

CONF'

Xavier Martinet

AI Research Engineer @ Meta



“Les dessous de Llama 3”



Le **18/09/2024** à **19h**

Au Palace, 4 rue Voltaire, Nantes



icilundi

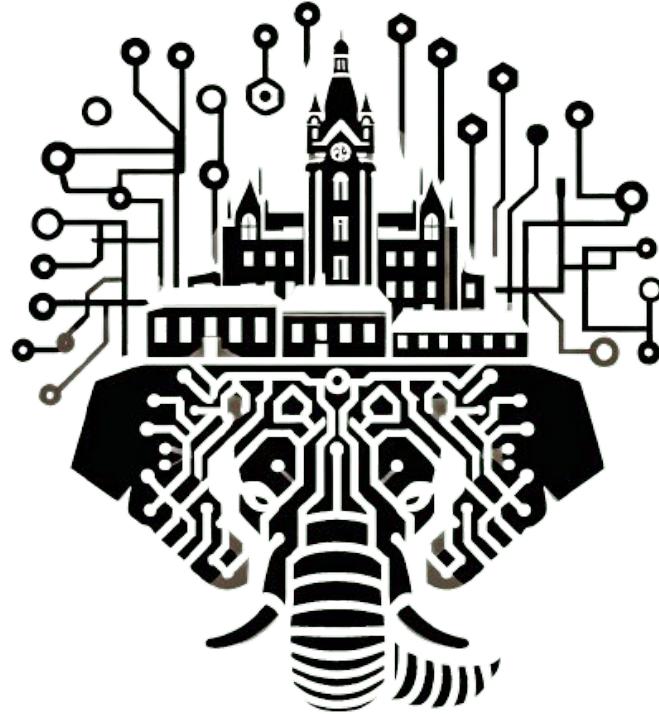


[sfΞir] **lonestone**

Qui n'est jamais venu au meetup Gen AI Nantes ?

GenAI Nantes

- Microphone icon: 15 événements / an
- Pirate flag icon: 1 hackathon
- Worker icon: 2 workshops
- Heart icon: 1 communauté de 400p

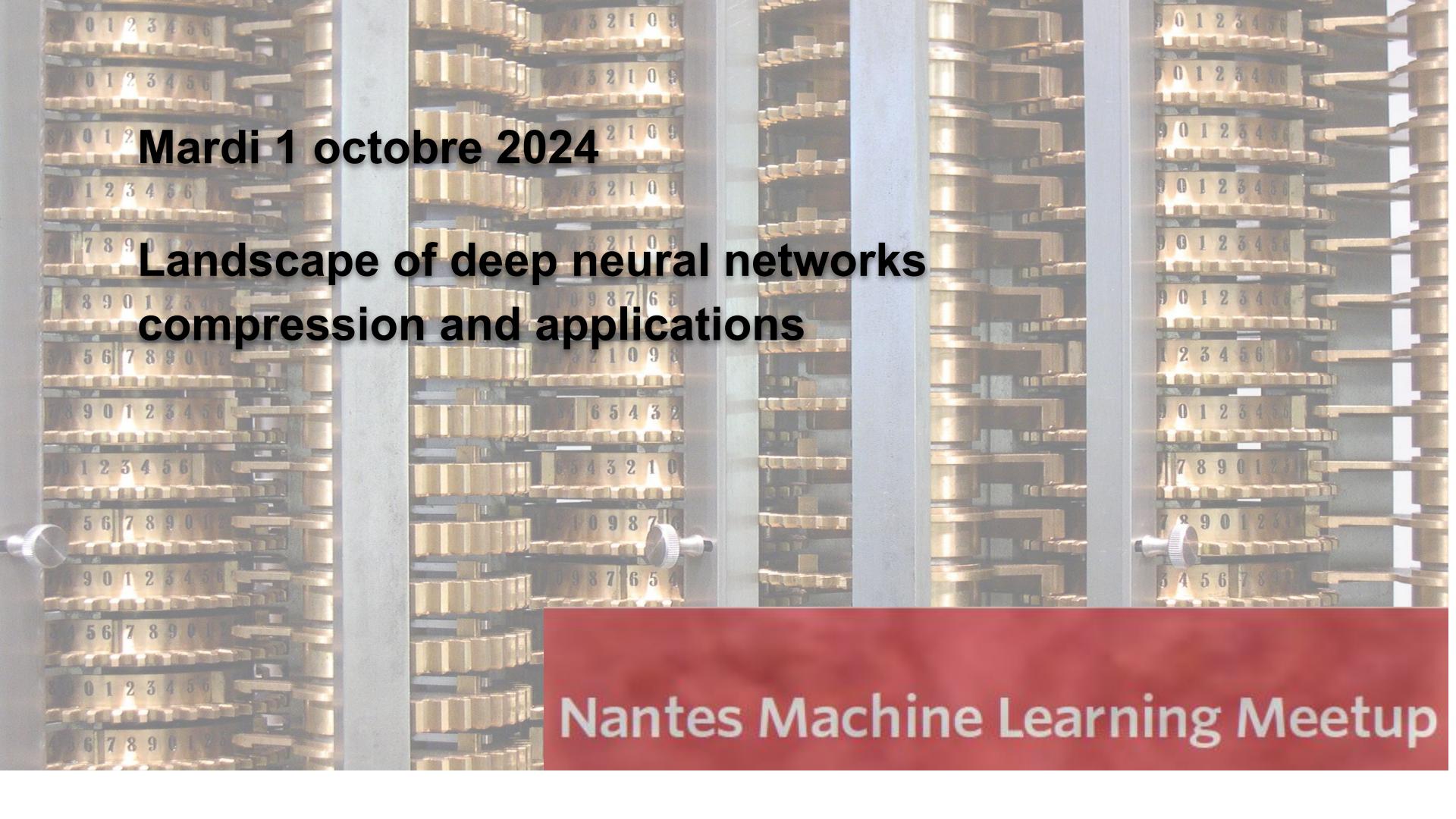




Save the date!

Shift revient en 2025,
du 28 au 30 mars.





Mardi 1 octobre 2024

Landscape of deep neural networks compression and applications

Nantes Machine Learning Meetup

Qui fait des projets GenAI ?



Qui chevauche des Llama* ?

*aucun Llama ne sera blessé pendant la soirée



Qui chevauche des Llama~~405B~~ mammouth ?

*aucun Llama ne sera blessé pendant la soirée

Schedule



1- News summer 2024



2- A multimodal Llama



3- Enjoy

The Llama 3 Herd of Models

July 23, 2024

Abstract

Modern artificial intelligence (AI) systems are powered by foundation models. This paper presents a new set of foundation models, called Llama 3. It is a herd of language models that natively support multilinguality, coding, reasoning, and tool usage. Our largest model is a dense Transformer with 405B parameters and a context window of up to 128K tokens. This paper presents an extensive empirical evaluation of Llama 3. We find that Llama 3 delivers comparable quality to leading language models such as GPT-4 on a plethora of tasks. We publicly release Llama 3, including pre-trained and post-trained versions of the 405B parameter language model and our Llama Guard 3 model for input and output safety. The paper also presents the results of experiments in which we integrate image, video, and speech capabilities into Llama 3 via a compositional approach. We observe this approach performs competitively with the state-of-the-art on image, video, and speech recognition tasks. The resulting models are not yet being broadly released as they are still under development.

 Download the Paper

Paper Llama3.1

- Datasets
- Tools
- Function calling
- Multimodal
- Security
- Long context
- Coding
- Reasoning
- ...

<https://ai.meta.com/research/publications/the-llama-3-herd-of-models/>

Focus on long context

Very large context ?

Use case:

- Q/A on large document
- RAG

LLM:

- Llama3.1: 128k
- Gemini 1.5 pro: **1m**
- GPT-4o: 128k
- Claude Sonet 3.5: **200k**
- Mistral Large 2: 128k

Very large context ?

Cost is quadratic

- Training
- Inference

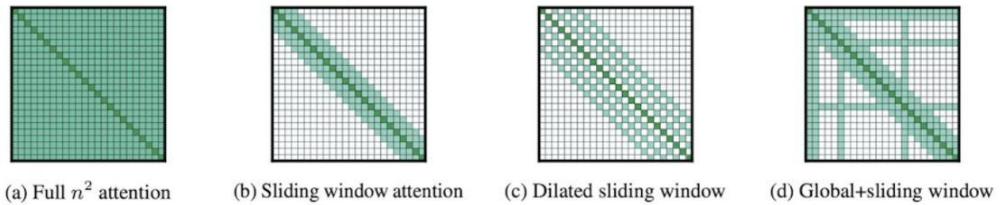


Figure 2: Comparing the full self-attention pattern and the configuration of attention patterns in our Longformer.

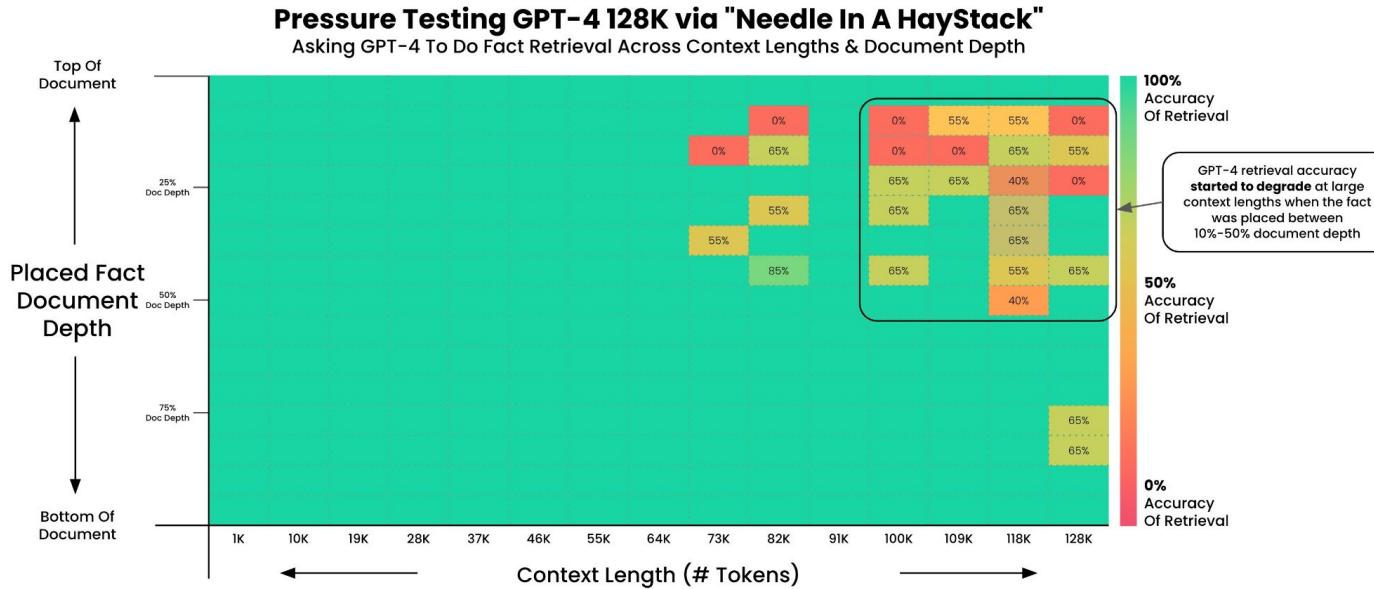
Large context training for Llama3

- * Trained with 8k token windows (14.000 B tokens)
- * Gradually increased to 128k for the last (800 B tokens)

The “needle in a haystack” benchmark (NIAH)

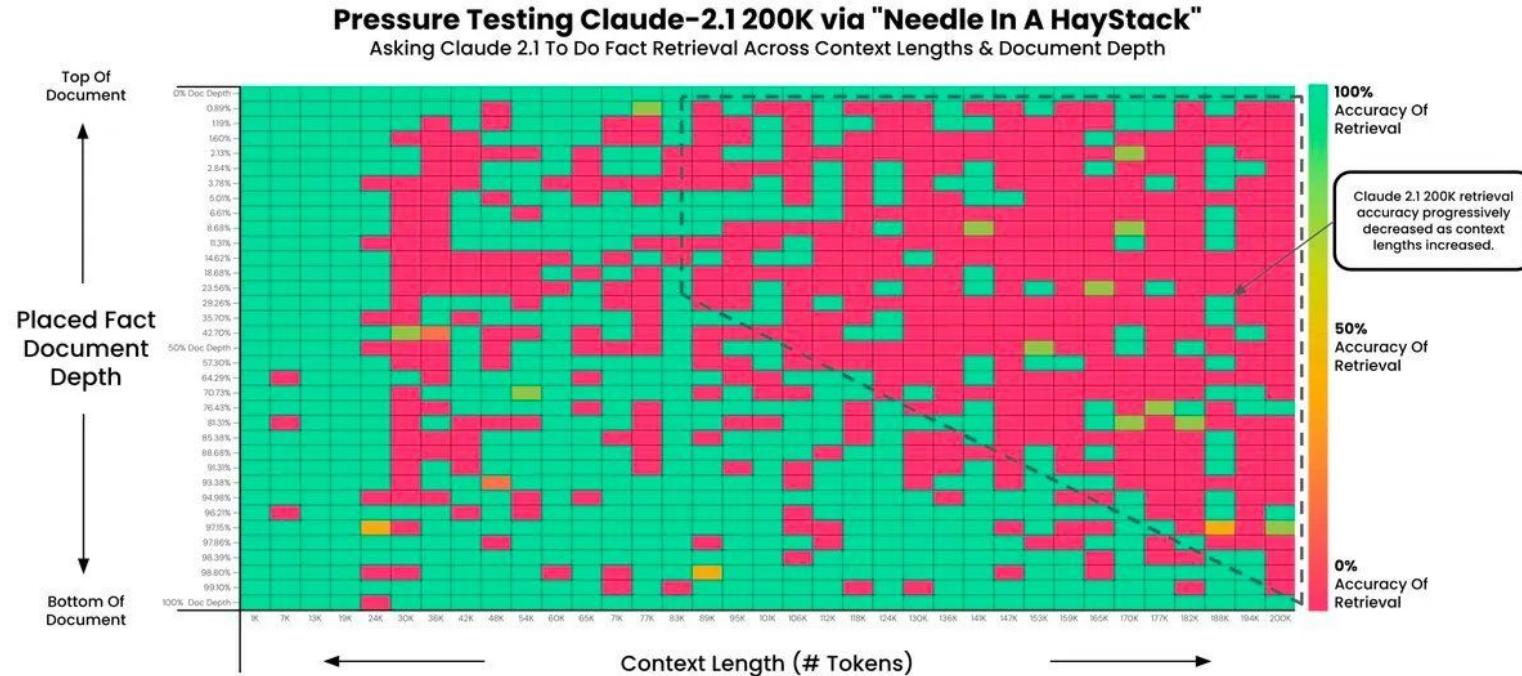
Let's benchmark a simple **retrieval** task

Single retrieval - GPT4-128k



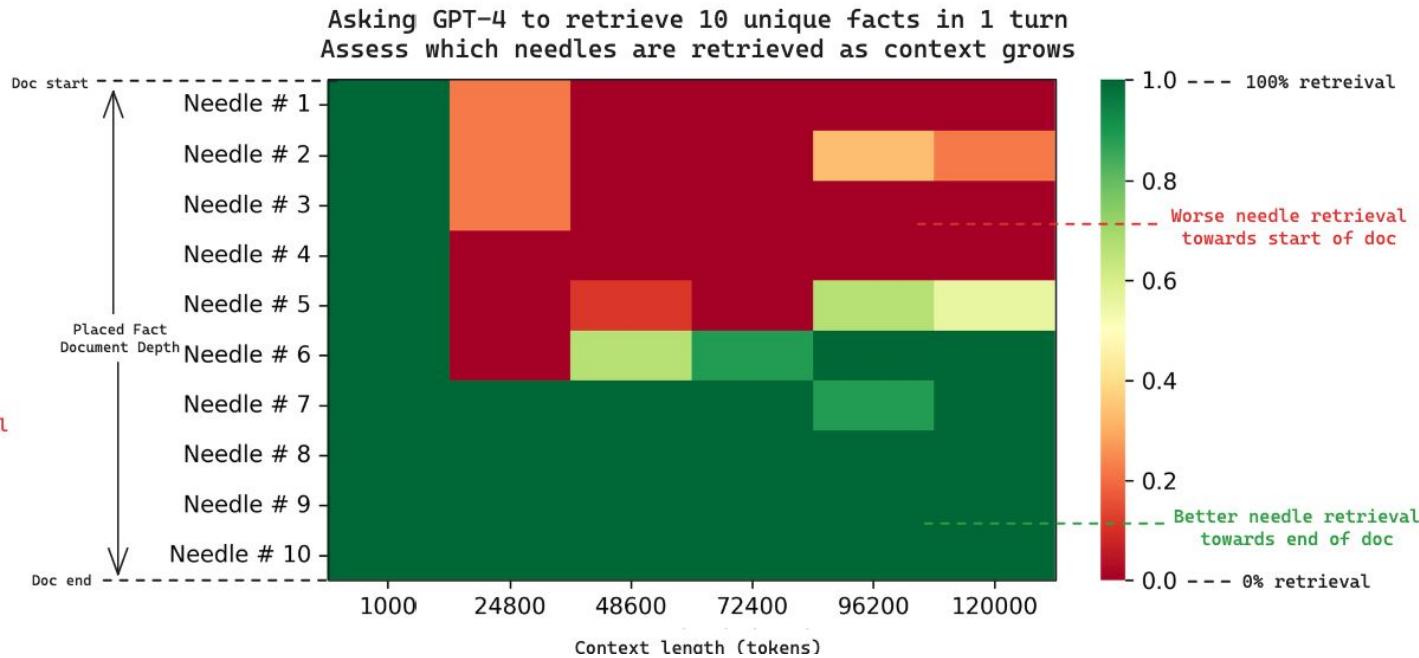
Goal: Test GPT-4 Ability To Retrieve Information From Large Context Windows
A fact was placed within a document. GPT-4 (1106-preview) was then asked to retrieve it. The output was evaluated for accuracy.
This test was run at 15 different document depths (top > bottom) and 15 different context lengths (1K > 128K tokens).
2x tests were run for larger contexts for a larger sample size.

Single retrieval - Claude 2.1 200k



Goal: Test Claude 2.1 Ability To Retrieve Information From Large Context Windows
A fact was placed within a document. Claude 2.1 (200K) was then asked to retrieve it. The output was evaluated (with GPT-4) for accuracy.
This test was run at 35 different document depths (top > bottom) and 35 different context lengths (1K > 200K tokens).
Document Depths followed a sigmoid distribution

Multi-needle GPT4-128k



The RULER benchmark

RULER: What's the Real Context Size of Your Long-Context Language Models?

Cheng-Ping Hsieh*, Simeng Sun*, Samuel Kriman, Shantanu Acharya
Dima Rekesh, Fei Jia, Yang Zhang, Boris Ginsburg

NVIDIA

{chsieh,simengs}@nvidia.com

Abstract

The needle-in-a-haystack (NIAH) test, which examines the ability to retrieve a piece of information (the “needle”) from long distractor texts (the “haystack”), has been widely adopted to evaluate long-context language models (LMs). However, this simple retrieval-based test is indicative of only a superficial form of long-context understanding. To provide a more comprehensive evaluation of long-context LMs, we create a new synthetic benchmark RULER with flexible configurations for customized sequence length and task complexity. RULER expands upon the vanilla NIAH test to encompass variations with diverse types and quantities of needles. Moreover, RULER introduces new task categories *multi-hop tracing* and *aggregation* to test behaviors beyond searching from context. We evaluate 17 long-context LMs with 13 representative tasks in RULER. Despite achieving nearly perfect accuracy in the vanilla NIAH test, almost all models exhibit large performance drops as the context length increases. While these models all claim context sizes of 32K tokens or greater, only half of them can maintain satisfactory performance at the length of 32K. Our analysis of Yi-34B, which supports context length of 200K, reveals large room for improvement as we increase input length and task complexity. We open source RULER to spur comprehensive evaluation of long-context LMs.

- Multiple needle retrieval
- Multi-hop needle
- Q/A
- More advanced tasks

<https://arxiv.org/pdf/2404.06654>

The RULER benchmark

Models	Claimed Length	Effective Length	4K	8K	16K	32K	64K	128K	Avg.	wAvg. (inc)	wAvg. (dec)
Llama2 (7B)	4K	-	85.6								
Gemini-1.5-Pro	1M	>128K	96.7	95.8	96.0	95.9	95.9	94.4	95.8	95.5 _(1st)	96.1 _(1st)
GPT-4	128K	64K	96.6	96.3	95.2	93.2	87.0	81.2	91.6	89.0 _(2nd)	94.1 _(2nd)
Llama3.1 (70B)	128K	64K	96.5	95.8	95.4	94.8	88.4	66.6	89.6	85.5 _(4th)	93.7 _(3rd)
Qwen2 (72B)	128K	32K	96.9	96.1	94.9	94.1	79.8	53.7	85.9	79.6 _(9th)	92.3 _(4th)
Command-R-plus (104B)	128K	32K	95.6	95.2	94.2	92.0	84.3	63.1	87.4	82.7 _(7th)	92.1 _(5th)
GLM4 (9B)	1M	64K	94.7	92.8	92.1	89.9	86.7	83.1	89.9	88.0 _(3rd)	91.7 _(6th)
Llama3.1 (8B)	128K	32K	95.5	93.8	91.6	87.4	84.7	77.0	88.3	85.4 _(5th)	91.3 _(7th)
GradientAI/Llama3 (70B)	1M	16K	95.1	94.4	90.8	85.4	80.9	72.1	86.5	82.6 _(8th)	90.3 _(8th)
Mixtral-8x22B (39B/141B)	64K	32K	95.6	94.9	93.4	90.9	84.7	31.7	81.9	73.5 _(11th)	90.3 _(9th)
Yi (34B)	200K	32K	93.3	92.2	91.3	87.5	83.2	77.3	87.5	84.8 _(6th)	90.1 _(10th)
Phi3-medium (14B)	128K	32K	93.3	93.2	91.1	86.8	78.6	46.1	81.5	74.8 _(10th)	88.3 _(11th)
Mistral-v0.2 (7B)	32K	16K	93.6	91.2	87.2	75.4	49.0	13.8	68.4	55.6 _(13th)	81.2 _(12th)
LWM (7B)	1M	<4K	82.3	78.4	73.7	69.1	68.1	65.0	72.8	69.9 _(12th)	75.7 _(13th)
DBRX (36B/132B)	32K	8K	95.1	93.8	83.6	63.1	2.4	0.0	56.3	38.0 _(14th)	74.7 _(14th)
Together (7B)	32K	4K	88.2	81.1	69.4	63.0	0.0	0.0	50.3	33.8 _(15th)	66.7 _(15th)
LongChat (7B)	32K	<4K	84.7	79.9	70.8	59.3	0.0	0.0	49.1	33.1 _(16th)	65.2 _(16th)
LongAlpaca (13B)	32K	<4K	60.6	57.0	56.6	43.6	0.0	0.0	36.3	24.7 _(17th)	47.9 _(17th)

So how to support long context?

- Map/Reduce
- Mamba2 architecture
- Ask Xavier to work harder 😅

Mamba2 architecture (June '24)

- Linear complexity
- Good for long (infinite?) context
- Transformers are stronger for short context
- Optimized memory usage
- Codestral Mamba:
<https://mistral.ai/news/codestral-mamba/>
- 100m tokens context window:
<https://magic.dev/blog/100m-token-context-windows>

Context caching

- Offload the attention state and K/V cache
- Process attention once and read many times
 - Set a TTL on each query

Context caching

INPUT PRICING	\$7.00 / 1 million tokens
OUTPUT PRICING	\$21.00 / 1 million tokens
CONTEXT CACHING	\$1.75 / 1 million tokens
<hr/>	
CONTEXT CACHING (STORAGE)	\$4.50 / 1 million tokens per hour

Claude 3.5 Sonnet

- Our most intelligent model to date
- 200K context window

Input

- \$3 / MTok

Prompt caching

- \$3.75 / MTok - Cache write
- \$0.30 / MTok - Cache read

Output

- \$15 / MTok

Multiple/parallel function calling

tool_choice

Users can use `tool_choice` to specify how tools are used:

- "auto": default mode. Model decides if it uses the tool or not.
- "any": forces tool use.
- "none": prevents tool use.

```
import os
from mistralai import Mistral

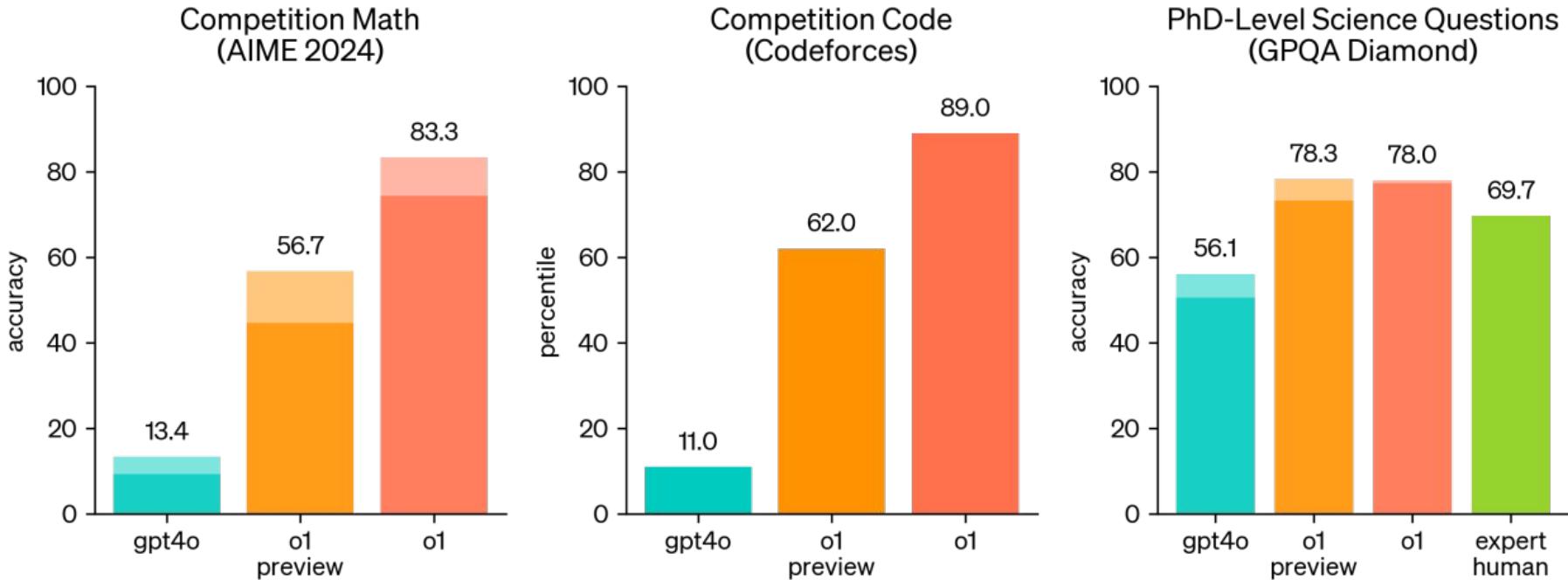
api_key = os.environ["MISTRAL_API_KEY"]
model = "mistral-large-latest"

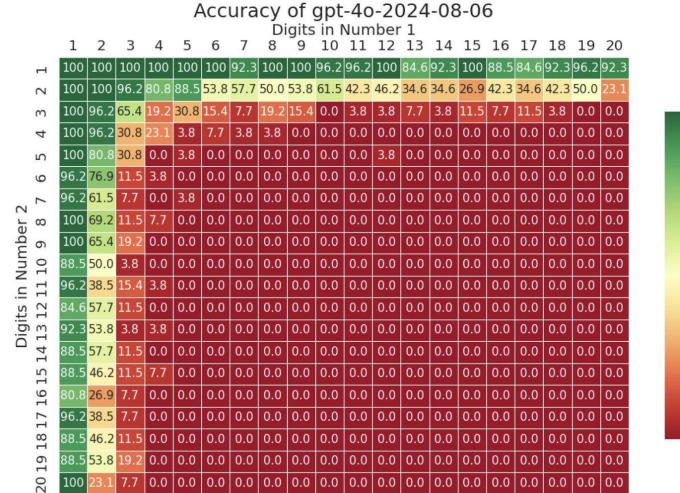
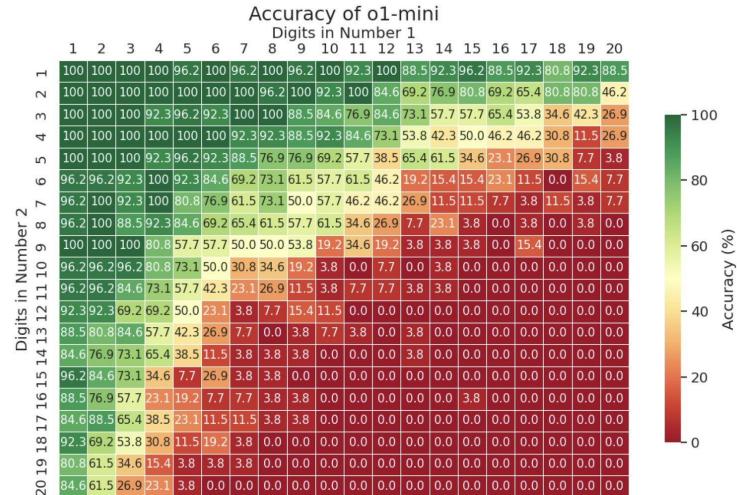
client = Mistral(api_key=api_key)
response = client.chat.complete(
    model=model,
    messages=messages,
    tools=tools,
    tool_choice="any",
)
response
```

OpenAI

Reasoning







<https://x.com/dallaslones/status/1834313445677605256>

SearchGPT●

Q What are you looking for?



<https://x.com/OpenAI/status/1816536290822881780>

iPhone 16 Pro

Hello, Apple Intelligence.





Reduce Interruptions

in Focus



Genmoji

Create a Memory Movie

Describe a Memory...

Create a Memory Movie

Create a Memory Movie

Create a Memory Movie



A more personal Siri

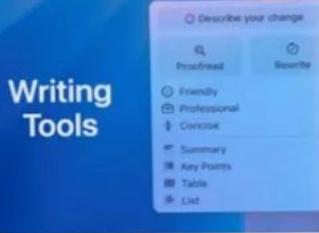
Private Cloud Compute



Clean Up in Photos

Summaries in Messages

Writing Tools



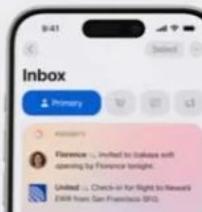
To:
Cc:
Subject:
From:

Dear Ms. H,
It was great to my letter
cover letter.

Thanks,
Jenny Fitt
Dept. of Jax

Apple Intelligence

Priority messages in Mail

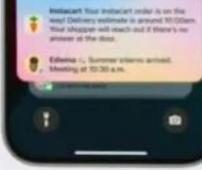


Inbox

Priority

Florence ... Invited an Online with opening to Florence tonight.
United ... Check-in for flight to Newark 2018 from San Francisco (SFO).

Priority notifications



Priority notifications

Image Playground



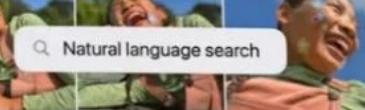
Image Wand



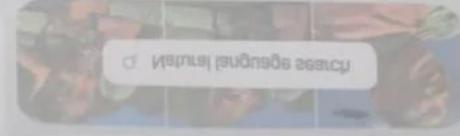
Audio recording



summaries



Natural language search



<https://x.com/joffrey/status/1829081299576914091>



<https://x.com/GuangyuRobert/status/1831006762184646829>



sakana.ai

首页

探索

资产

个人主页

AI 工具

图片生成

智能画布

视频生成

故事创作 Beta

消息中心

反馈

邀请

AI 作图

轻松实现创意图片

图片生成

智能画布

图片

视频

热门

摄影

插画

设计

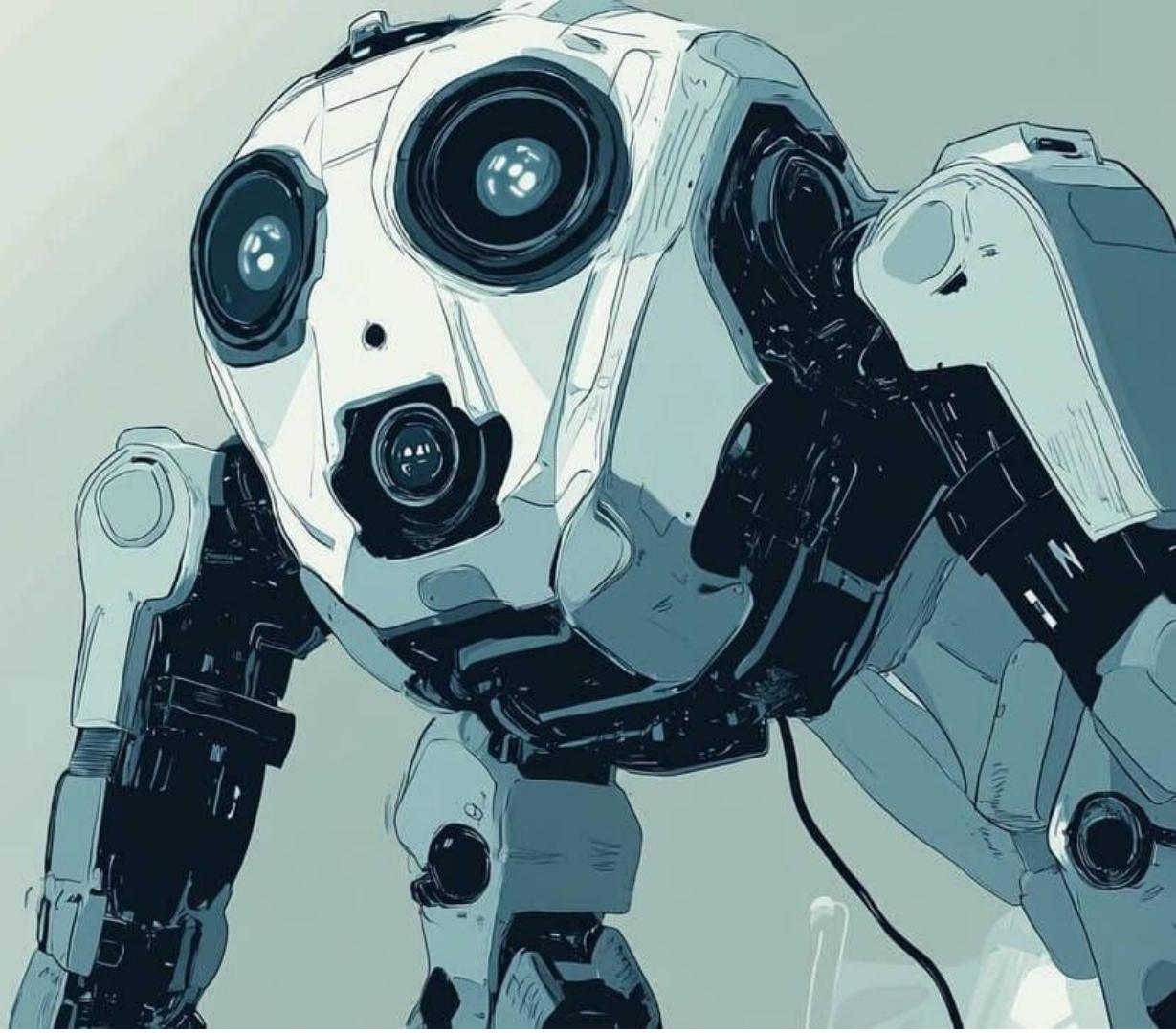


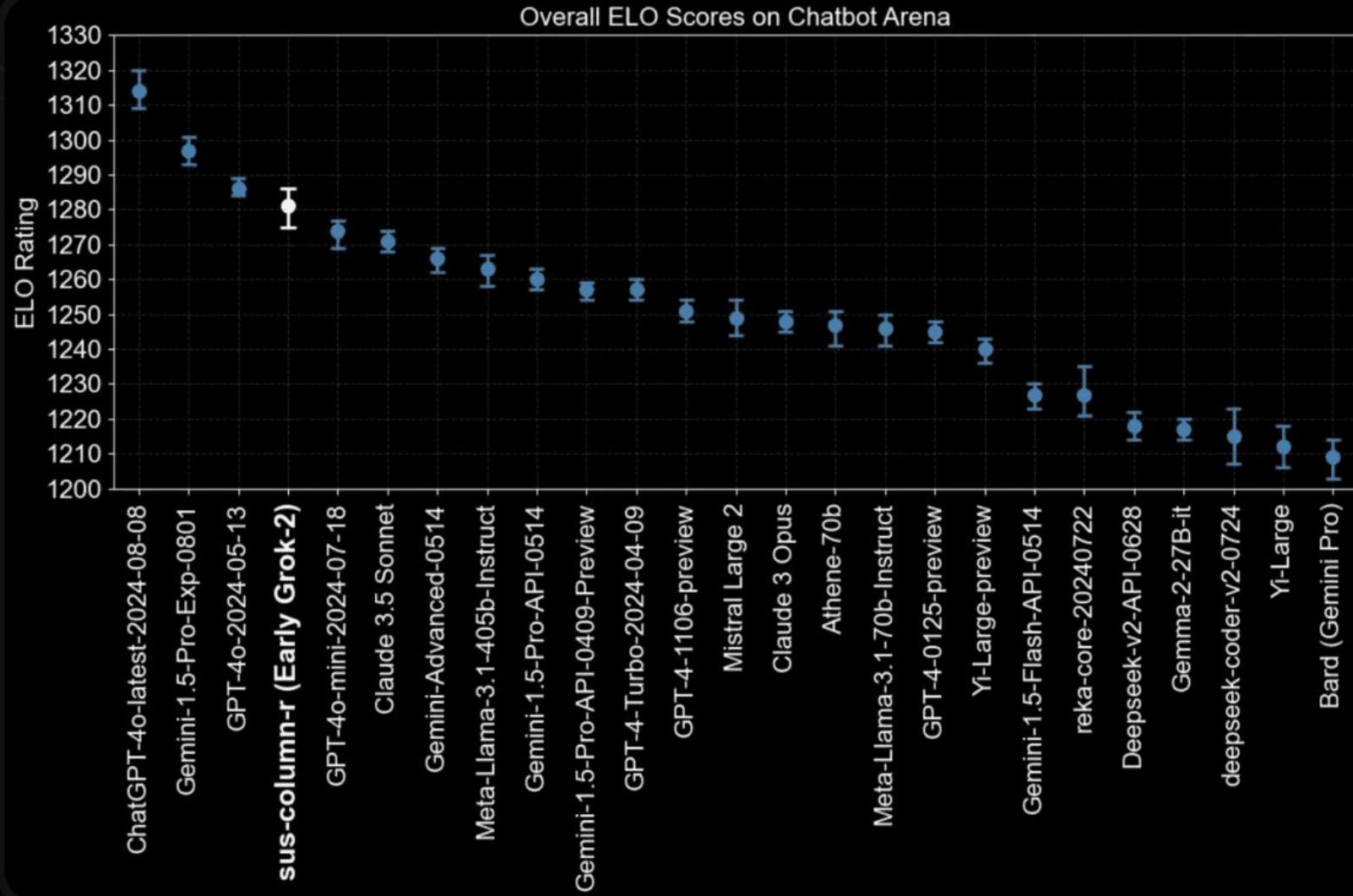
AI 视频

让你的创意动起来

视频生成

<https://www.youtube.com/shorts/pqTX4GI4uGc?app=desktop>







Llama 3.1



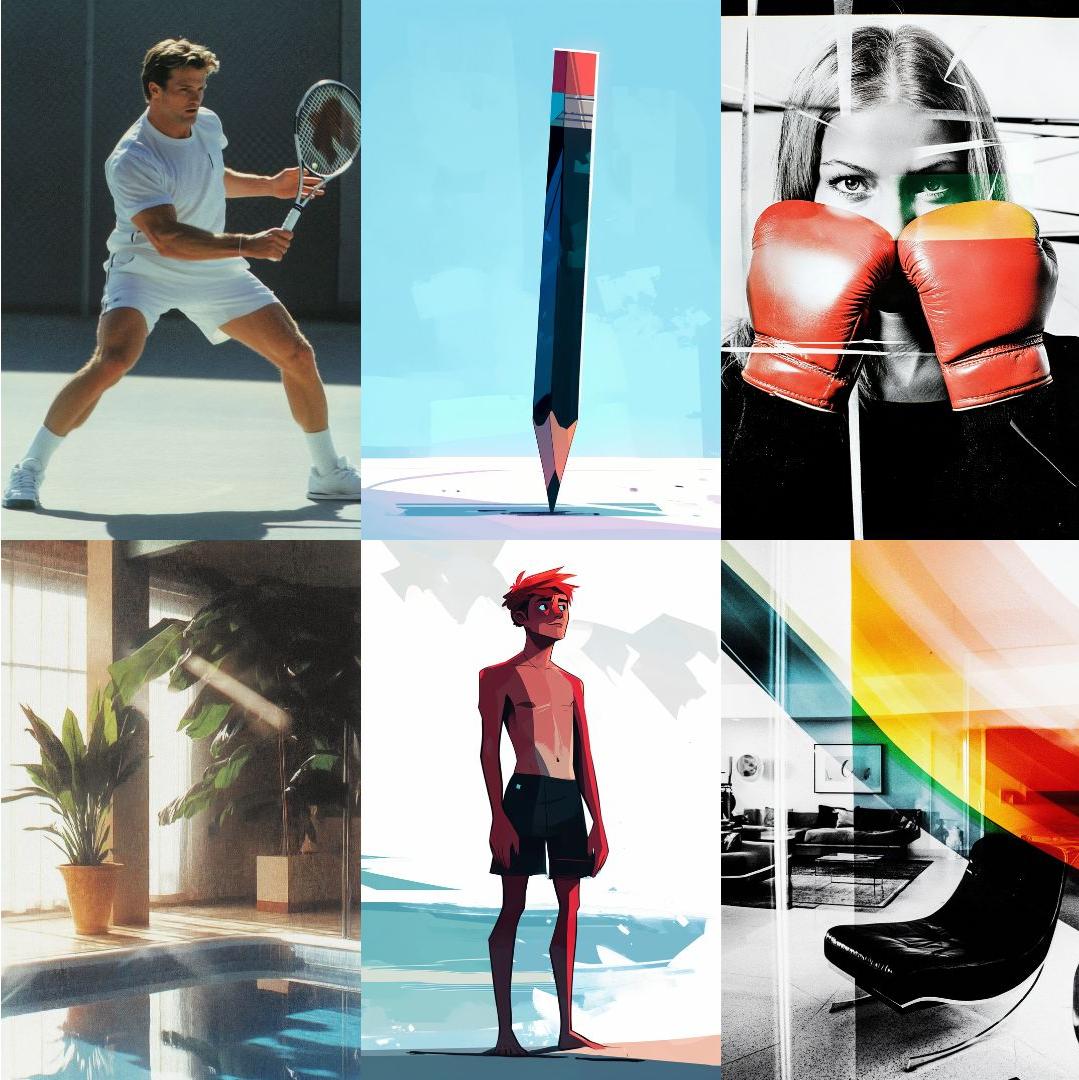
FLUX

PAR BLACK FOREST



RESSOURCES MIDJOURNEY

CODES SREF





CHIC AVANT-GARDE EDITORIAL PHOTO
--SREF 3752650426

PROMPT : PHOTO OF [SUJET], BLACK AND
WHITE, WHITE --SREF 3752650426 --AR 2:3













📄 --sref 44243



📄 --sref 44242



📄 --sref 33464



📄 --sref 88019



📄 --sref 87023



📄 --sref 85169



📄 --sref 56321



📄 --sref 65982



📄 --sref 69145





CONF'

Xavier Martinet

AI Research Engineer @ Meta



“Les dessous de Llama 3”



Le **18/09/2024** à **19h**

Au Palace, 4 rue Voltaire, Nantes



icilundi



[sfΞir] **lonestone**