**UET**
Since 2004
**ĐẠI HỌC CÔNG NGHỆ, ĐHQGHN**
VNU-University of Engineering and Technology

**VNU**
Since 1906
**ĐẠI HỌC QUỐC GIA HÀ NỘI**
Vietnam National University, Hanoi

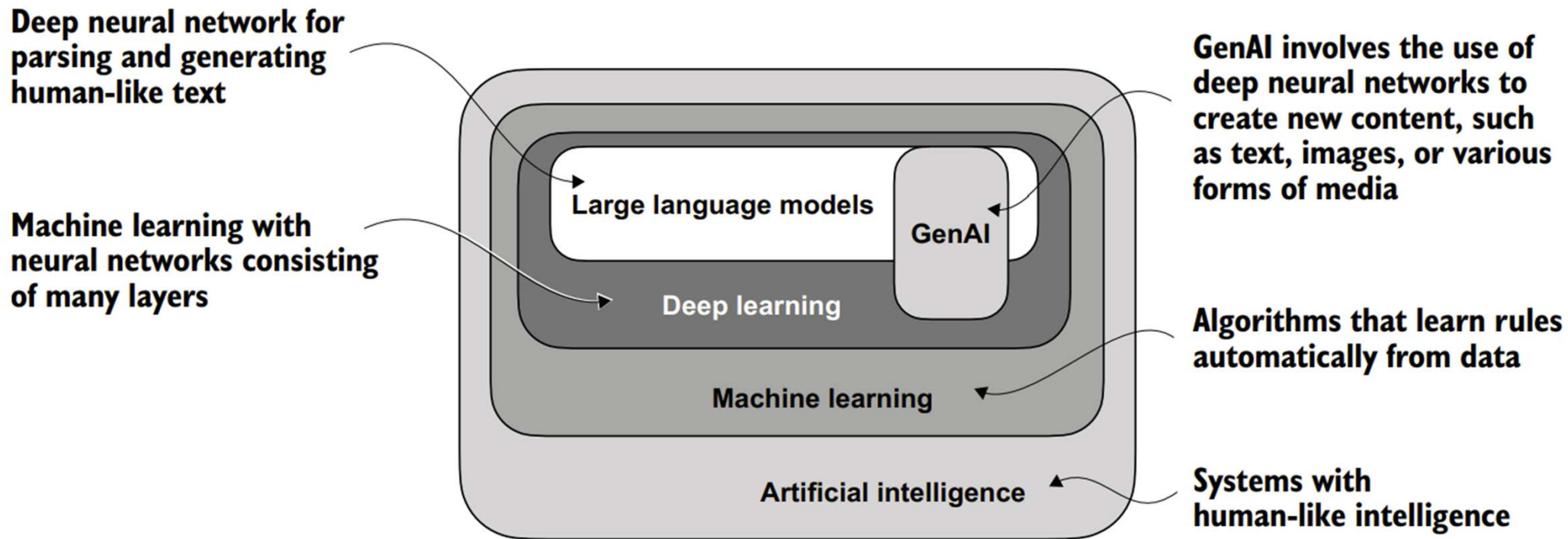# INT3121 – Topics in Computer Science
## Lecture 1: Foundation of LLMs

**Hanoi, 09/2025**

# Course Introduction

- Introduction to Large Language Models (LLMs)

- Word embedding

- Tokenization

- Self-attention

- Muli-headed self-attention

- Vision Transformer

# From AI to Large Language Models

**Deep neural network for parsing and generating human-like text**

**Machine learning with neural networks consisting of many layers**

Large language models

GenAI

Deep learning

Machine learning

Artificial intelligence

**GenAI involves the use of deep neural networks to create new content, such as text, images, or various forms of media**

**Algorithms that learn rules automatically from data**

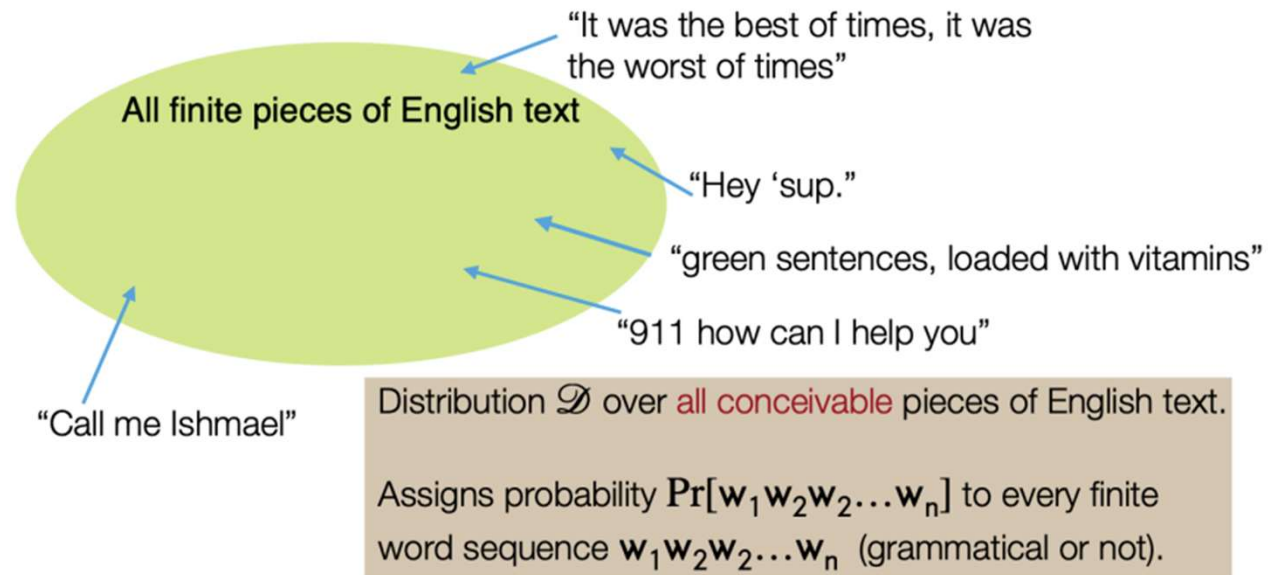**Systems with human-like intelligence**

Source: COS 597G

# Application of LLMs

- Chatbot (ChatGPT, Gemini,…)
- Machine translation
- Generation of novel texts
- Sentiment analysis
- Text summarization
- Language-instructed image generation
- Domain-specific LLM-based applications
  - Code generation (Copilot)
  - Task/project magagement (Notion)
  - Multimedia generation/editing (Canva)

Source: BLM

# Language Models: Narrow Sense

- A probabilistic model that assigns a probability $P[w_1, w_2, \ldots, w_n]$ to every finite sequence $w_1, \ldots, w_n$ (grammatical or not)

All finite pieces of English text

"It was the best of times, it was the worst of times"

"Hey 'sup."

"green sentences, loaded with vitamins"

"911 how can I help you"

"Call me Ishmael"

Distribution $\mathcal{D}$ over all conceivable pieces of English text.

Assigns probability $\Pr[w_1 w_2 w_2 \ldots w_n]$ to every finite word sequence $w_1 w_2 w_2 \ldots w_n$ (grammatical or not).

Source: COS 324

# Language Models: Narrow Sense

Conditional probability

$$p(w_1, w_2, w_3, \ldots, w_N) =$$
$$p(w_1)\, p(w_2|w_1)\, p(w_3|w_1, w_2) \times \ldots \times p(w_N|w_1, w_2, \ldots w_{N-1})$$
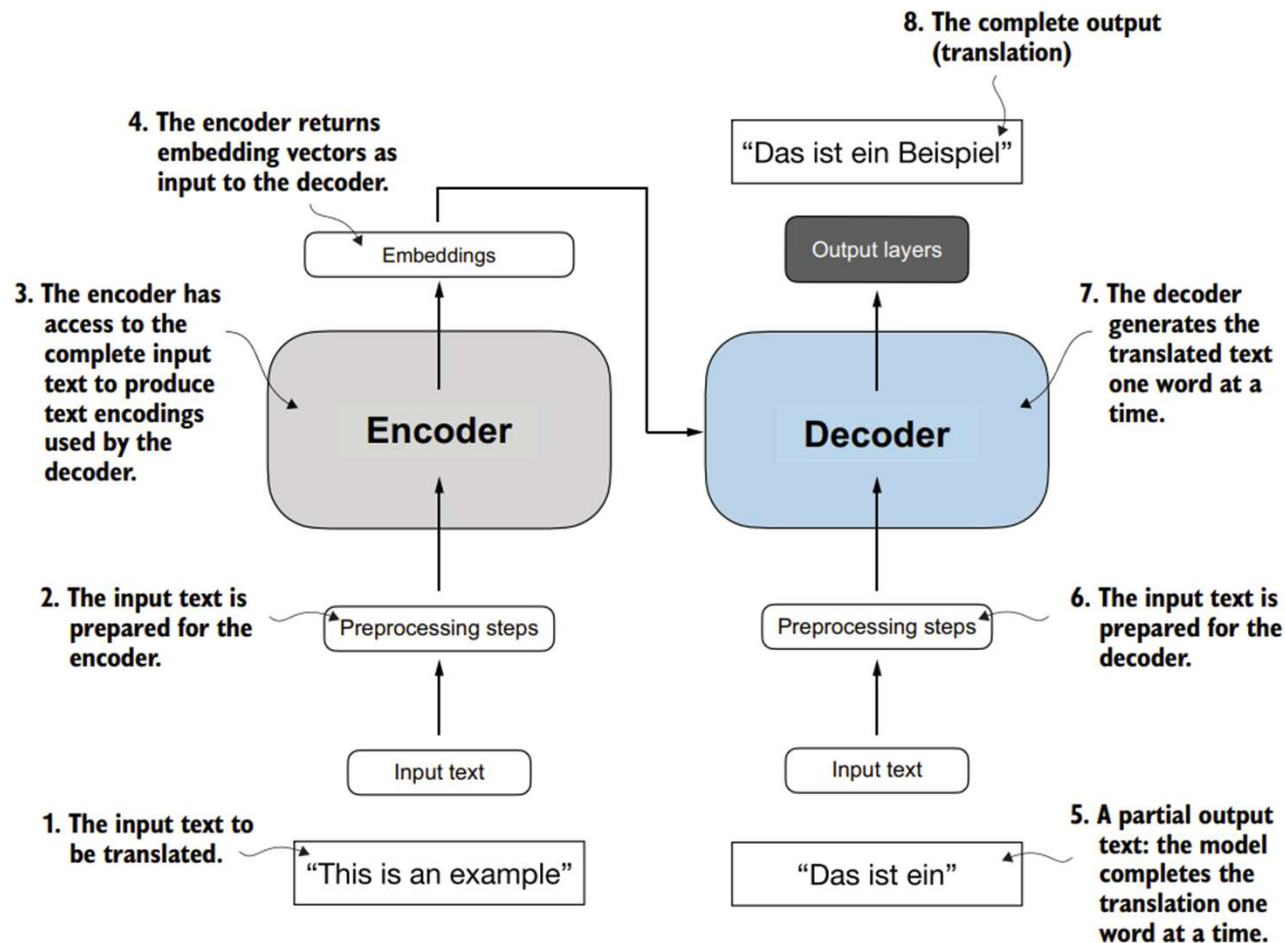
Sentence: "the cat sat on the mat"

$$P(\text{the cat sat on the mat}) = P(\text{the}) * P(\text{cat}|\text{the}) * P(\text{sat}|\text{the cat})$$
$$* P(\text{on}|\text{the cat sat}) * P(\text{the}|\text{the cat sat on})$$
$$* P(\text{mat}|\text{the cat sat on the})$$
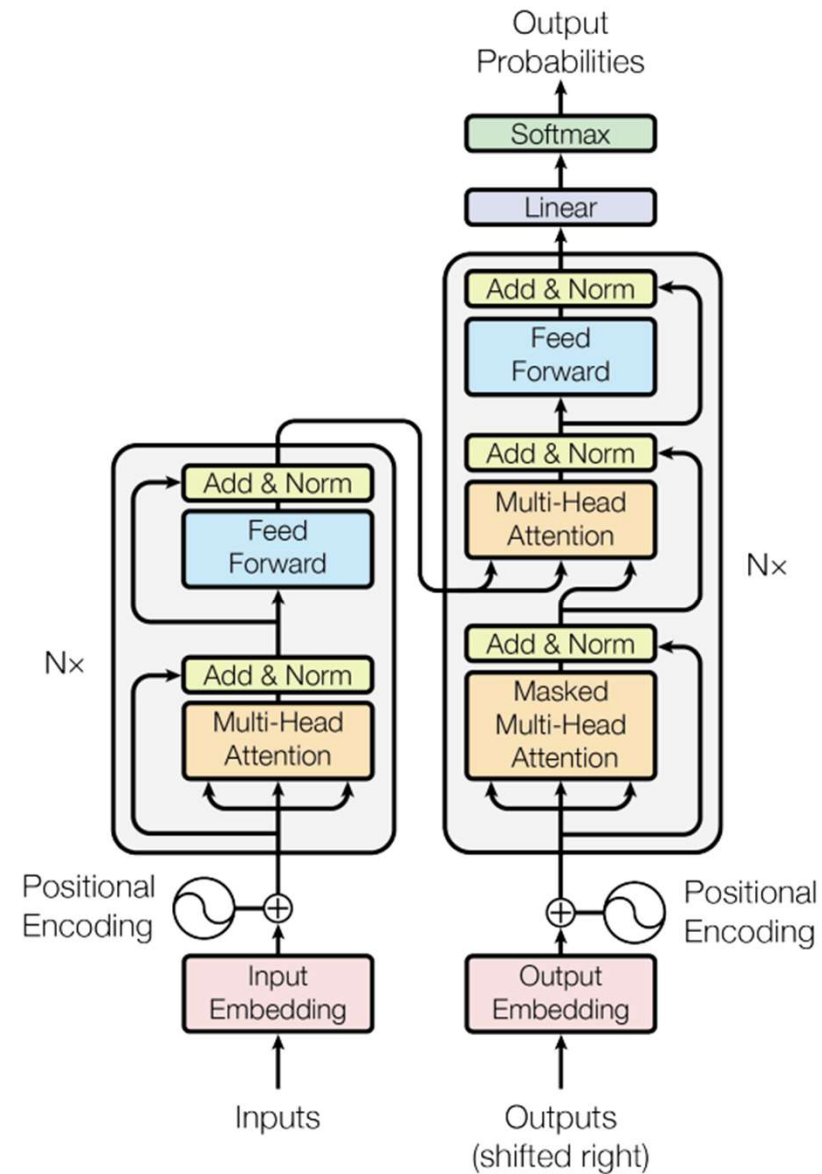
Implicit order

Source: COS 484

GPT-3 still acts in this way but the model is implemented as a very large neural network of 175-billion parameters!
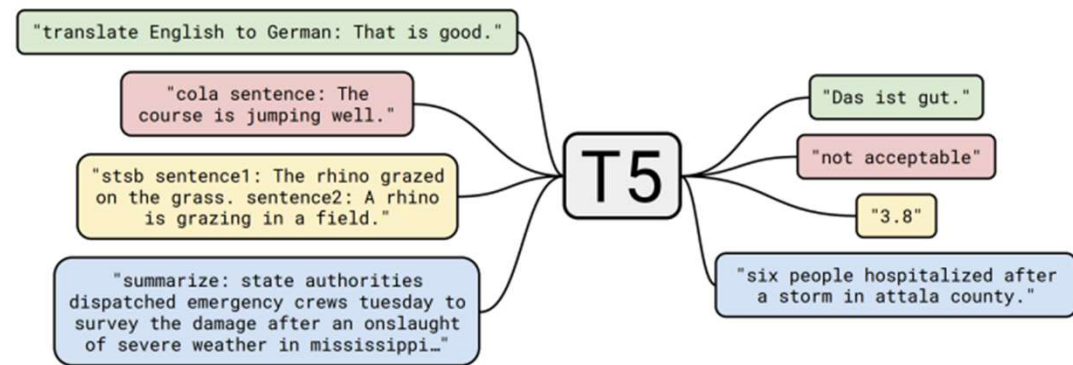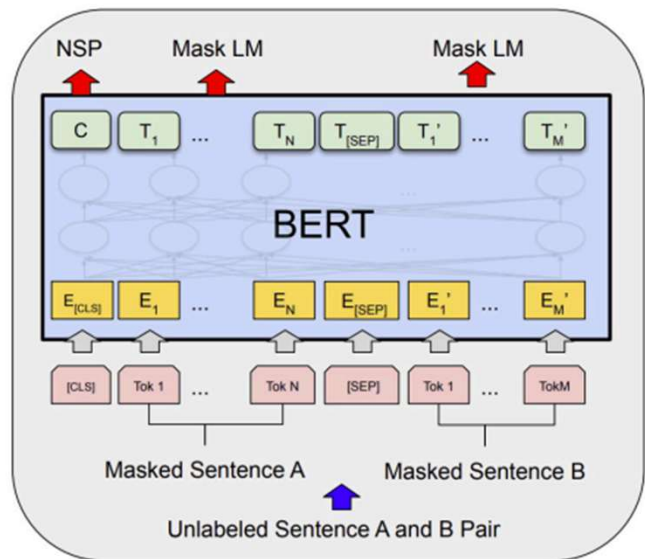
# The Transformer Architecture: Overview



8. The complete output (translation)

4. The encoder returns embedding vectors as input to the decoder.

3. The encoder has access to the complete input text to produce text encodings used by the decoder.

"Das ist ein Beispiel"

Embeddings

Output layers

Encoder

Decoder

7. The decoder generates the translated text one word at a time.

2. The input text is prepared for the encoder.

Preprocessing steps

Preprocessing steps

6. The input text is prepared for the decoder.

Input text

Input text

1. The input text to be translated.

"This is an example"

"Das ist ein"

5. A partial output text: the model completes the translation one word at a time.

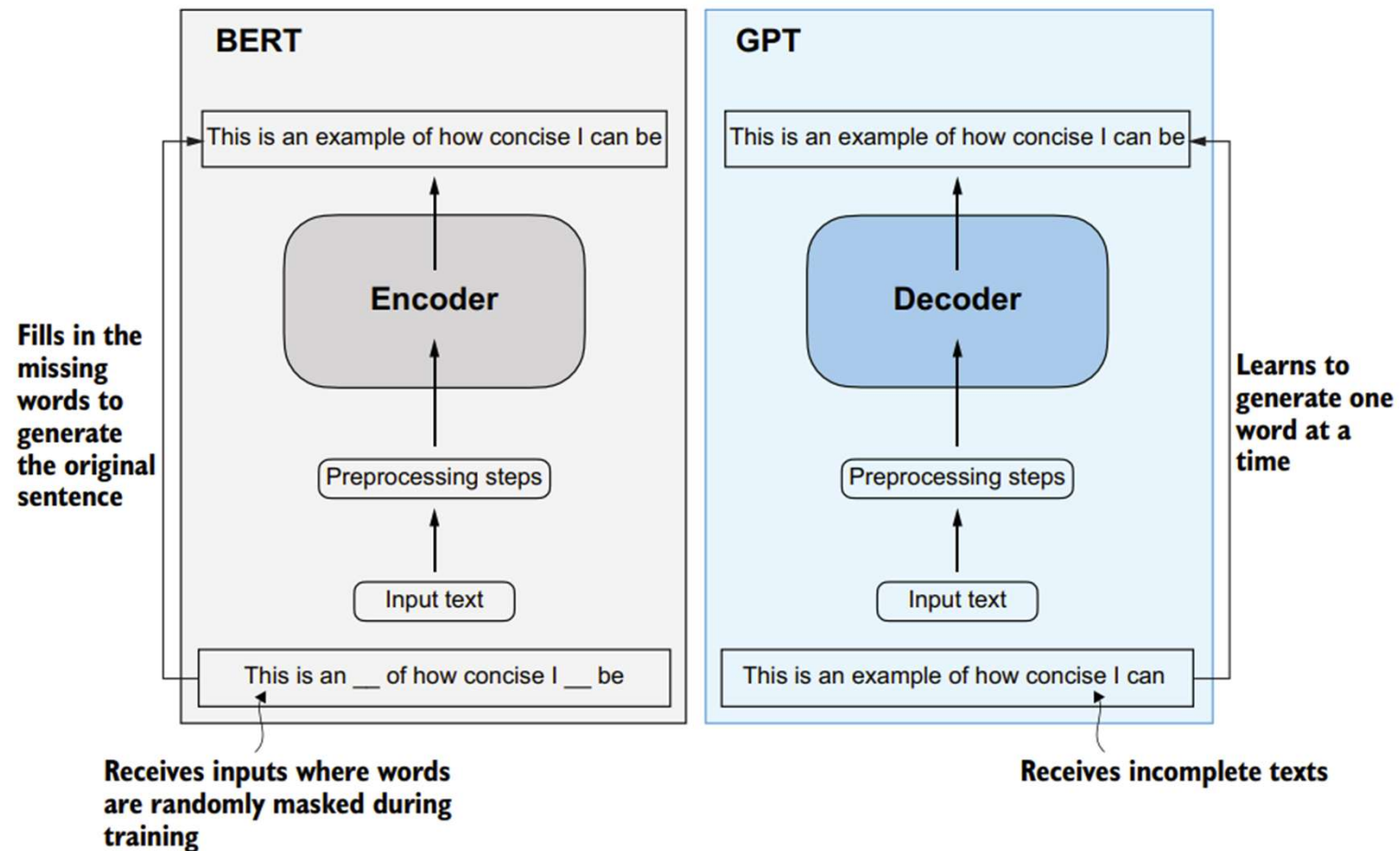# The Transformer Architecture: Detail

# Language Models: Broad Sense

- Decoder-only models (GPT-x models)
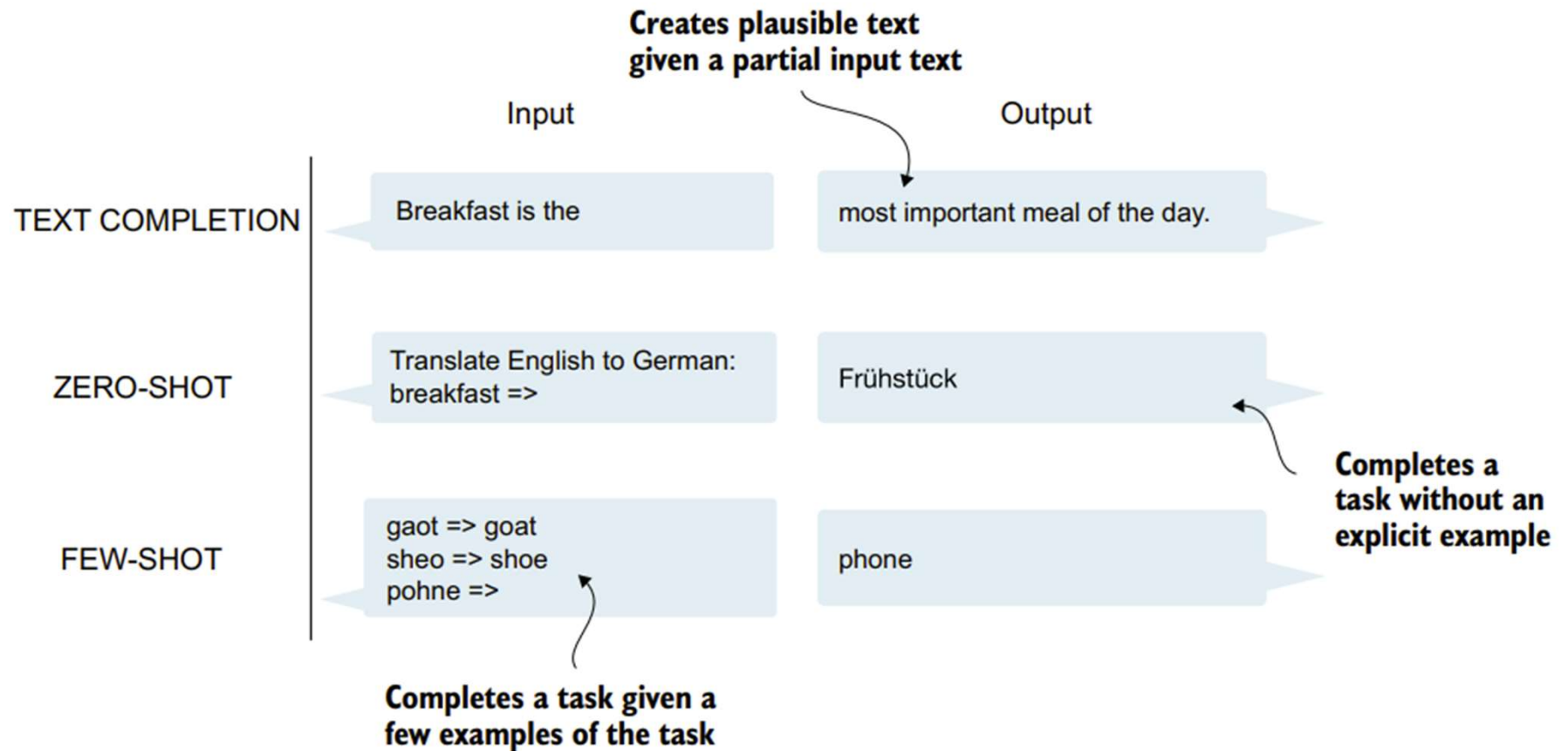- Encoder-only models (BERT, RoBERTa, ELECTRA)
- Encoder-decoder models (T5, BART)



Source: COS 597G

# GPT vs BERT



**BERT**

This is an example of how concise I can be

**Encoder**

Preprocessing steps

Input text

This is an __ of how concise I __ be

**Fills in the missing words to generate the original sentence**

**Receives inputs where words are randomly masked during training**

**GPT**

This is an example of how concise I can be

**Decoder**

Preprocessing steps

Input text

This is an example of how concise I can

**Learns to generate one word at a time**

**Receives incomplete texts**

# GPT-based Model Capabilities



**Creates plausible text given a partial input text**

| | Input | Output |
|---|---|---|
| TEXT COMPLETION | Breakfast is the | most important meal of the day. |
| ZERO-SHOT | Translate English to German: breakfast => | Frühstück |
| FEW-SHOT | gaot => goat<br>sheo => shoe<br>pohne => | phone |

**Completes a task without an explicit example**

**Completes a task given a few examples of the task**

# Stages of Building and Using LLMs

An LLM is pretrained on unlabeled text data.

The LLM has a few basic capabilities after pretraining.

- Text completion
- Few-shot capabilities

- Internet texts
- Books
- Wikipedia
- Research articles

Raw, unlabeled text (trillions of words)

Train

Pretrained LLM (foundation model)

Train

- Classification
- Summarization
- Translation
- Personal assistant
- ...

Fine-tuned LLM

A pretrained LLM can be further trained on a labeled dataset to obtain a fine-tuned LLM for specific tasks.
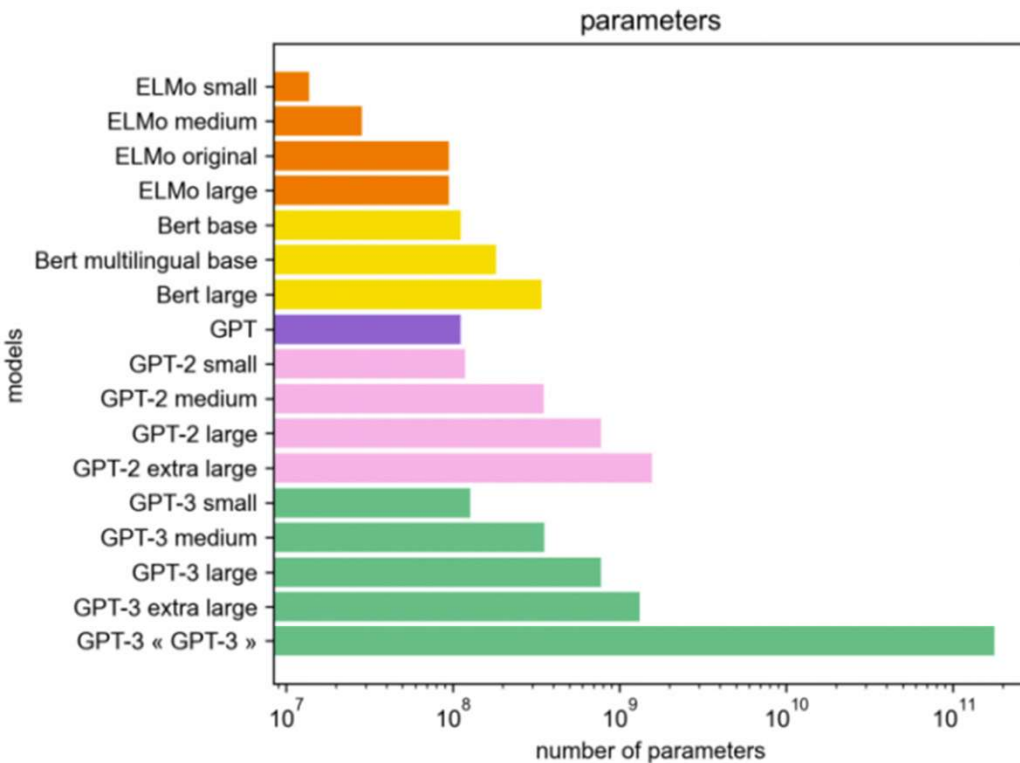
Labeled dataset

# GPT 3 Training Dataset

Table 1.1   The pretraining dataset of the popular GPT-3 LLM

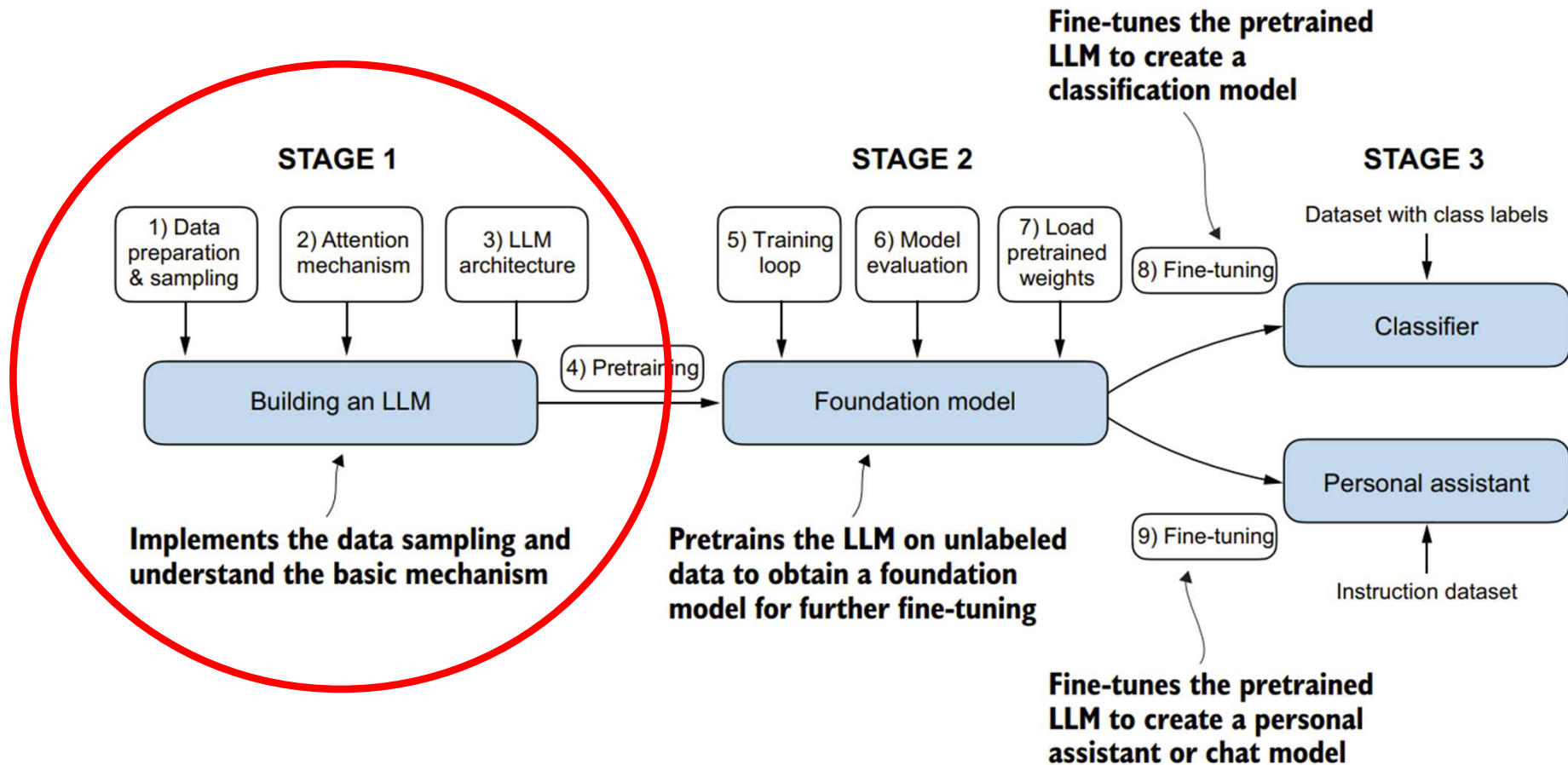| Dataset name | Dataset description | Number of tokens | Proportion in training data |
|---|---|---|---|
| CommonCrawl (filtered) | Web crawl data | 410 billion | 60% |
| WebText2 | Web crawl data | 19 billion | 22% |
| Books1 | Internet-based book corpus | 12 billion | 8% |
| Books2 | Internet-based book corpus | 55 billion | 8% |
| Wikipedia | High-quality text | 3 billion | 3% |

**GPT-3 has 175 billion parameters**

# LLMs Comparison



parameters — number of parameters (models: ELMo small, ELMo medium, ELMo original, ELMo large, Bert base, Bert multilingual base, Bert large, GPT, GPT-2 small, GPT-2 medium, GPT-2 large, GPT-2 extra large, GPT-3 small, GPT-3 medium, GPT-3 large, GPT-3 extra large, GPT-3 « GPT-3 »)

corpus size — tokens (models: ELMo, ELMo large, Bert, Bert multilingual, GPT, GPT-2, GPT-3)

More recent models: PaLM (540B), OPT (175B), BLOOM (176B)…

Image source: https://hellofuture.orange.com/en/the-gpt-3-language-model-revolution-or-evolution/

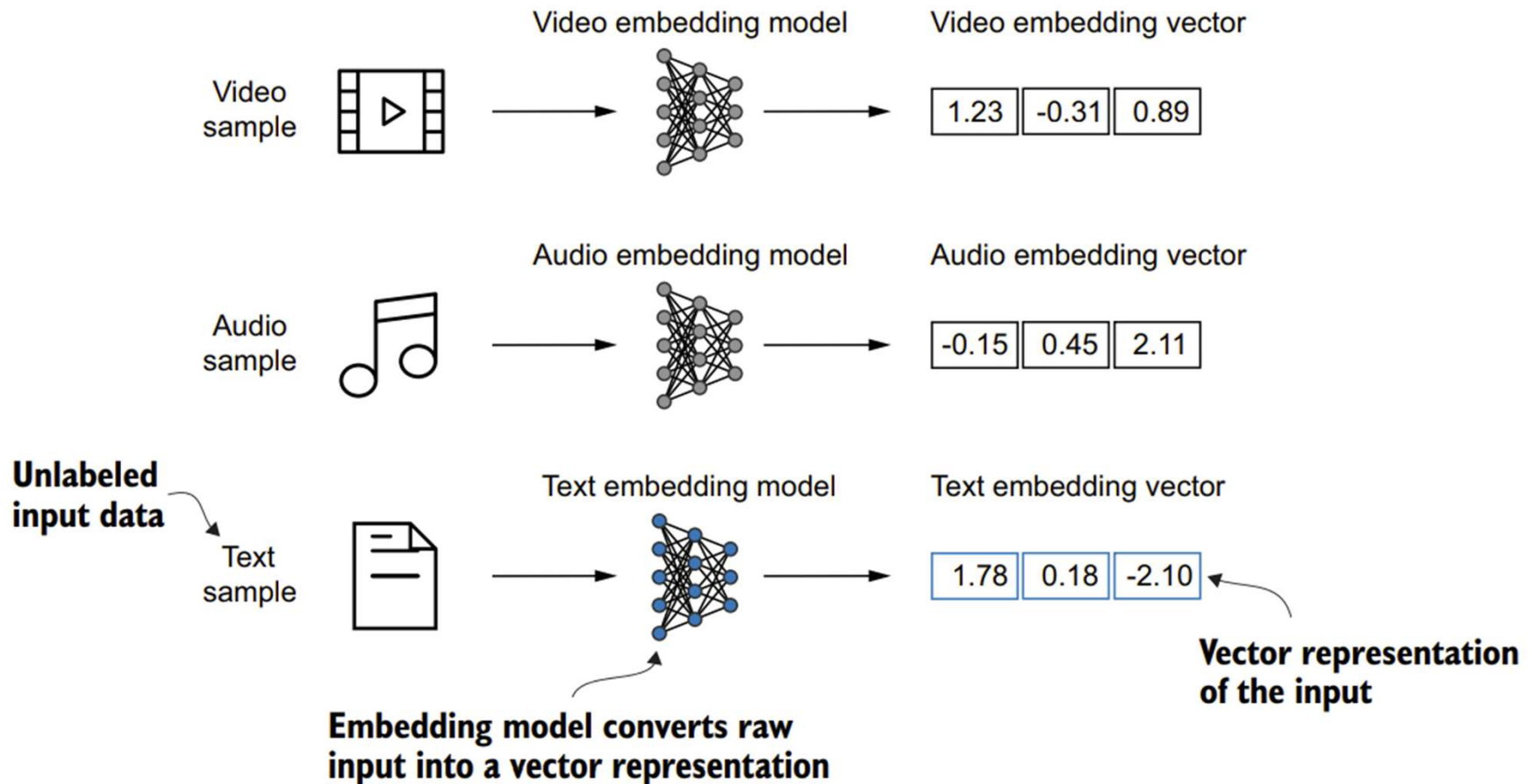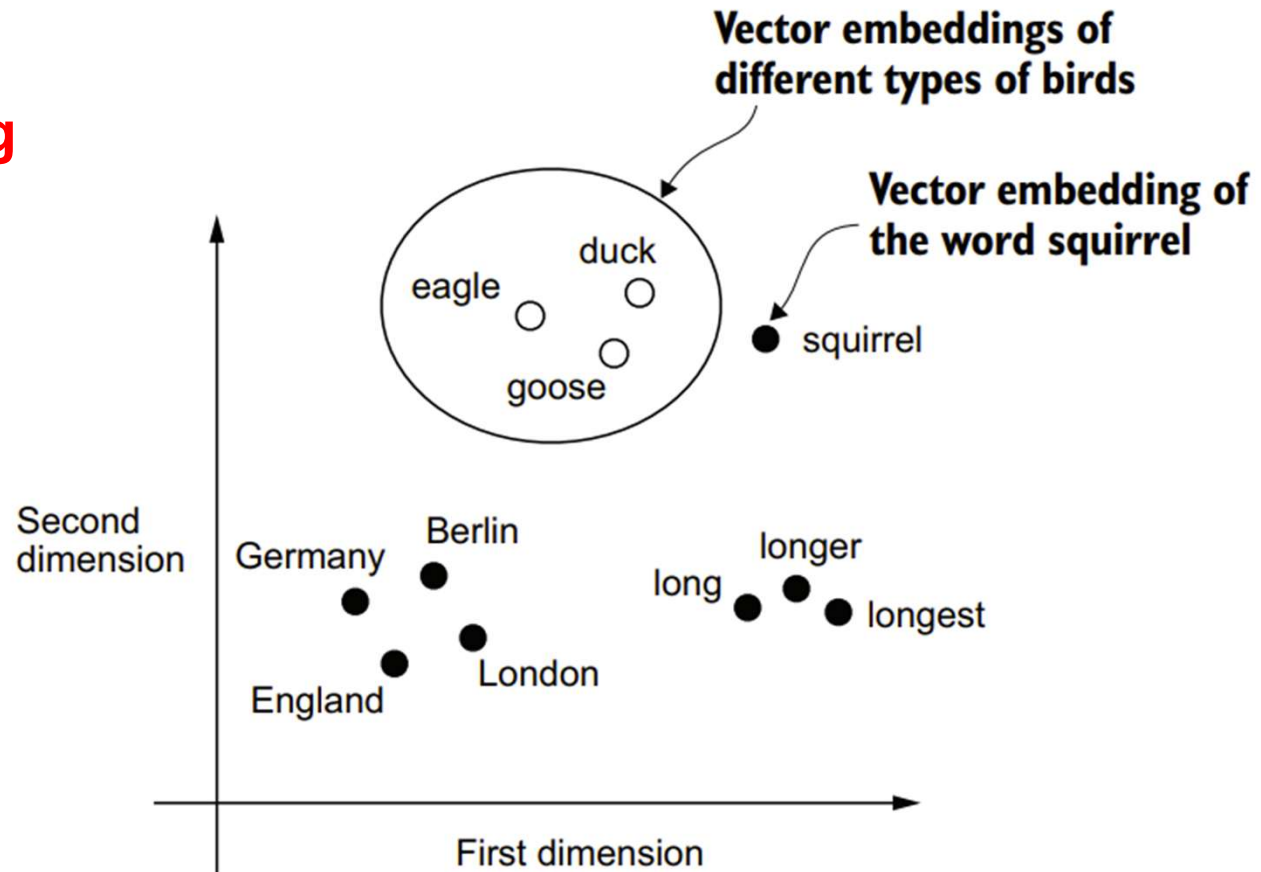# A Closer Look at the GPT Architecture

# Implementing LLMs

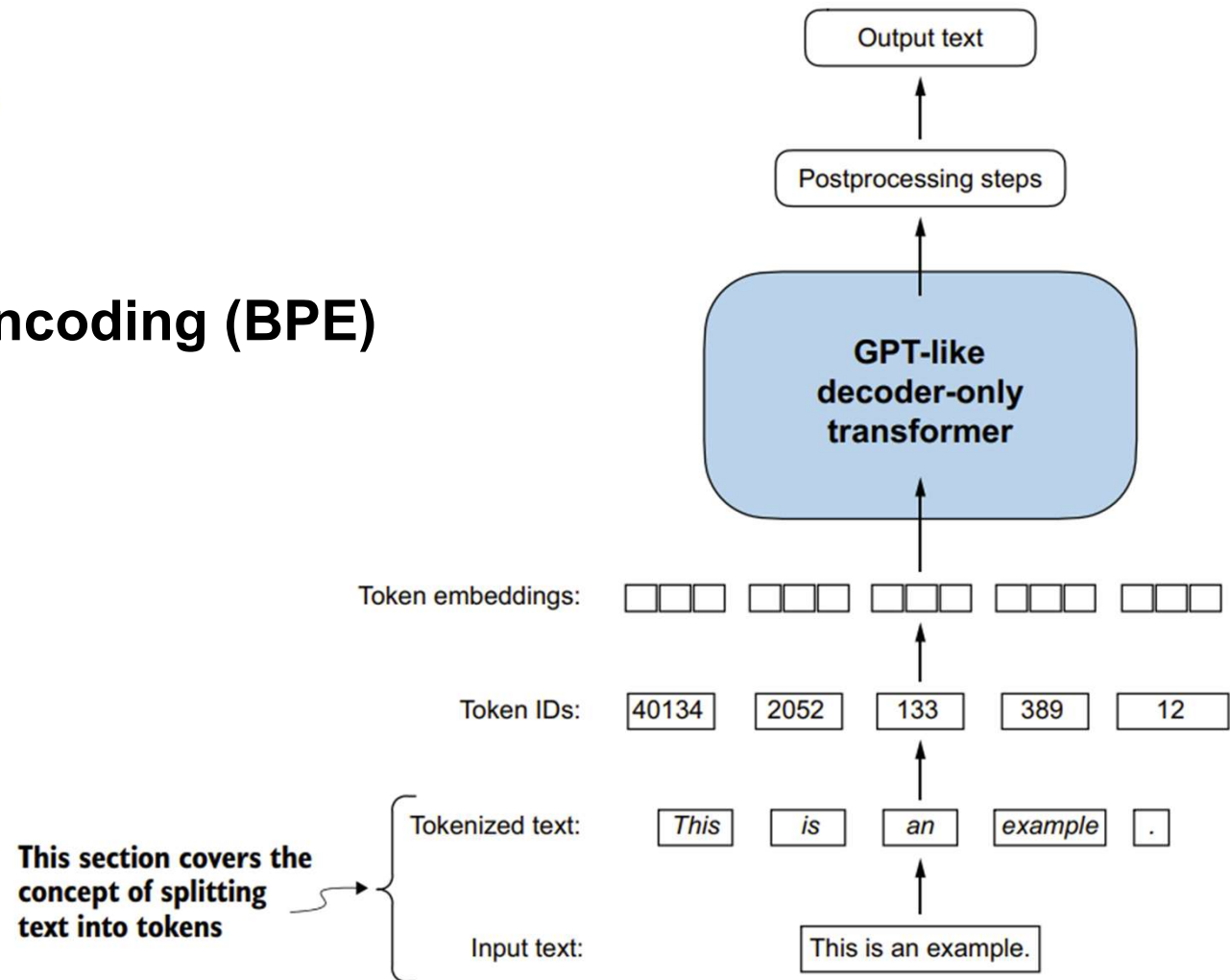# Word embedding − Modelling word in computer

# Word embedding − Modelling word in computer

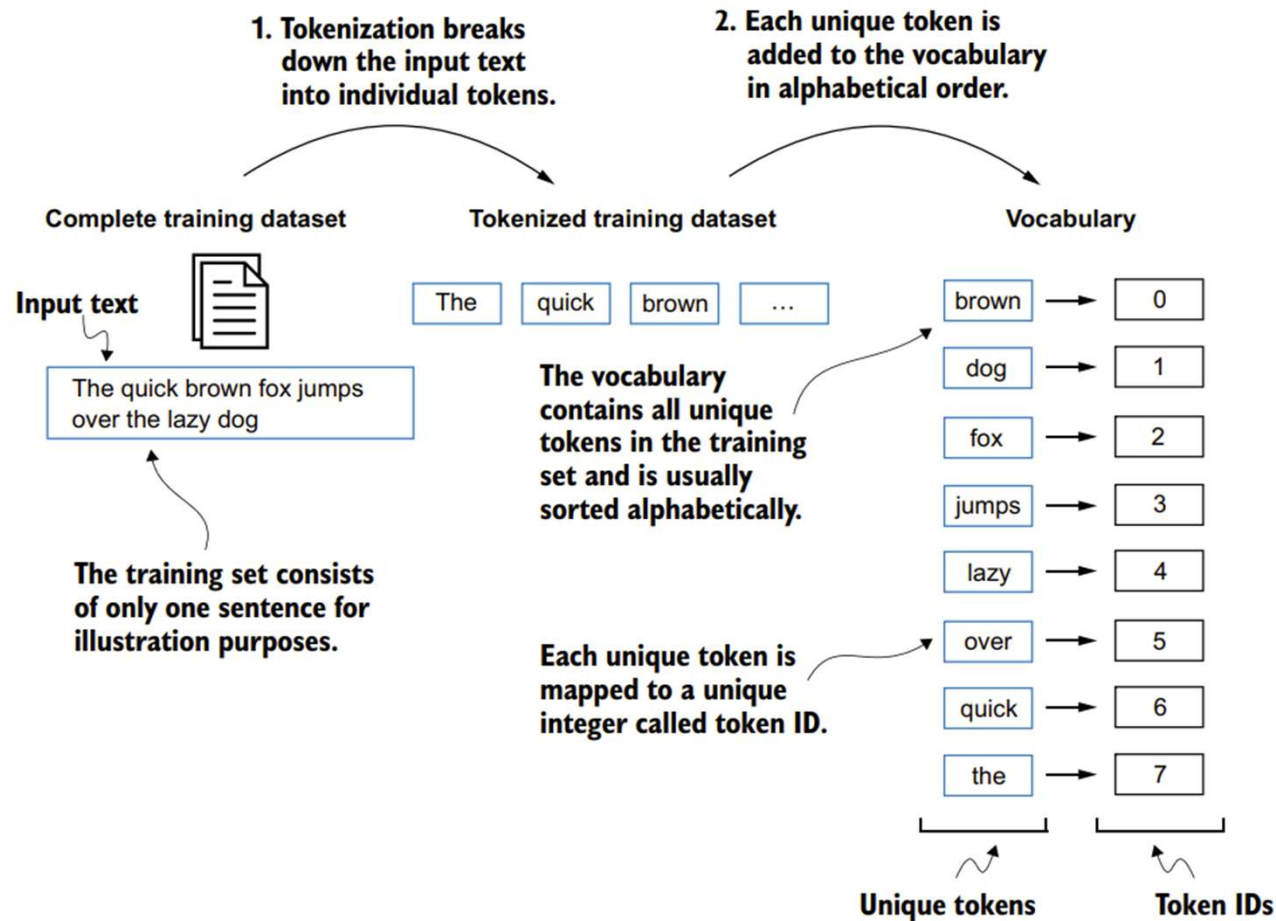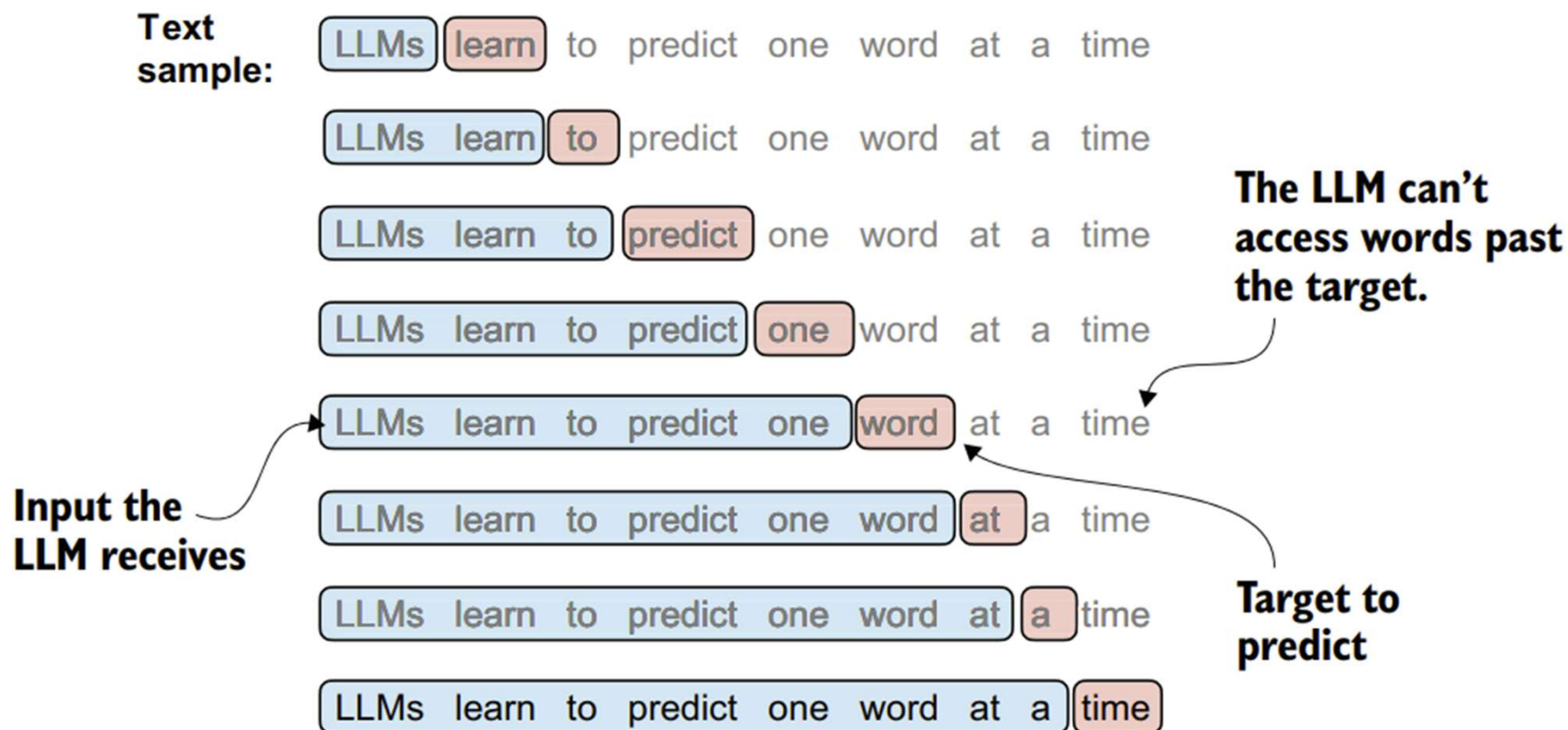**Bringing words that have closer meaning closer in the embedding space**



Vector embeddings of different types of birds

Vector embedding of the word squirrel

# Tokenization

**ChatGPT: Byte Pair Encoding (BPE)**

# Tokenization in action



1. Tokenization breaks down the input text into individual tokens.

2. Each unique token is added to the vocabulary in alphabetical order.

**Complete training dataset**

Input text

The quick brown fox jumps over the lazy dog

The training set consists of only one sentence for illustration purposes.

**Tokenized training dataset**

| The | quick | brown | … |

The vocabulary contains all unique tokens in the training set and is usually sorted alphabetically.

Each unique token is mapped to a unique integer called token ID.

**Vocabulary**

| brown | → | 0 |
| dog | → | 1 |
| fox | → | 2 |
| jumps | → | 3 |
| lazy | → | 4 |
| over | → | 5 |
| quick | → | 6 |
| the | → | 7 |

Unique tokens        Token IDs

# Input to LLM

# Self-Attention mechanism

The last step is multiplying each value vector with its respective attention weight and then summing them to obtain the context vector

# Self-Attention Overall Picture



The embedding vector $x^{(2)}$ of the second input token

$d_{in} = 3$

Inputs $X$

$n = 6$

Weight matrix $W_q$

Weight matrix $W_k$

Weight matrix $W_v$

$d_{in} = 3$

$d_{out} = 2$

Embedded queries, where the second row is the query vector $q^{(2)}$ corresponding to the second input token $x^{(2)}$

Queries $Q$

Keys $K$

Values $V$

We multiply the inputs $X$ with weight matrix $W_v$ to get the value matrix $V$.

**Causal attention**

Attention weight matrix containing the attention scores for each pair of inputs

| | Your | journey | starts | with | one | step |
|---|---|---|---|---|---|---|
| Your | 1.0 | | | | | |
| journey | 0.55 | 0.44 | | | | |
| starts | 0.38 | 0.30 | 0.31 | | | |
| with | 0.27 | 0.24 | 0.24 | 0.23 | | |
| one | 0.21 | 0.19 | 0.19 | 0.18 | 0.19 | |
| step | 0.19 | 0.16 | 0.16 | 0.15 | 0.16 | 0.15 |

Masked out future tokens for the "Your" token

$n = 6$

| | Your | journey | starts | with | one | step |
|---|---|---|---|---|---|---|
| Your | 0.19 | 0.16 | 0.16 | 0.15 | 0.17 | 0.15 |
| journey | 0.20 | 0.16 | 0.16 | 0.14 | 0.16 | 0.14 |
| starts | 0.20 | 0.16 | 0.16 | 0.14 | 0.16 | 0.14 |
| with | 0.18 | 0.16 | 0.16 | 0.15 | 0.16 | 0.15 |
| one | 0.18 | 0.16 | 0.16 | 0.15 | 0.16 | 0.15 |
| step | 0.19 | 0.16 | 0.16 | 0.15 | 0.16 | 0.15 |

$n = 6$

Context vector corresponding to the second input token

Context vectors $Z$
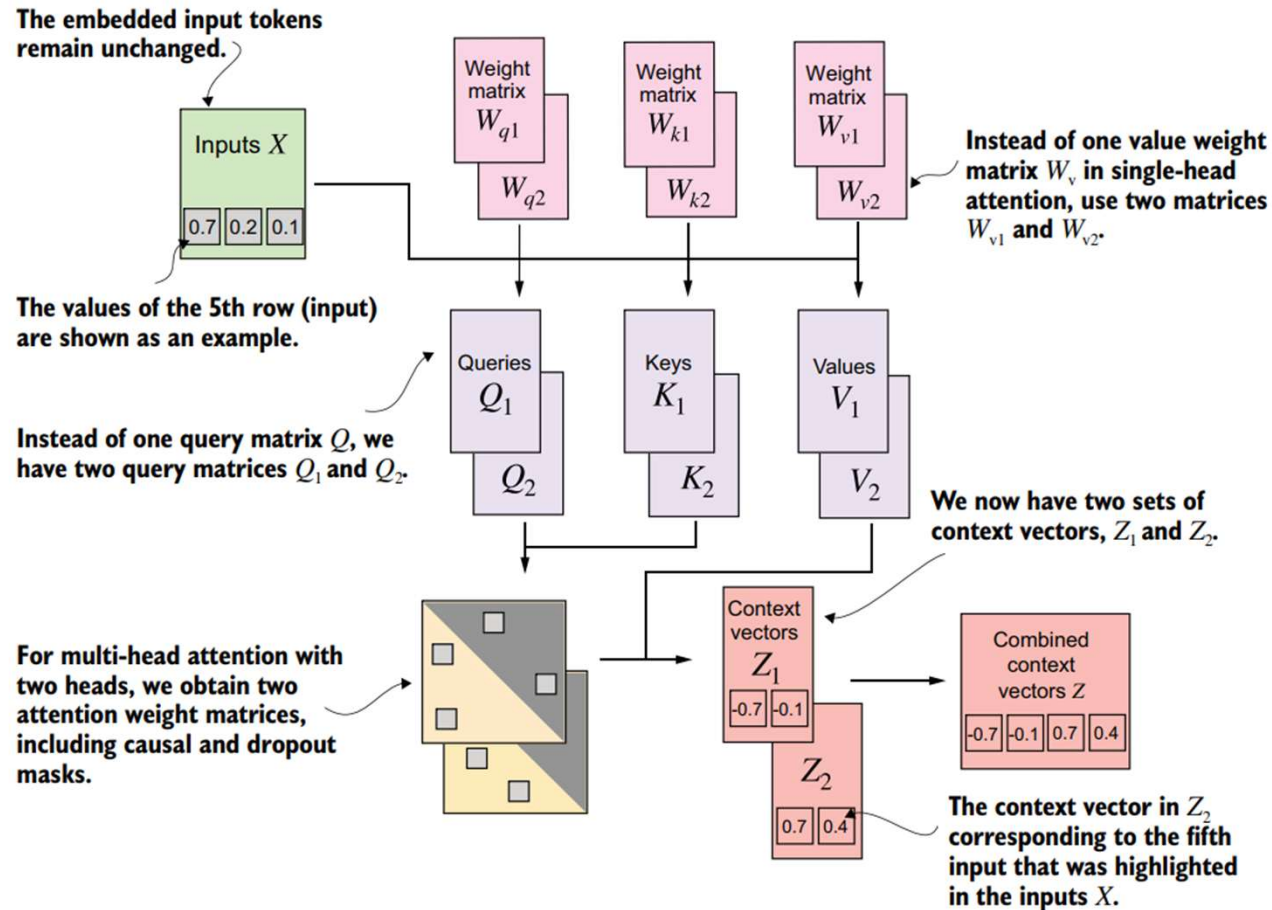
$n = 6$

$d_{out} = 2$

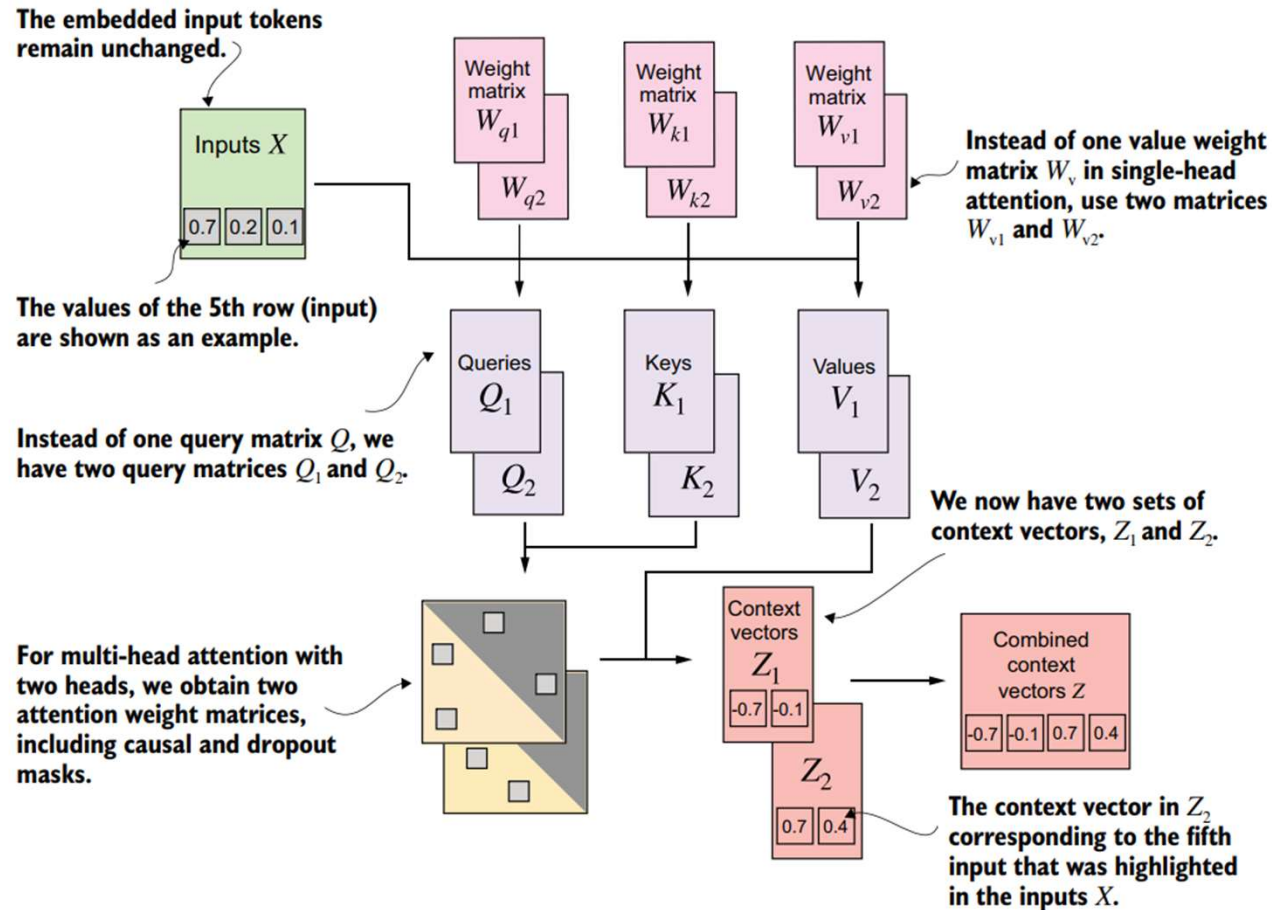# Self-Attention mechanism with dropout

**Help prevent overfitting**

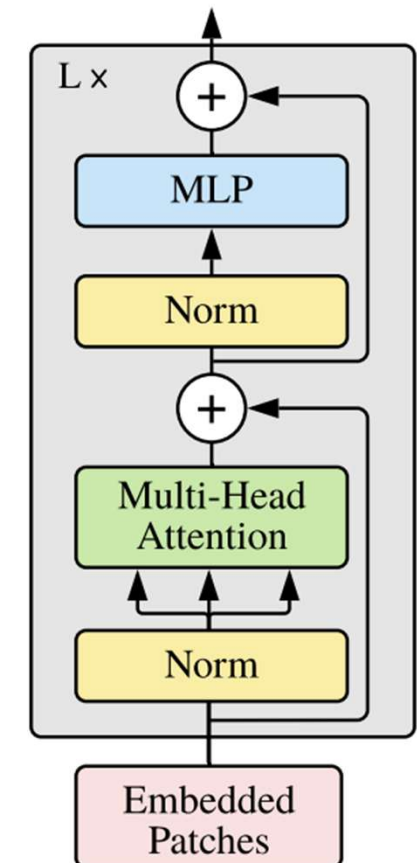**Only used during training**

# Multi-headed self-attention mechanism

# Multi-headed self-attention mechanism



The embedded input tokens remain unchanged.

Inputs $X$

0.7 0.2 0.1

The values of the 5th row (input) are shown as an example.

Weight matrix $W_{q1}$   $W_{q2}$

Weight matrix $W_{k1}$   $W_{k2}$

Weight matrix $W_{v1}$   $W_{v2}$

Instead of one value weight matrix $W_v$ in single-head attention, use two matrices $W_{v1}$ and $W_{v2}$.

Instead of one query matrix $Q$, we have two query matrices $Q_1$ and $Q_2$.

Queries $Q_1$   $Q_2$

Keys $K_1$   $K_2$

Values $V_1$   $V_2$

We now have two sets of context vectors, $Z_1$ and $Z_2$.

For multi-head attention with two heads, we obtain two attention weight matrices, including causal and dropout masks.

Context vectors $Z_1$

-0.7 -0.1

$Z_2$

0.7 0.4

Combined context vectors $Z$

-0.7 -0.1 0.7 0.4

The context vector in $Z_2$ corresponding to the fifth input that was highlighted in the inputs $X$.
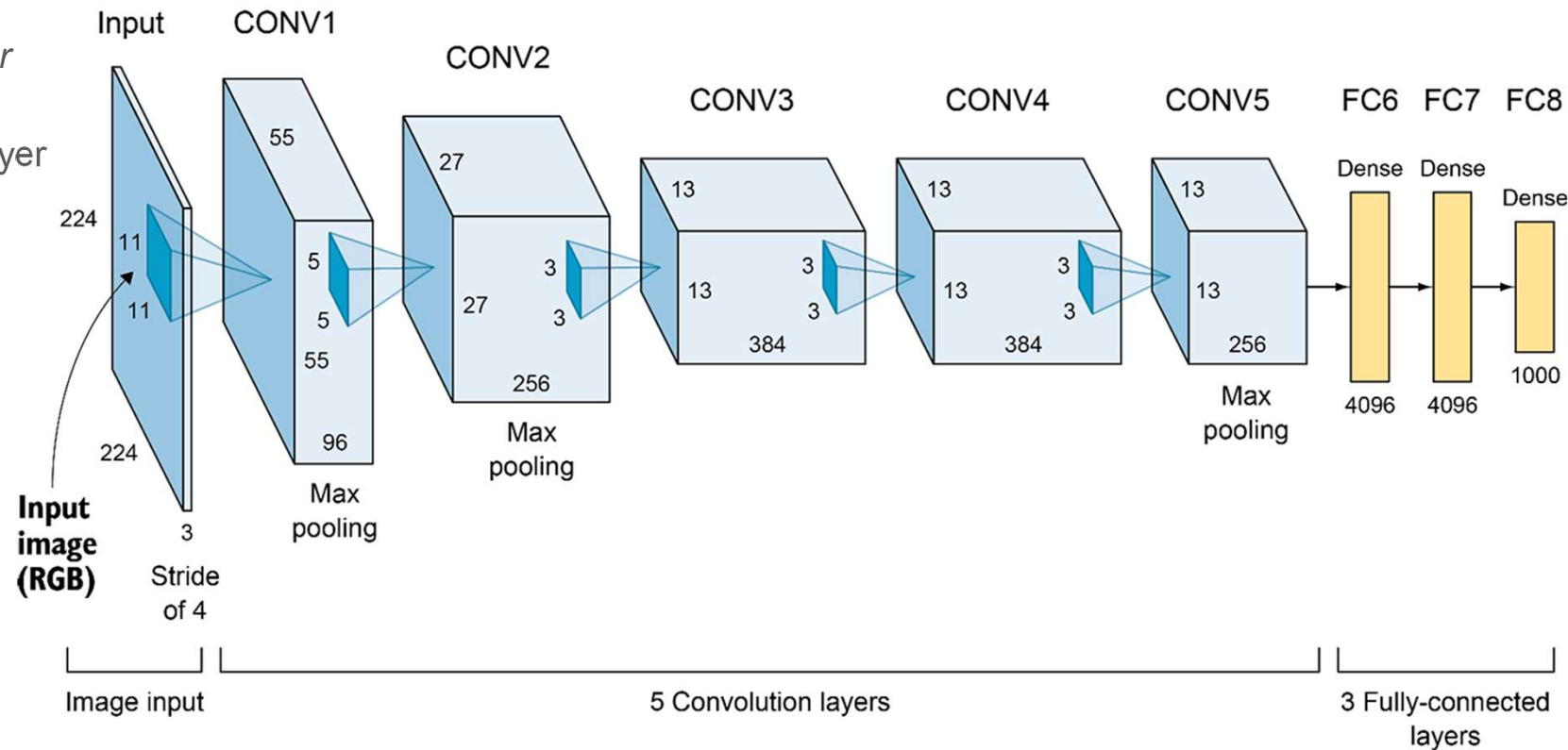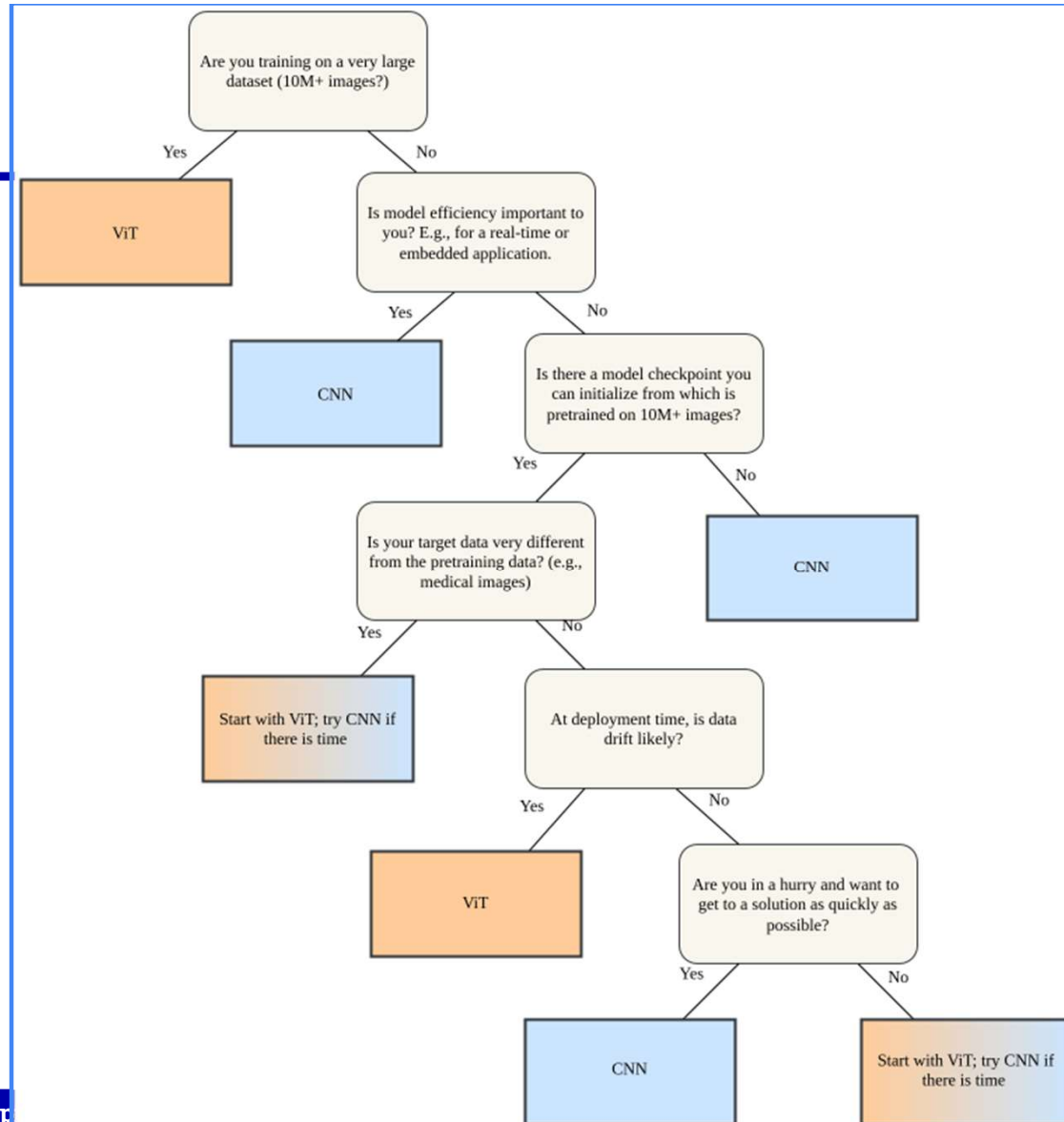
# Vision Transformer (ViT)

# Convolutional neural networks (CNN)

Layer types:
- *Convolutional layer*
- Pooling layer
- Fully-connected layer

# ViT vs CNN for Vision



https://tobiasvanderwerff.com/2024/05/15/cnn-vs-vit.html

# Further reading

- Attention Is All You Need (https://arxiv.org/pdf/1706.03762)
- The Illustrated Transformer (https://jalammar.github.io/illustrated-transformer/)
- Build a Large Language Model (From Scratch) (https://github.com/rasbt/LLMs-from-scratch)