# Demystifying ChatGPT
## A shallow deep-dive into LLMs

SKYLYZE | DSA

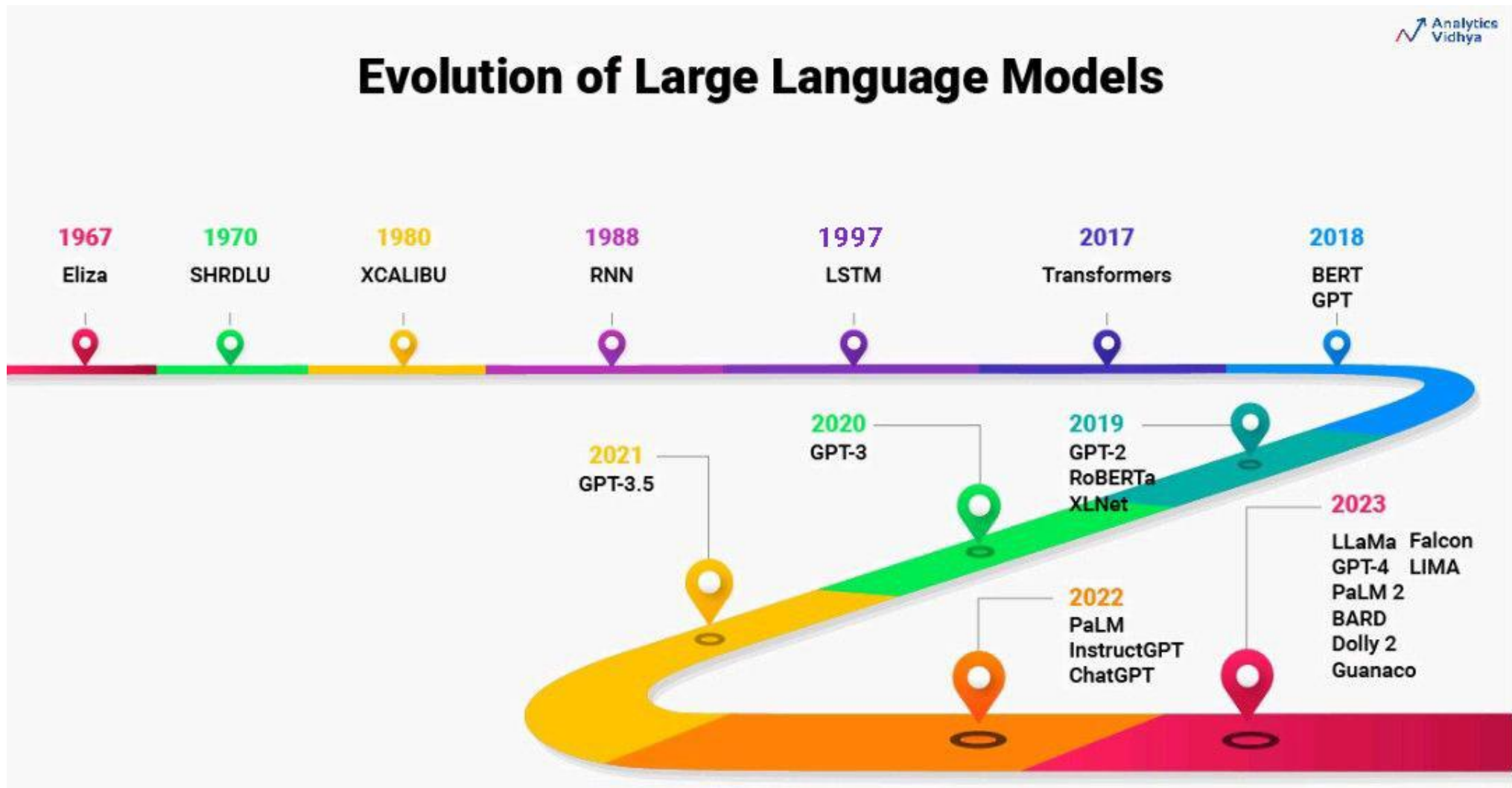**DSA Daten- und Systemtechnik GmbH | SKYLYZE | Roland Stoffel | Pascal Rößner**

# AGENDA

- Introduction ChatGPT aka LLM (Lecture: 15min)
  - Defintion and Use-Cases
  - Theoretical Background
  - Project Lifecycle
- From Scratch (Hands-On Session: 35-45min)
  - Train your own model
- Challenges of LLMS and Generative AI (Lecture: 10min)
  - Infrastructure and Costs
  - Transfer to Image Generation
- Experimentation (Hands-On Session: 20min)
  - Have Fun with Stable Diffusion
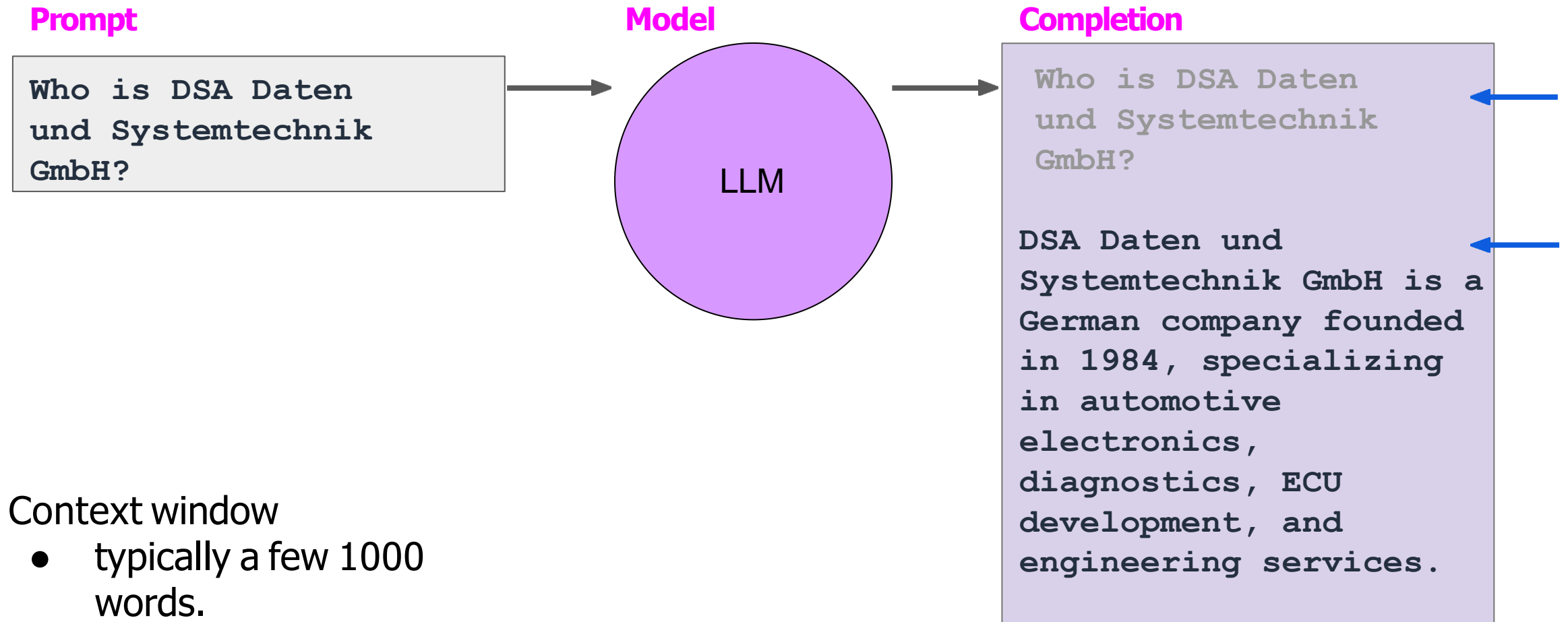
# WHAT IS A LLM

- ChatGPT is a so-called Large Language Model (LLM)

- Large Language Models are deep learning models trained on huge datasets to perform NLP tasks.

- Its core objective is to learn and understand human languages precisely.

- Large Language Models enable the machines to interpret languages just like the way we, as humans, interpret them.

| Name | Release Date | Parameter Size |
|------|--------------|----------------|
| GPT-4 | 2023 | 4.6 trillion |
| LLaMA | 2022 | 1.5 trillion |
| FLAN UL2 | 2022 | 1.3 trillion |
| BLOOM | 2022 | 176 billion |
| LaMDA | 2021 | 173 billion |
| MT-NLG | 2020 | 530 billion |
| GPT-3 | 2020 | 175 billion |
| GPT-2 | 2019 | 1.5 billion |
| BERT | 2018 | 340 million |
| ELMo | 2017 | 94 million |

# HISTORICAL EVOLUTION

# PROMPTS AND COMPLETIONS

**Prompt**

```
Who is DSA Daten
und Systemtechnik
GmbH?
```

**Model**

LLM

**Completion**

```
Who is DSA Daten
und Systemtechnik
GmbH?

DSA Daten und
Systemtechnik GmbH is a
German company founded
in 1984, specializing
in automotive
electronics,
diagnostics, ECU
development, and
engineering services.
```

Context window
- typically a few 1000 words.

The teacher's book?

The teacher taught the student with the book.

The student's book?

# TRANSFORMER ARCHITECTURE

## Attention Is All You Need

**Ashish Vaswani***
Google Brain
avaswani@google.com

**Noam Shazeer***
Google Brain
noam@google.com

**Niki Parmar***
Google Research
nikip@google.com

**Jakob Uszkoreit***
Google Research
usz@google.com

**Llion Jones***
Google Research
llion@google.com

**Aidan N. Gomez*** †
University of Toronto
aidan@cs.toronto.edu

**Łukasz Kaiser***
Google Brain
lukaszkaiser@google.com

**Illia Polosukhin*** ‡
illia.polosukhin@gmail.com

### Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to

Output
Probabilities

Softmax

Linear

Add & Norm
Feed
Forward

Add & Norm
Multi-Head
Attention

Nx

Add & Norm
Feed
Forward

Nx

Add & Norm
Multi-Head
Attention

Add & Norm
Masked
Multi-Head
Attention

Positional
Encoding

Positional
Encoding

Input
Embedding

Output
Embedding

Inputs

Outputs
(shifted right)

# TRANSFORMER ARCHITECTURE

## Attention Is All You Need

**Ashish Vaswani***
Google Brain
avaswani@google.com

**Noam Shazeer***
Google Brain
noam@google.com

**Niki Parmar***
Google Research
nikip@google.com

**Jakob Uszkoreit***
Google Research
usz@google.com

**Llion Jones***
Google Research
llion@google.com

**Aidan N. Gomez*** [†]
University of Toronto
aidan@cs.toronto.edu

**Łukasz Kaiser***
Google Brain
lukaszkaiser@google.com

**Illia Polosukhin*** [‡]
illia.polosukhin@gmail.com

### Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to

- Scale efficiently
- Parallelprocess
- Attention to input meaning

# ATTENTION IS ALL YOU NEED

The teacher taught the student with the book.

# ATTENTION IS ALL YOU NEED

The teacher taught the student with the book.

# ATTENTION IS ALL YOU NEED

The teacher taught the student with the book.

# SELF-ATTENTION

The
teacher
taught
the
student
with
a
book
.

The
teacher
taught
the
student
with
a
book
.

# SELF-ATTENTION

The
teacher
taught
the
student
with
a
book
.

The
teacher
taught
the
student
with
a
book
.

# SELF-ATTENTION

The
teacher
taught
the
student
with
a
book
.

The
teacher
taught
the
student
with
a
book
.

# TRANSFORMERS

# TRANSFORMERS

# TRANSFORMERS



**Encoder**
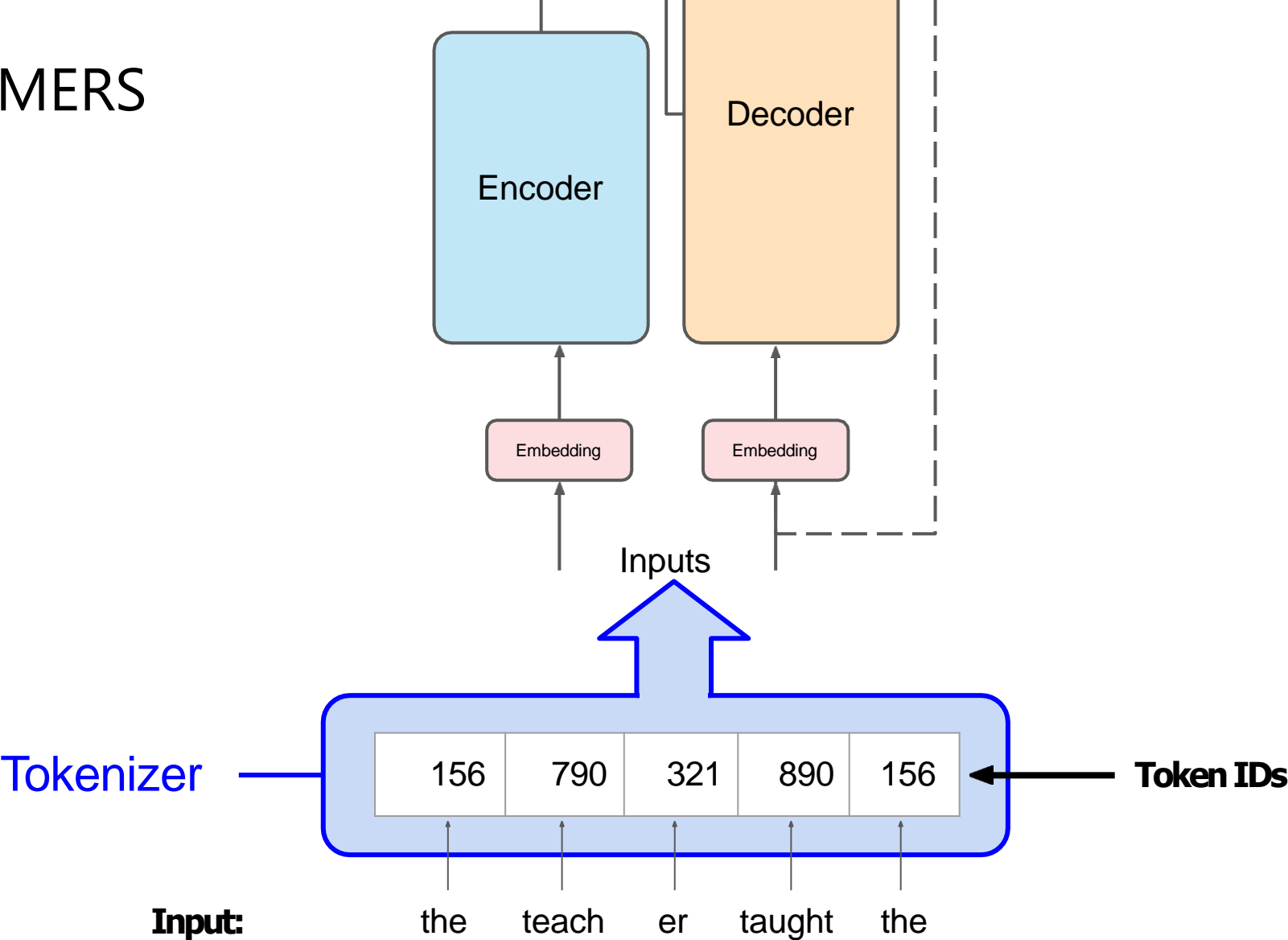Encodes inputs ("prompts") with contextual understanding and produces one vector per input token.
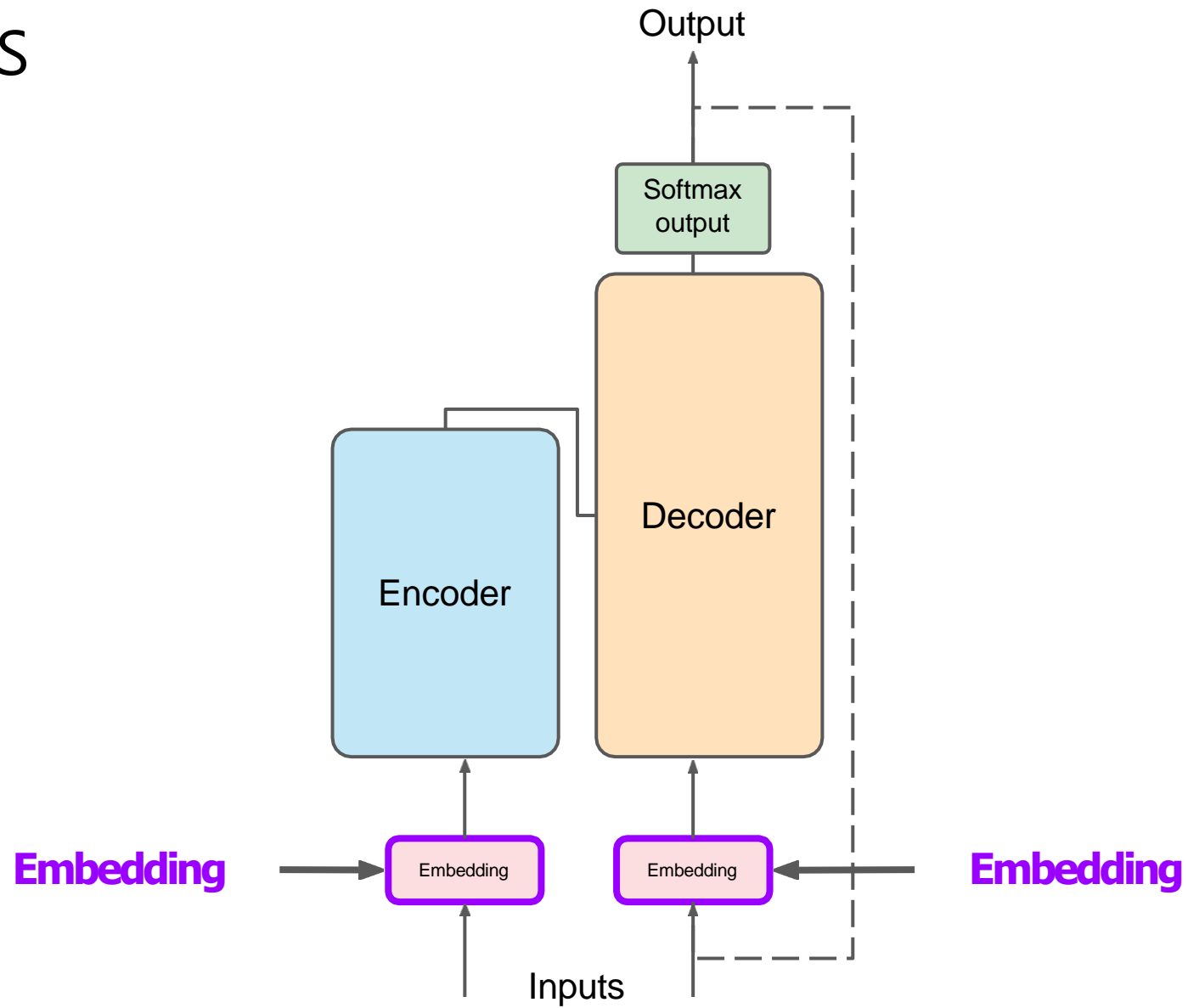
**Decoder**
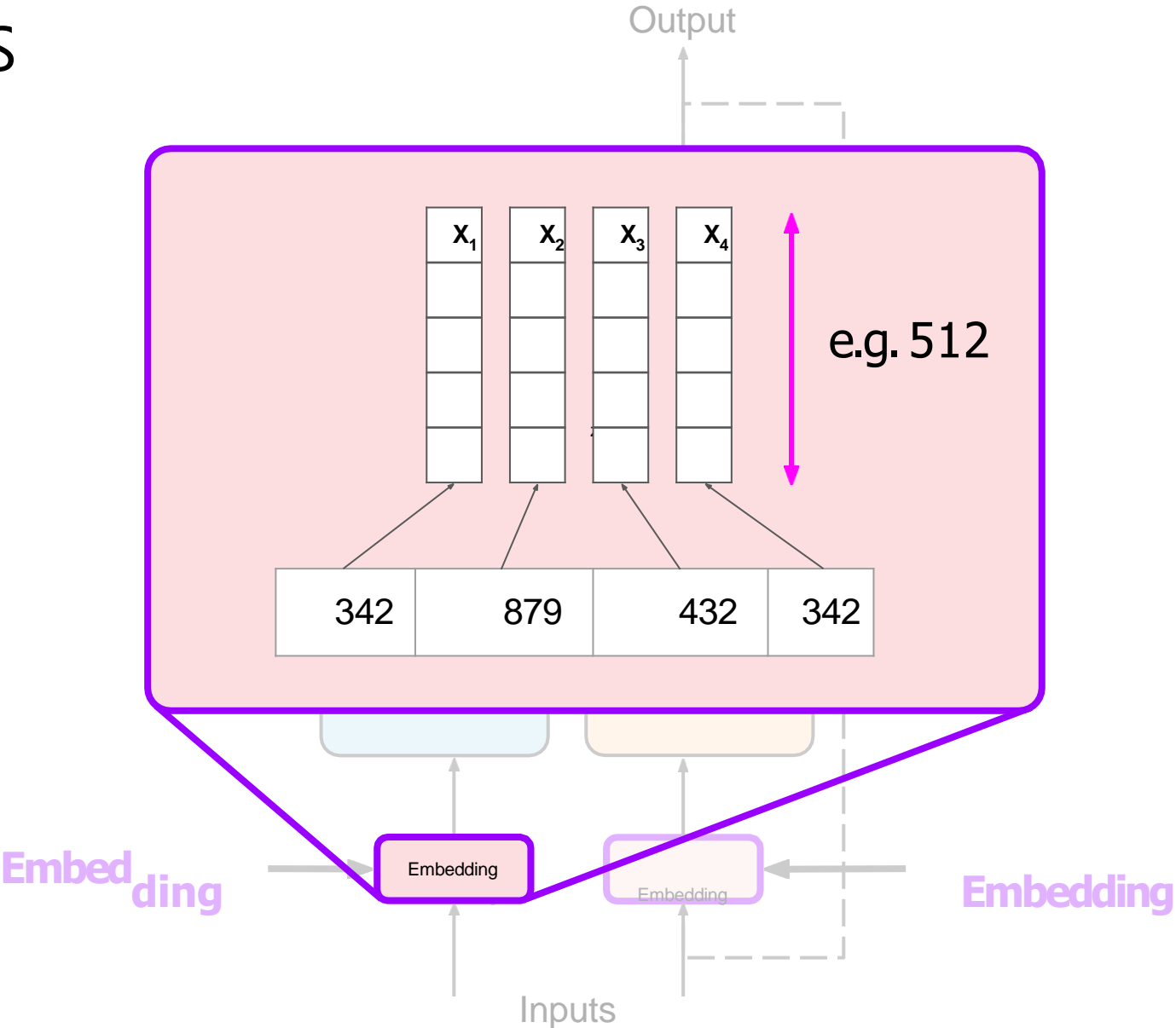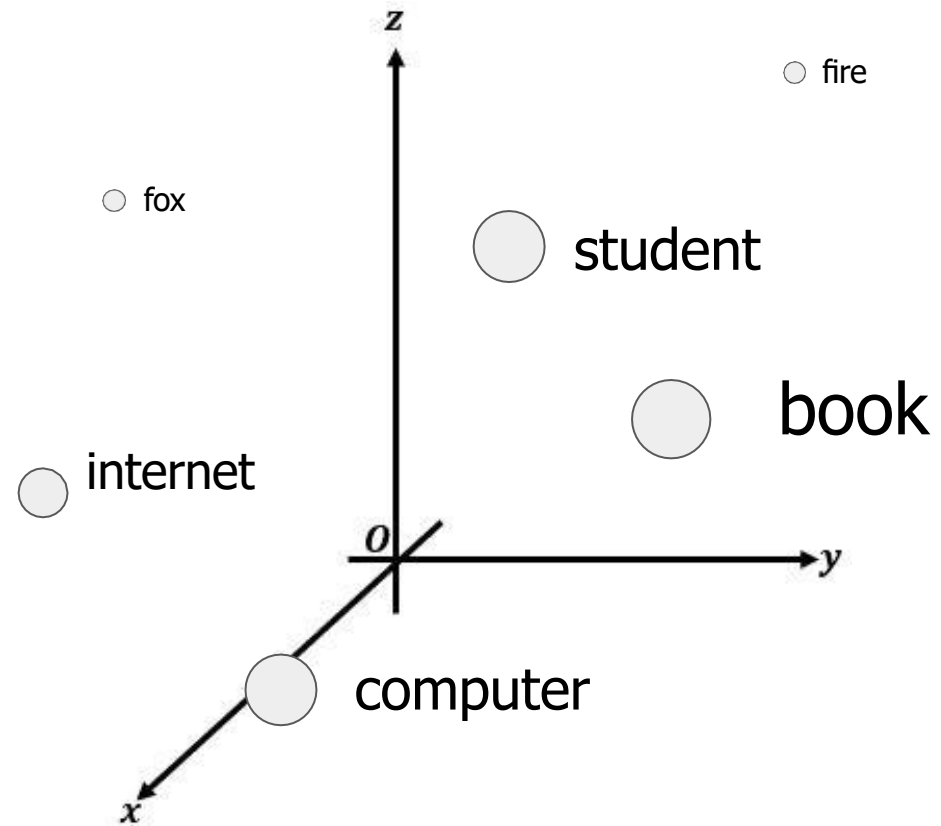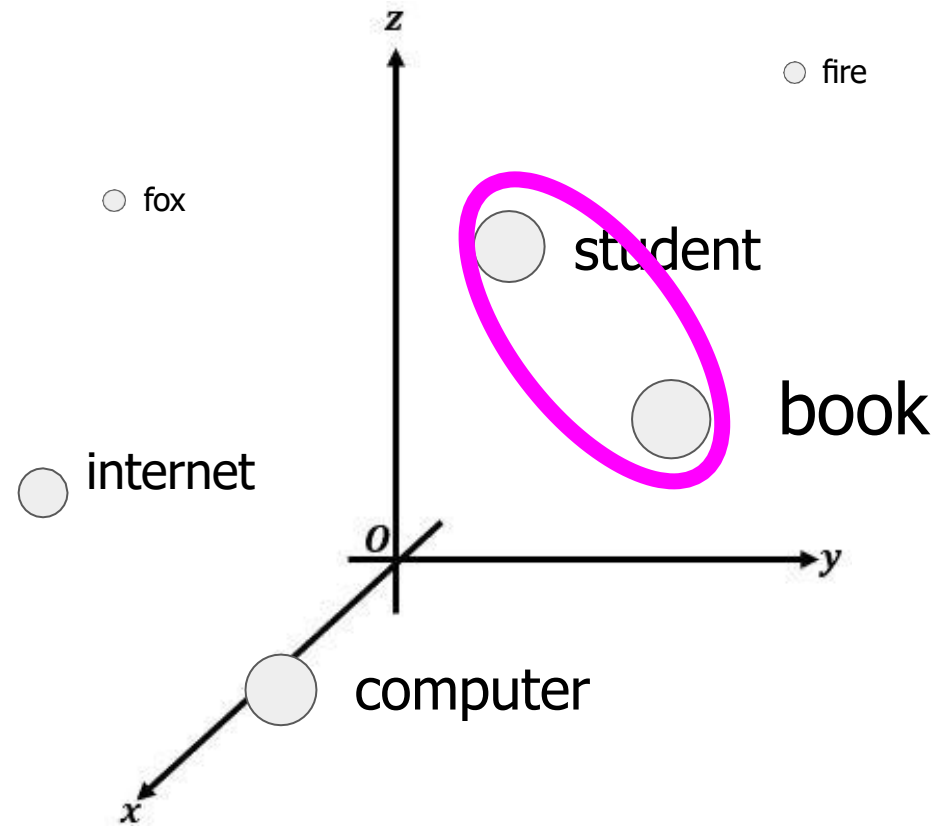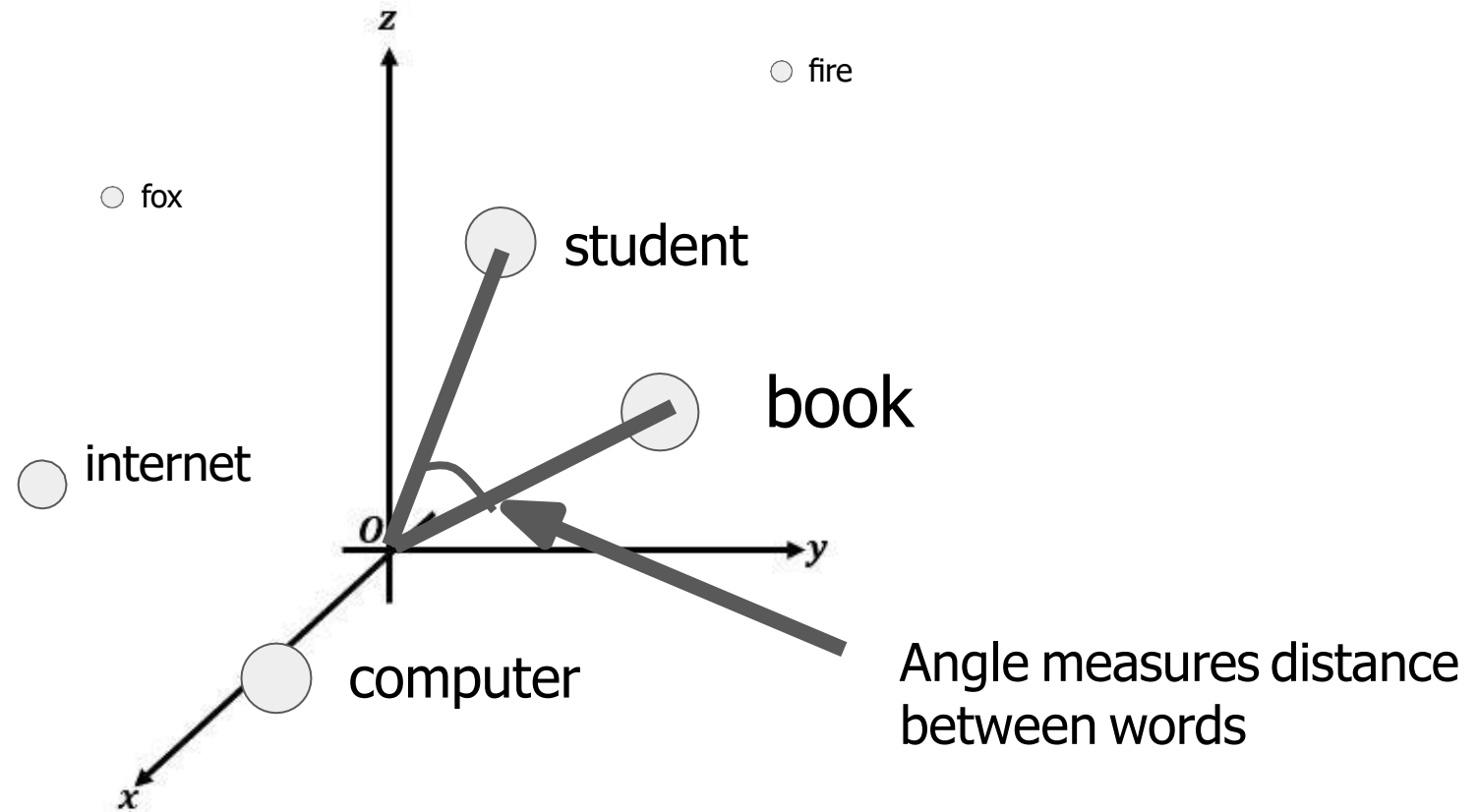Accepts input tokens and generates new tokens.

# TRANSFORMERS



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

**Decoder**
Accepts input tokens and
generates new tokens.

Output

Softmax output

Decoder

Encoder

Embedding

Embedding

Inputs

# TRANSFORMERS



Encoder

Decoder

Embedding          Embedding

Inputs

Tokenizer — | 342 | 879 | 432 | 342 | ← **Token IDs**

**Input:**     the     teacher     taught     the

# TRANSFORMERS

Decoder

Encoder

Embedding          Embedding

Inputs

Tokenizer

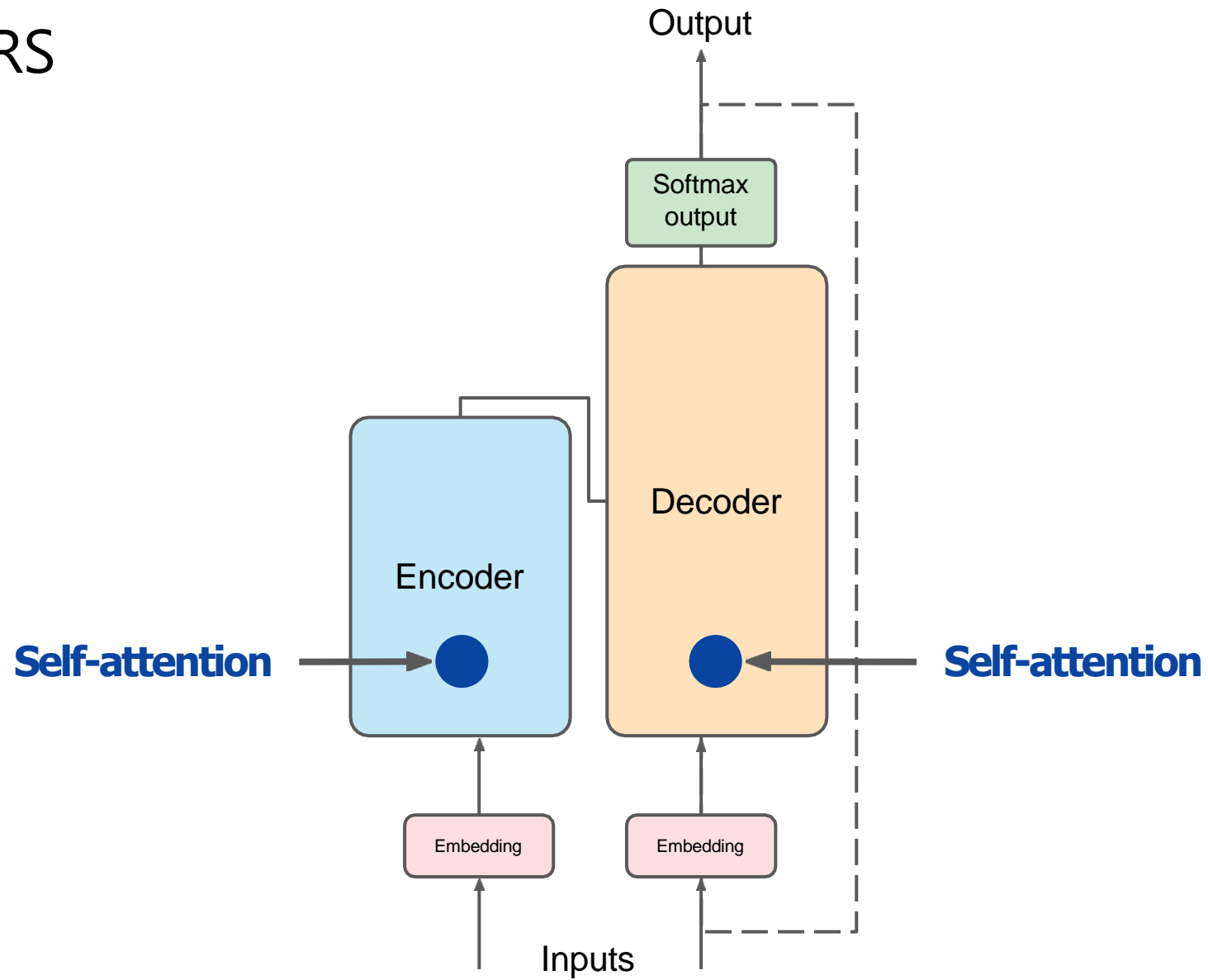| 156 | 790 | 321 | 890 | 156 |

**Token IDs**

**Input:**   the   teach   er   taught   the

# TRANSFORMERS

# TRANSFORMERS



Output

e.g. 512

Embedding

Embedding

Inputs

Demystifying ChatGPT

# TRANSFORMERS

# TRANSFORMERS

z

○ fire

○ fox

student

book

internet

O

y

Angle measures distance
between words

computer

x

# TRANSFORMERS



Output

Softmax output

Decoder

Encoder

**Self-attention**

**Self-attention**
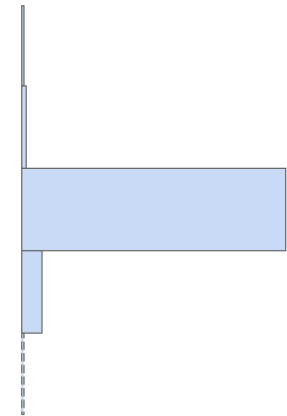
Embedding

Embedding

Inputs

# TRANSFORMERS

# TRANSFORMERS

# TRANSFORMERS



Output

Softmax output

**Softmax output**

### Cooler temperature (e.g <1)

| prob | word |
|------|------|
| 0.001 | apple |
| 0.002 | banana |
| 0.400 | cake |
| 0.012 | donut |
| … | … |

**Strongly peaked probability distribution**

### Higher temperature (>1)

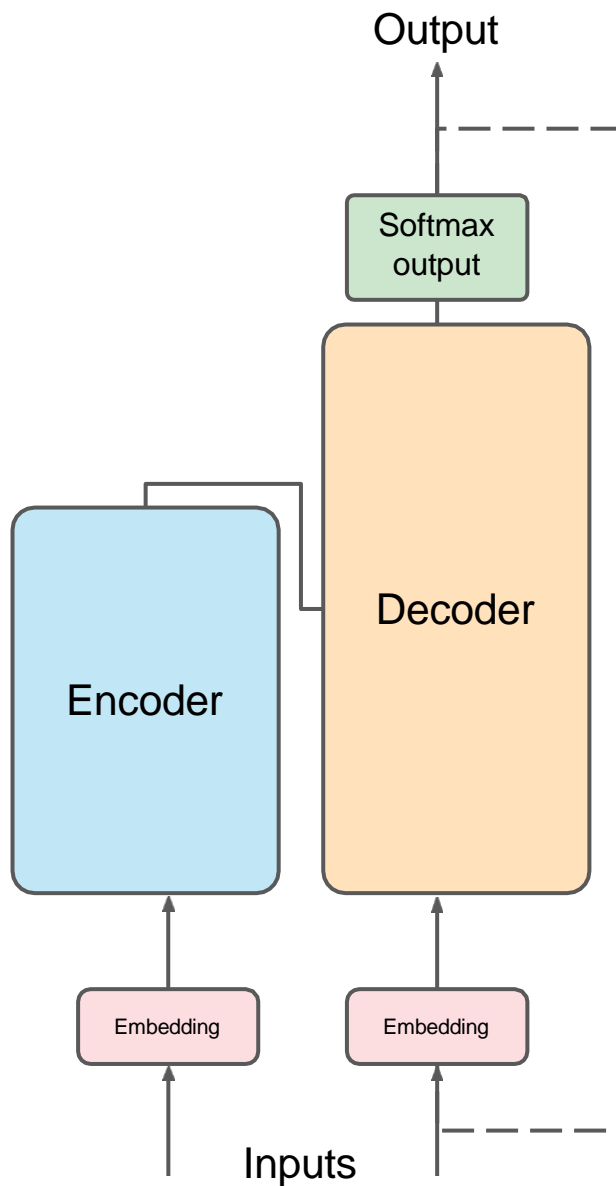| prob | word |
|------|------|
| 0.040 | apple |
| 0.080 | banana |
| 0.150 | cake |
| 0.120 | donut |
| … | … |

**Broader, flatter probability distribution**

# TRANSFORMERS

## Translation:
## sequence-to-sequence task

Milujem seminár DSA



Output

Softmax output

Decoder

Encoder

Embedding

Embedding
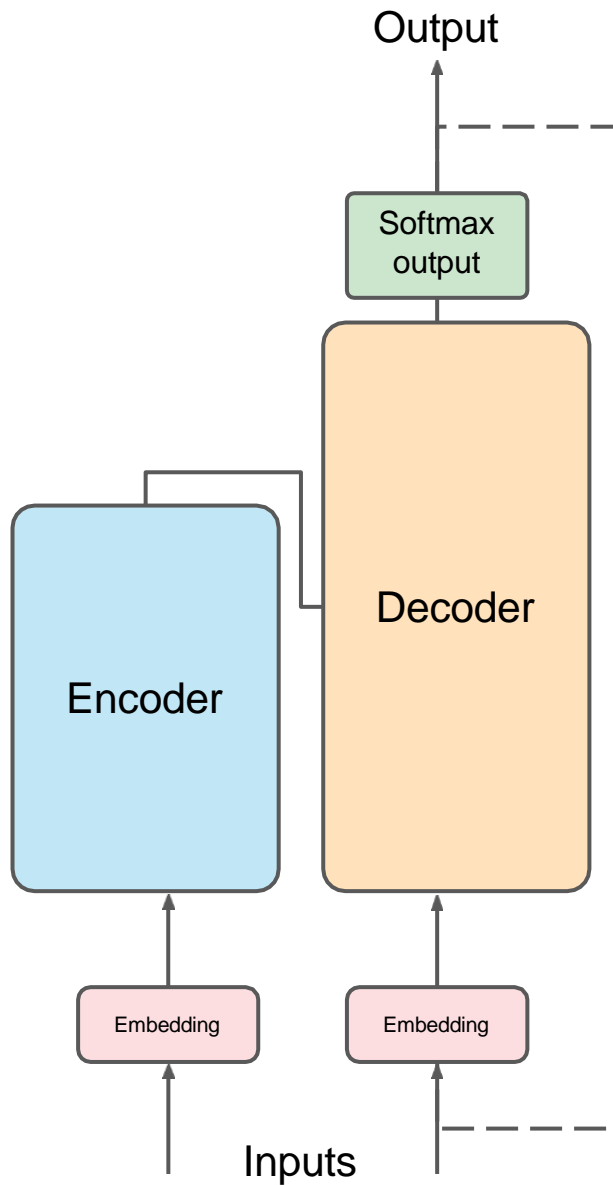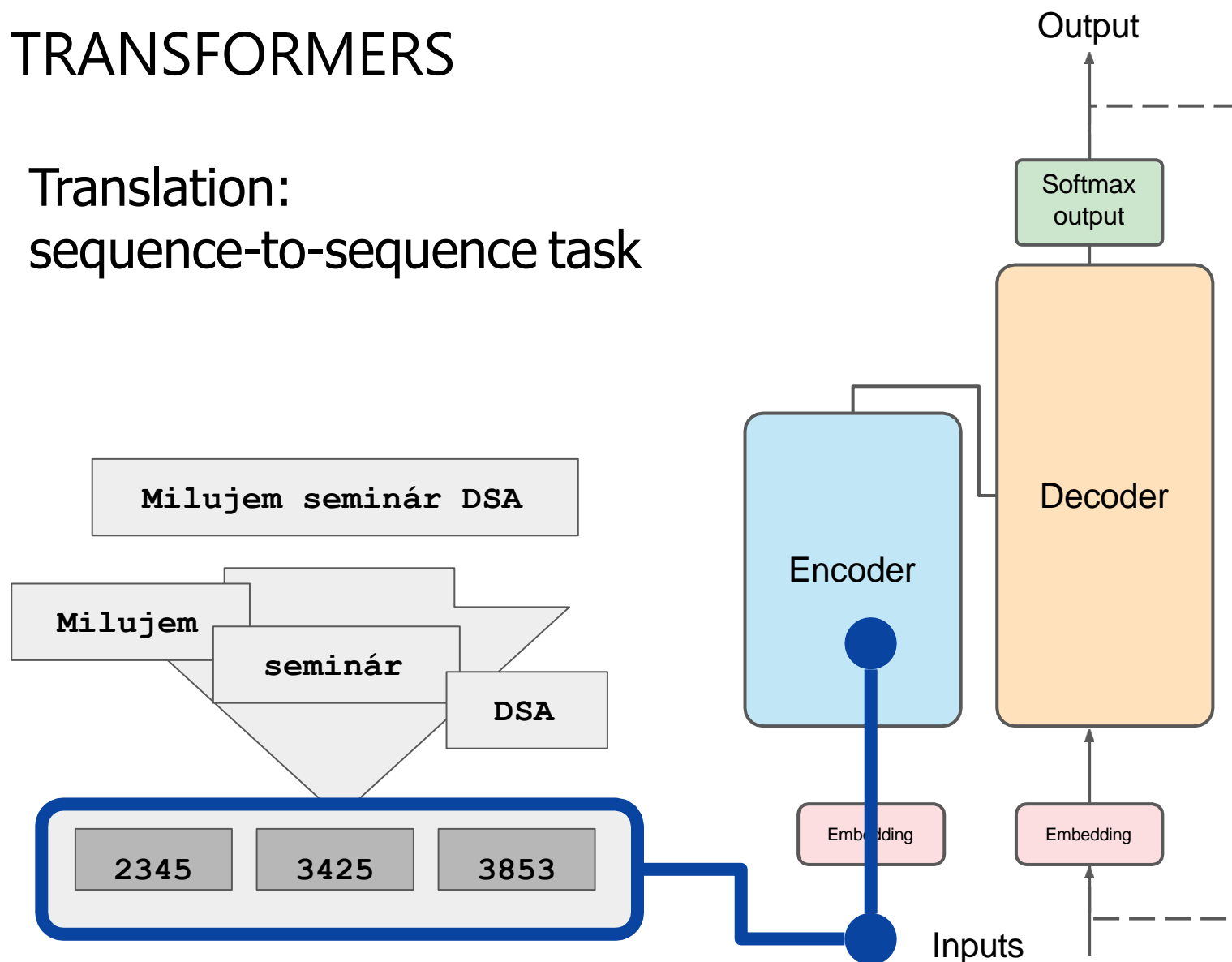
Inputs

# TRANSFORMERS

Translation:
sequence-to-sequence task

Milujem seminár DSA

Milujem

seminár

DSA

| 2345 | 3425 | 3853 |

Output

Softmax
output

Encoder

Decoder

Embedding

Embedding
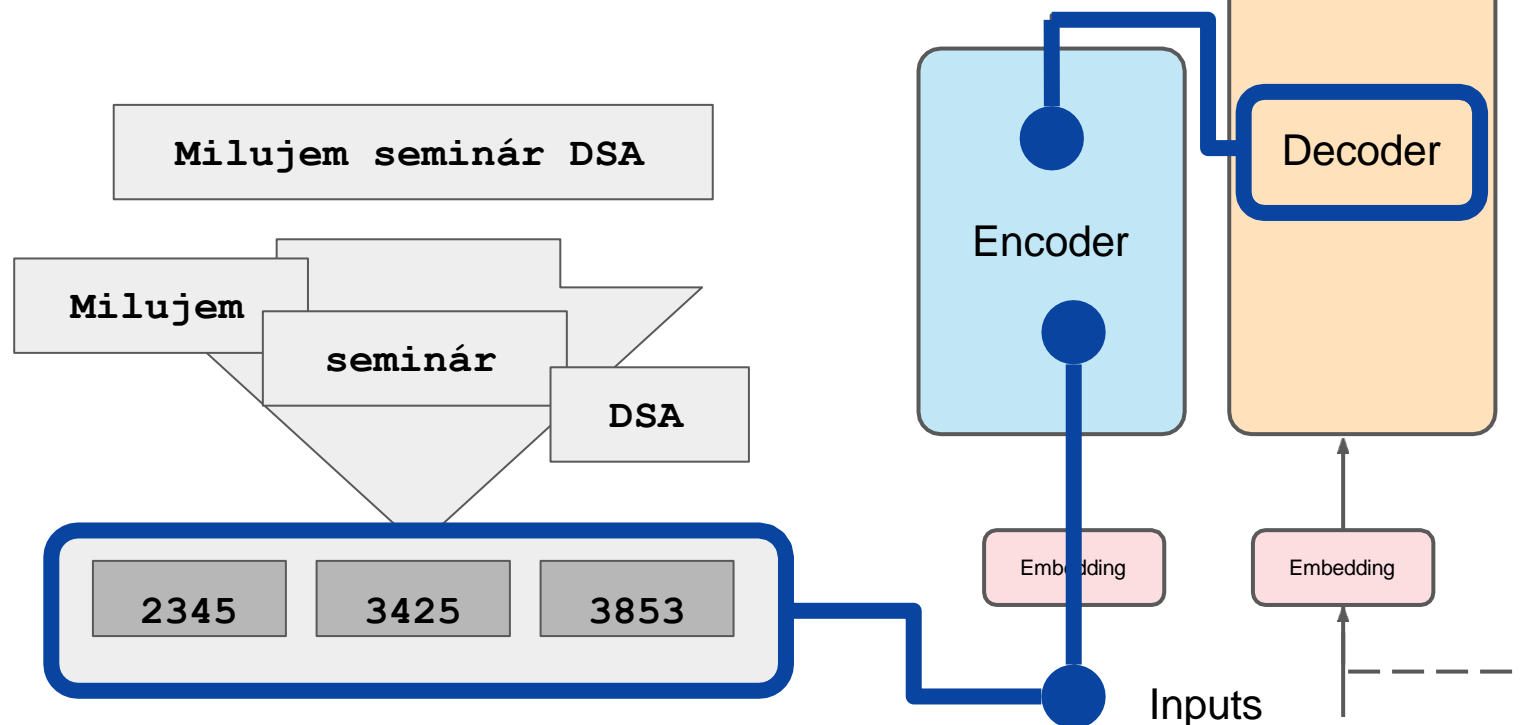
Inputs

# TRANSFORMERS

Translation:
sequence-to-sequence task

Milujem seminár DSA

Milujem

seminár

DSA

2345    3425    3853

Output

Softmax
output

Decoder

Encoder

Embedding

Embedding

Inputs

# TRANSFORMERS

Translation:
sequence-to-sequence task

Output

Softmax
output

Decoder

Encoder

**Milujem seminár DSA**

**Milujem**

**seminár**

**DSA**

| 2345 | 3425 | 3853 |

Embedding
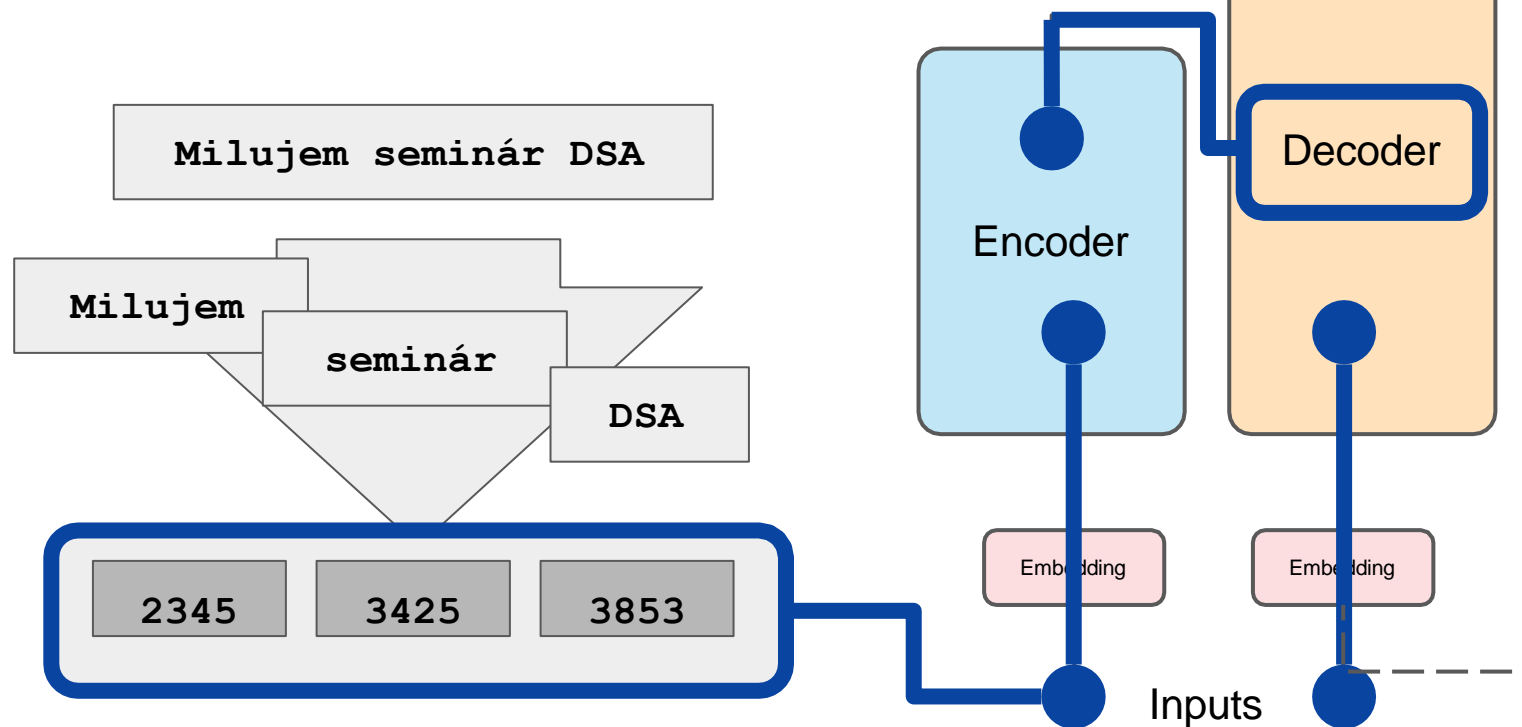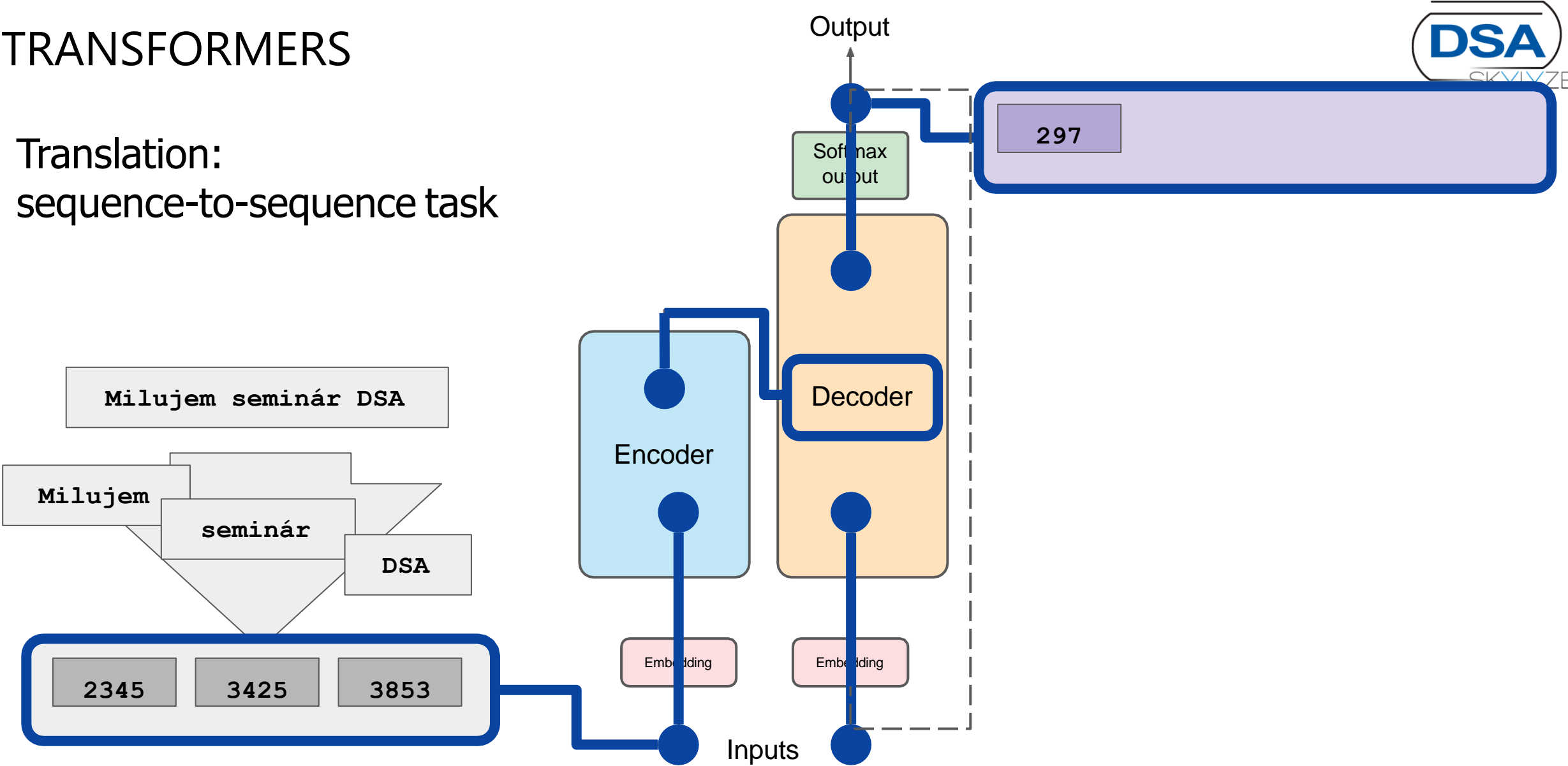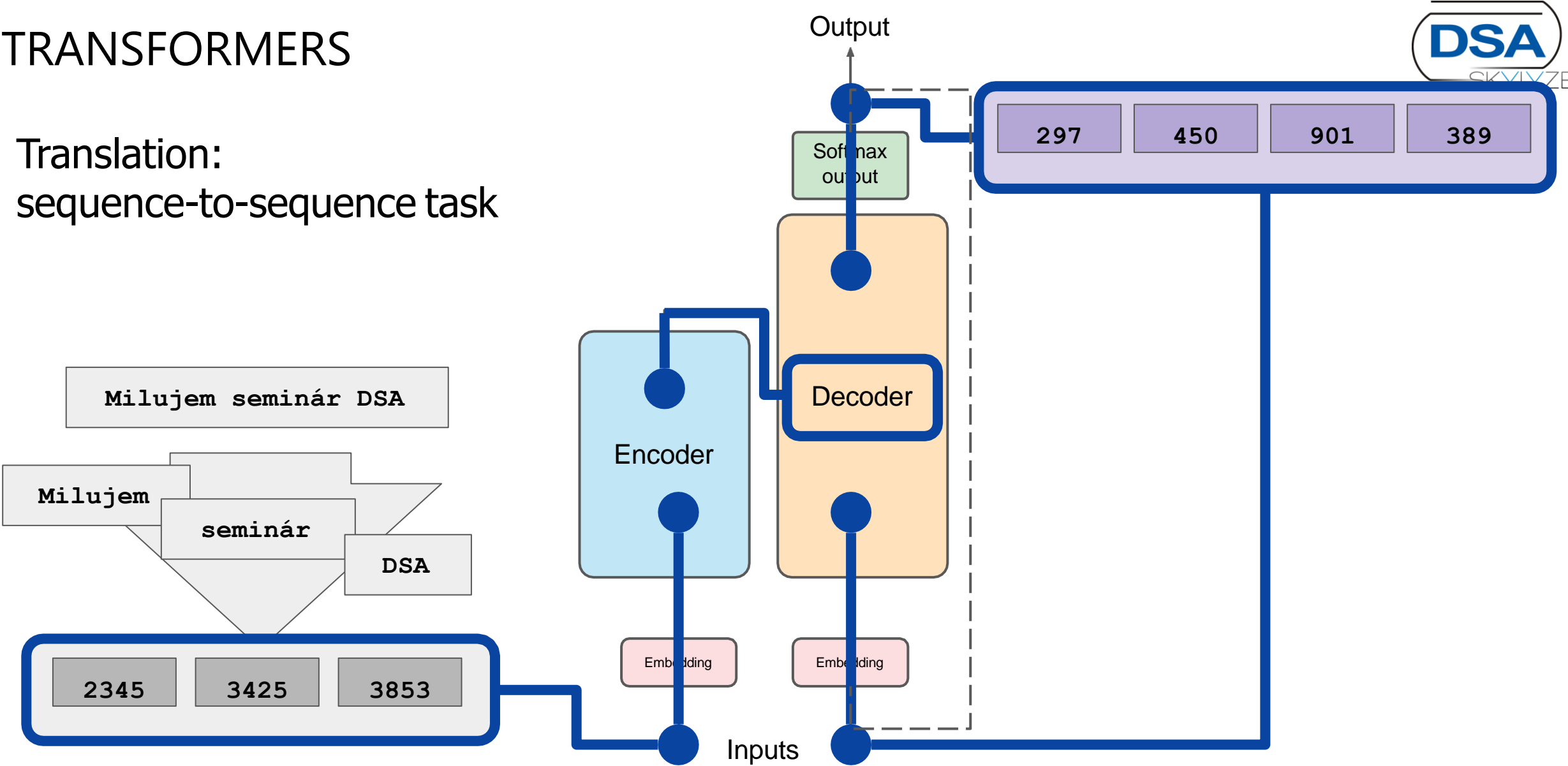
Embedding

Inputs

# TRANSFORMERS



Translation:
sequence-to-sequence task

# TRANSFORMERS

## Translation: sequence-to-sequence task

# TRANSFORMERS

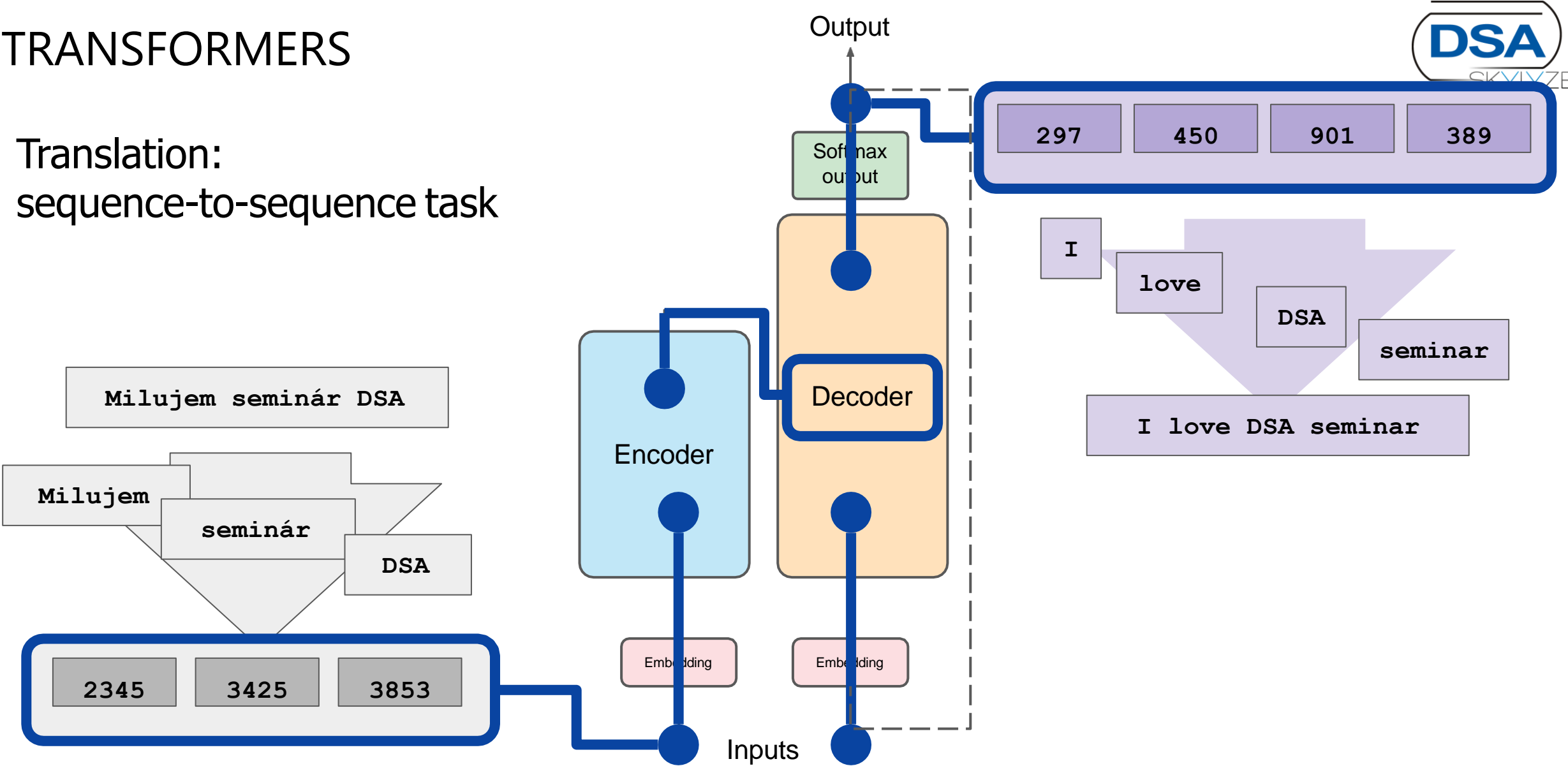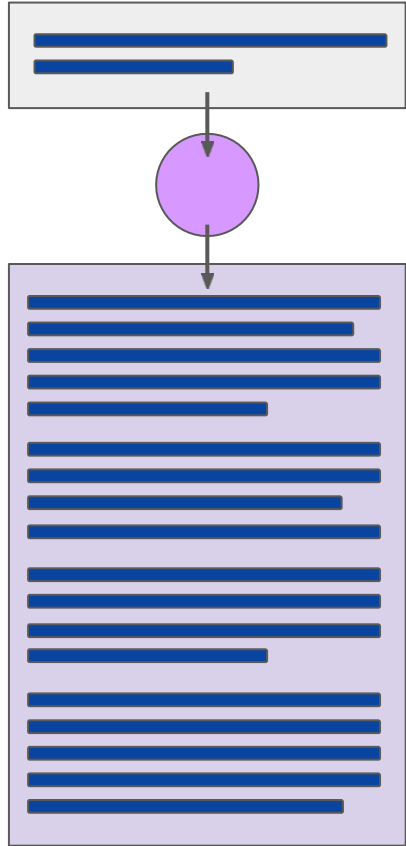Translation:
sequence-to-sequence task

Output

Milujem seminár DSA

Milujem

seminár

DSA

| 2345 | 3425 | 3853 |

| 297 | 450 | 901 | 389 |

Softmax output

Decoder

Encoder

Embedding

Embedding

Inputs

# TRANSFORMERS

Translation:
sequence-to-sequence task
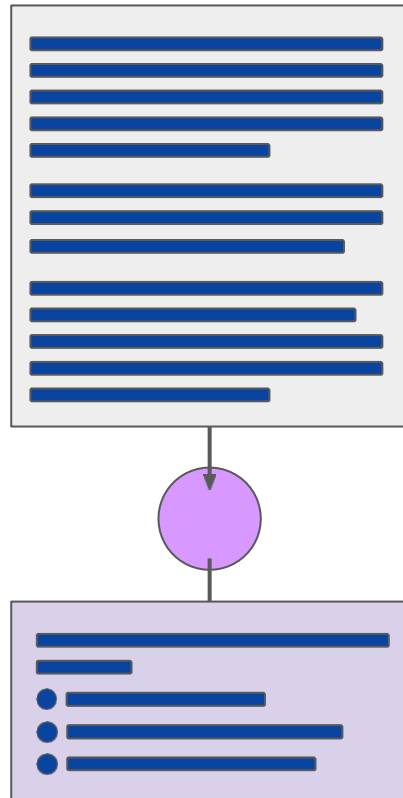
# LLM USE CASES & TASKS
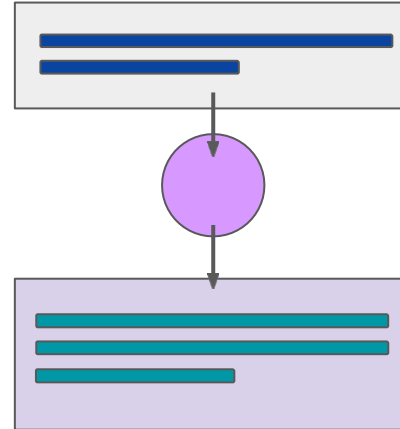
Essay Writing

Summarization
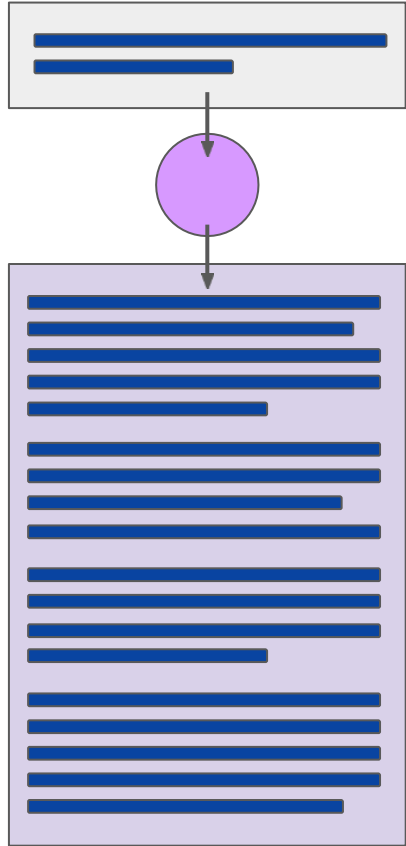
Translation

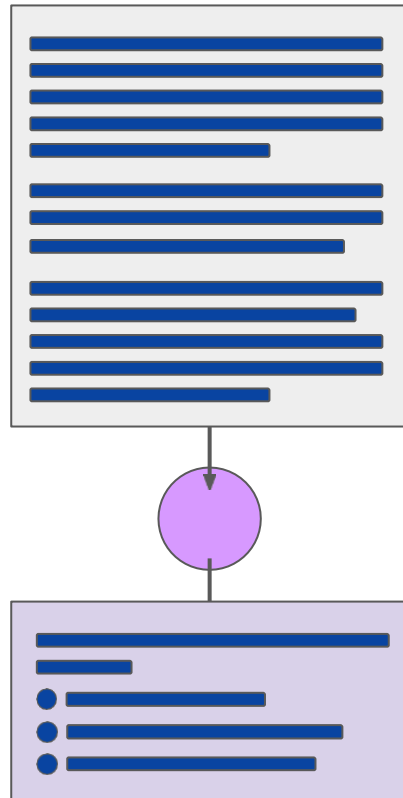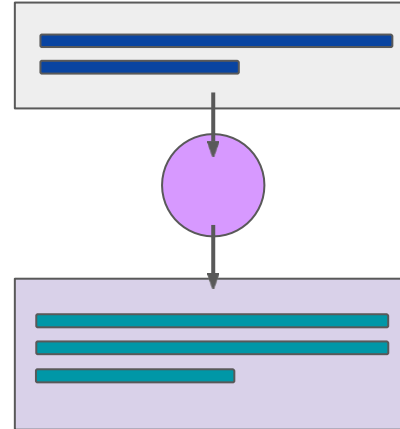# LLM USE CASES & TASKS

Essay Writing

Summarization

Translation

Information retrieval

Invoke APIs and actions

Action call

External Applications

# 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 |

# HANDS-ON

# 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 |

# CHALLENGES OF TRAINING LLMS

Infrastructure and costs

- Trained on massive corpus of text data

- Billions to trillions of parameter
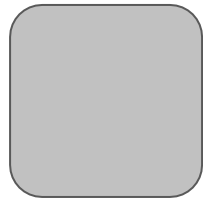
- GPU-Hardware/Infrastructure needed

- Approximate GPU RAM needed to store 1B parameters
  - 1 parameter = 4 bytes (32-bit float)
  - 1B parameters = 4 x 1e-09 bytes = 4GB

- Training needs ~20 extra bytes per parameter

# APPROXIMATE GPU RAM NEEDED TO TRAIN 1B-PARAMS

Memory needed to store model

Memory needed to train model

**4GB @ 32-bit
full precision**

**80GB @ 32-bit
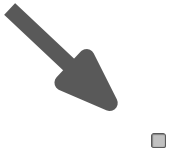full precision**

# GPU RAM NEEDED TO TRAIN LARGER MODELS

**1B param model**

**175B param model**

(GPT-3)

**500B param model**

(PaML)

**40,000 GB @ 32-bit full precision**

**14,000 GB @ 32-bit full precision**

# CHALLENGES OF TRAINING LLMS

Infrastructure and costs

- It would take 288 years to train GPT-3 (175B) on a single NVIDIA Tesla V100 GPU.
  - 1Y = 8760h
  - 275W * 8760h = 2409000Wh = 2490kWh
  - 2 Person-Household : 2000 – 2800kWh

- Researchers calculated that OpenAI could have trained GPT-3 in as little as 34 days on 1,024x A100 GPUs

- It is estimated that GPT-3 cost around $4.6 million dollars to train from scratch

- Cost of this workshop ~90$ (1$/min)



- **Jun 2017**
- **16GB RAM**
- **250-300 Watt**

- **Jun 2021**
- **80GB RAM**
- **300-400 Watt**

# GENERATIVE AI

Demystifying ChatGPT

# GENERATIVE AI

DSA
SKYLYZE

| NLP | Computer Vision | Audio | Multimodal |
|---|---|---|---|
| Text Generation | Image Classification | Text-to-Speech | Feature Extraction |
| Text-to-Text | Object Detection | Audio Classification | Text-to-Image |
| Translation | Image Segmentation | Speech recognition | **Stable Diffusion** |
| Summarization | Depth Estimation | Audio-to-Audio | |
| Text Classification | Image-to-Image | | Document Question Answering |
| Question Answering | | | |

Demystifying ChatGPT

# STABLE DIFFUSION

Demystifying ChatGPT

# ACKNOWLEDGEMENT AND COPYRIGHT NOTICE