# EPFL



# Training multilingual **Sovereign LLMs in the** Swiss context

#### Martin Jaggi

19 June 2025, Davos Cyber Alp Retreat









# Generative AI is the fourth wave of the computing revolution





# **EFFL** Al is Changing the World

### Significant and accelerating **improvements** in AI capabilities



Select AI Index technical performance benchmarks vs. human performance

Source: Al Index, 2025 | Chart: 2025 Al Index report

### **Increasing adoption** of Al

## **Impact today:**

- economy:
  - ~36% of occupations using
  - Al in substantial way
- education:
  - majority of students use Al



# **ETH** But is it for the better?



Φ sAler NCCR



# **ETH** Who gets to decide the future?

#### TECH / META

Meta won't release its multimodal Llama Al model in the EU / Meta says the European regulatory environment is too 'unpredictable.'

by Jess Weatherbed Jul 18, 2024, 8:07 PM GMT+8

**Canada Joins Russia And China In Being** Excluded From Google's Al Chatbot, Bard



Growing Conflict

**ChatGPT banned in Italy over privacy** concerns

French competition watchdog hits Google with 250 million euro fine

**US bans Chinese AI LLM DeepSeek from** government devices

**Brazilian** regulator bans Meta from using data for Al training









## Al systems that are developed, deployed, and governed by trusted national institutions, minimizing reliance on foreign entities.



### EPFL **EH Swiss Sovereign Al Research**

### Science of Al

Foundational scientific advances in Al





## Sovereign Al

a sovereign Al ecosystem for Switzerland that is open, trustworthy, and safe to deploy in society

## sAlence

collaboration across technical and societal disciplines





### Context:

- Al has unprecedented speed and depth of societal impact
- 2/3rds of world population left out (+1) by leading AI models focusing on English and Chinese
- Lack of cultural and democratic representation
- Risks not well understood due science lagging behind commercial providers

### Unique and timely opportunity

- worlds largest public Al supercomputer, Sept 2024 (~11K GH200 GPUs)
- world-class talented researchers
   Swiss Al Initiative: 80+ research groups, 1000+ researchers, 12+ institutions



# **EPFL** Case Study: The Swiss Al Initiative LLM

- First release of new open-data open-weights 8B and 70B LLM planned in July 2025
  - fully transparent and compliant training data (CH,EU law)
    multilingual >1000 languages

### Highlights:

- architecture & recipe innovations
- rigorous data compliance & ethics
- improved data quality



# **Training New 8B and 70B Parameter LLMs**

Efficiency



Pretraining LLMs of varying sizes (in particular 8/70B!) concurrently across 4000 GPUs with engineering and scientific **innovations** in architectures, training recipes, data curation and evaluation.

**Live Training Loss** 



EPFL ML<sub>s</sub>O

Scale

#### Engineering

# **Data: Breaking Barriers with Multilingual Data**

#### **Multilingual**

Our multilingual datasets cover 15T tokens in almost 2000 languages and scripts and comply with training opt-outs.





EPFL ML<sub>s</sub>O

Open

Fair

#### **FineWeb2-HQ**

# **Data: Breaking Barriers with Multilingual Data**

#### Multilingual

Our multilingual datasets cover 15T tokens in almost 2000 languages and scripts and comply with training opt-outs.



#### Open

#### Fair

# **EPFL Ethical Training Data / Compliance**

- Web crawling
- Robots.txt : of today or at crawl time?
- Compliance Gap (downstream performance)



ning & Optimization Laboratory Lear Machine



# **EPFL Ethical Training Data / Compliance**





Machine Learning & Optimization Laboratory



# **EPFL Ethical Training Data / Compliance**



Machine Learning & Optimization Laboratory

# **EPFL Multilingual Data Curation**

FineWeb2 (>1000 langs)
 FineWeb2-HQ (classifier based filtering by instruct similarity)

Ours
1.8333
0.3550
0.6670
0.3870
0.6040
0.3400
0.3860
0.7510
0.5720
0.0820

Machine Learning & Optimization Laboratory

DCLM*	FW-Edu*	FW*
2.3889	2.4444	3.3333
0.3530	0.3850	0.3010
0.6470	0.6970	0.5880
0.4100	0.3770	0.3850
0.5960	0.5700	0.5930
0.3160	0.3470	0.3030
0.3840	0.4180	0.3560
0.7510	0.7410	0.7620
0.5610	0.5660	0.5550
0.1240	0.0320	0.0370

13

# **EPFL Multilingual Data Curation**

FineWeb2 (>1000 langs)
 FineWeb2-HQ (classifier based filtering by instruct similarity)

Dataset	Ours
Average Rank	1.8333
ARC (Challenge)	0.3550
ARC (Easy)	0.6670
CommonsenseQA	0.3870
HellaSwag	0.6040
MMLU	0.3400
OpenBookQA	0.3860
PIQA	0.7510
WinoGrande	0.5720
TriviaQA	0.0820

Machine Learning & Optimization Laboratory





# **EPFL** Mitigating the Curse of Multilinguality

- Quality filtering helps to turn a curse to a benefit
- 1B model, seen equal tokens in the target language





### average rank benchmark (french)

ng data		filtered	original	
19B rman	119B arabic	119B danish	1.83	3.11
		Garmon	2.06	3.00





#### 

# Not Just Llama: Model and Training Innovations

We leverage our expertise and ongoing research from machine learning, optimization and deep learning architectures to find stronger models and more efficient training recipes (and test them at scale!).



#### **Optimizers comparison** EPFL



ning & Optimization Laboratory Machine Lear



# **EPFL Constant LR & Cheap Scaling Laws**

constant learning rate for LLM training followed by cooldown

cost of scaling law



ning & Optimization Laboratory Ũ Machine







#### EPFL ML<sub>s</sub>O

# Not Just ML: Low-level software optimizations!

Continuous effort together with CSCS has allowed the large-scale deployment.

e 3	Efficient Backend	Optimized Megatron Framework	High- Performance File System
	Node Interconnect	Custom CUDA Kernels	Compute and Communicatio Overlap
	Optimal DP, TP, PP Size	Optimal Batch Size	









# **Ongoing Evaluation**

Massive

Multitask

#### Good English Performance; Great Multilingual Capabilities



#### Multilingual Language Understanding





#### CSCS ntro Svizzero di Calcolo Scientifico Swiss National Supercomputing Centre



# Thank you 🙏 EPFL



Antoine **Bosselut** 



Martin Jaggi



Imanol Schlag



Antoni-Joan Solergibert



**Bettina** Messmer



Allen Huang



Maximilian Böther



Simin Fan



lhor Protsenko



Matteo Pagliardini



Kyle Matoba



Roman Macháček



Jan Deriu



Andrei Semenov



Theofilos Manitaras



Auguste Poiroux



Stefano Schuppli



Lukas Drescher



Henrique Mendonça









Negar Foroutan



Romanou





Alexander Hägele



Alejandro Hernández









Anna Sotnikova



Zeming Chen



Badr **AIKhamissi** 



**Syrielle** Montariol





Paul Teiletche







Dhia

Garbaya

Joost VandeVondele



Dongyang Fan



Fawzi Roberto Mohamed



Ayush Kumar Tarun



Matin Ansaripour



Ilia Badanin





Tiancheng Chen





Eduard

Durech





#### Marfurt





Machine Learning & Optimization Laboratory



