# Large Language Models

# Prompts and completions

**Prompt**

> Where is Ganymede located in the solar system?
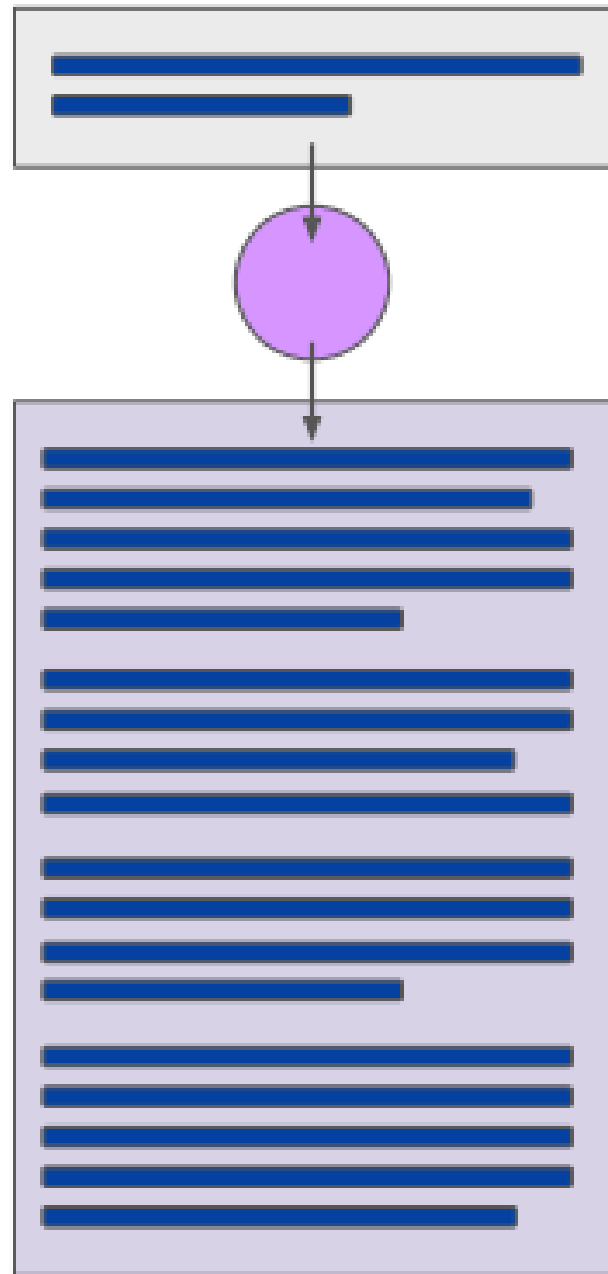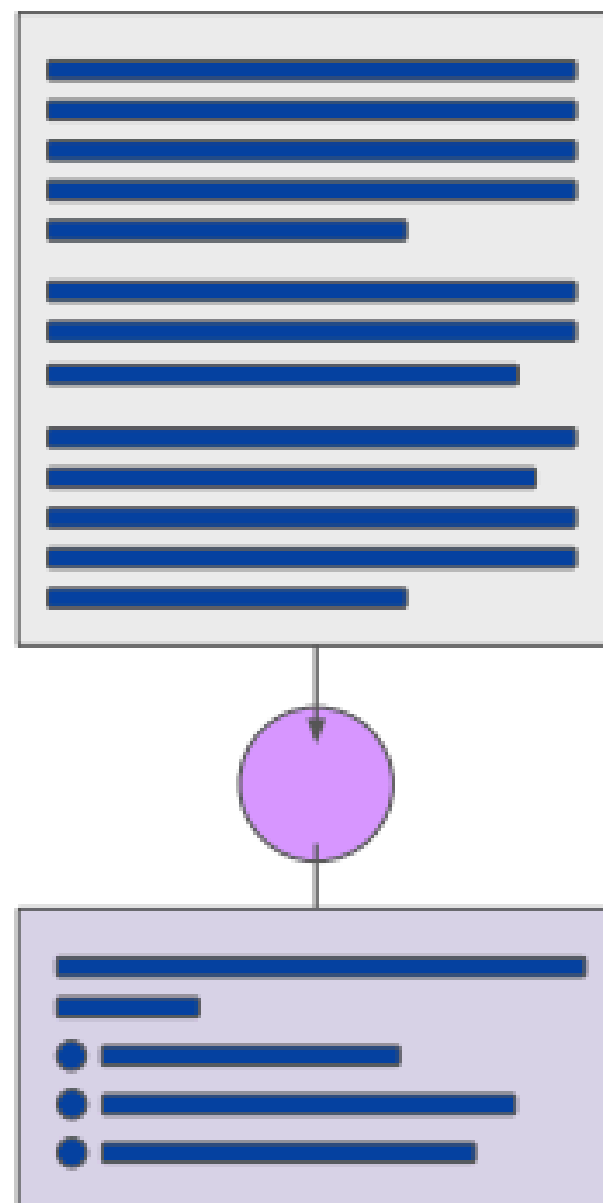
**Model**

LLM

**Completion**

> Where is Ganymede located in the solar system?
>
> Ganymede is a moon of Jupiter and is located in the solar system within Jupiter's orbit.

Context window
- typically a few 1000 words.

DeepLearning.AI

aws

# LLM use cases & tasks

Essay Writing

Summarization

Translation

Information retrieval

Invoke APIs and actions

Action call

External Applications

# The significance of scale: language understanding

BERT*
110M

BLOOM
176B

*Bert-base

DeepLearning.AI

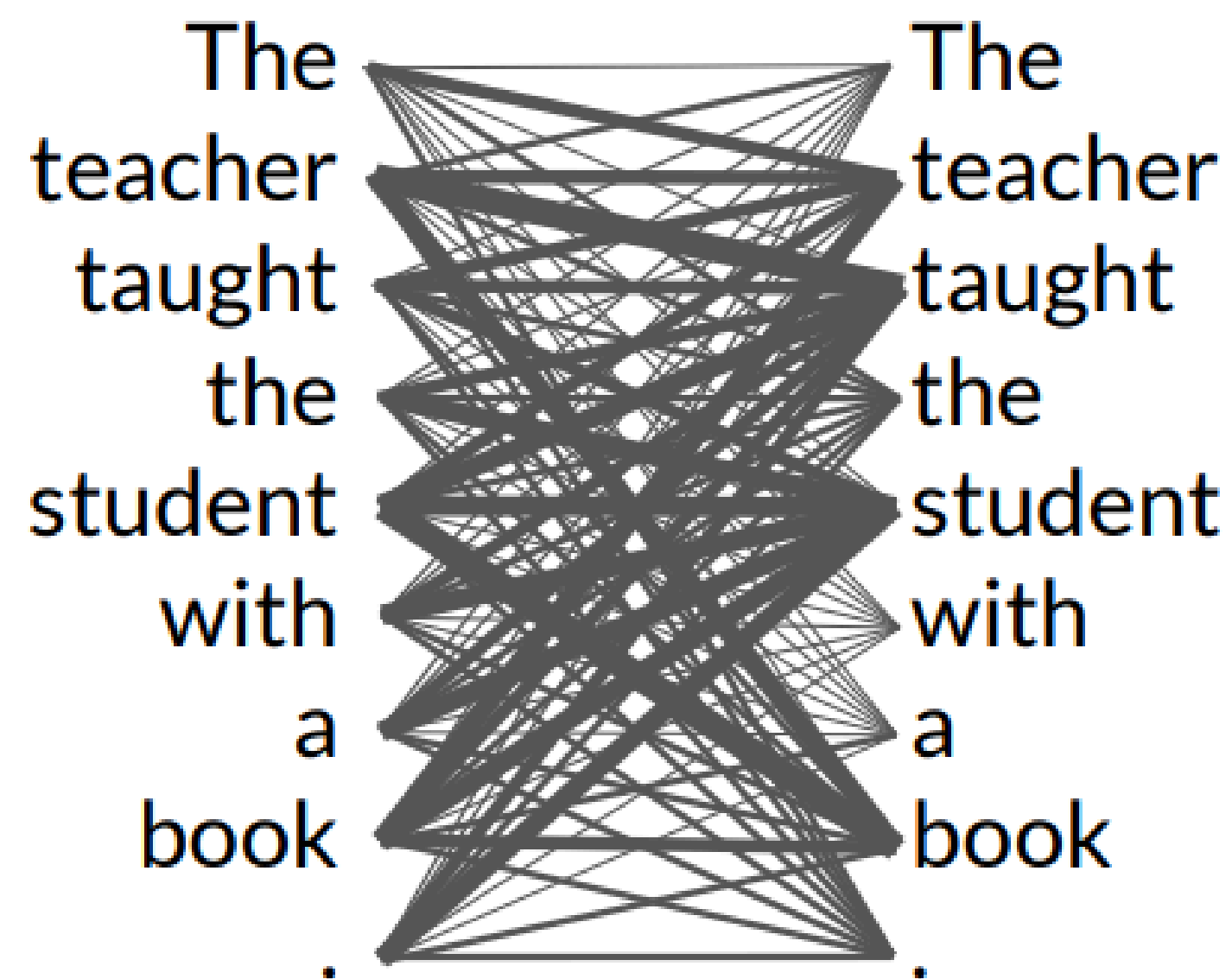How LLMs work -
Transformers architecture

# Transformers



The teacher taught the student with the book.

# Self-attention

# Self-attention

The
teacher
taught
the
student
with
a
book
.

The
teacher
taught
the
student
with
a
book
.

# Transformers

# Transformers

Encoder

Decoder

Embedding          Embedding

Inputs

Tokenizer          | 342 | 879 | 432 | 342 |          Token IDs

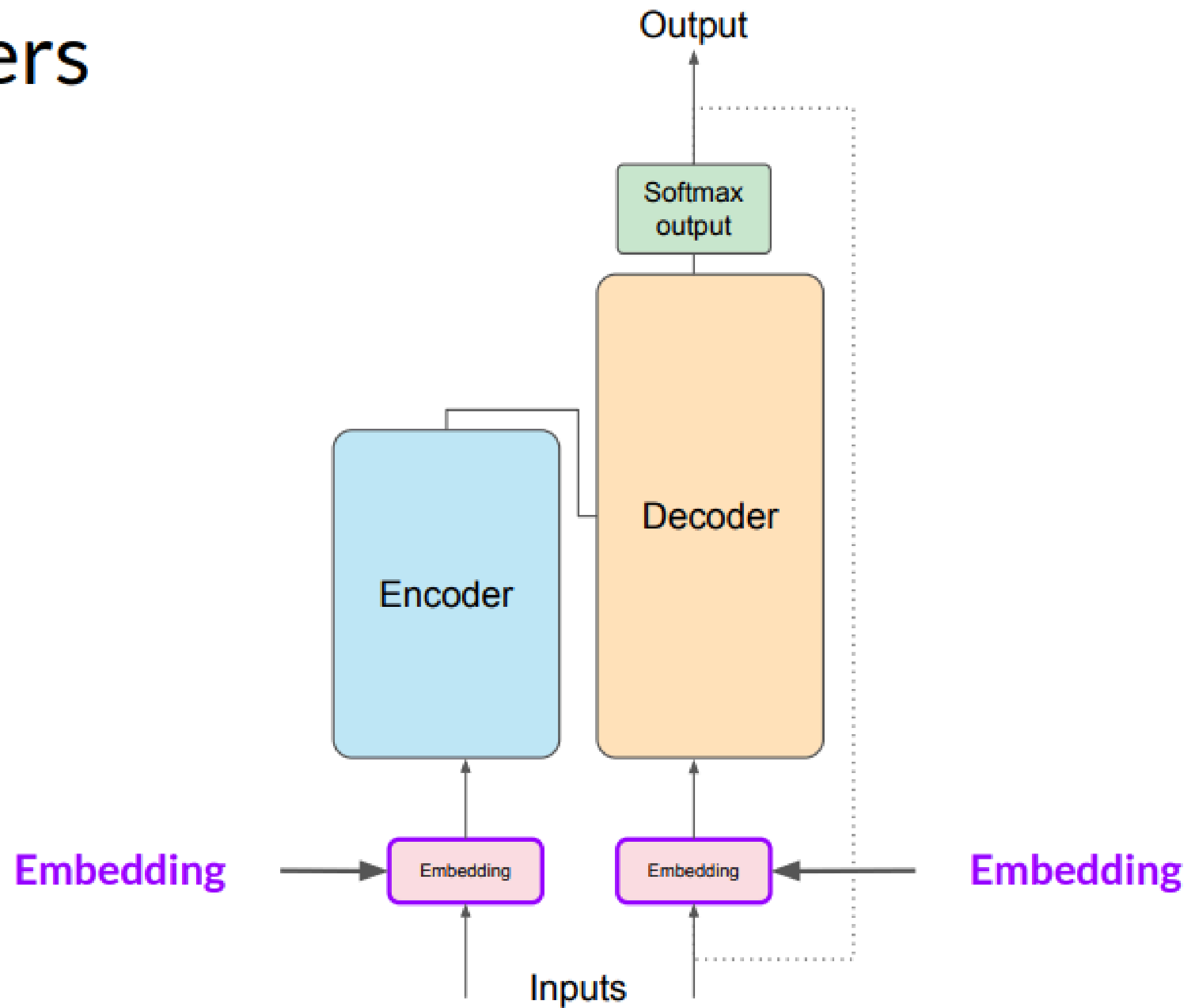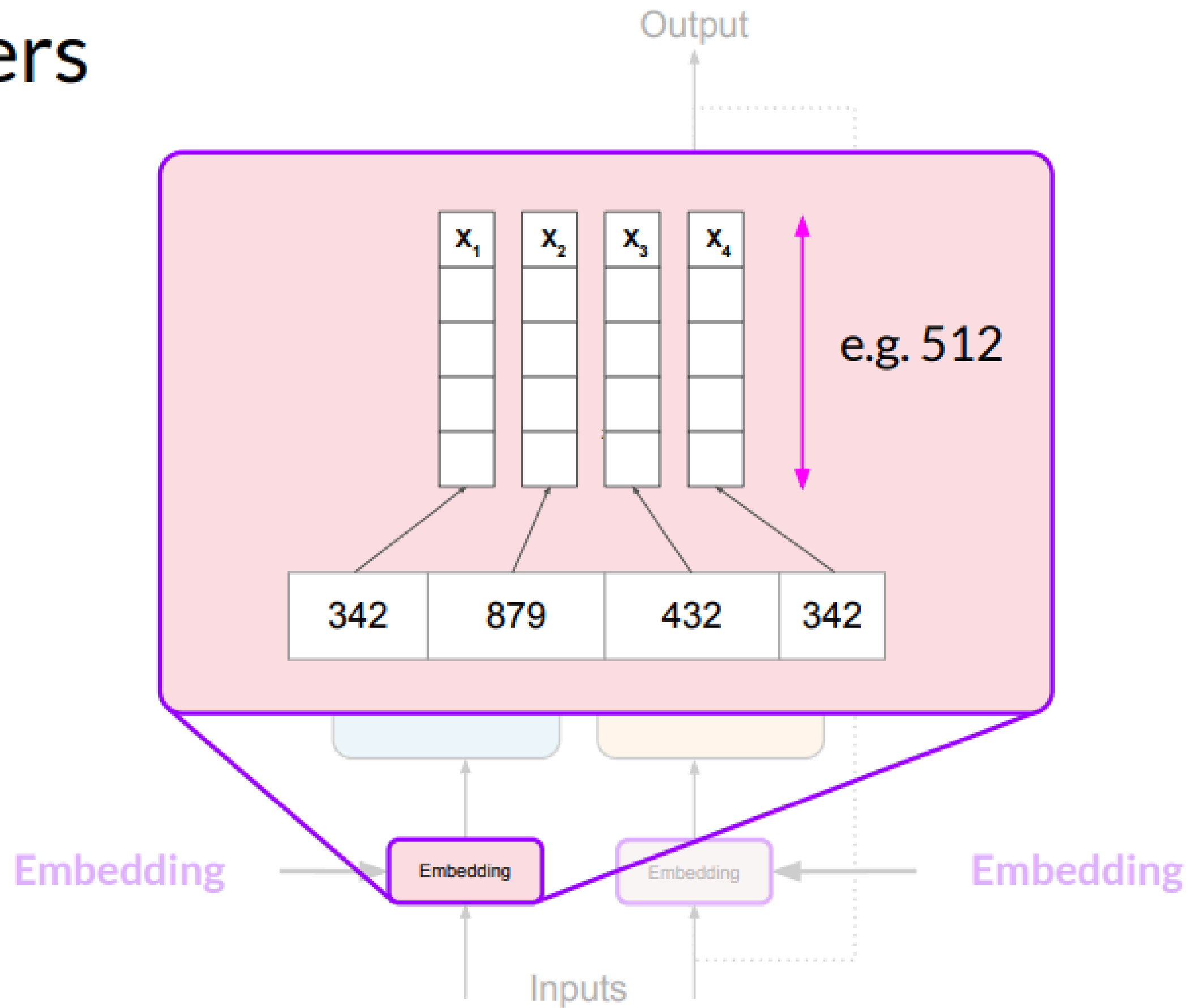Input:          the     teacher     taught     the

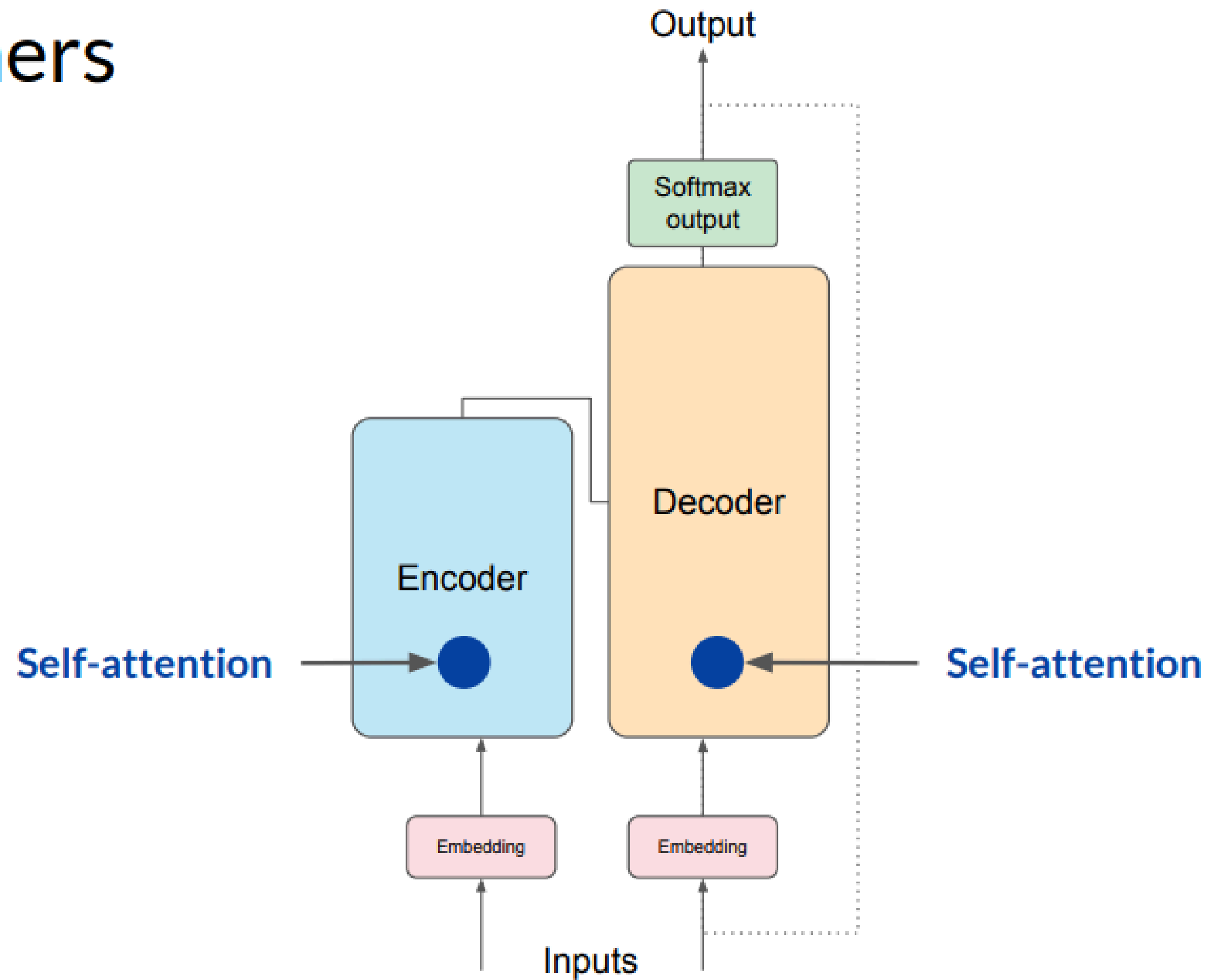DeepLearning.AI          aws

# Transformers

# Transformers

# Transformers
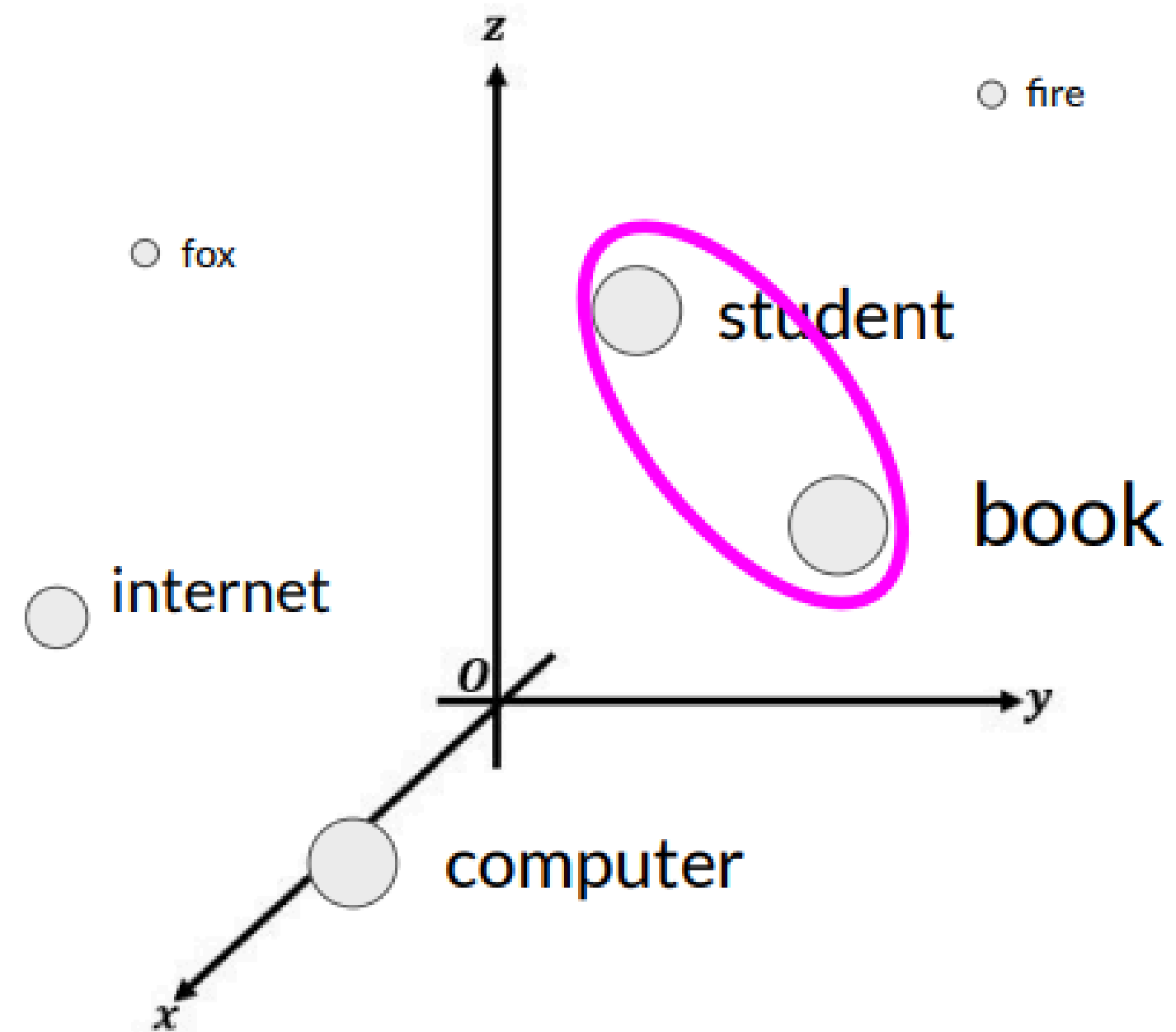
# Transformers

# Transformers

# Transformers

# Transformers



Output

Softmax output

| P1 | P2 | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | Pn |

Encoder

Embedding

Embedding

Inputs

# Transformers

Translation:
sequence-to-sequence task

# Transformers

Translation:
sequence-to-sequence task

# Transformers



**Encoder**

Encodes inputs ("prompts") with contextual understanding and produces one vector per input token.

**Decoder**

Accepts input tokens and generates new tokens.

# Transformers

# Summary of in-context learning (ICL)

**Prompt** // Zero Shot

```
Classify this review:
I loved this movie!
Sentiment:
```
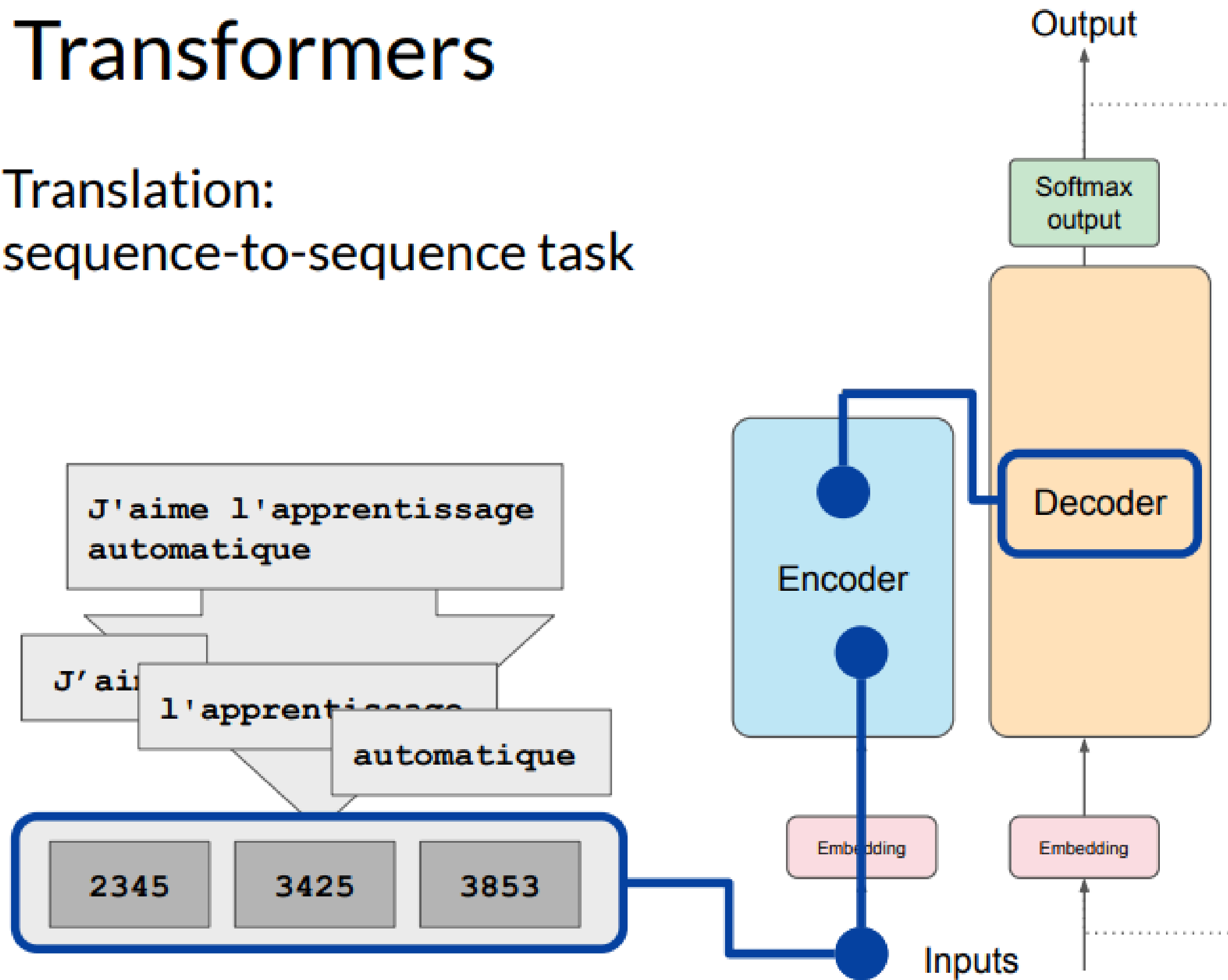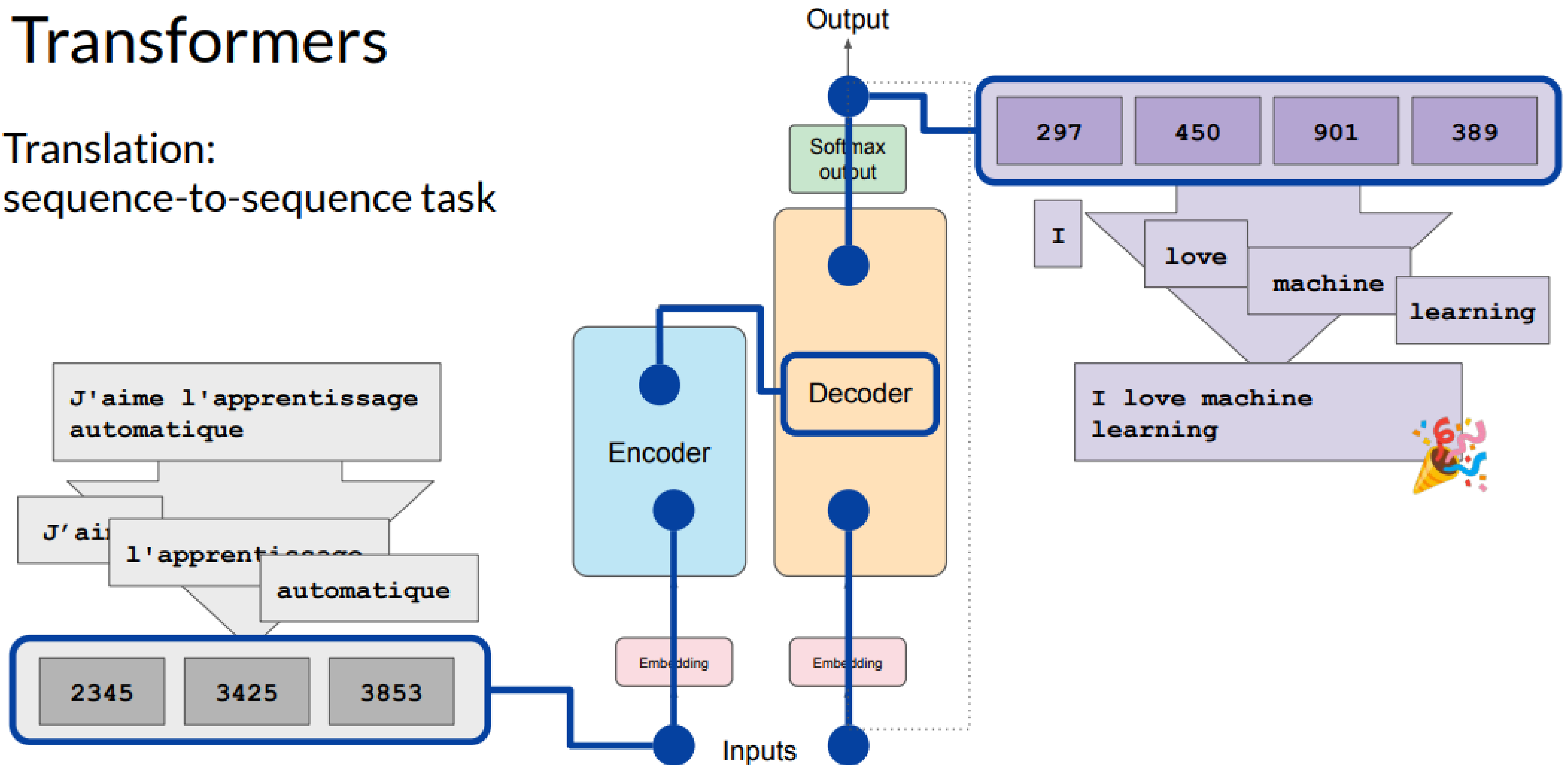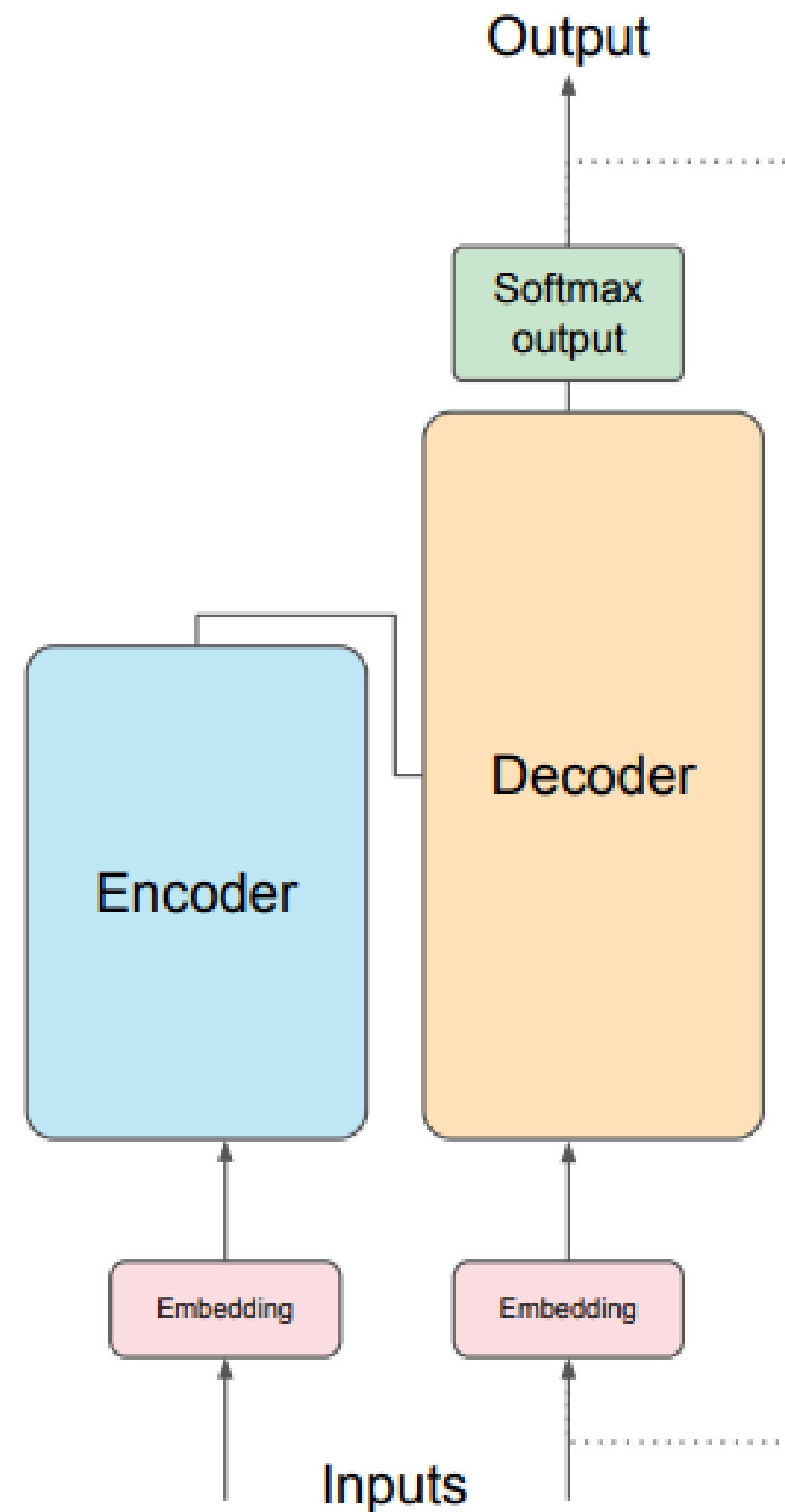
**Context Window**
(few thousand words)

**Prompt** // One Shot

```
Classify this review:
I loved this movie!
Sentiment: Positive

Classify this review:
I don't like this
chair.
Sentiment:
```

**Prompt** // Few Shot    >5 or 6 examples

```
Classify this review:
I loved this movie!
Sentiment: Positive

Classify this review:
I don't like this
chair.
Sentiment: Negative

Classify this review:
Who would use this
product?
Sentiment:
```

# The significance of scale: task ability

BERT*
110M

BLOOM
176B →

*Bert-base

# Generative AI project lifecycle



| Scope | Select | Adapt and align model | | Application integration | |
|---|---|---|---|---|---|
| Define the use case | Choose an existing model or pretrain your own | Prompt engineering<br><br>Fine-tuning<br><br>Align with human feedback | Evaluate | Optimize and deploy model for inference | Augment model and build LLM-powered applications |

DeepLearning.AI

aws

# Model size vs. time



BERT-L
340M

GPT-2
1.5B

GPT-3
175B

PaLM
540B

Growth powered by:
- Introduction of transformer
- Access to massive datasets
- More powerful compute resources

2018          2022          2023

DeepLearning.AI          aws

# Model size vs. time



BERT-L
340M

GPT-2
1.5B

GPT-3
175B

PaLM
540B

increase?

Trillion(s)

2018          2022      2023

DeepLearning.AI      aws