# Lumigator:
# evaluating LLMs made simple

Davide Eynard

FOSDEM, 2 Feb 2025

"I look at HuggingFace leaderboard and try a bunch of those… I feel overwhelmed and don't know which model to trust"

Engineer from a leading technology company

"I tried a couple of models… eventually I just went with GPT4 because it seemed like the safest bet"

Data scientist from customer engagement company

# You...

- code
- work with AI and LLMs
- ... but are not necessarily publishing on NeurIPS!
- need to choose LLMs and spend as little time as possible

# We are...

**Contributors** 18



+ 4 contributors

Mozilla.ai's goal:

- Make AI easier and more transparent
- Understand what is broken and fix it
- Integrate existing tools
- Only build what is missing

# Building Lumigator

Lumigator is a **platform** that **guides** users through the process of selecting the right language model for their **specific** needs.

# Key features

- Infrastructure agnostic (local/cluster/cloud)
- Relying on existing tools / standards for interoperability
- API + SDK + UI
- Extensible: by the community for the community
- First use case: evaluation of summarization models

# Lumigator

- Datasets
- Experiments

## Experiments

Experiments are a logical sequence of inference and inference tasks that run sequentially to evaluate an LLM.

**+ Create Experiment**

| EXPERIMENT TITLE | CREATED | STATUS | OPTIONS |
|---|---|---|---|
| ⌄ OpenAI vs more OSS models | 30 Jan 2025, 08:50 | RUNNING | ⋯ |
| oai://gpt-4o-mini | 30 Jan 2025, 08:50 | SUCCEEDED | |
| hf://Falconsai/text_summarizat... | 30 Jan 2025, 08:50 | RUNNING | |
| oai://gpt-4o | 30 Jan 2025, 08:50 | SUCCEEDED | |
| hf://facebook/bart-large-cnn | 30 Jan 2025, 08:50 | SUCCEEDED | |
| mistral://open-mistral-7b | 30 Jan 2025, 08:50 | PENDING | |
| ⌄ OpenAI vs OSS models | 30 Jan 2025, 08:47 | INCOMPLETE | ⋯ |
| ⌄ test_eval_mistral | 29 Jan 2025, 16:36 | SUCCEEDED | ⋯ |

GitHub ↗

Documentation ↗

## Experiment Details

✕

TITLE
OpenAI vs more OSS models

DESCRIPTION
Another experiment comparing more OAI models with more OSS models

STATUS                    SUCCEEDED    ⧉ Logs

JOB ID
3c216ce3-f5e8-4cba-8199-6ed5bbec92f3    ⧉

DATASET
dialogsum.csv

USE-CASE
summarization

MODEL
oai://gpt-4o-mini

CREATED
30 Jan 2025, 08:50

RUN TIME
00: 00: 28

SAMPLES LIMIT
5

TOP-P
0.5

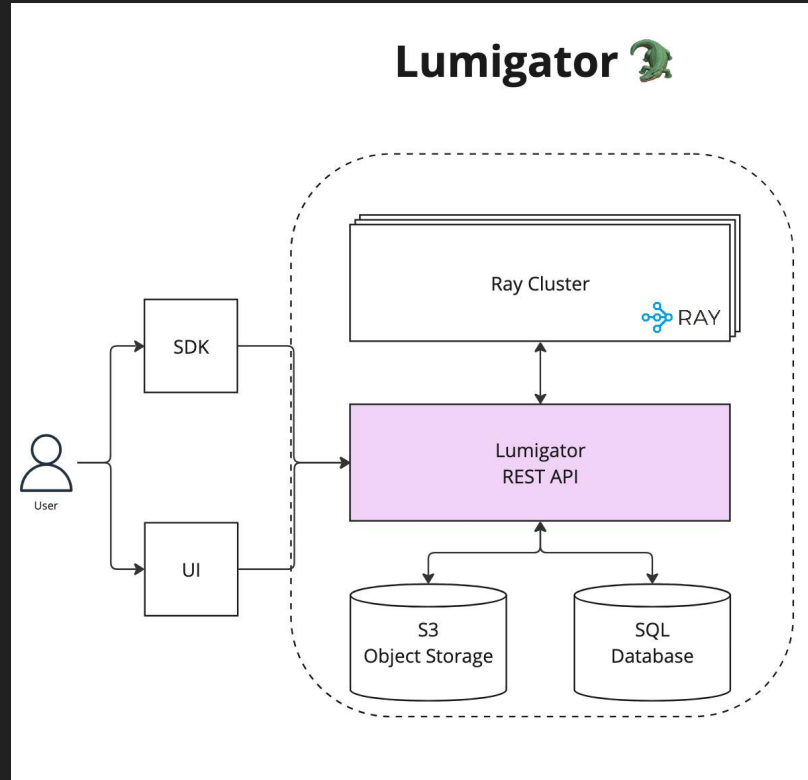↗ View Results          ⬇ Download Results

| EXAMPLES | GROUND TRUTH | PREDICTIONS | ROUGE-1 ↑↓ | ROUGE-2 ↑↓ | ROUGE-L ↑↓ | METEOR ↑↓ | BERT P ↑↓ | BERT F1 ↓ |
|---|---|---|---|---|---|---|---|---|
| 1 #Person1#: Did you go to school today? #Person2#: Of course. Did you? #Person1#: I didn't want to, so I didn't. #Person2#: That's sad, but have you gone to the movies recently? #Person1#: That's a switch. #Person2#: I'm serious, have you? #Person1#: No, I haven't. Why? #Person2#: I really want to go to the movies this weekend. #Person1#: So go then. #Person2#: I really don't want to go by myself. #Person1#: Well anyway, do you plan on going to school tomorrow? #Person2#: No, I think I'm going to go to the movies. | #Person1# didn't go to school today. #Person2# wants to skip class tomorrow to go to the movies. | Person1 and Person2 discuss their decision to skip school, with Person1 admitting not wanting to attend and Person2 expressing their desire to go to the movies. Person2 ultimately decides not to go to school the next day and prefers to watch a movie instead, but is hesitant to go alone. | 0.35 | 0.21 | 0.26 | 0.28 | 0.89 | 0.90 |
| 2 #Person1#: Hello, how are you doing today? #Person2#: I ' Ve been having trouble breathing lately. #Person1#: Have you had any type of cold lately? #Person2#: No, I haven ' t had a cold. I just have a heavy feeling in my chest when I try to breathe. #Person1#: Do you have any allergies that you know of? #Person2#: No, I don ' t have any allergies that I know of. #Person1#: Does this happen all the time or mostly when you are active? #Person2#: It happens a lot when I work out. #Person1#: I am going to send you to a pulmonary specialist who can run tests on you for asthma. #Person2#: Thank you for your help, doctor. | #Person2# has trouble breathing. The doctor asks #Person2# about it and will send #Person2# to a pulmonary specialist. | Person2 is experiencing trouble breathing, particularly a heavy feeling in the chest during physical activity, and has not had a cold or known allergies. Person1, likely a doctor, suggests referring Person2 to a pulmonary specialist for asthma testing. | 0.43 | 0.19 | 0.36 | 0.41 | 0.89 | 0.89 |
| 3 #Person1#: Hey Jimmy. Let's go workout later today. #Person2#: Sure. What time do you want to go? #Person1#: How about at 3:30? #Person2#: That sounds good. Today we work on Legs and forearm. #Person1#: Hey. I just played basketball earlier, so my legs are a little sore. Let's | #Person1# invites Jimmy to go workout and persuades him into working out on arms and stomach. | Person1 suggests going to work out at 3:30, but after mentioning soreness in their legs, they propose focusing on arms and stomach instead. Despite initial hesitation, Person2 agrees to this change and confirms the workout time. | 0.34 | 0.12 | 0.26 | 0.40 | 0.87 | 0.89 |

# OpenAI vs more OSS models

| MODEL | METEOR | BERT P | BERT R | BERT F1 | ROUGE-1 | ROUGE-2 | ROUGE-L | MODEL SIZE | PARAMETERS | RUN TIME |
|---|---|---|---|---|---|---|---|---|---|---|
| > gpt-4o-mini | 0.32 | 0.88 | 0.89 | 0.88 | 0.37 | 0.17 | 0.30 | | | 00: 00: 28 |
| > facebook/bart-large-cnn | 0.39 | 0.88 | 0.87 | 0.87 | 0.33 | 0.18 | 0.27 | 1.63 GB | 406 M | 00: 04: 01 |
| > open-mistral-7b | 0.28 | 0.88 | 0.87 | 0.88 | 0.30 | 0.08 | 0.24 | | | 00: 07: 38 |
| > gpt-4o | 0.30 | 0.88 | 0.88 | 0.88 | 0.28 | 0.11 | 0.24 | | | 00: 00: 51 |
| > Falconsai/text_summarization | 0.45 | 0.87 | 0.89 | 0.88 | 0.32 | 0.12 | 0.25 | 242 MB | 60.5 M | 00: 07: 21 |

# Architecture

Mozilla.ai

# Web UI

# REST API

# Ray backend

# Object Store

# SDK

# Choose your own component

- Different classes of supported models
  - HuggingFace Transformers
  - APIs
  - Local models (e.g. llamafile, vLLM) through OpenAI-compatible APIs
- Different type of jobs
  - Simple (e.g. annotation) as well as composite (e.g. inference+evaluation)
- Different access
  - UI, Python SDK, CURL calls
- Different deployments
  - Local, cluster, cloud

# What next?

- 

**Mozilla.ai**

Lumigator Public Roadmap H1 2025

Board | New view

label:initiative

## 🕵️ In Discovery 2

**lumigator #628**
Expand LLM Evaluation Use Case - Translation

**lumigator #716**
Ease of Use and Ease of Contribution

## 🔄 In Progress 1
This is actively being worked on

**lumigator #629**
Extend Lumigator Evaluation Metrics

# Wanna know more?
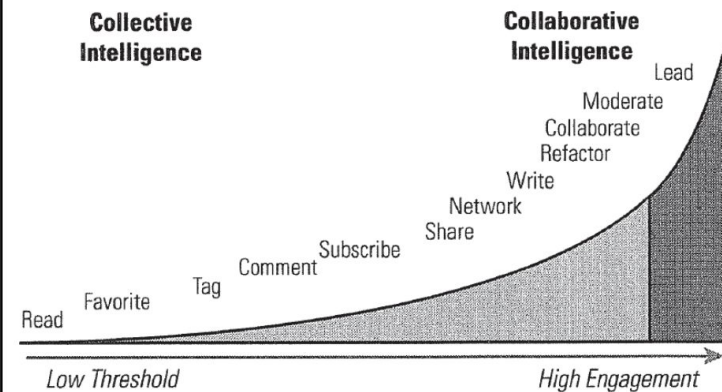
- [This talk](#)
- [Code](#)
- [Documentation](#)
- [Roadmap](#)

# Wanna contribute?

- [Try it out](#)
- Submit a [bug](#)/[feature request](#)
- Check our [documentation](#)
- Pick up a [good first issue](#)
- Add a new [feature](#) :-)

See our [CONTRIBUTING.md](#)



**Figure 2.2** The power law of participation

Collective Intelligence

Collaborative Intelligence

Lead
Moderate
Collaborate
Refactor
Write
Network
Share
Subscribe
Comment
Tag
Favorite
Read

Low Threshold

High Engagement

Source: Mayfield 2006

# Thank you!