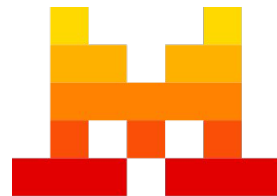# Choosing the Right Model

●●●

Krzysztof Nowak (Institute of Polish Language, Kraków)

# Outline

- model types
- architecture
- cost
- computational resources
- performance

———

LLM Ecosystem

# Focus, Licensing, Cost

| Model | Provider | Type | Cost Range ($/1M) | Access |
|-------|----------|------|-------------------|--------|
| Claude 3.7 Sonnet | Anthropic | Proprietary | $1.5-7.5 | API/Web |
| GPT-4.1 | OpenAI | Proprietary | $2.0-8.0 | API/ChatGPT |
| Gemini Ultra | Google | Proprietary | $1.0-2.0 | AI Studio |
| Qwen 2.5 72B | Alibaba | Open/Hosted | $0.5-1.5 | API/Local |
| Llama 3.1 405B | Meta | Open Source | Free | Local only |
| GPT-3.5 Turbo | OpenAI | Proprietary | $0.5-1.0 | API/Web |
| Mistral Large 2 | Mistral AI | Proprietary | $1.0-3.0 | API/Web |
| DeepSeek V3 | DeepSeek | Open/Hosted | $0.3-1.0 | API/Local |

## Corpus Linguistics

| Research Task | Primary Model | Alternative | Budget Option | Why |
|---|---|---|---|---|
| Large-scale classification | DeepSeek V3 | GPT-3.5 | GPT-3.5 | Cost efficiency, batch processing |
| Sentiment analysis | Claude 3.7 | Mistral Large 2 | GPT-3.5 | Nuanced reasoning needed |
| POS tagging | Qwen 2.5 | DeepSeek V3 | GPT-3.5 | Multilingual support |
| Named entity recognition | GPT-4.1 | Claude 3.7 | Qwen 2.5 | Accuracy across domains |

## Theoretical Linguistics

| Research Task | Primary Model | Alternative | Budget Option | Why |
|---|---|---|---|---|
| Syntax analysis | Claude 3.7 | GPT-4.1 | DeepSeek V3 | Complex reasoning required |
| Semantic annotation | GPT-4.1 | Claude 3.7 | Mistral Large 2 | Precision in meaning analysis |
| Pragmatic analysis | Claude 3.7 | GPT-4.1 | Qwen 2.5 | Context understanding |
| Formal semantics | DeepSeek V3 | Claude 3.7 | GPT-3.5 | Mathematical notation handling |

**SUPPORT POOL OF EXPERTS PROGRAMME**

AI Privacy Risks & Mitigations

**Large Language Models (LLMs)**

By Isabel BARBERÁ

Data sensitivity

- confidential interviews
- sensitive information
- proprietary corpora

GDPR compliance

## Data Sensitivity & Privacy

| Requirement | Recommended Models | Avoid |
|---|---|---|
| Confidential data | Llama 3.1, DeepSeek V3 (local) | All proprietary cloud APIs |
| GDPR compliance | Mistral Large 2, Llama 3.1 | US-based models |
| Institutional policy | Llama 3.1, Qwen 2.5 (local) | Varies by institution |

# Do I really need an LLM?



cost

speed

interpretability

specific tasks

# Practical Considerations
## Access and Running Models Locally

# Web Interface

Classify the following sentence as positive, negative, or neutral:
"Quo usque tandem abutere Catilina patientia nostra?"

The sentence "Quo usque tandem abutere Catilina patientia nostra?" translates roughly from Latin to: "How long, Catiline, will you abuse our patience?"

Classification:
Negative

Reason: The sentence is accusatory and expresses frustration and exasperation toward Catiline, indicating a negative attitude.

# API

```python
from openai import OpenAI
from google.colab import userdata
api_key = userdata.get('OPENAI_API_KEY') # read

client = OpenAI(api_key=api_key)

# contact the API by specifying the model you wa
# the response will be stored in a variable call
response = client.responses.create(
    model="gpt-4.1-nano",
    input="Write one sentence about the use of
)
```

# Running Models Locally

Servers



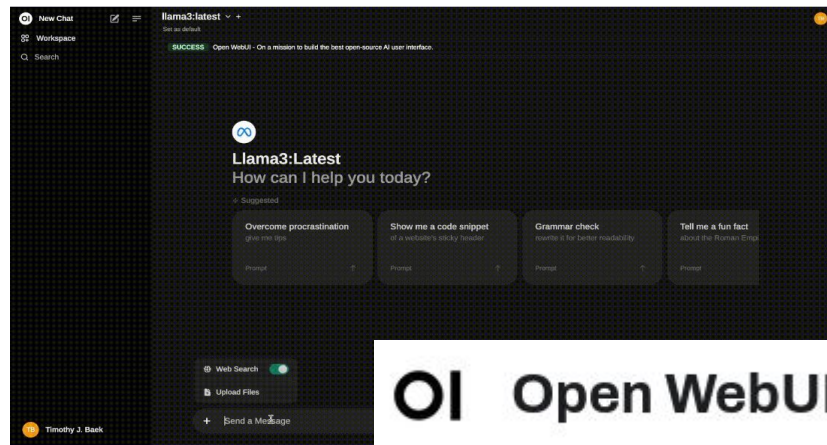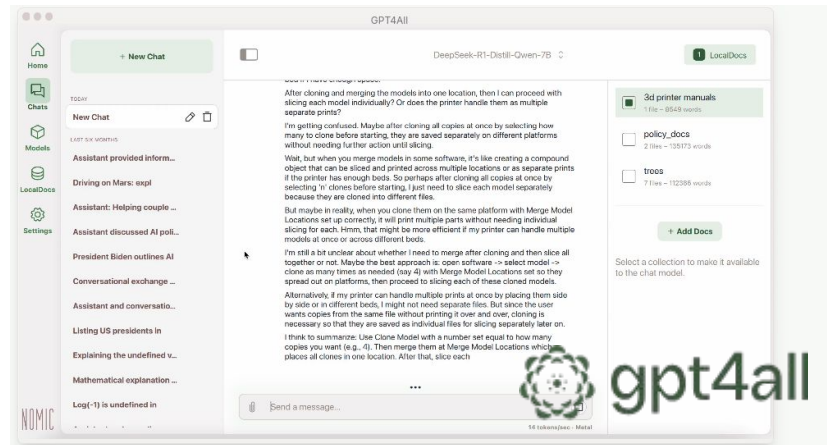LLaMA C++



ollama



LM Studio

# Running Models Locally

Clients

# Running Models Locally

Editors
Libraries





More: **https://github.com/ollama/ollama**

# Practical Considerations
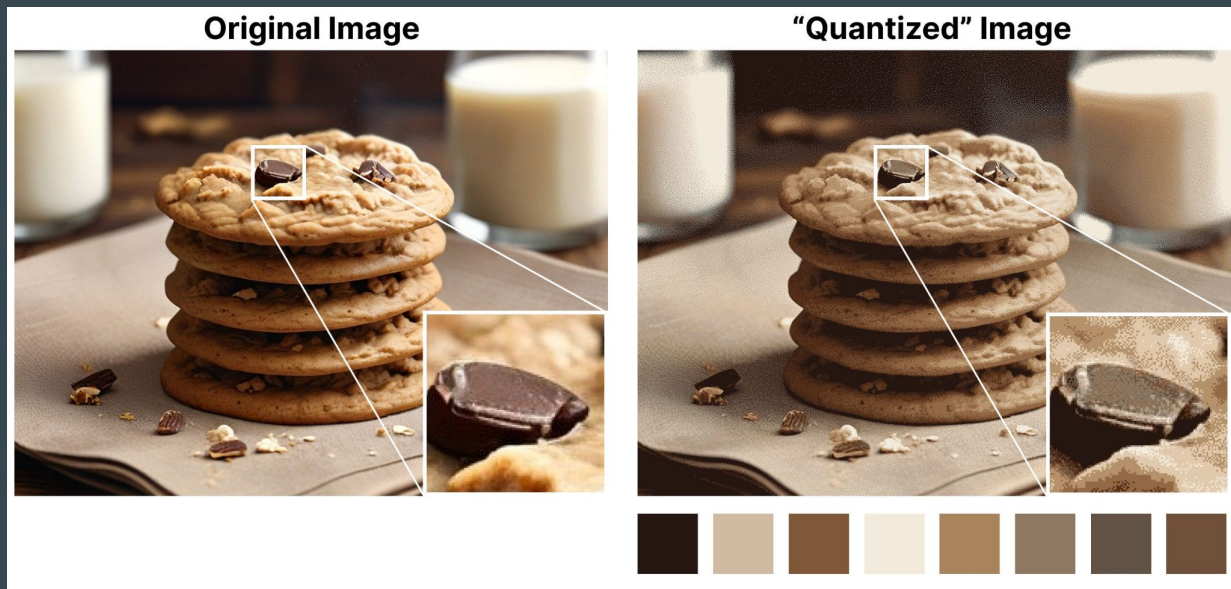## Computational Resources

# Hardware Requirements



7B model → 14GB VRAM

70B model → 140GB VRAM

Consumer GPUs (RTX 4090)

- quantized models
- quality compromises

# Quantization The Quality vs. Accessibility Trade-off

# Practical Considerations
## Cost

# Cost

Counting tokens

Cost Management

Context Window Limits

Latency

Performance Optimization

# Why Token Counting Matters

- **Subword units that LLMs process (not always = words)**

```
"unbelievable" → [un|believ|able]
    1 word          3 tokens
```

- **Punctuation, spaces, and special characters**

```
Sentence: "Hello, world!"
Tokens: [Hello][,][ world][!]
          (4 tokens)
```



**https://platform.openai.com/tokenizer**

Input: ~32 tokens

(18 + 14 tokens)

Output: ~95-105 tokens

For each of the sentences decide if its is positive, negative, or neutral: Quō ūsque tandem abūtere, Catilīna, patientiā nostrā?"

The sentence **"Quō ūsque tandem abūtere, Catilīna, patientiā nostrā?"** translates to **"How long, Catiline, will you abuse our patience?"**

This is a famous line from Cicero's First Catilinarian Oration, and it expresses frustration and accusation toward Catiline.

**Sentiment: Negative**

**Reason:** The sentence accuses Catiline of abusing the patience of the speaker and others, expressing exasperation and criticism.

| | **1 token** | **Request** | **1000 Requests** |
|---|---|---|---|
| Input | $0.000002 | $0.000064 | $0.064 |
| Output | $0.000008 | $0.0008 | $0.8 |
| **Total** | | **$0.000864** | **$0.864** |

# Model Selection

**Flagship chat models**　Our versatile, high-intelligence flagship models.

**GPT-4.1**
Flagship GPT model for complex tasks

**Cost-optimized models**　Smaller, faster models that cost less to run.

**GPT-4.1 mini**
Balanced for intelligence, speed, and cost

**GPT-4.1 nano**
Fastest, most cost-effective GPT-4.1 model

**Reasoning models**　o-series models that excel at complex, multi-step tasks.

**o3**
Our most powerful reasoning model

**o4-mini**
Faster, more affordable reasoning model

# Prompt Engineering

For each of the sentences decide if its is positive, negative, or neutral: Quō ūsque tandem abūtere, Catilīna, patientiā nostrā?"

~~The sentence "Quō ūsque tandem abūtere, Catilīna, patientiā nostrā?" translates to "How long, Catiline, will you abuse our patience?"~~

~~This is a famous line from Cicero's First Catilinarian Oration, and it expresses frustration and accusation toward Catiline.~~

**Sentiment: Negative**

~~Reason: The sentence accuses Catiline of abusing the patience of the speaker and others, expressing exasperation and criticism.~~

# Strategies

Classify sentiment (pos/neg/neu):

System (cached): "You are a sentiment classifier. Respond only with: positive, negative, or neutral."
User input: Quō ūsque tandem abūtere, Catilīna, patientiā nostrā?

Classify each (pos/neg/neu):
1. Quō ūsque tandem abūtere, Catilīna, patientiā nostrā?
2. Quam diū etiam furor iste tuus nōs ēlūdet?
3. Quem ad fīnem sēsē effrēnāta iactābit audācia?

System: Respond with only one word: positive, negative, or neutral
Output: negative

# Resources

- Model documentation
- Papers
- Blogs, Newsletters, Reddit
- Substack

## Interconnects

The cutting edge of AI, from inside the frontier AI labs, minus the hype. The border between high-level and technical thinking. Read by leading engineers, researchers, and investors on Wednesday mornings.

r/LocalLLaMA

AHEAD OF AI

By Sebastian Raschka, PhD

# Benchmarks, Leaderboards

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

# Chatbot Arena ([https://lmarena.ai](https://lmarena.ai))

⚔ Arena (battle)   ⚔ Arena (side-by-side)   💬 Direct Chat   🏆 Leaderboard
🎯 Prompt-to-Leaderboard   🔍 Arena Explorer   ℹ️ About Us

| 💬 Language | 🎲 Overview | 🔺 Price Analysis | 🏯 WebDev Arena | 👀 Vision | 🎨 Text-to-Image | 🎮 Copilot Arena | 🌐 Search | Arena-Hard-Auto |

Total #models: **235**.   Total #votes: **2,920,528**.   Last updated: 2025-05-11.

Code to recreate leaderboard tables and plots in this notebook. You can contribute your vote at lmarena.ai!

| Category | Apply filter | Overall Questions |
|---|---|---|
| Overall ▼ | ☐ Style Control ☐ Show Deprecated | #models: 235 (100%)   #votes: 2,920,528 (100%) |

| Rank* (UB) ▲ | Rank (StyleCtrl) ▲ | Model ▲ | Arena Score ▲ | 95% CI ▲ | Votes ▲ | Organization ▲ | License ▲ |
|---|---|---|---|---|---|---|---|
| 1 | 1 | Gemini-2.5-Pro-Preview-05-06 | 1446 | +8/-9 | 4500 | Google | Proprietary |
| 2 | 1 | o3-2025-04-16 | 1413 | +8/-7 | 6689 | OpenAI | Proprietary |
| 2 | 3 | ChatGPT-4o-latest (2025-03-26) | 1408 | +6/-6 | 10290 | OpenAI | Proprietary |
| 3 | 5 | Grok-3-Preview-02-24 | 1403 | +4/-4 | 14843 | xAI | Proprietary |

# Chatbot Arena (https://lmarena.ai)

## 📜 How It Works

○ **Blind Test**: Ask any question to two anonymous AI chatbots (ChatGPT, Gemini, Claude, Llama, and more).

○ **Vote for the Best**: Choose the best response. You can keep chatting until you find a winner.

○ **Play Fair**: If AI identity reveals, your vote won't count.

○ **NEW features**: **Upload an image** 🖼️ and chat. Use 🌐 **Search** for online LLMs. Use 🎨 **Text-to-Image** models like DALL-E 3, Flux, Ideogram to generate images! Use 🦀 **RepoChat** tab to chat with Github repos.

## 🏆 Chatbot Arena LLM Leaderboard

○ Backed by over **1,000,000+** community votes, our platform ranks the best LLM and AI chatbots. Explore the top AI models on our LLM leaderboard!

## 👇 Chat now!

🔍 Expand to see the descriptions of 101 models                                                                                            ◀

| 💬 Model A | 💬 Model B |
| --- | --- |

---

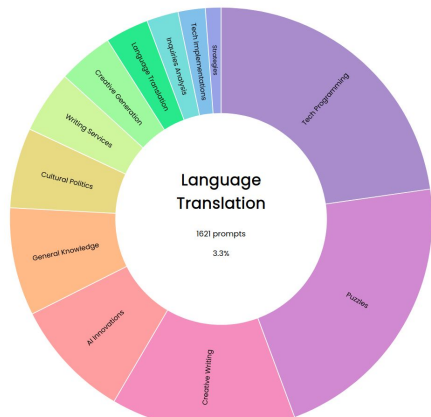| Chat | 🌐 Search | 🦀 RepoChat | 🎨 Text2Image |
| --- | --- | --- | --- |

For each of the sentences decide if its is positive, negative, or neutral: Quō ūsque tandem abūtere, Catilīna, patientiā nostrā? Quam diū etiam furor iste tuus nōs ēlūdet? Quem ad fīnem sēsē effrēnāta iactābit audācia?

**Send**

# Chatbot Arena ([https://lmarena.ai](https://lmarena.ai))

# Advice? Experiment!