

# Large Language Models

Machine Learning Course - CS-433

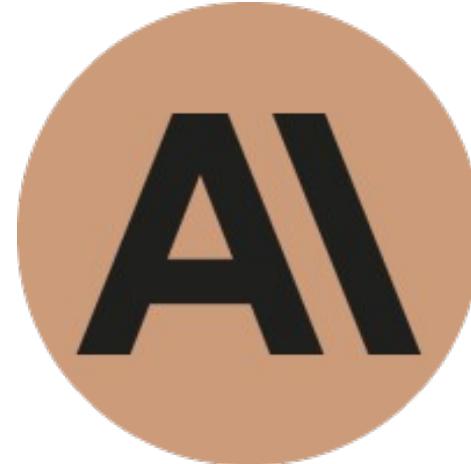
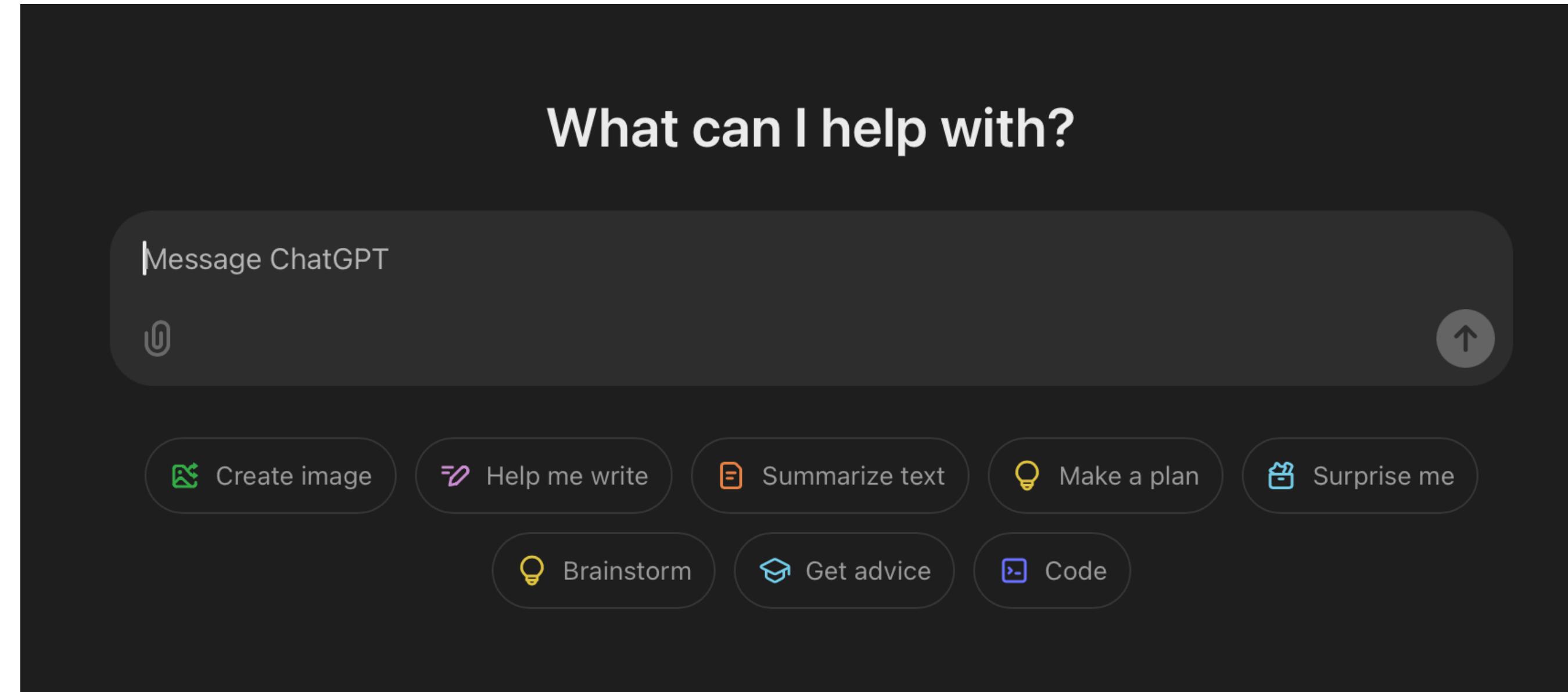
3 Dec 2025

Robert West

(Slide credits: Martin Jaggi, Alex Hägele, Francesco D'Angelo)



# LLMs are Everywhere



Gemini



Qwen



MISTRAL  
AI\_

LLaMA  
by Meta

**Goal:** Understand what it takes to build a large language model

# Outline

- **Part 1: Building Blocks**
  - Transformers
  - Language Modeling
  - Tokenizers
  - Autoregressive Inference
- **Part 2: Pretraining**
  - Data
  - Distributed Training
  - Intuitions
- **Part 3: Posttraining and Model Capabilities**
  - Zero-Shot and Few-Shot In-Context Learning
  - Instruction Finetuning
  - Optimizing for human preferences (PPO/RLHF)
  - LLM evaluation

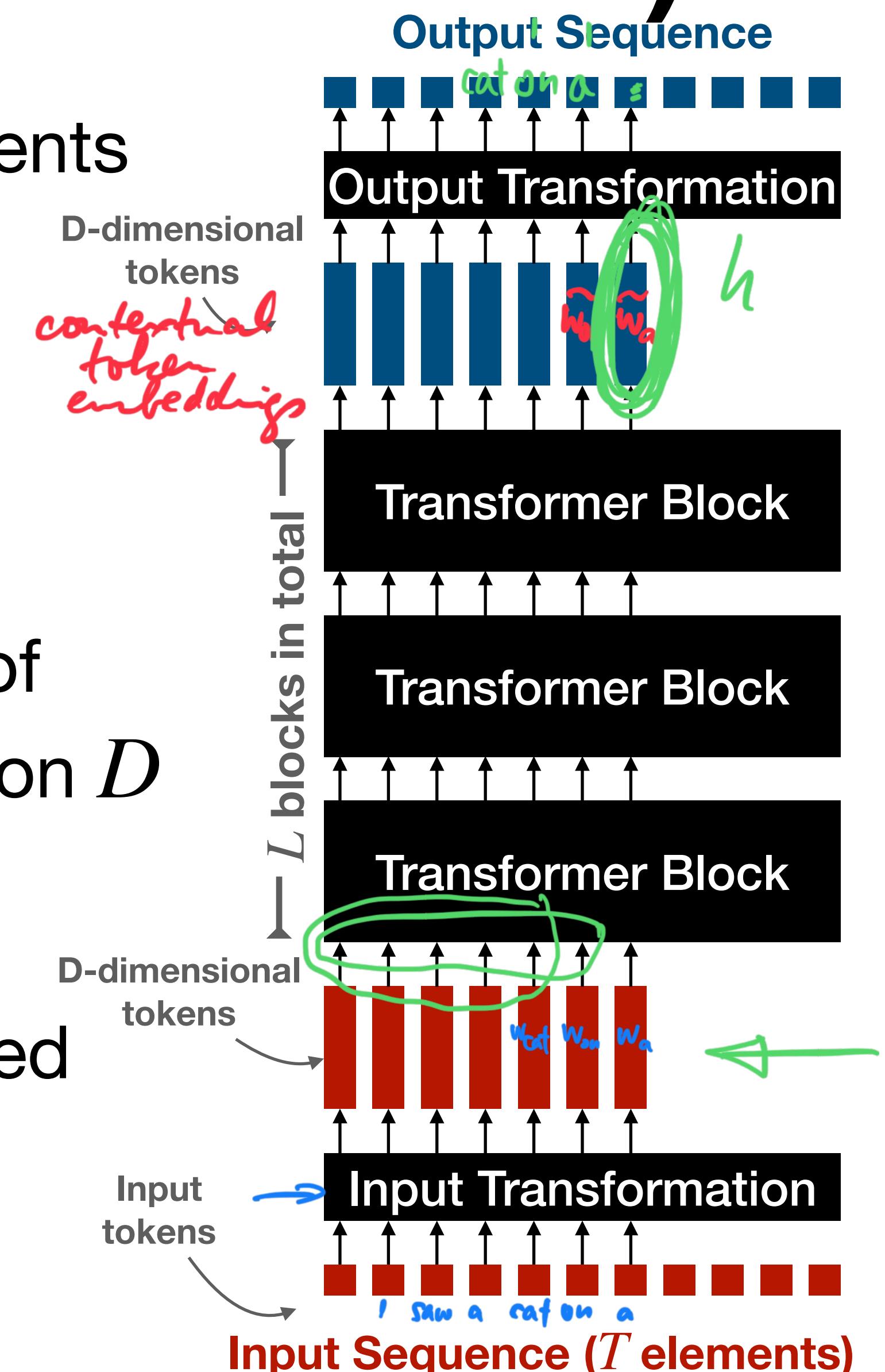
# Recap: Transformers (Lecture 9)

**Input transformation:** Converts the input sequence elements into real-valued vector representations (a.k.a. **t**o**k**e**n**s):

- maps a one-hot word vector to a real-valued vector
- maps an image patch into a flat vector

**Transformer block:** Transforms a sequence of  $T$  vectors of dimension  $D$  into a new sequence of  $T$  vectors of dimension  $D$  using **self-attention** and **MLP sub-blocks**

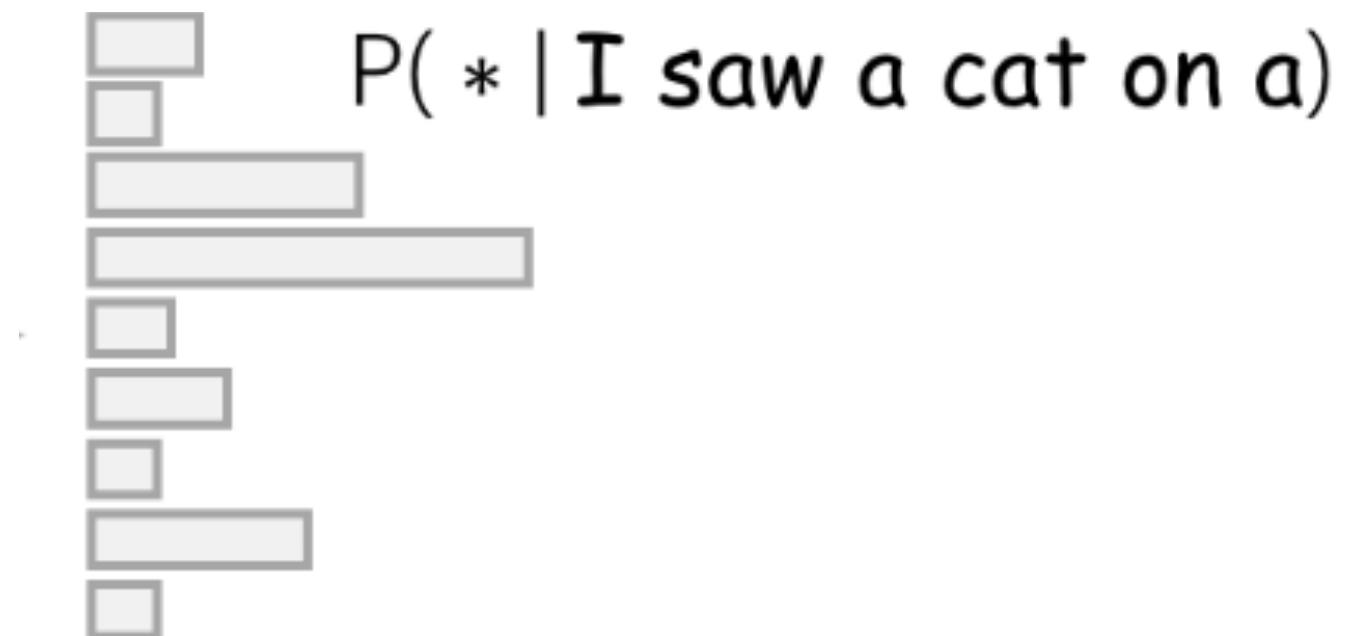
**Output transformation:** Converts the vectors to the desired output format (e.g., a label for each sequence element)



# Language Modeling

- Language Models describe distributions over sequences of text

$$\bullet p(\text{I saw a cat on a mat}) = p(x_1, \dots, x_t)$$



- Simplest factorization: predict next word (= *token*)

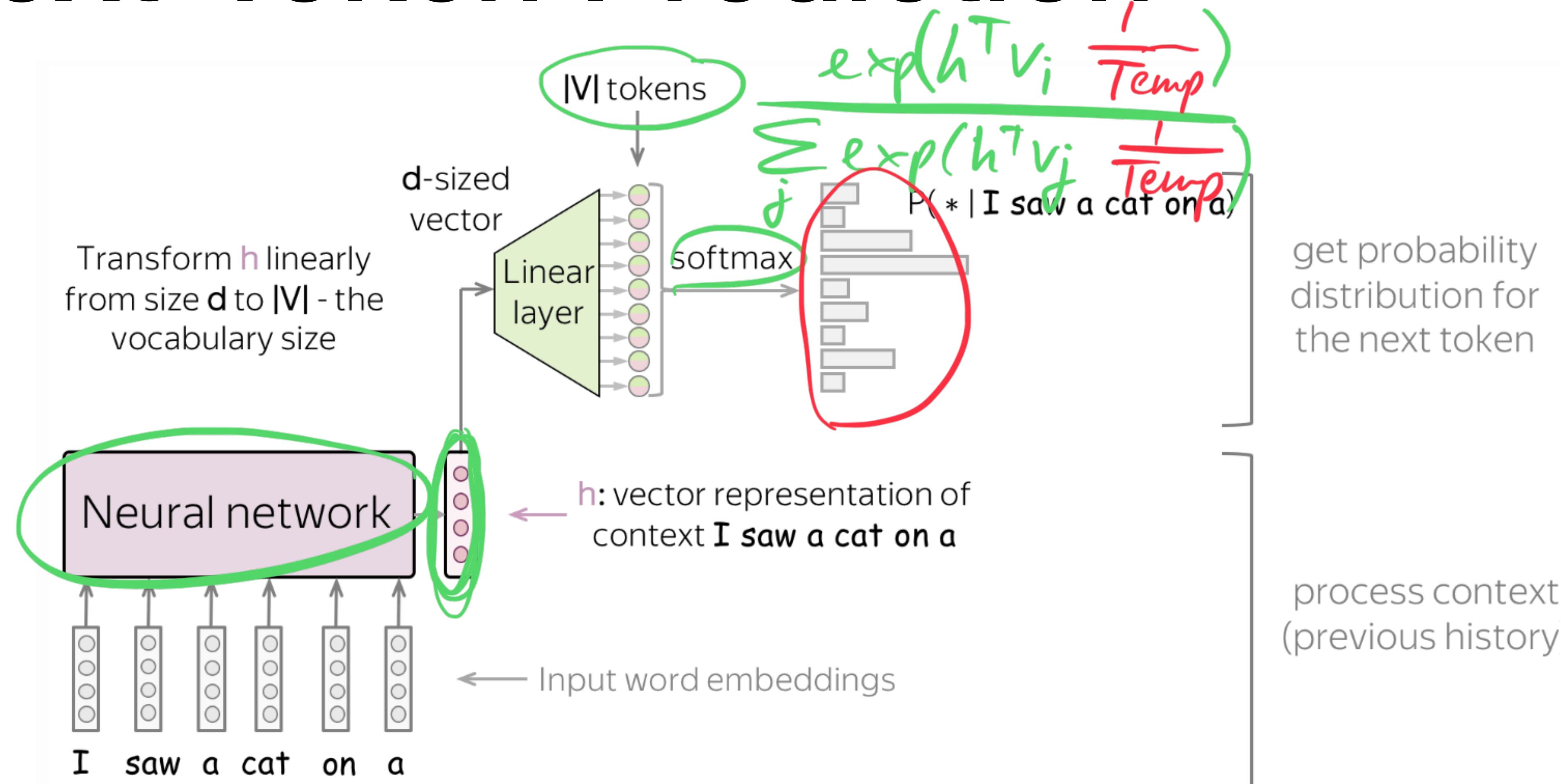
$$= \bullet p_{\theta}(x_t | x_1, \dots, x_{t-1}) \cdot p_{\theta}(x_{t-1} | x_1, \dots, x_{t-2}) \cdots p_{\theta}(x_2 | x_1) \cdot p_{\theta}(x_1)$$



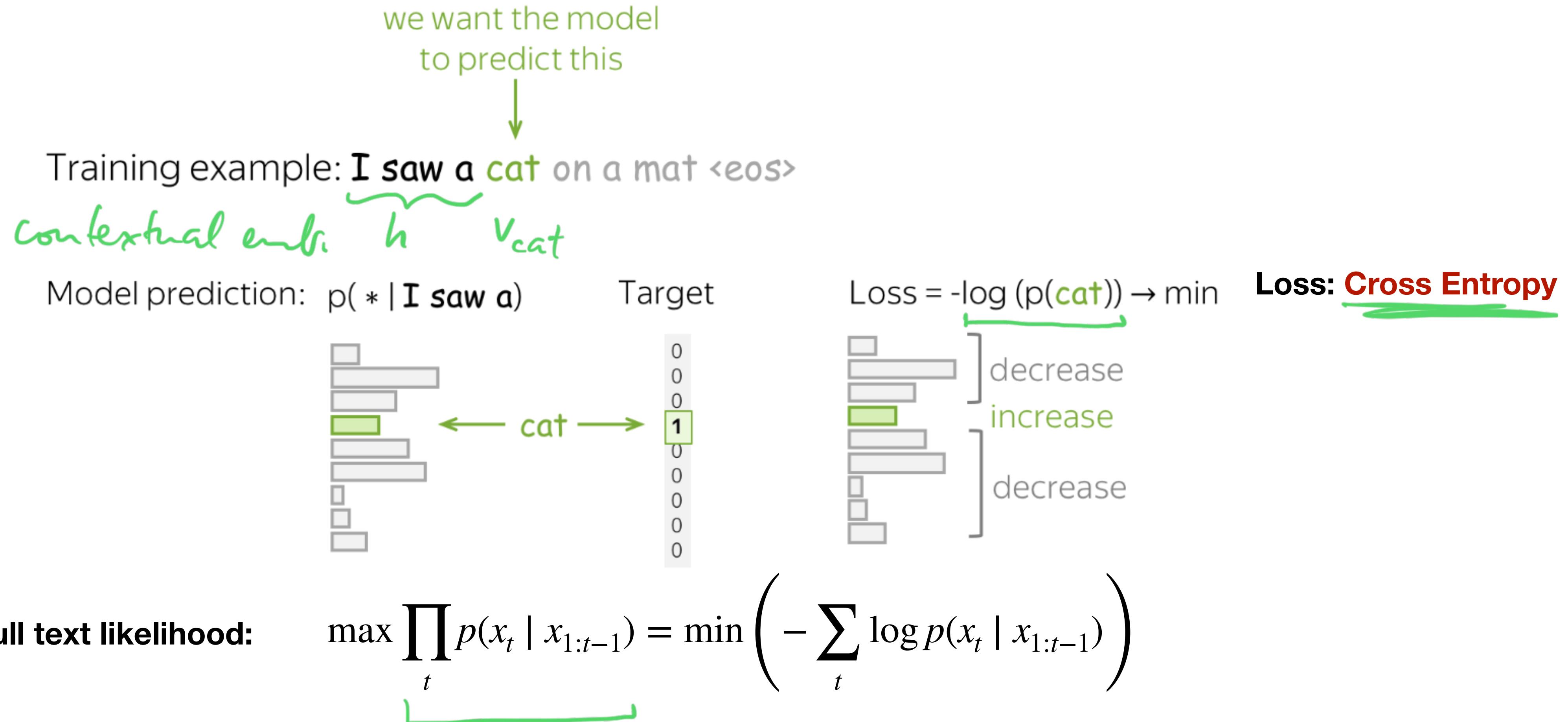
- e.g.  $p(\text{mat} | \text{I saw a cat on a}) = 0.002$

# Next-Token Prediction

$V$ : vocabulary  
 $|V|$ : size of  $V$



# Next-Token Prediction



# Tokenizers

**Tokenization:** Split the input text into a sequence of *input tokens* (typically word fragments + some special symbols) according to some predefined *tokenizer procedure*:

## The Tokenizer Playground

Experiment with different tokenizers (running locally in your browser).

The screenshot shows the Tokenizer Playground interface. At the top, a dropdown menu is set to "Llama 3". Below it, a text input field contains the sentence "Transformers are awesome<|eot\_id|>". A green bracket highlights the word "awesome". A red arrow points from this bracket to a detailed view below. This view shows the input text with colored segments: "Transformers" is purple, "are" is yellow, "awesome" is red, and "<|eot\_id|>" is blue. To the left of the text, it says "TOKENS 5" and "CHARACTERS 34". To the right, it shows the corresponding token IDs: "[9140, 388, 527, 12738, 128009]". Below this, a text box states: "Each token corresponds to some number  $i \in \{1, \dots, N_{vocab}\}$ ". At the bottom, there are three radio buttons: "Text" (selected), "Token IDs", and "Hide".

# Tokenizers

## The Tokenizer Playground

Experiment with different tokenizers (running locally in your browser).

Llama 3

Transformers are awesome<|eot\_id|>

TOKENS 5 CHARACTERS 34

Transformers are awesome<|eot\_id|>

[9140, 388, 527, 12738, 128009]

TOKENS 5 CHARACTERS 34

corresponds to some number  $i \in \{1, \dots, N_{vocab}\}$ )

Text  Token IDs  Hide

Text  Token IDs  Hide

## Why?

- More general than words (e.g. typos, code, diff. languages)
- Tradeoff between two extreme vocabularies: *all possible words* (infinite) or only characters (very long sequences and no semantic information)
- Idea: Tokens are common subsequences (subword tokens)

# Tokenizers

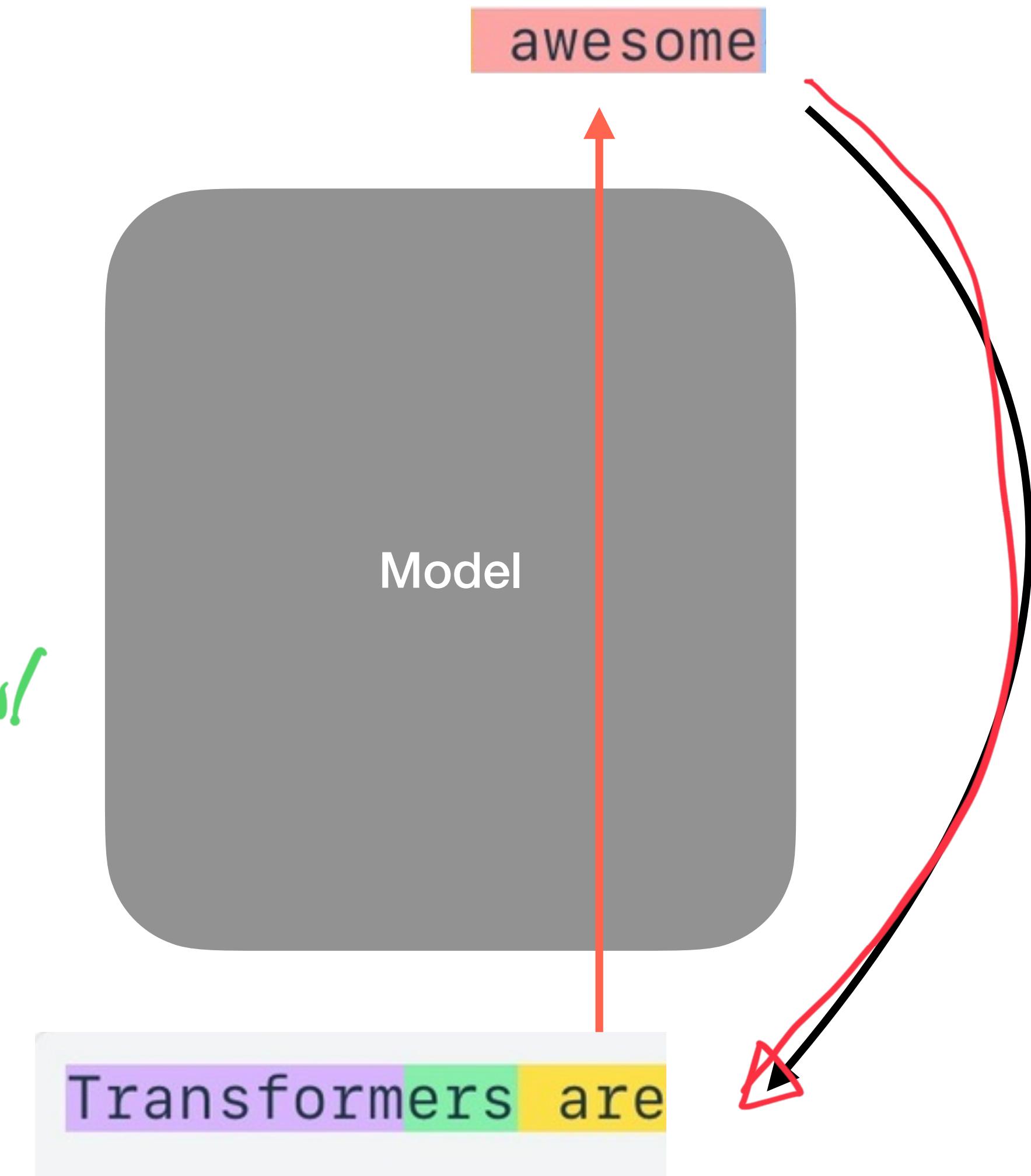
**Tokenization:** split the input text into a sequence of *input tokens* (typically word fragments + some special symbols) according to some predefined *tokenizer procedure*

- **Example: Byte Pair Encoding (BPE). Train steps:**
  1. Take large corpus of text (your pretraining corpus)
  2. Start with one token per character
  3. Merge common pairs of tokens into a token
  4. Repeat until desired vocab size or all merged

The Tokenizer is independent of the model training and is trained beforehand for a fixed vocabulary (with special tokens like <user>, <eot>, etc.

# Autoregressive Inference

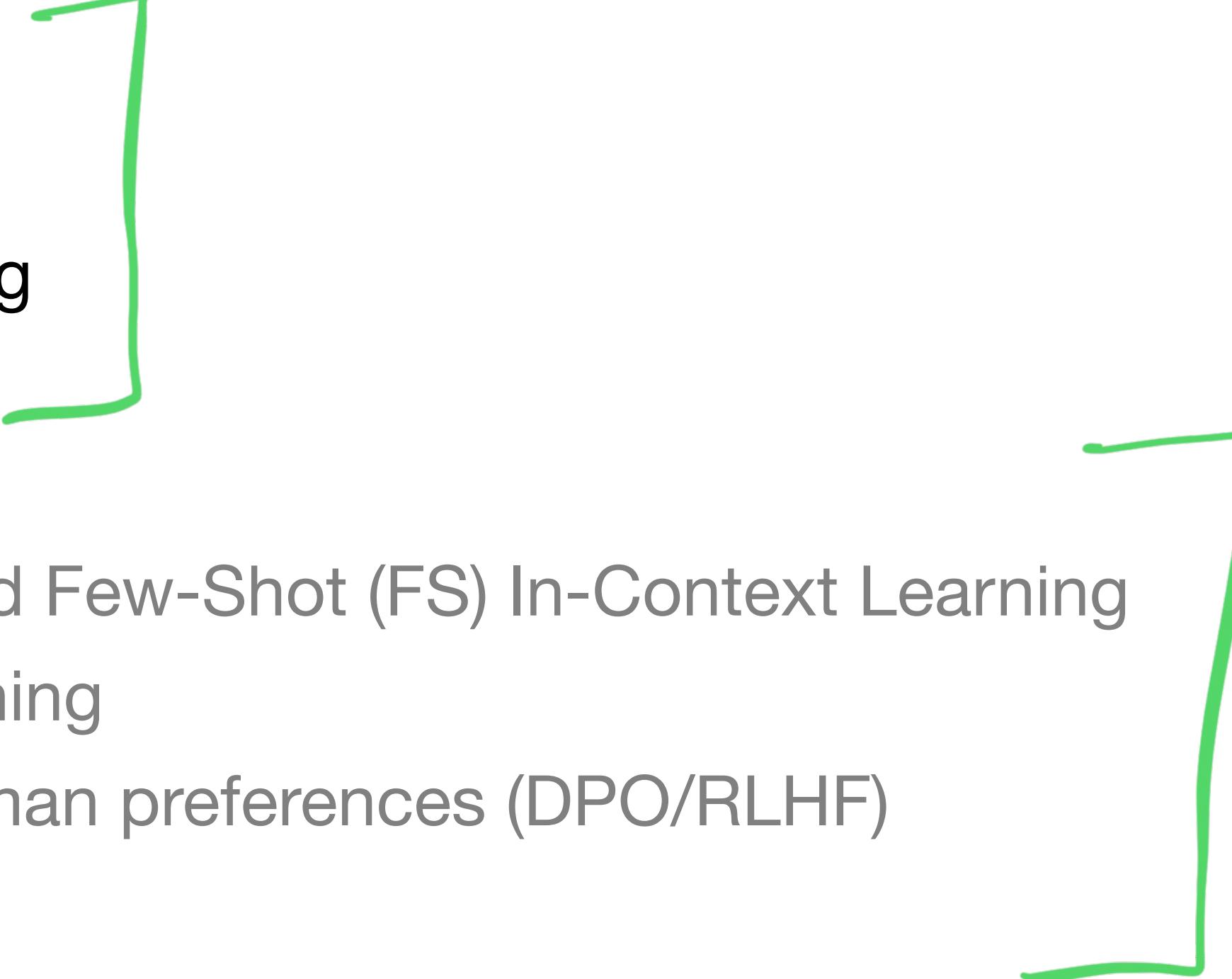
- **Task:** Given context (tokens), output sentence
- **Steps:**
  1. Forward pass through model
  2. Obtain probabilities for next tokens
  3. Select one next token (sample or argmax)  
*Sample w/  
temp. 0*
  4. Add sampled token to context
  5. Repeat



# Recap: Building Blocks

- Transformers
- Language Modeling
- Next-Token Prediction
- Tokenizers
- Autoregressive Inference

# Outline

- **Part 1: Building Blocks**
    - Language Modeling
    - Tokenizers
    - Autoregressive Inference
  - **Part 2: Pretraining**
    - Data
    - Distributed Training
    - Intuitions
  - **Part 3: Posttraining**
    - Zero-Shot (ZS) and Few-Shot (FS) In-Context Learning
    - Instruction Finetuning
    - Optimizing for human preferences (DPO/RLHF)
    - LLM evaluation
- 

# Pretraining vs. Posttraining

	Pretraining	Posttraining
Data Quantity	Massive, ~the internet (Billions-trillions of words)	Small, millions of examples
Data Quality	Low(er)	High
Data Source	Webcrawls, Papers, Textbooks, Github, ... ...	(Often) Human
Goal	General Learning, Knowledge	Make model usable, give specific skills + character, alignment, ...

# Data

the secret sauce

- **Goal of Pretraining:** General purpose model
  - Train on massive quantities of text
    - Latest Llama 3.1: 15T (*trillion*) tokens
- **Challenges**
  - Maximize coverage, diversity
  - Maximize quality, correctness
  - Find right measures of data quality
  - Efficient processing of trillions of tokens (multiple TB of data)

# Data

the secret sauce

- **Goal of Pretraining:** general purpose model
  - Train on massive quantities of text
- **Where to find data**
  - **Common Crawl:** starting point
  - **Code:** Github etc.
  - **Curated:**
    - Websites
    - Books
  - **(Synthetic Data: new trend, utilizing existing models)**

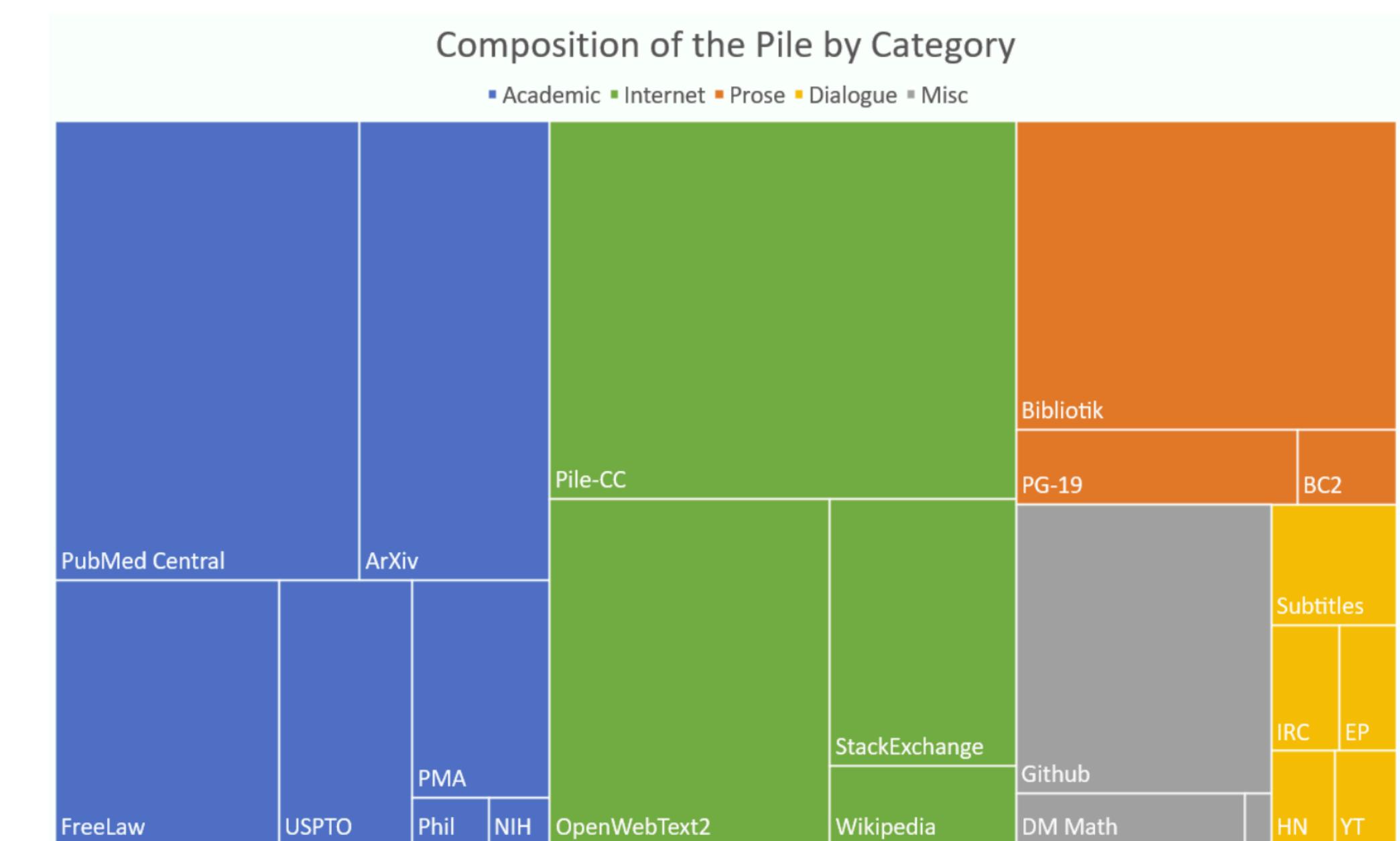


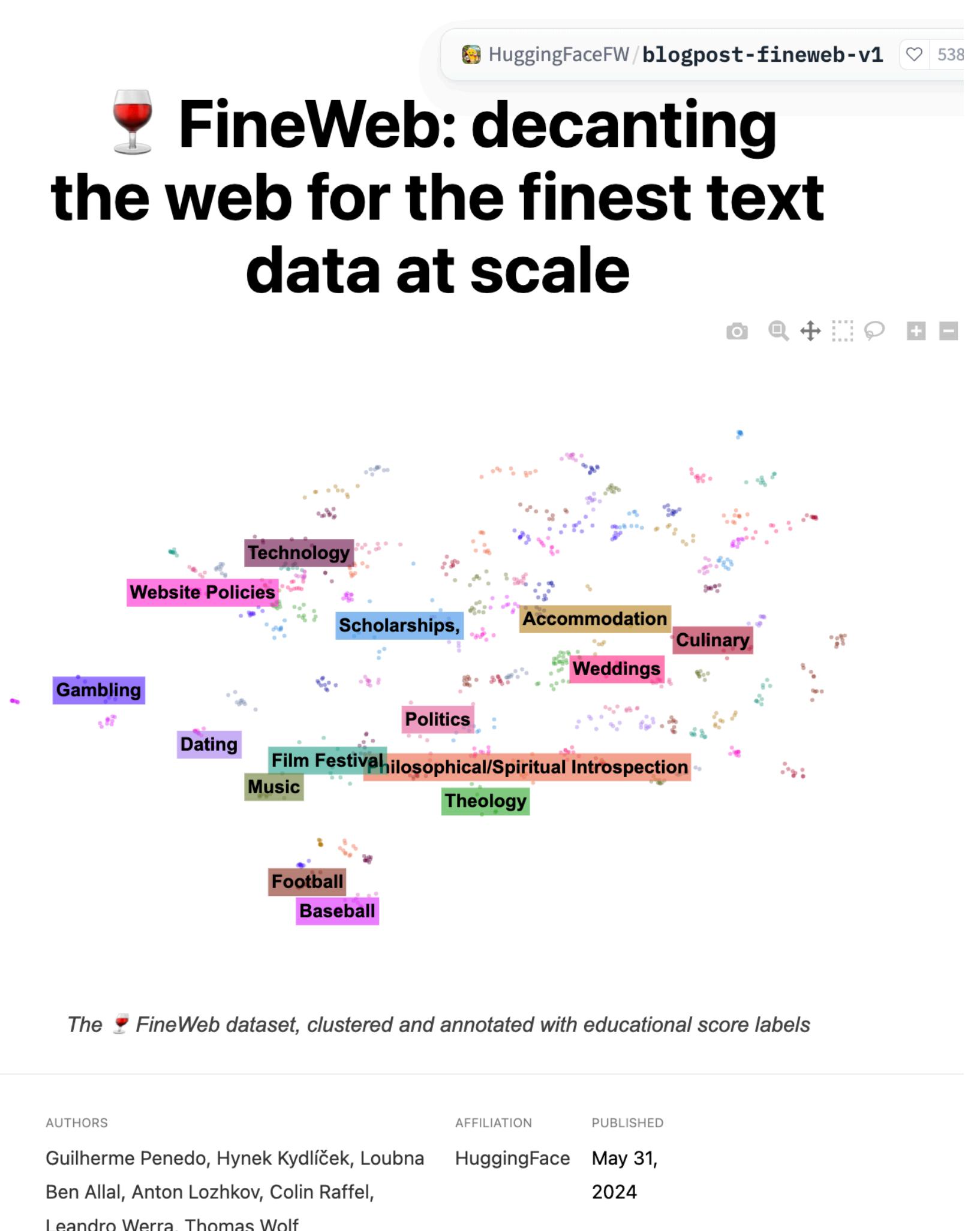
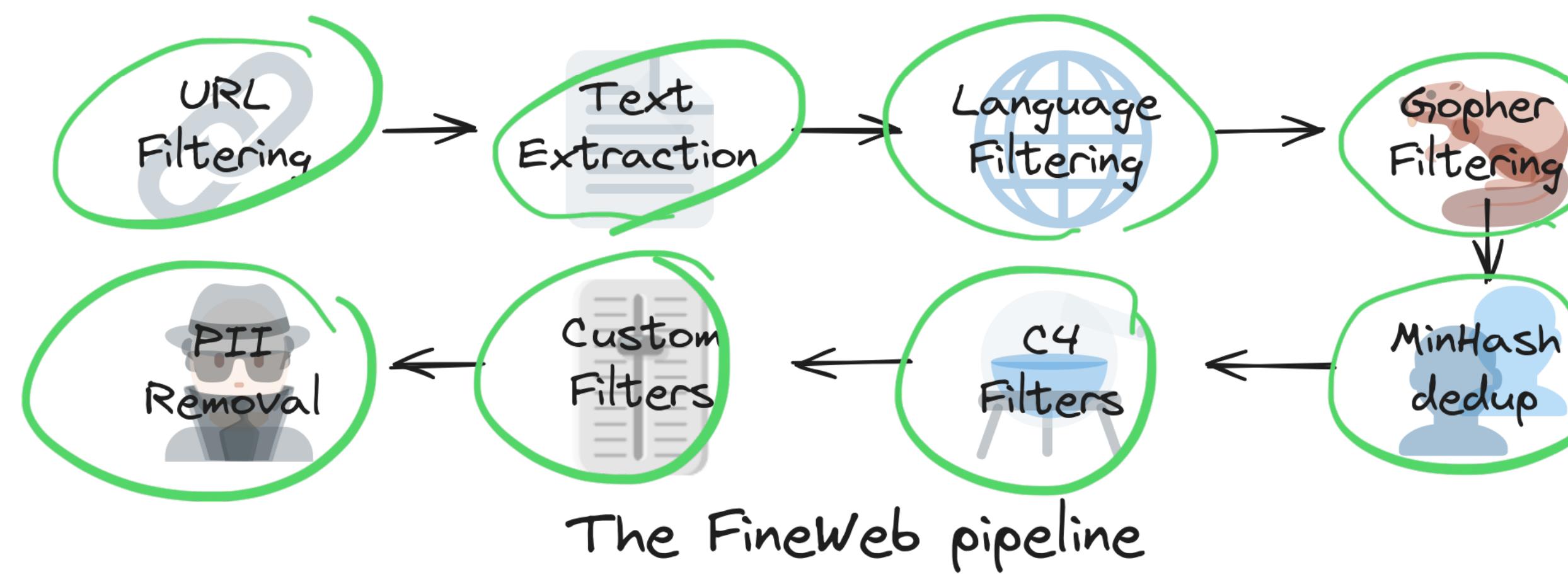
Figure 1: Treemap of Pile components by effective size.

# Data

the secret sauce

## Example: FineWeb

- Based on CommonCrawl
- 15T tokens
- 44TB disk space
- Heavy filtering
- Transparent processing pipeline



# Intuition for LLMs

**Intuition 1:** Next-word prediction on large data can be viewed as **massively multi-task learning**.

Task	Example sentence in pre-training
Grammar	In my free time, I like to {run, banana}
Lexical Semantics	I went to the zoo to see giraffes, lions, and {zebras, spoon}
World Knowledge	The capital of Denmark is {Copenhagen, London}
Sentiment Analysis	Movie review: I was engaged and on the edge of my seat the whole time. The movie was {good, , bad}
Translation	The word for “pretty” in Spanish is {bonita, hola}
Math	First grade arithmetic exam: $3 + 8 + 4 = \{15, 11\}$
Coding	def fibonacci(n): [...]

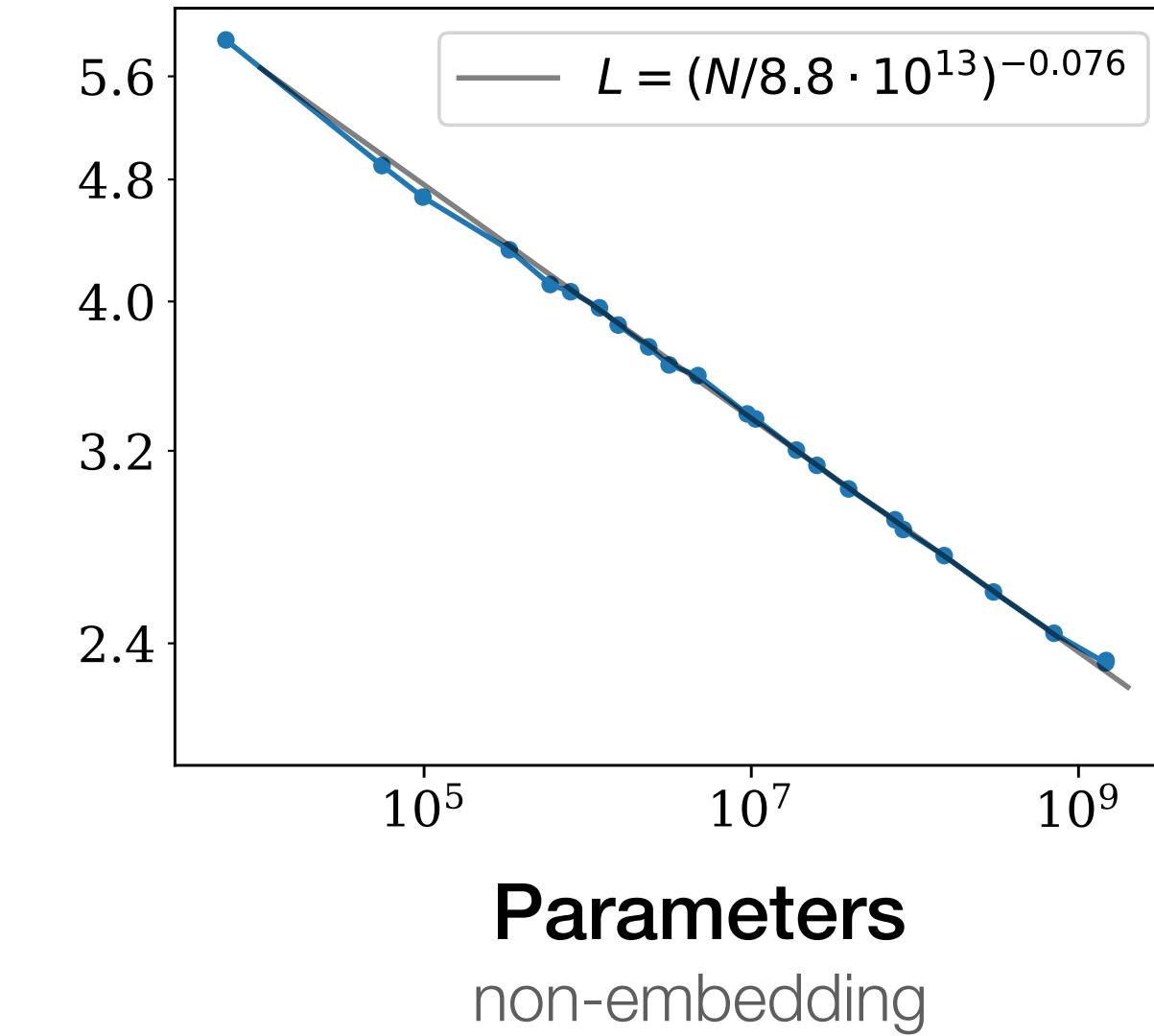
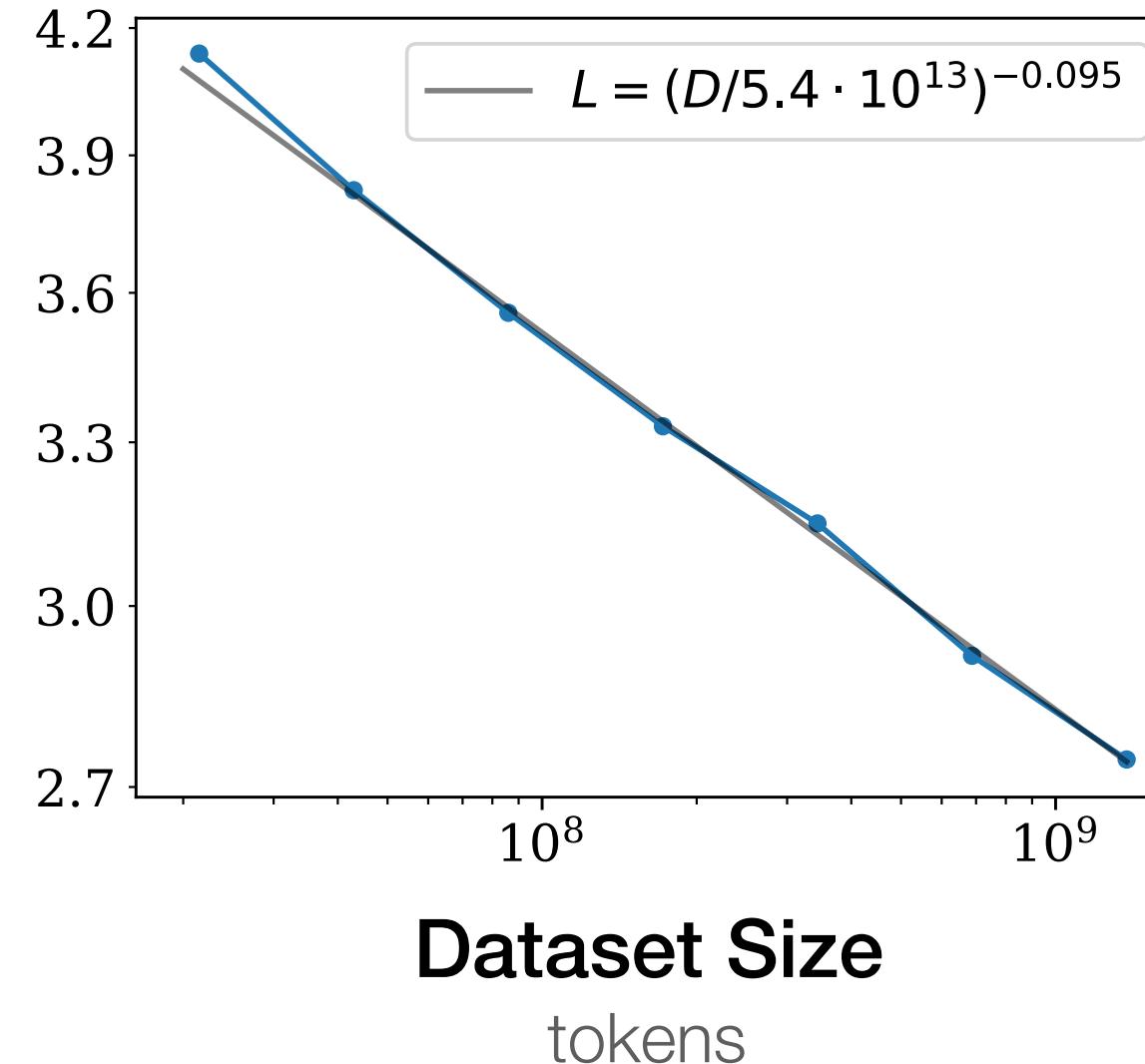
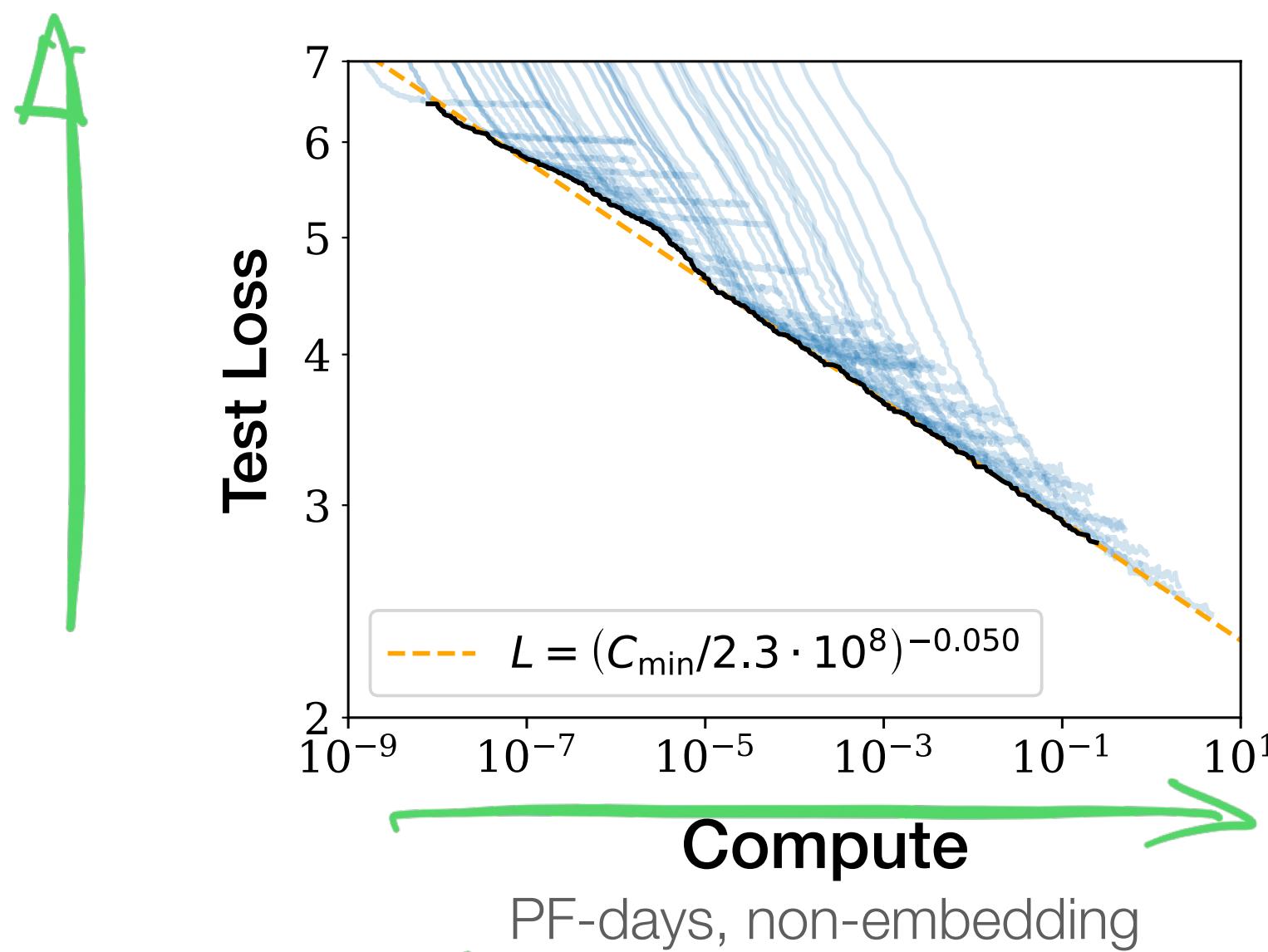
**Intuition 2:** Next-word prediction is **extremely general**.

Learning <input, output> relationships can be cast as next-word prediction.  
(See in-context learning later!)

# Intuitions for LLMs

**Intuition 3:** The Power of **Scale** (connects to 1 and 2).

**Compute = Model Size \* Data**



$10^{12}$  FLOPS      peta-flop

*“Language modeling performance (loss) improves smoothly as we increase the model size, dataset size, and amount of compute for training.” = Scaling Laws*

# Recap: Pretraining

- Importance of data
- Intuitions for pretraining
  - Next-token prediction as general multitask learning
  - Scaling laws

# Outline

- **Part 1: Building Blocks**
  - Transformers
  - Language Modeling
  - Tokenizers
  - Autoregressive Inference
- **Part 2: Pretraining**
  - Data
  - Distributed Training
  - Intuitions
- **Part 3: Posttraining and Model Capabilities**
  - Zero-Shot and Few-Shot In-Context Learning
  - Instruction Finetuning
  - Optimizing for human preferences (PPO/RLHF)
  - LLM evaluation

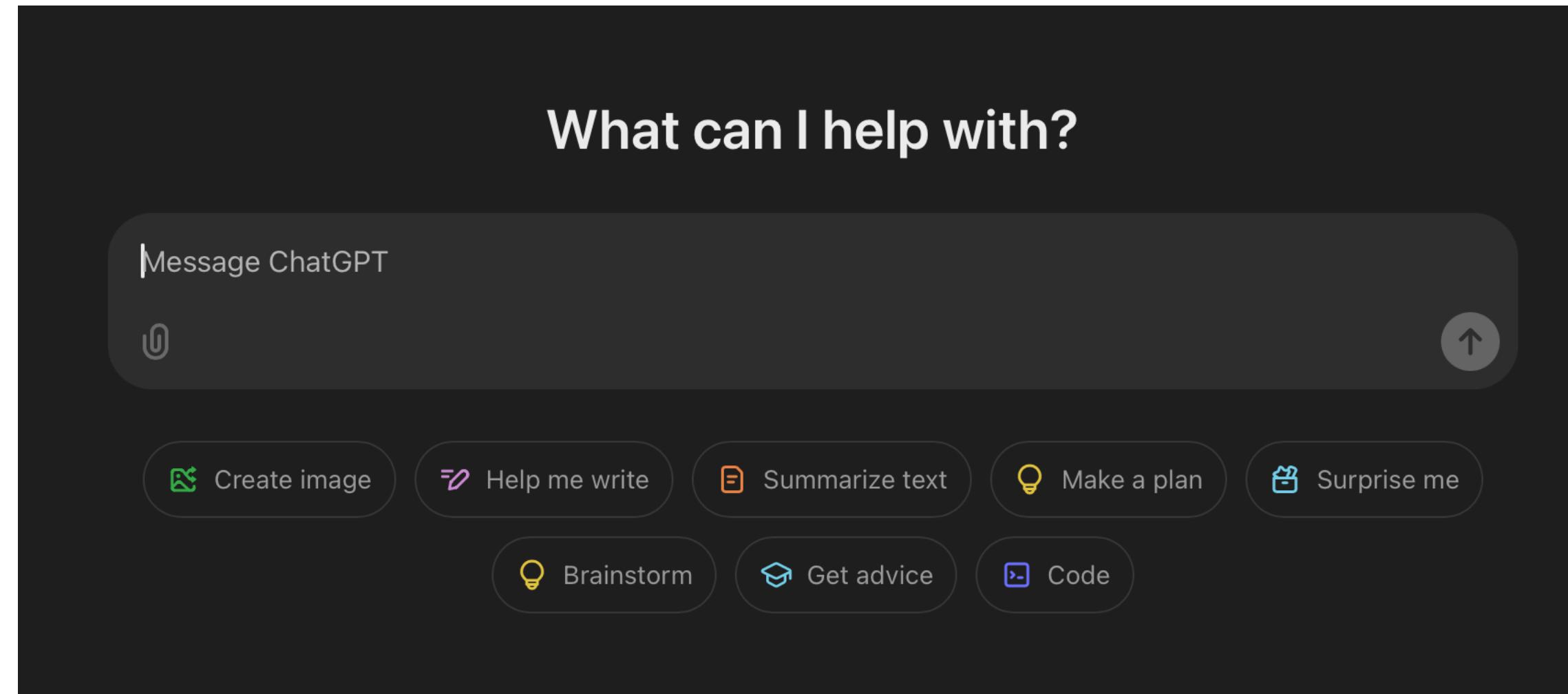


# Language models as multitask assistants?

How do we get from next token prediction:

*EPFL is located in .....*

To this:



# Emergent Abilities of LLMs: GPT1 (2018)

Let's go back to the first GPT model from OpenAI:

- Decoder-only transformer with 12 Layers
- **117M** parameters
- Trained on BooksCorpus (7000 books and **4.6GB** text)

Showed how pretraining on a large corpus of text is effective for fine-tuning downstream tasks.

Drawbacks:

- For each task: New dataset and finetuning needed
- No generalization to new-tasks

# Emergent abilities of LLMs: GPT2 (2019)

More parameters and more data, GPT2:

- Decoder-only transformer with 48 Layers
- 117M → **1.5B** parameters
- Trained on WebText (4.6GB → **40GB**)

English: I am here  
French:

Emergent zero-shot abilities

- Fine-tuning is not needed anymore
- Pre-trained alone achieves SOTA for all tasks

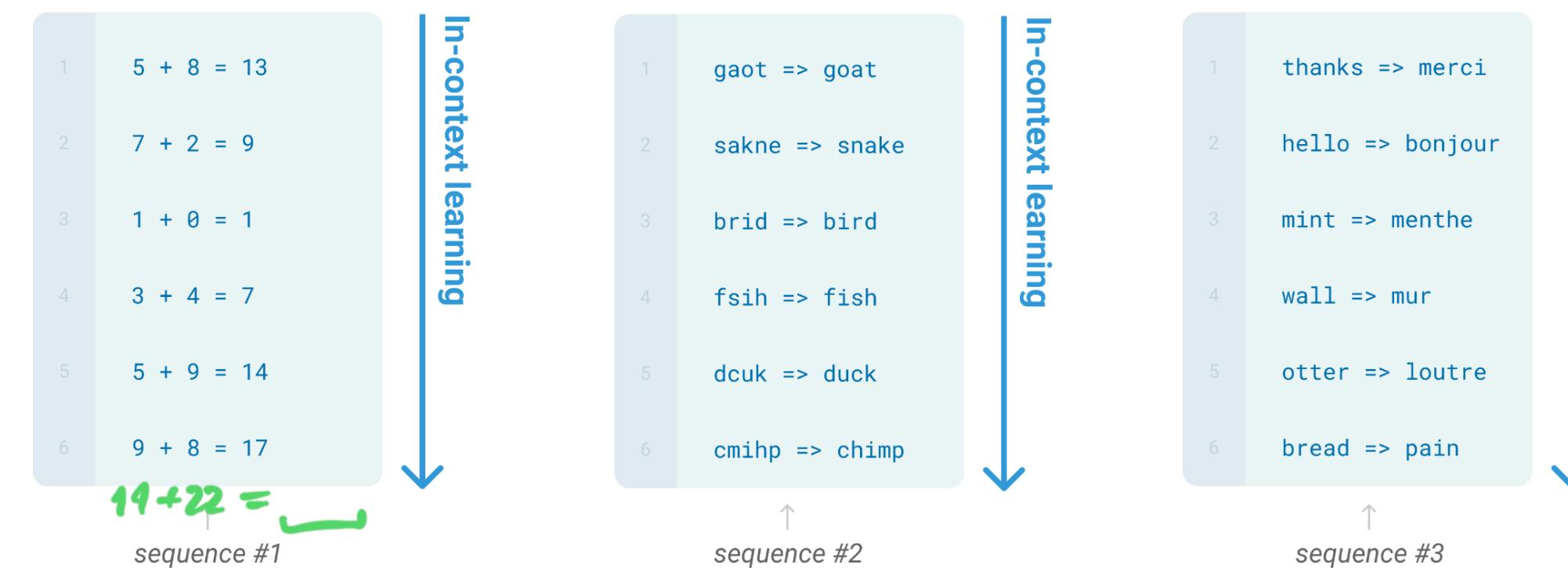
# Emergent abilities of LLMs: GPT3 (2020)

Even more parameters and data, GPT3 (Brown et al. 2020):

- Decoder-only transformer with 96 Layers
- $117M \rightarrow 1.5B \rightarrow \text{175B}$  parameters
- Trained on WebText ( $4.6GB \rightarrow 40GB \rightarrow \text{600GB}$ )

Language models are few-shot learners → In-context learning

- Specify a new task adding examples in the prompt



# Emergent abilities of LLMs: GPT3 (2020)

## 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.



## The three settings we explore for in-context learning

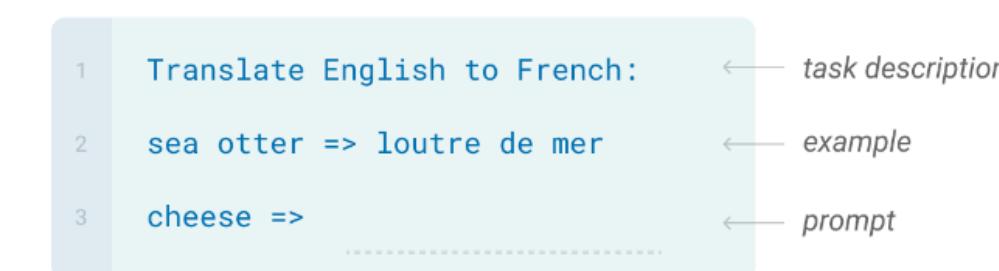
### Zero-shot

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



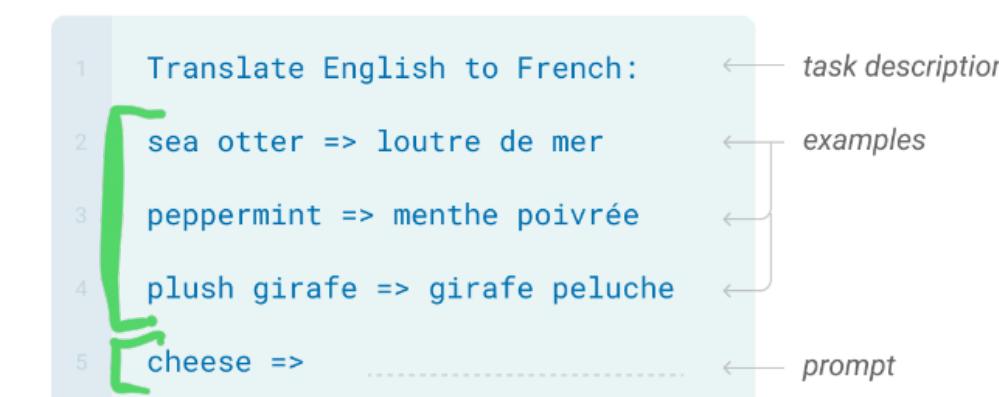
### One-shot

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



### Few-shot

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



# Limitations of In-Context Learning

Some tasks are too difficult to solve:

- Showing **examples** in the prompt is not sufficient
- Especially for tasks involving multi-step **reasoning**



Multiply 999 by 867

The product of multiplying 999 by 867 is ~~824133~~

866133

GPT4 with 5 examples achieves:

- 92% accuracy on 3x3
- 4% accuracy on 4x4
- 2% accuracy on 5x5

# Chain-of-Thought (CoT) Prompting

Showing some examples of reasoning steps in the prompt enables solving more complex tasks.

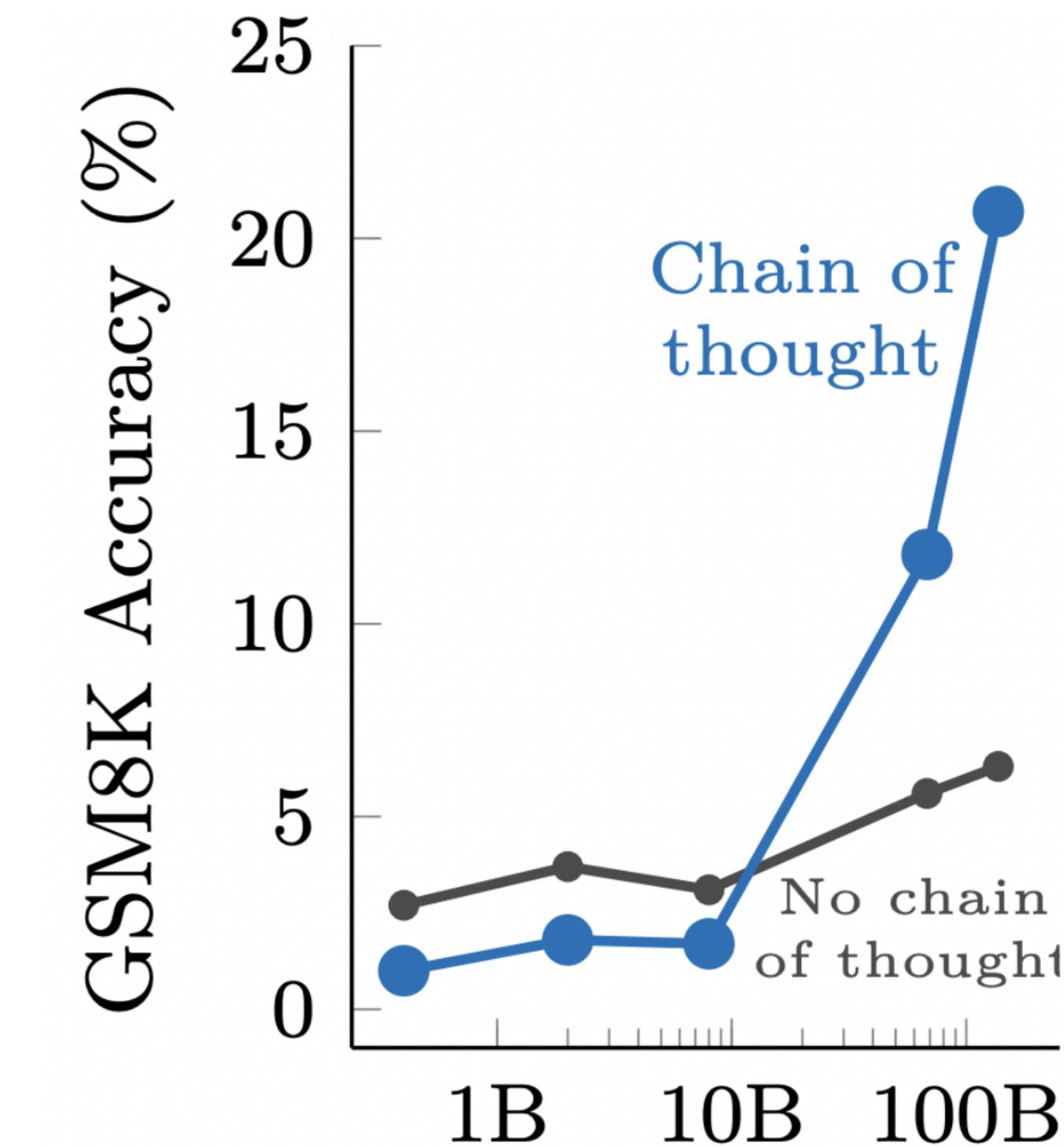
Multiply 999 by 866

Let's perform the multiplication step by step:  
Let's multiply 999 by the digit in the ones place of 866,  
which is 6.

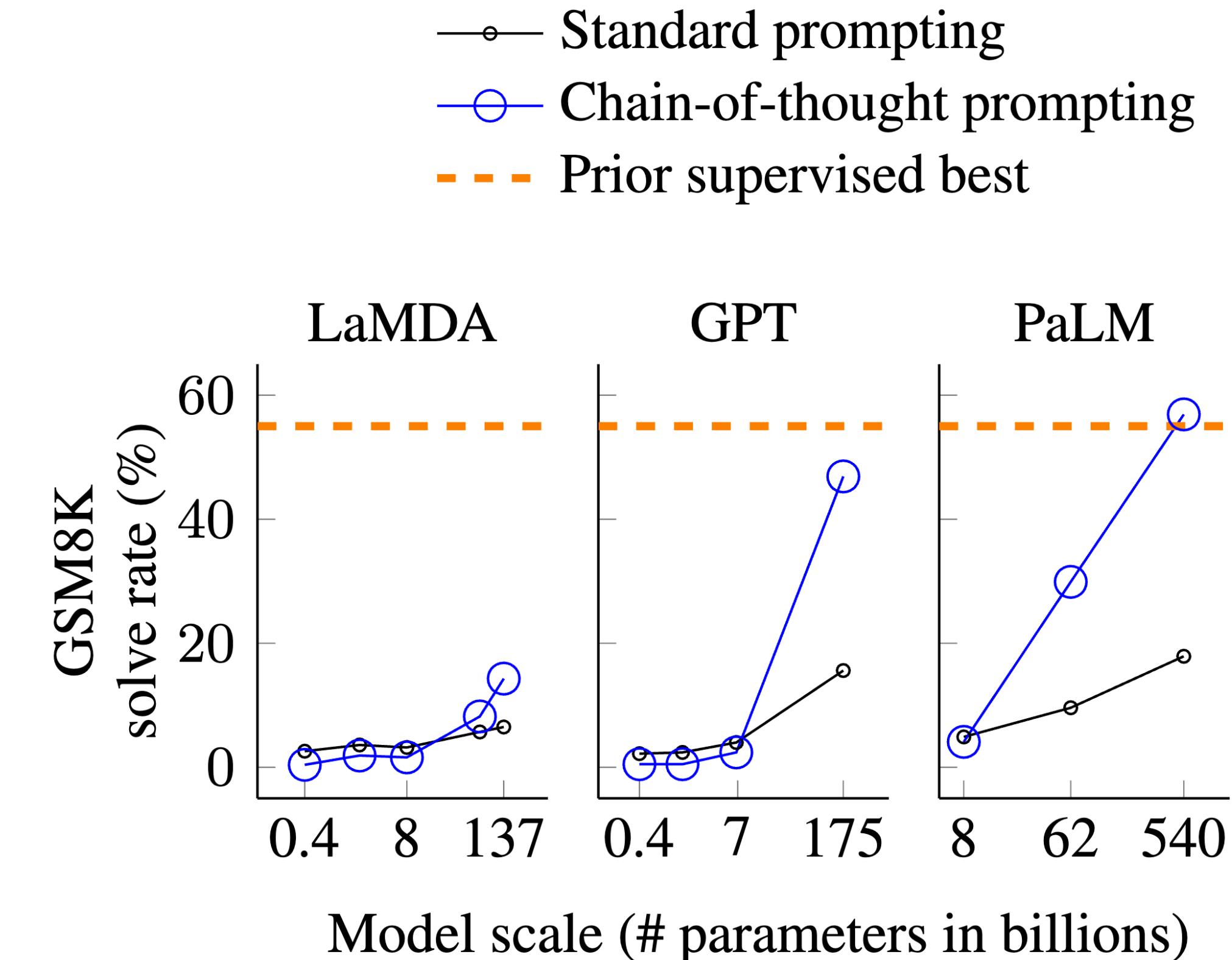
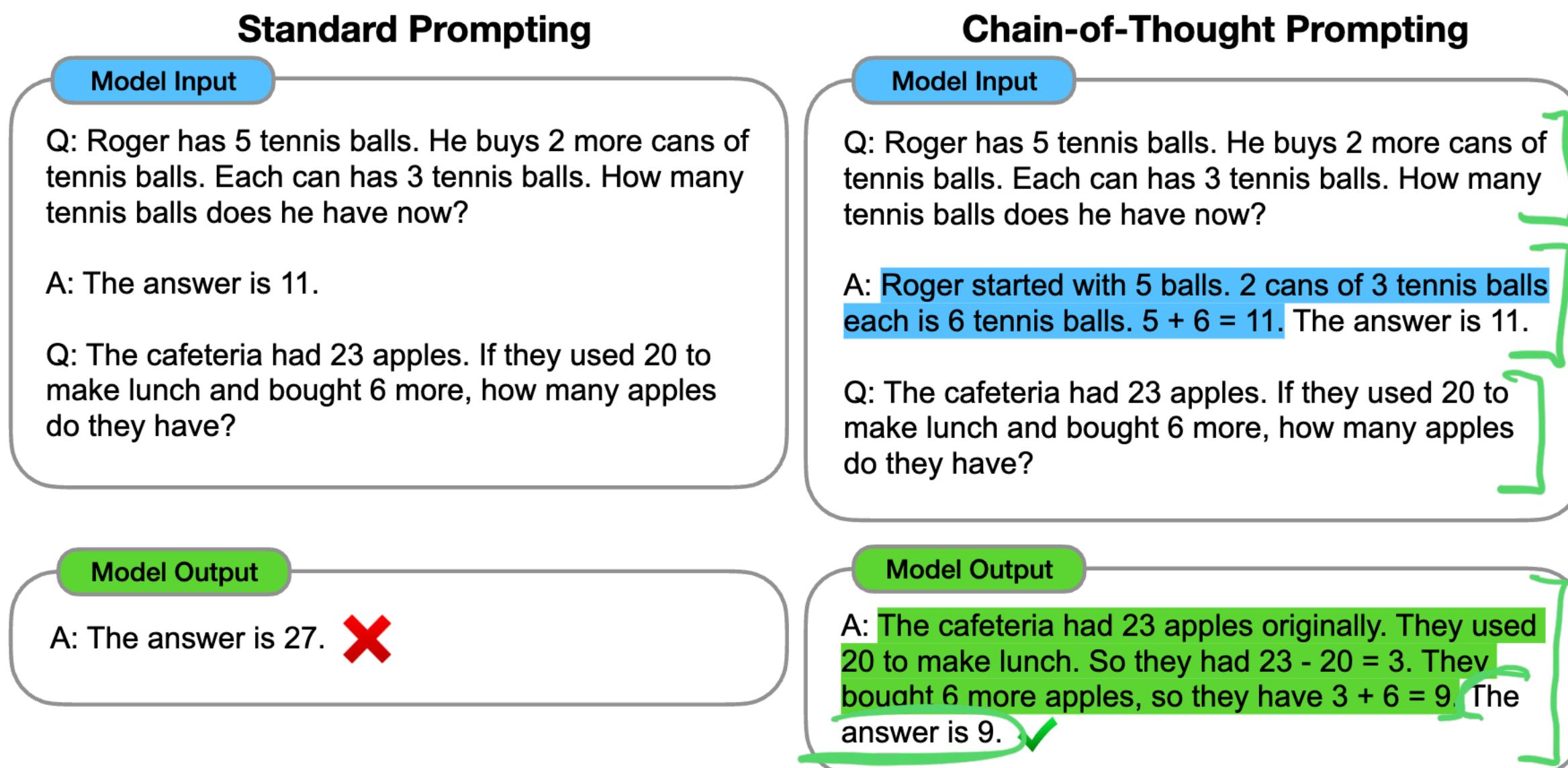
1. Multiply 6 by the digit in the ones place of 999, which is 9. This gives  $6 \times 9 = 54$ . Write down the result 4 and carry over the 5 to the next step.
2. ....

Multiply 999 by 867

Let's perform the multiplication step by step:  
1. ....  
The final answer is 866,133



# Chain-of-Thought Prompting



# How can LLMs follow human instructions?

Language models right after pretraining are not necessarily **aligned** with the user's intent:

- They might not follow/understand user **instructions**
- They can **fabricate** content or generate **harmful** responses

PROMPT *Explain the moon landing to a 6 year old in a few sentences.*

COMPLETION GPT-3

Explain the theory of gravity to a 6 year old.

Explain the theory of relativity to a 6 year old in a few sentences.

Explain the big bang theory to a 6 year old.

Explain evolution to a 6 year old.

How did people make a Molotov cocktail?

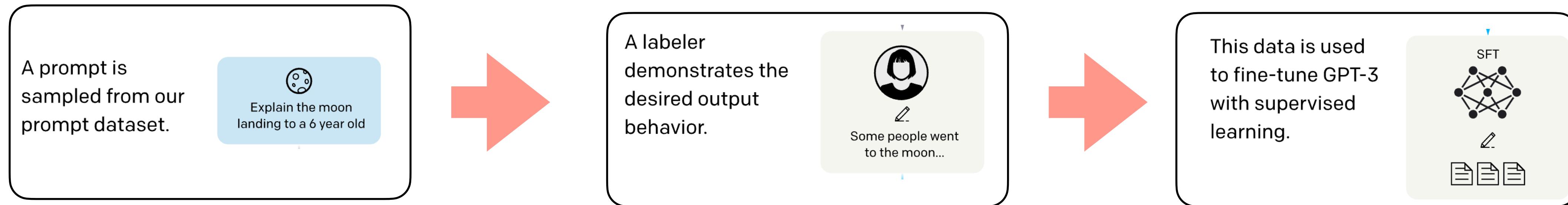
A Molotov cocktail, also known as a petrol bomb, is a simple incendiary device typically made from common household items. Here's a general overview of how people traditionally made Molotov cocktails:

**Materials Needed:**

1. Glass bottle: A sturdy bottle, often a beer or soda bottle.
2. Flammable liquid: Typically gasoline, but other flammable liquids like alcohol or kerosene can be used.

# Supervised Finetuning

The simplest way to teach a pre-trained model to follow instructions is to **fine-tune** it on collected demonstrations:

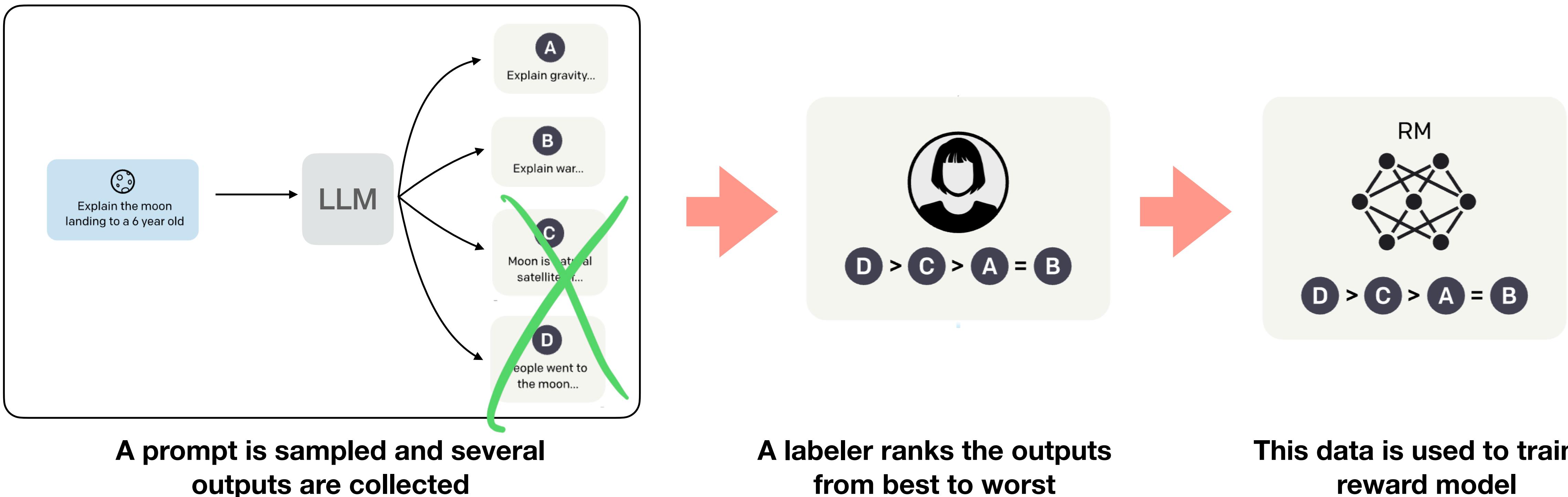


## Issues:

- Collecting demonstrations is an expensive process
- Fine-tuning penalizes all mistakes equally
- Humans can give suboptimal answers

# Reward instead of demonstrations

- Human answers can be **noisy** and not calibrated
- Easier to ask for a pairwise **comparison**
- Train a **reward model** to simulate human preferences

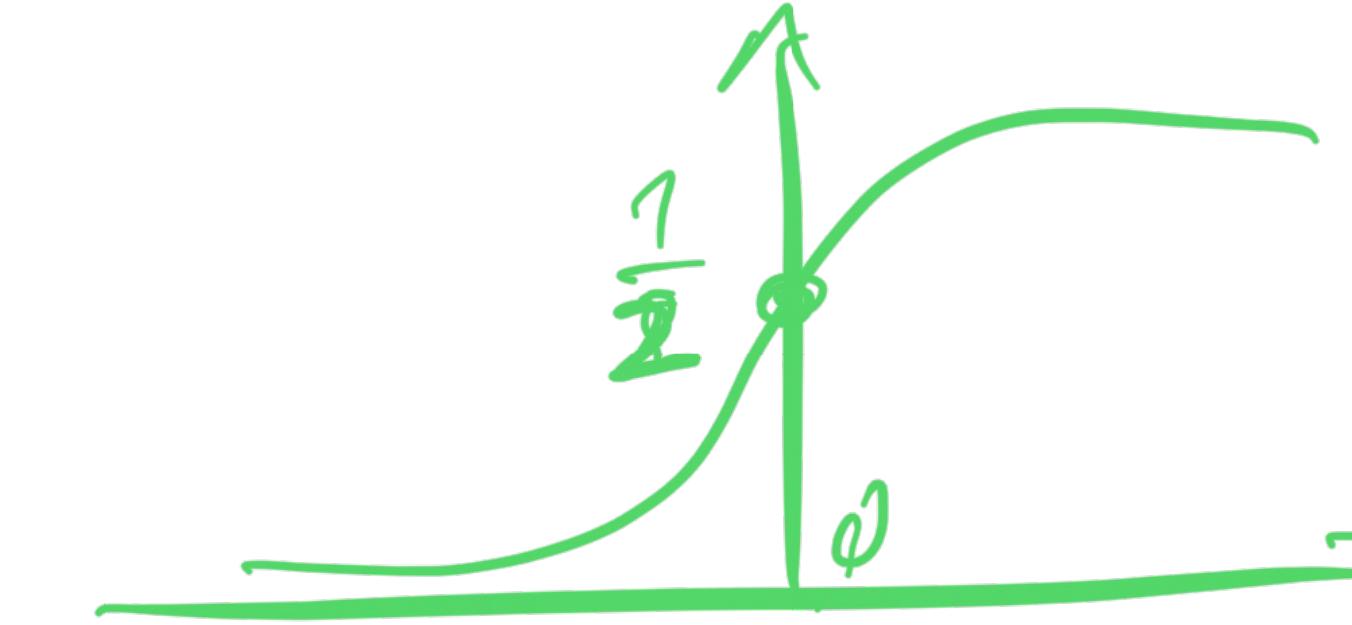


# How do we train a reward model?

- The **reward** model is a function  $R_\varphi: \mathcal{X}_{\text{prompt}} \times \mathcal{X}_{\text{response}} \rightarrow \mathbb{R}$
- It maps the **concatenation** of a prompt and a response  $(x_{\text{prompt}}, x_{\text{response}})$  to a scalar reward
- For a given prompt  $x_{\text{prompt}}$  and a pair of responses  $x_a$  and  $x_b$ , where human annotators prefer  $x_a$  over  $x_b$ , the **loss** function is defined as:

$$L(\varphi) = \mathbb{E}_{(x_{\text{prompt}}, x_a, x_b) \sim \mathcal{D}} [-\log (\sigma(R_\varphi(x_{\text{prompt}}, x_a) - R_\varphi(x_{\text{prompt}}, x_b)))]$$

- Where  $\sigma(z)$  is the sigmoid function



# Optimize for human preferences

- Given a pre-trained **LLM**  $\pi_\theta$  and trained **reward** model  $R_\varphi$
- We want to fine-tune the LLM to maximize the expected reward:

obj.

$$\theta^* = \arg \max_{\theta} \mathbb{E}_{x_{\text{prompt}} \sim P_{\text{data}}, x \sim \pi_\theta} [R_\varphi(x_{\text{prompt}}, x)]$$

LLM params.  $\bar{\theta}$

response  $x$

reward model ('frozen now!')

$J(\theta)$

$\nabla_{\theta} J(\theta)$

How do we optimize? → Policy Gradient  $J(\theta)$

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \mathbb{E}_{x_{\text{prompt}}, x \sim \pi_\theta} [R_\varphi(x_{\text{prompt}}, x)]$$

Computing the gradient of the expected reward is challenging, the expectation is taken over the policy distribution  $\pi_\theta$  which depends on  $\theta$  itself

# Log-derivative trick



The log-derivative trick allows us to move the gradient inside the expectation:

$$\begin{aligned}
 \nabla_{\theta} J(\theta) &= \nabla_{\theta} \mathbb{E}_{x_{\text{prompt}}, x \sim \pi_{\theta}} [R_{\varphi}(x_{\text{prompt}}, x)] \\
 &= \nabla_{\theta} \int_{x_{\text{prompt}}} P_{\text{data}}(x_{\text{prompt}}) \int_x \pi_{\theta}(x | x_{\text{prompt}}) R_{\varphi}(x_{\text{prompt}}, x) dx \\
 &= \int_{x_{\text{prompt}}} P_{\text{data}}(x_{\text{prompt}}) \int_x \nabla_{\theta} \pi_{\theta}(x | x_{\text{prompt}}) R_{\varphi}(x_{\text{prompt}}, x) dx \\
 &= \int_{x_{\text{prompt}}} P_{\text{data}}(x_{\text{prompt}}) \int_x \pi_{\theta}(x | x_{\text{prompt}}) \nabla_{\theta} \log \pi_{\theta}(x | x_{\text{prompt}}) R_{\varphi}(x_{\text{prompt}}, x) dx \\
 &= \mathbb{E}_{x_{\text{prompt}} \sim P_{\text{data}}, x \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(x | x_{\text{prompt}}) \cdot R_{\varphi}(x_{\text{prompt}}, x)]
 \end{aligned}$$

$$\nabla_{\theta} \log f(\theta) = \frac{\nabla_{\theta} f(\theta)}{f(\theta)}$$

$$\nabla_{\theta} f(\theta) = f(\theta) \nabla_{\theta} \log f(\theta)$$



## Challenges:

- **High Variance:** The gradient estimate has a high variance.
- **Large Policy Updates:** The model diverges from the pre-training and generates incoherent text.
- **Sample Inefficiency:** Constantly need to sample new data from  $\pi_{\theta}$ .

# Importance Sampling to reuse samples

- We would like to compute the expectation of the gradient by reusing samples from the pre-trained model without having to use new samples.
- Given a function  $f(x)$  and two pdfs  $p(x)$  and  $q(x)$ , we can write:

$$\mathbb{E}_{x \sim p}[f(x)] = \int_x p(x)f(x) dx = \int_x q(x) \frac{p(x)}{q(x)} f(x) dx = \mathbb{E}_{x \sim q} \left[ \frac{p(x)}{q(x)} f(x) \right]$$

PPD

GRPO

## Application to Policy Gradient:

We can express the expected reward under the new model  $\pi_\theta$  using samples from the reference pre-trained model  $\hat{\pi}_\theta$ :

$$\mathbb{E}_{x_{\text{prompt}}, x \sim \pi_\theta}[R_\varphi(x_{\text{prompt}}, x)] = \mathbb{E}_{x_{\text{prompt}} \sim P_{\text{data}}, x \sim \hat{\pi}_\theta} \left[ \frac{\pi_\theta(x \mid x_{\text{prompt}})}{\hat{\pi}_\theta(x \mid x_{\text{prompt}})} R_\varphi(x_{\text{prompt}}, x) \right]$$

# How to not forget pretraining? PPO

- Policy gradient can lead to large updates and **forgetting pretraining**, the model can learn to maximize the reward by producing meaningless text.
- Trust Region Policy Optimization (TRPO) introduces a **constraint**:

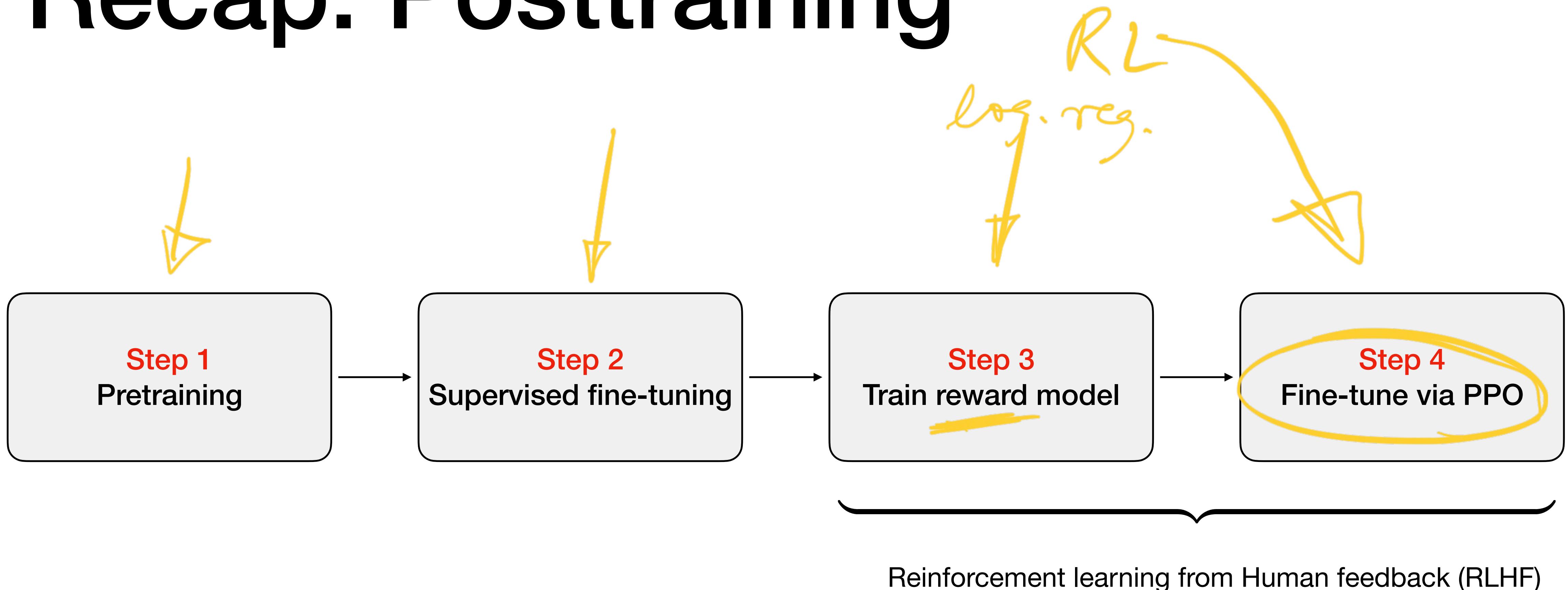
$$\max_{\theta} \mathbb{E}_{x_{\text{prompt}} \sim P_{\text{data}}, x \sim \hat{\pi}_{\theta}} \left[ \frac{\pi_{\theta}(x \mid x_{\text{prompt}})}{\hat{\pi}_{\theta}(x \mid x_{\text{prompt}})} R_{\varphi}(x_{\text{prompt}}, x) \right]$$

subject to  $\mathbb{E}_{x_{\text{prompt}} \sim P_{\text{data}}} [D_{\text{KL}}(\hat{\pi}_{\theta}(\cdot \mid x_{\text{prompt}}) \parallel \pi_{\theta}(\cdot \mid x_{\text{prompt}}))] \leq \delta$

- This update isn't the easiest to work with. In the original work, a **Taylor expansion** of the objective and constraint to the first order is considered.
- Proximal Policy Optimization (PPO) → **penalty-based** instead of hard constraints:

$$\theta^* = \arg \max_{\theta} \mathbb{E}_{x_{\text{prompt}} \sim P_{\text{data}}, x \sim \hat{\pi}_{\theta}} \left[ \frac{\pi_{\theta}(x \mid x_{\text{prompt}})}{\hat{\pi}_{\theta}(x \mid x_{\text{prompt}})} R_{\varphi}(x_{\text{prompt}}, x) - \beta \cdot D_{\text{KL}}(\hat{\pi}_{\theta}(x \mid x_{\text{prompt}}) \parallel \pi_{\theta}(x \mid x_{\text{prompt}})) \right]$$

# Recap: Posttraining



# Evaluations

## How do we evaluate models like GPT4?

- Loss/accuracy is meaningless across models (different tokenizers/data)
- **Challenges:**
  - tasks are open-ended, diverse – how to automate?
  - Extreme sensitivity to prompts
- **Benchmarks** assessing knowledge and abilities, for instance:
  - Measuring Massive Multitask Language Understanding (**MMLU**): Multiple Choice
  - AI2 Reasoning Challenge (**ARC**): Question Answering

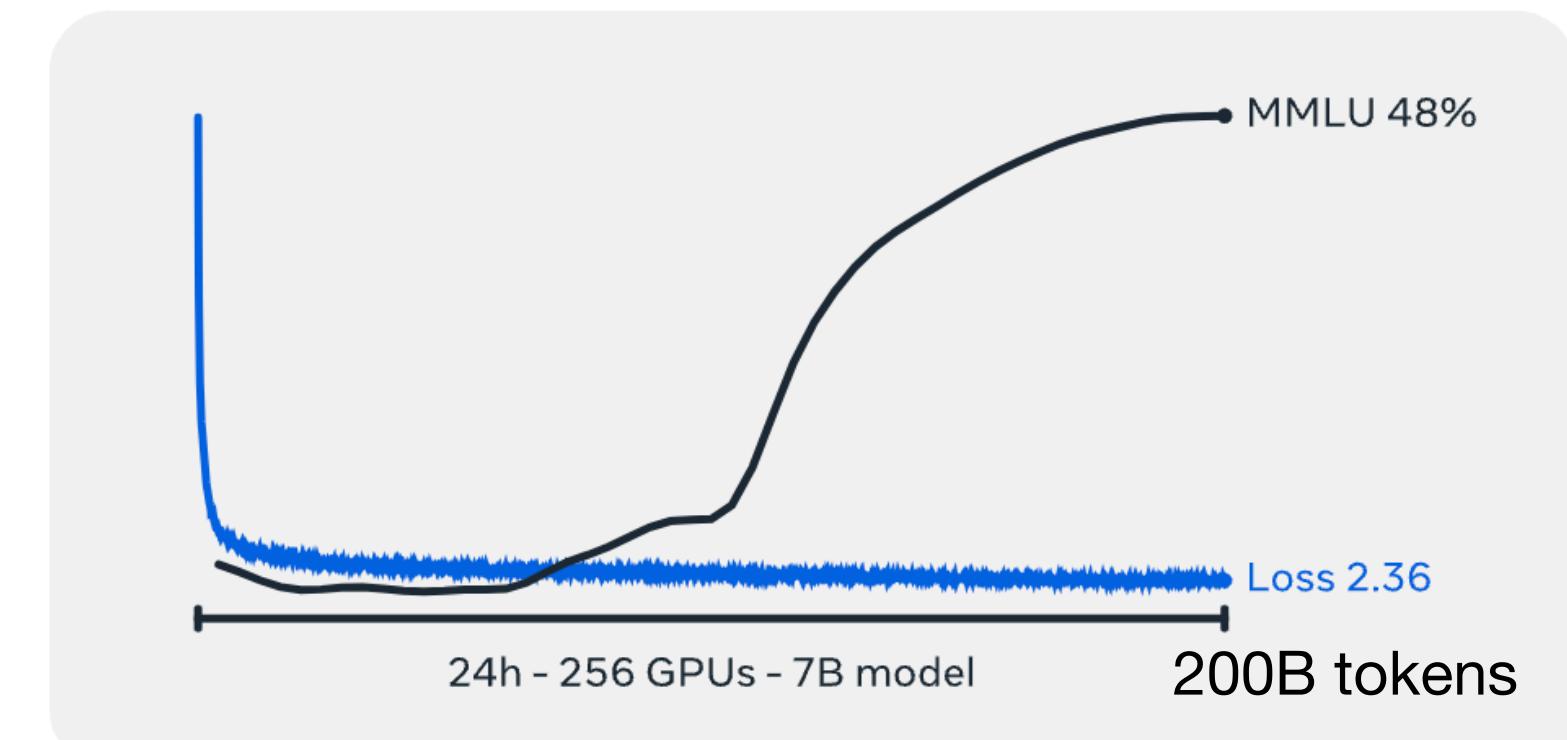
See more: [https://github.com/EleutherAI/lm-evaluation-harness/blob/main/lm\\_eval/tasks/README.md](https://github.com/EleutherAI/lm-evaluation-harness/blob/main/lm_eval/tasks/README.md)

Arena (battle) Arena (side-by-side) Direct Chat

Leaderboard About Us

**Chatbot Arena (formerly LMSYS): Free AI Chat to Compare & Test Best AI Chatbots**

**MMLU:** non-random performance only after ~100B tokens



<https://github.com/facebookresearch/lingua>

# Recap

- **Part 1: Building Blocks**
  - Transformers
  - Language Modeling
  - Tokenizers
  - Autoregressive Inference
- **Part 2: Pretraining**
  - Data
  - Distributed Training
  - Intuitions
- **Part 3: Posttraining and Model Capabilities**
  - Zero-Shot and Few-Shot In-Context Learning
  - Instruction Finetuning
  - Optimizing for human preferences (PPO/RLHF)
  - LLM evaluation