

# LLM Infrastructure on HPC: Workflows, Constraints, and Solutions (Text Lab)

Ahmad Alhineidi  
Data Science Lab (DSL)

29.06.2026

PASC 26

MS2C – Serving Inference: Leveraging  
HPCs in the Age of Generative AI

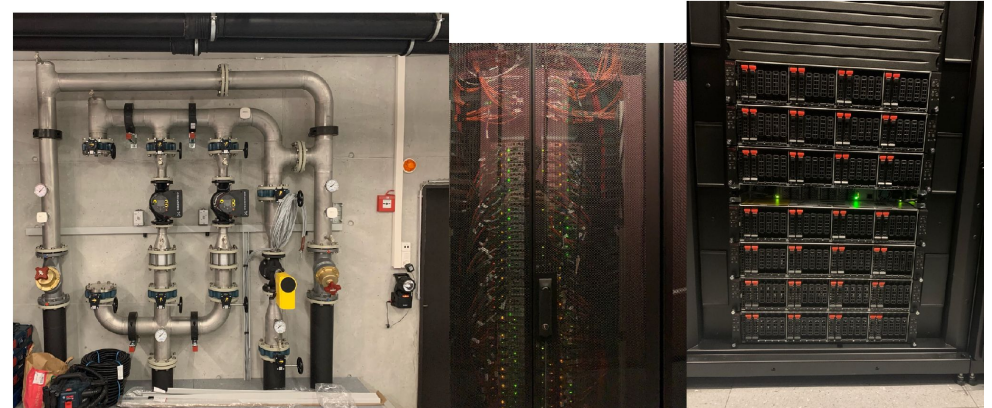


# Content

- University of Bern HPC setup
- Text Lab overview & usage (Live demo)
- GPUs and LLMs choice

# UniBe HPC - UBELIX

- GPU partition
- 21 GPU nodes · 174 GPUs total
- RTX3090, RTX 4090, A100, H100, H200, RTX PRO 6000 Blackwell
- Most GPU nodes: 128 CPU cores and ~740 GB RAM (rtx4090/h100/h200)



# UniBe HPC - UBELIX

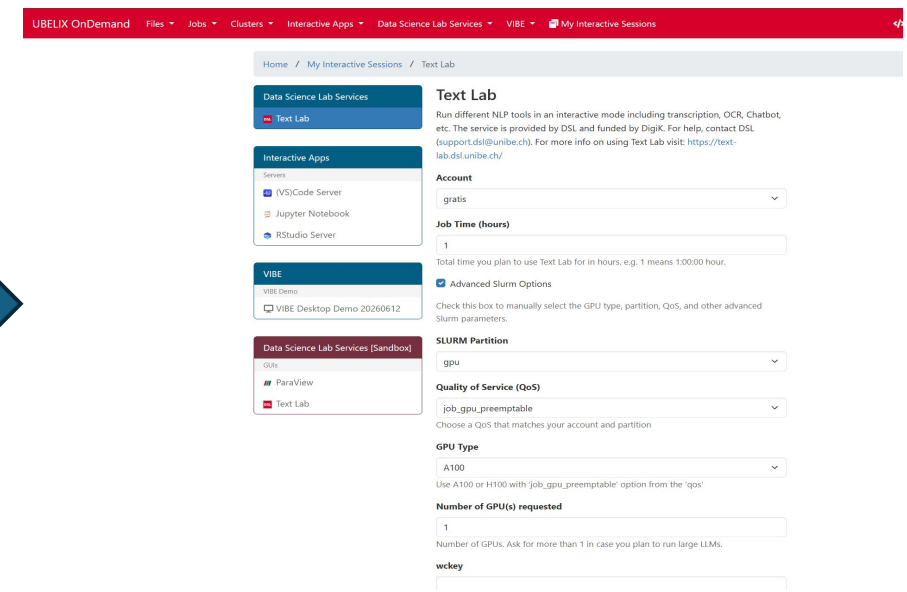
- Open OnDemand deployed since 2025
- Improves HPC accessibility, useability
- Allows to easily build custom interactive apps

```
Rocky 9.7 Blue Onyx
-----
FQDN:      submit02.ubelix.unibe.ch (10.4.129.22, 10.1.129.22)
Processor: 128x AMD EPYC 7742 64-Core Processor
Kernel:    5.14.0-611.55.1.el9_7.x86_64
Memory:    128.229 GB
-----
          **
         ****
        *****
       *****
      *****
     *****
    *****
   *****
  *****
 *****
*****
*****

Welcome to the UBELIX HPC
Please report any issues at:
https://serviceportal.unibe.ch

https://ubelix.hpc.unibe.ch - for more information

-----
Important changes to the UBELIX usage model in December 2025:
The University of Bern is introducing a pricing scheme for UBELIX.
Please refer to https://hpc-unibe.ch/github.io/costs/overview/ for details.
```



# Text Lab - a bit of History

- Whisper container on UBELIX for audio transcription
- Very long Manuel
- Includes command-line interface
- 😞 confusion, more support, where is my data? How do I transfer my data to the cluster? Where is the output?
- Why not just make a User interface?  Text Lab

b. If you do not see your audio file listed you can retrieve the path to it as such:

```
find $HOME -name "<audio_file.mp3>"
```

(Replace <audio\_file.mp3> with the name of your audio file)

The complete command should look as such:

```
/storage/research/dsl_shared/solutions/whisper/scripts/transcribe.sh  
/storage/homefs/<your_user_name>/<path/to/your/audio_file.mp3>
```

(Replace <your\_user\_name> with your UBELIX user name, e.g.: ab12c345; and replace the <path/to/your/audio\_file.mp3> with the result from the 'find' command)

Press *Enter* to send your command

3. If the script has launched your Whisper job successfully it will report a job number. Example:

```
Submitted batch job <1234567>
```

Check the status of your job on the cluster:

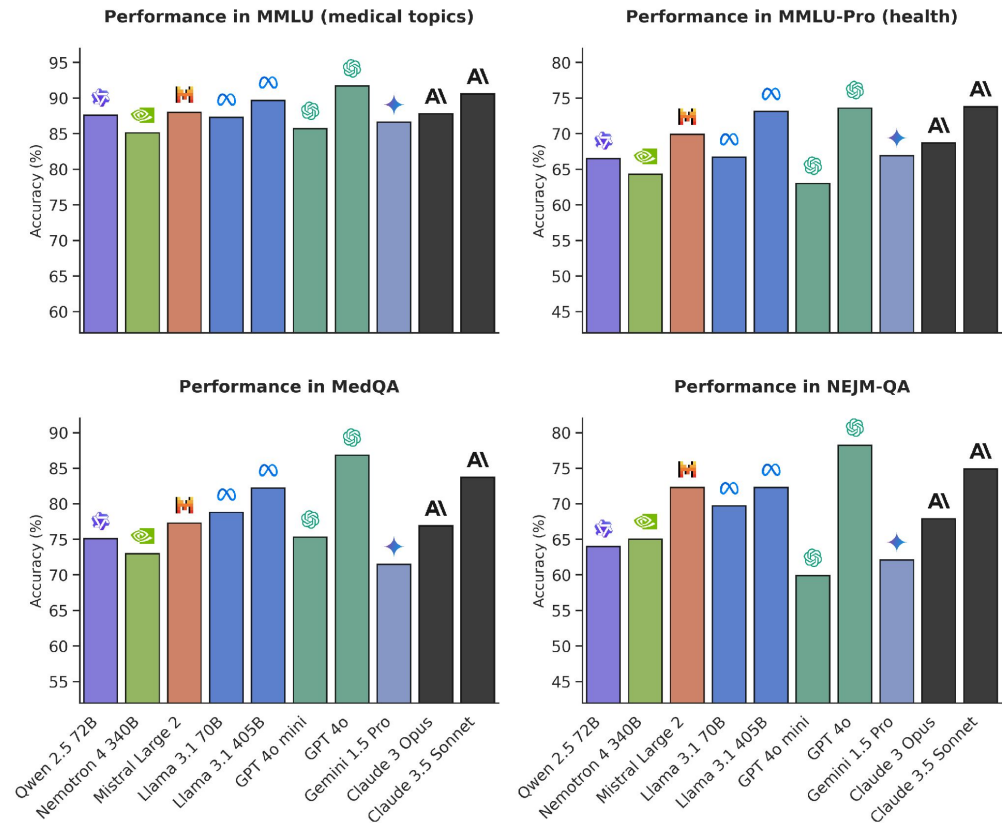
```
squeue --me
```

The output looks like this:

| JOBID   | PARTITION | NAME    | USER     | ST | TIME | NODES | NODELIST(REASON) |
|---------|-----------|---------|----------|----|------|-------|------------------|
| 1234567 | gpu       | whisper | ld23y671 | R  | 0:07 | 1     | gnode18          |

# Open-source LLMs

- Managed to match/come close to performance of closed models
- Can download and run locally
- Still requires expensive hardware to serve efficiently
- Quantization and smaller models are cheaper and can perform well in certain tasks



# Text Lab - motivation

- Make different “SOTA” models available & easily accessible == zero code
- Data privacy (Runs locally, University infrastructure)
- Experiment with different models / HPC deployment

# What is Text Lab?

- Framework with different NLP/Text Processing features
- Runs within the University internal network and infrastructure

## What can you do with Text Lab?

- Chat, transcription, OCR, Data visualization, knowledge graphs
- Live Demo
- User's guide: <https://text-lab.dsl.unibe.ch/>

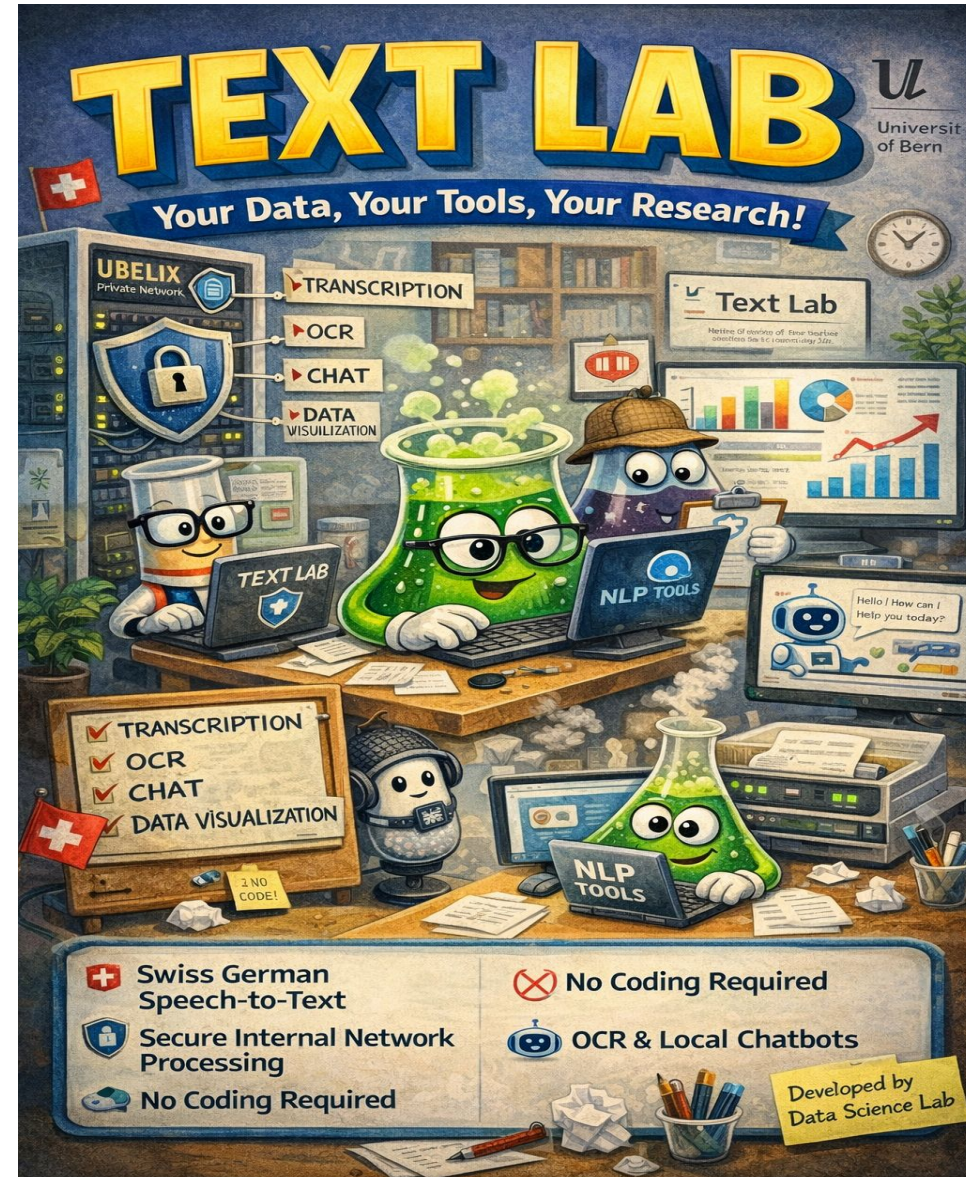
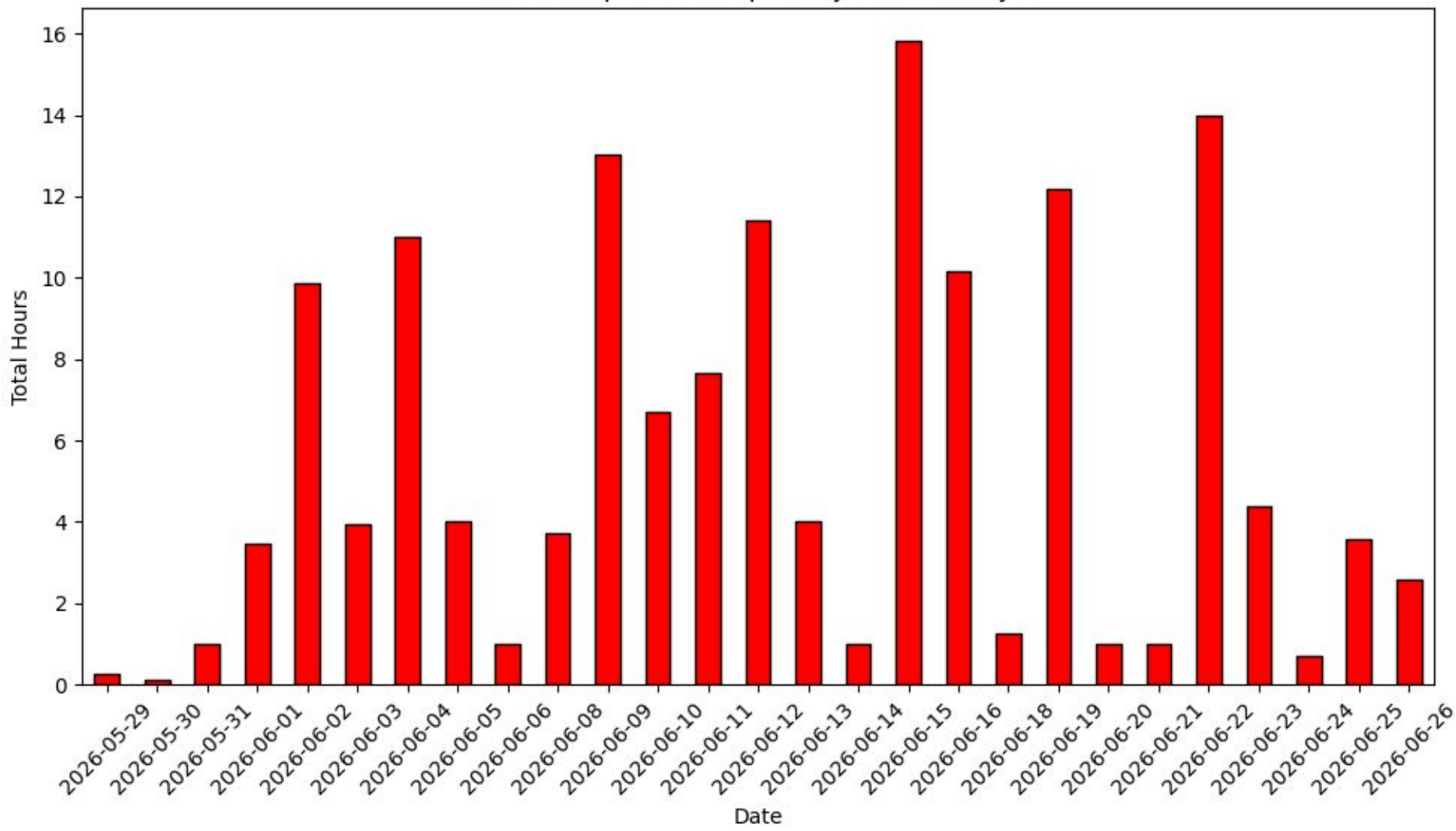


Image source: AI generated, openai

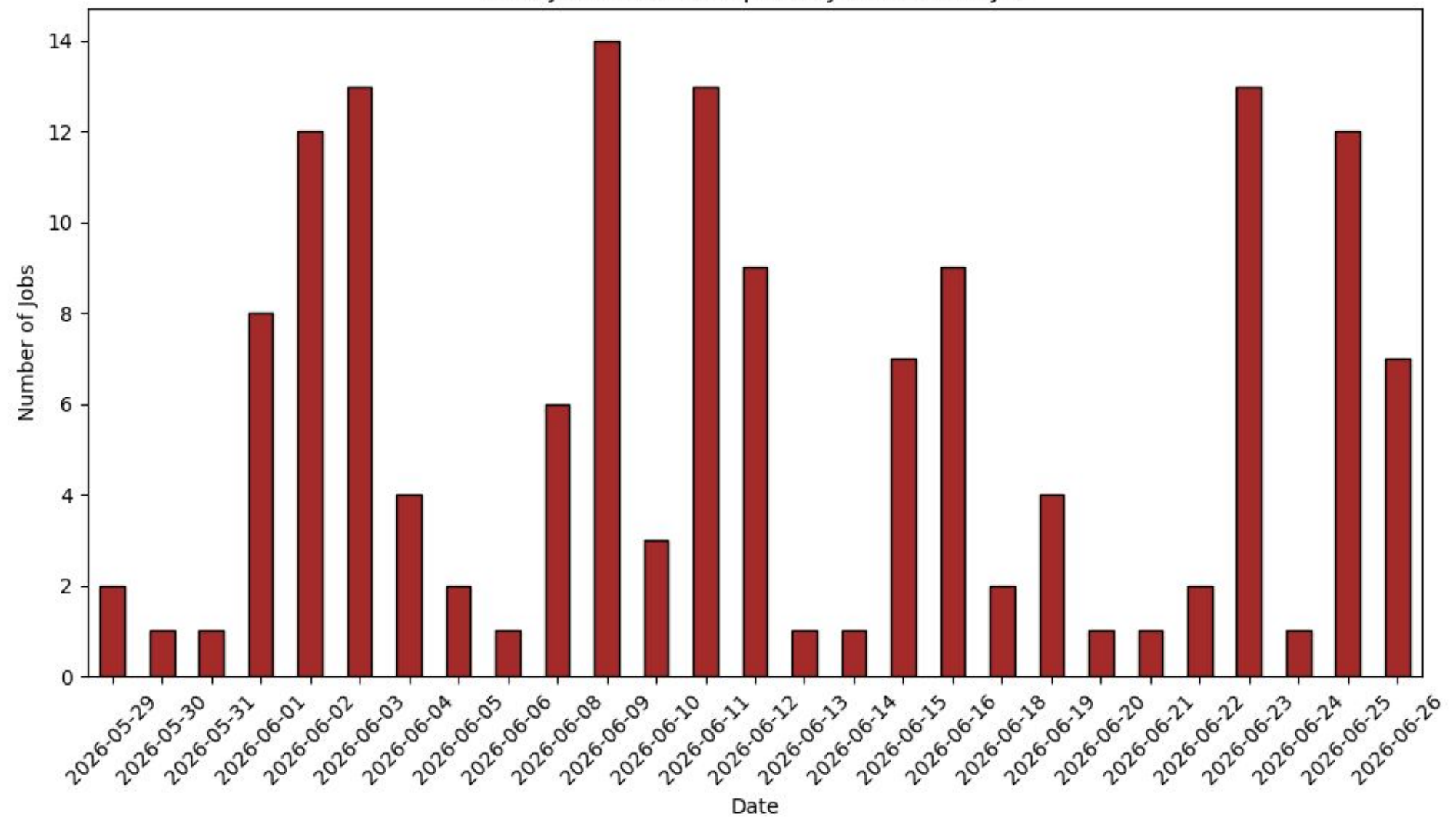
# Usage

- Anyone really using Text Lab?
- We're moving from prototype to useful, functional tools
- Next step to know suitable users  communicate the tool

Total Compute Hours per Day (Last 30 Days)



Total Jobs Submitted per Day (Last 30 Days)

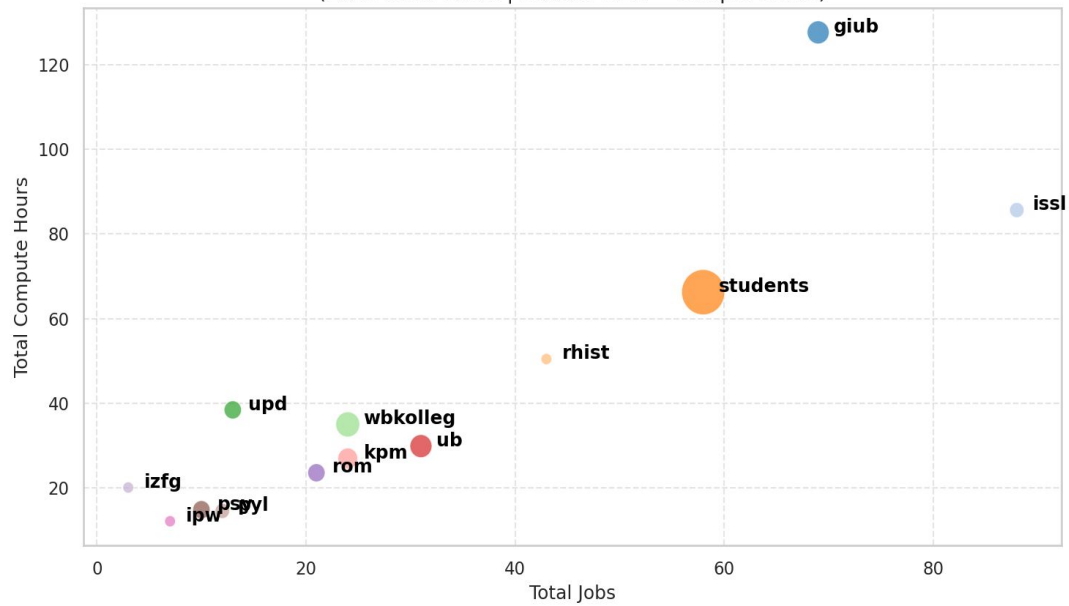


# Usage - last 30 days

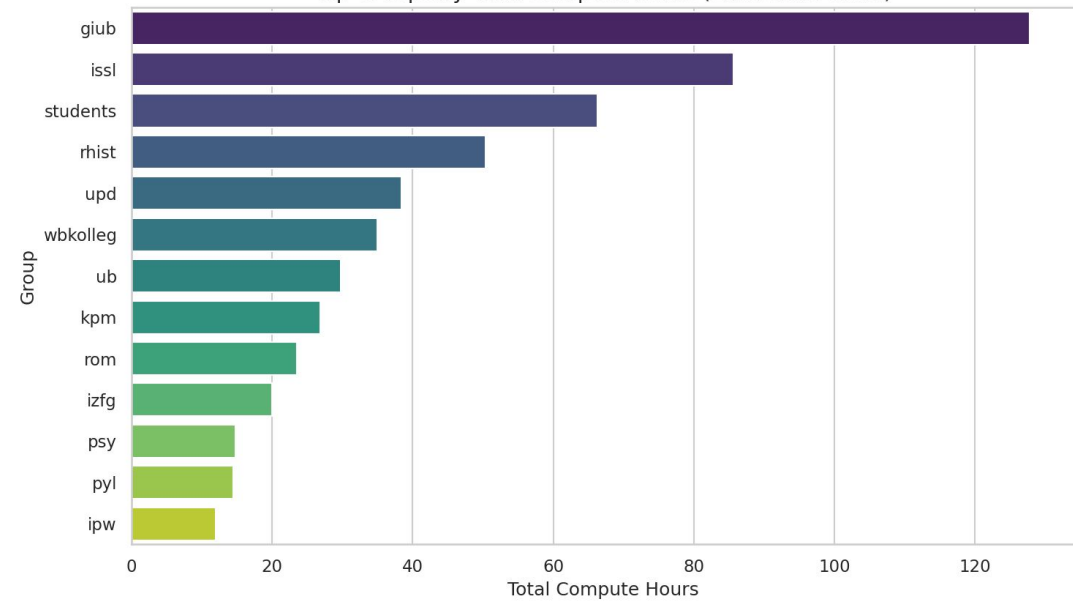
- Average Compute Hours per Day: 4.96
- Average Jobs per Day: 5.63
- Total Jobs Run: 169 (2026-05-29 to 2026-06-26)
- Total Unique Users: 47
- Total Unique Groups: 24

# Usage

Group Profiles: Total Jobs vs. Total Hours  
(Excl. math & dsl | Bubble Size = Unique Users)

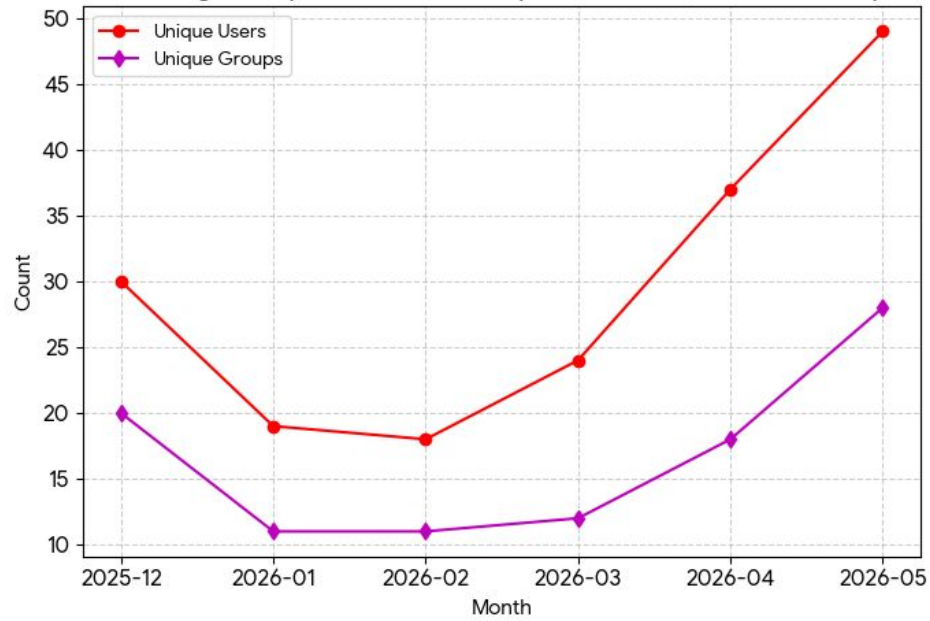


Top Groups by Total Compute Hours (Excl. math & dsl)

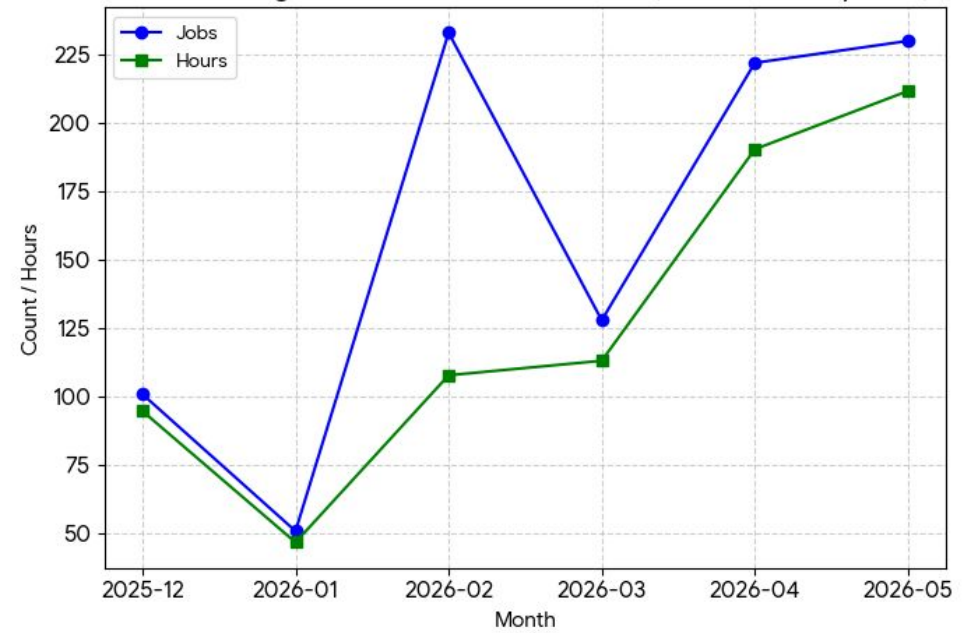


# Usage

Text Lab Usage: Unique Users and Groups Over Time (Dec 2025 - May 2026)

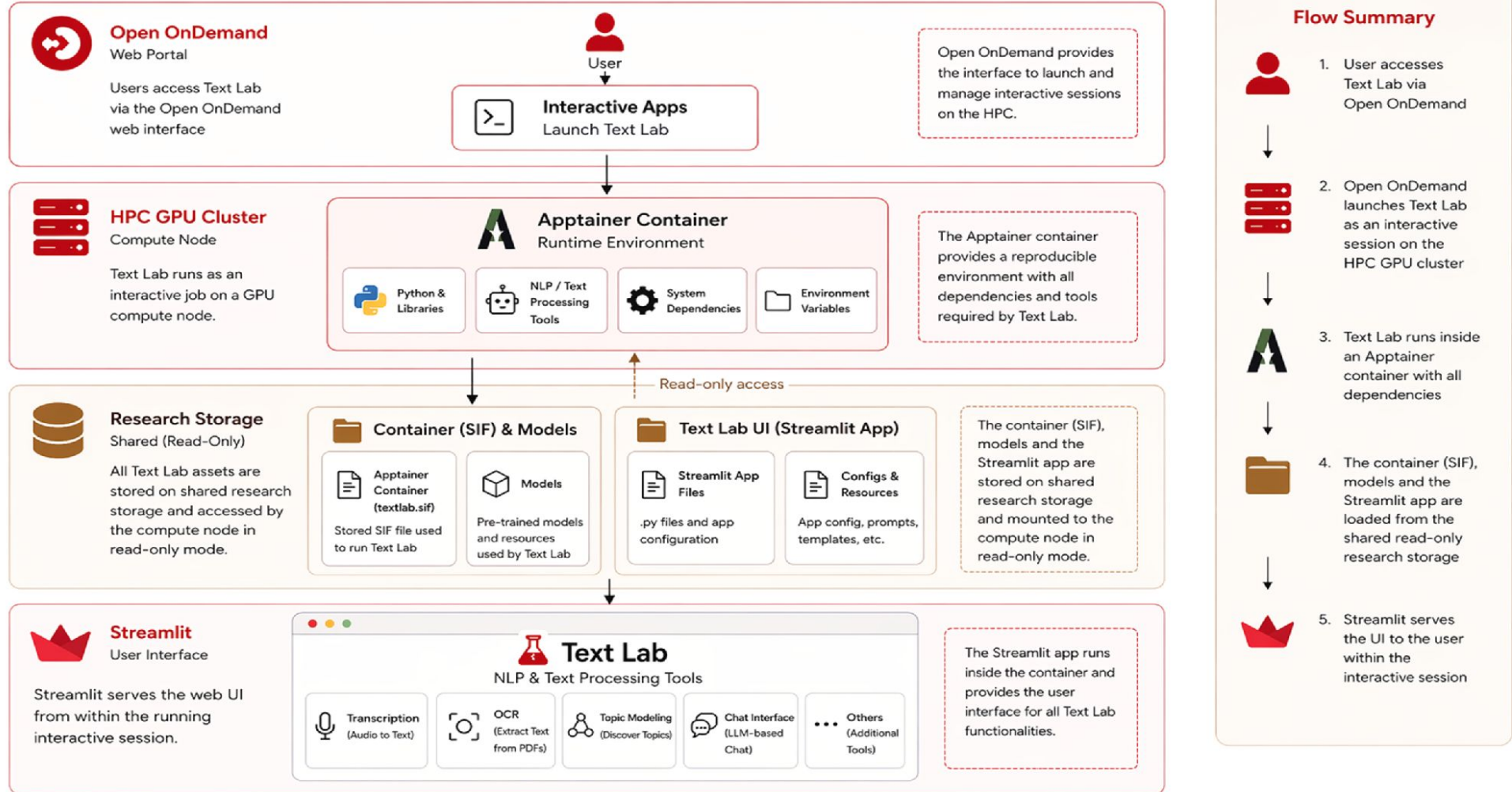


Text Lab Usage: Jobs and Hours Over Time (Dec 2025 - May 2026)



## Text Lab Architecture

Text Lab runs as an Open OnDemand Interactive app on the HPC GPU cluster



# LLMs Choice

- Open-source models range from less than 1 B □ more than 600B (~404GB model weights)
- **Hardware requirement example 1:** Qwen 3.5 122B (81GB), **4** RTX 4090, OR **2** A100, OR **2** H100
- **Hardware requirement example 2:** Qwen-coder 480B (290 GB) **13** RTX4090 !?, OR **4** A100, OR **4** H100
- Choice between acceptable performance and Hardware restrains
- Model used between ~9 – 30B parameters. Can run on single GPU
- Perform reasonable with MCP, given good environment
- Less problematic with other transformer models for transcription (whisper large ~ 3GB), OCR, etc

# Local LLMs (The Storage Problem)

- Initial approach: each user download model on ~HOME in first run
- **Models are huge.** A single modern open-weight LLM is 5–40 GB
- Many users, one cluster. If every researcher downloads their own copy into their home directory, storage cost multiplies by the number of users —  $N$  users  $\times$   $M$  models =  $N \times M$  copies of the same bytes
- **Models update frequently** (Gemma1, 2, 3, 4, etc). Regular updates on Text Lab
- If Text Lab becomes popular  $\square$  storage problem

# Local LLMs (The Storage Problem)

- The intuitive first approach was to embed the models inside the Apptainer SIF image itself. One file to rule them all!
- First problem, build time is very long for every update.
- Second problem, introduce latency loading the model

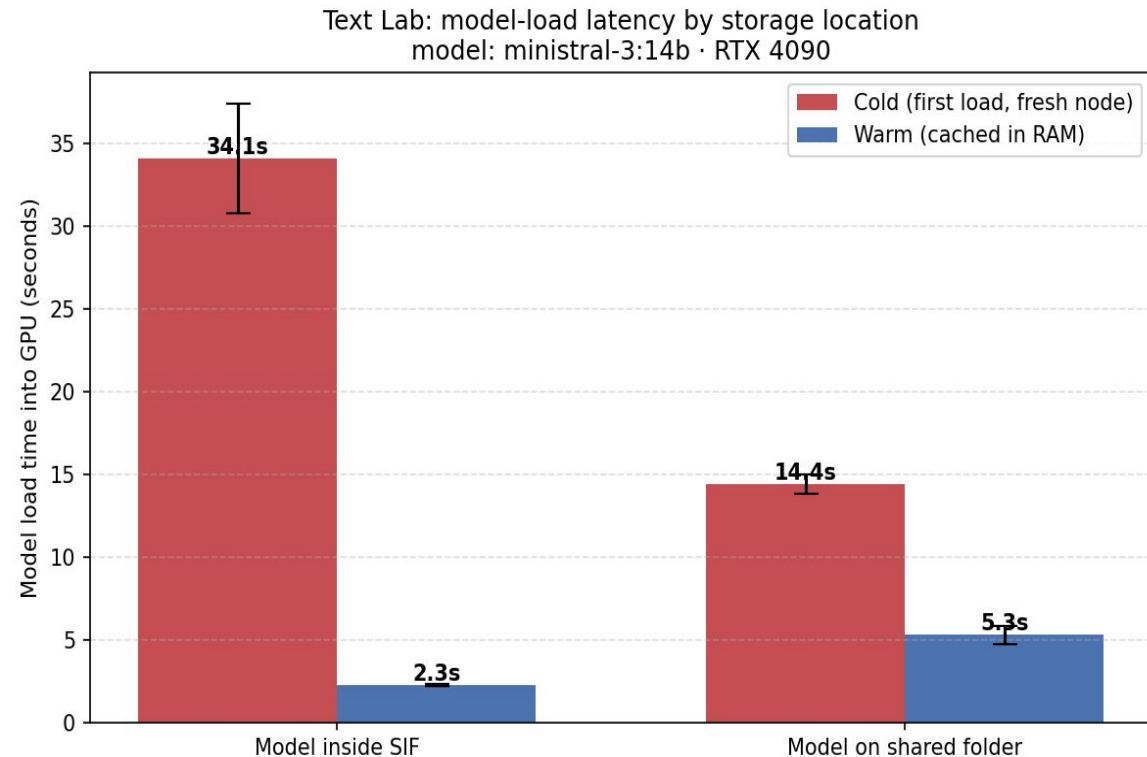
# Local LLMs (The Storage Problem)

- Current solution: One Shared, Read-Only Cache
- Store each model once, in a central read-only folder on cluster storage, shared by all Text Lab users.
- **Read-only = safe & consistent.** No user can corrupt or accidentally delete the weights; everyone runs the exact same, validated model.
- **Zero footprint in user space.** Researchers consume *no* home-directory quota for model weights. Their quota stays for their own data.
- **De-duplication at scale.** Storage cost is M models (one copy each), not  $N \times M$ . As user count grows, the savings grow with it.

# LLM inside SIF vs. Shared Folder Benchmark

- Cold start from inside the SIF is  $\sim 2.4\times$  slower than from the shared folder.
- The larger the model, the longer the wait for the first load
- Ollama reports load\_duration (nanoseconds) in every /api/generate response
- Model: minstral-3:14b, ollama version 0.24.0, GPU: RTX 4090

| Storage location       | Cold load (mean±std) | Warm load (mean±std) | Cold runs |
|------------------------|----------------------|----------------------|-----------|
| Model inside SIF       | 34.1 ± 3.301 s       | 2.268 ± 0.084 s      | 3         |
| Model on shared folder | 14.392 ± 0.587 s     | 5.289 ± 0.567 s      | 3         |



## Local LLMs (The Storage Problem) – Storage saving

- Shared read-only cache: ~393 GB on disk, 14 models, stored exactly once. (Some older models clean-up)

| Text Lab users | Storage saved |
|----------------|---------------|
| 10             | ~ 3.8 TB      |
| 25             | ~9.6 TB       |
| 50             | ~19.2 TB      |
| 100            | ~38.4 TB      |
| 200            | ~76.8 TB      |

# Which GPU should I use?

- RTX4090, A100, H100, H200??
- The larger VRAM, the better?
- Certain GPUs are less available than others

# Text Lab - tokens per second GPU benchmark

- Measures how fast a local LLM generates text (**tokens per second**) on different UBELIX GPU type

# Text Lab - tokens per second GPU benchmark

- Experiment set up:

**Model:** minstral-3:14b

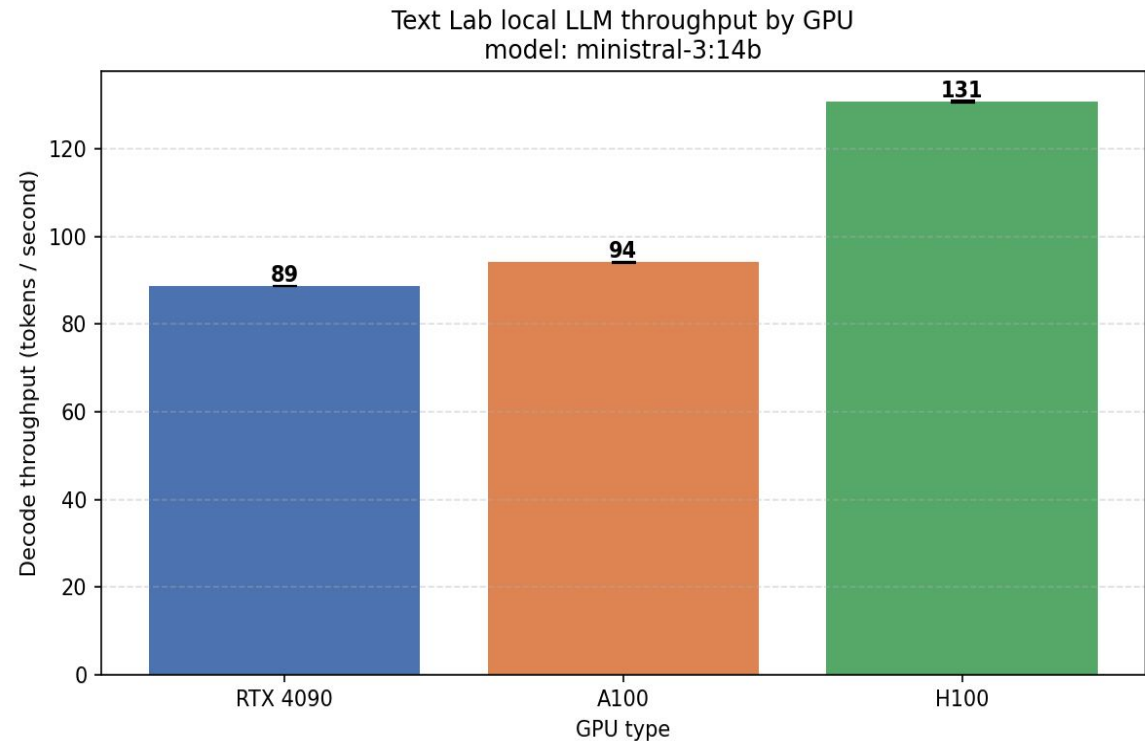
**Ollama version:** 0.24.0

**Flash attention:** true

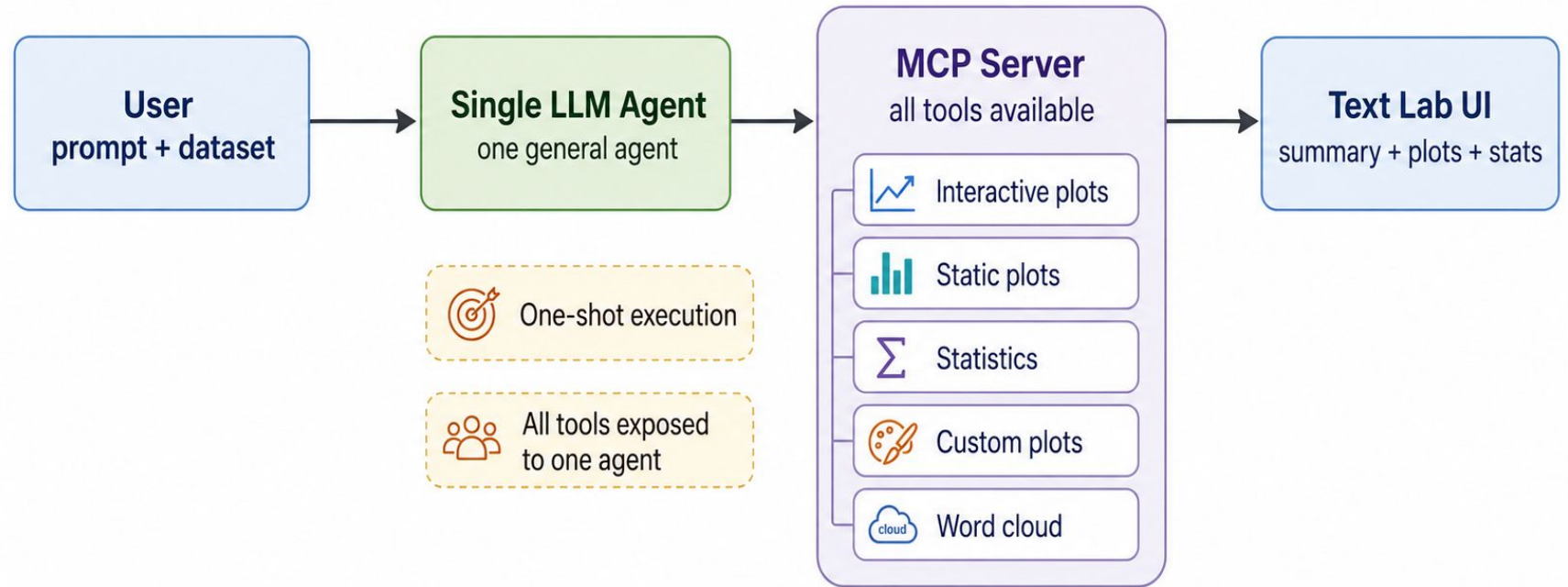
**Sampling options:** {"num\_predict": 256, "temperature": 0, "seed": 42, "top\_p": 1.0, "top\_k": 0, "num\_ctx": 4096}

**tokens\_per\_sec = eval\_count / eval\_duration \* 1e9**

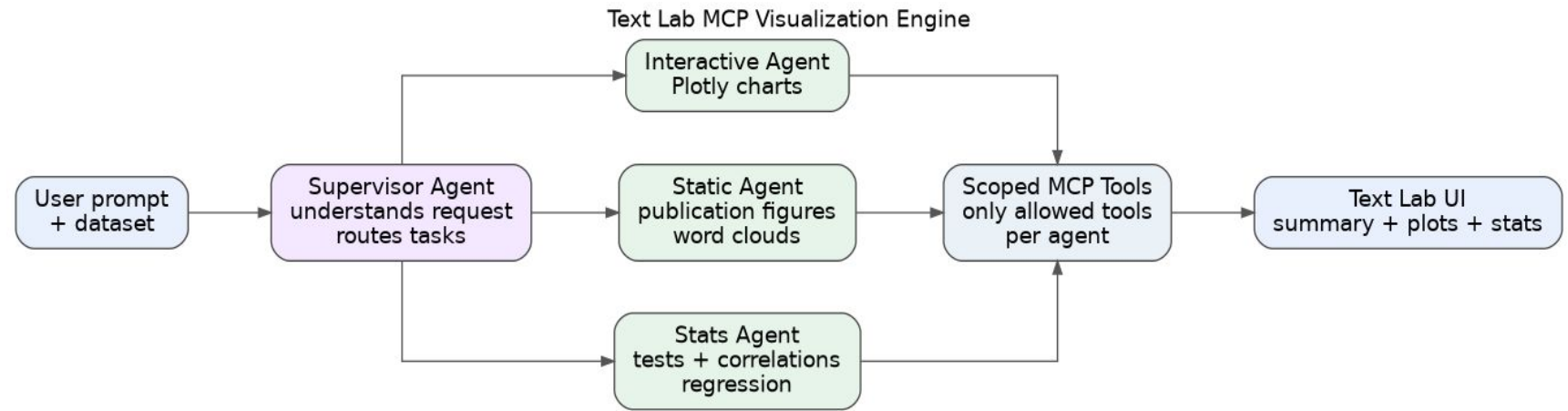
**eval\_count** = tokens generated, **eval\_duration** = nanoseconds spent generating them, and **1e9** converts nanoseconds to seconds.

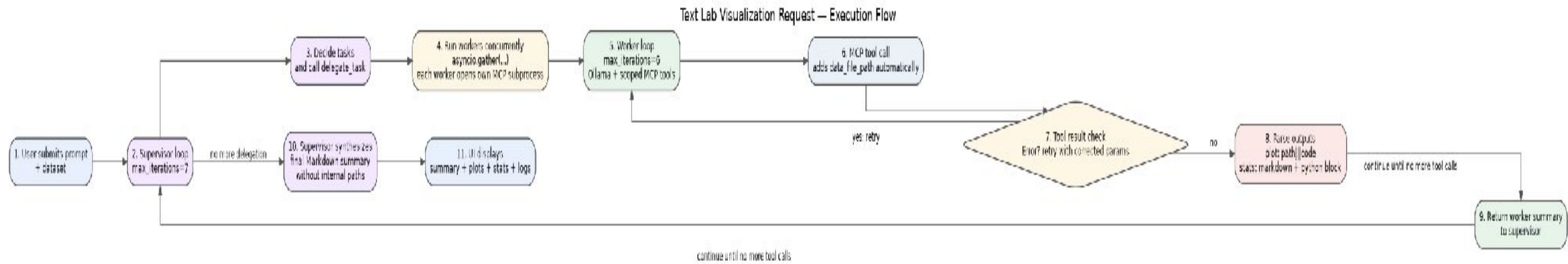


- Initial approach
- 1 agent (LLM)
- Access to all tools
- Failure when scaling



- Multi-agent architecture
- 1 supervisor agent, 3 specialized agents
- Each LLM has access to own tools only
- Integration when failed





# What's Next?

- Improve reach out
- More features (Machine Translation, Metadata extraction)
- Maybe communicate with other university with similar setup (HPC + Open OnDemand)

# Thank you!

- Questions
- Try Text Lab on: <https://ondemand.hpc.unibe.ch/>