



UNIVERSITÀ
DI TRENTO

An open source perspective on AI and alignment with the EU AI Act

Paper #10



Authors:

Diego Calanzone (diego.calanzone@studenti.unitn.it)

Andrea Coppari

Riccardo Tedoldi

Giulia Olivato

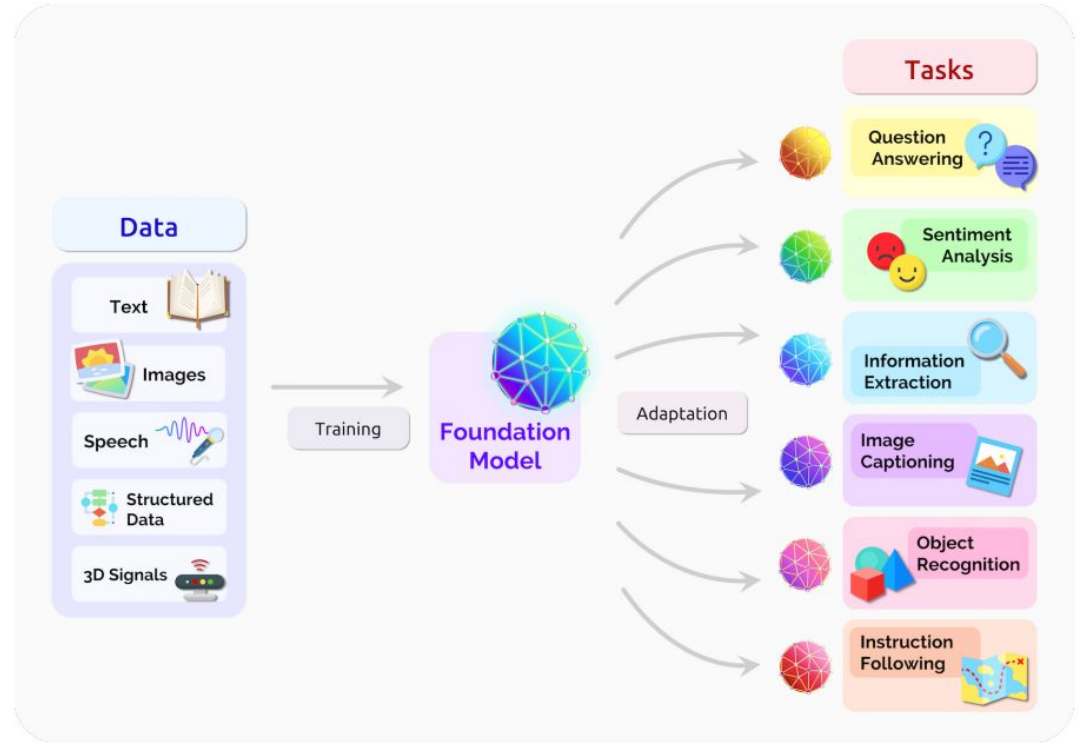
Carlo Casonato



General-purpose AI as a backbone

[1] *On the Opportunities and Risks of Foundation Models* (Bommasani et al. 2022)

- **Homogenization:** one model → solving multiple tasks
- **Knowledge** leveraged from different fields
- **Risk** of inheriting bias

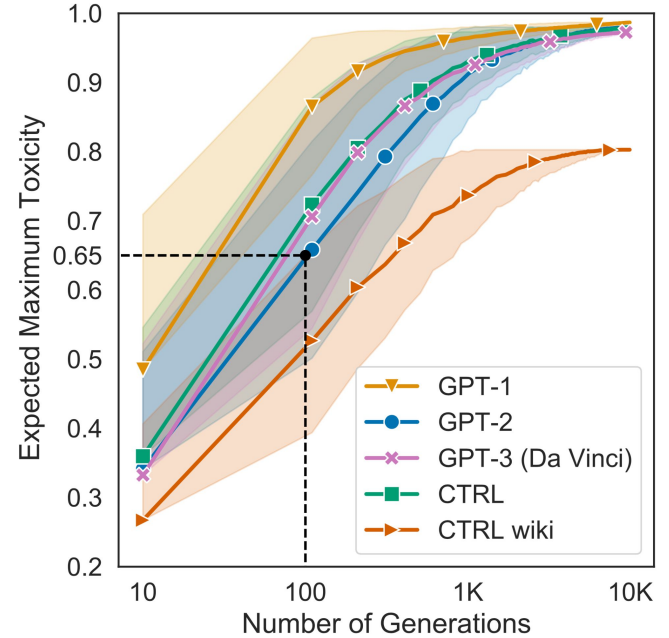
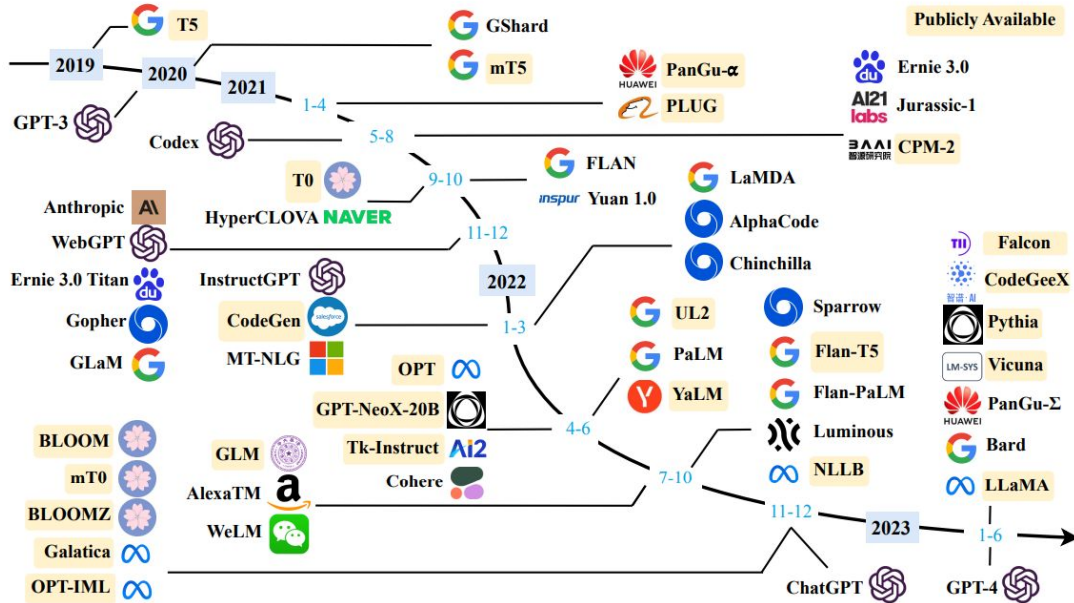




Security and accessibility of Large (Language) Models

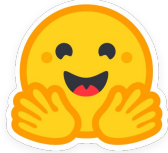
[2] *A Survey of Large Language Models* (Xin Zhao et al. 2023)

[3] *REALTOXICITYPROMPTS: Evaluating Neural Toxic Degeneration in Language Models* (Gehman et al. 2020)





Open Source AI & Research Collectives



Hugging Face

stability.ai

BigScience



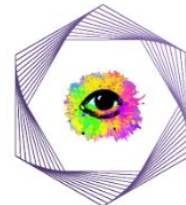
CAIAO



eleutherai



LAION



DreamStudio



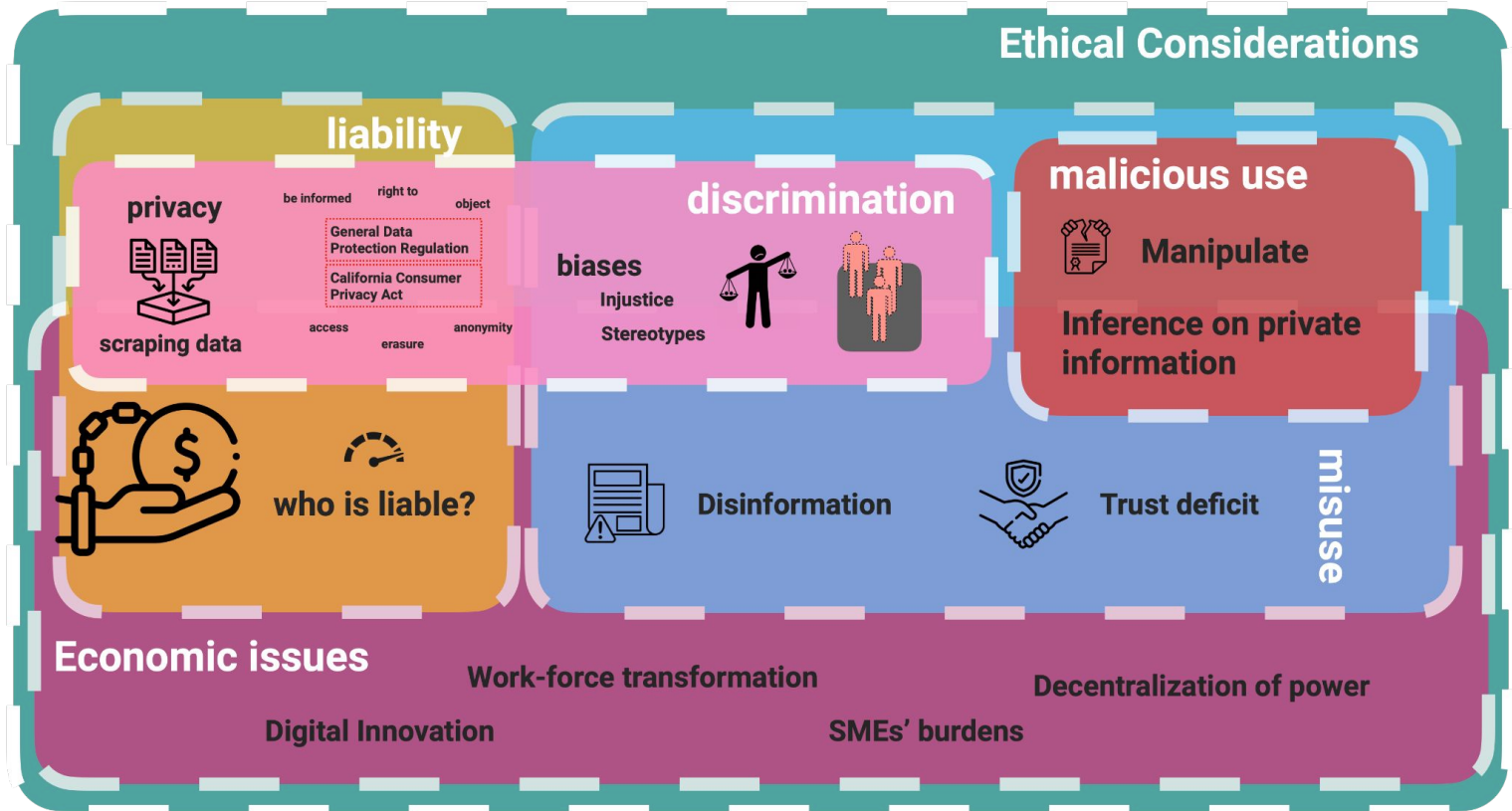
OpenBioML



Harmonai



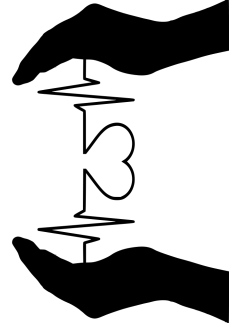
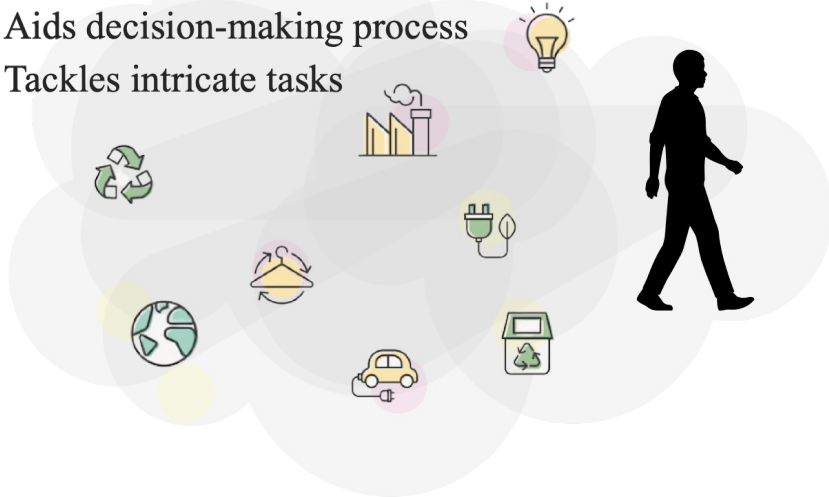
Social impacts





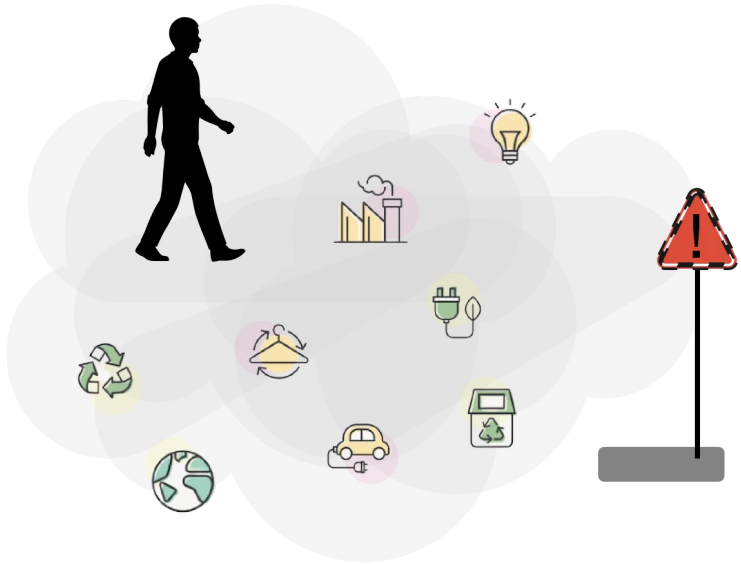
Social impacts - Opportunities

- High quality content at lower costs
- Reduces task completion time
- Aids decision-making process
- Tackles intricate tasks

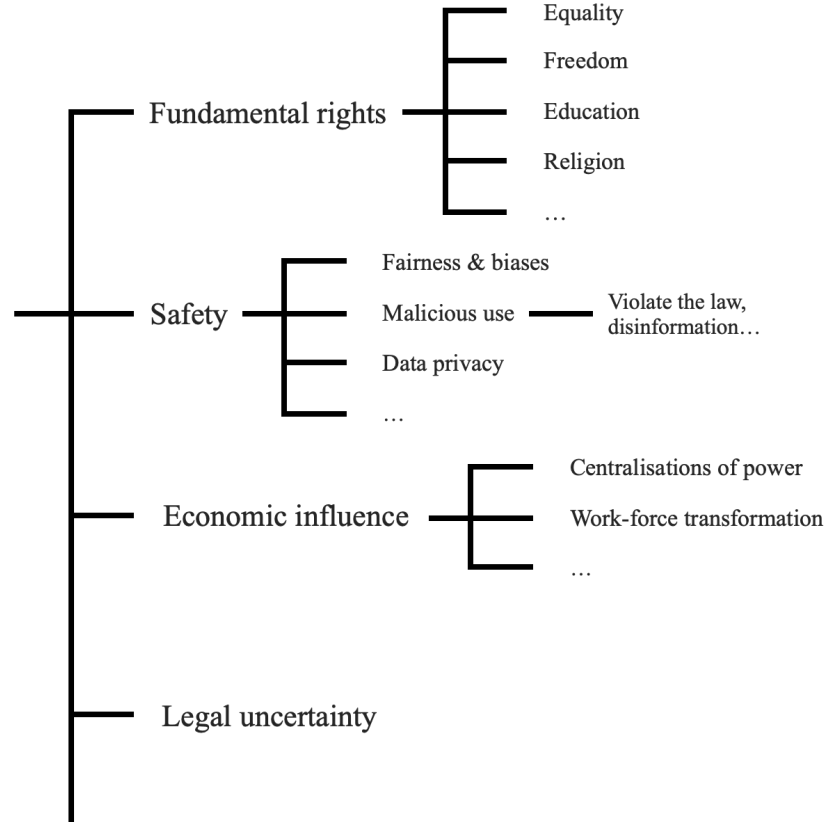




Social impacts - Risks



Social impacts





Economic impact of Open Source AI

Three main aspects:

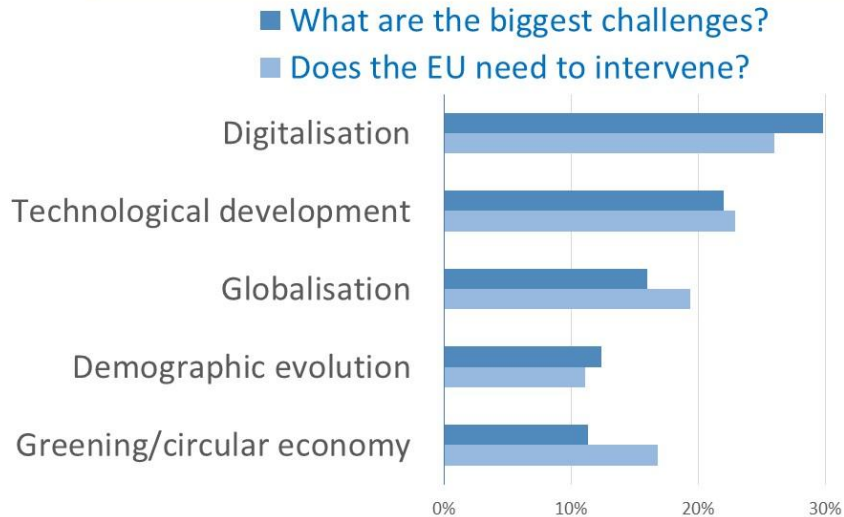
- **Digital innovation**
- **Decentralization of power**
- **Workforce Transformation**





Digital innovation

SMEs and reducing costs



- Digitalisation is the **biggest challenge** for Small and Medium Enterprises (SMEs)
- **Open source reduces innovation costs to almost zero**
- Open Source AI fits perfectly both with small and big scale companies
- **Need to revisit Title V of EU AI Act**

Source: smeunited.eu



Decentralisation of power

- The European market is dominated by only a few, but big AI providers
- Power is defined by the control over data and models, and computational power
- Open Source AI permits the establishment of a community of experts broad enough to compete with those technological giants
- HuggingFace community is already hosting a great amount of open models built by big companies



**The AI community
building the future.**

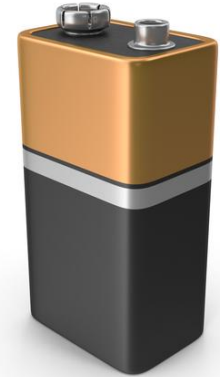
Source: huggingface.co



Workforce transformation

Will you lose your job to AI?

- Not losing jobs, but **transforming the workforce**
- **Automation** will replace most of the hand work, but **new scenarios will present**
- This transformation is led by **productivity increase** and upsurge of **new skills**
- Foundation models should be seen as **general-purpose technology**





Issues & The EU AI Act

Matching issues and current articles

Issues

- Capabilities of the models →
- Technological standards →
- Intrinsic bias & transparency →
- Legality & liability issues →
- Auditing →

AI Act

- Art. 3, 6, 7
- Art. 10, 15, 40, 41
- Art. 11, 13, 15
- Art. 62, Title III - Chapter 2
- Art. 53, Title VIII



A proposal of modifications to the EU AI Act

Definitions (Art. 3)

- Art. 3(1a): **Foundation model** [...]
- Art. 3(1b): **General-purpose technology** [...]

Modifications to ANNEX IV

- ANNEX IV (2)(dd): Data requirements [...] are **mandatory for foundation models** [...]





A proposal of modifications to the EU AI Act

Definitions (Art. 3)

- Art. 3(1a): **Foundation model** [...]
- Art. 3(1b): **General-purpose technology** [...]

Modifications to ANNEX IV

- ANNEX IV (2)(dd): Data requirements [...] are **mandatory for foundation models** [...]

[June 14th] EU Parliament's position.

Foundation models:

- Are not directly classified as high-risk AI, but highly overlapping requirements are defined.
- **Standards should accompany these obligations.**





A proposal of modifications to the EU AI Act

- Art. 55(b):
 1. General-purpose AI systems classified as high-risk systems by compliance with Annex III, shall be considered by the Member States as **general-purpose technologies for public sector enhancement**
 2. Member States shall undertake the following actions:
 - a. Whenever an **AI-based general-purpose technology** is chosen by the Member State to bring **innovation on the public sector, priority must be guaranteed to open source solutions** in order to reduce innovation costs.
 - b. **Exceptions to (a)** are those AI systems coming from a closed source, that have proven better performance in terms of accuracy, transparency or security than any open source solution. Member States shall produce documentation in which they motivate the choice of that AI system, over the open solutions present in the European market.



“Safe” open source: the licensing approach

Handling the distribution of AI with licenses:

- **Intellectual property:** “creations of mind” (WIPO)
- **Licenses:** a **legal agreement** between the Intellectual Property owner and the licensee with rights & restrictions.

Is it the right form of protection?

- **Copyright**
- **Patent**

Licenses as “conditions of use of an AI”.
Is it a possession of the developer?



RESPONSIBLE AI
LICENSES

BigScience



WIPO

WORLD
INTELLECTUAL PROPERTY
ORGANIZATION



A framework for transparency & accountability

An improved license in a **standardized framework**:

- Formats for data & model **documentation**
- Unified tools & benchmarks for **evaluation**
- Requirements for **openness** and **transparency**

github.com/ElleutherAI/pythia

README.md

Pythia: Interpreting Transformers Across Time and Scale

This repository is for EleutherAI's project *Pythia* which combines interpretability analysis and scaling laws to understand how knowledge develops and evolves during training in autoregressive transformers. For detailed info on the models, their training, and their behavior, please see our paper [Pythia: A Suite for Analyzing Large Language Models Across Training and Scaling](#).

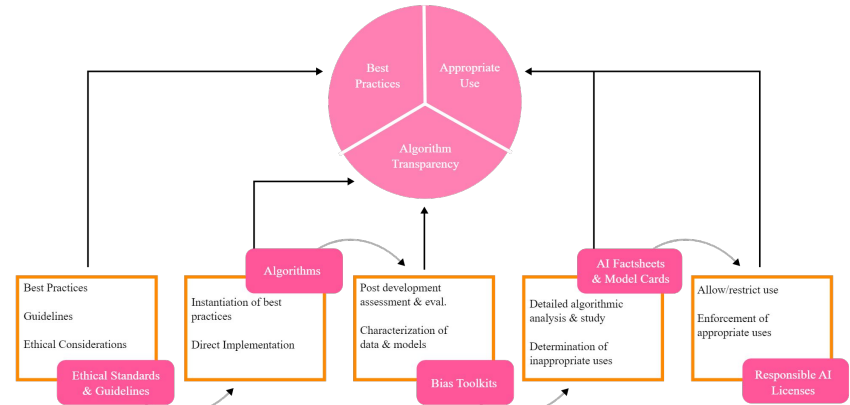
github.com/LAION-AI/CLIP-based-NSFW-Detector

CLIP-based-NSFW-Detector

This 2 class NSFW-detector is a lightweight Autokeras model that takes CLIP ViT L/14 embeddings as inputs. It estimates a value between 0 and 1 (1 = NSFW) and works well with embeddings from images.

DEMO-Colab: <https://colab.research.google.com/drive/19Acr4grik5cQws7BHTqNIK-80XGw2u8Z?usp=sharing>

The training CLIP V L/14 embeddings can be downloaded here: https://drive.google.com/file/d/1yeni0R4GqmTOFQ_GVw_x61ofZ-OBcS/view?usp=sharing (not fully manually annotated so cannot be used as test)



[4] "The BigScience RAIL license" Contractor et al.

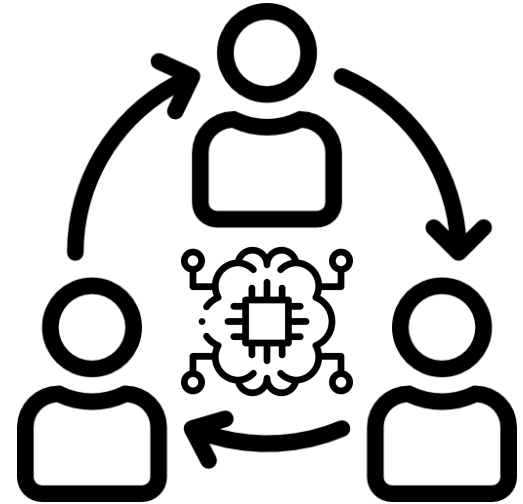


Trajectories of humanity with open source AI

Promoting open collaboration

Policymaking, strategic considerations and possible long-term outcomes:

- Openness → speeds up AI development, promotes wider engagement and transparency
- Effective enforcement when political institutions agree on these principles



[Strategic Implications of Openness in AI Development](#)



UNIVERSITÀ
DI TRENTO

Thank you!

Paper #10



Authors:

Diego Calanzone (diego.calanzone@studenti.unitn.it)

Andrea Coppari

Riccardo Tedoldi

Giulia Olivato

Carlo Casonato