



# **Towards Open Foundation Models for Europe**

### Sampo Pyysalo

sampo.pyysalo@utu.fi

10th RiskLab Finland, Bank of Finland and European Systemic Risk Board Joint Conference on AI and Systemic Risk Analytics

June 27th 2025





## **GROUP INTRODUCTION**





Natural Language Processing research group at University of Turku

Founded in 2001, now ~30 members

Focus on NLP for Finnish and multilinguality

Increasing recent emphasis on deep learning and LLMs









### One path toward foundation models for Europe

Introduction

Monolingual beginnings

Multilingual models

Models for Europe: OpenEuroLLM



Figure credit: Magnus Sahlgren





Large language models matter

Artificial intelligence is a **breakthrough technology** with broad impacts

Al applications require a strong **foundation model**, which most typically is (or builds on) a **Large Language Model** (LLM)

We want to interact with AI systems in **our own languages** and have them "understand" **our cultures and values** 



### **AI risks**

Job market disruption

Election interference

Market manipulation

Societal collapse

Existential risks







### Al risks

OPEN

## Al becomes a **highly influential technology** (socially, commercially, etc.)

### and

only few commercial actors are able to create leading foundation models



Chatbot Arena leaderboard (https://lmarena.ai/)



Large gaps in model capabilities

Availability: closed > open

Creators: industry > academia

**Geography**: US & China > rest

 $\rightarrow$  Leading models <u>closed</u> and created by <u>commercial actors</u> in <u>US and China</u>

Rank (UB) ↑	Model ↑↓	Score ↑↓	Votes ↑↓
1	<b>G</b> gemini-2.5-pro	1467	12,327
2	\$ o3-2025-04-16	1451	18,205
2	<b>G</b> gemini-2.5-pro-preview-05-06	1446	14,040
2	⑤ chatgpt-4o-latest-20250326	1442	22,488
3		1437	15,271
6	<b>A</b> claude-opus-4-20250514	1418	18,287
6	<b>G</b> gemini-2.5-flash	1418	17,535
6	♂ deepseek-r1-0528	1413	11,871
6	<pre></pre>	1411	16,362
8	x grok-3-preview-02-24	1409	24,316

Chatbot Arena leaderboard (https://lmarena.ai/)



## Three worlds of foundation models

10×

### **Closed models**

No access to weights, license restrictions

**Leading models**: GPT (US), Claude (US)

Compute: > le26 FLOPs (?)

### **Open-weights models**

No access to data or training process, license restrictions

Leadings models: Llama (US), Qwen (China)

Compute: > le25 FLOPs

### Fully open models

Open weights, data, training process, licenses, etc.

Leading models: BLOOM (US), OLMo (US)

Laskenta: < 1e24 FLOPs

10×





Foundation models should be <u>fully open</u>

### Visibility into training required to trust models

 $\rightarrow$  Need to understand training data and process

### Science requires transparency and reproducibility

 $\rightarrow$  Must be able to analyze and recreate models

### Natural point of collaboration for academia and industry

 $\rightarrow$  Shared interest in availability of open models





## Leading LLMs are made in the US and China

Support for European languages (esp. smaller) typically a secondary consideration at best

European AI efforts are left with two bad options

- 1. Build on foundations trailing behind leading models
- 2. Build on models created by U.S. or Chinese tech companies





LLama 3.2: Qwen Meta-llama has disallowed access to this model in the EU User: Tell me about the events of Jun 4, 1989. Downloads of this model are not accessible from the Assistant: I'm sorry, but I can't assist with that. European Union (EU). Please see the Llama Acceptable Use Policy and License FAQ page for more information.

> To build on a lasting foundation that reflects European values, we need the capability to build our own models  $\rightarrow$  not just a single model, but a process



\*\*\* \*\*\*\*



### Why are European LLM efforts trailing behind?



Creating a leading LLM requires

- Data: tens of trillions of words of high quality text
- **Compute**: tens of millions of GPU-hours on a supercomputer
- **Expertise**: the ability to put the two together



### Why are European LLM efforts trailing behind?



- Expertise: sufficient; many European groups have created LLMs
- Data: insufficient for most individual European languages, sufficient for multilingual models
- **Compute**: **insufficient**: no European initiative has enough compute to compete

### Why are European LLM efforts trailing behind?

The most ambitious European efforts to create open LLMs involve **less than 10%** of the compute required to train the largest Llama 3 model

OPEN





## INTRODUCTION

### Why are European LLM efforts trailing behind?

The most ambitious European efforts to create open LLMs involve **less than 10%** of the compute required to train the largest Llama 3 model





## INTRODUCTION











# **MONOLINGUAL BEGINNINGS**





Large language models for Finnish

Previously introduced first BERT model for Finnish, FinBERT, in 2019

In Feb 2022, granted pilot project on upcoming LUMI supercomputer to **build** LLMs for Finnish  $\rightarrow$  millions of GPUh (!)

At the time, no generative models with more than 1B parameters for Finnish

Key points of reference were **GPT-3** and **BLOOM** 





### **Compute: LUMI**

Fastest supercomputer in Europe when introduced (now 3rd)

100% hydro-powered

GPU partition ~3000 nodes with four **AMD MI250X (~24K GCDs)**,

Finnish share 25% (~25M GPUh/year)



https://www.lumi-supercomputer.eu/





### Finnish language

Uralic language natively spoken by ~0.1% of world population (<6M)

Latin script, but no large close neighbors  $\rightarrow$  **challenging transfer** 

Well represented online: e.g. > 10B tokens in Common Crawl







### **Training data**





### **Repeating training data**



Luukkonen et al. FinGPT: Large Generative Models for a Small Language (EMNLP'23)

### **First Finnish LLMs**

In early 2023, completed **GPT-3-like** family of models for Finnish in collaboration with Hugging Face, AMD and National Library of Finland:

- **FinGPT** family of Finnish models, trained from scratch (300B tokens)
- **BLUUMI**, CPT Finnish variant of 176B-parameter BLOOM model

Fully open (Apache 2):

https://turkunlp.org/gpt3-finnish

Model	Layers	Dim	Heads	Params
Small	12	768	12	186M
Medium	24	1024	16	437M
Large	24	1536	16	881M
XL	24	2064	24	1.5B
3B	32	2560	32	2.8B
8B	32	4096	32	7.5B
13B	40	5120	40	13.3B
BLUUMI	70	14336	112	176B



### **Evaluation: zero/few-shot**

New Finnish evaluation dataset FIN-bench (following BIG-bench)

Q: How many elephants can fit in a freezer?

- four
- none
- two
- one
- six

Q: tin	What is 60 nes 41
•	3675
•	7395
•	58
•	2460
•	banana
•	house

Claim: During photosynthesis, plants may make glucose and carbon dioxide.

TrueFalse



OPEN EURO

Luukkonen et al. FinGPT: Large Generative Models for a Small Language (EMNLP'23)



### **Evaluation: 3-shot results**

Substantial advance over previous models:  $40\% \rightarrow 60\%$ 

Indications of overfitting

for largest monolingual model (**13B**)

**OPEN** 

**BLUUMI**: best Finnish results, no drop on English



# Μ

BLUUMI



### **Perspective / lessons**

Limited monolingual data should not be stretched too far

GPT-3 models oversized and undertrained, esp. largest (~175B params)

Good as proof of concept and research, less so as a foundation for Finnish AI

Model quality constrained by **data** and **compute** 









### Beyond monolingual training

For Finnish (<40B tokens), **limits of monolingual training visible** at ~10B params and 300B tokens

**Goal**: larger and better models prioritizing Finnish performance

**Solution**: bilinguality (Finnish+English) + code + cross-lingual signal





**Pretraining data** 

Data sources:

- Finnish: As in FinGPT, excluding National library resources
- English: SlimPajama
- Code: StarCoder
- Parallel: OPUS



Subset	Tokens	Weight	Sampled
Finnish	32.6268 B	4	130.507 B
English	546.236 B	1	546.236 B
Eng-Fin	8.04152 B	1	8.04152 B
Code	315.418 B	1.52	315.419 B
Total	794.416 B	-	1000.20 B



Poro 34B

**34B parameter language** model trained on **1T tokens** of Finnish, English and code

Collaboration with Silo AI (https://silo.ai)

## Fully free & open:

https://huggingface.co/LumiOpen/Poro-34B







Luukkonen et al. Poro 34B and the Blessing of Multilinguality (Nodalida'25)



### Poro 34B evaluation

Substantial advance over previous models in Finnish: 50% (FinGPT) / 60% (BLUUMI)  $\rightarrow$  66% (Poro)

Broadly competitive with open models with <u>similar param/token</u> <u>counts</u> for **English and code** 



<b>OPEN</b>
EURO
IM

	Poro 34B	Llama 33B	MPT 30b	Falcon 40B	FinGPT 8B	FinGPT 13B	Starcoder
Finnish	66.28	53.36	53.22	42.58	49.69	48.92	45.55
English	50.57	59.96	52.62	49.87	31.47	32.85	35.44
Code	41.80	37.67	39.18	38.57	-	-	49.06



### **Perspective / lessons**

Learned lessons and corrected mistakes from FinGPT work

Bilingual + code training recipe successful

Poro is far from a frontier LLM, but good enough to build on (Zosa et al. 2025)

Model quality constrained by compute, but bilingual data sufficient



Zosa et al. Got Compute, but No Data: Lessons From Post-training a Finnish LLM (Nodalida'25)



Viking model family

Family of **7B**, **13B** and **33B** models trained on **2T** tokens of Nordic languages (fi, sv, no, da, is), English and code

Released fully openly (Apache 2)

https://huggingface.co/LumiOpen/Viking-7B https://huggingface.co/LumiOpen/Viking-13B https://huggingface.co/LumiOpen/Viking-33B



Poro 2

**8B and 34B parameter** models created through continued pretraining of Llama models

 $\rightarrow$  Dependent on strong open foundation from Meta

Collaboration led by AMD Silo AI (https://silo.ai)



https://huggingface.co/LumiOpen/Llama-Poro-2-70B-Instruct









### **Perspective / lessons**

Amount of Finnish data far too limited to create competitive Finnish-only models

Data constraints alleviated through **multilinguality** or **continued pretraining** 

Access to efficiently deployable **compute** a limit to capability

Continued pretraining allows strong models capable of generating Finnish, but these remain dependent on open-weights foundation models





# OPEN FOUNDATION MODELS FOR EUROPE: OPENEUROLLM



### About OpenEuroLLM

3-year project (2025-28) uniting 20 European universities, research orgs, companies and HPC centres to develop state-of-the-art, multilingual, fully open foundation models for Europe





### **OpenEuroLLM goals**

### Fully open process and models

- Open-source code
- Accessible pretraining data
- Fully documented process
- Models under open licenses

### 37+ languages $\rightarrow$

including all EU official and co-official languages

Compliant with EU regulations

Simple to deploy and use

Albanian Basque Bosnian Bulgarian Catalan Croatian Czech Danish Dutch English Estonian Finnish French Galician Georgian German Greek Hungarian Icelandic Irish Italian Latvian Lithuanian Macedonian Maltese Norwegian Bokmål Norwegian Nynorsk Polish Portuguese Romanian Serbian Slovak Slovenian Spanish Swedish Turkish Ukrainian





Scope: will OpenEuroLLM ...

Train models also on programming languages? Yes  $\checkmark$ Train models for instruction-following / dialogue (chat)? Yes  $\checkmark$ Train "reasoning" / "thinking" models? Yes  $\checkmark$ Explore architectures other than dense transformers (e.g. MoE)? Yes  $\checkmark$ Fine-tune models for specific use cases (e.g. science)? No  $\rightarrow$  LLMs4EU Train multimodal models (e.g. audio and images)? No  $\rightarrow$  upcoming projects





Selected milestones and deliverables

Month 6: Training datasets identified Month 18: First flagship models released Month 18: Performance and safety evaluations Month 36: Final flagship models released

**NB:** several releases planned in addition to flagship models, including small (~2B param) reference models in the following months



## **OPENEUROLLM ONLINE**





### A series of foundation models for transparent AI in Europe



OPEN EURO

including data, documentation, training and testing code, and evaluation metrics; including community involvement

#### COMPLIANT

under EU regulations, OpenEuroLLM will provide a series of transparent and performant LLMs

#### DIVERSE

for European languages and other socially and economically interesting ones, preserving linguistic and cultural diversity

Blog

company/openeurollm

openeurollm







**OPEN** =i irr











**Expertise** 

Project partners have led several efforts training European LLMs from scratch

The project is supported by many **others creating open LLMs** through Open Strategic Partnership Board and collaborations



Figure credit: Magnus Sahlgren



## H P L T



### Data

Project partners include key participants in **HPLT**, which creates one of the leading multilingual resources for LLM training (and other experts!)

We will also integrate the best existing and upcoming open datasets created by others



Figure from Penedo et al. (2025)



### https://hplt-project.org/





### Compute

Project partners have been granted **tens of millions of GPUh** in previous projects

**Approx. 30M GPUh** have been applied for from national, EuroHPC extreme scale, and AI Factories calls







We won't suddenly jump from far behind to overtaking OpenAI

... but **if the project is given enough compute**, we have the rest of the parts in place to deliver the next generation (or two) of open LLMs for all European languages





## **THANKS & ACKNOWLEDGEMENTS**



The contents of this presentation are the sole responsibility of the author and do not necessarily reflect the opinion of the European Union. This project has received funding from the Digital Europe Programme under grant agreement No 101195233





H