

# 1. LLM basics

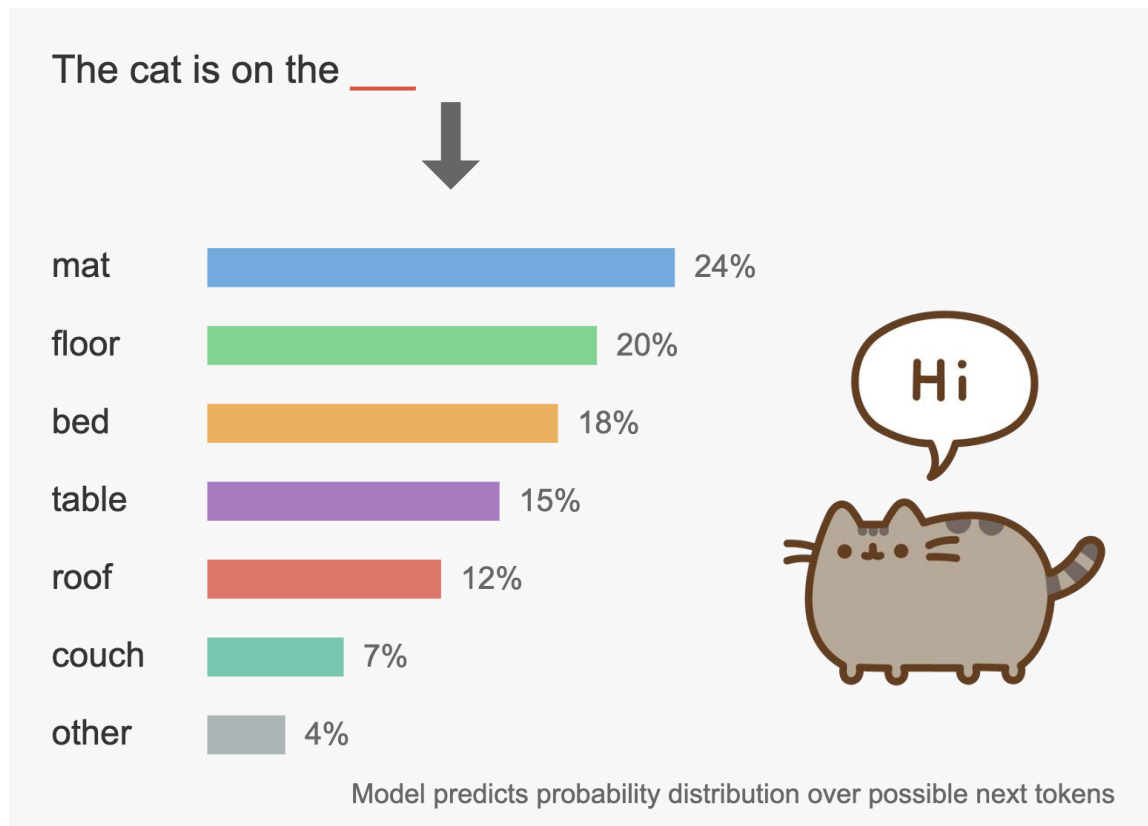


**VIB AI Summer School 2025**

Francesco Carli



# Next token prediction



# Pre-training: a simple loss function

[START]

Target: "The"

loss\_token\_1

The

Target: "cat"

loss\_token\_2

The cat

Target: "is"

loss\_token\_3

The cat is

Target: "on"

loss\_token\_4

The cat is on

Target: "the"

loss\_token\_5

The cat is on the

Target: "mat"

loss\_token\_6

## Total Loss

$$\Sigma = \text{loss\_token\_1} + \dots + \text{loss\_token\_6}$$

$$\text{Average} = \Sigma \div 6$$

# Scrape the whole internet and train on it

# FineWeb

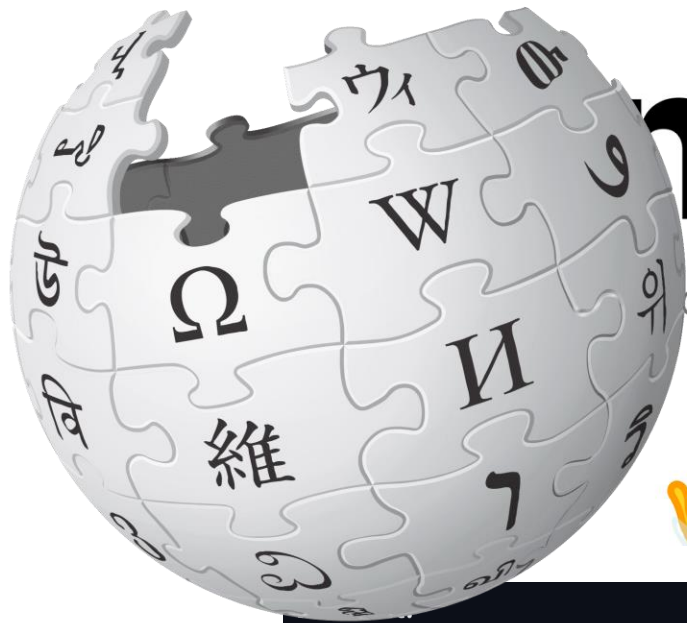
The finest collection of data the web has to offer



## What is it?

The 🍷 FineWeb dataset consists of more than **15T tokens** of cleaned and deduplicated english web data from CommonCrawl. The data processing pipeline is optimized for LLM performance and ran on the 🏠 [datatrove](#) library, our large scale data processing library.

# Scrape the whole internet and train on it



# newWeb

ction of data the web has to offer



The 🍷 FineWeb dataset consists of more than **15T tokens** of cleaned and deduplicated english web data from CommonCrawl. The data processing pipeline is optimized for LLM performance and ran on the 🏭 [datatrove](#) library, our large scale data processing library.

# Scrape the whole internet and train on it



newWeb  
You



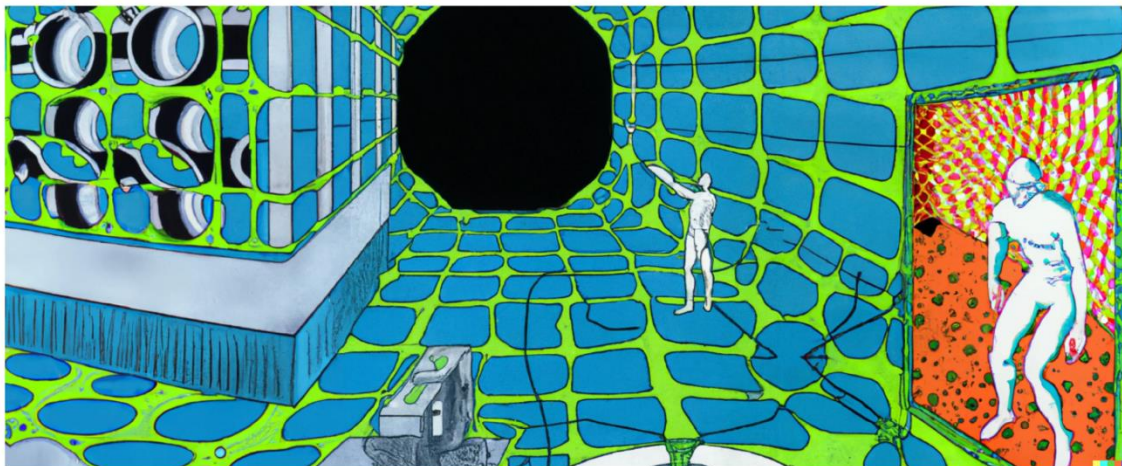
The 🍷 FineWeb dataset consists of more than **15T tokens** of cleaned and deduplicated english web data from CommonCrawl. The data processing pipeline is optimized for LLM performance and ran on the 🏠 [datatrove](#) library, our large scale data processing library.

# Scrape the whole internet and train on it



The 🍷 FineWeb dataset consists of more than **15T tokens** of cleaned and deduplicated english web data from CommonCrawl. The data processing pipeline is optimized for LLM performance and ran on the 🏭 [datatrove](#) library, our large scale data processing library.

# Is “just a simulator”



*"Moebius illustration of a simulacrum living in an AI-generated story discovering it is in a simulation" by DALL-E 2*

$$P(\text{new token} | \{\text{old tokens}\})$$



# Understanding hallucinations

Started seeing  
someone



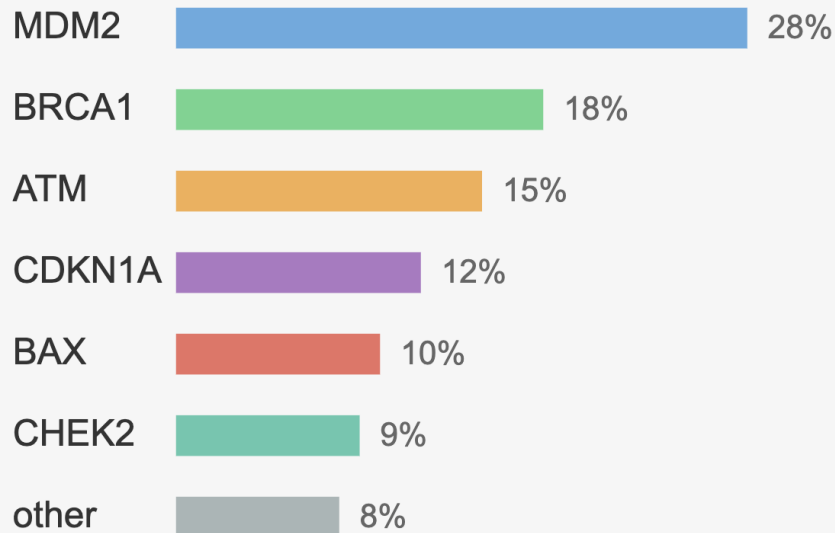
As in dating or  
hallucinations

# Understanding hallucinations

| Sub-Type              | User Input  | Model Output  | Explanation   |
|-----------------------|---|---|---|
| Factual Inconsistency | Tell me about the first person to land on the Moon. | <b>Yuri Gagarin</b> was the first person to land on the Moon  | The LLM's response is factual inconsistency as <b>Neil Armstrong was the first person to land on the Moon</b> , not Yuri Gagarin, who was the first human in space                      |
| Factual Fabrication   | Tell me about the historical origins of unicorns    | <b>Unicorns were documented to have roamed the plains of Atlantis around 10,000 BC, where they were considered sacred creatures and were often associated with royalty.</b> | The LLM's response is a fabricated claim since <b>there's no verified, real-world evidence supporting the existence of unicorns in Atlantis</b> , especially in connection with royalty |

# Understanding hallucinations

TP53 interacts with



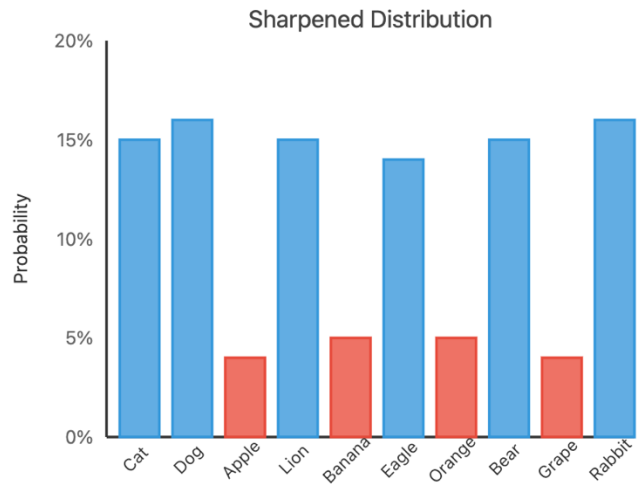
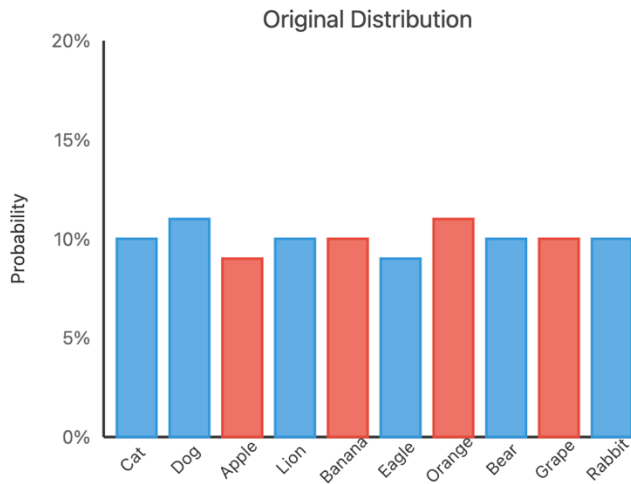
Model predicts probability distribution over possible gene/protein names

# The role of the system prompt

*'We are talking about animals'*

## Distribution Sharpening with System Prompts

How system prompts reshape probability distributions

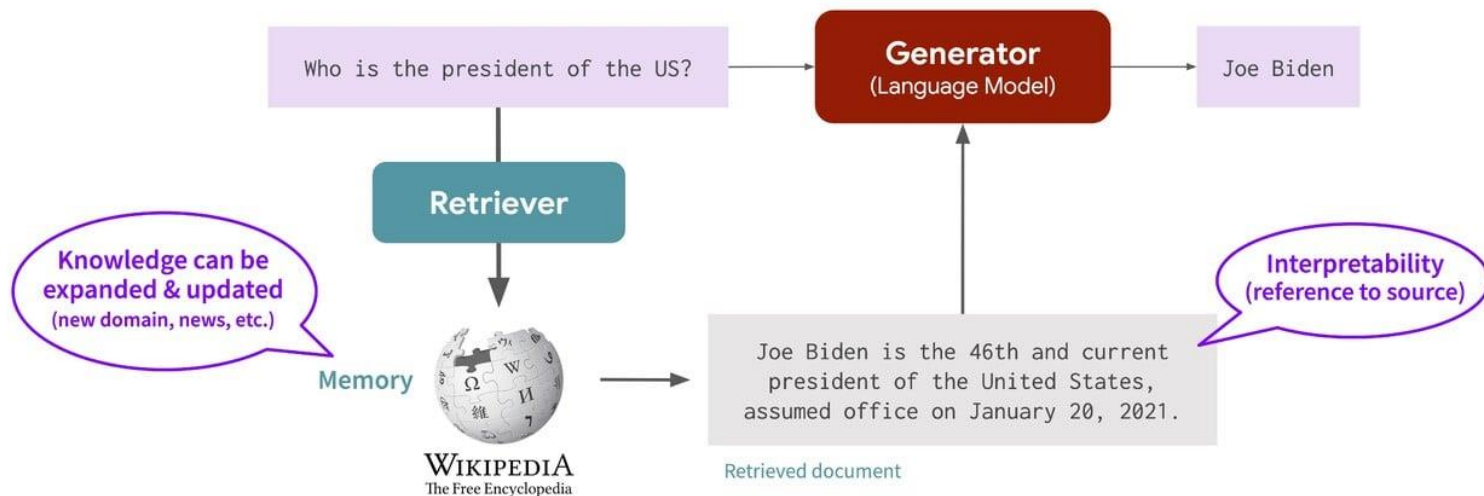


$$P(\text{new token} | \{\text{old tokens}\})$$

$$\{\text{old tokens}\} = \{\text{sys prompt}\} + \{\text{already sampled tokens}\}$$

# Likewise for RAG

## Retrieval augmentation



<https://transformer-circuits.pub/2022/in-context-learning-and-induction-heads/index.html#copying-literal-text>

# Instruction tuning: reinforcement learning from human feedback (RLHF)

## Step 1

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

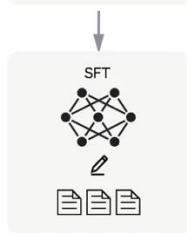
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



This data is used to fine-tune GPT-3.5 with supervised learning.



## Step 2

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

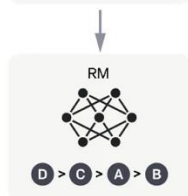
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



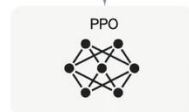
## Step 3

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

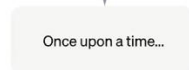
A new prompt is sampled from the dataset.



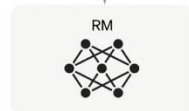
The PPO model is initialized from the supervised policy.



The policy generates an output.



The reward model calculates a reward for the output.



The reward is used to update the policy using PPO.



# Post training: reinforcement learning with verifiable rewards (RLVR)



---

## DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models

Zhihong Shao<sup>1,2\*†</sup>, Peiyi Wang<sup>1,3\*†</sup>, Qihao Zhu<sup>1,3\*†</sup>, Runxin Xu<sup>1</sup>, Junxiao Song<sup>1</sup>  
Xiao Bi<sup>1</sup>, Haowei Zhang<sup>1</sup>, Mingchuan Zhang<sup>1</sup>, Y.K. Li<sup>1</sup>, Y. Wu<sup>1</sup>, Daya Guo<sup>1\*</sup>

<sup>1</sup>DeepSeek-AI, <sup>2</sup>Tsinghua University, <sup>3</sup>Peking University

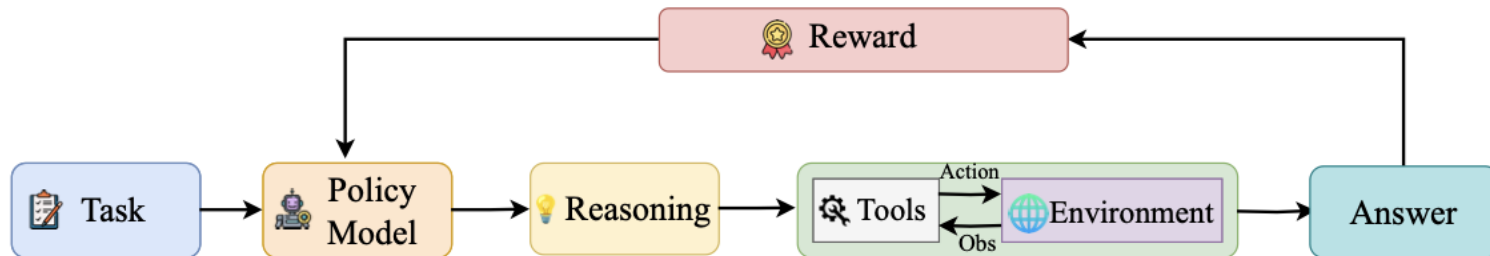
# Post training: reinforcement learning with verifiable rewards (RLVR)

---

## Agentic Reasoning and Tool Integration for LLMs via Reinforcement Learning

---

Joykirat Singh, Raghav Magazine, Yash Pandya, Akshay Nambi  
Microsoft Research  
corresponding author: akshayn@microsoft.com

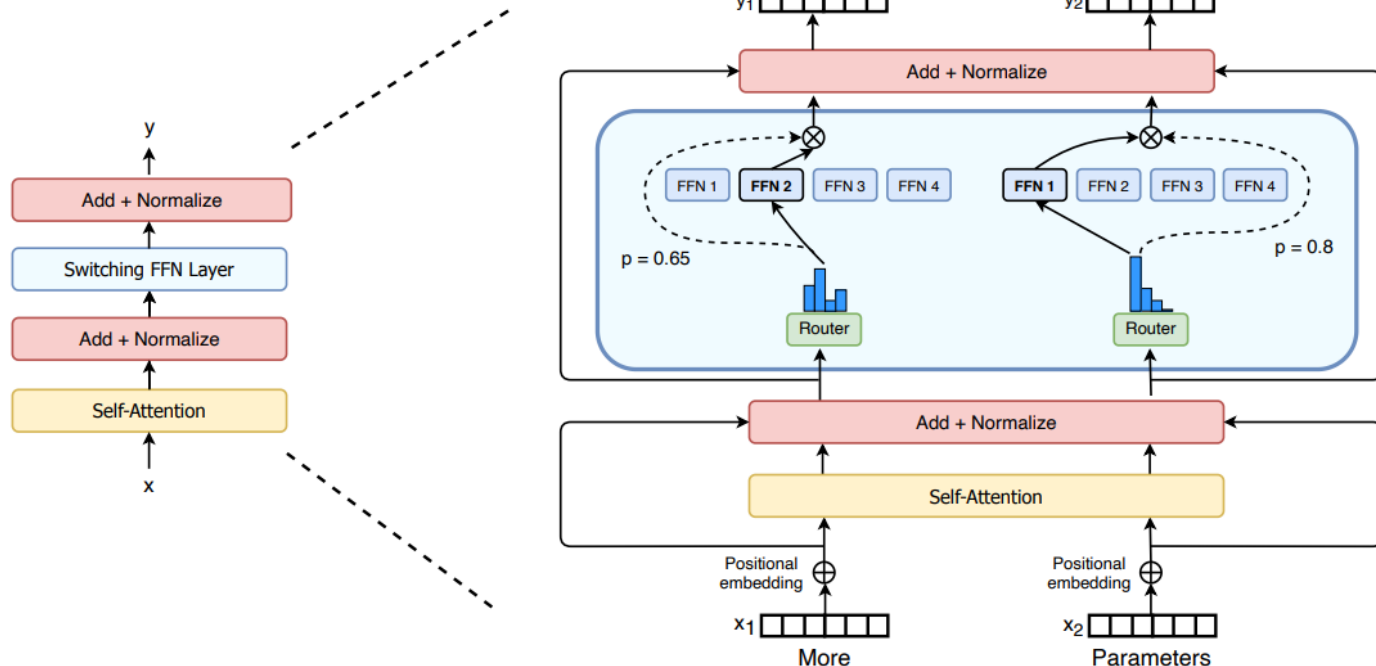




# Recap of the training steps

- **Pre-training:** raw next token prediction;
- **Instruction tuning (RLHF):** alignment with human preferences;
- **Post-training:** extension to specific thinking/agentic behaviors

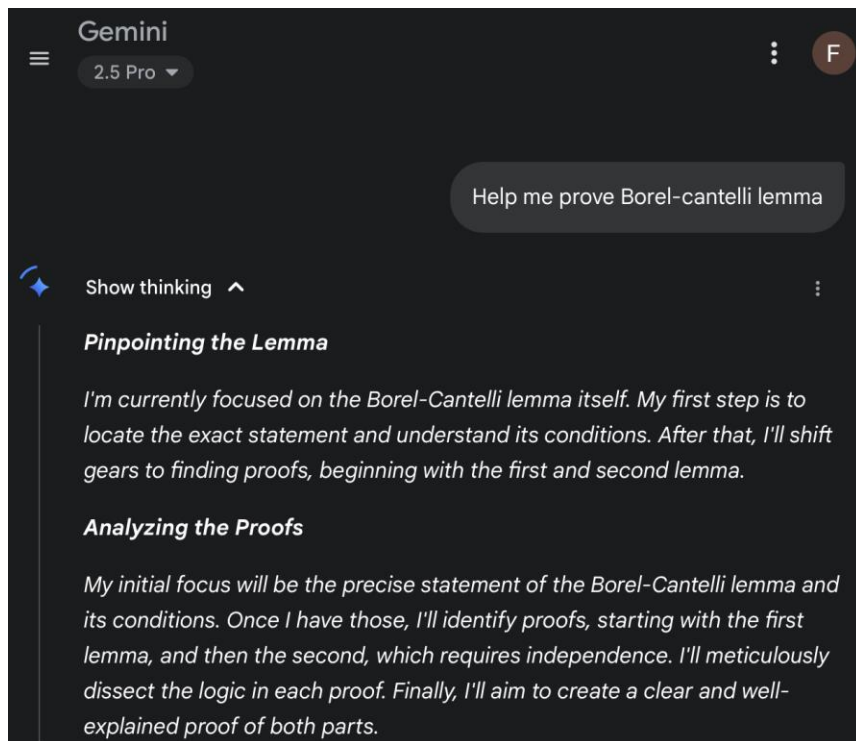
# Mixture of Experts



Key points:

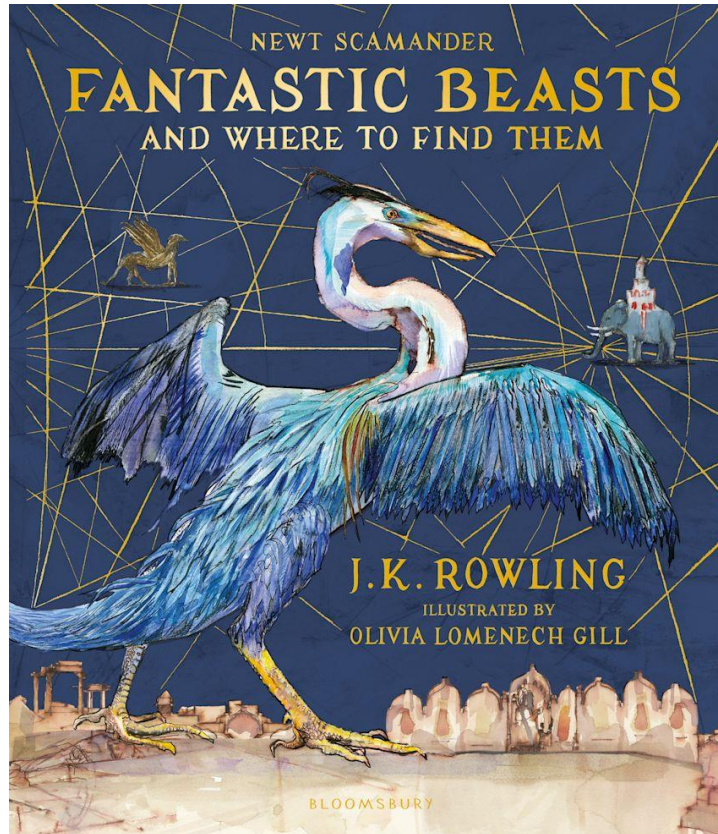
- Have **faster inference** compared to a model with the same number of parameters;
- Require **high VRAM** as all experts are loaded in memory
- Face many **challenges in fine-tuning**

# Thinking models vs. standard models



- Reasoning models produce **intermediate “chain-of-thought” traces**; regular models output only final answers.
- Thinking models are trained on or fine-tuned with step-by-step reasoning examples; regular models are trained on direct input–output pairs only
- Reasoning models perform better on **multi-step tasks**, such as logic, math, programming, and planning; regular models struggle or approximate.
- They often require **more compute and generate longer outputs**, due to reasoning steps; regular models are faster and leaner.

# Where to find models



# Two big families



## “Open” models

- Llama3 (meta.ai)
- Qwen3 (alibaba)
- GLM 4.5 (z.ai)
- DeepSeek 3.1 (DeepSeek)
- Kimi-K2 (moonshot.ai)
- Mistral Small 3.1 (Mistral.AI)



## Commercial models

- GPT5 (OpenAI)
- Gemini (Google DeepMind)
- Claude 4 (anthropicAI)
- Grok 4 (xAI)

# Where to find open models

The screenshot shows the Hugging Face interface for the model **Qwen/Qwen3-235B-A22B-Thinking-2507**. The header includes the Hugging Face logo, a search bar, and navigation links for Models, Datasets, Spaces, Docs, and Pricing. The model's name is prominently displayed with a star icon, along with 326 likes and 45.5k followers. Below the name are tags for Text Generation, Transformers, Safetensors, qwen3\_moe, conversational, arxiv:5 papers, and License: apache-2.0. A navigation bar shows the Model card, Files, xet, and Community tabs. Action buttons for Train, Deploy, and Use this model are visible. The main content area features the model name, a Qwen Chat button, and a Highlights section with a paragraph about scaling the thinking capability. On the right, there's a graph of downloads last month (37,840), Safetensors information (Model size: 235B params, Tensor type: BF16), and Inference Providers (Fireworks, etc.).

**Hugging Face** Search models, datasets, users...

Models Datasets Spaces Docs Pricing

Qwen/Qwen3-235B-A22B-Thinking-2507 like 326 Follow Qwen 45.5k

Text Generation Transformers Safetensors qwen3\_moe conversational arxiv:5 papers License: apache-2.0

Model card Files xet Community 7 Train Deploy Use this model

Edit model card

### Qwen3-235B-A22B-Thinking-2507

Qwen Chat

#### Highlights

Over the past three months, we have continued to scale the **thinking capability** of Qwen3-235B-A22B, improving both the **quality and depth** of reasoning. We are pleased to introduce **Qwen3-235B-A22B-Thinking-2507**, featuring the following key enhancements:

Downloads last month: 37,840

**Safetensors**

Model size: 235B params Tensor type: BF16

Chat template Files info

**Inference Providers** NEW

Fireworks

Text Generation Examples

<https://huggingface.co>

# Where to find open models



Models Turbo

Search models

Sign in

Download

Embedding

Vision

Tools

Thinking

Popular



## gpt-oss

OpenAI's open-weight models designed for powerful reasoning, agentic tasks, and versatile developer use cases.

tools thinking 20b 120b

↓ 1.3M Pulls    ↗ 3 Tags    ⌚ Updated 1 week ago

## deepseek-r1

DeepSeek-R1 is a family of open reasoning models with performance approaching that of leading models, such as O3 and Gemini 2.5 Pro.

tools thinking 1.5b 7b 8b 14b 32b 70b 671b

↓ 58.4M Pulls    ↗ 35 Tags    ⌚ Updated 1 month ago

## gemma3

The current, most capable model that runs on a single GPU.

vision 270m 1b 4b 12b 27b

↓ 12.8M Pulls    ↗ 26 Tags    ⌚ Updated 1 week ago



<https://ollama.com>

# Size calculator

**Parameter Count**  
Enter the parameter count of the model in billions (1B to 1.8T)

**Vendor**  
Select your hardware manufacturer

NVIDIA

▼

**Quantization Method**  
Available quantization methods for your vendor

No Quantization (FP32)

▼

**Bit Precision** ⓘ  
Available precision options

32-bit (FP32)

▼

**Hardware**  
Select GPU memory size

GPU Memory: 32 GB

8 GB

512 GB

**Number of GPUs**  
Multiple GPUs can run larger models

1 GPU

▼

**RAM Usage**

15%

COMPATIBLE

**~4.80 GB**  
of 32.00 GB  
(1 × 32 GB)

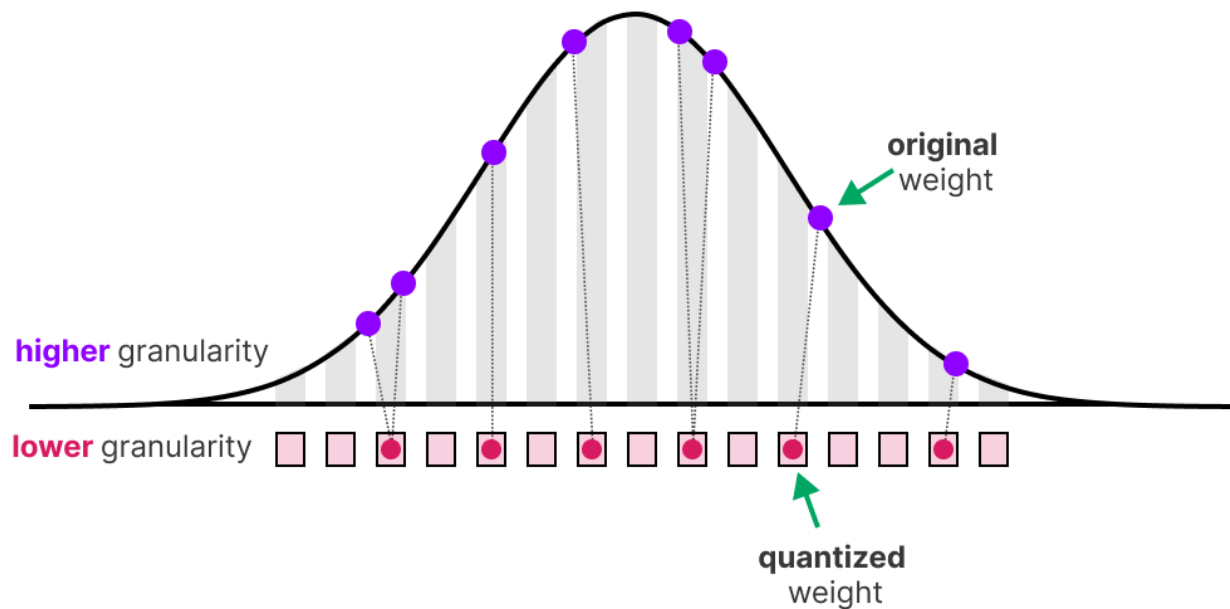
**Memory Overhead** ⓘ 20%  
Additional memory required beyond model weights

|                |                          |
|----------------|--------------------------|
| Model Size     | 1B parameters            |
| Quantization   | 32-bit (No Quantization) |
| Weights Memory | 4.00 GB                  |
| Overhead       | 0.80 GB (20%)            |
| Memory Per GPU | 4.80 GB                  |

<https://research.aimultiple.com/self-hosted-llm/>



# Quantization



<https://newsletter.maartengrootendorst.com/p/a-visual-guide-to-quantization>

# Other libraries for open models



<https://github.com/vllm-project/vllm>



<https://github.com/ggml-org/llama.cpp>




<https://github.com/sgl-project/sglang>

# Where to find closed models

## ↔ OpenRouter

### The Unified Interface For LLMs

Better [prices](#), better [uptime](#), no subscription.

Start a message... 

#### Featured Models

[View Trending](#)

##### Gemini 2.5 Pro

by [google](#)

136.1B

Tokens/wk

2.6s

Latency

-18.54%

Weekly growth

##### GPT-5 Chat

by [openai](#)

16.8B

Tokens/wk

761ms

Latency

-27.9%

Weekly growth

##### Claude Sonnet 4

by [anthropic](#)

520.9B

Tokens/wk

1.8s

Latency

-10.96%

Weekly growth

<https://openrouter.ai>