# LARGE LANGUAGE MODELS

Continual Learning

# WHO AM I?

➢ Bachelor & Master Computer Science (2016-2022)

➢ PhD at the Artificial Intelligence & Machine Learning Lab (since 2022)

➢ Research Scientist at the German Research Center for AI (DFKI) (since 2023)

➢ Visiting Research Scientist at Adobe (2023 & 2024)

➢ Co-founder of OcciGlot Research Initiative for European Language Models (since 2024)

# AGENDA

**1** Introduction

➢ Context

➢ Natural Language Processing

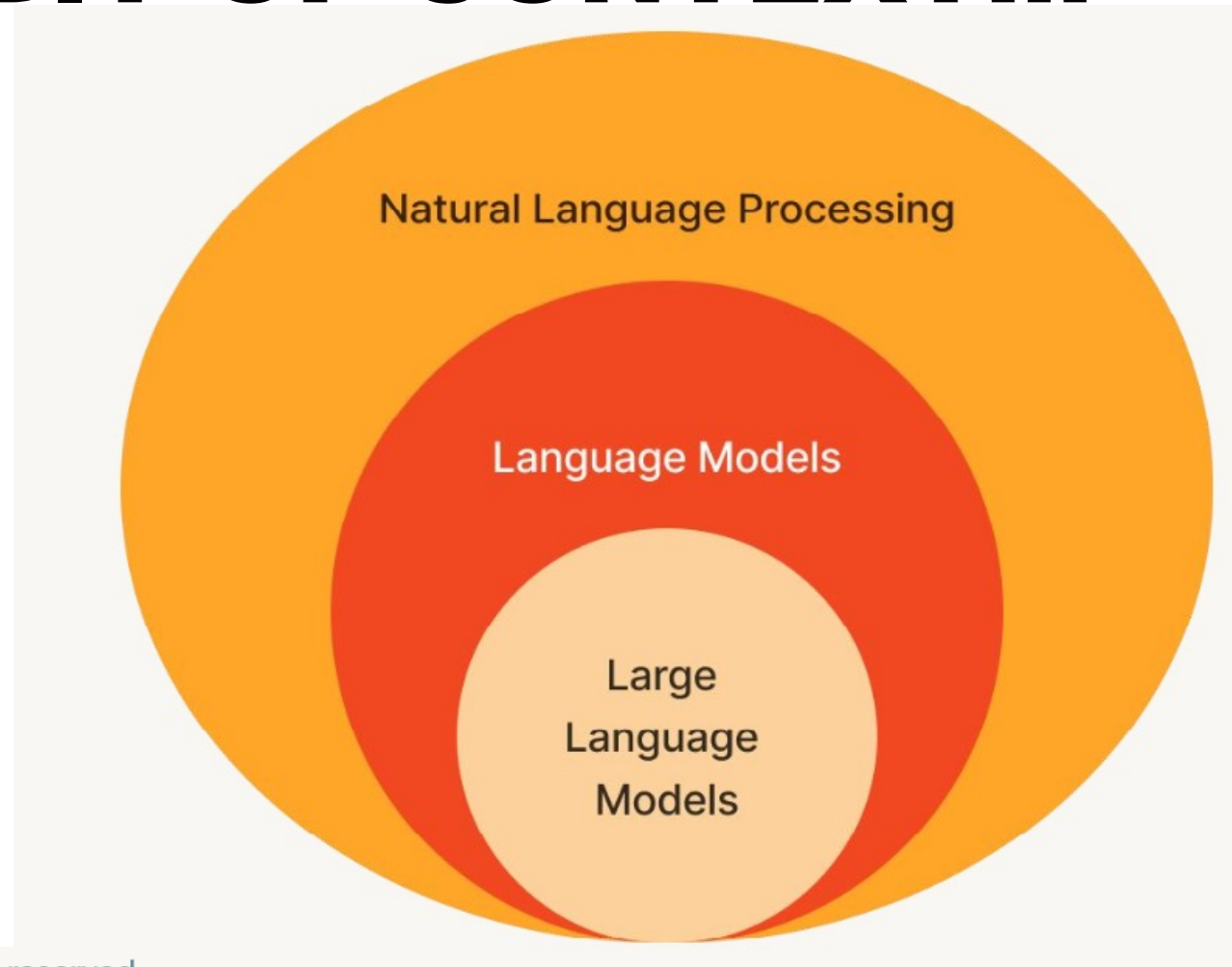➢ Language Modelling

➢ Components of LLMs

**2** Building LLMs

➢ Stages of Training

➢ Pre-Training

➢ Instruction Tuning

➢ Preference Tuning

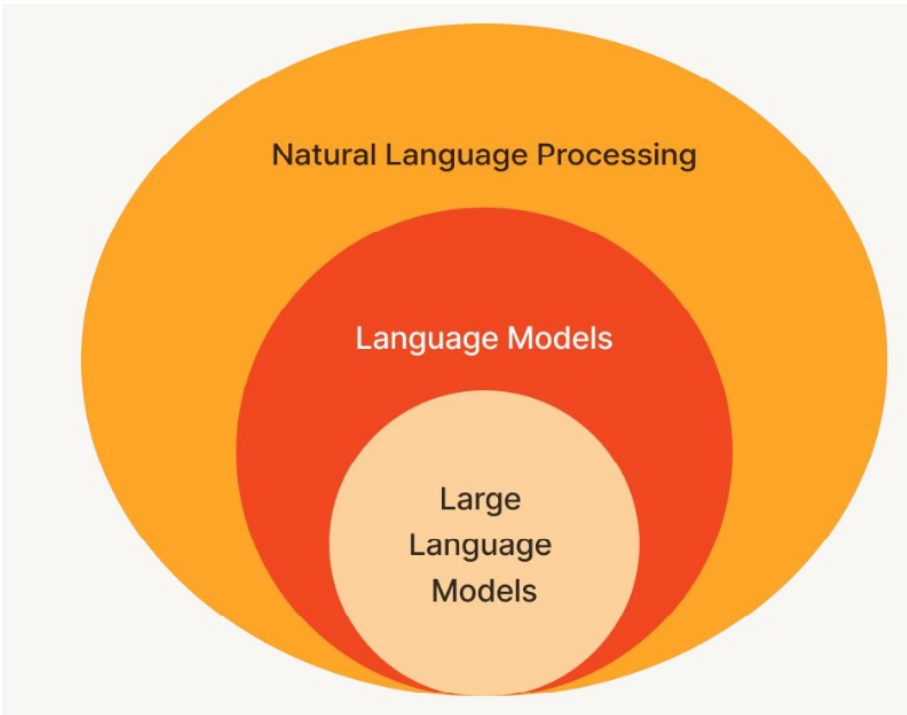**3** Applications & Challenges

➢ Domain Adaptation

➢ Unique Challenges
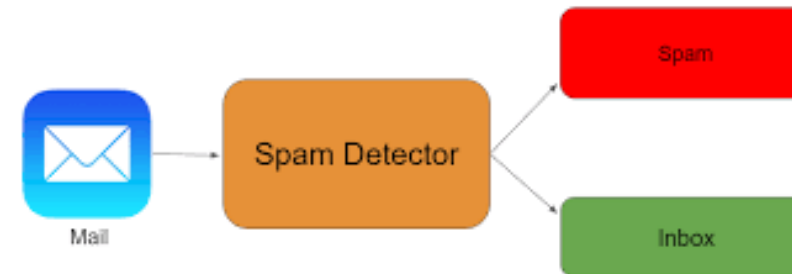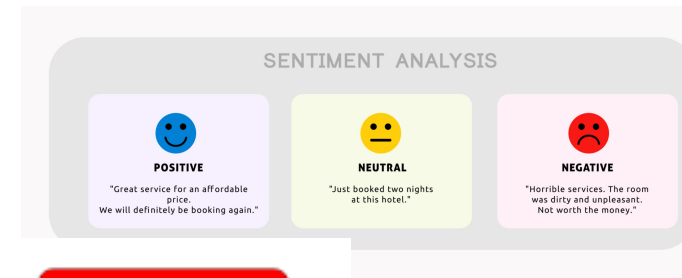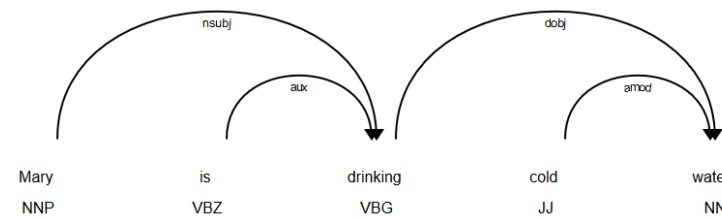
# LITTLE BIT OF CONTEXT...

# NATURAL LANGUAGE PROCESSING (NLP)

➢ Over 50 years old field

➢ Originated from linguistics

*"giving computers the ability to support and manipulate human language"*

Figures from:
https://bhatnagar91.medium.com/analyzing-youtube-fans-feelings-uncovering-the-power-of-sentiment-analysis-5d840909ac58
https://towardsdatascience.com/spam-detection-in-emails-de0398ea3b48
https://suttipong-kull.medium.com/how-to-extract-subject-verb-and-object-by-nlp-4149323a7d7d

# LANGUAGE MODELING

➤ Probabilistic of sequences of "words"

"*dsfh hjaiorpghh fdhjol adhjj auezoijh*"         ➡         **Low probability**
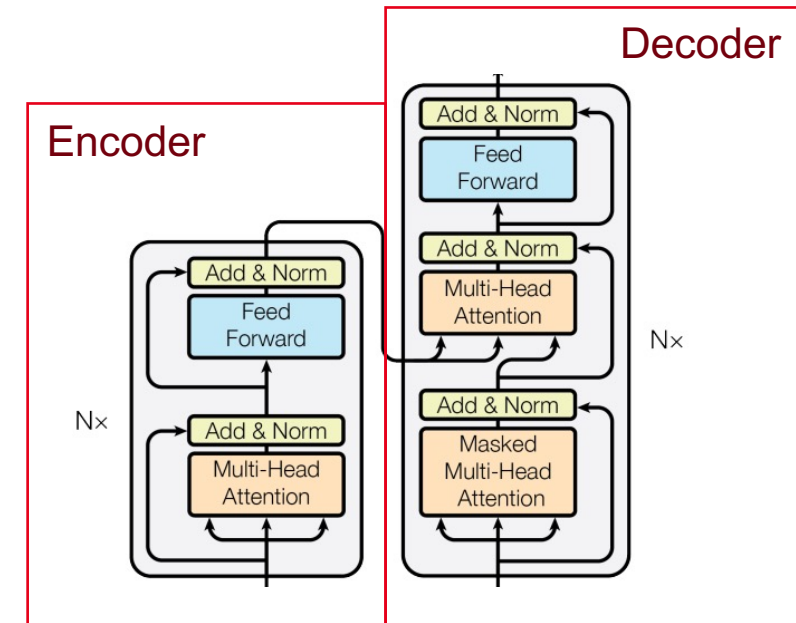
"*language modelling is fun*"         ➡         **High probability**

➤ **Assumption:** *The probability of the next word in a sequence only depends on the previous ones*

$$\mathcal{P}(w_n|w_{n-1} \ldots w_1) = \prod \mathcal{P}(w_i|w_{i-1})$$
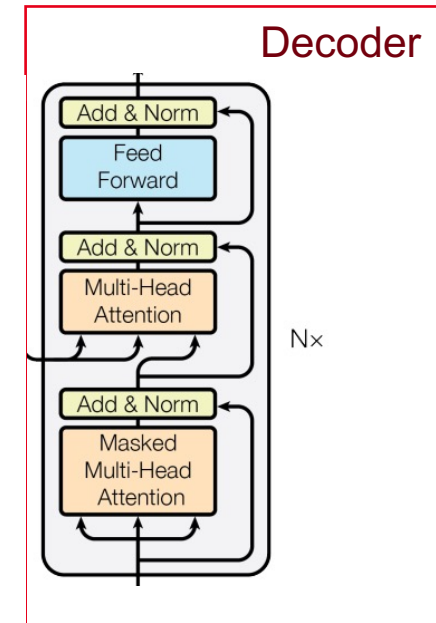
# LARGE LANGUAGE MODELS

➤ (Large) neural networks that model text probabilistics

   ➤ BERT (2018) up to **340M** parameters & **2.5B** words

   ➤ GPT-4 (2023) probably 16x111B MoE = **1.7T** parameters

   ➤ Llama-3 (2024) **15T** words



Vaswani et al. "Attention is all you need." *NeurIPS* (2017).

How do we need to setup our neural network

to perform language modelling on text?

$$\mathcal{P}(w_n|w_{n-1} \dots w_1) = \prod \mathcal{P}(w_i|w_{i-1})$$

Decoder



Vaswani et al. "Attention is all you need." *NeurIPS* (2017).

# LLMS IN A NUTSHELL

➢ Build a large vocabulary of (sub-)words (30k-250k)

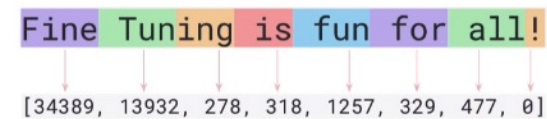➢ Tokenize input text into sequence of vocabulary IDs

> *"language modelling is fun"*→ ['<s>', '⎵language', '⎵mod', 'elling', '⎵is', '⎵fun'] → [1, 3842, 968, 3572, 349, 746]

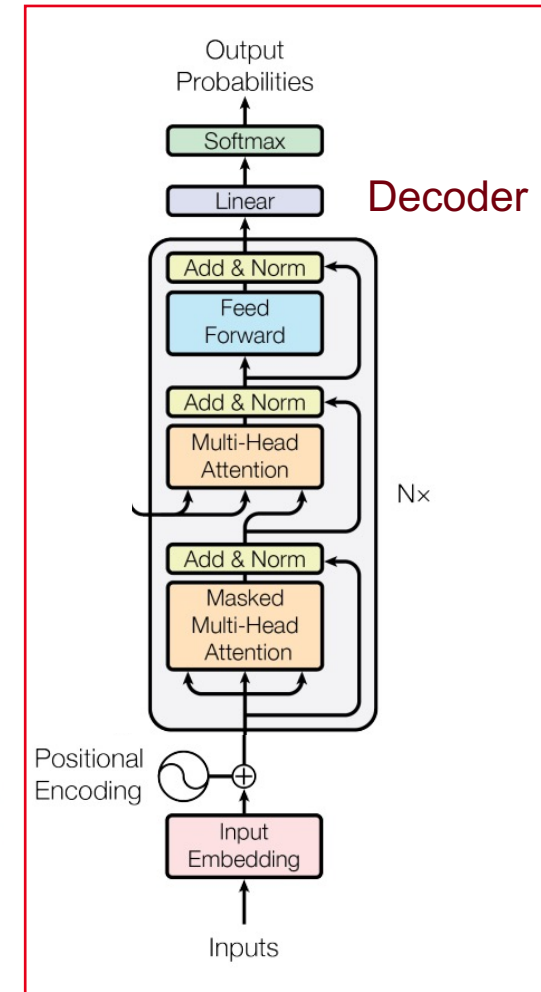➢ Each token is embedded to a learned representation

➢ Make forward pass through the model

➢ Sample next word from LM layer → classification over entire vocabulary

# *What do we need to also generate new text?*



Decoder

Nx

https://teetracker.medium.com/llm-fine-tuning-step-tokenizing-caebb280cfc2

Vaswani et al. "Attention is all you need." *NeurIPS* (2017).

# AUTOREGRESSIVE SAMPLING



Shanahan, Murray, Kyle McDonell, and Laria Reynolds. "Role play with large language models." *Nature* (2023)

# TRAINING LLMS

➢ LLMs for chat applications are usually trained in 3 stages

**1**    Pre-Training

➢ Build capable base-model

➢ General purpose

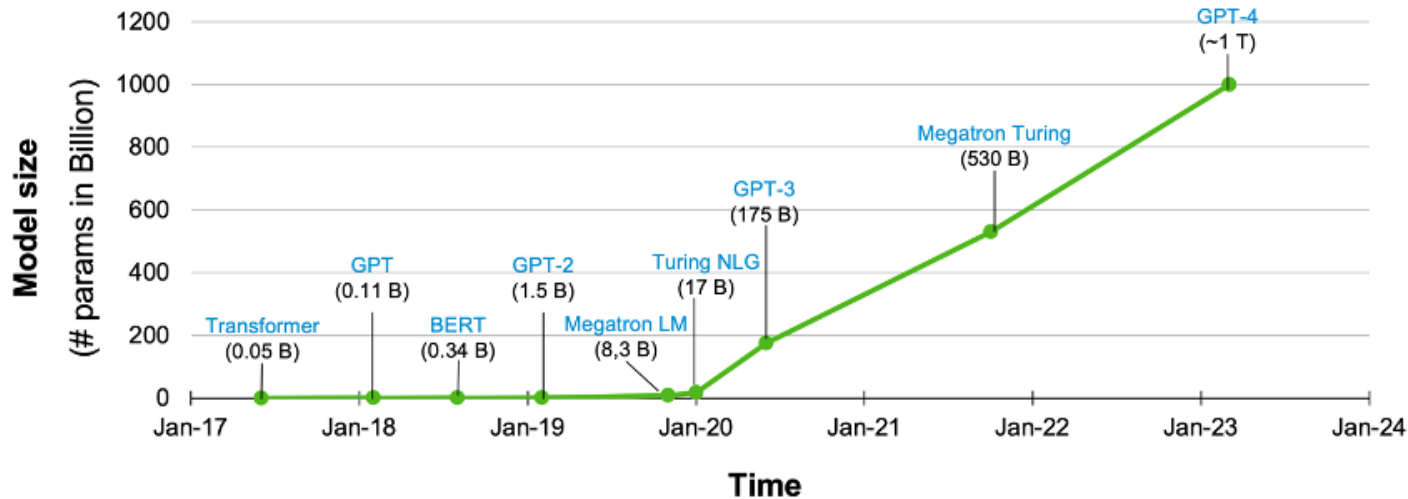➢ Application-specific models are built on top

**2**    Instruction Tuning

➢ Train capability to follow instructions

➢ For example for a chat model

**3**    Preference Tuning

➢ *Align* model to some preferred „*behavior*"

# PRE-TRAINING



https://medium.com/@gladabhi/optimize-cost-to-host-llm-with-sagemaker-async-endpoints-1a6755e458c5

➢ LLM sizes have grown immensely

➢ Significant jump in last 3-4 years

➢ Growing computational requirements

# PRE-TRAINING

➢ Enabled through rapid hardware improvements

➢ Hardware performance/throughput

➢ Similar results at lower precision

➢ Energy optimization



NVIDIA Computex 2024 Keynote

# DATA IS EVERYTHING
# TODO: SPLIT UP



Dataset comparisons

- Limited architectural changes in recent years

- Quality of data far more important

- Large amounts of data on the web

- Filtering & Curation is key

https://huggingface.co/spaces/HuggingFaceFW/blo
gpost-fineweb-v1

# DATA IS EVERYTHING
# TODO: SPLIT UP



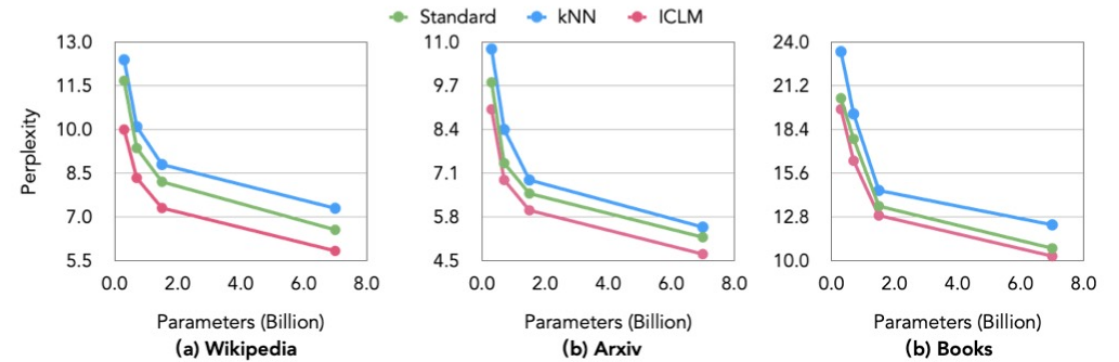Shi, Weijia, et al. "In-Context Pretraining: Language Modeling Beyond Document Boundaries." *arXiv:2310.10638* (2023).

➢ Design choices go beyond what data

➢ How should data be presented during training?

➢ Packing, Context Grouping, *"Curriculum"*

  *can have significant influences*

| Task | Naive Packing | Fewer Truncations Packing | Percentage Increase |
|------|---------------|---------------------------|---------------------|
| truthfulqa_mc | 0.452648 | 0.467687 | 3.32% |
| arc_challenge | 0.517918 | 0.528157 | 1.98% |
| truthful_qa_de | 0.485529 | 0.492979 | 1.53% |
| arc_challenge_de | 0.480375 | 0.493174 | 2.66% |
| hellaswag | 0.776041 | 0.773352 | -0.35% |
| hellaswag_de | 0.655248 | 0.653356 | -0.29% |
| MMLU | 0.573719 | 0.579802 | 1.06% |
| MMLU-DE | 0.504509 | 0.503863 | -0.13% |

https://occiglot.eu/posts/llama-3-german-8b/

> " In contrast to OLMo 1.0, we trained OLMo 1.7 with a two-stage curriculum:
>
> • In the first stage, we train the model from scratch on the Dolma 1.7 dataset. We set a cosine learning rate schedule with a warmup of 2500 steps, a peak learning rate of 3e-4, and a cosine decay to 3e-5 after 3T tokens. We cut off this stage after 2T tokens, when the learning rate is still high.
>
> • At this point we switch to the second stage, in which we train on a curated subset of Dolma 1.7 for another 50B tokens, while linearly decaying the learning rate to 0. We curate this high-quality subset by (1) using all available Wikipedia, OpenWebMath and Flan data, (2) removing Dolma CC, CC News, and Megawika, and (3) rebalancing remaining sources to achieve approximately equal proportions of each. See exact token counts and relative proportions of this second stage mix below.

# How does this definition of "curriculum learning" algin with last weeks lecture?

Olmo 1.7 blogpost
https://blog.allenai.org/olmo-1-7-7b-a-24-point-improvement-on-mmlu-92b43f7d269d

# CURRICULA IN LLM PRE-TRAINING

➢ Very informal/flexible definition

➢ Any change to data-mixture or training setup

➢ Usually hard → easy or

 noisy data → clean/high-quality data

> " In contrast to OLMo 1.0, we trained OLMo 1.7 with a two-stage curriculum:
>
> • In the first stage, we train the model from scratch on the Dolma 1.7 dataset. We set a cosine learning rate schedule with a warmup of 2500 steps, a peak learning rate of 3e-4, and a cosine decay to 3e-5 after 3T tokens. We cut off this stage after 2T tokens, when the learning rate is still high.
>
> • At this point we switch to the second stage, in which we train on a curated subset of Dolma 1.7 for another 50B tokens, while linearly decaying the learning rate to 0. We curate this high-quality subset by (1) using all available Wikipedia, OpenWebMath and Flan data, (2) removing Dolma CC, CC News, and Megawika, and (3) rebalancing remaining sources to achieve approximately equal proportions of each. See exact token counts and relative proportions of this second stage mix below.

Olmo 1.7 blogpost
https://blog.allenai.org/olmo-1-7-7b-a-24-point-improvement-on-mmlu-92b43f7d269d

> In contrast to OLMo 1.0, we trained OLMo 1.7 with a two-stage curriculum:
>
> • In the first stage, we train the model from scratch on the Dolma 1.7 dataset. We set a cosine learning rate schedule with a warmup of 2500 steps, a peak learning rate of 3e-4, and a cosine decay to 3e-5 after 3T tokens. We cut off this stage after 2T tokens, when the learning rate is still high.
>
> • At this point we switch to the second stage, in which we train on a curated subset of Dolma 1.7 for another 50B tokens, while linearly decaying the learning rate to 0. We curate this high-quality subset by (1) using all available Wikipedia, OpenWebMath and Flan data, (2) removing Dolma CC, CC News, and Megawika, and (3) rebalancing remaining sources to achieve approximately equal proportions of each. See exact token counts and relative proportions of this second stage mix below.

## *What is the scoring and pacing function?*

Olmo 1.7 blogpost
https://blog.allenai.org/olmo-1-7-7b-a-24-point-improvement-on-mmlu-92b43f7d269d

# CURRICULA IN LLM PRE-TRAINING

➤ Very informal/flexible definition

➤ Any change to data-mixture or training setup

➤ Usually hard → easy or

   noisy data → clean/high-quality data

➤ Setup often chosen based on best-practices

   (manual decision of pacing)

➤ Scoring based on some quality heuristic

> In contrast to OLMo 1.0, we trained OLMo 1.7 with a two-stage curriculum:
>
> • In the first stage, we train the model from scratch on the Dolma 1.7 dataset. We set a cosine learning rate schedule with a warmup of 2500 steps, a peak learning rate of 3e-4, and a cosine decay to 3e-5 after 3T tokens. We cut off this stage after 2T tokens, when the learning rate is still high.
>
> • At this point we switch to the second stage, in which we train on a curated subset of Dolma 1.7 for another 50B tokens, while linearly decaying the learning rate to 0. We curate this high-quality subset by (1) using all available Wikipedia, OpenWebMath and Flan data, (2) removing Dolma CC, CC News, and Megawika, and (3) rebalancing remaining sources to achieve approximately equal proportions of each. See exact token counts and relative proportions of this second stage mix below.

Olmo 1.7 blogpost
https://blog.allenai.org/olmo-1-7-7b-a-24-point-improvement-on-mmlu-92b43f7d269d

# INSTRUCTION TUNING

➢ Goal: Capabilities beyond completing sentences

    Input: "*Teach me how to bake bread*"    →    Output: "*in an oven.*"

➢ Additional training on structured conversations

➢ General instruction following

➢ Elicit in-context learning (few-shot prompting)

Chat with
**DiscoResearch/Llama3-DiscoLeo-Instruct-8B-v0.1**

# INSTRUCTION TUNING

<|begin_of_text|><|start_header_id|>system<|end_header_id|>
You are a helpful, smart, kind, and efficient AI assistant. You always fulfill the user\'s requests to the best of your ability.<|eot_id|>
<|start_header_id|>human<|end_header_id|>
Teach me how to bake bread<|eot_id|>
<|start_header_id|>assistant<|end_header_id|>
Baking bread can be a rewarding and enjoyable experience, especially when using fresh ingredients. To bake your own bread, you'll need to know a few basic steps and have access to a few kitchen tools.
Here's how you can make bread:
1.Gather your ingredients: [...]
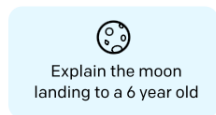2.Prepare your dough: [...]
3.Knead the dough: [...]
[...]
With practice, you'll get the hang of it and can experiment with different flavors and shapes<|eot_id|>
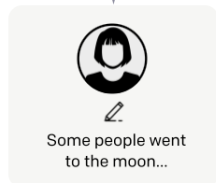
# PREFERENCE TUNING

**Step 1**

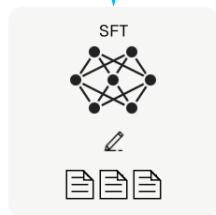**Collect demonstration data, and train a supervised policy.**

A prompt is sampled from our prompt dataset.
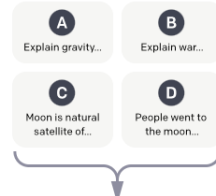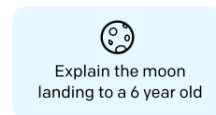
A labeler demonstrates the desired output behavior.

This data is used to fine-tune GPT-3 with supervised learning.

**Step 2**

**Collect comparison data, and train a reward model.**

A prompt and several model outputs are sampled.

A labeler ranks the outputs from best to worst.

This data is used to train our reward model.

**Step 3**

**Optimize a policy against the reward model using reinforcement learning.**

A new prompt is sampled from the dataset.

The policy generates an output.

The reward model calculates a reward for the output.

The reward is used to update the policy using PPO.

Direct Preference Optimization (DPO)

x: "write me a poem about the history of jazz"

preference data    maximum likelihood    final LM

Rafailov, Rafael, et al. "Direct preference optimization: Your language model is secretly a reward model."NeurIPS (2023)

Ouyang, Long, et al. "Training language models to follow instructions with human feedback." *NeurIPS* (2022)

# PREFERENCE TUNING

*How does preference tuning relate to*

*domain adaptation & transfer learning?*

# PREFERENCE TUNING

➢ Goal: Instill (human) preference on outputs for same input

➢ No changes to the underlying task

➢ Still adaptation of the target domain

➢ We change the distribution of model outputs/predictions

# DOMAIN ADAPTATION

➢ Again: Terminology is blurry between different fields

➢ LLM domain adaption focuses on use cases (domain)

➢ E.g. Adapt LLM for a different language

# CONTINUAL PRE-TRAINING

➢ Take pre-trained Language Model

➢ Continue pre-training stage on new data/other language

*Why would we want to use a pre-trained*

*model instead of training from scratch?*



occiglot-eu5-7b-v0.1 ✎                                    updated Mar 7

First release of 7B LLMs models for the 5 biggest European languages. All models initialised from
mistral-7b-v0.1.

occiglot/occiglot-7b-eu5
Text Generation • Updated 10 days ago • ↓ 11k • ♥ 25

occiglot/occiglot-7b-eu5-instruct
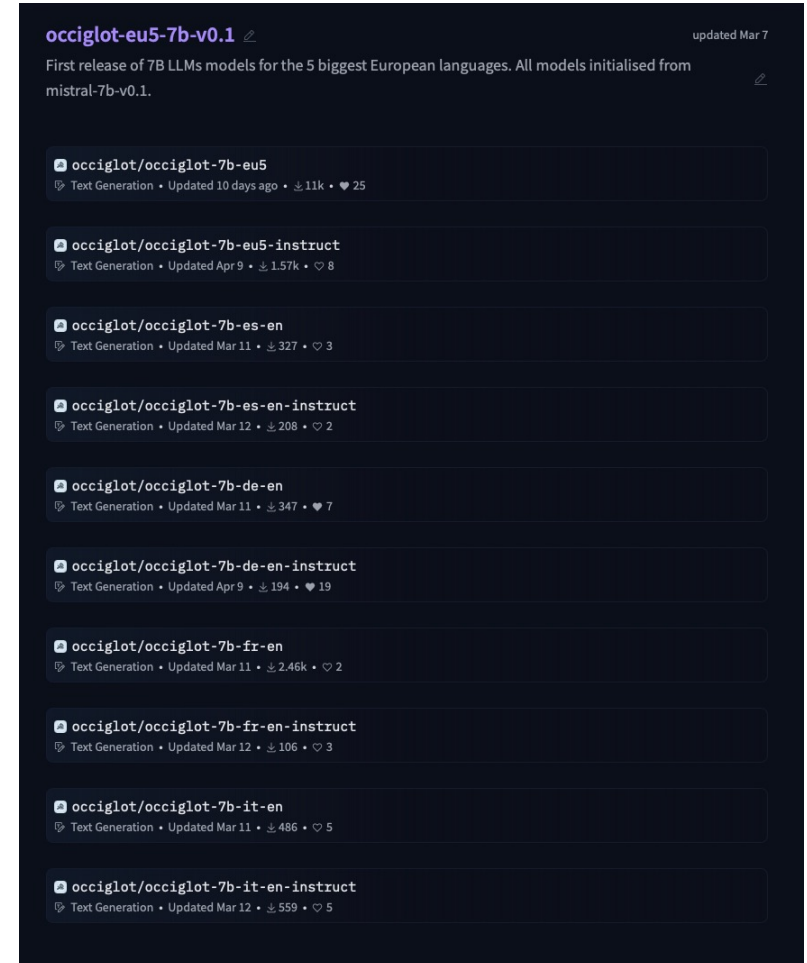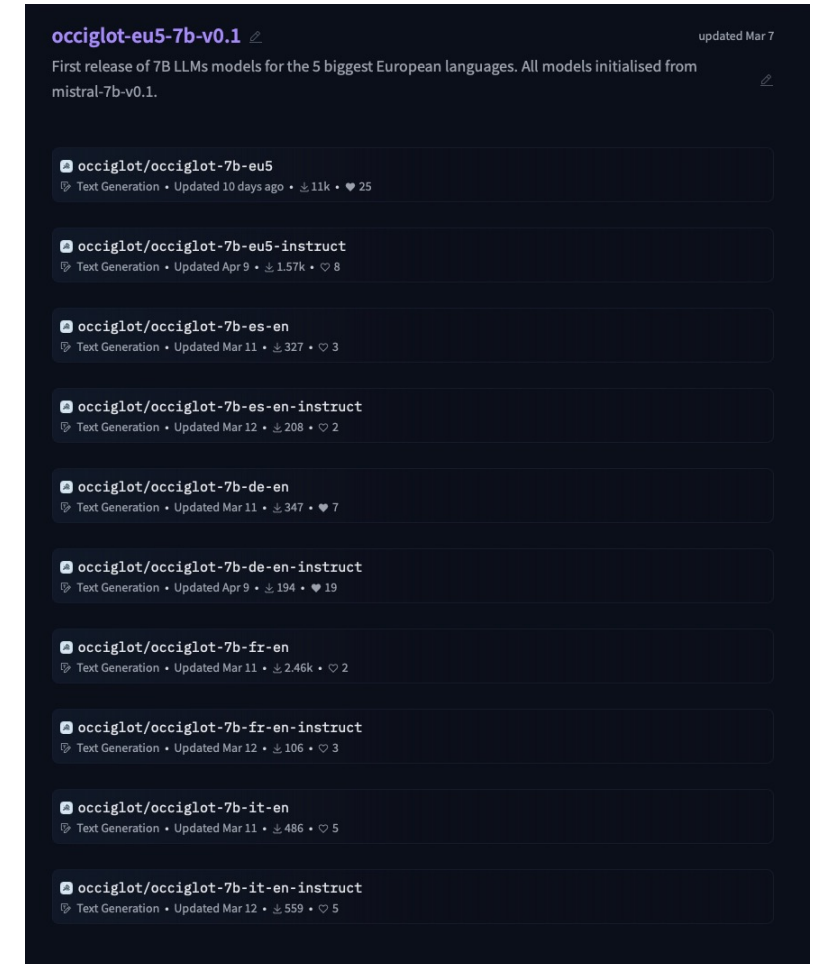Text Generation • Updated Apr 9 • ↓ 1.57k • ♡ 8

occiglot/occiglot-7b-es-en
Text Generation • Updated Mar 11 • ↓ 327 • ♡ 3

occiglot/occiglot-7b-es-en-instruct
Text Generation • Updated Mar 12 • ↓ 208 • ♡ 2

occiglot/occiglot-7b-de-en
Text Generation • Updated Mar 11 • ↓ 347 • ♥ 7

occiglot/occiglot-7b-de-en-instruct
Text Generation • Updated Apr 9 • ↓ 194 • ♥ 19

occiglot/occiglot-7b-fr-en
Text Generation • Updated Mar 11 • ↓ 2.46k • ♡ 2

occiglot/occiglot-7b-fr-en-instruct
Text Generation • Updated Mar 12 • ↓ 106 • ♡ 3

occiglot/occiglot-7b-it-en
Text Generation • Updated Mar 11 • ↓ 486 • ♡ 5

occiglot/occiglot-7b-it-en-instruct
Text Generation • Updated Mar 12 • ↓ 559 • ♡ 5

# CONTINUAL PRE-TRAINING

➢ Take pre-trained Language Model

➢ Continue pre-training stage on new data/other language

➢ LLM training is costly & hard

➢ Language Modelling can transfer between languages

➢ Maybe we want to retain performance on original language

➢ Often not enough mono-lingual data for training from scratch



occiglot-eu5-7b-v0.1 ✎                                                    updated Mar 7

First release of 7B LLMs models for the 5 biggest European languages. All models initialised from mistral-7b-v0.1.

🖼 occiglot/occiglot-7b-eu5
🗐 Text Generation · Updated 10 days ago · ⬇ 11k · ❤ 25

🖼 occiglot/occiglot-7b-eu5-instruct
🗐 Text Generation · Updated Apr 9 · ⬇ 1.57k · ♡ 8

🖼 occiglot/occiglot-7b-es-en
🗐 Text Generation · Updated Mar 11 · ⬇ 327 · ♡ 3

🖼 occiglot/occiglot-7b-es-en-instruct
🗐 Text Generation · Updated Mar 12 · ⬇ 208 · ♡ 2

🖼 occiglot/occiglot-7b-de-en
🗐 Text Generation · Updated Mar 11 · ⬇ 347 · ❤ 7

🖼 occiglot/occiglot-7b-de-en-instruct
🗐 Text Generation · Updated Apr 9 · ⬇ 194 · ❤ 19

🖼 occiglot/occiglot-7b-fr-en
🗐 Text Generation · Updated Mar 11 · ⬇ 2.46k · ♡ 2

🖼 occiglot/occiglot-7b-fr-en-instruct
🗐 Text Generation · Updated Mar 12 · ⬇ 106 · ♡ 3

🖼 occiglot/occiglot-7b-it-en
🗐 Text Generation · Updated Mar 11 · ⬇ 486 · ♡ 5

🖼 occiglot/occiglot-7b-it-en-instruct
🗐 Text Generation · Updated Mar 12 · ⬇ 559 · ♡ 5

# CHALLENGES W/ CLOSED MODELS

➤ The best models come from large companies

➤ Very sparse/limited information on data & training

*What challenges does that*

*entail for continual pre-training*



https://www.analyticsvidhya.com/blog/2024/04/meta-llama-a-breakthrough-in-open-ai-models/



https://de.wikipedia.org/

# CHALLENGES W/ CLOSED MODELS

➢ The best models come from large companies

➢ Very sparse/limited information on data & training

➢ Memory buffers for replay are hard to construct

➢ What data may have been seen already?

➢ Careful rewarming/choice of learning rate



https://www.analyticsvidhya.com/blog/2024/04/meta-llama-a-breakthrough-in-open-ai-models/



https://de.wikipedia.org/

# CHALLENGES

➢ Large-scale model training is very different than previous examples in this lecture

*What are some distinct challenges/differences*

*arising for LLM training?*

# CHALLENGES

➢ Large-scale model training is very different than previous examples in this lecture

➢ A lot of machine learning engineering required

➢ Usually only trained for one epoch

➢ You may only be able to effort one training run at full scale

   ➢ No extensive hyper-parameter ablations

   ➢ Instead try to establish scaling laws from smaller models

➢ Scale of the data may lead to unintended inclusion of problematic material

# EVALUATION

➢ How to evaluate a (pre-trained) LLM

➢ Naïve idea: Based on the Language Modeling objective

$$\mathcal{P}(w_n|w_{n-1}\dots w_1) = \prod \mathcal{P}(w_i|w_{i-1})$$

➢ Model the (un-)certainty of the LLM for a given sequence (entropy) → inverse

$$\frac{1}{\prod \mathcal{P}(w_i|w_{i-1})}$$

➢ Normalize by word length (geometric average) & write in log-scale for numeric stability

$$exp\left(\frac{1}{n}\sum log\big(\mathcal{P}(w_i|w_{i-1})\big)\right)$$

# EVALUATION

➤ Problems with Perplexity

 ➤ How to you chose representative test set?

 ➤ Usefulness of comparison between different approaches might be limited

 ➤ Does not capture desired probabilities of chat bots very well

# EVALUATION

## AI2 Reasoning Challenge (ARC)

Which property of a mineral can be determined just by looking at it?

- (A) luster
- (B) mass
- (C) weight
- (D) hardness

Clark, Peter, et al. "Think you have solved question answering? try arc, the ai2 reasoning challenge." *arXiv:1803.05457* (2018)

## HellaSwag

**Category**: Shaving (ActivityNet; In-domain)
A bearded man is seen speaking to the camera and making several faces. the man

a) then switches off and shows himself via the washer and dryer rolling down a towel and scrubbing the floor. (0.0%)
b) then rubs and wipes down an individual's face and leads into another man playing another person's flute. (0.0%)
c) is then seen eating food on a ladder while still speaking. (0.0%)
d) then holds up a razor and begins shaving his face. (100.0%)

Zellers, Rowan, et al. "Hellaswag: Can a machine really finish your sentence?." ACL (2019(

## Massive Multitask Language Understanding (MMLU)

**Microeconomics**
One of the reasons that the government discourages and regulates monopolies is that
(A) producer surplus is lost and consumer surplus is gained. ✗
(B) monopoly prices ensure productive efficiency but cost society allocative efficiency. ✗
(C) monopoly firms do not engage in significant research and development. ✗
(D) consumer surplus is lost with higher prices and lower levels of output. ✓

**Conceptual Physics**
When you drop a ball from rest it accelerates downward at 9.8 m/s². If you instead throw it downward assuming no air resistance its acceleration immediately after leaving your hand is
(A) 9.8 m/s² ✓
(B) more than 9.8 m/s² ✗
(C) less than 9.8 m/s² ✗
(D) Cannot say unless the speed of throw is given. ✗

Hendrycks, Dan, et al. "Measuring massive multitask language understanding." ICLR (2021).

## GMS8K

**Problem**: Beth bakes 4, 2 dozen batches of cookies in a week. If these cookies are shared amongst 16 people equally, how many cookies does each person consume?
**Solution**: Beth bakes 4 2 dozen batches of cookies for a total of 4*2 = <<4*2=8>>8 dozen cookies
There are 12 cookies in a dozen and she makes 8 dozen cookies for a total of 12*8 = <<12*8=96>>96 cookies
She splits the 96 cookies equally amongst 16 people so they each eat 96/16 = <<96/16=6>>6 cookies
**Final Answer**: 6

Cobbe, Karl, et al. "Training verifiers to solve math word problems." *arXiv:2110.14168* (2021).

# EVALUATION

https://huggingface.co/spaces/open-llm-leaderboard/open_llm_leaderboard
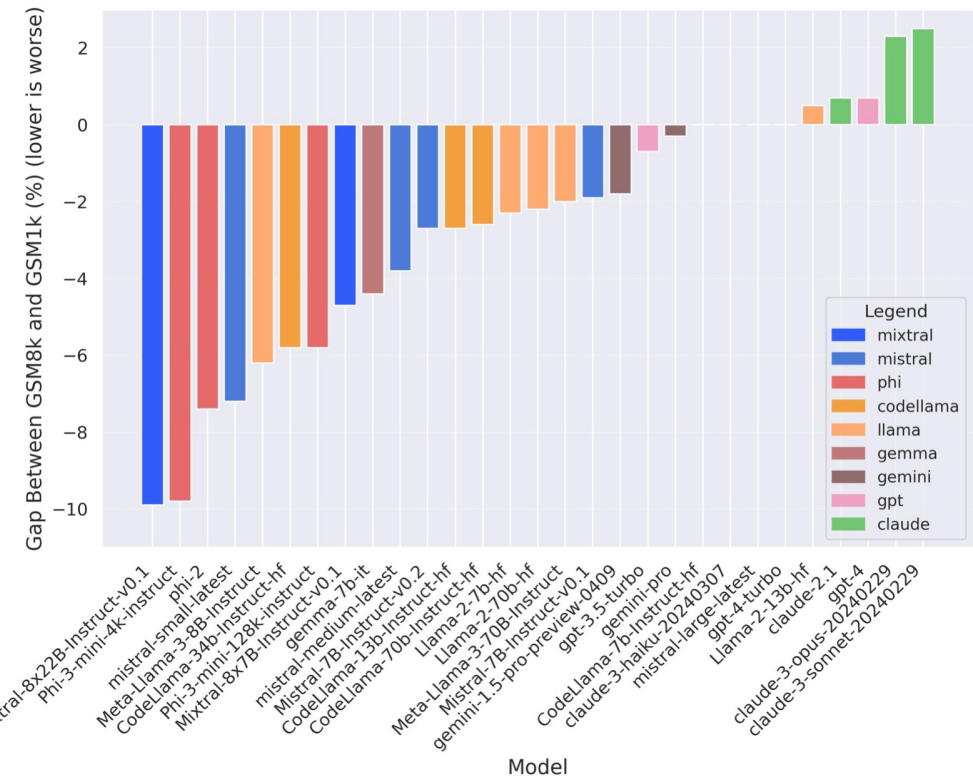
# EVALUATION

*What are potential problems of these open benchmarks?*

# EVALUATION

➢ Leaderboards can be gamed

➢ Developers might overfit to a certain benchmark

➢ Webcrawling leads to (unintentional)

   training data contamination





Zhang, Hugh, et al. "A careful examination of large language model performance on grade school arithmetic." *arXiv:2405.00332* (2024).

https://x.com/cHHillee/status/1635790330854526981

# CONTRIBUTE @ OCCIGLOT

manuel@occiglot.org

Discord

# QUESTIONS?



Generated with Dalle-3/ChatGPT