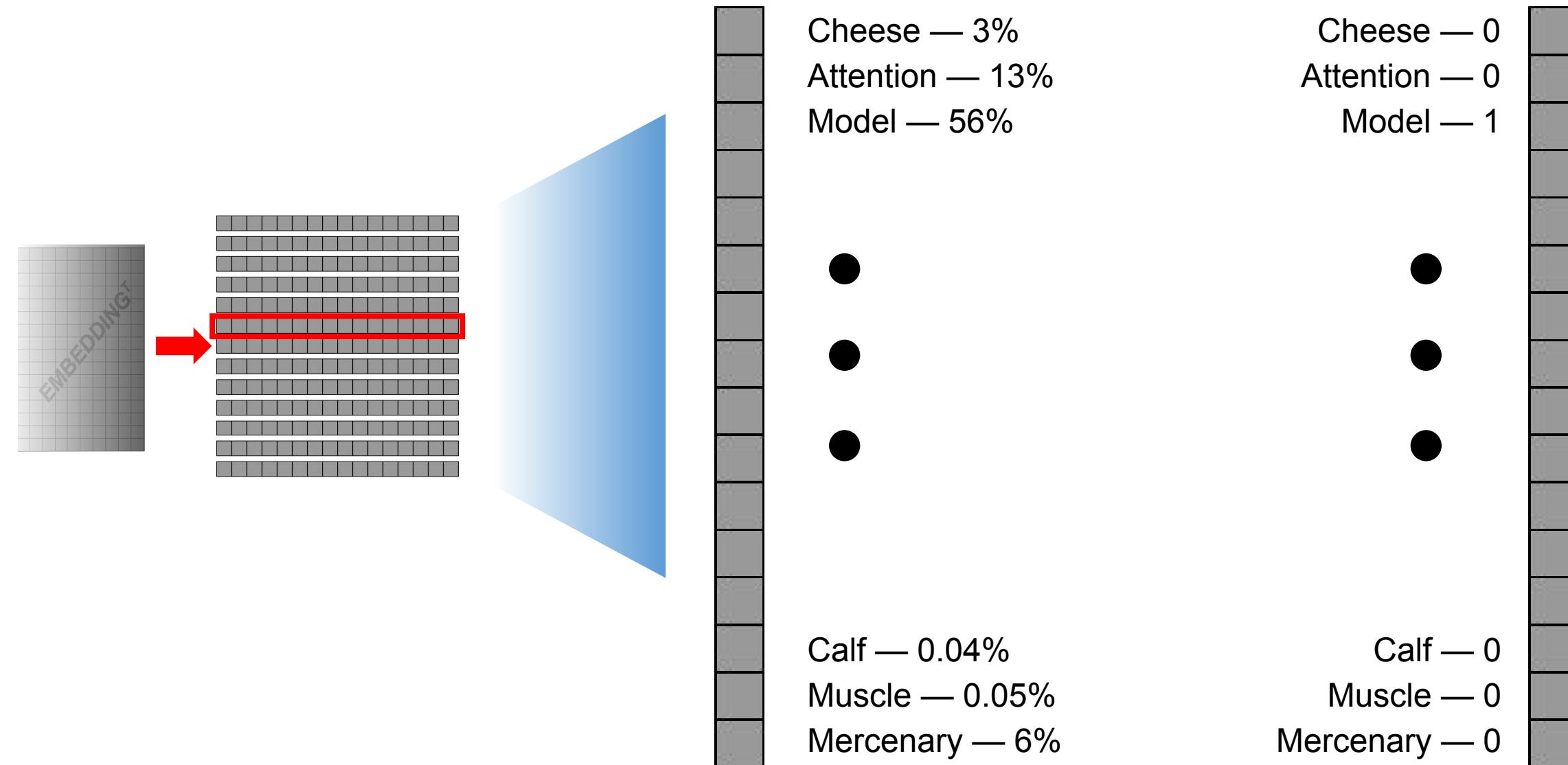


Pretraining a GPT-style transformer



Pretreaining a GPT-style transformer

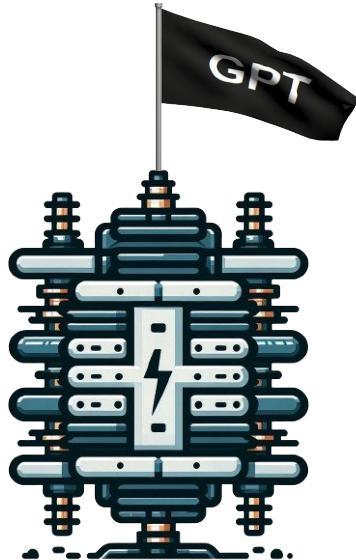
Classification target



Pretraining a GPT-style transformer

Input

A
transformer
is
a
deep
learning
model
that
adopts
the
mechanism
of
self



Target

transformer
is
a
deep
learning
model
that
adopts
the
mechanism
of
self
attention



Next word prediction

Pretraining a BERT-style transformer

Sample

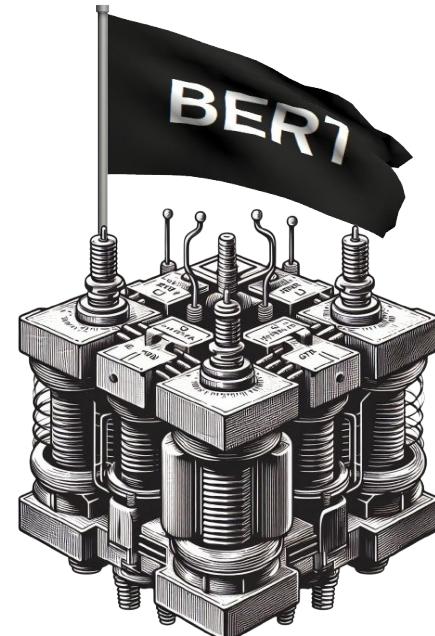
A
transformer
is
a
deep
learning
model
that
adopts
the
mechanism
of
self
attention



~ 15%

Input

[CLS]
A
[MASK]
is
a
deep
learning
model
that
[MASK]
the
mechanism
of
self
[MASK]



Target

/
/
transformer
/
/
/
/
/
adopts
/
/
/
/
attention

Masked words prediction

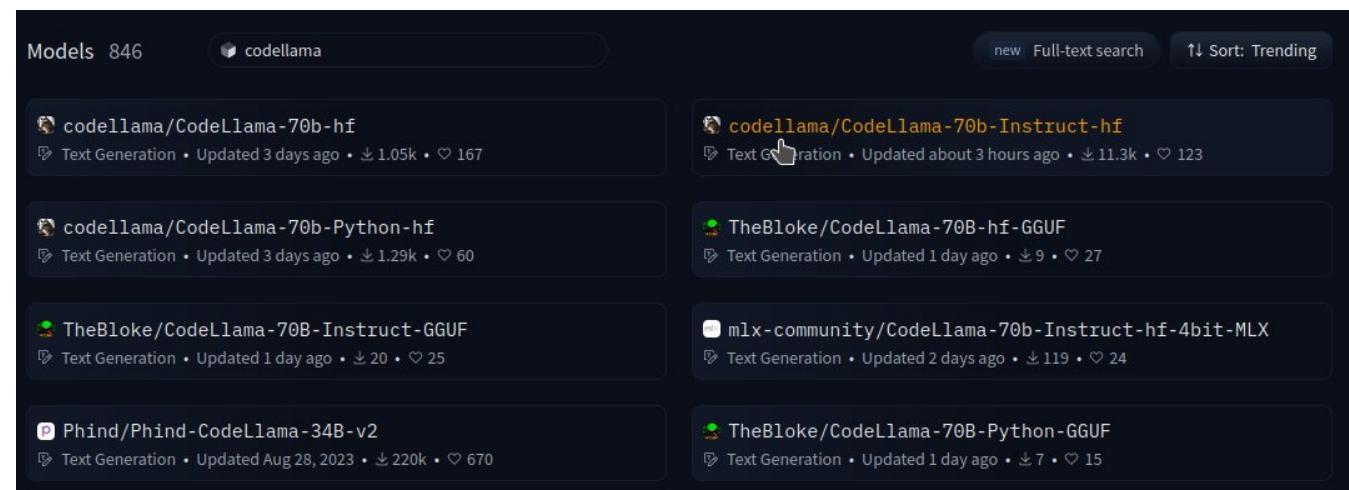
Finding a pretrained model

The largest free pretrained transformer models database

1 — Go to <https://huggingface.co/>



2 — Look for a model



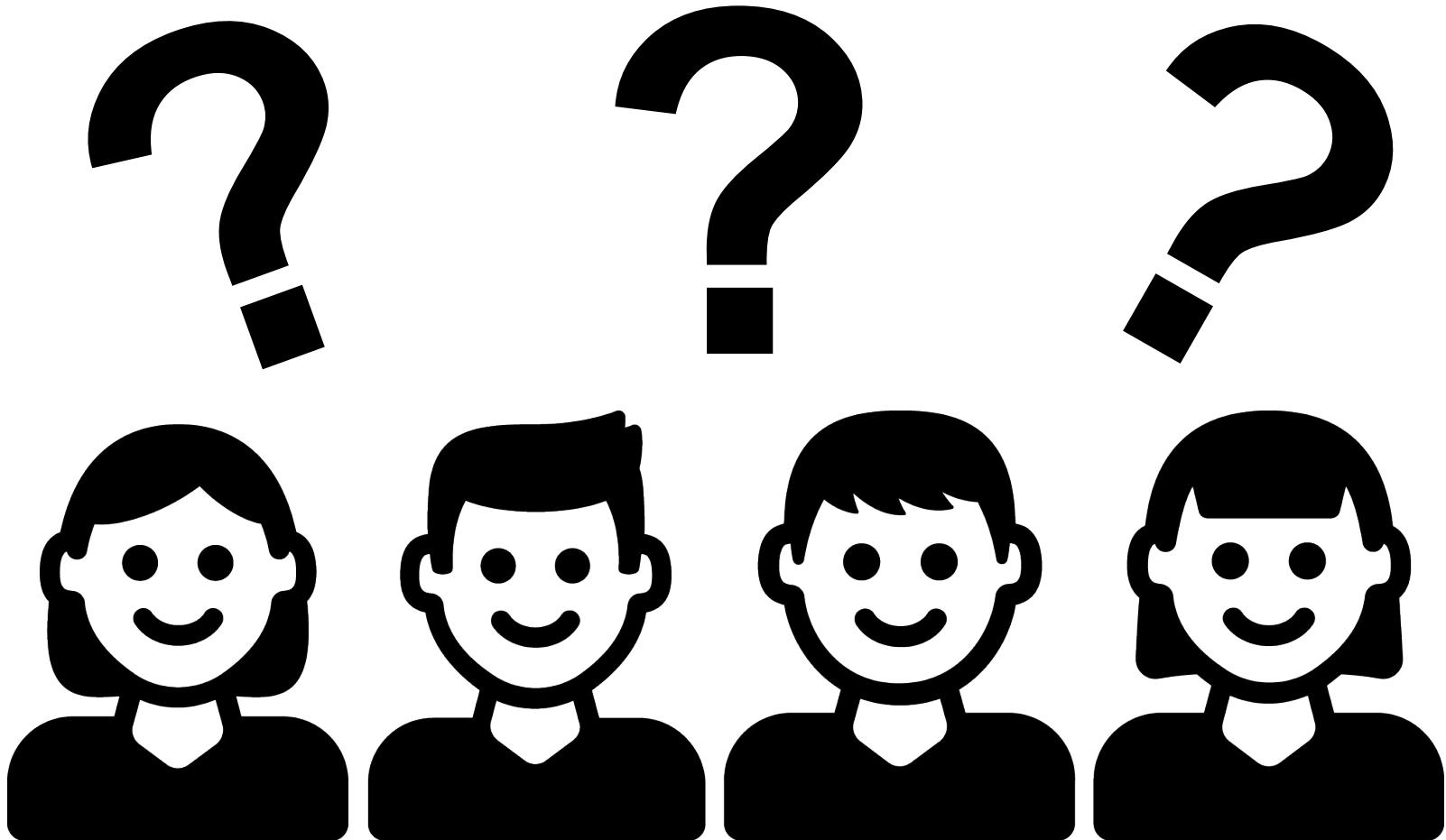
3 — Download the model

```
ncassereau@jean-zay:/fidle/transformers$ git clone https://huggingface.co/codellama/CodeLlama-70b-Instruct-hf
```

4 — Enjoy

```
1  from transformers import AutoTokenizer, AutoModelForCausalLM
2
3  path = "/fidle/transformers/CodeLlama-70b-Instruct-hf"
4  tokenizer = AutoTokenizer.from_pretrained(path)
5  model = AutoModelForCausalLM.from_pretrained(path)
6
7  input = tokenizer(["I love Jean Zay", "Hatim truly is the best"], padding=True)
8  output = model(**input)
9
```

Question break #4



9



«Attention is
All You Need»

Transformers



1

What is a **Transformer**? The magic
of the **Attention Mechanism**

2

The different **Transformers**
architectures

3

**Pre-training and Foundation
Models**

4

**Specialization of Foundation
Models (especially LLM)**

5

Example: IMDB Reviews

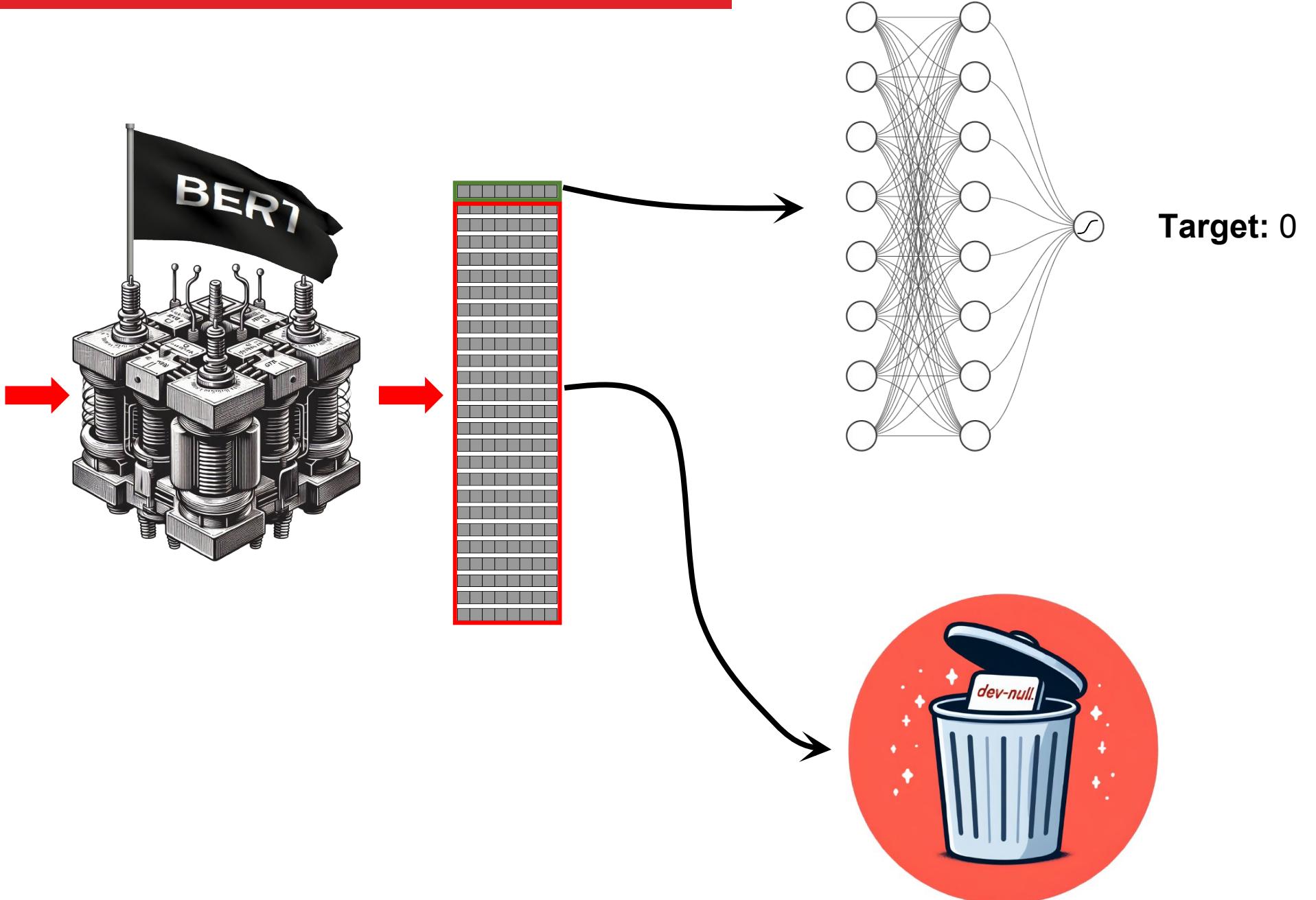


Fine-tuning of language models



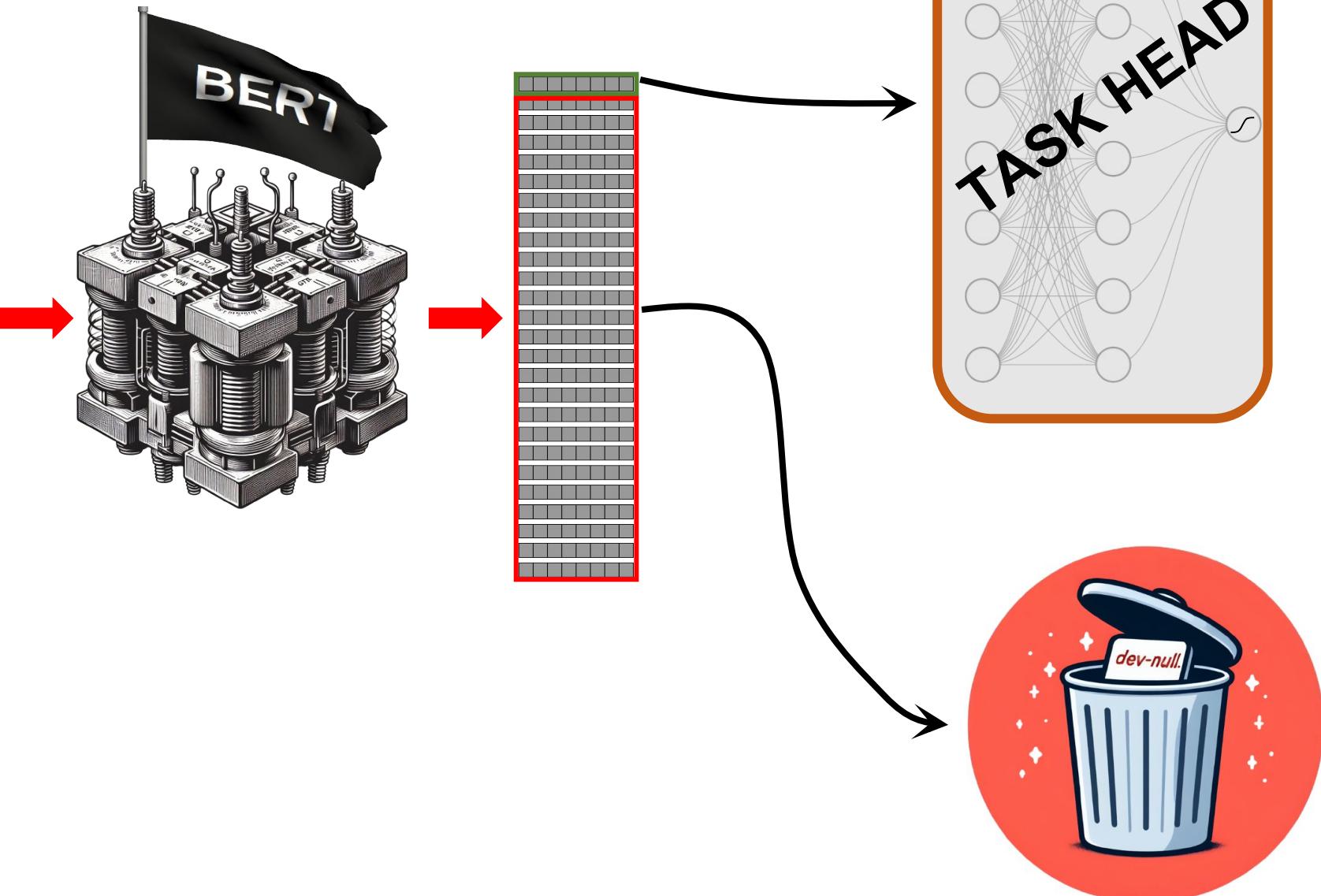
Fine-tuning a BERT — Sentence classification

[CLS] Star Wars
The Last Jedi is
really bad!



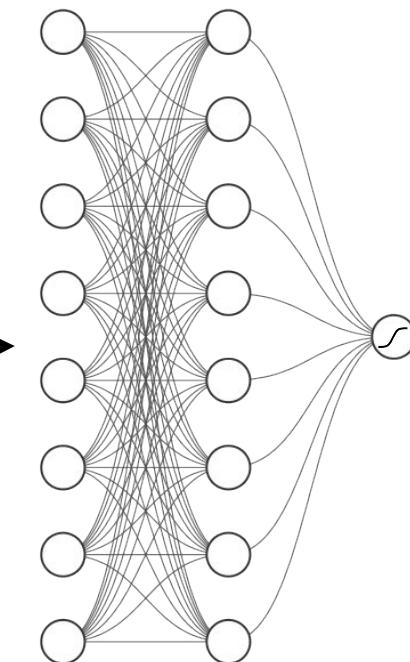
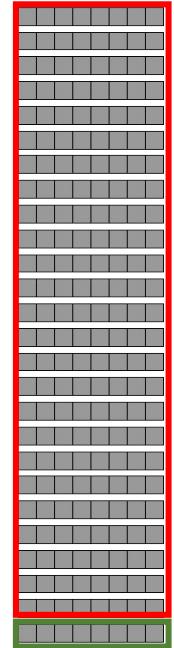
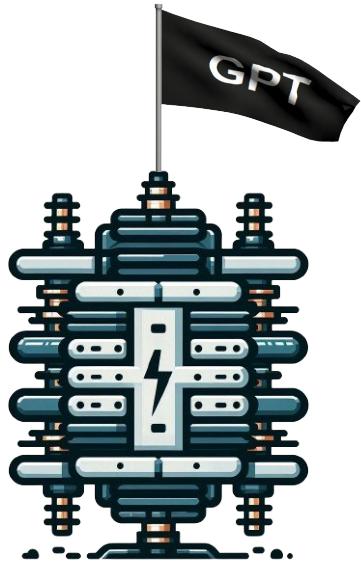
Fine-tuning a BERT — Sentence classification

[CLS] Star Wars
The Last Jedi is
really bad!



Fine-tuning a GPT — Sentence classification

Star Wars The
Last Jedi is really
bad!

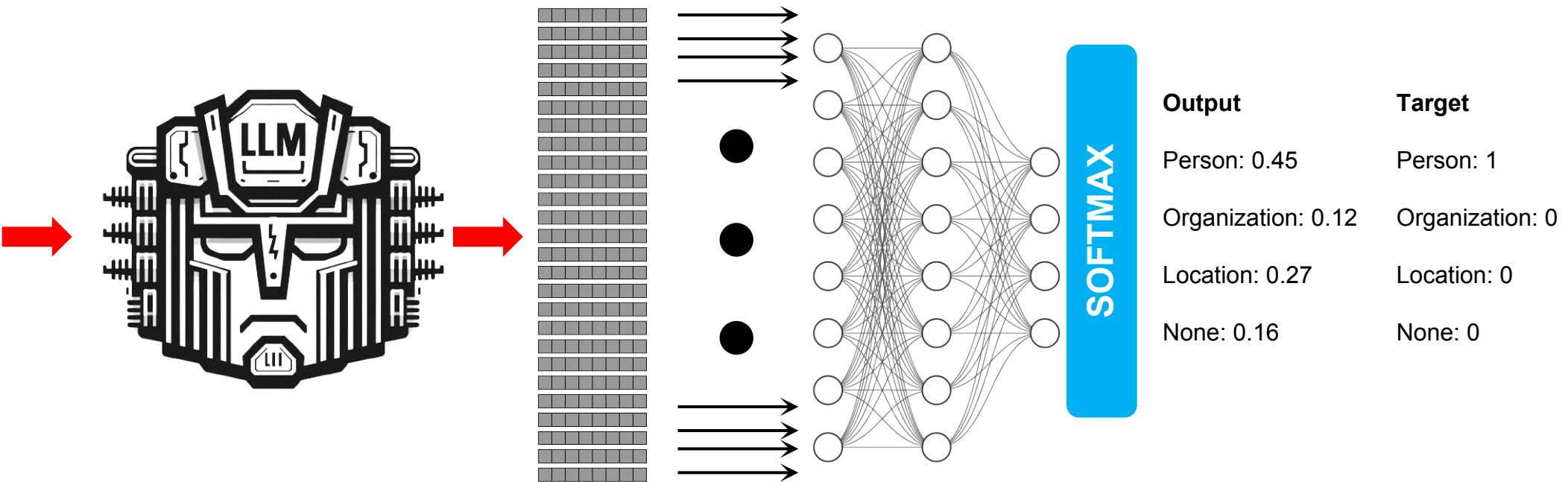


Target: 0



Fine-tuning a transformer — Named Entity Recognition

Zidane came to
Paris to watch
the Paris Saint-
Germain



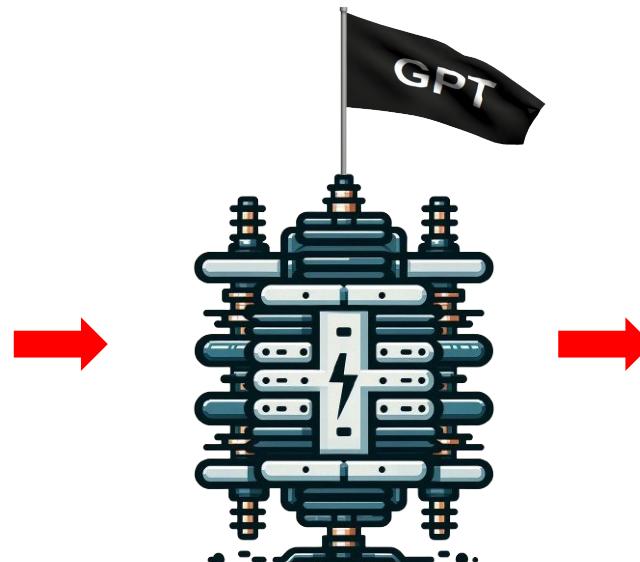
Fine-tuning a GPT with prompting

Review

This film is really trash!

Template

{{ REVIEW }} This review is (positive, negative or neutral):



Output

Positive: 0.18
Negative: 0.44
Neutral: 0.38

Target

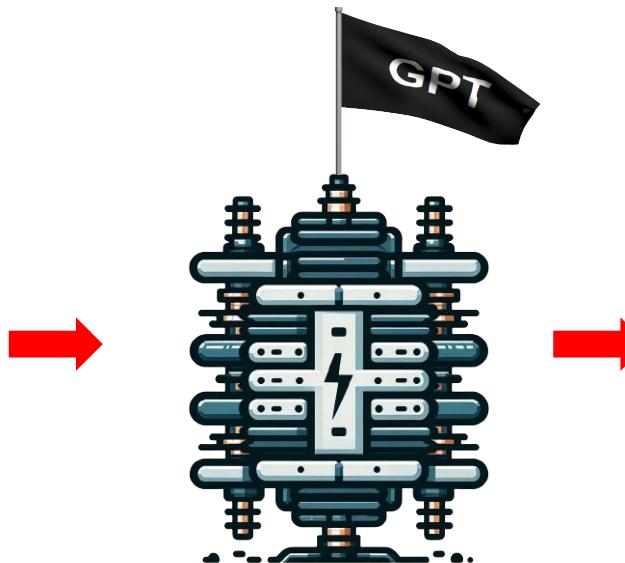
Positive: 0
Negative: 1
Neutral: 0

Fine-tuning a GPT — Example of summarization

Input

Template

{{ INPUT }} TL;DR:



GENERATION HEAD

Target

She dumped me
because I took a sip
before toasting.

Fine-tuning a GPT with templates

Type	Task	Input ([x])	Template	Answer ([z])
Text CLS	Sentiment	I love this movie.	[X] The movie is [Z].	great fantastic ...
	Topics	He prompted the LM.	[X] The text is about [Z].	sports science ...
	Intention	What is taxi fare to Denver?	[X] The question is about [Z].	quantity city ...
Text-span CLS	Aspect Sentiment	Poor service but good food.	[X] What about service? [Z].	Bad Terrible ...
	Text-pair CLS	[X1]: An old man with ...		Yes
		[X2]: A man walks ...	[X1] ? [Z], [X2]	No ...
Tagging	NER	[X1]: Mike went to Paris.		organization
		[X2]: Paris	[X1] [X2] is a [Z] entity.	location ...
				The victim ... A woman
Text Generation	Summarization	Las Vegas police ...	[X] TL;DR: [Z]	I love you. A woman
	Translation	Je vous aime.	French: [X] English: [Z]	I fancy you. ...

Orienting the style with templates

System

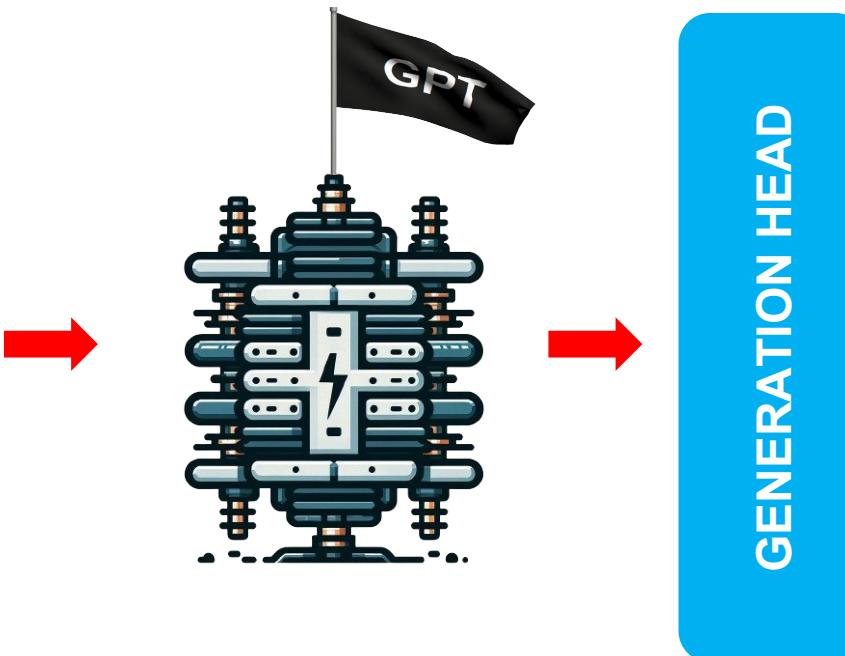
Speak as if you were a pirate.

Input

What is the Cayley-Hamilton theorem ?

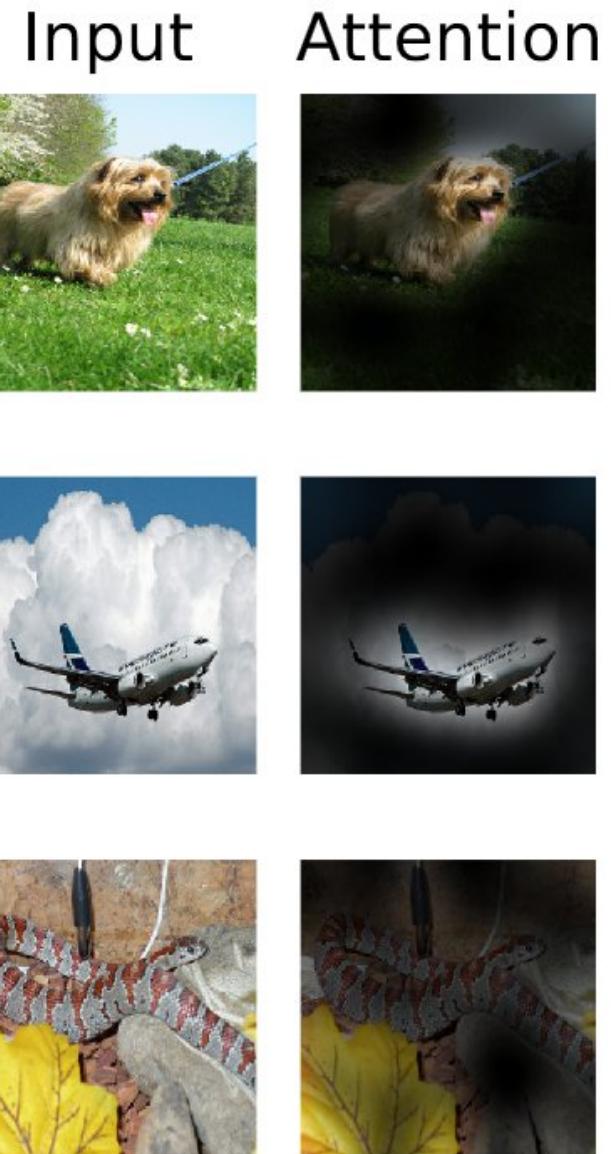
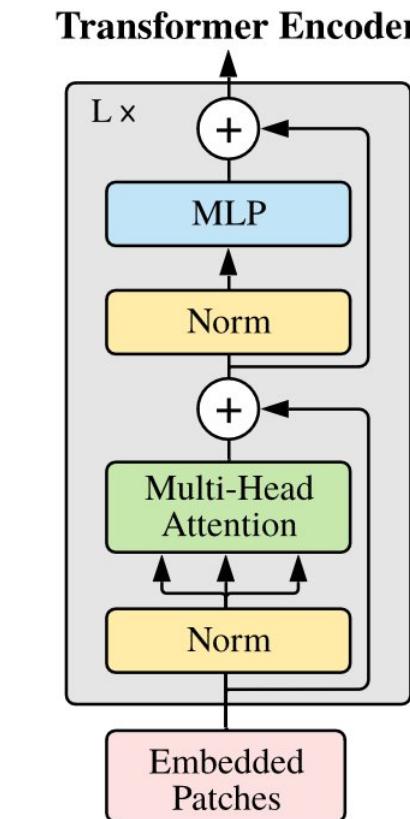
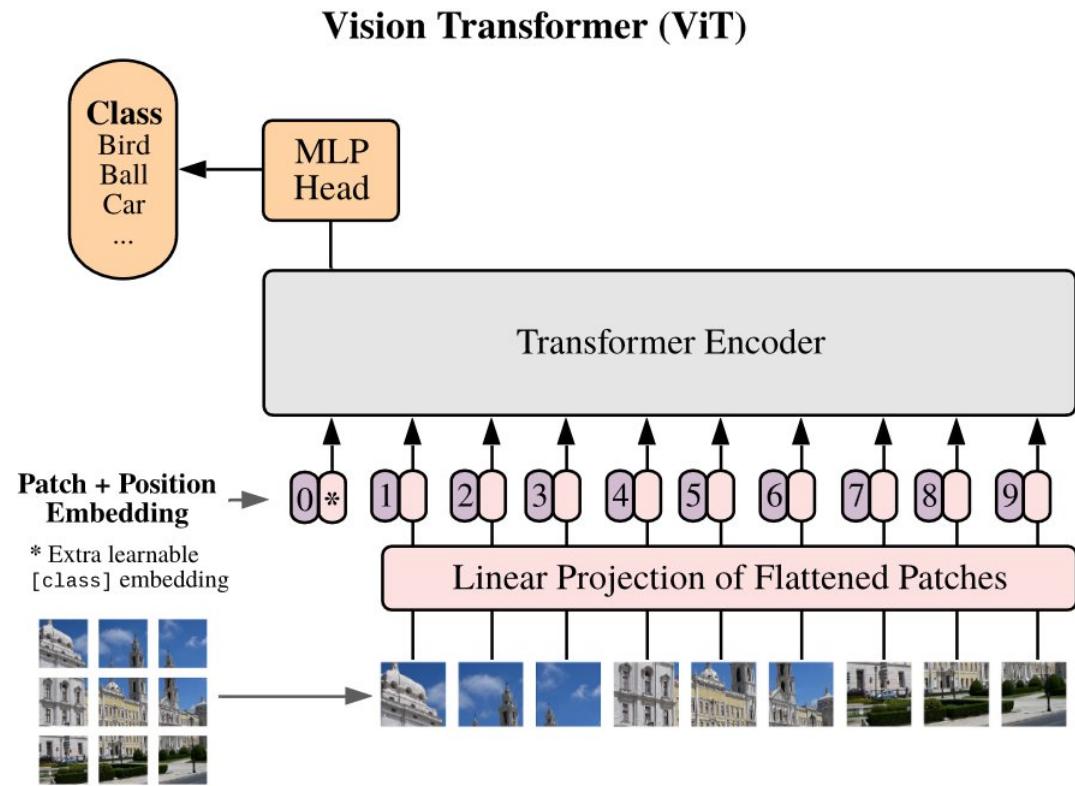
Template

```
<|system|>
{{ SYSTEM }}
<|user|>
{{ INPUT }}
<|assistant|>
```



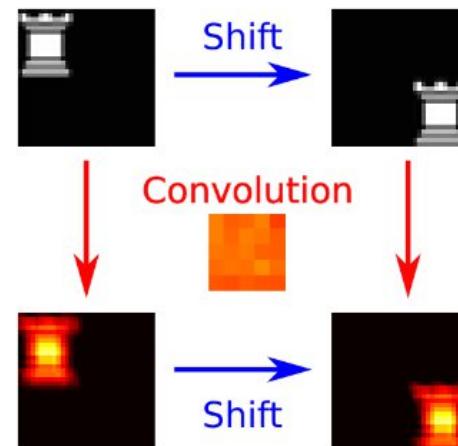
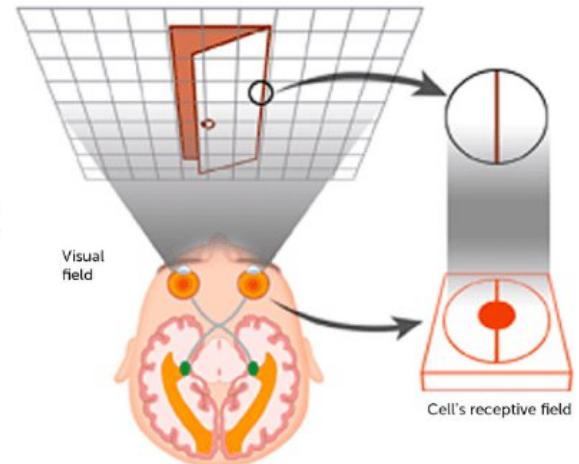
Ahoy, me hearties! I be hearin' 'bout a mathematical concept called the Cayley-Hamilton theorem. It's a fancy way o' sayin' that every square matrix is a root of its own characteristic equation, which is a polynomial equation that describes the matrix's properties. This theorem has important implications in linear algebra and matrix theory, as it allows for the calculation of matrix powers and inverses in certain cases. So, if you ever find yourself in a sea of matrices, remember the Cayley-Hamilton theorem and you'll be sailing smoothly!

Vision Transformer

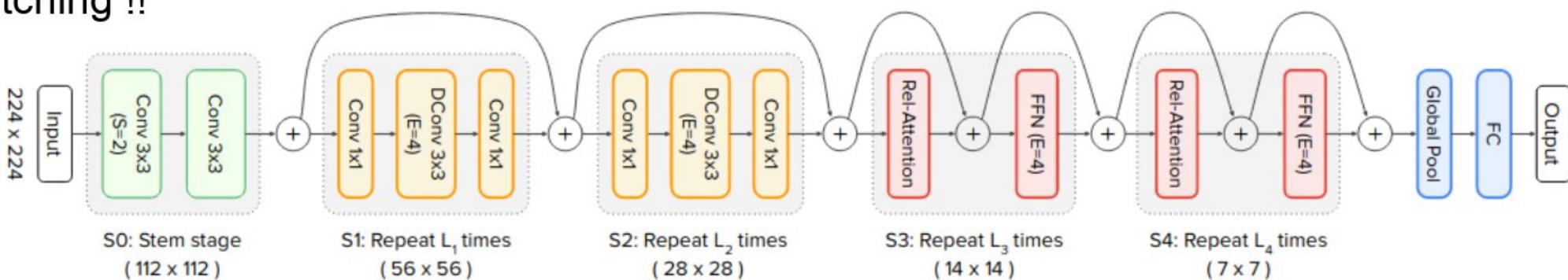


ConvNet

Translation Equivariance

**Self-Attention Net**Global Receptive Field
Context comprehension

No patching !!



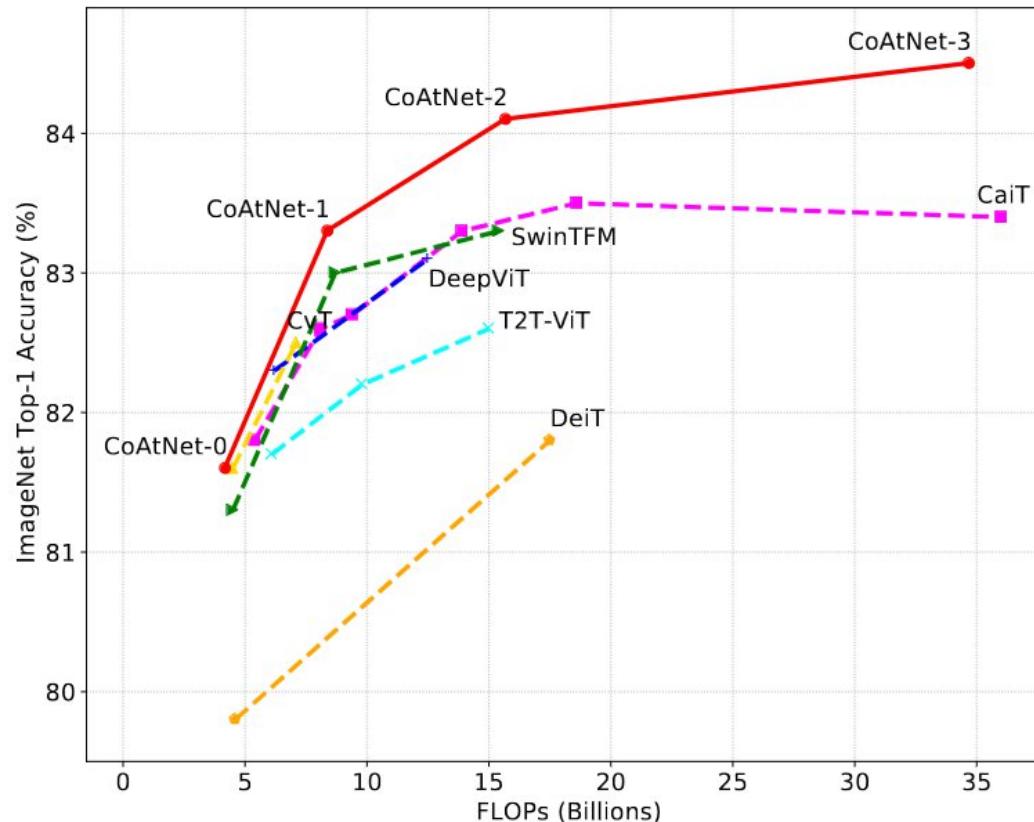


Figure 2: Accuracy-to-FLOPs scaling curve under ImageNet-1K only setting at 224x224.

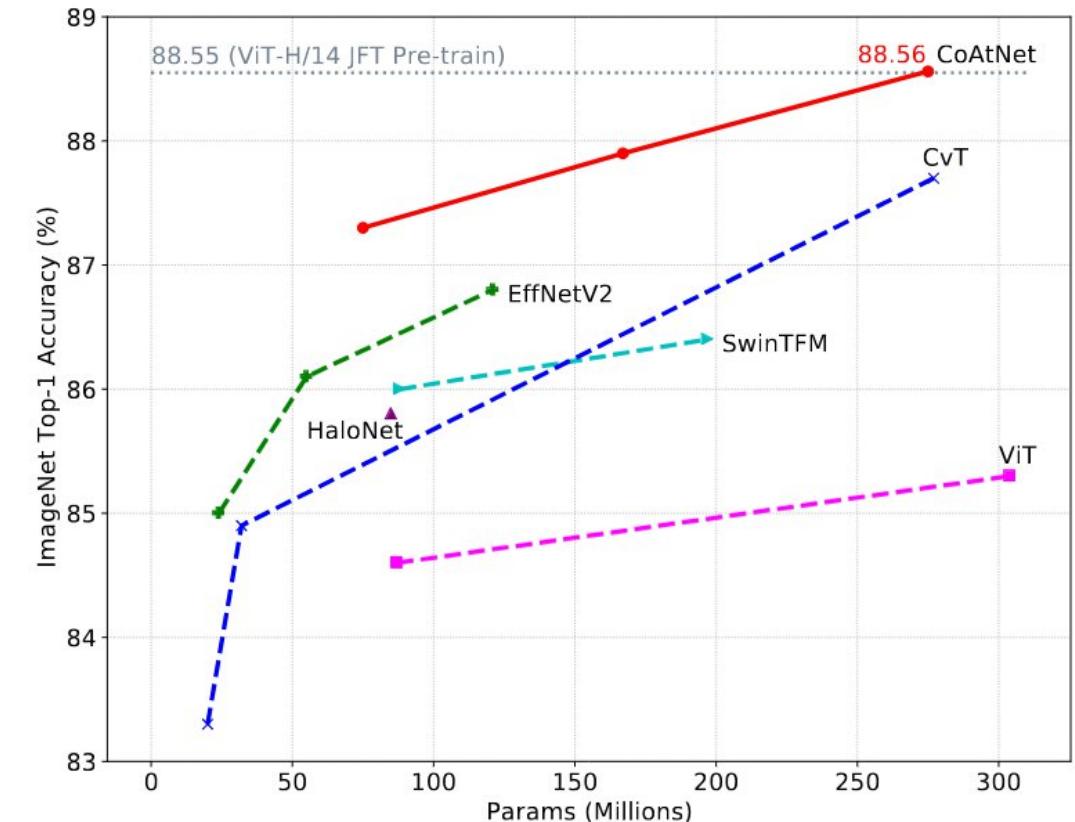
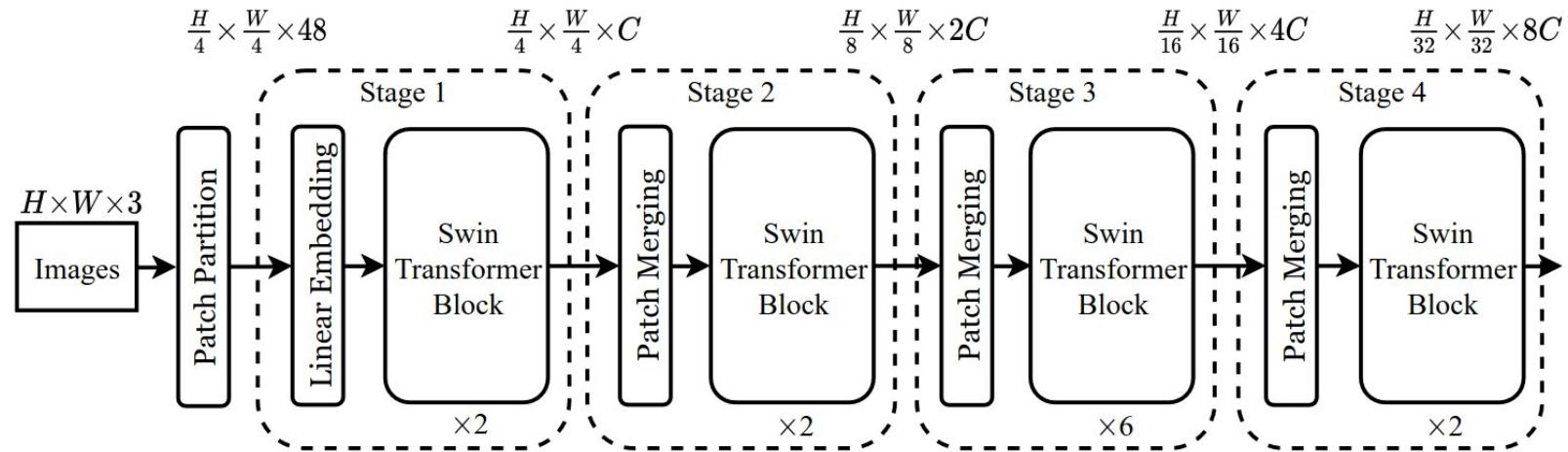
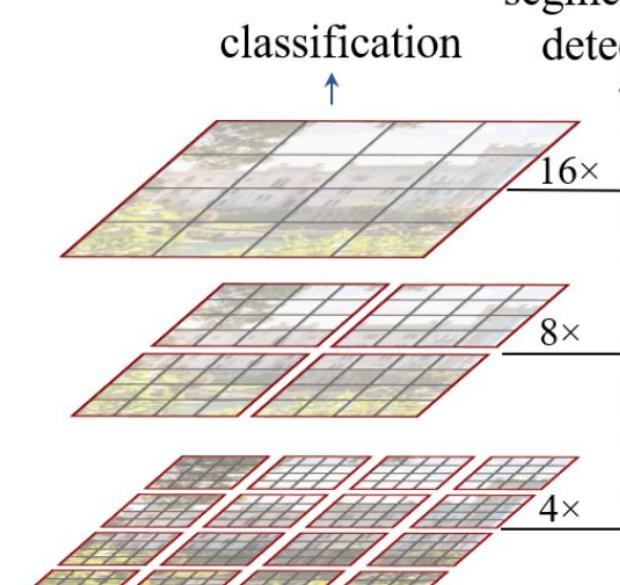


Figure 3: Accuracy-to-Params scaling curve under ImageNet-21K ⇒ ImageNet-1K setting.

SwinTransformer

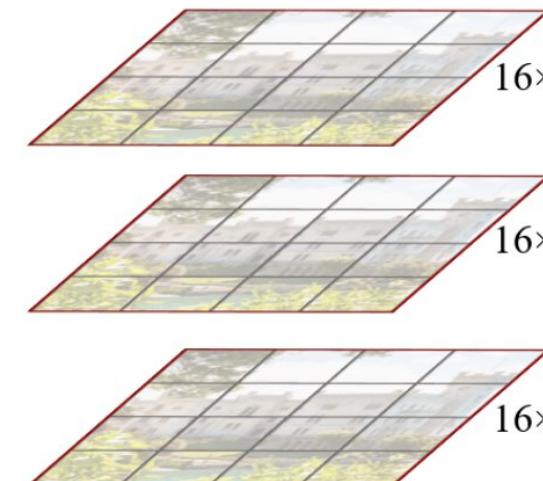


classification
segmentation
detection ...



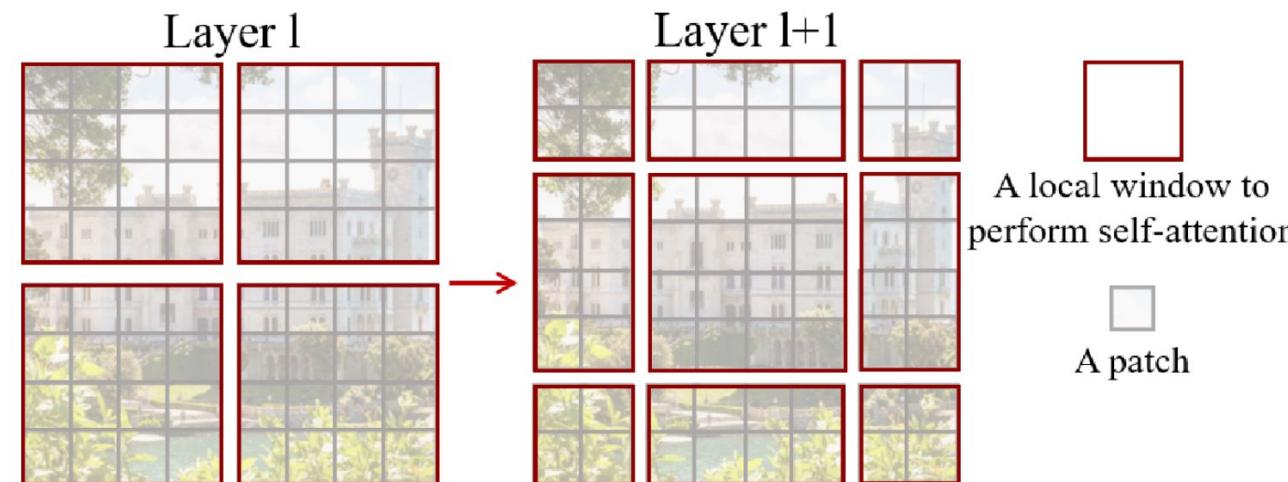
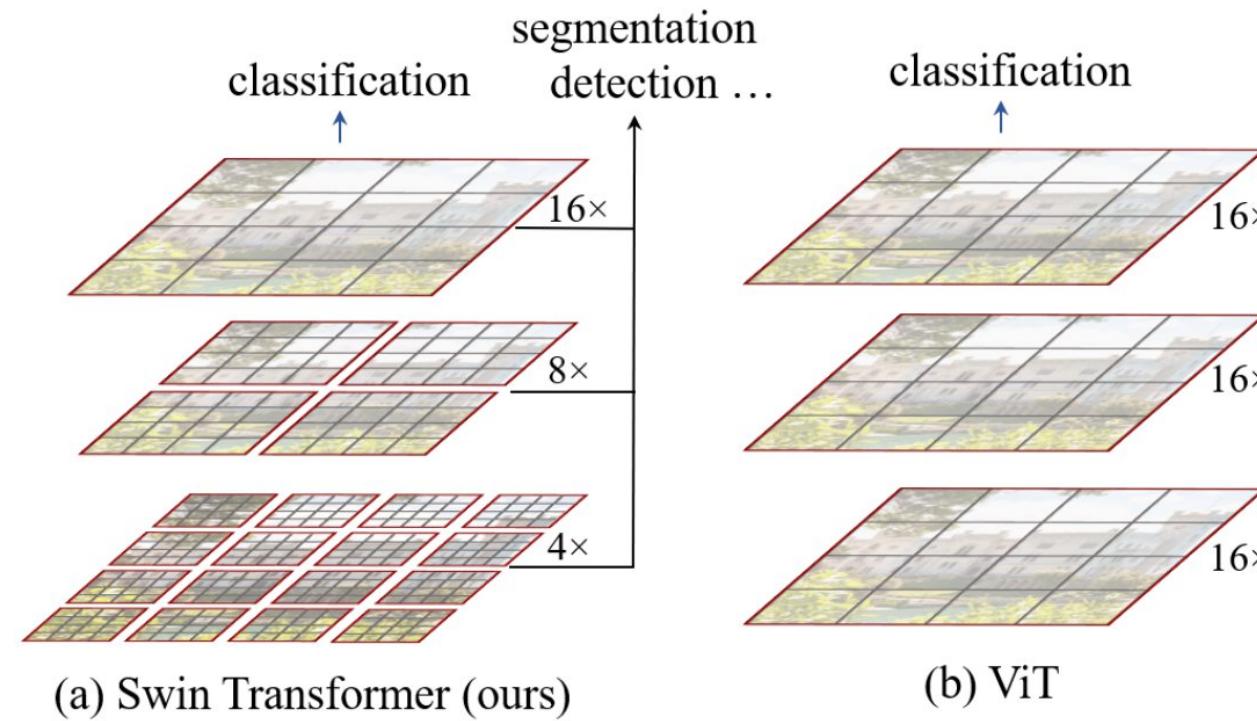
(a) Swin Transformer (ours)

classification



(b) ViT

SwinTransformer



9



«Attention is
All You Need»

Transformers



1

What is a **Transformer**? The magic
of the **Attention Mechanism**

2

The different **Transformers**
architectures

3

Pre-training and Foundation
Models

4

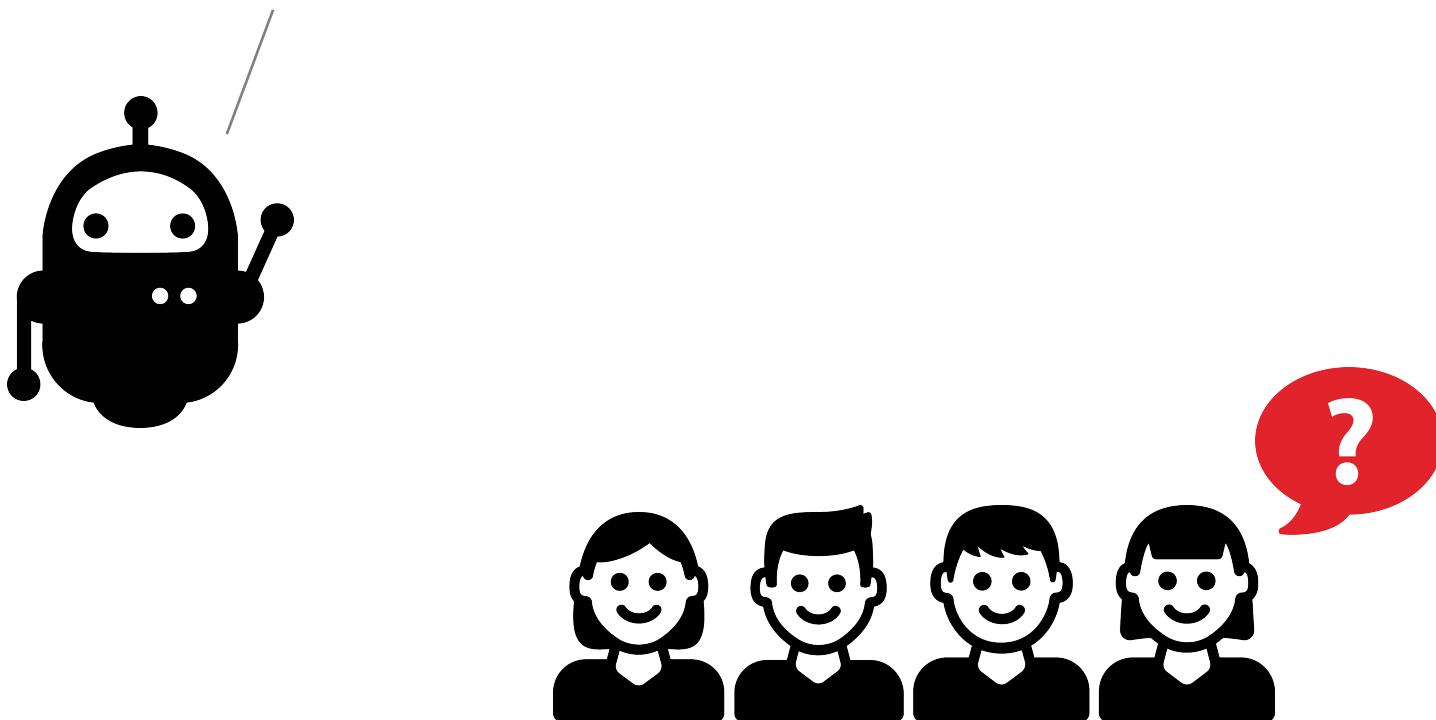
Specialization of Foundation
Models (especially LLM)

5

Example: IMDB Reviews



Merci beaucoup !!
Quelques questions ?





FIDLE

<https://youtube.com/@CNRS-FIDLE>

<https://fidle.cnrs.fr>

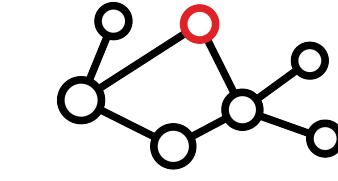


Attribution-NonCommercial-NoDerivatives 4.0 International (CC BY-NC-ND 4.0)
<https://creativecommons.org/licenses/by-nc-nd/4.0/>

Next, on Fidle :



6



Graph Neural
Network

GNN

Jeudi
15
Février
à 14h00