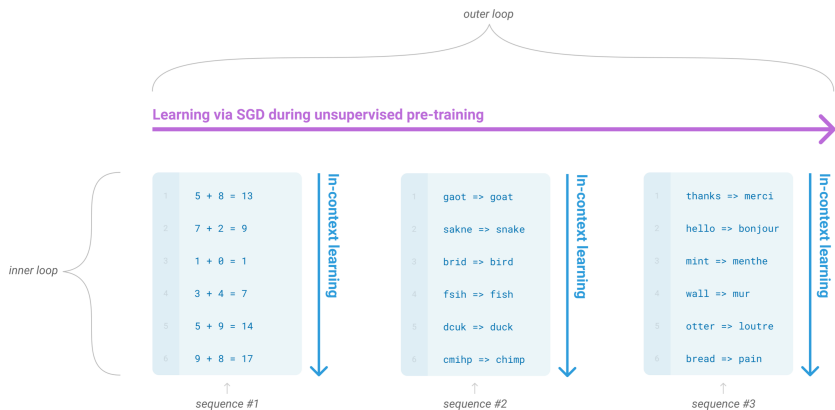# Intro to GPT & X-shot learning

# GPT & Benchmarks

**Learning goals**

- Recap GPT and the ideas behind standard language modelling

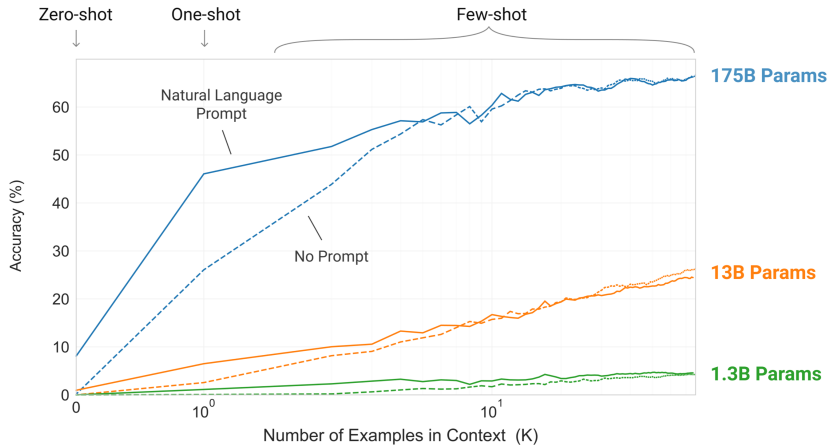- Understand the difference between fine-tuning and X-shot learning

# GPT

- Like BERT, GPT is a language model.
- But not MLM, but a conventional language model: it predicts the next word (or subword).
- Like BERT, GPT is trained on a huge corpus, actually an even huger corpus.
- Like BERT, GPT is a transformer architecture.
- Difference 1: GPT is a single model that aims to solve all tasks.
    - It can also switch back and forth between tasks and solve tasks within tasks, another human capability that is important in practice. "fluidity"
- Difference 2: GPT leverages task descriptions.
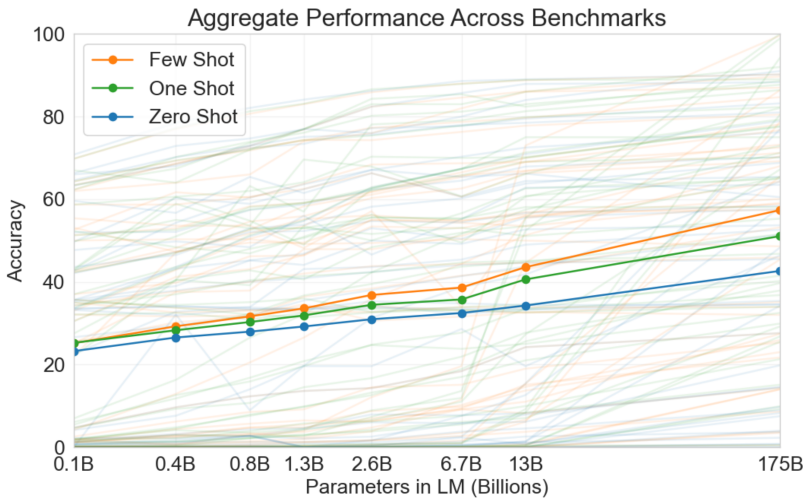- Difference 3: GPT is effective at few-shot learning.

# GPT: TWO TYPES OF LEARNING



outer loop

Learning via SGD during unsupervised pre-training

inner loop

In-context learning

| | |
|---|---|
| 1 | 5 + 8 = 13 |
| 2 | 7 + 2 = 9 |
| 3 | 1 + 0 = 1 |
| 4 | 3 + 4 = 7 |
| 5 | 5 + 9 = 14 |
| 6 | 9 + 8 = 17 |

↑
sequence #1

| | |
|---|---|
| 1 | gaot => goat |
| 2 | sakne => snake |
| 3 | brid => bird |
| 4 | fsih => fish |
| 5 | dcuk => duck |
| 6 | cmihp => chimp |

↑
sequence #2

| | |
|---|---|
| 1 | thanks => merci |
| 2 | hello => bonjour |
| 3 | mint => menthe |
| 4 | wall => mur |
| 5 | otter => loutre |
| 6 | bread => pain |

↑
sequence #3

# GPT: EFFECTIVE IN-CONTEXT LEARNING

# X-SHOT COMPARISON AND EFFECT OF LARGER CORPORA



Aggregate Performance Across Benchmarks

# FINE-TUNING (NOT USED BY GPT)



Traditional fine-tuning (not used for GPT-3)

**Fine-tuning**

The model is trained via repeated gradient updates using a large corpus of example tasks.

```
1   sea otter => loutre de mer        ← example #1
```

**gradient update**

```
1   peppermint => menthe poivrée      ← example #2
```

**gradient update**

```
1   plush giraffe => girafe peluche   ← example #N
```

**gradient update**

```
1   cheese => ..............................   ← prompt
```

# ZERO-SHOT (NO GRADIENT UPDATE)

**Zero-shot**

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.

```
1    Translate English to French:          ←  task description

2    cheese =>                             ←  prompt
```

# ONE-SHOT (NO GRADIENT UPDATE)

**One-shot**

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.

```
1   Translate English to French:        ←  task description

2   sea otter => loutre de mer          ←  example

3   cheese =>                           ←  prompt
```

# FEW-SHOT (NO GRADIENT UPDATE)

**Few-shot**

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.

```
1    Translate English to French:        ←——— task description

2    sea otter => loutre de mer          ←——— examples

3    peppermint => menthe poivrée        ←

4    plush girafe => girafe peluche      ←

5    cheese =>  ..........................  ←——— prompt
```

# ARCHITECTURE

| Model Name | $n_{\mathrm{params}}$ | $n_{\mathrm{layers}}$ | $d_{\mathrm{model}}$ | $n_{\mathrm{heads}}$ | $d_{\mathrm{head}}$ | Batch Size | Learning Rate |
|---|---|---|---|---|---|---|---|
| GPT-3 Small | 125M | 12 | 768 | 12 | 64 | 0.5M | $6.0 \times 10^{-4}$ |
| GPT-3 Medium | 350M | 24 | 1024 | 16 | 64 | 0.5M | $3.0 \times 10^{-4}$ |
| GPT-3 Large | 760M | 24 | 1536 | 16 | 96 | 0.5M | $2.5 \times 10^{-4}$ |
| GPT-3 XL | 1.3B | 24 | 2048 | 24 | 128 | 1M | $2.0 \times 10^{-4}$ |
| GPT-3 2.7B | 2.7B | 32 | 2560 | 32 | 80 | 1M | $1.6 \times 10^{-4}$ |
| GPT-3 6.7B | 6.7B | 32 | 4096 | 32 | 128 | 2M | $1.2 \times 10^{-4}$ |
| GPT-3 13B | 13.0B | 40 | 5140 | 40 | 128 | 2M | $1.0 \times 10^{-4}$ |
| GPT-3 175B or "GPT-3" | 175.0B | 96 | 12288 | 96 | 128 | 3.2M | $0.6 \times 10^{-4}$ |

# TRAINING CORPUS

| Dataset | Quantity (tokens) | Weight in training mix | Epochs elapsed when training for 300B tokens |
|---------|-------------------|------------------------|----------------------------------------------|
| Common Crawl (filtered) | 410 billion | 60% | 0.44 |
| WebText2 | 19 billion | 22% | 2.9 |
| Books1 | 12 billion | 8% | 1.9 |
| Books2 | 55 billion | 8% | 0.43 |
| Wikipedia | 3 billion | 3% | 3.4 |

# LOSS AS A FUNCTION OF COMPUTE