STAGE 1

A building block of an LLM that we implemented in the previous chapter

In this chapter, we implement the other parts of the LLM

In the next chapter, we add the training loop and pretrain the LLM

1) Data preparation & sampling → 2) Attention mechanism → 3) LLM Architecture

Building an LLM

4) Pretraining →

STAGE 2

5) Training loop → 6) Model evaluation → 7) Load pretrained weights

Foundation model

8) Finetuning →

STAGE 3

Dataset with class labels

Classifier

9) Finetuning →

Personal assistant

Instruction dataset

After learning about the attention mechanism, let us learn about the LLM architecture now.



"Every effort moves you **forward**"

The goal is to generate new text one word at a time

**GPT model**

Output layers

In this chapter, we implement a GPT model including all of its subcomponents

**Transformer block**

Transformer blocks are a key component of GPT-like LLMs

Masked multi-head attention

We implemented the attention module in the previous chapter

Embedding layers

Embedding layers and tokenization were covered in Chapter 2

Tokenized text

"Every effort moves you"

# Zoom Into the Transformer block



Outputs have the same form and dimensions as the inputs

The transformer block

The input tokens to be embedded

| | |
|---|---|
| Every | [[0.2961, ..., 0.4604], |
| effort | [0.2238, ..., 0.7598], |
| moves | [0.6945, ..., 0.5963], |
| you | [0.0890, ..., 0.5833]] |

A view into the "Feed forward" block

Shortcut connection

This tensor represents an embedded text sample that serves as input to the transformer block

Each row is a 768-dimensional vector representing an embedded input token

③ What we are yet to learn:

ⓐ Transformer blocks

④ We will scale upto the size of a small GPT-2 model → 124 million parameters

**Language Models are Unsupervised Multitask Learners**

Alec Radford [*][1]  Jeffrey Wu [*][1]  Rewon Child [1]  David Luan [1]  Dario Amodei [**][1]  Ilya Sutskever [**][1]

**Abstract**

Natural language processing tasks, such as question answering, machine translation, reading comprehension, and summarization, are typically approached with supervised learning on task-specific datasets. We demonstrate that language models begin to learn these tasks without any explicit supervision when trained on a new dataset of millions of webpages called WebText. When conditioned on a document plus questions, the answers generated by the language model reach 55 F1 on the CoQA dataset - matching or exceeding the performance of 3 out of 4 baseline systems without using the 127,000+ training examples. The capacity of the language model is essential to the success of zero-shot task transfer and increasing it improves performance in a log-linear fashion across tasks. Our largest model, GPT-2, is a 1.5B parameter Transformer that achieves state of the art results on 7 out of 8 tested language modeling datasets in a zero-shot setting

competent generalists. We would like to move towards more general systems which can perform many tasks – eventually without the need to manually create and label a training dataset for each one.

The dominant approach to creating ML systems is to collect a dataset of training examples demonstrating correct behavior for a desired task, train a system to imitate these behaviors, and then test its performance on independent and identically distributed (IID) held-out examples. This has served well to make progress on narrow experts. But the often erratic behavior of captioning models (Lake et al., 2017), reading comprehension systems (Jia & Liang, 2017), and image classifiers (Alcorn et al., 2018) on the diversity and variety of possible inputs highlights some of the shortcomings of this approach.

Our suspicion is that the prevalence of single task training on single domain datasets is a major contributor to the lack of generalization observed in current systems. Progress towards robust systems with current architectures is likely to require training and measuring performance on a wide range of domains and tasks. Recently, several benchmarks

| Parameters | Layers | $d_{model}$ |
|------------|--------|-------------|
| 117M | 12 | 768 |
| 345M | 24 | 1024 |
| 762M | 36 | 1280 |
| 1542M | 48 | 1600 |

⑤ Open AI has made GPT-2 weights public.
GPT-3, 4 weights have not yet been made public.

We will use these parameters

```
GPT_CONFIG_124M = {
    "vocab_size": 50257,       # Vocabulary size
    "context_length": 1024,    # Context length
    "emb_dim": 768,            # Embedding dimension
    "n_heads": 12,             # Number of attention heads
    "n_layers": 12,            # Number of layers
    "drop_rate": 0.1,          # Dropout rate
    "qkv_bias": False          # Query-Key-Value bias
}
```

Context length: how many maximum words are used to predict the next word(here word means tokens)

Ervery token in the vocabulary will we have projected in vector space(here we use 768) and the embedding shoud be such that the so meaning should be captured
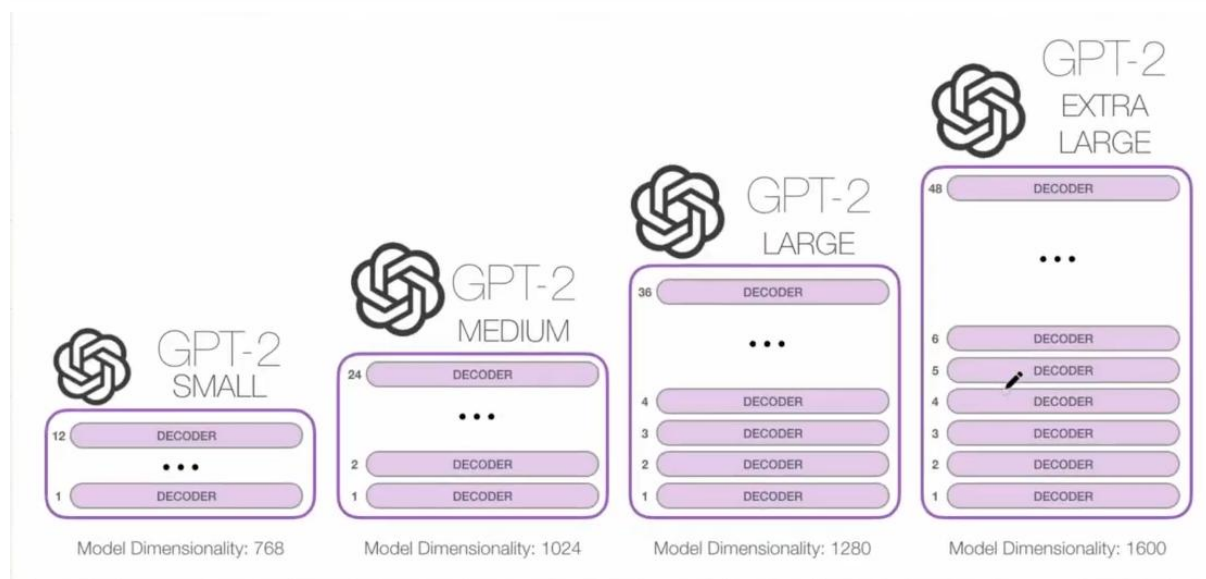
Number of heads: Number of attention head

# Number of layers: number of transformer blocks

```
[[ 0.0256, ..., 0.6890],
 [-0.0178, ..., 0.7431],
 [ 0.4558, ..., 0.7814],
 [ 0.0702, ..., 0.7134]]
```

**Outputs have the same form and dimensions as the inputs**

**The transformer block**

Dropout

Feed forward

LayerNorm 2

Dropout

Masked multi-head attention

LayerNorm 1

Linear layer

GELU activation

Linear layer

**A view into the "Feed forward" block**

**The input tokens to be embedded**

**Shortcut connection**

```
Every    [[0.2961, ..., 0.4604],
effort    [0.2238, ..., 0.7598],
moves     [0.6945, ..., 0.5963],
you       [0.0890, ..., 0.5833]]
```

**This tensor represents an embedded text sample that serves as input to the transformer block**

**Each row is a 768-dimensional vector representing an embedded input token**

GPT-2
EXTRA
LARGE

GPT-2
LARGE

GPT-2
MEDIUM

GPT-2
SMALL

| 48 | DECODER |
| ... |
| 6 | DECODER |
| 5 | DECODER |
| 4 | DECODER |
| 3 | DECODER |
| 2 | DECODER |
| 1 | DECODER |

| 36 | DECODER |
| ... |
| 4 | DECODER |
| 3 | DECODER |
| 2 | DECODER |
| 1 | DECODER |

| 24 | DECODER |
| ... |
| 2 | DECODER |
| 1 | DECODER |

| 12 | DECODER |
| ... |
| 1 | DECODER |

Model Dimensionality: 768    Model Dimensionality: 1024    Model Dimensionality: 1280    Model Dimensionality: 1600

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 $W_{v2}$

Queries $Q_1$
$Q_2$

Keys $K_1$
$K_2$

Values $V_1$
$V_2$

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

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 $X$

Using this configuration, we will build a GPT placeholder architecture (Dummy GPT model)

This will give us a birds eye view of how everything fits together.

At the end of this chapter, we use multiple transformer blocks to implement the GPT model that we will train in the upcoming chapter

7) Final GPT architecture
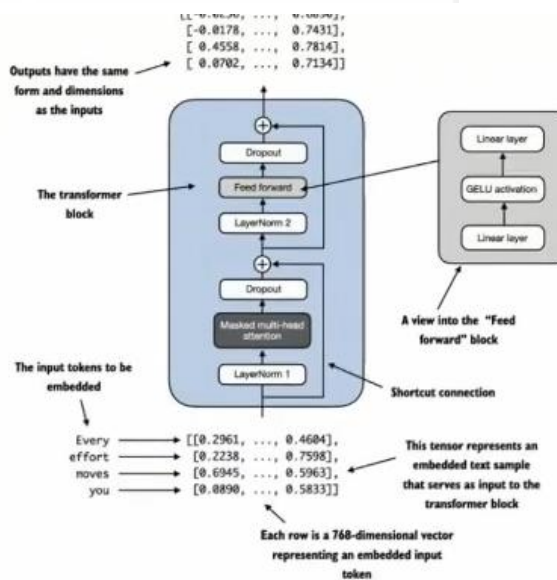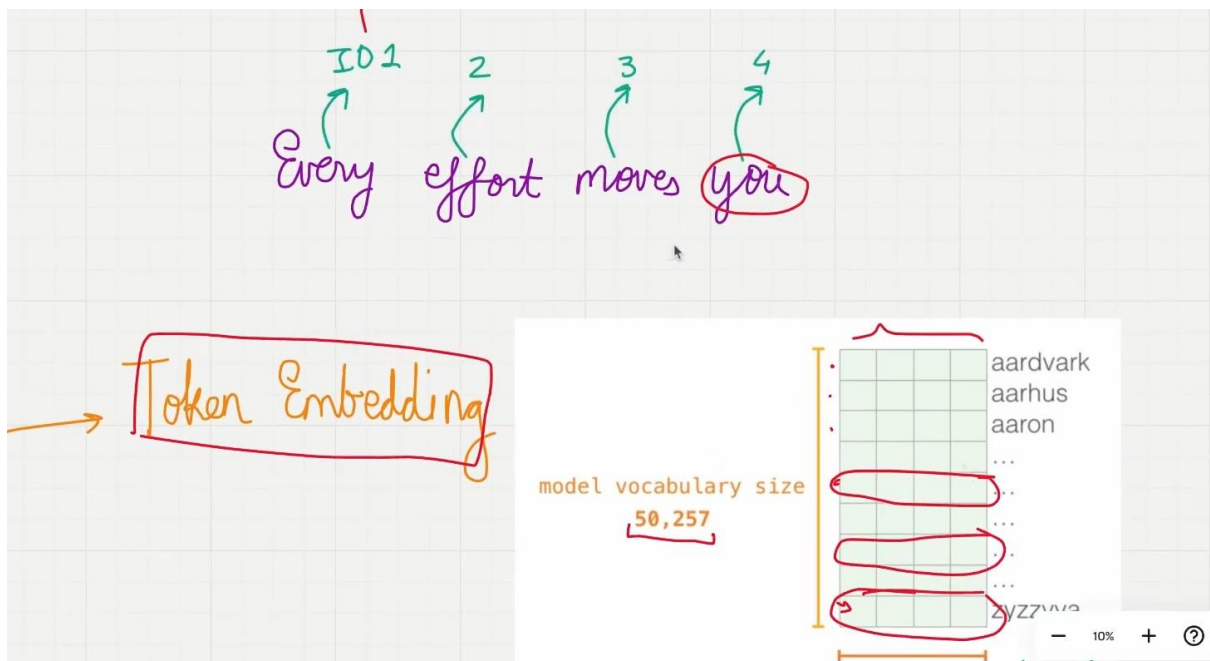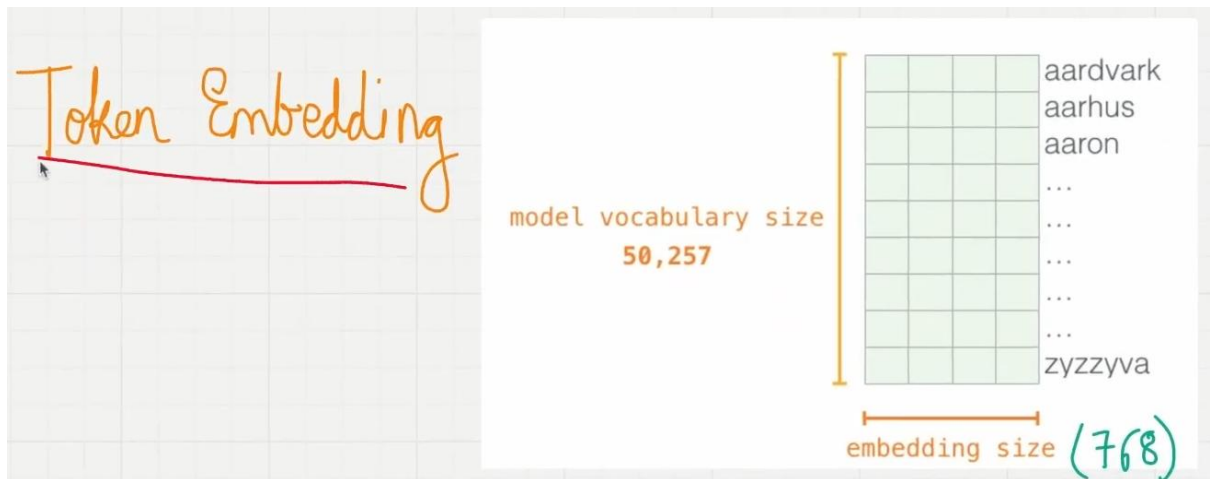
We combine the building blocks 2-5, including the multi-head attention module covered in the previous chapter, into a transformer block

6) Transformer block

In the upcoming sections, we implement the different building blocks 2-5 that are used inside a GPT model

2) Layer normalization

3) GELU activation

4) Feed forward network

5) Shortcut connections

We start this chapter by implementing a GPT placeholder model to see the overall structure of the model

1) GPT backbone

→ Mental model in which we will code the GPT architecture

Outputs have the same form and dimensions as the inputs

[[-0.0230, ..., 0.0850],
 [-0.0178, ..., 0.7431],
 [ 0.4558, ..., 0.7814],
 [ 0.0702, ..., 0.7134]]

The transformer block

Linear layer

GELU activation

Linear layer

Dropout
Feed forward
LayerNorm 2
Dropout
Masked multi-head attention
LayerNorm 1

A view into the "Feed forward" block

Shortcut connection

The input tokens to be embedded

Every → [[0.2961, ..., 0.4604],
effort → [0.2238, ..., 0.7598],
moves → [0.6945, ..., 0.5963],
you → [0.0890, ..., 0.5833]]

This tensor represents an embedded text sample that serves as input to the transformer block

Each row is a 768-dimensional vector representing an embedded input token

768

ID 1    2    3    4

Every effort moves you

---

Token Embedding



model vocabulary size
50,257

aardvark
aarhus
aaron
...
...
...
...
...
zyzzyva

embedding size (768)

---

ID 1    2    3    4

Every effort moves (you)

Token Embedding



model vocabulary size
,50,257

aardvark
aarhus
aaron
...
...
...
...
...
zyzzyva

# Positional Embedding

Context size
**1024**

1
2
3
...
...
...
...
1024

embedding size (768)

→ Output logits for 1 test sample:

Token 1 — — — — — — ○ — —

Token 2 — — — — — — — — — —

Token 3 — — — — — — — — —

Token 4 — — — — — — — — —

Each element here corresponds to the probability of it being the next token.

# Vocabulary size (50257)

At the end of this chapter, we use multiple transformer blocks to implement the GPT model that we will train in the upcoming chapter

**7) Final GPT architecture**

We combine the building blocks 2-5, including the multi-head attention module covered in the previous chapter, into a transformer block

**6) Transformer block**

In the upcoming sections, we implement the different building blocks 2-5 that are used inside a GPT model

**2) Layer normalization**

**3) GELU activation**

**4) Feed forward network**

**5) Shortcut connections**

**1) GPT backbone**

We start this chapter by implementing a GPT placeholder model to see the overall structure of the model