

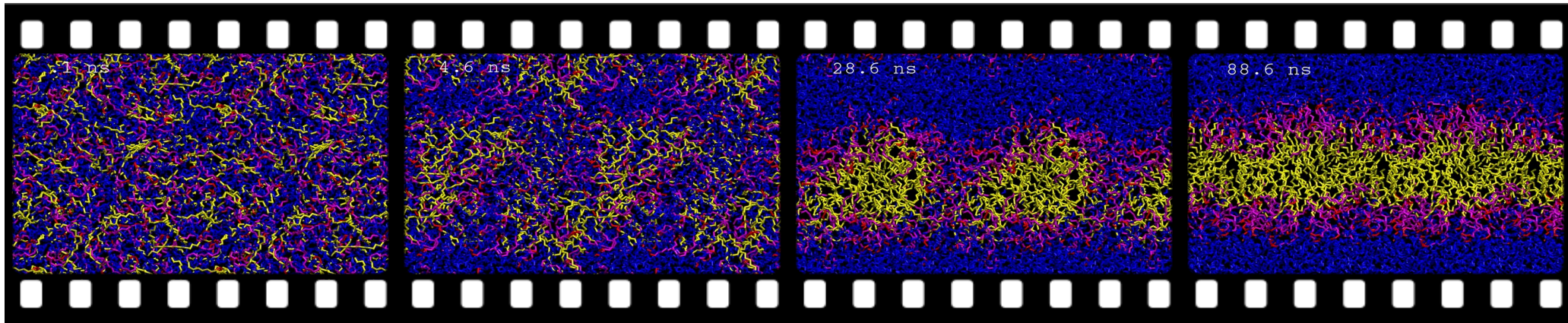


Science ouverte, données et IA : automatiser l'annotation scientifique

Pierre Poulain

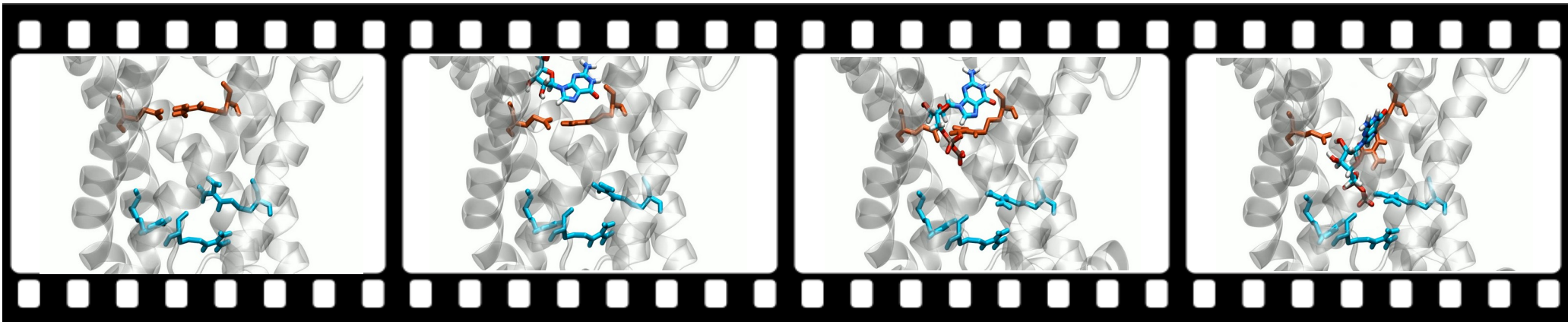
Laboratoire de Biochimie Théorique
UMR 8266 CNRS & Université Paris Cité

What is molecular dynamics (MD)?



water + detergent

Senac et al, Langmuir, 2017. Movie by Patrick Fuchs.



Gagelin et al, Nature communications, 2023.

Supercomputers → high cost



Source : Photothèque CNRS/Cyril Frésillon (droits réservés)

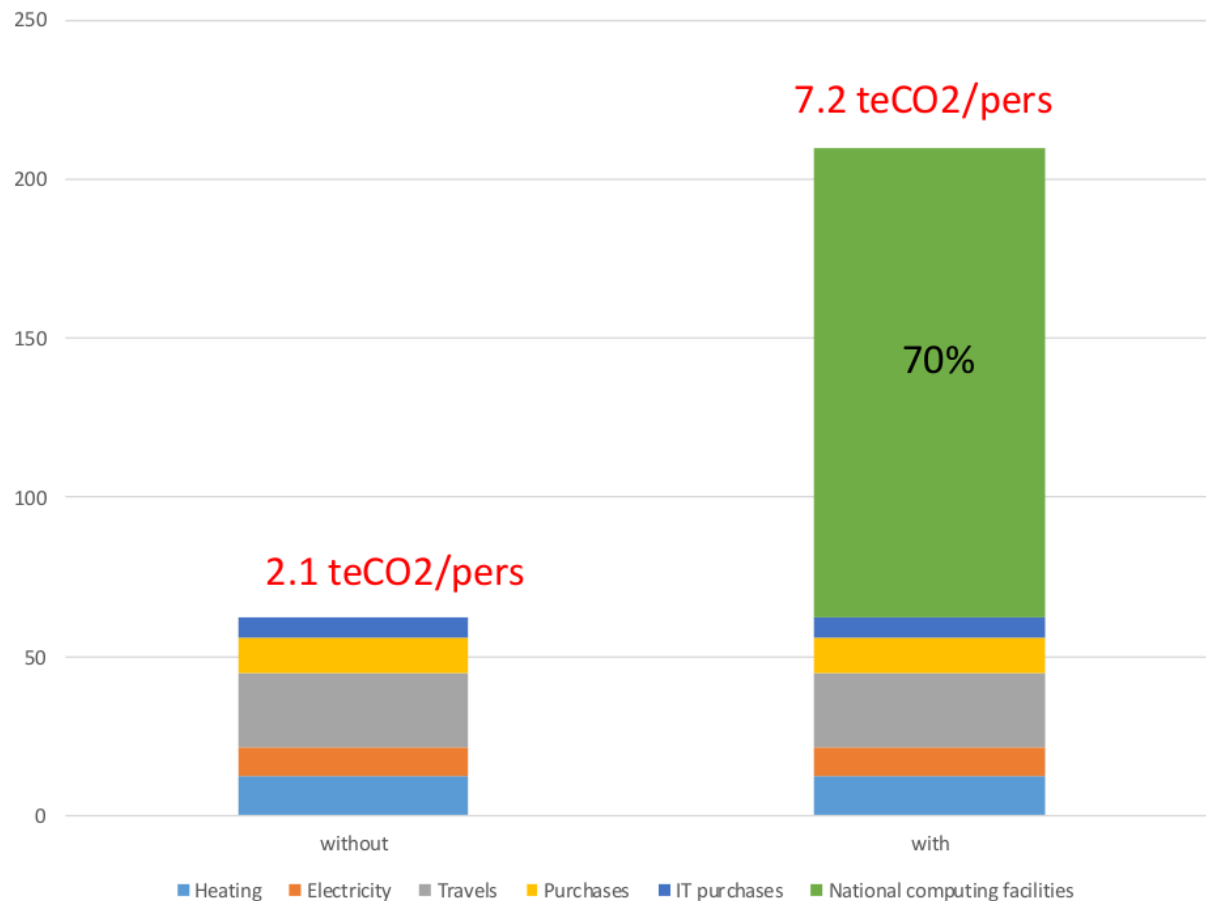


In 2025:

- CT 7 (simulation in biology)
- 132 Mh CPU
- 8 Mh GPU
- **Total cost ~8 M€**

Patrick Fuchs, 2026.

Supercomputers → high environmental cost



Welcome to the Dark Matter of MD

Data that is **technically accessible**,
but neither **indexed, curated**,
or easily **searchable**.



eLife RESEARCH ARTICLE

MDverse, shedding light on the dark matter of molecular dynamics simulations

Johanna KS Tiemann^{1*}, Magdalena Szczuka², Lisa Bouarroudj³, Mohamed Oussaren⁴, Steven Garcia⁵, Rebecca J Howard⁶, Lucie Delemotte⁶, Erik Lindahl^{6,8}, Marc Baaden⁷, Kresten Lindorff-Larsen⁹, Matthieu Chavent^{2,*,} Pierre Poulain⁸

¹Linderstrøm-Lang Centre for Protein Science, Department of Biology, University of Copenhagen, Copenhagen, Denmark; ²Institut de Pharmacologie et Biologie Structurale, CNRS, Université de Toulouse, Toulouse, France; ³Université Paris Cité, CNRS, Institut Jacques Monod, Paris, France; ⁴Independent researcher, Amsterdam, Netherlands; ⁵Department of Biochemistry and Biophysics, Science for Life Laboratory, Stockholm University, Stockholm, Sweden; ⁶Department of applied physics, Science for Life Laboratory, KTH Royal Institute of Technology, Stockholm, Sweden; ⁷Laboratoire de Biochimie Théorique, CNRS, Université Paris Cité, Paris, France

*For correspondence: johanna.tiemann@gmail.com (JKST); matthieu.chavent@ipbs.fr (MC); pierre.poulain@u-paris.fr (PP)

Present address: ⁸NovozymesA/S, Lyngby, Denmark

Competing interest: See page 16

Funding: See page 16

Preprint posted: 02 May 2023

Sent for Review: 28 June 2023

Reviewed preprint posted: 20 September 2023

Reviewed preprint revised: 07 July 2024

Version of Record published: 30 August 2024

Reviewing Editor: Shoaib Haider, University College London, United Kingdom

© Copyright Tiemann et al. This article is distributed under the terms of the [Creative Commons Attribution License](#), which permits unrestricted use and redistribution provided that the original author and source are credited.

Abstract The rise of open science and the absence of a global dedicated data repository for molecular dynamics (MD) simulations has led to the accumulation of MD files in generalist data repositories, constituting the *dark matter of MD*— data that is technically accessible, but neither indexed, curated, or easily searchable. Leveraging an original search strategy, we found and indexed about 250,000 files and 2000 datasets from Zenodo, Figshare and Open Science Framework. With a focus on files produced by the Gromacs MD software, we illustrate the potential offered by the mining of publicly available MD data. We identified systems with specific molecular composition and were able to characterize essential parameters of MD simulation such as temperature and simulation length, and could identify model resolution, such as all-atom and coarse-grain. Based on this analysis, we inferred metadata to propose a search engine prototype to explore the MD data. To continue in this direction, we call on the community to pursue the effort of sharing MD data, and to report and standardize metadata to reuse this valuable matter.

eLife assessment The study presents a **valuable** tool for searching molecular dynamics simulation data, making such datasets accessible for open science. The authors provide **convincing** evidence that it is possible to identify noteworthy molecular dynamics simulation datasets and that their analysis can produce information of value to the community.

Introduction The volume of data available in biology has increased tremendously (Marx, 2013; Stephens et al., 2015), through the emergence of high-throughput experimental technologies, often referred to as -omics, and the development of efficient computational techniques, associated with high-performance computing resources. The Open Access (OA) movement to make research results free and available to anyone (including e.g. the Budapest Open Access Initiative and the Berlin declaration on Open Access to Knowledge) has led to an explosive growth of research data made available by scientists (Wilson

Tiemann et al. eLife 2023;12:RP90061. DOI: <https://doi.org/10.7554/eLife.90061> 1 of 22

A data catalogue for MD simulation data (and more)

MDverse

MDverse Search

I'm Feeling Lucky

WHY? Enable, facilitate & foster the reuse of MD data

MD data availability

Data repository	First dataset	# datasets	# files	Size (GB)
Zenodo	2014	2 369	504 891	26 686
Figshare	2012	1 373	222 703	1 177
NOMAD	2021	16 119	555 990	2 189
MDposit CINECA (MDDDB)		310	251 581	52
MDposit INRIA (MDDDB)	2025	719	186 293	1 022
MDposit MMB (MDDDB)		4 050	411 854	15 177
GPCRmd	2017	830	2 414	624
ATLAS	2023	1 938	69 768	10 944
Total		27 708	2 205 494	57 809

Reuse MD data → **Extract & normalize metadata**

Where to find metadata?

MD simulations of P-glycoprotein started from three different crystal structures: 3G5U, 4M1M and 4KSB

> <https://doi.org/10.6084/m9.figshare.4806544>

Cite

Download all (2.55 GB)

Share

Embed

+ Collect

Version 4 ▾

Dataset posted on 2018-01-10, 07:17 authored by [Karmen Condic-Jurkic](#)

Molecular simulations of mouse P-glycoprotein started from three different crystal structures with PDB IDs corresponding to 3G5U, 4KSB and 4M1M. The protein was transferred to POPC/cholesterol bilayer, followed by energy minimization and 10 ns equilibration period. All the simulations were performed in Gromacs 3.3.3 in conjunction with GROMOS54a7 force field at 300 K and pressure of 1 bar. Each system was run in triplicates for 200 ns.

USAGE METRICS

2341

3328

1

VIEWS

DOWNLOADS

CITATIONS 

[Learn more about usage metrics](#) 



Source: Figshare

Named Entity Recognition (NER)

MD simulations of P-glycoprotein MOL started from three different crystal structures: 3G5U MOL , 4M1M MOL and 4KSB MOL .
Molecular simulations of mouse P-glycoprotein MOL started from three different crystal structures with PDB IDs corresponding to 3G5U MOL , 4KSB MOL and 4M1M MOL . The protein was transferred to POPC MOL / cholesterol MOL bilayer, followed by energy minimization and 10 ns equilibration period. All the simulations were performed in Gromacs SOFTNAME 3.3.3 SOFTVERS in conjunction with GROMOS54a7 FFM force field at 300 K TEMP and pressure of 1 bar. Each system was run in triplicates for 200 ns STIME .

Strategies based on language models:

- Encoder-like models (BERT): ~ 100 millions parameters (SLM)
→ *fine tuning* (many manually annotated examples)
- Decoder-like models (GPT): ~ 10-1000 billions parameters (LLM)
→ ***harness*** (instructions + constraints) / **access to models?**

Manual approach

OpenAI, GPT-5.5, 10/06/2026

In the following text, extract all molecule names, software versions, software names, forcefield or model names, simulation temperatures, and simulation times (duration):

MD simulations of P-glycoprotein started from three different crystal structures: 3G5U, 4M1M and 4KSB

Molecular simulations of mouse P-glycoprotein started from three different crystal structures with PDB IDs corresponding to 3G5U, 4KSB and 4M1M. The protein was transferred to POPC/cholesterol bilayer, followed by energy minimization and 10 ns equilibration period. All the simulations were performed in Gromacs 3.3.3 in conjunction with GROMOS54a7 force field at 300 K and pressure of 1 bar. Each system was run in triplicates for 200 ns.

Afficher moins ^

Manual approach

OpenAI, GPT-5.5, 10/06/2026

In the following text, extract all molecule names, software versions, software names, forcefield or model names, simulation temperatures, and simulation times (duration):

MD simulations of P-glycoprotein started from three different crystal structures: 3G5U, 4M1M and 4KSB

Molecular simulations of mouse P-glycoprotein started from three different crystal structures with PDB IDs corresponding to 3G5U, 4KSB and 4M1M. The protein was transferred to POPC/cholesterol bilayer, followed by energy minimization and 10 ns equilibration period. All the simulations were performed in Gromacs 3.3.3 in conjunction with GROMOS54a7 force field at 300 K and pressure of 1 bar. Each system was run in triplicates for 200 ns.

Afficher moins ^

Extracted entities from the text:

Category	Extracted Value
Molecule names	P-glycoprotein (mouse P-glycoprotein), POPC, cholesterol
Software names	Gromacs
Software versions	Gromacs 3.3.3
Forcefield / model names	GROMOS54a7 force field
Simulation temperatures	300 K
Simulation times (duration)	10 ns equilibration period; 200 ns production simulation (each system, run in triplicate)

Additional identifiers present (not requested but mentioned):

- Crystal structure / PDB IDs: 3G5U, 4KSB, 4M1M
- Pressure: 1 bar

Manual approach

OpenAI, GPT-5.5, 10/06/2026

- Automation?
- Model selection?
- Detailed instructions?
- Output format?
- Quality control?



Extracted entities from the text:

Category	Extracted Value
Molecule names	P-glycoprotein (mouse P-glycoprotein), POPC, cholesterol
Software names	Gromacs
Software versions	Gromacs 3.3.3
Forcefield / model names	GROMOS54a7 force field
Simulation temperatures	300 K
Simulation times (duration)	10 ns equilibration period; 200 ns production simulation (each system, run in triplicate)

Additional identifiers present (not requested but mentioned):

- Crystal structure / PDB IDs: 3G5U, 4KSB, 4M1M
- Pressure: 1 bar

Automation & model selection: OpenRouter

The Unified Interface For LLMs

Better prices, better uptime, no subscriptions.

[Get API Key](#) [Explore Models](#)

100T Monthly Tokens **8M+** Global Users **60+** Providers **400+** Models

One API for Any Model

Access all major models through a single, unified interface. OpenAI SDK works out of the box.

[Browse all](#)

Higher Availability

anthropic/claude-opus-4.8

Reliable AI models via our distributed infrastructure. Fall back to other providers when one goes down.

[Learn more](#)

Price and Performance

Throughput

Latency

Keep costs in check without sacrificing speed. OpenRouter runs at the edge for minimal latency between your users and their inference.

[Learn more](#)

Custom Data Policies

Protect your organization with fine grained data policies. Ensure prompts only go to the models and providers you trust.

[View docs](#)

<https://openrouter.ai/>

Automation & model selection: OpenRouter

The screenshot displays the OpenRouter AI model selection interface, showing several AI models with their specifications and pricing. The models are presented in a grid-like layout, with each model card containing a title, description, and a table of key metrics.

Model	Provider	Modalities	IN / OUT PRICE	Avg	CONTEXT	RELEASED
Anthropic: Claude Sonnet 4.6	anthropic/claude-sonnet-4.6	T, I, A, V, T	\$3 / \$15			
Qwen: Qwen3.6 27B	qwen/qwen3.6-27b	T, I, A, V, T	\$0,289 / \$2,40 per 1M		262K	Apr 27, 2026
Google: Gemini 3.1 Pro Preview	google/gemini-3.1-pro-preview					
Google: Gemma 4 31B	google/gemma-4-31b-it	T, I, A, V, T	\$0,12 / \$0,35 per 1M	Low	262K	
Anthropic: Claude Fable 5	anthropic/claude-fable-5	T, I, A, V, T	\$10 / \$50 per 1M	High	1M	Jun 9, 2026
DeepSeek: DeepSeek V4 Pro	deepseek/deepseek-v4-pro	T, I, A, V, T	\$0,435 / \$0,87 per 1M	Low	1M	

API

```
import requests
import json

# First API call with reasoning
response = requests.post(
    url="https://openrouter.ai/api/v1/chat/completions",
    headers={
        "Authorization": "Bearer <OPENROUTER_API_KEY>",
        "Content-Type": "application/json",
    },
    data=json.dumps({
        "model": "qwen/qwen3.7-plus",
        "messages": [
            {
                "role": "user",
                "content": "How many r's are in the word 'strawberry'?"
            }
        ],
        "reasoning": {"enabled": True}
    })
)

# Extract the assistant message with reasoning_details
response = response.json()
response = response['choices'][0]['message']
```


Automation, output format & quality control

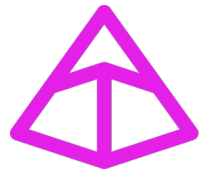
Data model  Pydantic

```
class Entity(BaseModel):
    category: str = Field(..., description="Category identifying the entity type.")
    text: str = Field(..., description="Extracted text content.")

class Molecule(Entity):
    category: Literal["MOL"] = Field(
        "MOL", description="Category for molecule entities."
    )
```

Automation, output format & quality control

Data model



Pydantic

Forces LLM to output
the data model



Instructor

Results (best of)

Model name	Release date	Parameters (B)	Precision	Unitary cost (\$)	Unitary inf. time (s)
anthropic/claude-sonnet-4.6	2026/02	?	0.89	0.002	0.5
z-ai/glm-5.1	2026/04	754	0.88	0.001	5.6
qwen/qwen3.6-27b	2026/04	27	0.87	0.002	8.0
google/gemini-3.1-pro-preview	2026/02	?	0.85	0.004	2.5
minimax/minimax-m2.7	2026/03	229	0.85	> 0.001	3.1
google/gemma-4-31b-it	2026/04	31	0.84	> 0.001	1.1
openai/gpt-5.5	2026/04	?	0.84	0.003	1.5

Results (worst of)

Model name	Release date	Parameters (B)	Precision	Unitary cost (\$)	Unitary inf. time (s)
deepseek/deepseek-chat	2024/12	685	0.74	< 0.001	2.2
mistralai/mistral-small-3.2-24b-instruct	2025/06	24	0.73	< 0.001	0.2
mistralai/mixtral-8x22b-instruct	2024/04	141	0.68	0.001	0.2
mistralai/ministral-8b-2512	2025/12	8	0.67	< 0.001	0.4
meta-llama/llama-3.1-8b-instruct	2024/07	8	0.56	< 0.001	1.0

Simulation time normalization



MD simulations of P-glycoprotein MOL started from three different crystal structures: 3G5U MOL , 4M1M MOL and 4KSB MOL . Molecular simulations of mouse P-glycoprotein MOL started from three different crystal structures with PDB IDs corresponding to 3G5U MOL , 4KSB MOL and 4M1M MOL . The protein was transferred to POPC MOL / cholesterol MOL bilayer, followed by energy minimization and 10 ns equilibration period. All the simulations were performed in Gromacs SOFTNAME 3.3.3 SOFTVERS in conjunction with GROMOS54a7 FFM force field at 300 K TEMP and pressure of 1 bar. Each system was run in triplicates for 200 ns STIME .

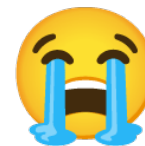
- 200ns
- 350 ns
- 50 ns
- 6μs

Regular expressions (Regex)

`([0-9]+) (\.\?[0-9]+)? *(ps|ns|μs|ms|s)`



- 5-microsecond
- 200-300ns
- one hundred nanosecond
- 1.5 micro-sec



Simulation time normalization



Model: DeepSeek V4 Pro

Prompt:

You are a unit normalization assistant for molecular dynamics simulation times.

Your tasks:

- Convert all time units to standard time abbreviations (ps, ns, μ s, ms, s)
- Separate numerical values from time units

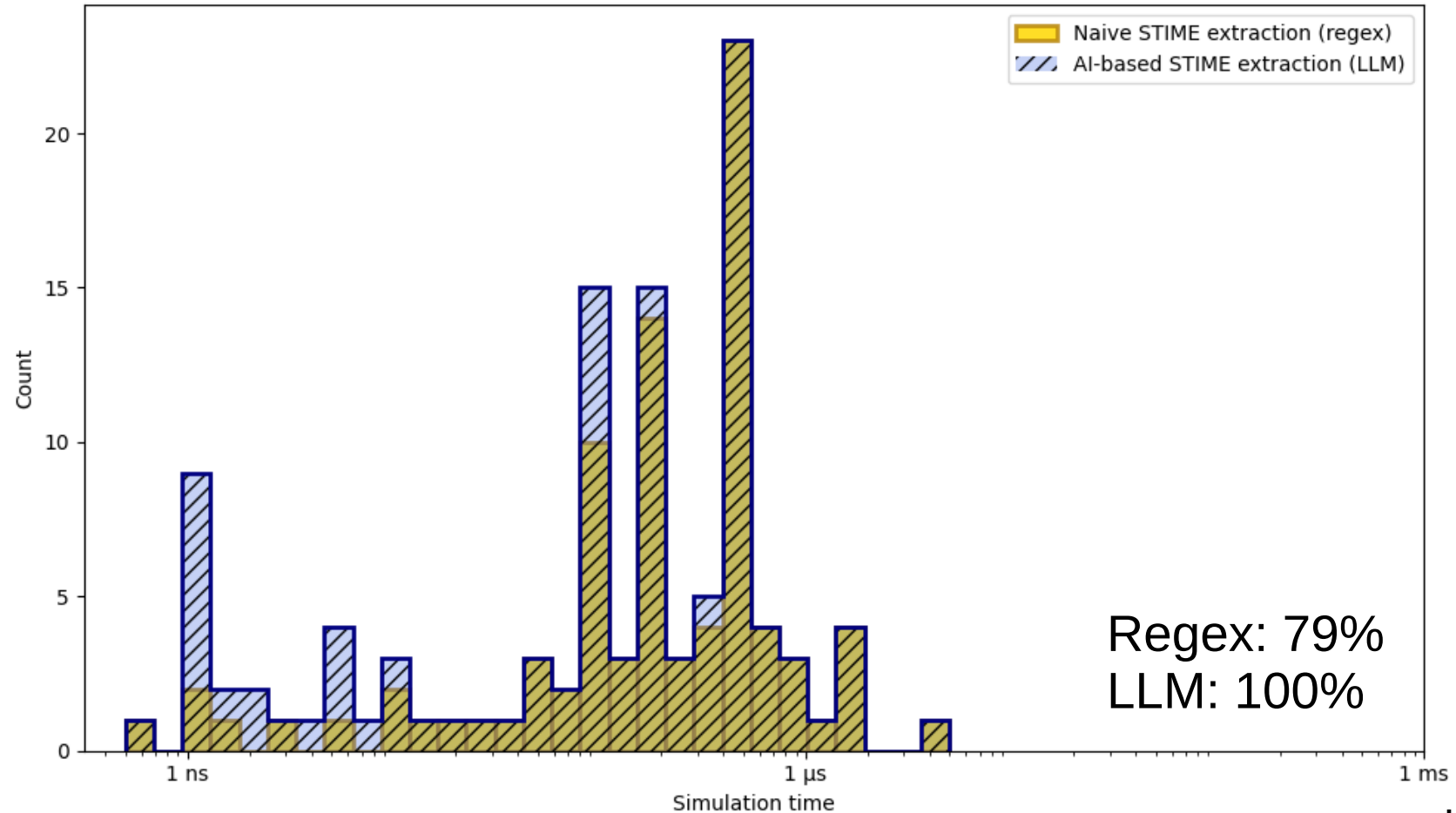
Rules:

- No markdown, no explanation
- Use only standard time units: ps (picoseconds), ns (nanoseconds), μ s (microseconds), ms (milliseconds), s (seconds)
- Always separate value and unit (e.g. "500ns" \rightarrow value: 500, unit: "ns")
- Take in consideration values written in letter (e.g. "one hundred"), and convert it to numeric value
- If the simulation time is an interval, separate each simulation time in the interval.
- If the unit is missing or the unit is not a time unit, output the normalized unit to "None"
- If the numerical value is missing, output the normalized value to "None"

Simulation time normalization



Distribution of simulation times



Conclusion

Metadata extraction and normalization are key for the reuse of open research data (and publications).

Large language models (e.g., ChatGPT) can be used as research tools, provided they are **equipped with sufficient harness**.

The field of LLMs is evolving rapidly:

- Most small models released last year or two-years ago perform significantly worse than current SOTA models.
- Access to large, up-to-date models (**proprietary and open-weight**) is important.

Discussion

Access to large, up-to-date models (proprietary and open-weight)

What I don't want to do:

- Read benchmarks (mostly biased and not corresponding to my use case)
- Find the “best” model
- Find how much VRAM (GPU) is needed to run the model (whichllm)
- Find how to run the model (vLLM)
- Connect to a computer cluster:
 - Wait to get access to GPU
 - Download an (open-weight) model
 - Wait...
 - Run inferences
- Repeat

Local 'computer'

```
(base) pierre@banana 📄 ~ $ uvx whichllm@latest plan "deepseek v4 pro"
```

Model Info

```
Model: deepseek-ai/DeepSeek-V4-Pro  
Params: 1600.0B (49.0B active) | Arch: deepseek | Context: unknown  
License: mit
```

VRAM Required (context: 4096)

Quant	VRAM	Quality Loss
Q2_K	473.0 GB	-25%
Q3_K_M	659.2 GB	-8%
Q4_K_M *	845.5 GB	-5%
Q5_K_M	1031.8 GB	-3%
Q6_K	1218.0 GB	-2%
Q8_0	1590.6 GB	-1%
F16	2987.5 GB	0%

LDLC PRO
GARANTIE 3 ANS INCLUSE MENU

NOS UNIVERS ▼ Chercher un produit, une marque, une caté 🔍

VOTRE PANIER :

DÉSIGNATION	PRIX	QUANTITÉ	SOUS-TOTAL
VOS PRODUITS			
 PNY NVIDIA H200 NVL + DE 15 JOURS	29 999€96	- 6 +	179 999€75

Vous avez un code promo ?
Entrez-le ici

TOTAL
179 999€75 HT
Voir le panier TTC

PASSER COMMANDE



API Hugging Face: open-weight only

 Inference Providers · Metrics for top trending models

[Browse all models](#)

[Learn more](#)

Filter by model or provider...

Model	Provider	Input \$/1M	Output \$/1M	Context	Latency(s)	Throughput(t/s)	Tools	Structured
 nvidia/NVIDIA-Nemotron-3-Ultra-550B-A55...	together	\$0.60	\$3.60	512,288	0.31	106	Yes	No
 deepseek-ai/DeepSeek-V4-Pro	novita 	\$1.60	\$3.38	1,048,576	0.40	37	Yes	No
 deepseek-ai/DeepSeek-V4-Pro	together	\$1.74	\$3.48	512,000	0.60	50	Yes	Yes
 deepseek-ai/DeepSeek-V4-Pro	fireworks-ai 	-	-	1,048,576	1.05	66	Yes	No
 deepseek-ai/DeepSeek-V4-Pro	deepinfra	\$1.74	\$3.48	65,536	1.43	29	Yes	Yes
 Qwen/Qwen3.6-35B-A3B	deepinfra	\$0.15	\$0.95	262,144	0.33	131	Yes	Yes
 google/gemma-4-31B-it	novita 	\$0.14	\$0.40	262,144	0.65	79	Yes	Yes
 google/gemma-4-31B-it	together	\$0.39	\$0.97	262,144	0.29	41	Yes	Yes
 google/gemma-4-31B-it	deepinfra 	\$0.13	\$0.38	262,144	0.39	21	Yes	Yes
 deepseek-ai/DeepSeek-V4-Flash	novita  	\$0.14	\$0.28	1,048,576	0.75	98	Yes	No
 deepseek-ai/DeepSeek-V4-Flash	fireworks-ai	-	-	-	1.00	78	Yes	No
 deepseek-ai/DeepSeek-V4-Flash	deepinfra	\$0.14	\$0.28	1,048,576	1.09	25	Yes	Yes
 Qwen/Qwen3.6-27B	ovhcloud	-	-	-	0.36	73	Yes	Yes
 google/gemma-4-26B-A4B-it	novita	\$0.13	\$0.40	262,144	0.66	31	No	Yes
 google/gemma-4-26B-A4B-it	deepinfra  	\$0.07	\$0.34	262,144	0.50	39	Yes	Yes
 moonshotai/Kimi-K2.6	novita 	\$0.80	\$3.40	262,144	1.17	25	Yes	No
 moonshotai/Kimi-K2.6	together	\$1.20	\$4.50	262,144	0.39	51	Yes	Yes
 moonshotai/Kimi-K2.6	fireworks-ai 	-	-	262,144	0.35	99	No	No
 moonshotai/Kimi-K2.6	deepinfra	\$0.75	\$3.50	262,144	0.64	18	Yes	No

<https://huggingface.co/inference/models>

API Albert

Des modèles pour tous vos usages

Albert API offre l'accès à une large gamme de modèle fondation open source d'IA générative. **Nous mettons à jour régulièrement les modèles disponibles pour vous permettre d'utiliser les modèles qui font l'état de l'art.**

D'où proviennent les modèles disponibles ?

Ces modèles ne sont pas propres à l'administration, il s'agit de modèle mis à disposition par des tiers comme Mistral ou Meta sur la plateforme 🤖 [HuggingFace](#) 🌐. Hébergé sur nos serveurs, aucune de vos données n'est envoyée à ces fournisseurs de modèles.

120B, 08/2025



APACHE 2.0 OPENWEIGHT-LARGE

[openai/gpt-oss-120b](#)

Modèle de chat pour des tâches complexes





24B, 06/2025



APACHE 2.0 OPENWEIGHT-MEDIUM

[mistralai/Mistral-Small-3.2-24B-Instruct-2506](#)

Modèle de chat pour des tâches d'une complexité modérée et d'analyse d'images





8B, 12/2025



APACHE 2.0 OPENWEIGHT-SMALL

[mistralai/Mistral-3-8B-Instruct-2512](#)

Modèle de chat pour des tâches simples





<https://ia.numerique.gouv.fr/outils-ia/albert-api/>

API Scaleway



Generative APIs - Serverless

Use the latest AI models via API, pay by thousand tokens.

Try out new models with our [free tier](#): 1 million tokens and 60 minutes of audio transcription.

All requests performed using [Batches API](#) are priced with a -50% discount.

Select Region

 Paris ▼

Name ↕	Tasks ↕	Input tokens ↕	Output tokens ↕	
qwen3.5-397b-a17b	Chat and code	€0.60 / MILLION TOKENS	€3.60 / MILLION TOKENS	Try
qwen3.6-35b-a3b	Chat and Vision	€0.25 / MILLION TOKENS	€1.50 / MILLION TOKENS	Try
gemma-4-26b-a4b-it	Chat and Vision	€0.25 / MILLION TOKENS	€0.50 / MILLION TOKENS	Try
gpt-oss-120b	Chat	€0.15 / MILLION TOKENS	€0.60 / MILLION TOKENS	Try
mistral-medium-3.5-128b	Chat and Vision	€1.50 / MILLION TOKENS	€7.50 / MILLION TOKENS	Try
whisper-large-v3	Audio transcription	€0.003 / AUDIO MINUTE	Free	Try
llama-3.3-70b-instruct	Chat	€0.90 / MILLION TOKENS	€0.90 / MILLION TOKENS	Try
qwen3-235b-a22b-instruct-2507	Chat	€0.75 / MILLION TOKENS	€2.25 / MILLION TOKENS	Try

<https://www.scaleway.com/en/pricing/model-as-a-service/>

Conclusion

Metadata extraction and normalization are key for the reuse of open research data (and publications).

Large language models (e.g., ChatGPT) can be used as research tools, provided they are **equipped with sufficient harness**.

The field of LLMs is evolving rapidly:

- Most small models released last year or two-years ago perform significantly worse than current SOTA models.
- Access to large, up-to-date models (**proprietary and open-weight**) is important.

Thanks



IJM, Paris, France

Lisa Bouarroudj
Mohamed Oussaren

IPBS, Toulouse, France

Magdalena Szczuka
Matthieu Chavent

LBT, Paris, France

Marc Baaden
Karine Duong
Giulia Di Gennaro
Benoist Laurent
Essmay Touami
Inès Zenati
Salahudin Sheikh

INRAE, Jouy-en-Josas, France

Arnaud Ferré

Amsterdam, Netherlands

Steven Garcia

Univ. Copenhagen, Denmark

Johanna K. S. Tiemann
Kresten Lindorff-Larsen

KTH Royal Inst. Tech. Stockholm, Sweden

Lucie Delemotte
Erik Lindahl

Stockholm Univ. , Sweden

Rebecca J. Howard



DET FRANSKE
KULTURCENTER
I DANMARK

