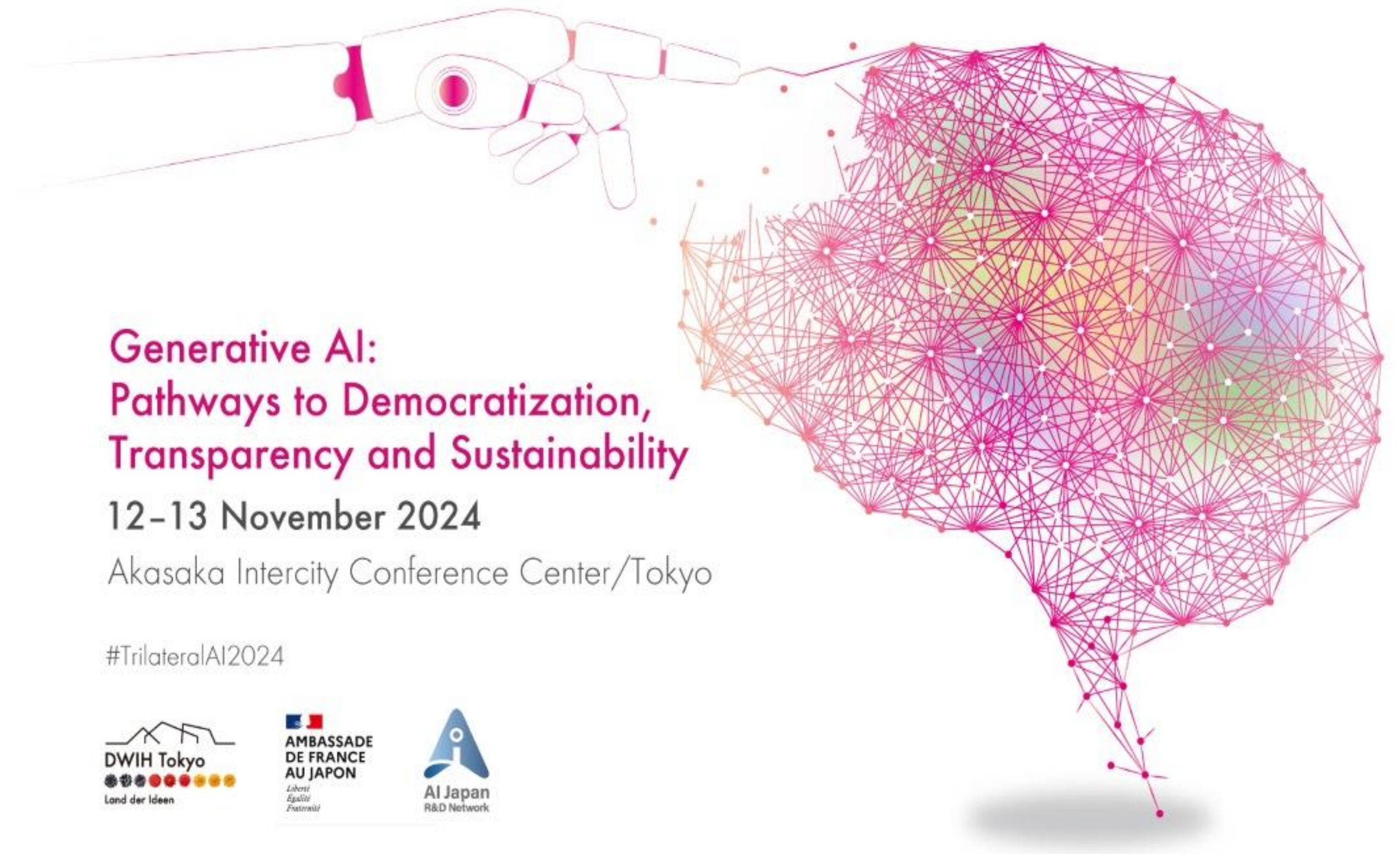
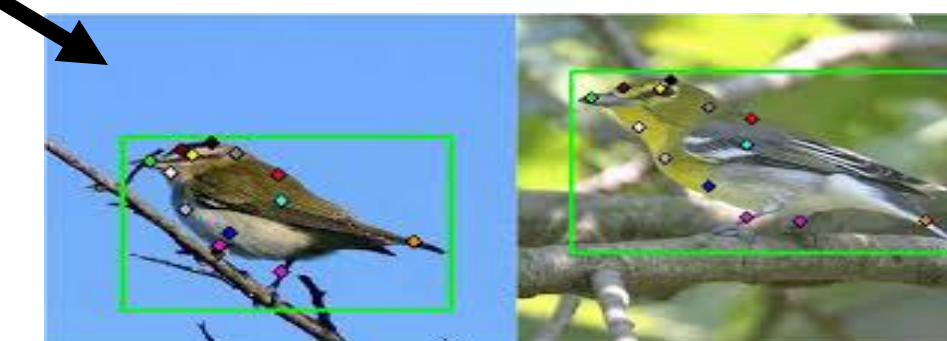
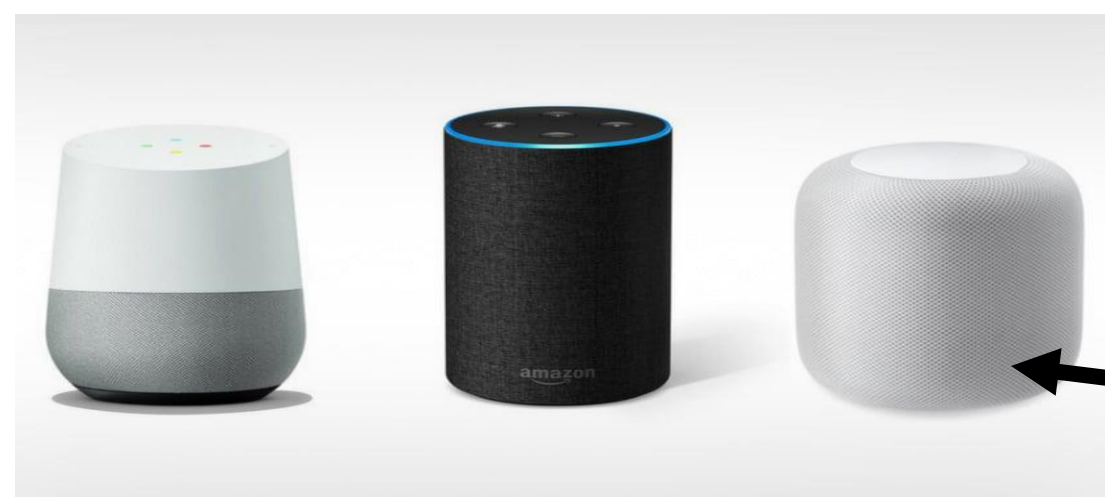
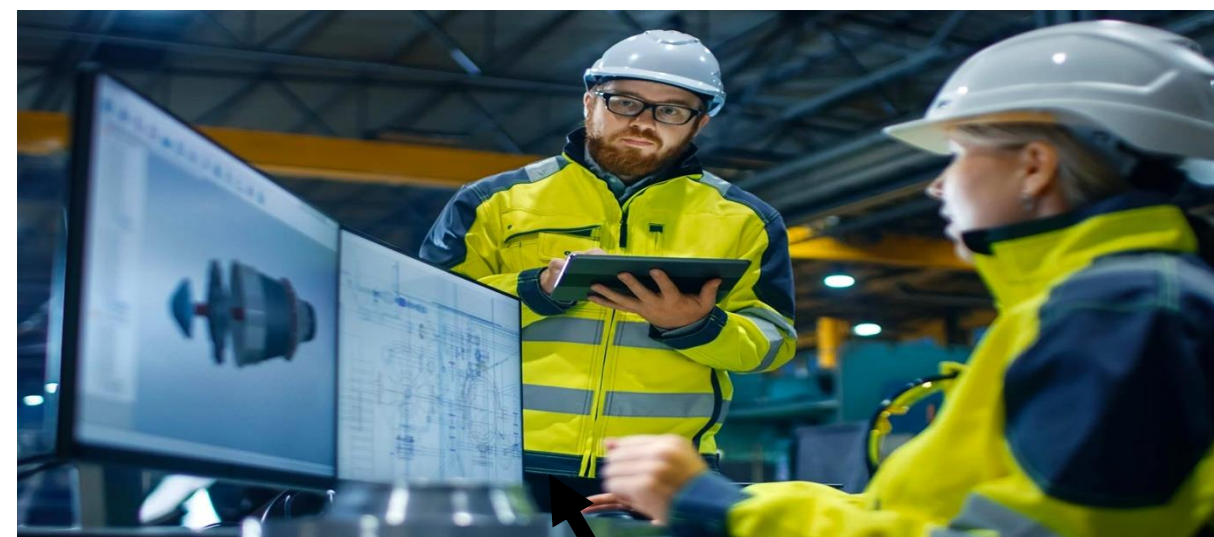


Generative and General-purpose AI : Strenghts and Weaknesses for AI Democratization

Céline Hudelot
Professor
MICS Laboratory
CentraleSupélec



AI is now everywhere !



Democratization of AI

Four meanings and goals [Seger et al, 23]

Democratization of AI Use

« Making it easier for a wide range of people to access and use the technology, without coding experience »

Democratization of AI Development

« Helping a wider range of people contribute to AI design and development processes »

Democratization of AI Profit

« Facilitating the broad and equitable distribution of value accrued to organisations that build and control advanced AI capabilities »

Democratization of AI Governance

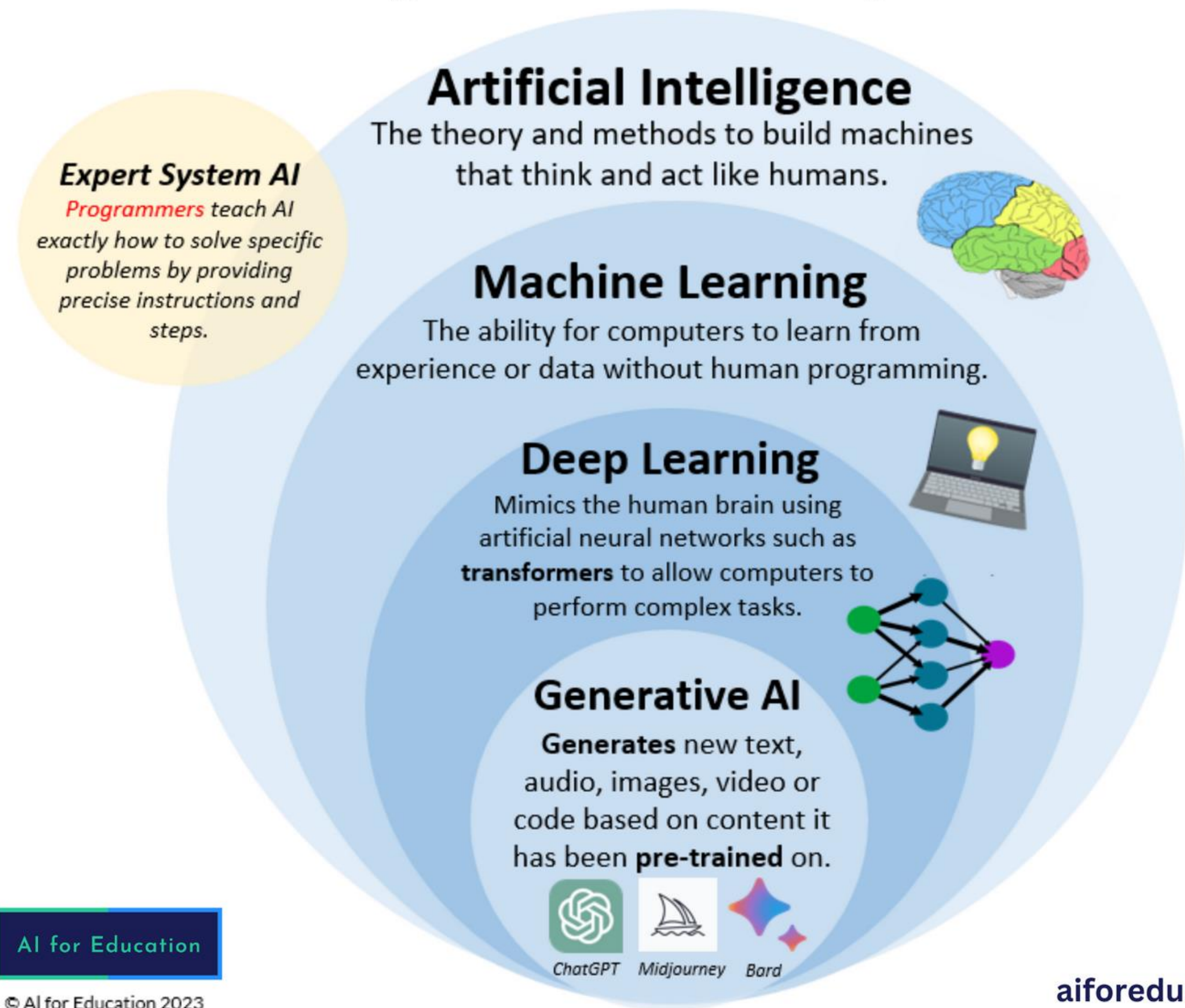
« Ensure that decisions around questions such as usage, development, and profits reflect the will and the preference of the society »

Why Generative AI and General-purpose AI can be considered as Pathways to the Democratization of AI Use and AI Development ?

The Rise of the Generative AI Paradigm

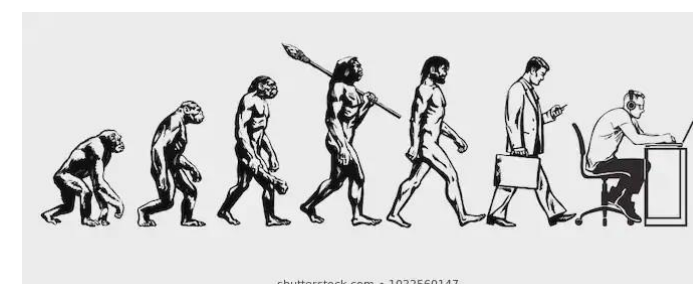
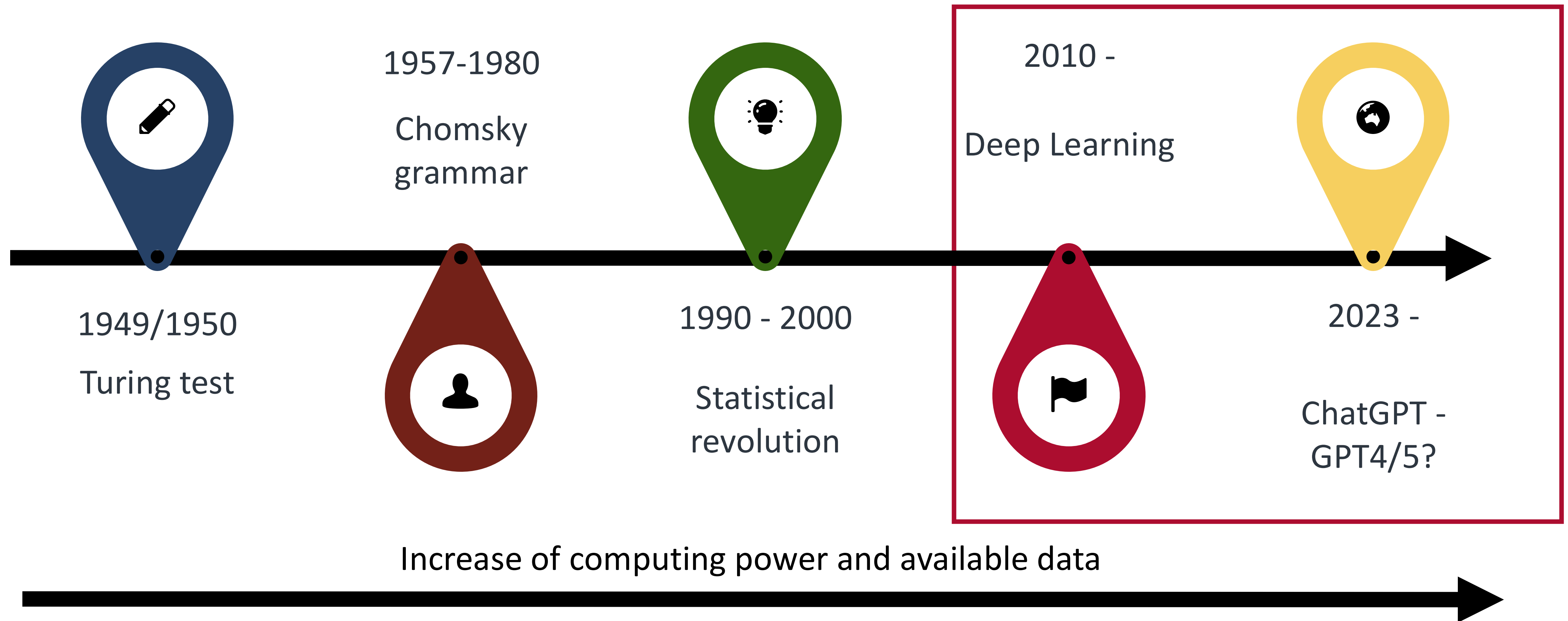
Defining Generative AI

To understand generative artificial intelligence (GenAI), we first need to understand how the technology builds from each of the AI subcategories listed below.



aiforeducation.io

The Rise of the Generative AI Paradigm

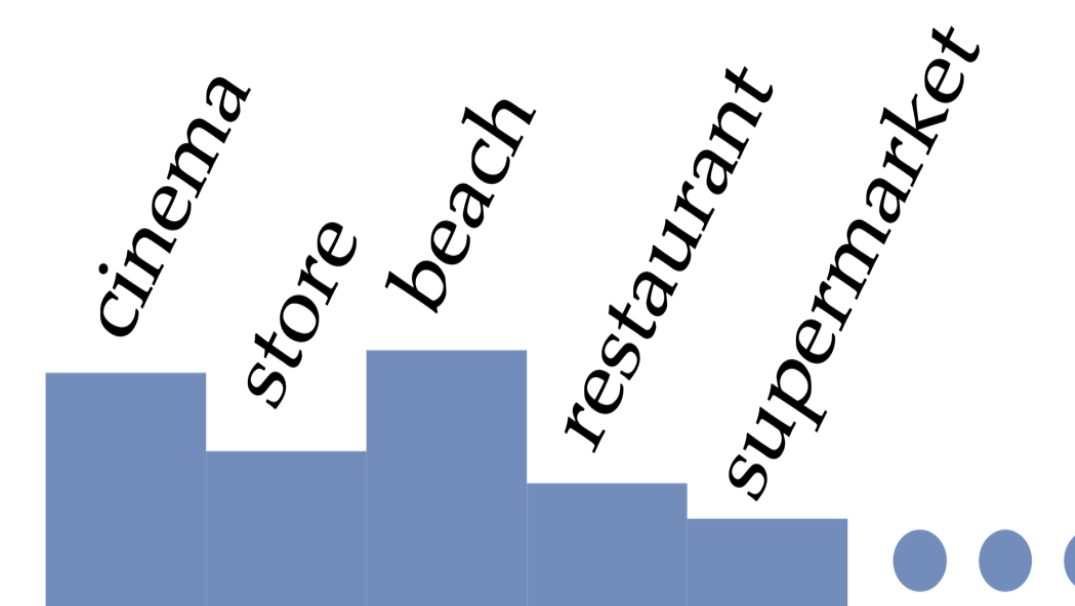
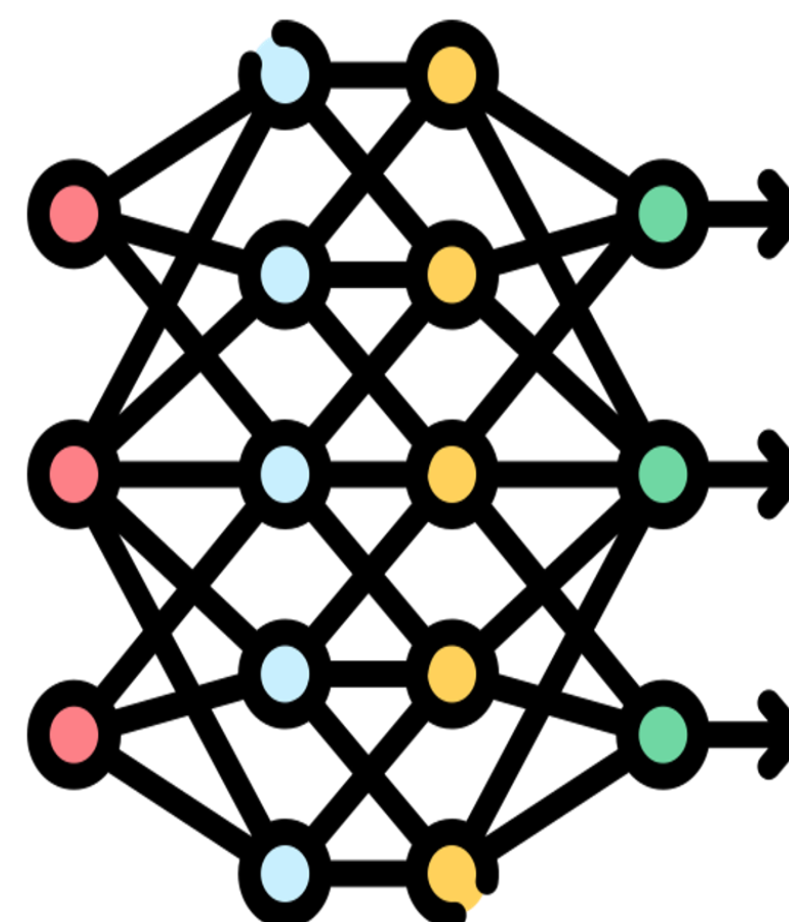


Generative AI : Large Language Models

What is a Large Language Model (LLM)?

A LLM is a neural network designed to process and generate human text.

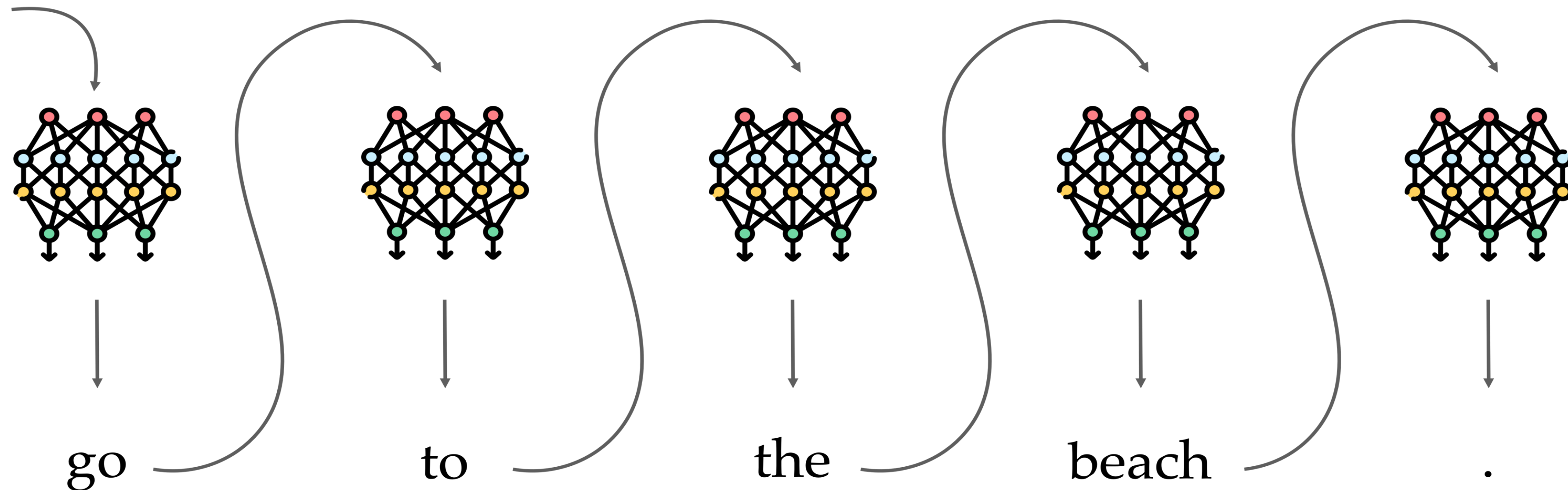
“I want to go to the”



Generative AI : Large Language Models

How does a LLM generate text?

I want to

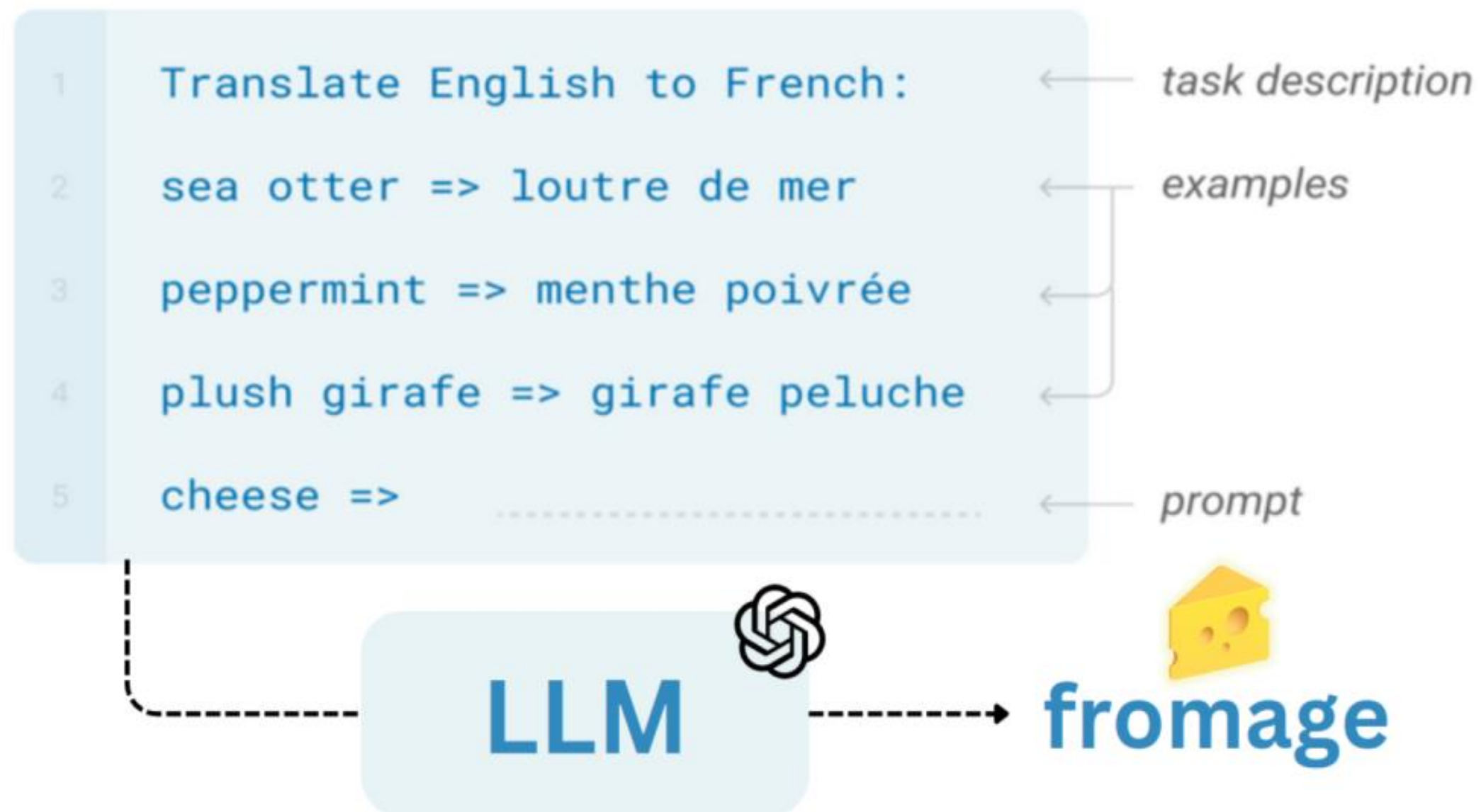


This is a very generic framework. How can we leverage this to perform multiple tasks?

Generative AI : LLMs

How can we leverage LLMs to perform many tasks?

Key ingredients : The era of In-context learning and prompt-engineering !



Prompt : input to a Generative AI model, that is used to guide its output.

In-context learning : paradigm that allows language models to learn tasks given only a few examples in the form of demonstration.

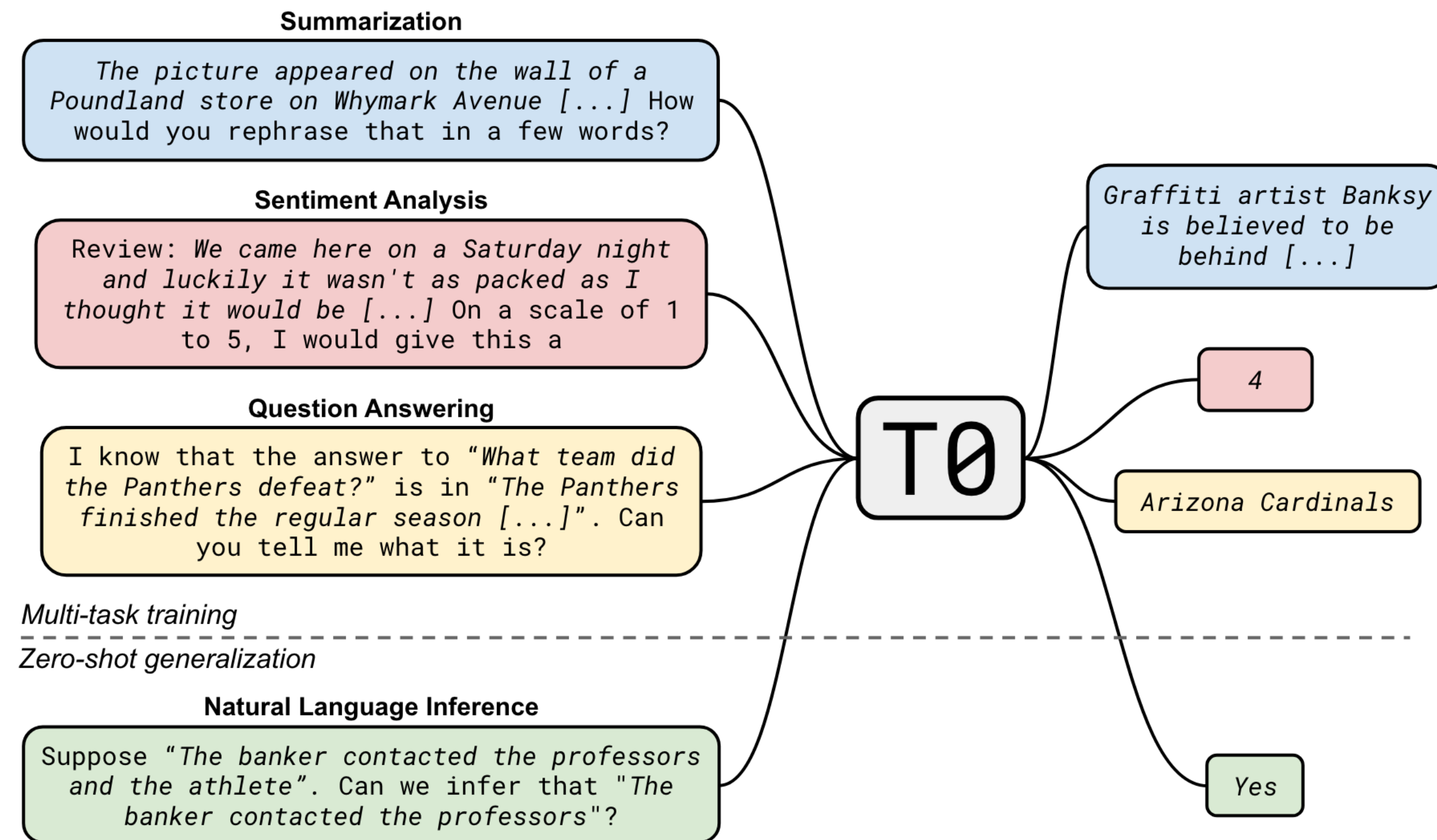
Prompt-engineering : developing and optimizing prompts to efficiently use language models (LMs) for a wide variety of applications and research topics.

Brown et al., "Language models are Few-Shot Learners"

Generative AI : LLMs

How can we leverage LLMs to perform many tasks?

We can prompt them, using **natural language prompts**, for a **large set of tasks**.
We can **dialog** with them !



Prompting and In-context learning enable users and developers to use AI easily, only by prompting and for a large set of tasks.

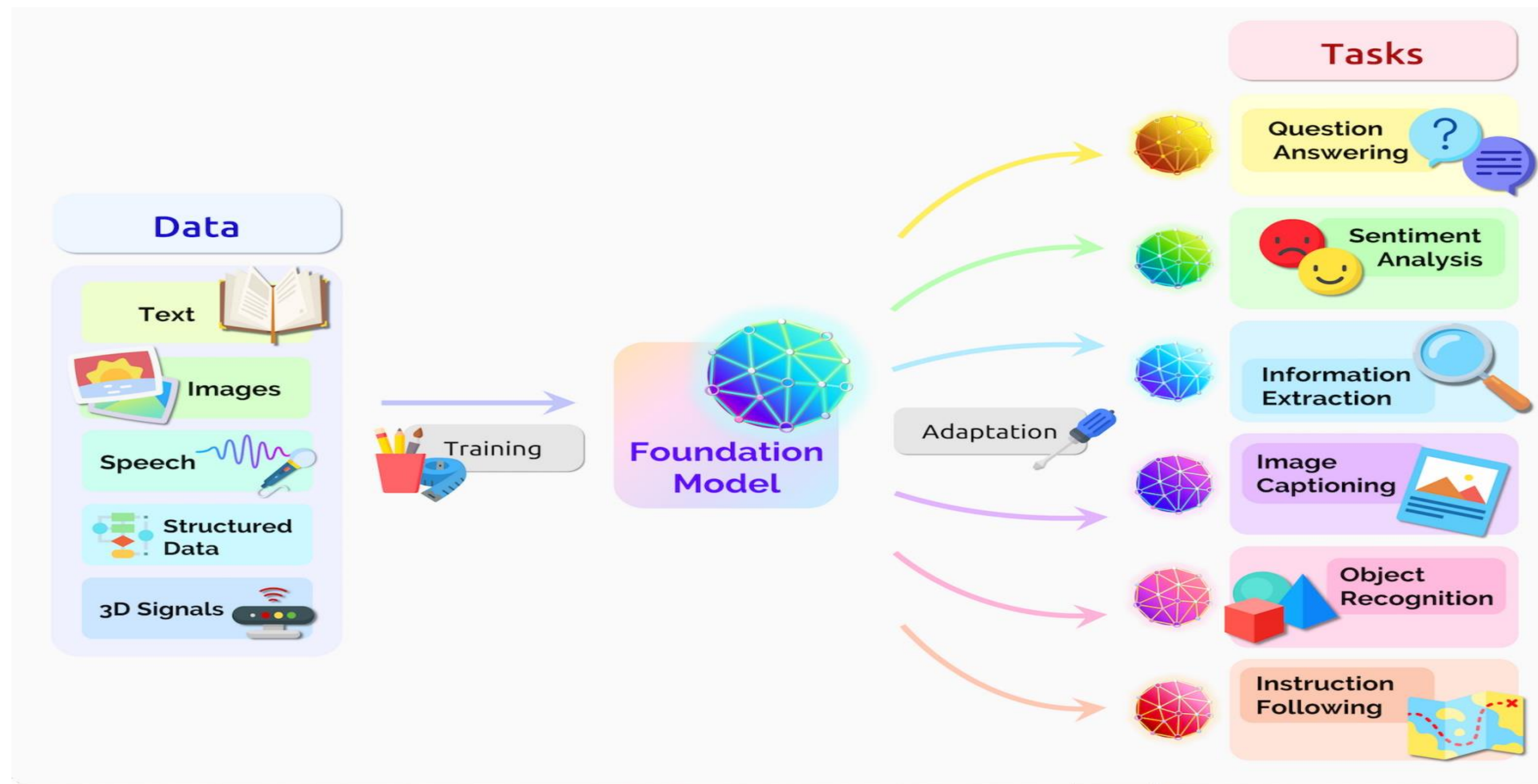
An important key step to the Democratization of AI Use and AI development.

Sahn et al., "Multitask Prompted Training Enable Zero-Shot Task Generalization"

Generative AI : the era of Foundations models

Not only for textual data. LLMs are just Foundations models

A new major paradigm for building AI system : train a model on broad data and adapt it to a wide range of downstream tasks.



Foundation models are **another key ingredients to the Democratization of AI**

Bommasani et al., "On the Opportunities and Risks of Foundation Models"

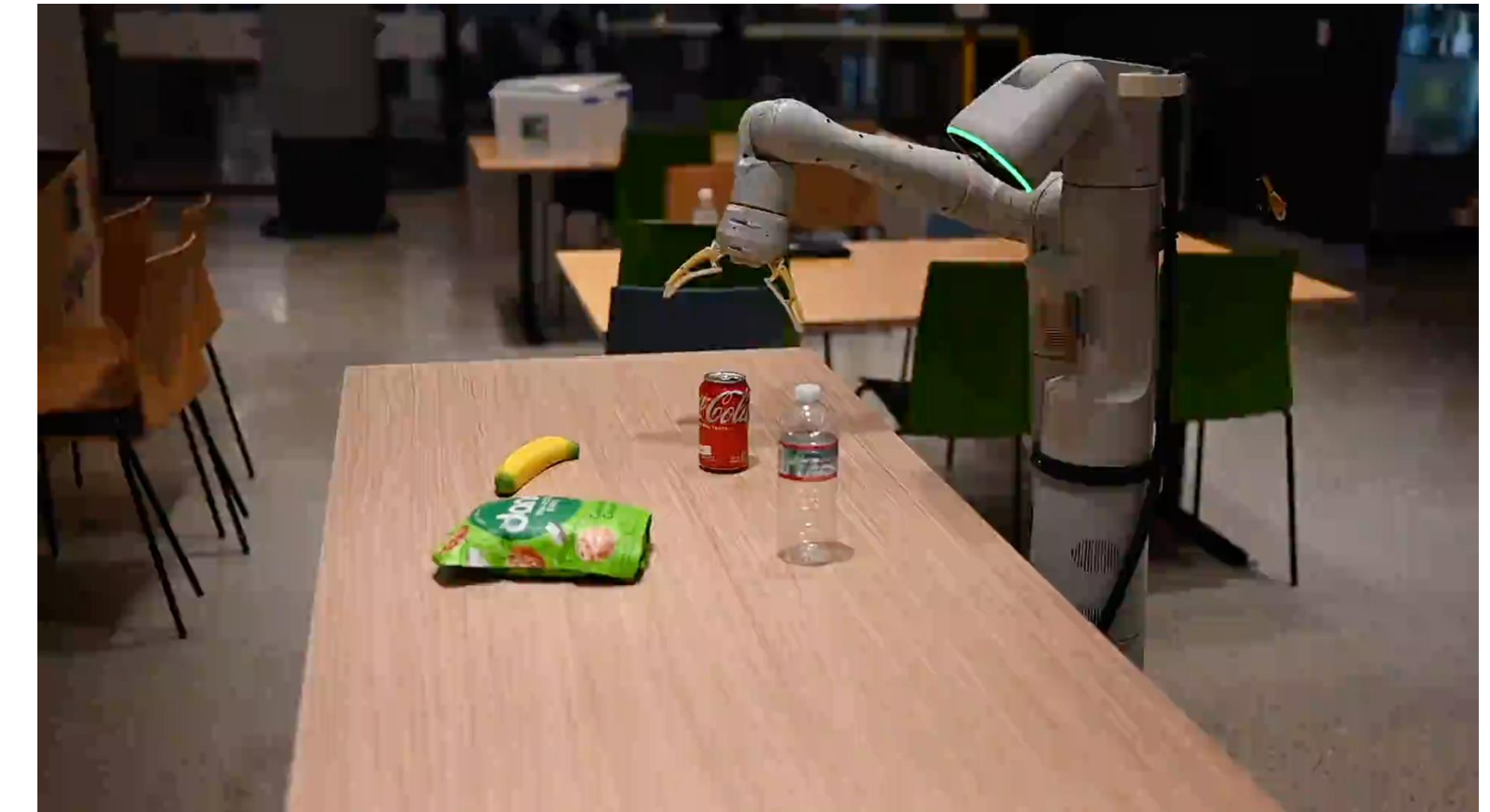
Generative AI : the era of Foundations models

Examples of well-known foundation models

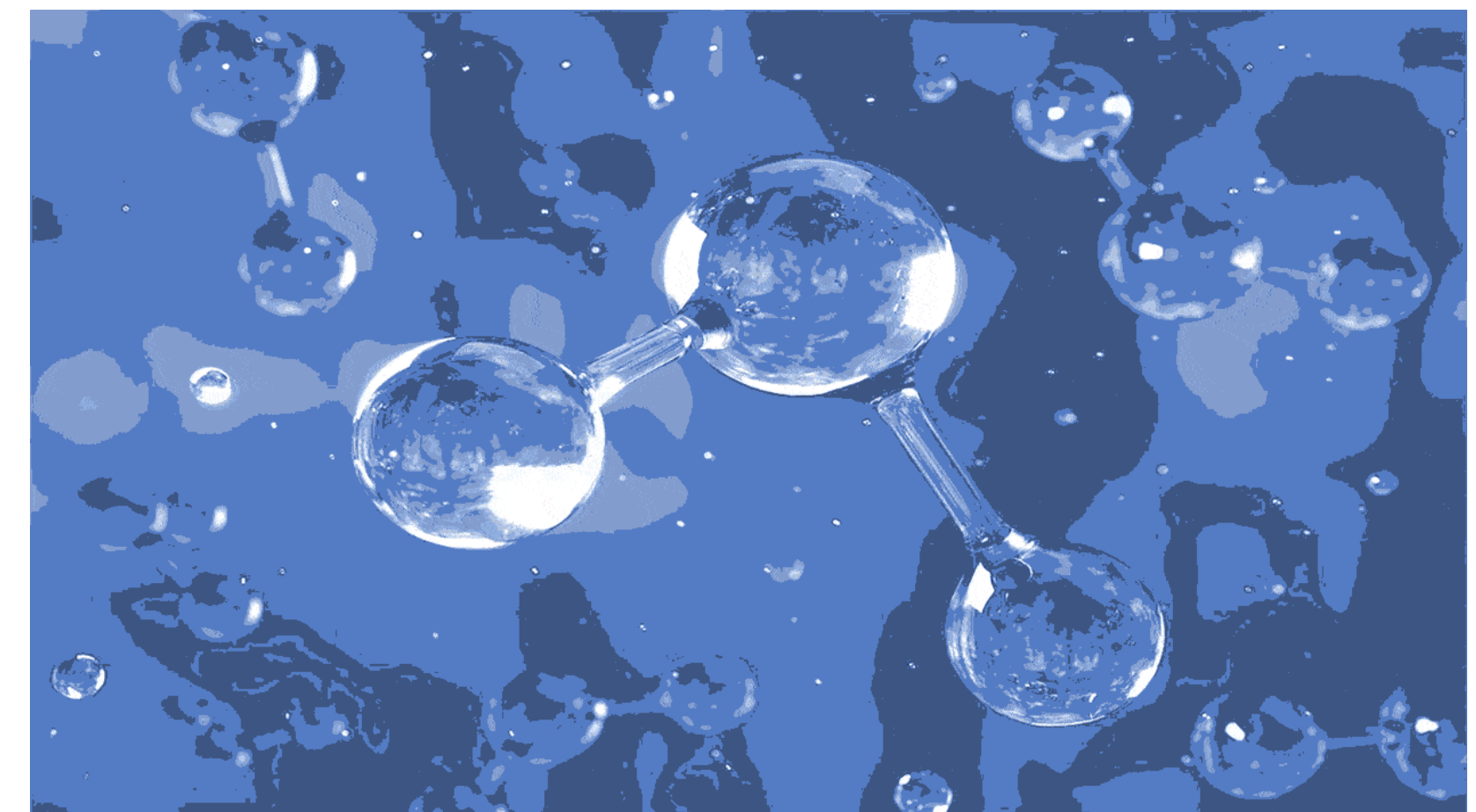
Segment Anything by Meta AI



RT-2: Vision-Language-Action Models by Google Deep Mind



MoLFormers, foundation models on chemicals by IBM

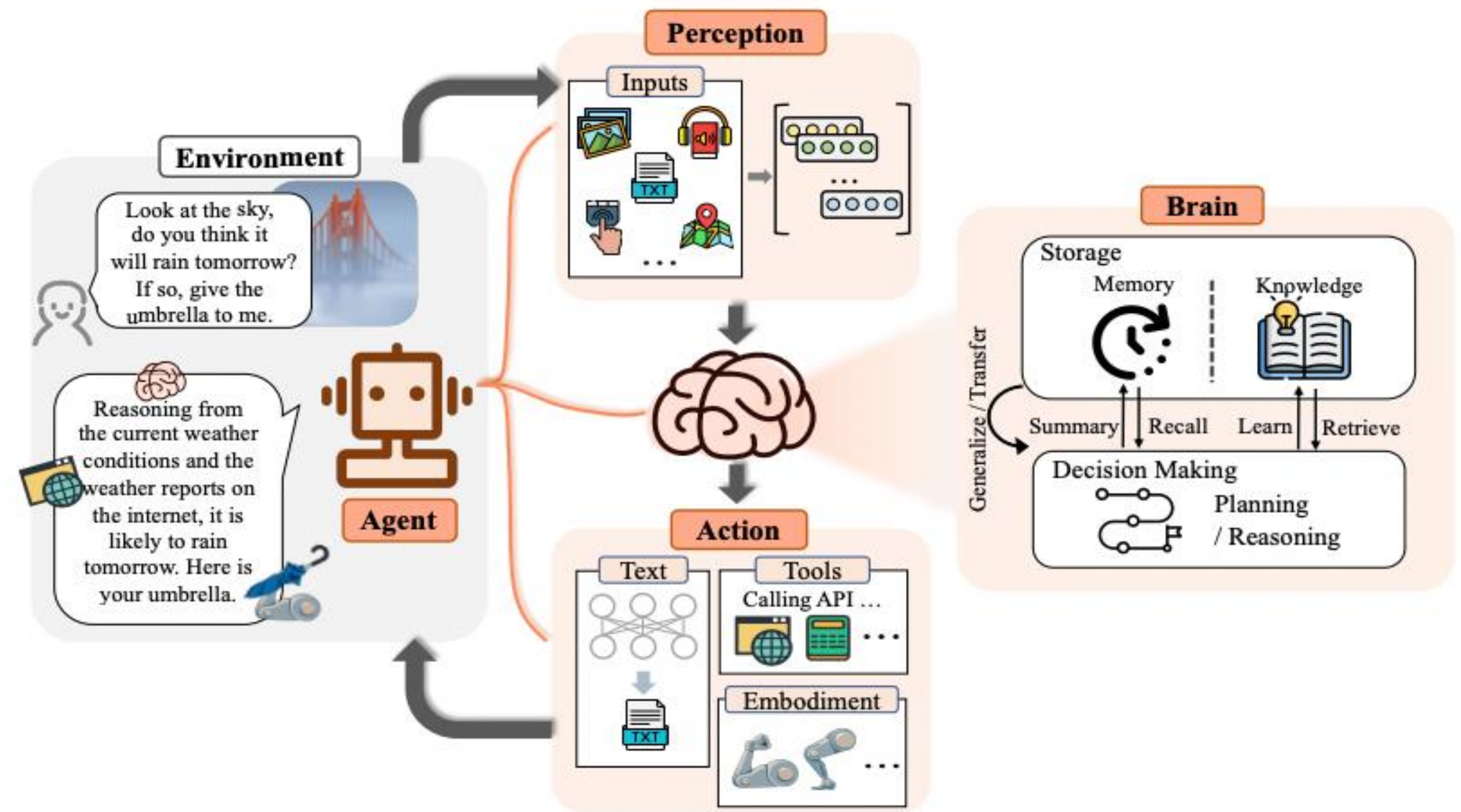


Generative AI : a new era for AI agents.

AI agent : can gather useful sensory information and interact with their environments in meaningful ways

Enable the orchestration of **complex pipelines** with LLM and foundations models applications executing complex tasks through the use of an architecture combining them with key modules like planning and memory.

AI Agent paradigm is also a key ingredient to **the Democratization of AI (use and development)**



Generative AI : a new era for AI agents.

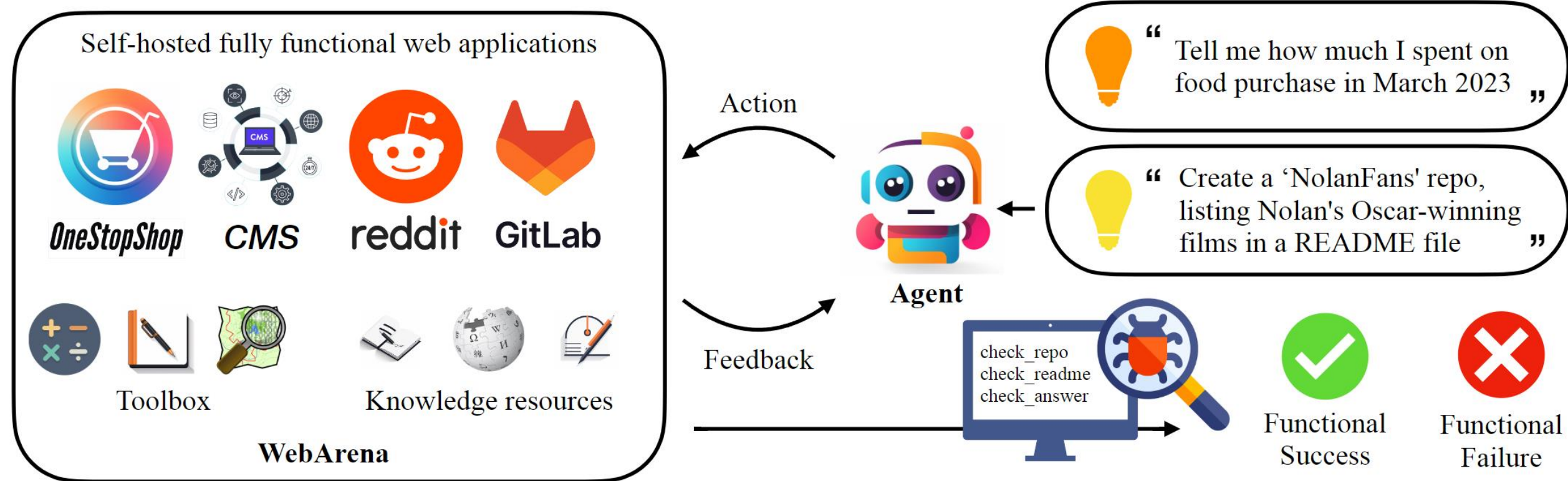
AI agent : A route towards Artificial General Intelligence (AGI)

Foundation Models for Decision Making @NeurIPS 2023

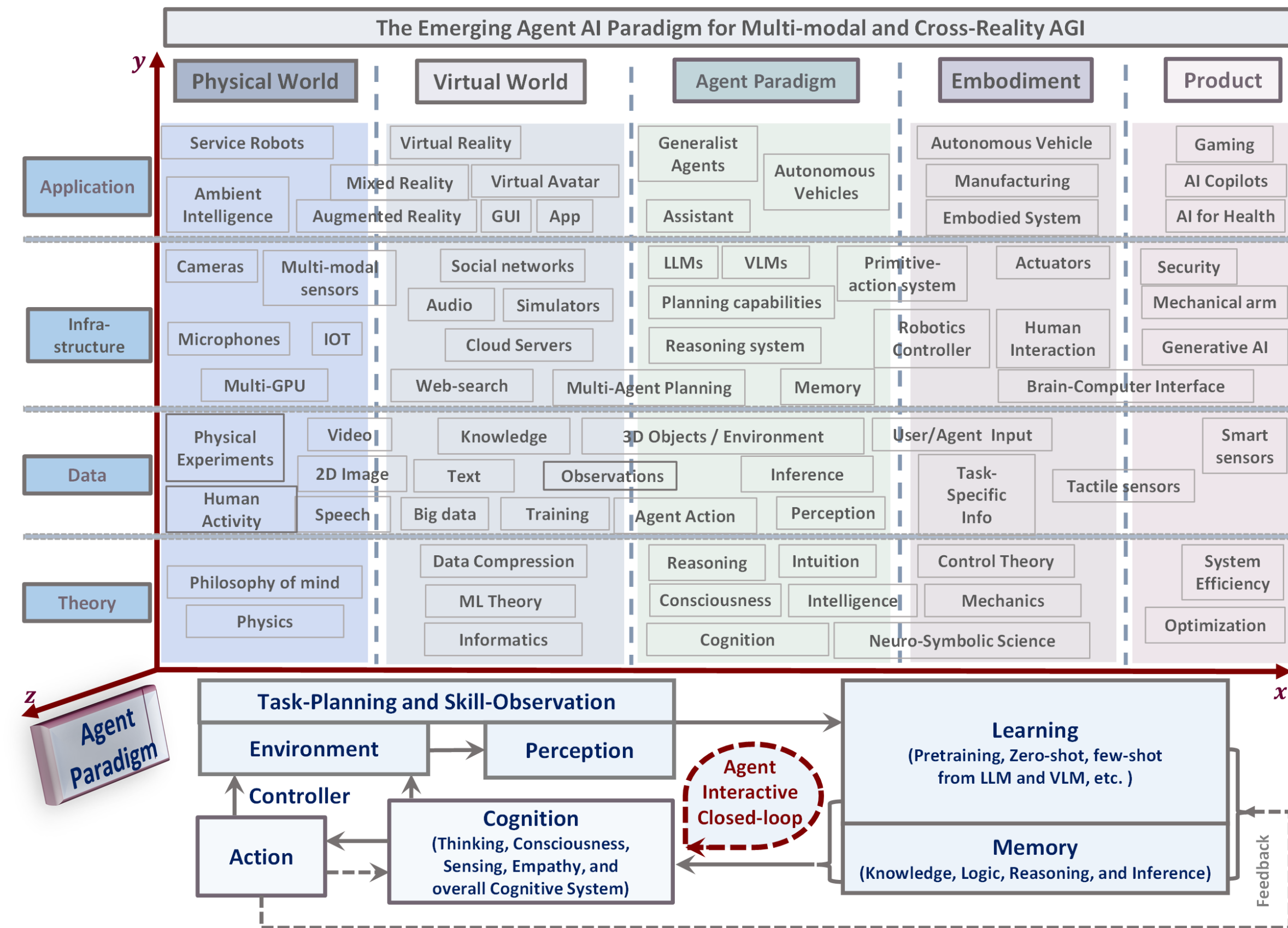
Hall E2, December 15, 2023 (Friday) @8:15 am

December 15, 2023 (Friday)

Foundation Models and Decision Making come together to solve complex tasks at scale.



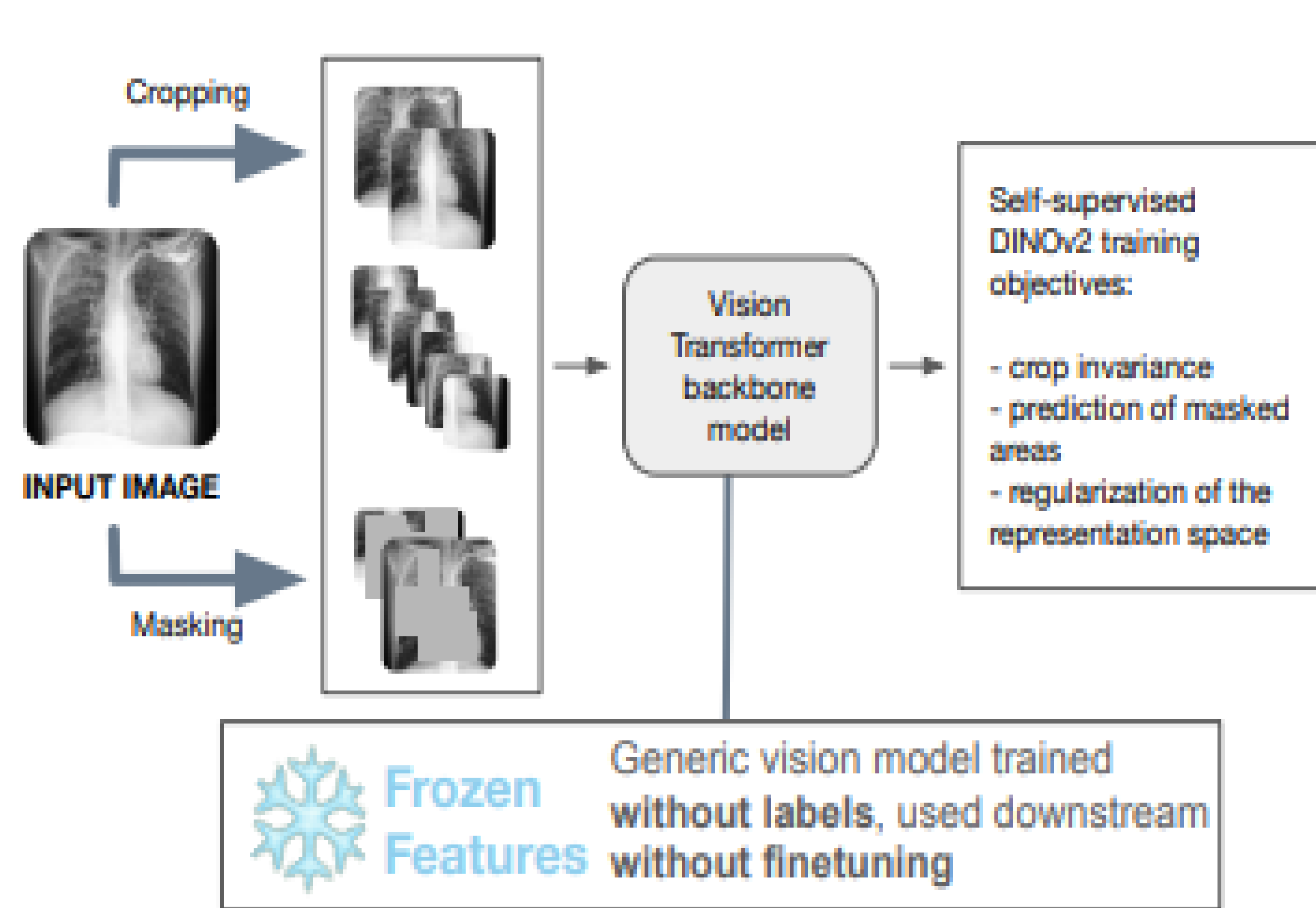
Zhou et al., "WebArena: A Realistic Web Environment for Building Autonomous Agents"



Durante et al., "Agent AI: Surveying the Horizons of Multimodal Interaction"

Use case study with a health foundation model

RayDINO, a X-ray foundation model

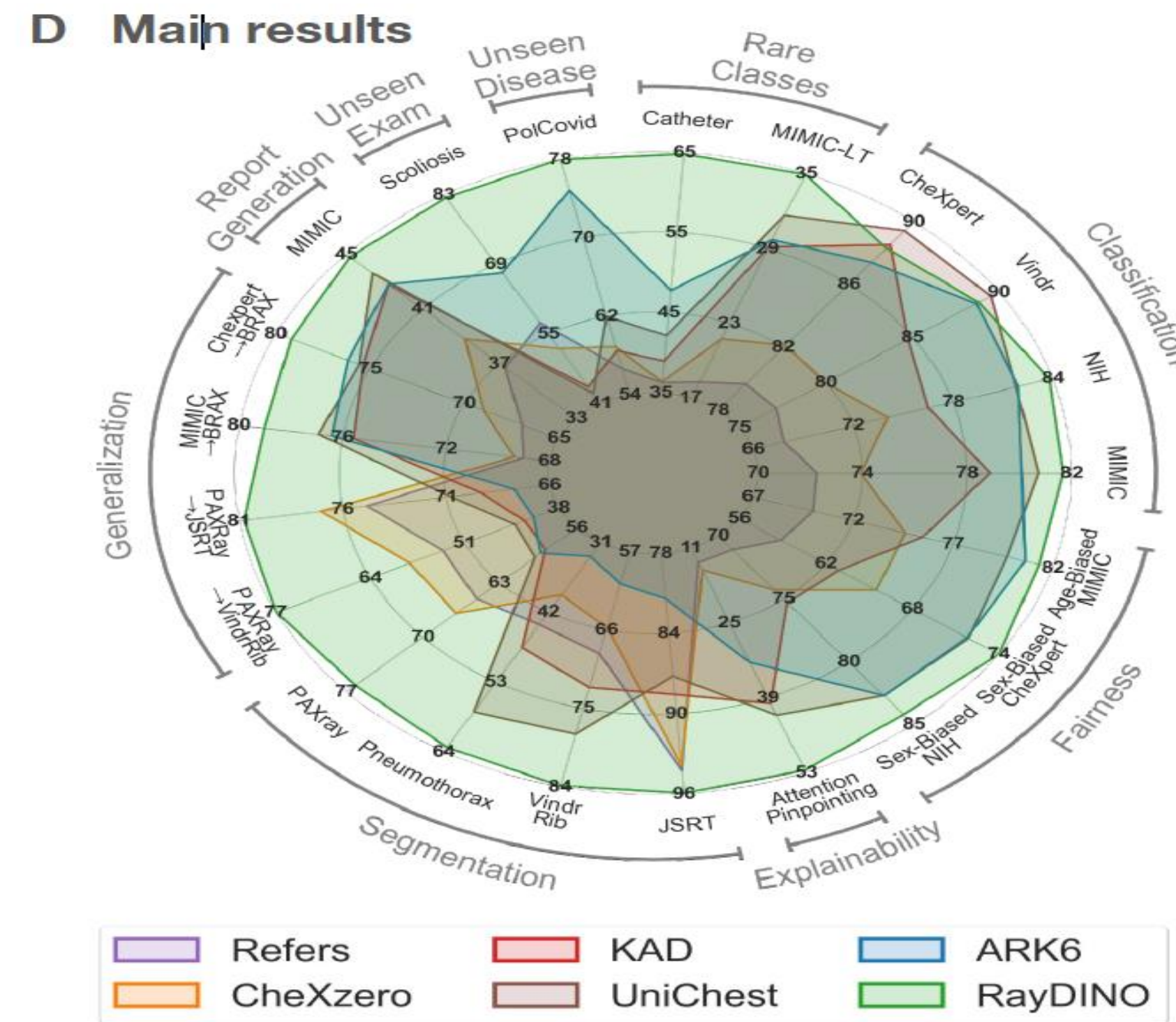
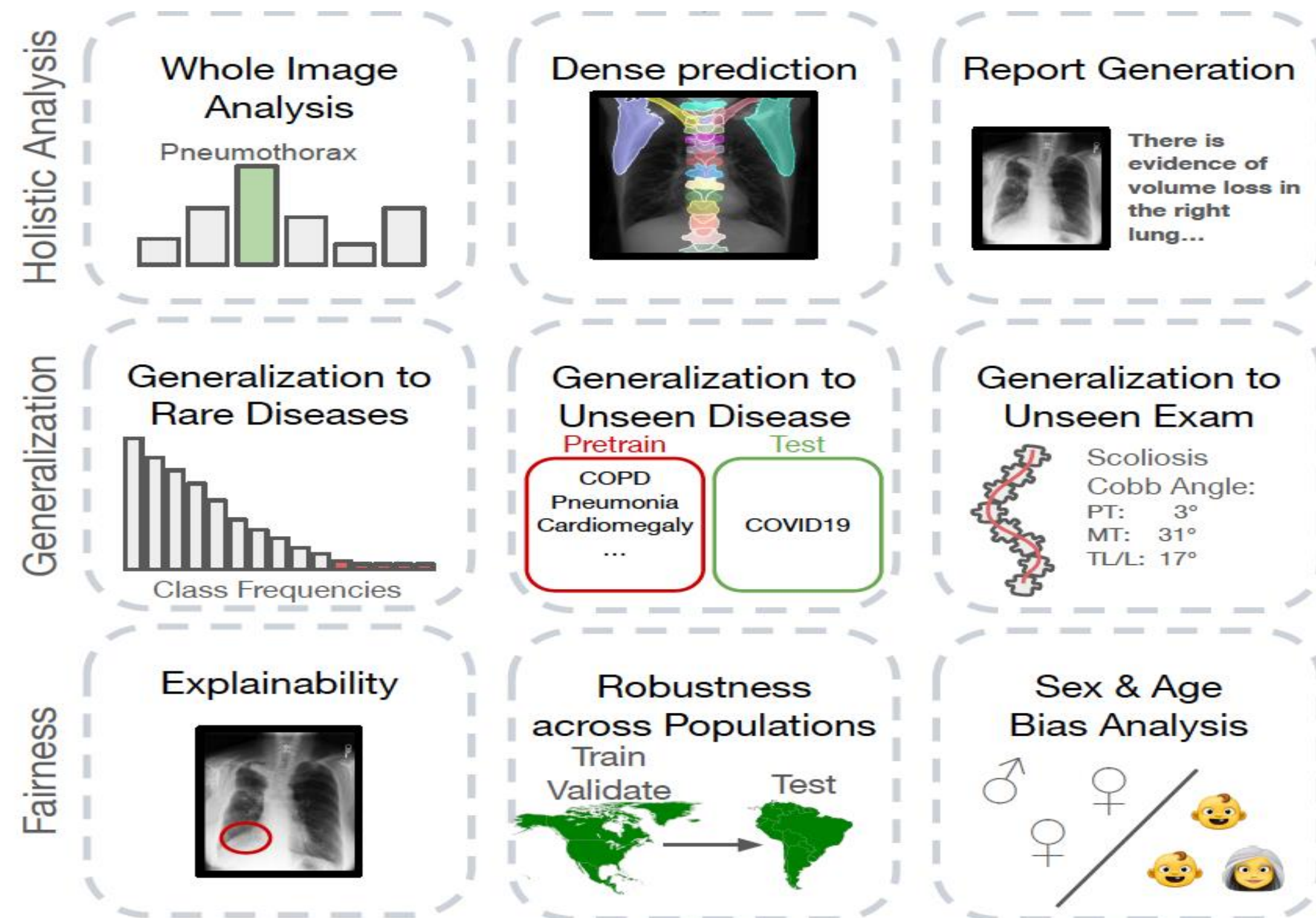


Pretraining / Evaluation			
MIMIC 377k imgs 31k patients 14 labels USA	CheXpert 224k imgs 65k patients 14 labels USA	NIH 112k imgs 31k patients 15 labels USA	PadChest 160k imgs 67k patients 19 labels Spain
Evaluation			
VinDr 18k imgs 28 labels Vietnam	PAXRay++ 15k img 157 labels Multiple	Scoliosis 609 img 3 labels Canada	Pneumothorax 12k img 1 label USA
PolCovid 4k imgs 3 labels Poland	JSRT 154 img 5 labels Japan	VinDr Rib 245 img 20 labels Vietnam	BRAX 41k imgs 19k patients 14 labels Brasil

Moutakanni et al., "Advancing human-centric AI for robust X-ray analysis through holistic self-supervised learning »

Use case study with a health foundation model

RayDINO, a X-ray foundation model



Moutakanni et al., "Advancing human-centric AI for robust X-ray analysis through holistic self-supervised learning »

Use case study with a health foundation model

RayDINO, a X-ray foundation model

- Unified Approach for Comprehensive Radiology Interpretation: Classification, Segmentation, and Report Generation.
 - ✓ Make easier AI use and development.
- Ability to deal with rare health conditions that cannot be annotated on thousands of images.
 - ✓ A crucial assessment for the practical utility in routine clinical practice.
 - ✓ Potential for large-scale deployment.
- Out-of-Domain Task Performance (POLCOVID Dataset)
 - ✓ Potential for Exploring New Diseases and Exams.
- Robustness to Patient Demographics and their social origins
 - ✓ AI for all
- We also perform gender bias analysis
 - ✓ Ensure fairness and reliability across diverse applications and populations.

Generative AI and General-purpose AI can be considered as Pathways to the Democratization of AI Use and AI Development

but

Accessibility and Ease of Use do not mean Adoption.

(a) Factuality Hallucination

Who was the first person to walk on the moon?

Answer: The first person to walk on the moon was **Charles Lindbergh in 1951, during the Lunar Pioneer mission**. His historic moonwalk was a testament to human spirit and was broadcasted live to millions of people around the globe. ❌

Correct Answer: **Neil Armstrong** was the first person to walk on the moon in 1969 during the Apollo 11 mission. ✅

(b) Faithfulness Hallucination

Please summarize the following news article:

Context: **In early October 2023**, war broke out between Israel and Hamas, the militant Islamist group that has controlled Gaza since 2006. Hamas fighters fired rockets ... civilians and taking dozens of hostages.

Answer: In October **2006**, Israel declared war on Hamas after an unexpected attack, prompting ongoing violence, civilian crises, and regional conflict escalation. ❌

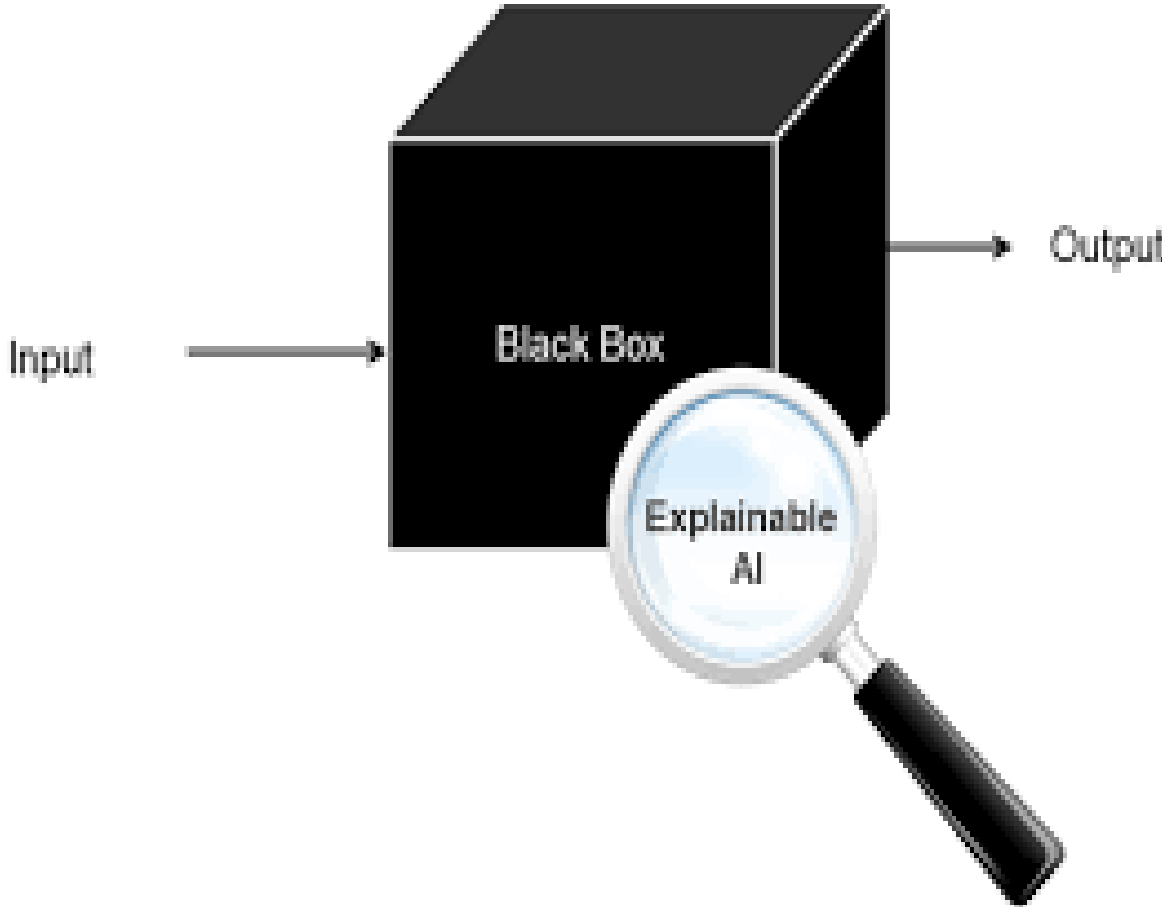
Hallucinations



Press release >

Generative AI: UNESCO study reveals alarming evidence of regressive gender stereotypes

Bias



Explainability, Transparency

Generative AI and General-purpose AI can be considered as Pathways to the Democratization of AI Use and AI Development

but

Accessibility and Ease of Use do not mean Adoption.



Trustworthy AI is needed

**Next session :
transparency and
explainability**

Others issues for the Democratization of AI Development

« Helping a wider range of people contribute to AI design and development processes »

Goals

- Accelerate AI Innovation and Progress
- Enabling more people to participate in AI design and development : cater diverse needs and interests
- External evaluation

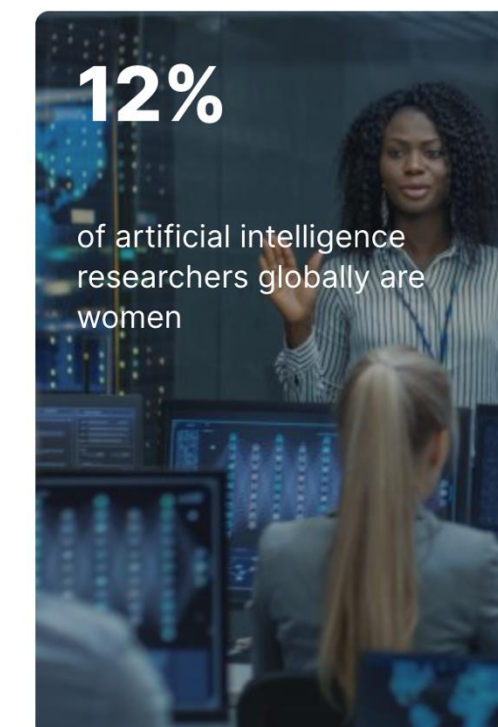
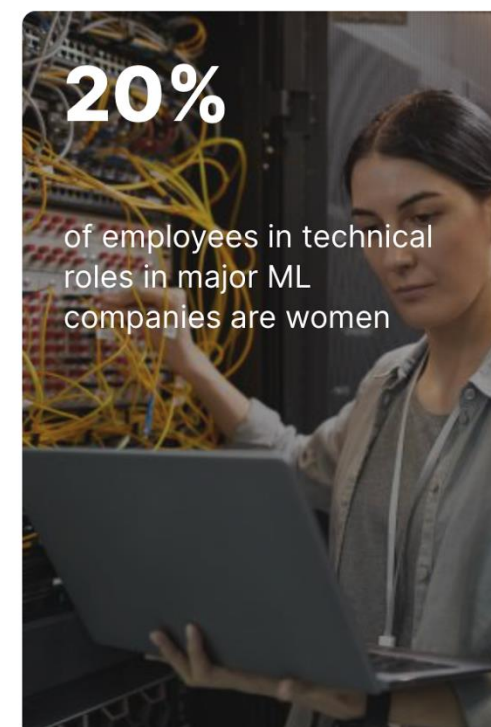
Seger et al., “Democratising AI: Multiple Meanings, Goals, and Methods”

Some issues



Big tech monopole

Figures of Women in AI



Source : Unesco

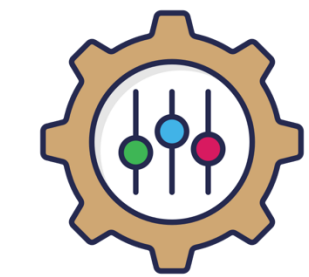
Huge models downsides



EFFICIENCY



COST

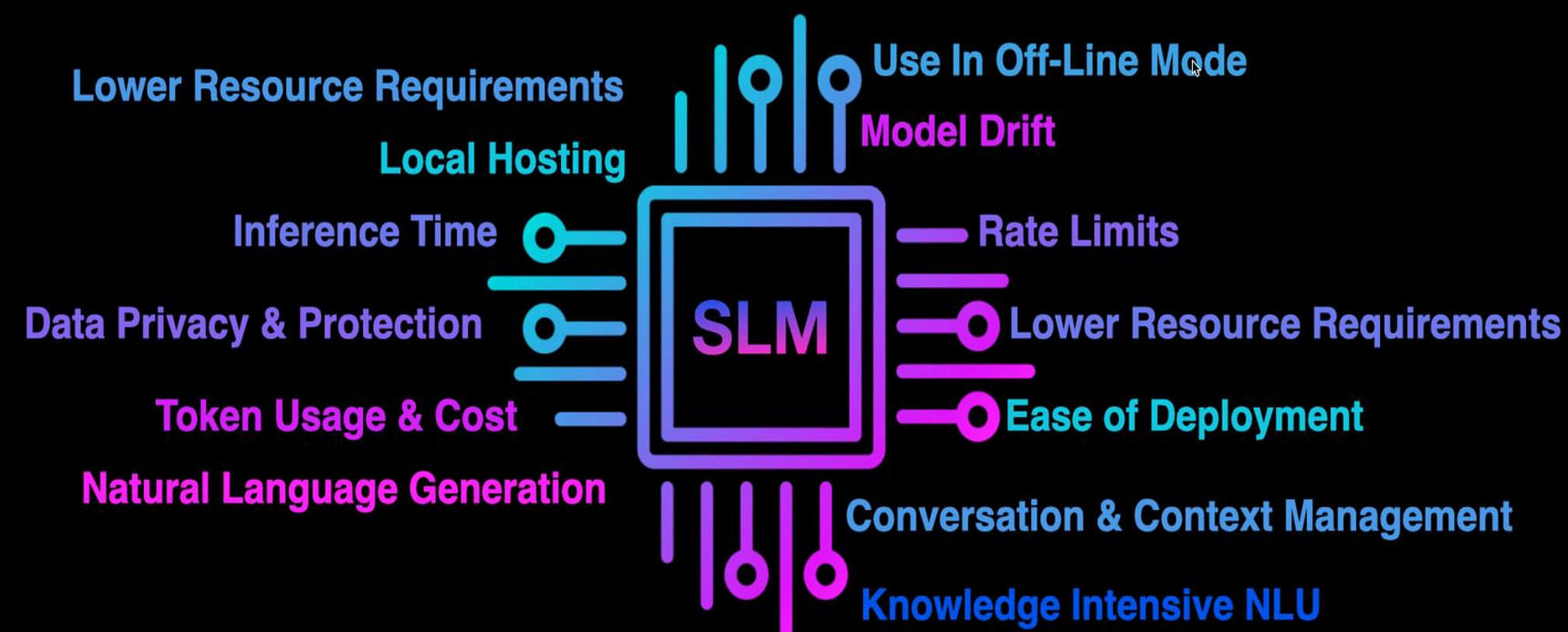


CUSTOMIZATION

Gender bias,
narrow
demographic
actors

Small Language Models (SLMs) at the rescue

SLM = Small Language Model



www.cobusgreyling.com

Small Languages Models (SLMs) :

- Less than 5 billions parameters
- Retain accuracy and/or adaptability of LLMs while being subject to constraints:
 - Training or inference hardware
 - Data availability
 - Bandwidth
 - Generation time

SLM example : CROISSANT LLM

An Industrial and Academic Partnership for a or a sovereign LLM



CroissantLLM: A Truly Bilingual French-English Language Model

Manuel Faysse^{1,5} Patrick Fernandes^{6,8,11} Nuno M. Guerreiro^{2,5,6,8}
António Loison¹ Duarte M. Alves^{6,8} Caio Corro⁹ Nicolas Boizard^{4,5}
João Alves² Ricardo Rei^{2,7,8} Pedro H. Martins² Antoni Bigata Casademunt¹⁰
François Yvon⁹ André F.T. Martins^{2,6,8} Gautier Viaud¹ Céline Hudelot⁵
Pierre Colombo^{3,5}



Manuel Faysse
CentraleSupélec, MICS
PhD student



Pierre Colombo
CentraleSupélec, MICS
Assistant Professor



Céline Hudelot
CentraleSupélec, MICS
Professor



Nicolas Boizard
CentraleSupélec, MICS
Diabolocom
PhD Student



A research projet with industrial aims



Research

CroissantLLM is a research project aiming to study how **bilingualism** impacts language model pretraining and performance.



Industry

The final model is designed to be **small enough** to run on local hardware, but **good enough** to run generative tasks that are often reserved to larger models (**inference-optimal** training). It is trained on **permissively licensed data** only.



Open-Source

This project is rooted in open-source, with openly released models, data, code bases and evaluation benchmarks, enabling **researchers and practitioners** to benefit from it.

The Model (Chinchilla tradeoffs)

Generative models are often transformer *decoders* (GPT, Mistral, LLaMa) which performance is closely related to (1) the # of model parameters, and (2) the # of training tokens.

Training the best model given a **fixed compute budget**

For a given compute budget, there is an optimal ratio between parameter count and training data size (~20 for **Chinchilla Scaling laws**)

OR

Training the best model of a **fixed size**

By training longer than the Chinchilla ratio, we continue to improve the model but performance gains are increasingly costly.

We decide to overtrain a small model (1,3B) → 2307 token:param ratio vs 20 Chinchilla-optimal.



Lighter



Faster



Capable



Training Cost

The data

Training language models requires huge amounts of data... In French and under permissive licenses, gathering sufficiently is even more challenging !

Corpus



Culture

Permissively licensed data source

Public domain books

Podcasts

Poetry

Song lyrics

Movie subtitles



Business

Industrial and Administrative data

Legal corpora

Parliamentary debates

Administrative Decisions

Public business documents



Knowledge

Scientific and Factual data

Encyclopedia

Textbooks

Theses abstracts

Scientific publications



Translations

English and French Parallel data

Huge quantity of translation pairs sourced from different domains

Filtered with SOTA quality estimation methods



Internet

Filtered web data

Web-scale data filtered to obtain high quality French and English texts

Github Code under open licenses

The data

Training language models requires huge amounts of data... In French and under permissive licenses, gathering sufficiently is even more challenging !

Corpus



Culture

Permissively licensed data source



Business

Industrial and Administrative data



Knowledge

Scientific and Factual data



Translations

English and French Parallel data



Internet

Filtered web data



Identification and scrapping

Deduplication

Filtering

Upsampling



License checks
Large scale download scripts
Data quality and diversity checks

The data

Training language models requires huge amounts of data... In French and under permissive licenses, gathering sufficiently is even more challenging !

Corpus



Culture

Permissively licensed data source



Business

Industrial and Administrative data



Knowledge

Scientific and Factual data



Translations

English and French Parallel data



Internet

Filtered web data



Identification and scrapping

Deduplication

Filtering

Upsampling



URL deduplication
Exact deduplication
Fuzzy deduplication
(MinHash LSH)

The data

Training language models requires huge amounts of data... In French and under permissive licenses, gathering sufficiently is even more challenging !

Corpus



Culture

Permissively licensed data source



Business

Industrial and Administrative data



Knowledge

Scientific and Factual data



Translations

English and French Parallel data



Internet

Filtered web data



Identification and scrapping

Deduplication

Filtering

Upsampling



Rule-based filtering
Filtering of Toxic, Violent or Political content
Perplexity Filtering
(Data Quality)

The data

Training language models requires huge amounts of data... In French and under permissive licenses, gathering sufficiently is even more challenging !

Corpus



Culture

Permissively licensed data source



Business

Industrial and Administrative data



Knowledge

Scientific and Factual data



Translations

English and French Parallel data



Internet

Filtered web data



Identification and scrapping

Deduplication

Filtering

Upsampling



Final data mix

Upsampling FR to obtain a balanced corpus FR / GB

Transparency & Open-Source

Project rooted in transparency, to serve as a useful resource for industrial practitioners and researchers !



Documented training process from beginning to end



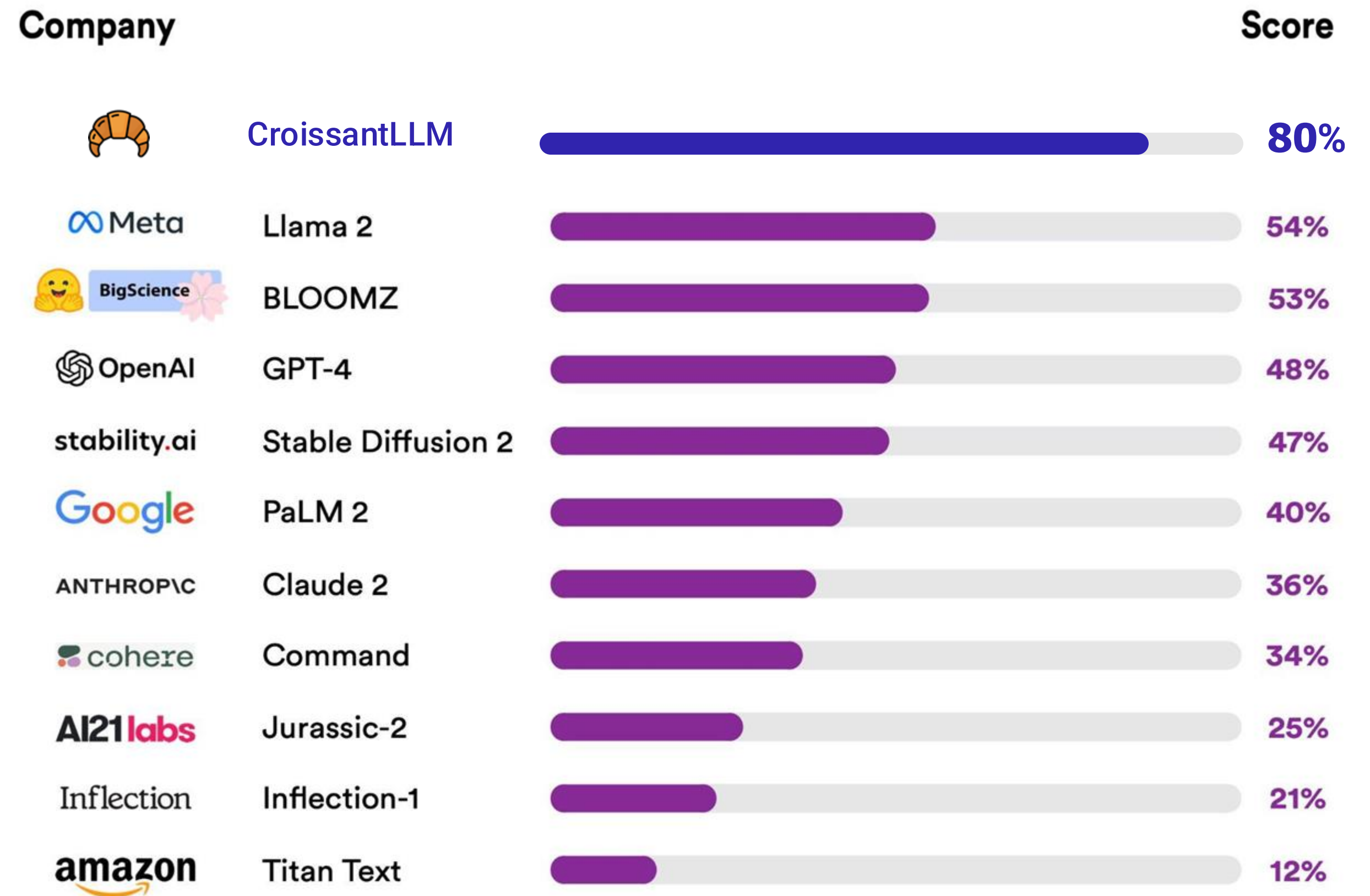
Openly available:

- Training corpus
- Model checkpoints
- Evaluation Benchmarks
- Code bases



No usage restrictions (MIT)

Foundation Model Transparency Index Total Scores

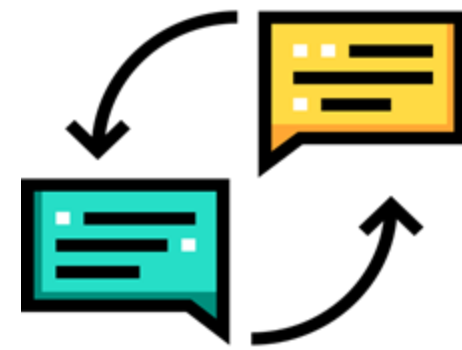


Croissant LLM applications

Specific tasks



- Writing assistance
- Summarization
- Orthographic correction
- Prompt Compression (RAG)
- Retrieval



Translation

CroissantLLM is the best decoder of its size in translation ^{FR} / ^{GB} matching the performance of Mistral and Llama models of 10 x the size.

Phone & CPU



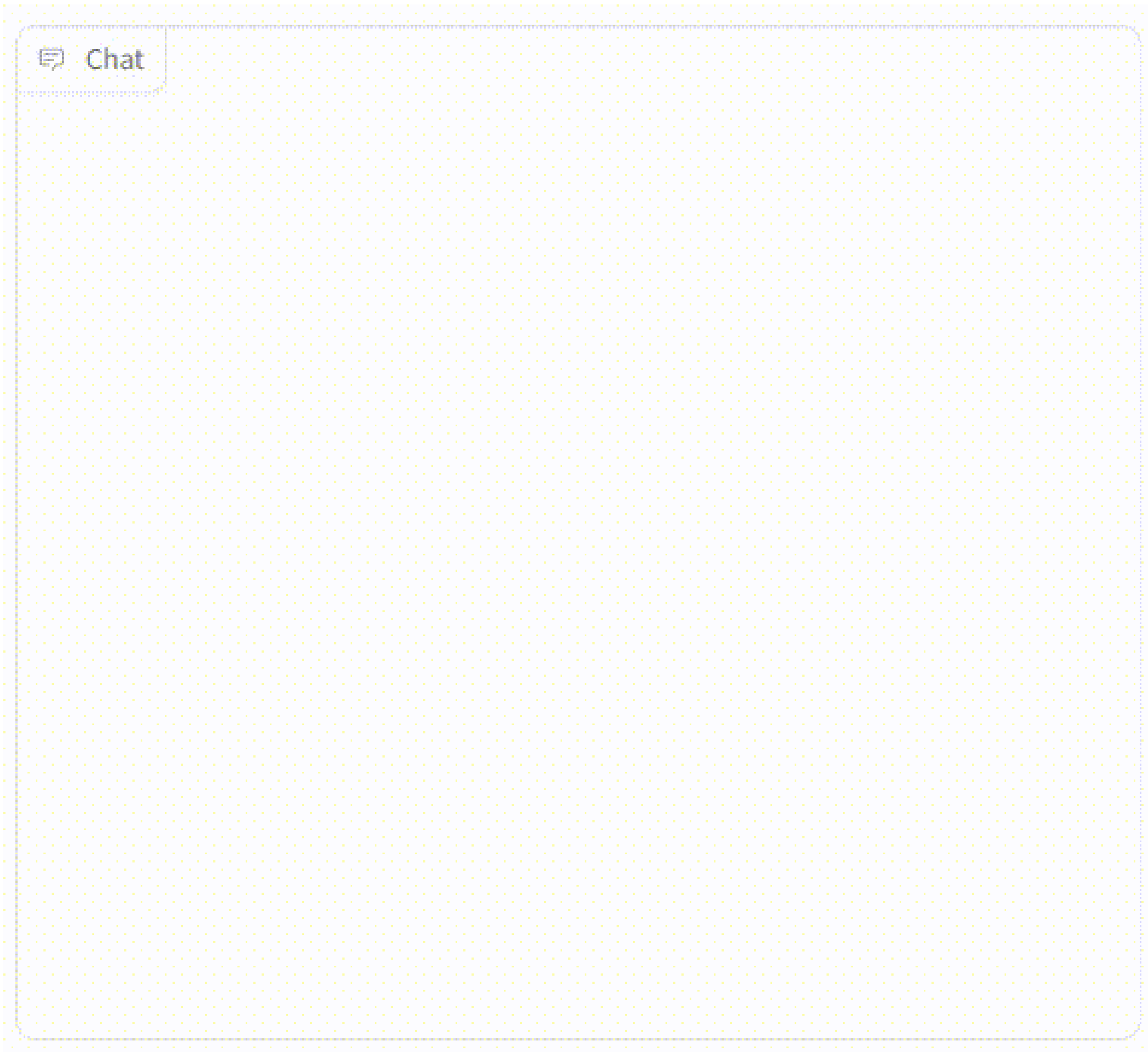
Unmatched performance in ^{FR} amongst models lightweight enough to run on phones and local hardware.



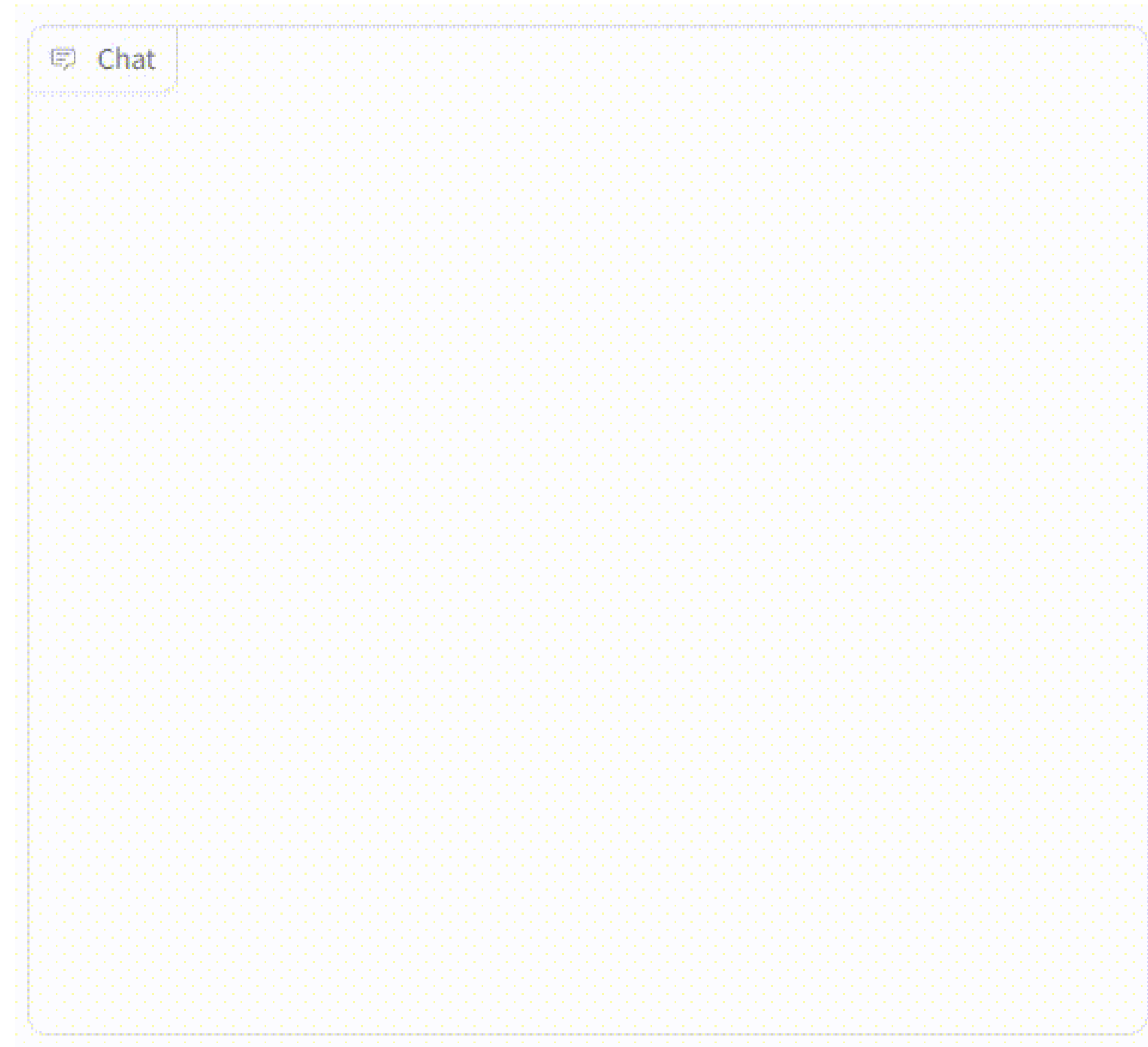
Frugality

The lightweight model reduces costs and energy requirements.

Generation speed



More than **30 tokens per second on CPU**



More than **120 tokens per second on lower end GPU (T4)**

A small recap

- ✓ Generative AI, foundations models and Agent-based AI are important enablers of the democratization of AI.
- ✓ They can be used, with high performance, in a wide range of applications and by a wide range of users (**through prompting or with small task-level adapters**)

BUT

Many obstacles need to be overcome to achieve true democratisation:

- **Trustworthy** Generative AI : **explainable, fair, interpretable, robust, transparent, safe and secure.**
- Encourage initiatives to avoid development being monopolised by a few large players (**sovereignty at stake**).
- LLMs and foundation models are not the only paths : SLMs, RAG, specialized models...
- Need of Governance and Education

The clusters

We presented projects involving a lot of compute

ADASTRA



JEANZAY



The Team



Pierre Colombo



Manuel Faysse



Nicolas Boizard



Hippolyte
Gisserot-Boukhlef



Maria
Vakalopoulou



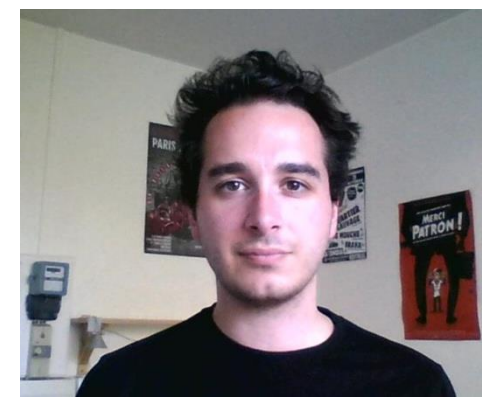
Theo
Moutakanni



Céline Hudelot



Duarte Alves



Caio Corro



Nuno Guerreiro



André Martins

And many others.....