# Introduction



**LINAGORA**

**Jean-Pierre Lorré**
**Directeur de recherche**
jplorre@linagora.com

## Agenda

- Context: why Lucie 7B ?

- Lucie 7B details: training & evaluation

- Future work

Lucie-7B LLM Context

# LINAGORA

## Leader in Open Source

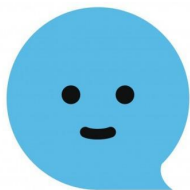**#GOODTECHFORGOOD**

**100% Open Source**

**GAFAM-free**

Founded in **2000**
**160** employees
**6** offices worldwide

## Research topics: NLP, Speech Recognition

LE VOICE LAB

ANITI

LinTO.ai

LinTO STUDIO

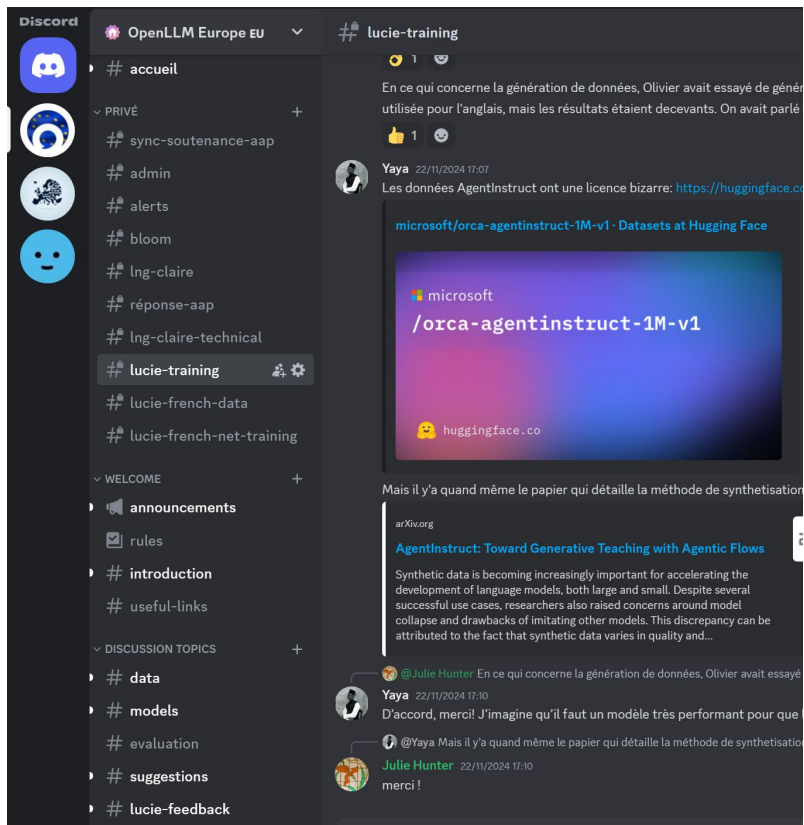**ASR & NLP platform**

**OpenLLM-Europe**

**Community for the development of sovereign, and truly Open Source LLM > 1 100 members**
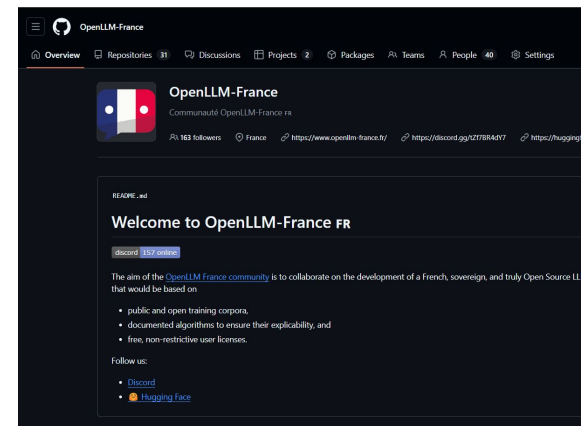
# The OpenLLM Community 🌍

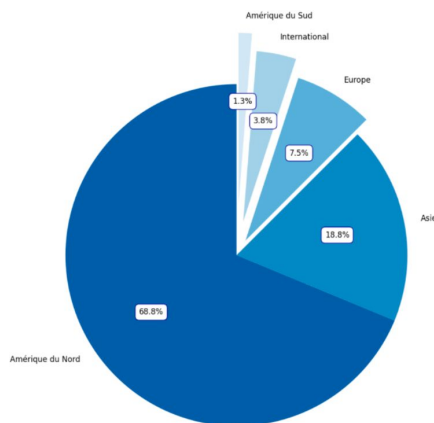Discord (1 100+ members)   Hugging Face (100+ members)   GitHub (40 dev, 100+ followers)



https://discord.gg/VFJxQnqrEU

# **Motivation –** Cultural Representation

## **Better representation of French and French-speaking communities**

Geographical distribution of LLMs with more than one billion parameters since 2018



LLAMA V2 : Language distribution in pretraining data with percentage

| Language | Percent | Language | Percent |
|----------|---------|----------|---------|
| en | 89.70% | uk | 0.07% |
| unknown | 8.38% | ko | 0.06% |
| de | 0.17% | ca | 0.04% |
| fr | 0.16% | sr | 0.04% |
| sv | 0.15% | id | 0.03% |
| zh | 0.13% | cs | 0.03% |
| es | 0.13% | fi | 0.03% |
| ru | 0.13% | hu | 0.03% |
| nl | 0.12% | no | 0.03% |
| it | 0.11% | ro | 0.03% |
| ja | 0.10% | bg | 0.02% |
| pl | 0.09% | da | 0.02% |
| pt | 0.09% | sl | 0.01% |
| vi | 0.08% | hr | 0.01% |

**Not just language:**
- History
- Politics
- Art
- Religion
- Social practices
- Cooking …

Llama 3: pretrained on 15T tokens and 5% non-English data
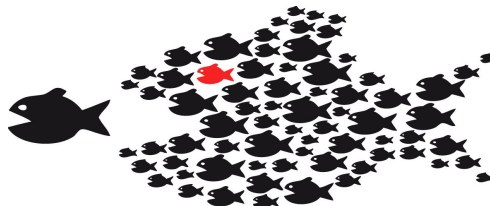
# Motivation – Truly Open Source AI

## Open Source AI Definition

| 4 Freedoms | Open Weights | Open Code | Data |
|---|---|---|---|
| Use<br>Study<br>Modify<br>Share | Model weights and (hyper-) parameters | Source code used to train the system<br>Source code used to create the dataset | The complete list of datasets used to train the system **and** the actual datasets when allowed |

A license that allows unrestricted usage
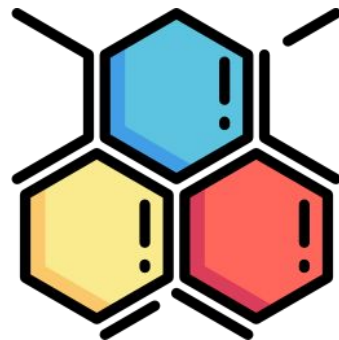
**Community impact**

## Motivation – Small Language Models

Small, specialised models can achieve comparable or even better performance than large, general-purpose models on given tasks

Smaller models are more resource efficient at both training and inference times

Small, specialised models can be combined with other AI or non-AI models to develop complex applications

# OpenLLM-France Project

T0: 01/09/24

2 years

10.5 M€

## Develop multimodal, voice and text LLM models that are trusted, controlled and transparent

Focus on the education application domain
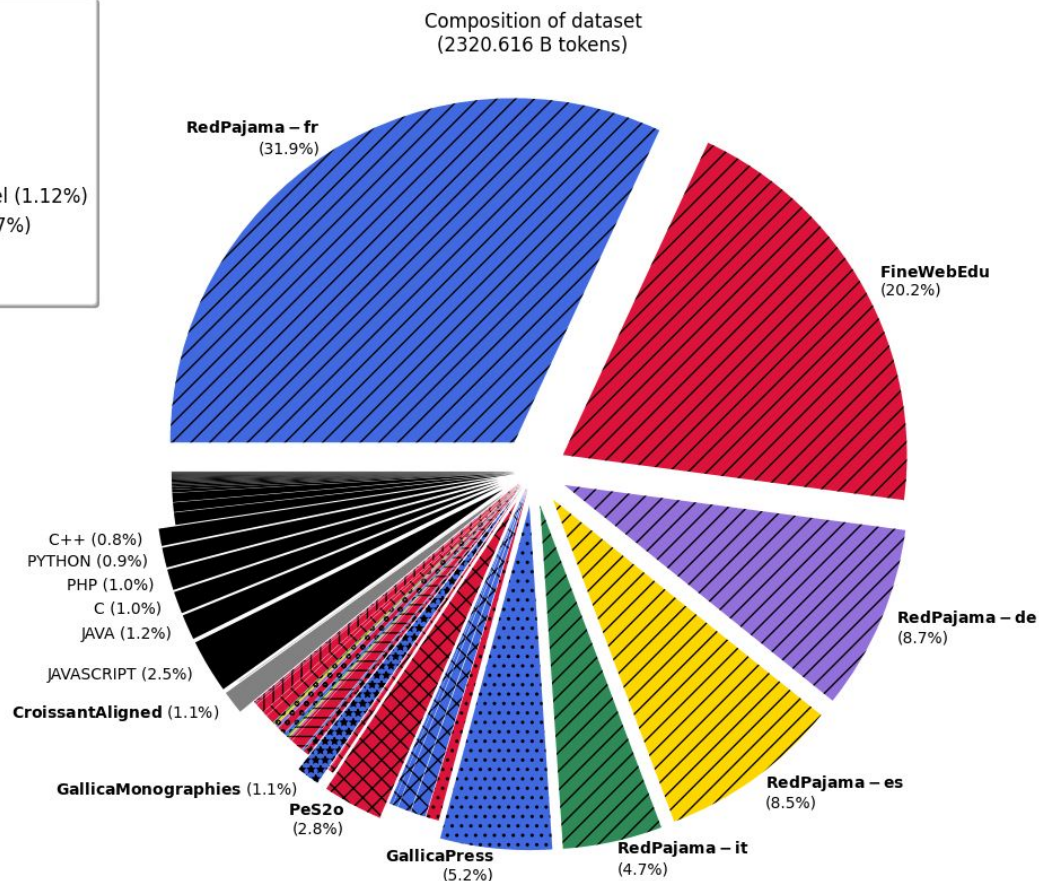Taking account of ethical, legal and environmental aspects

Lucie-7B LLM

🛢 **Lucie Training Dataset**   [OpenLLM-France/Lucie-Training-Dataset](OpenLLM-France/Lucie-Training-Dataset)

**Categories**
- ▨ Web (73.9%)
- ▦ Newspaper (5.86%)
- ◺ Technical (4.79%)
- ▦ Book (1.36%)
- ▤ Legislative (1.00%)
- ▦ Wiki (0.832%)
- ◹ Math (0.628%)
- ▥ Forum (0.536%)
- ▦ Dialogue (0.0779%)

**Languages**
- ■ French (40.3%)
- ■ English (26.4%)
- ■ German (8.90%)
- ■ Spanish (8.65%)
- ■ Italian (4.83%)
- ■ Multilingual Parallel (1.12%)
- ■ Programming (9.87%)

Composition of dataset
(2320.616 B tokens)

RedPajama – fr
(31.9%)

FineWebEdu
(20.2%)

RedPajama – de
(8.7%)

RedPajama – es
(8.5%)

RedPajama – it
(4.7%)

GallicaPress
(5.2%)

PeS2o
(2.8%)

GallicaMonographies (1.1%)

CroissantAligned (1.1%)

JAVASCRIPT (2.5%)

JAVA (1.2%)

C (1.0%)

PHP (1.0%)

PYTHON (0.9%)

C++ (0.8%)

# 🛢 Lucie Training Dataset - Data mix for Lucie pretraining

**Categories**
- Web (62.7%)
- Newspaper (4.60%)
- Technical (7.93%)
- Book (2.22%)
- Legislative (0.964%)
- Wiki (1.86%)
- Math (1.40%)
- Forum (1.00%)
- Dialogue (0.116%)

**Languages**
- French (32.1%)
- English (33.3%)
- German (6.93%)
- Spanish (6.65%)
- Italian (3.79%)
- Multilingual Parallel (2.50%)
- Programming (14.7%)

- Upsampling of English and higher quality data sets
- Final proportions:
  - French        40%  → **33%**
  - English       26%  → **33%**
  - Web data   74%  → **63%**
  - Code          10%  → **15%**

Composition of training dataset
(3110.235 B tokens)



- RedPajama – fr (23.8%)
- FineWebEdu (22.6%)
- RedPajama – de (6.5%)
- RedPajama – es (6.3%)
- RedPajama – it (3.5%)
- GallicaPress (3.9%)
- HAL (1.0%)
- Theses (0.9%)
- PeS2o (5.3%)
- GallicaMonographies (1.6%)
- MathPile (0.9%)
- Pile(StackExchange) (0.8%)
- CroissantAligned (2.4%)
- JAVASCRIPT (3.8%)
- JAVA (1.8%)
- C (1.5%)
- PHP (1.5%)
- PYTHON (1.4%)
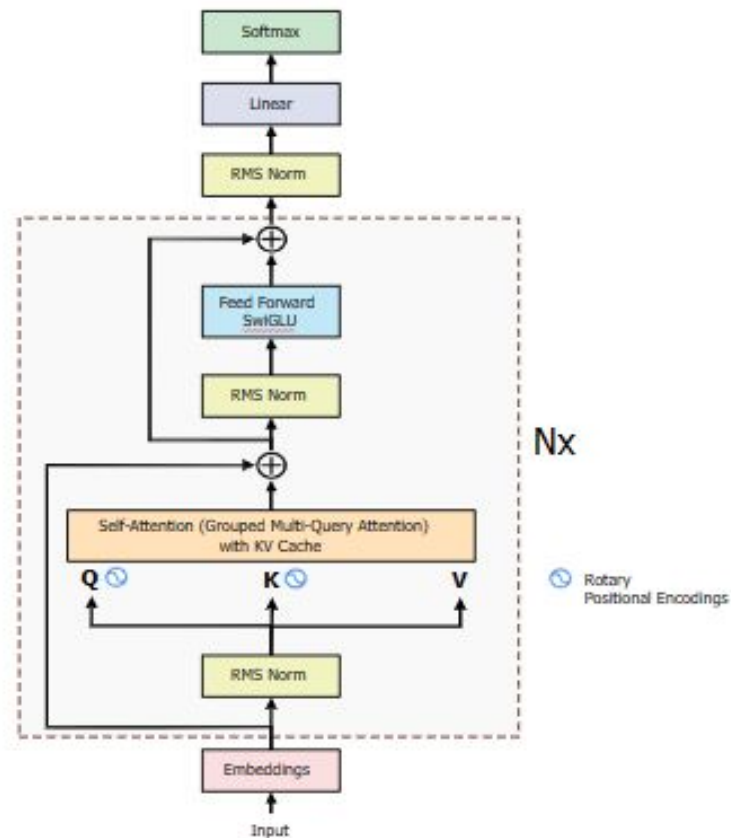- C++ (1.2%)
- C# (0.9%)

# 🤖 Architecture

- Causal decoder-only model. Next-token prediction.
- Llama 3 architecture
  - Group Query Attention

  - Rotary Positional Embedding (RoPE)

  - Configuration:
    - Vocabulary size: 65 k tokens
    - Layers:          32
    - Hidden size:     4096
    - Context length:  4096 (étendu ensuite à 32k)

- Most of the weights lie in the Feed-Forward Networks

  | | | | | | |
  |---|---|---|---|---|---|
  | a. | Embedding: | 266 M | (x2) | → | 0.53 B |
  | b. | Attention block: | 42 M | (x32) | → | 1.34 B |
  | c. | FFN block: | 151 M | (x32) | → | **4.83 B** |
  | | | | | Total | 6.70 B |

## Tokenization – Text Pre-Processing Options and Constraints on Tokens

| | Bloom, GPT, Falcon, OLMo | Gemma | Llama2, Mistral | Croissant | Lucie |
|---|---|---|---|---|---|
| Number of tokens | 65 – 250k | 256k | 32k | | 65k |
| Avoid OOV : byte-level / byte fallback | Byte-level BPE | | Unicode-level BPE with byte fallback | | |
| Unicode Normalization | (NFC for OLMo) | | | NFKC | NFC |
| Enforced split : isolate digits<br>abc12.3__4 → __abc \| 1 \| 2 \| . \| 3 \| __ \| 4 | ✔ | ✔ | ✔ | ✔ | ✔ |
| Enforced split : separate punctuation<br>abc. de.f... → __abc \| . \| __de \| . \| f \| ... | ✔ | ✔ | ✔ | | ✔ |
| Consecutive spaces<br>____\t\t\t\t\n\n → ____ \| \t\t\t\t\t \| \n\n | learned<br>(some only for __) | fixed<br>(max 30) | learned | | fixed<br>(max 8\|4\|2) |
| Prefix first words with space : at start / also after other kind of spaces | at start only | | | | also after<br>\n\t ( [ ' " « <- |

# Training Pipeline

Three pre-training phases:

- Main pre-training phase
  - 3.1T tokens
  - Knowledge of the world acquisition

- Context extension phase
  - 5B tokens
  - Extend the context length from 4096 to 32k tokens

- Annealing phase
  - 5B tokens
  - High-quality dataset with a focus on mathematical content



Composition of training dataset
(3110.235 B tokens)

Categories:
- Web (62.7%)
- Newspaper (4.60%)
- Technical (7.93%)
- Book (2.22%)
- Legislative (0.964%)
- Wiki (1.86%)
- Math (1.40%)
- Forum (1.00%)
- Dialogue (0.116%)

Languages:
- French (32.1%)
- English (33.3%)
- German (6.93%)
- Spanish (6.65%)
- Italian (3.79%)
- Multilingual Parallel (2.50%)
- Programming (14.7%)

# Parallel Training



Transformer layer #1 → Transformer layer #2

Lucie-7B was pre-trained on:

- 512 80GB-VRAM H100
- for about 500k GPU hours

The training code is based on a fork of Megatron-DeepSpeed

**3D Parallelism:**

- Data  Parallelism — 32
- Pipeline Parallelism — 4
- Tensor Parallelism — 4
- Batch size ~ 4M tokens



## Lucie Training

- Setup
  - Clone the repository
  - Environment setup
    - With python virtual environment (conda)
    - With Docker
  - Install Megatron-Deepspeed
- Train a model
  - 1. Pretraining (first main phase)
  - 2. Context Extension
  - 3. Annealing
  - 4. Instruct-Tuning and Finetuning
- Model conversion
  - From Megatron-Deepspeed to transformers
  - From LORA (PEFT) to full weights
  - Quantize models

# 📈 Learning Curves & Benchmark Evaluations



🇬🇧

**ARC Challenge (25-shot)**

**Hellaswag (10-shot)**

**MMLU Continuation**

**Winogrande (5-shot)**

**GSM8k (5-shot)**

**Truthfulqa MC2**

🇫🇷

**French Bench ARC Challenge (5-shot)**

**French Bench Hellaswag (5-shot)**

**French Bench Grammar (5-shot)**
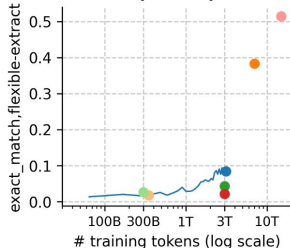
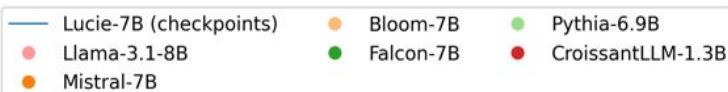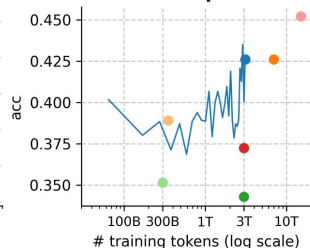**French Bench Vocab (5-shot)**

🌍

**ARC French (25-shot)**

**ARC Spanish (25-shot)**

**ARC German (25-shot)**

**ARC Italian (25-shot)**

Legend:
- Lucie-7B (checkpoints)
- Llama-3.1-8B
- Mistral-7B
- Bloom-7B
- Falcon-7B
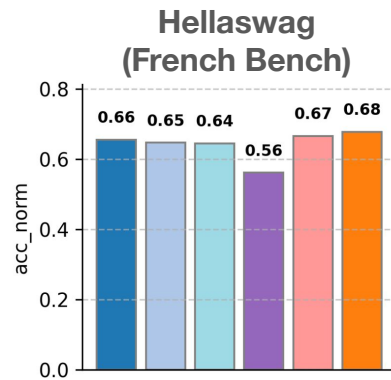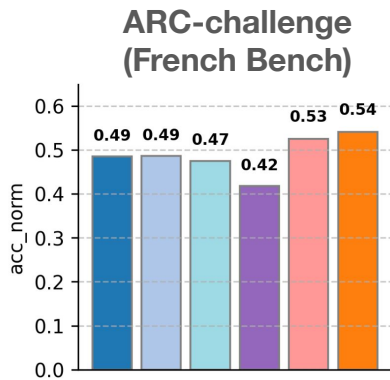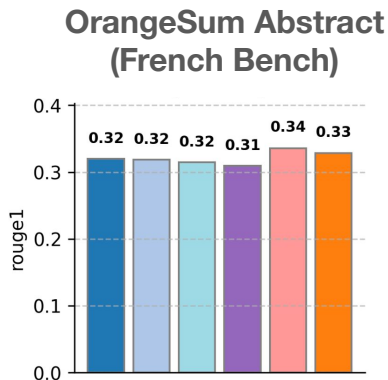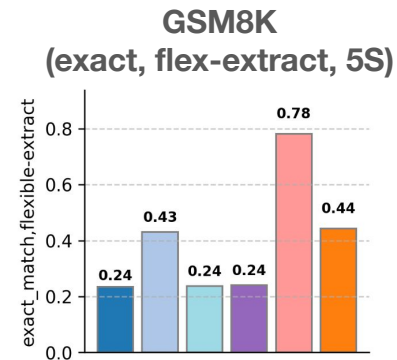- Pythia-6.9B
- CroissantLLM-1.3B

- Lucie state-of-the-art on French
- Right direction for German, Spanish and Italian
- OK in English, but not at the level of Llama-3 (8B) and Mistral-7B for Multitask Language Understanding (MMLU, Winogrande, GSM8k)

# 📈 Instruction Tuning (the start...)

**FQUAD v2 Gen (French Bench)**

- Math
  - Beats Llama 2, competitive with Mistral Instruct
  - Falls short of Llama 3.1
- Language-dependence
  - Competitive on French benchmarks
  - Less so when the benchmarks are in English

**GSM8K (exact, flex-extract, 5S)**

**OrangeSum Abstract (French Bench)**

**ARC-challenge (French Bench)**

**Hellaswag (French Bench)**

**ARC-challenge**

**Hellaswag**

Legend:
- Lucie-7B
- Lucie-7B-Instruct-v.1.1
- Lucie-7B-Instruct (human data)
- Llama-2-7B-Instruct
- Llama-3.1-8B-Instruct
- Mistral-7B-Instruct

🏛️ **Shared & Open Resources**   https://github.com/OpenLLM-France/Lucie-Training

**Future work**

## 🔮 Future Work

- Model Alignement
  - Model for Education (OpenLLM project just started – kick-off was 21/01/2025)
  - Propose open test platform to get community feedbacks
  - Start improvement loop with Reinforcement Learning (GRPO …)
- Reasoning
  - Function Calling (for math, physics, …) to calculators and API
  - **R**etrieval-**A**ugmented **G**eneration
- Multi-modality (Text prompt + Audio [+ Image/Video])
- Scale to more languages & alphabets (Greek, …), handle code switching and multilingual inputs
- Smaller model (1B) – Distillation and/or Training from scratch
- Data mix improvement (quantity, quality, nature)
- New Architectures
  - MAMBA (more linear, more efficient)
  - Hybrid Transformers(Attention) / RNN(LSTM) – TITAN

# LINAGORA

# MERCI

**Discord OpenLLM-Europe**

**Jean-Pierre Lorré**
jplorre@linagora.com
+33 6 88 34 63 85

Open LLM
Europe