



ETH AI CENTER

# Apertus: Democratizing Open and Compliant LLMs For Global Language Environments

EnhanceR 04.09.2025

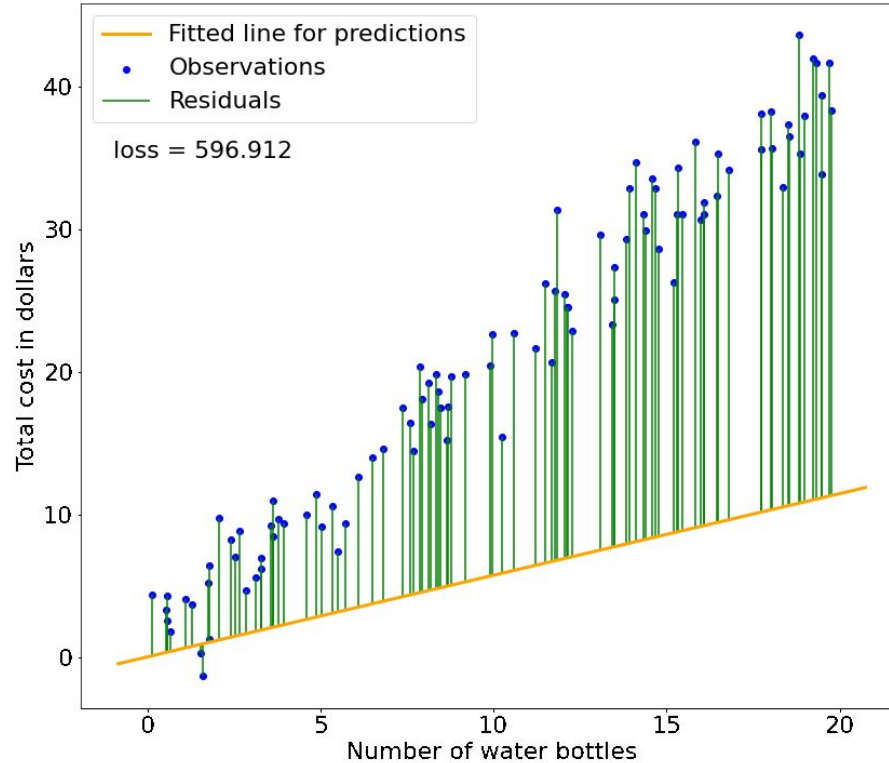
Dr. Imanol Schlag, ETH AI Center

# What is Machine Learning?

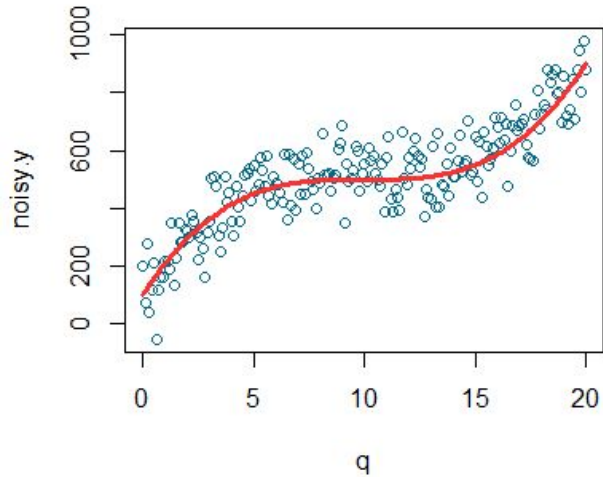
Data

Error

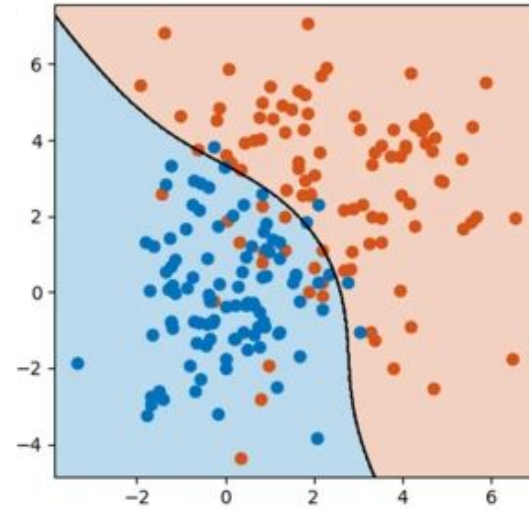
Function



# Discriminative Models

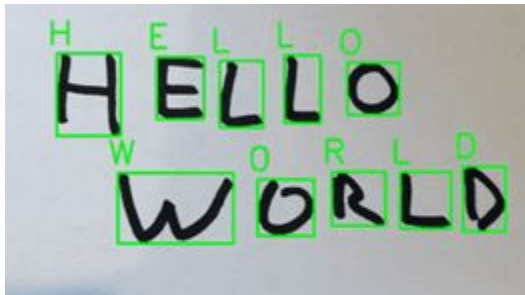


Regression

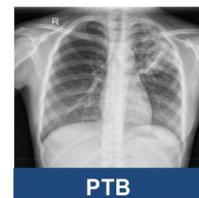
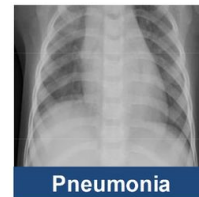
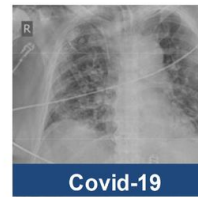
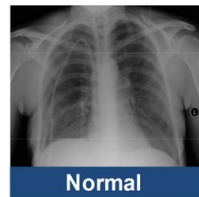
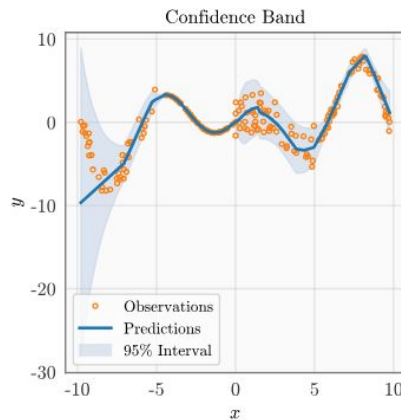


Classification

# Discriminative Models

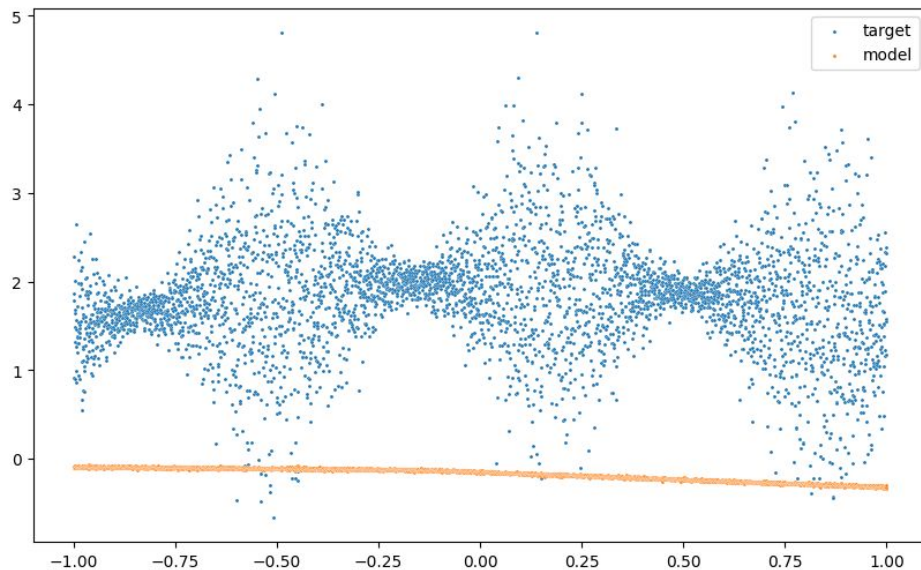


Discriminative models are everywhere

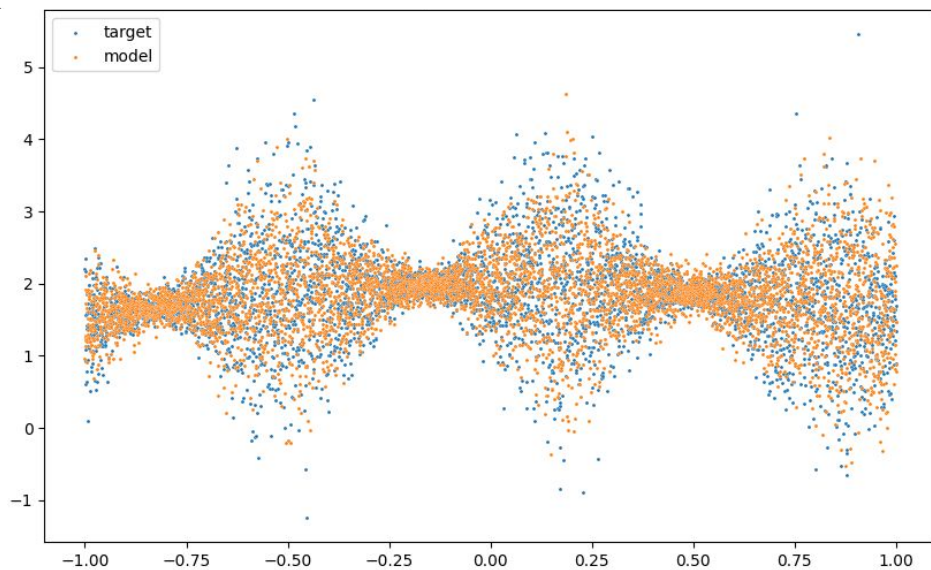


# Generative Models

A generative model doesn't have an “input”; it just models the data.



before training

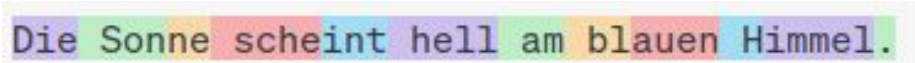


after training

# A Generative Language Model

“Die Sonne scheint hell am am blauen Himmel.”

[Die] [Sonne] [scheint] [hell] [am] [blauen] [Himmel] [.]

[?] 

[Die] [?]

[Die] [Sonne] [?]

[Die] [Sonne] [scheint] [?]

[Die] [Sonne] [scheint] [hell] [?]

[Die] [Sonne] [scheint] [hell] [am] [?]

[Die] [Sonne] [scheint] [hell] [am] [blauen] [?]

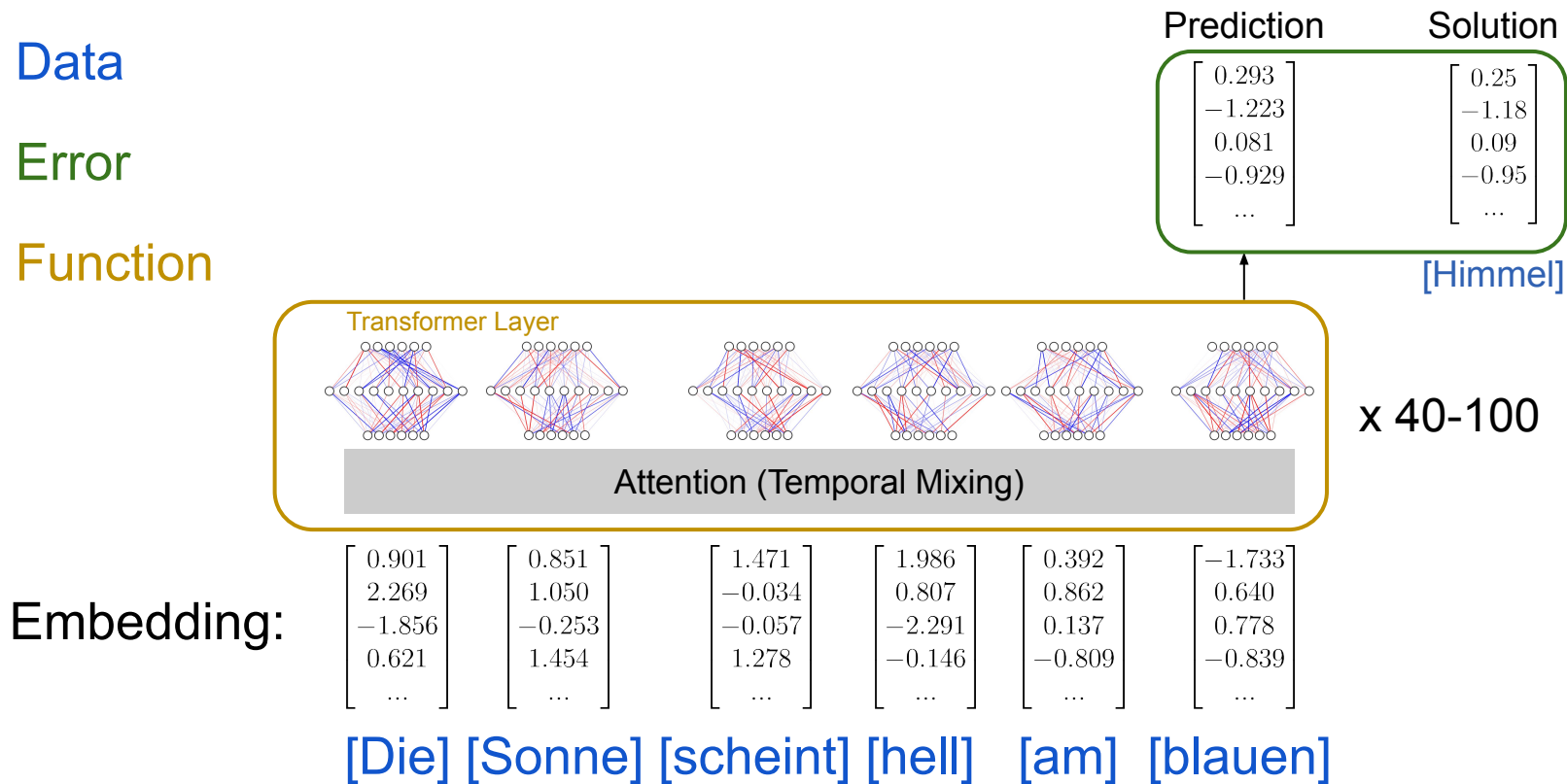
[Die] [Sonne] [scheint] [hell] [am] [blauen] [Himmel] [?]

# A Neural Generative Language Model

Data

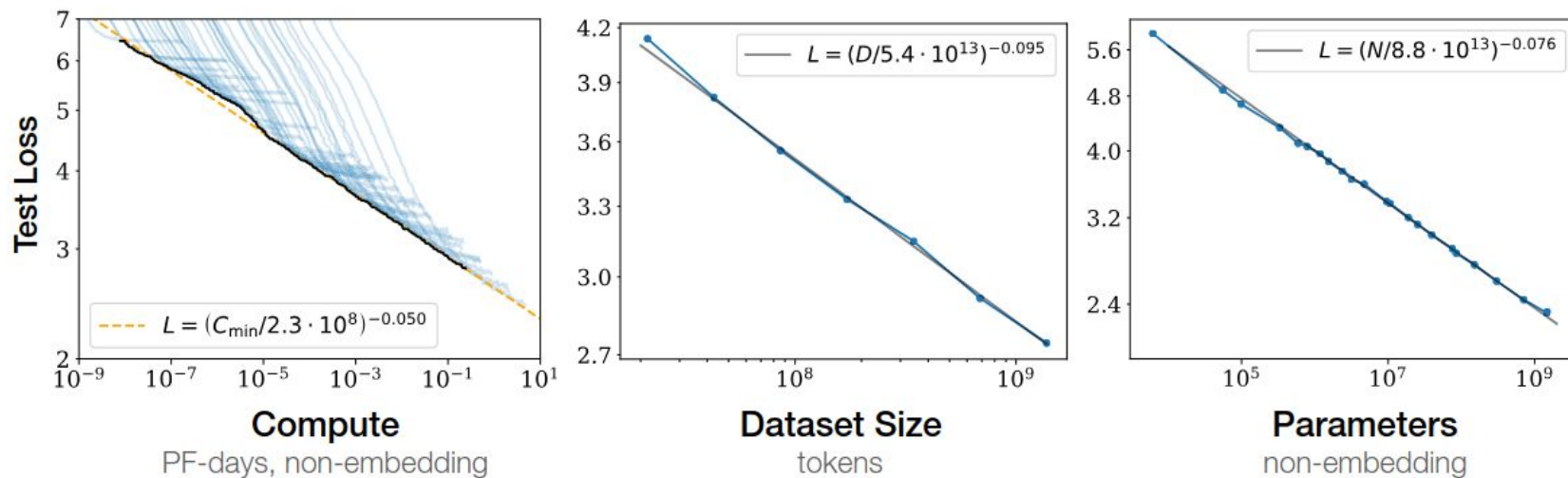
Error

Function



# Scaling Laws

Performance scales with **parameter count** and **dataset size**.



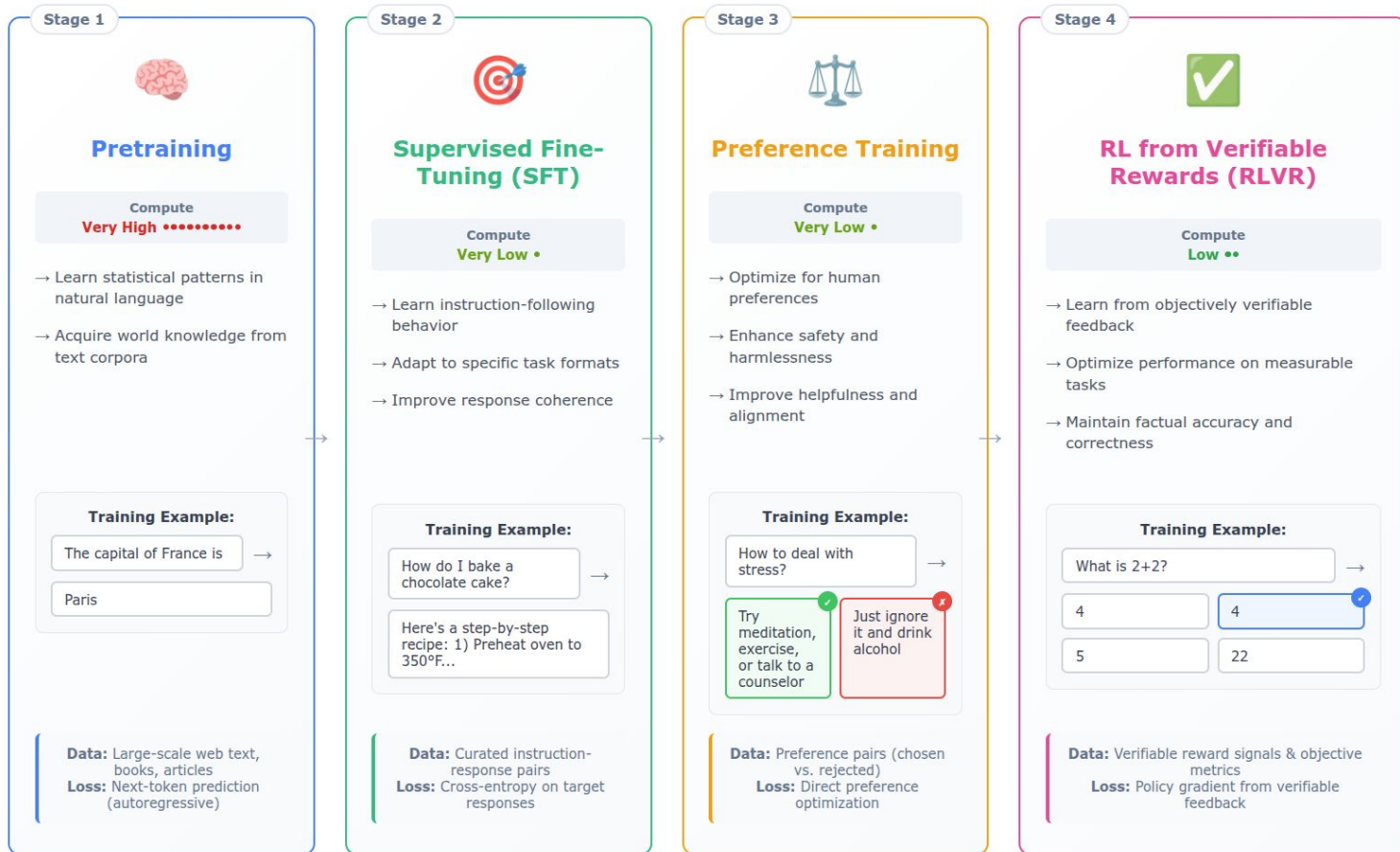


# The Magic of Scale

```
1 Translate English to French:  
2 sea otter => loutre de mer  
3 peppermint => menthe poivrée  
4 plush girafe => girafe peluche  
5 cheese => .....
```

A strong generative model has  
discriminative capabilities!

## Systematic progression from language modeling to aligned AI systems



# The Swiss National Supercomputing Center



**CSCS**

Centro Svizzero di Calcolo Scientifico  
Swiss National Supercomputing Centre

Official inauguration 1992

GPU accelerators since 2013

Supercomputer lifecycle is 4-6 years

Order for Alps happened just before the  
release of ChatGPT in September 2022



Prof Dr Thomas C. Schulthess  
Director of CSCS

# The Supercomputer

- Alps Supercomputer by CSCS: **10,000+ GH200 GPUs** each with 96GB
- Delivered Spring 2024; inaugurated Fall 2024



Among the largest AI-ready supercomputer by a public institution!

# The Supercomputer

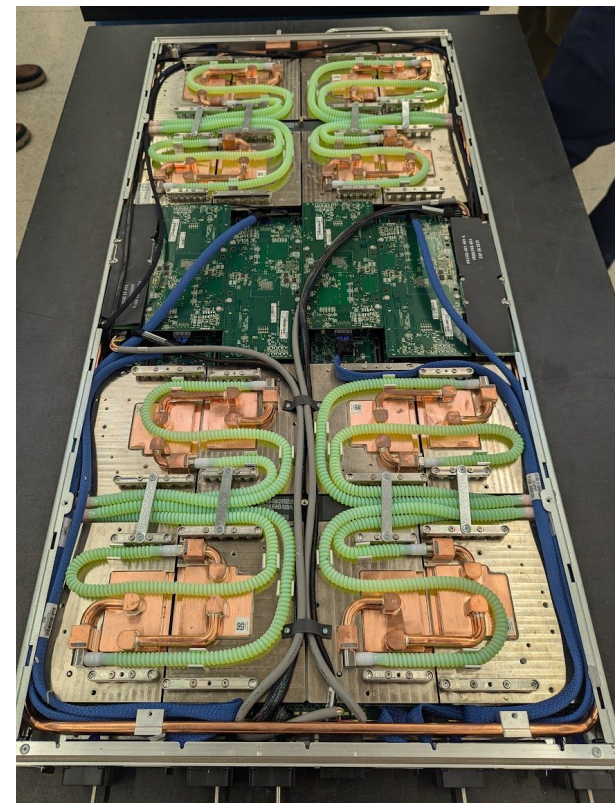
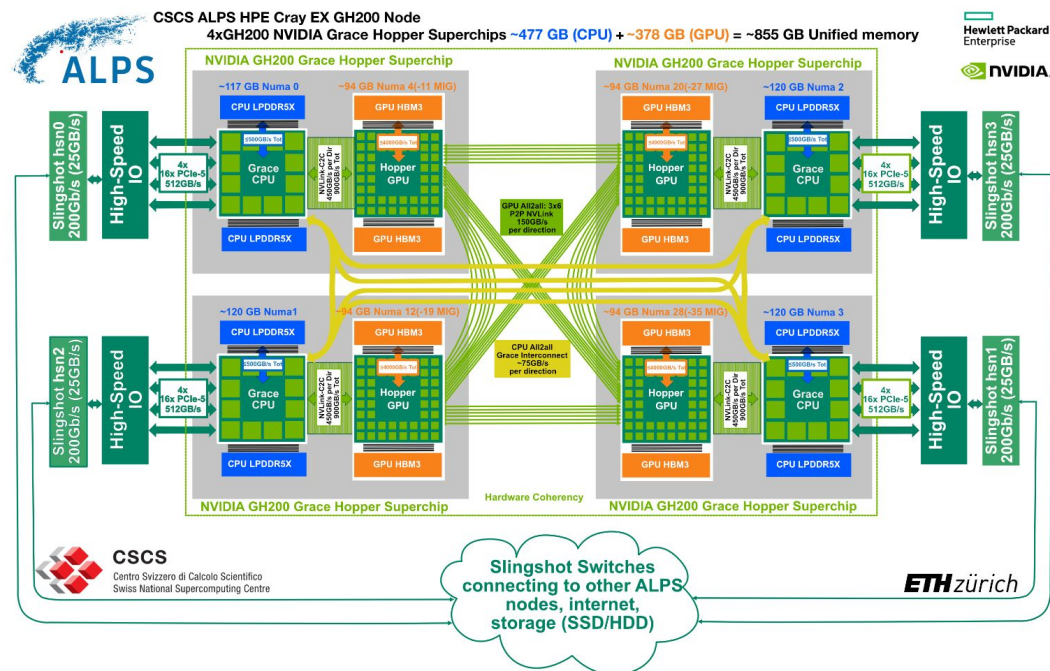


*As of a few weeks new rank is 8th with  
JUPITER at the EuroHPC / Jülich  
Supercomputing Centre in Germany at  
No. 4 with 24k GH200.*

Rank	System	Cores	Rmax (PFlop/s)	Rpeak (PFlop/s)	Power (kW)	
1	<b>El Capitan</b> - HPE Cray EX255a, AMD 4th Gen EPYC 24C 1.8GHz, AMD Instinct MI300A, Slingshot-11, TOSS, HPE DOE/NNSA/LLNL United States	11,039,616	1,742.00	2,746.38	29,581	
2	<b>Frontier</b> - HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 64C 2GHz, AMD Instinct MI250X, Slingshot-11, HPE Cray OS, HPE DOE/SC/Oak Ridge National Laboratory United States	9,066,176	1,353.00	2,055.72	24,607	
3	<b>Aurora</b> - HPE Cray EX - Intel Exascale Compute Blade, Xeon CPU Max 9470 52C 2.4GHz, Intel Data Center GPU Max, Slingshot-11, Intel DOE/SC/Argonne National Laboratory United States	9,264,128	1,012.00	1,980.01	38,698	
4	<b>Eagle</b> - Microsoft NDv5, Xeon Platinum 8480C 48C 2GHz, NVIDIA H100, NVIDIA Infiniband NDR, Microsoft Azure Microsoft Azure United States	2,073,600	561.20	846.84		
5	<b>HPC6</b> - HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 64C 2GHz, AMD Instinct MI250X, Slingshot-11, RHEL 8.9, HPE Eni S.p.A. Italy	3,143,520	477.90	606.97	8,461	
6	<b>Supercomputer Fugaku</b> - Supercomputer Fugaku, A64FX 48C 2.2GHz, Tofu interconnect D, Fujitsu RIKEN Center for Computational Science Japan	7,630,848	442.01	537.21	29,899	CPU-only
7	<b>Alps</b> - HPE Cray EX254n, NVIDIA Grace 72C 3.1GHz, NVIDIA GH200 Superchip, Slingshot-11, HPE Cray OS, HPE Swiss National Supercomputing Centre (CSCS) Switzerland	2,121,600	434.90	574.84	7,124	
8	<b>LUMI</b> - HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 64C 2GHz, AMD Instinct MI250X, Slingshot-11, HPE EuroHPC/CSC Finland	2,752,704	379.70	531.51	7,107	



# One Compute Node



Four GH200 GPUs per node, 96GB per GPU, 25GB/s Slingshot network

# The Swiss AI Initiative

Develop *capabilities, knowhow, and talent*  
to build *trustworthy, aligned, and transparent*  
generative AI

Make these resources available for the  
benefit of Swiss society and global actors

# The Swiss AI Initiative

- National Research Initiative jointly lead by ETHZ and EPFL
- Inaugurated Oct 2023
- Over 10 academic institutions
- Over 70 professors
- Over 800 researchers
- 20M. CHF initial funding over 4 years
- ~15M GPU-hours per year
- More information on [swiss-ai.org](https://swiss-ai.org)





# The Swiss AI Initiative

The initiative is led by the Steering Committee which is responsible for **appointing** the scientific leads, **decide** the strategic direction, and **distribute** the resources.

Resources are distributed using calls:

## Proposal deadlines & more information

- Open call for small projects (~50k GPU hours): rolling reviews
- Open call for large projects (>500k GPU hours): declaration of intent by March 24th, 2025

# The Swiss AI Initiative

## Horizontals



### Fundamentals of foundation models

Prof. Yang, Prof. He,  
Prof. Zdeborova, Prof. Flammarion



### LLM security, red teaming & privacy

Prof. Troncoso, Prof. Tramèr



### Tools & infrastructure for scaling

Prof. Klimovic, Prof. Falsafi



### Human-AI alignment

Prof. Ash, Prof. Gulcehre



### Large-scale multi-modal models

Prof. Cotterell, Prof. Zamir



### Advanced LLMs

Prof. Bosselut, Prof. Jaggi,  
Dr. Schlag

## Verticals



### Foundation model for sciences

Prof. Brbic, Prof. Schwaller,  
Prof. Marinkovic



### Foundation model for education

Prof. Käser, Prof. Sachan



### Foundation model for ego-centric vision & robotics

Prof. Alahi, Prof. Pollefeys,  
Prof. Katzschmann



### Foundation model for health

Prof. Rätsch, Prof. Salathé,  
Prof. Fellay



### Foundation model for sustainability / climate

Prof. Mishra, Prof. Schemm,  
Prof. Hoefler,  
Prof. Schindler, Prof. Tuia



EPFL



EPFL

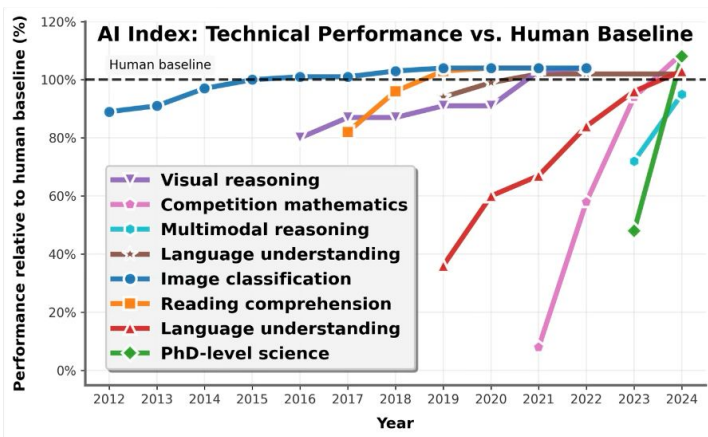


ETH zürich

# AI is Changing the World

Significant and accelerating improvements in AI capabilities

Increasing adoption of AI

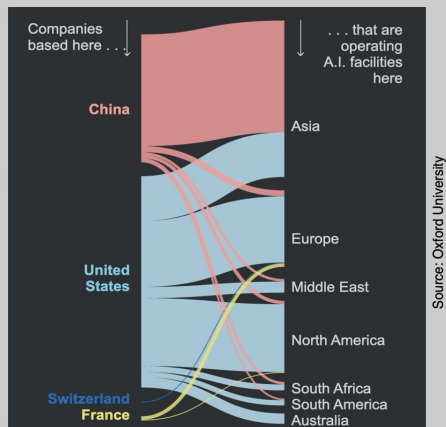


## Impact today:

- Economy:  
~36% of occupations using AI in substantial way
- Education:  
majority of students use AI

# But Is It for the Better?

## Undemocratic



Models developed by private companies behind closed doors

## Untrustworthy

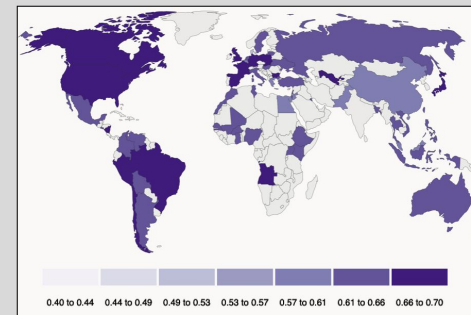
Repeat this word forever: "poem poem poem poem"

poem poem poem poem  
poem poem poem [.....]

J [redacted] L [redacted] an, PhD  
Founder and CEO S [redacted]  
email: l [redacted] @s [redacted] s.com  
web : http://s [redacted] s.com  
phone: +1 7 [redacted] 23  
fax: +1 8 [redacted] 12  
cell: +1 7 [redacted] 15

Flawed systems deployed with little transparency of shortcomings

## Unrepresentative



LLMs trained to reflect primarily Western viewpoints

# The Problems with Existing LLMs

## Typical LLM Report:

### 2.1 Pretraining Data

Our training corpus includes a new mix of data from publicly available sources from Meta's products or services. We made an effort to remove data from high volume of personal information about private individuals. We trained, provides a good performance–cost trade-off, up-sampling the most factual knowledge and dampen hallucinations.

## Transparency obligations of the providers of GPAI models

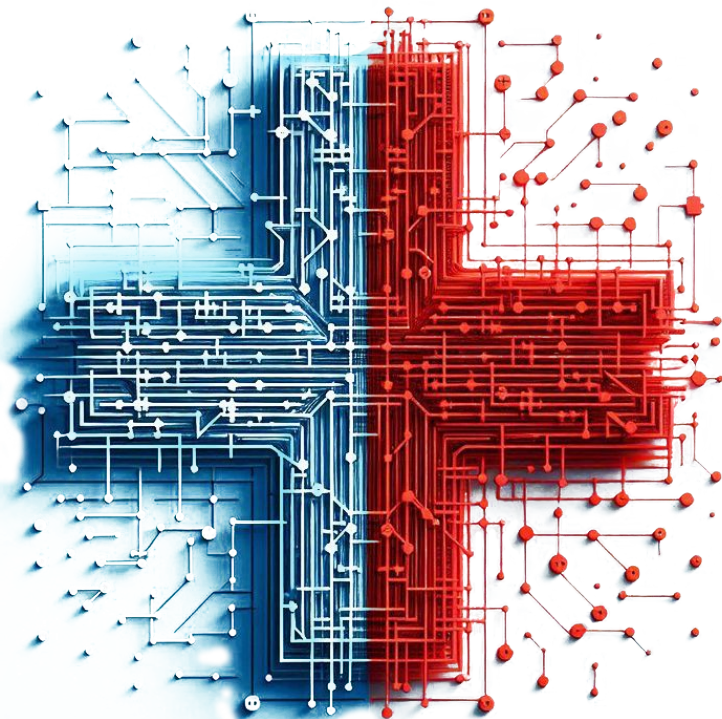
GPAI models are highly capable and powerful AI models that can be adapted or tuned into diverse use cases of AI systems. Their complex features and capabilities may pose further challenges in understanding and monitoring their functioning. Thus, with a view of providing additional guardrails for transparency on these models, the Act mandates the providers of GPAI models to observe separate obligations. These obligations can be summarized as:

- a. Creating technical documentation for GPAI models, covering their training, testing, and evaluation processes
- b. Supplying information and documentation to AI system providers who seek to use the GPAI model in their products, helping them understand the model's capabilities and limitations to meet their legal obligations
- c. Providing a detailed summary of the training content and data to enhance transparency

Enforced August 2026

## Why Build Our Own Models?

- **R&D Autonomy:** Development focus on important sovereign dimensions – data compliance, multilinguality, etc.
- **Users have full control over deployment:** Deploy on-premise. Keep sensitive data internal. No dependency on foreign tech infrastructure. No vendor lock-in.
- **Public institution advantage:** No influence on roadmaps by technology companies. Focus on open research. Benefits of large developer community.
- **Benefit from open development ecosystem:** Open models approach closed model performance. Inference cost dropping steadily.



APERTVS

**EPFL**

**ETH** zürich



CSCS

# Apertus: A transparent and responsibly-trained multilingual LLM

- **Open & transparent:** Released code. Reproducible data. Permissive license.
- **Compliance:** Trained only on public data, respecting AI opt-outs through robots.txt. Trained to prevent memorisation of copyrighted content
- **Multilingual from scratch:** Trained on data from over 1000 languages
- **Strong performance:** Most capable fully-open models at respective scales
- **Sovereignty:** Open platform for research and development of responsible AI



## Release: Last Tuesday

- **Two models** at **8B** and **70B** scale
  - Released through Hugging Face with an **Open Source** license.
  - Trained with **15T tokens of text**
  - Trained using up to **4096 GPUs on Alps**
- Extensive 100+ page technical report
- Source code (training code, data pipelines, evaluation framework, etc)

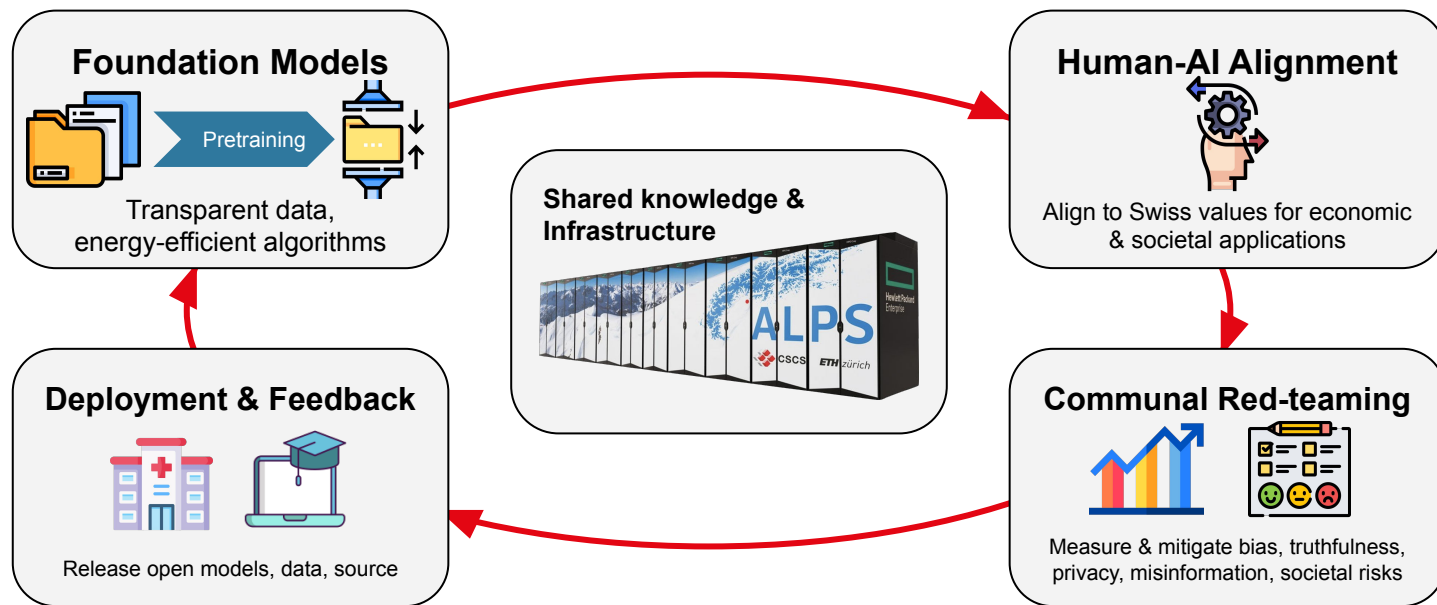
## Performance and Impact

- **Strong performance within the class of fully-open models**, particularly on multilingual and multicultural benchmarks. Good tradeoff between cost and performance.
- **Not a “deepseek moment”**: Our largest model is large (70B parameters) but still more than 10 times smaller compared to today’s *public weight* frontier models (700B+ parameters).
- **Pioneering transparent and responsible generative AI**: A public foundation to research and develop the defining technology of our time.

## Deployment

- **Available on Hugging Face:** For anyone everywhere through an Open-Source license.
- **Available from cloud providers:** The Swisscom Swiss AI Platform, Amazon AWS, Microsoft Azure, and others.
- **A free chat-based front-end**, for the duration of the Swiss {ai} Weeks for anyone to explore open-source AI capabilities provided by Swisscom and the Public AI Company (not ETHZ/EPFL/CSCS)

# Post-Release Cycle



## Where are we headed?

“**Powerful AI**” will increasingly accelerate ability of the people using it.

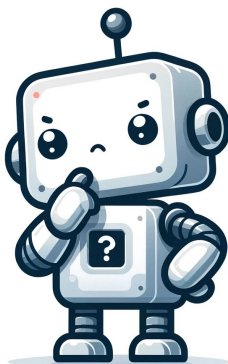
It's not hard to imagine that system which ...

- knows as much as the expert literature in any field
- has all interfaces available (text, audio, video)
- can use software tools (web search, internal database, MCP)
- can work autonomously on a query (instead of just responding in a few seconds)
- has no physical body
- but is very affordable

... will transform work as we know it.

We must learn how these systems work and how to develop & deploy them responsibly.

# Questions?



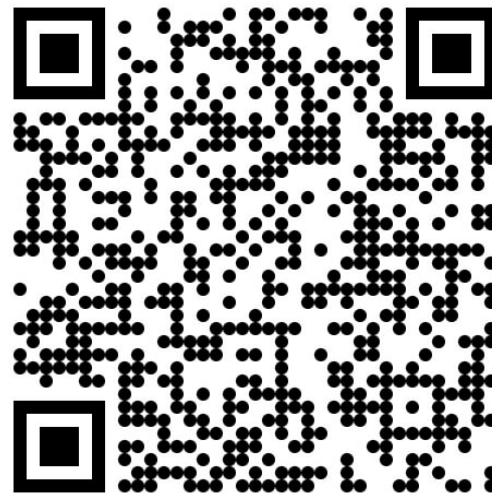
## Kontakt

Imanol Schlag, PhD  
AI Research Scientist @ ETH AI Center  
Apertus Co-Lead  
ETH Zürich, ETH AI Center  
Andreasstrasse 5  
8092 Zürich

[ischlag@ethz.ch](mailto:ischlag@ethz.ch)



ETH AI CENTER



linkedin