

Inteligență Artificială Generativă - LLMs

AI Curs 7 – 09.04.2025

Capitole

- Întroducere în Generative AI – Alexandru Manole
- Embeddings – Răzvan Petec
- Stable Diffusion?

Despre mine



- Student la doctorat în anul II la UBB FMI
- Domenii de interes: Computer Vision, Multitask Models, Image Generation
- alexandru.manole@ubbcluj.ro

Despre voi



- <https://www.menti.com/>
- cod: **5193 0229**

Introducere în Generative AI



ChatGPT

Gemini

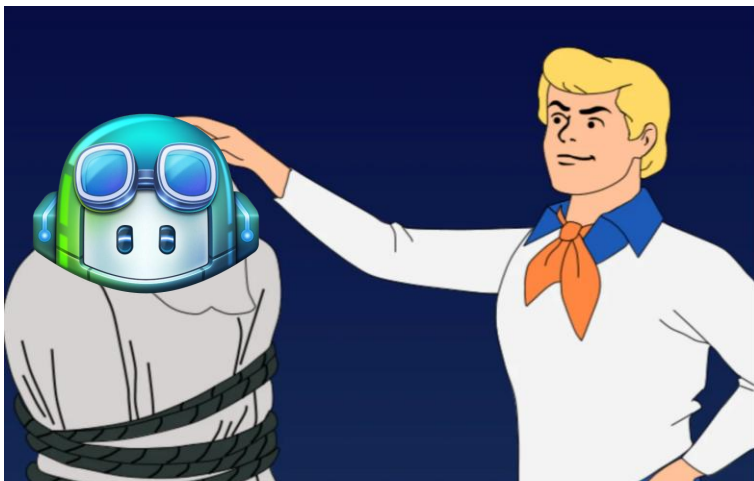


Claude

DALL·E

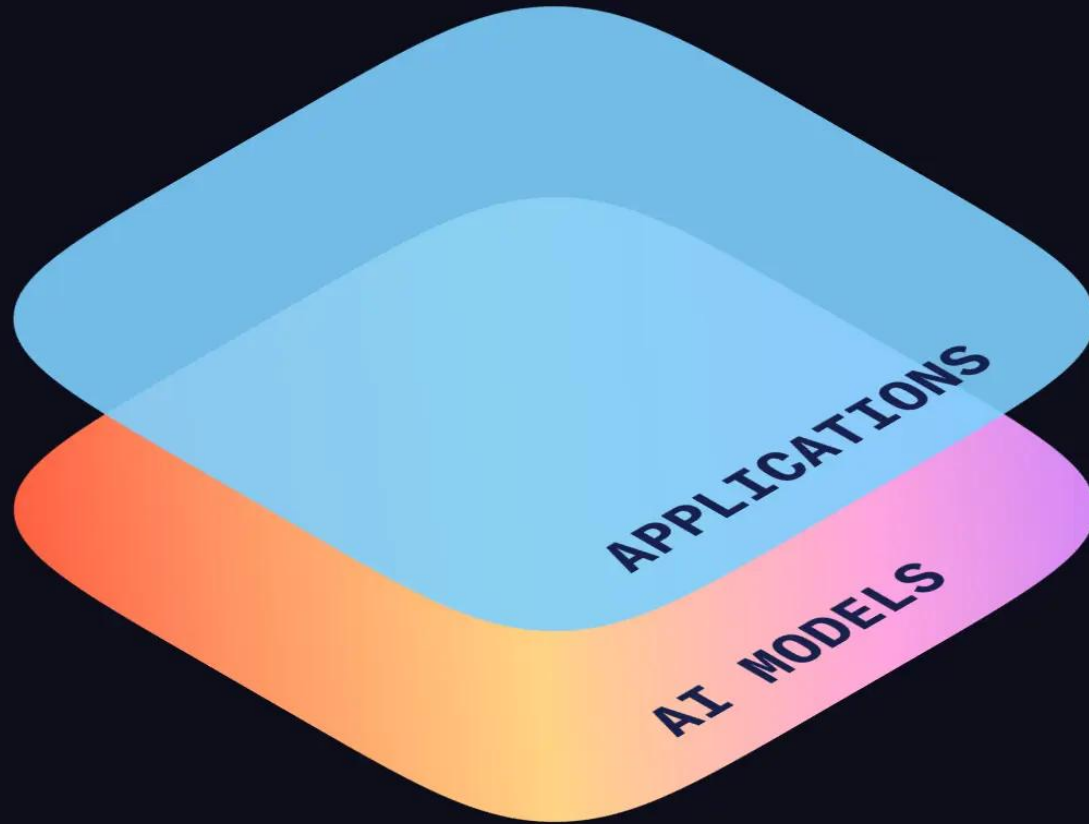


Introducere în Generative AI



- ChatGPT, Gemini, Copilot, Claude sunt sisteme software / tool-uri complexe
- Aceste au în spate modele inteligente generative (LLMs)

The Generative Tech Stack



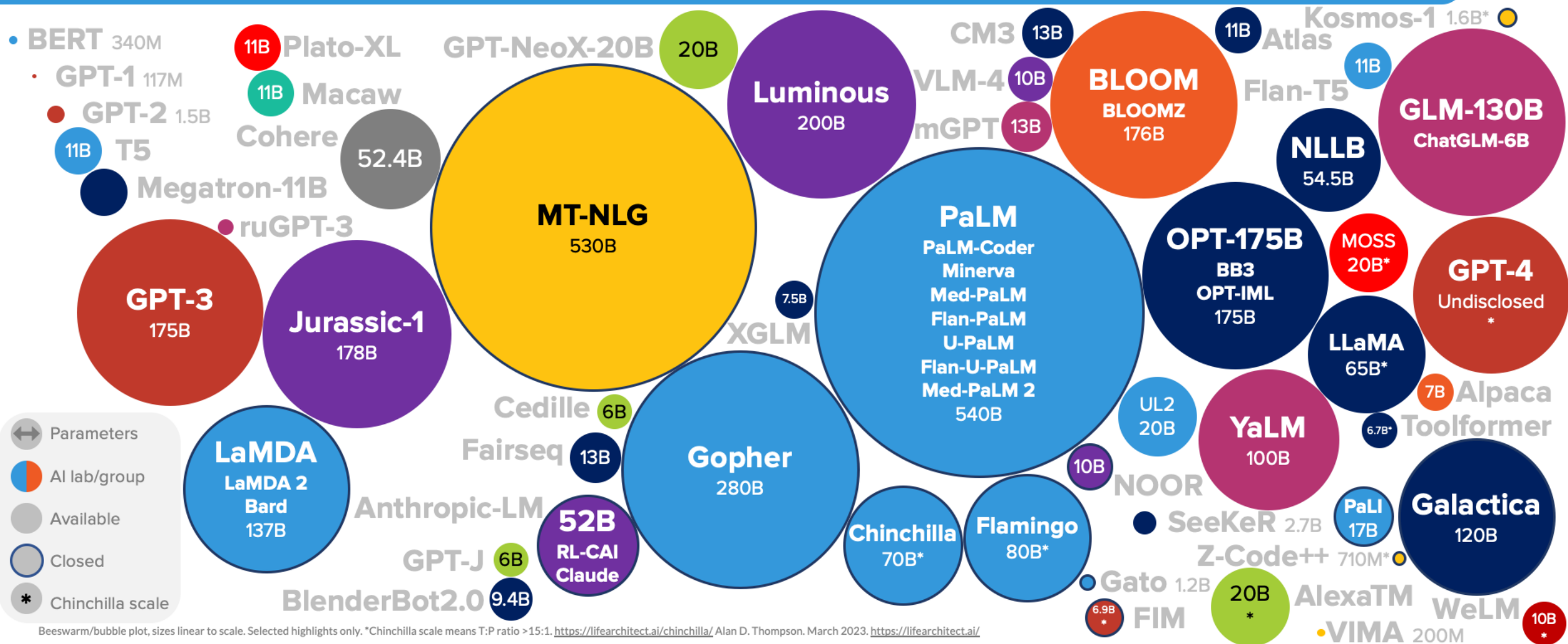
- **TOOLS THAT COLLABORATE WITH AI MODELS.**

WORKFLOWS, SECURITY,
NETWORK EFFECTS, PAYMENTS, ETC.
(10,000'S OF THESE)

- **MODELS GENERATE UNIQUE AND NOVEL OUTPUT.**

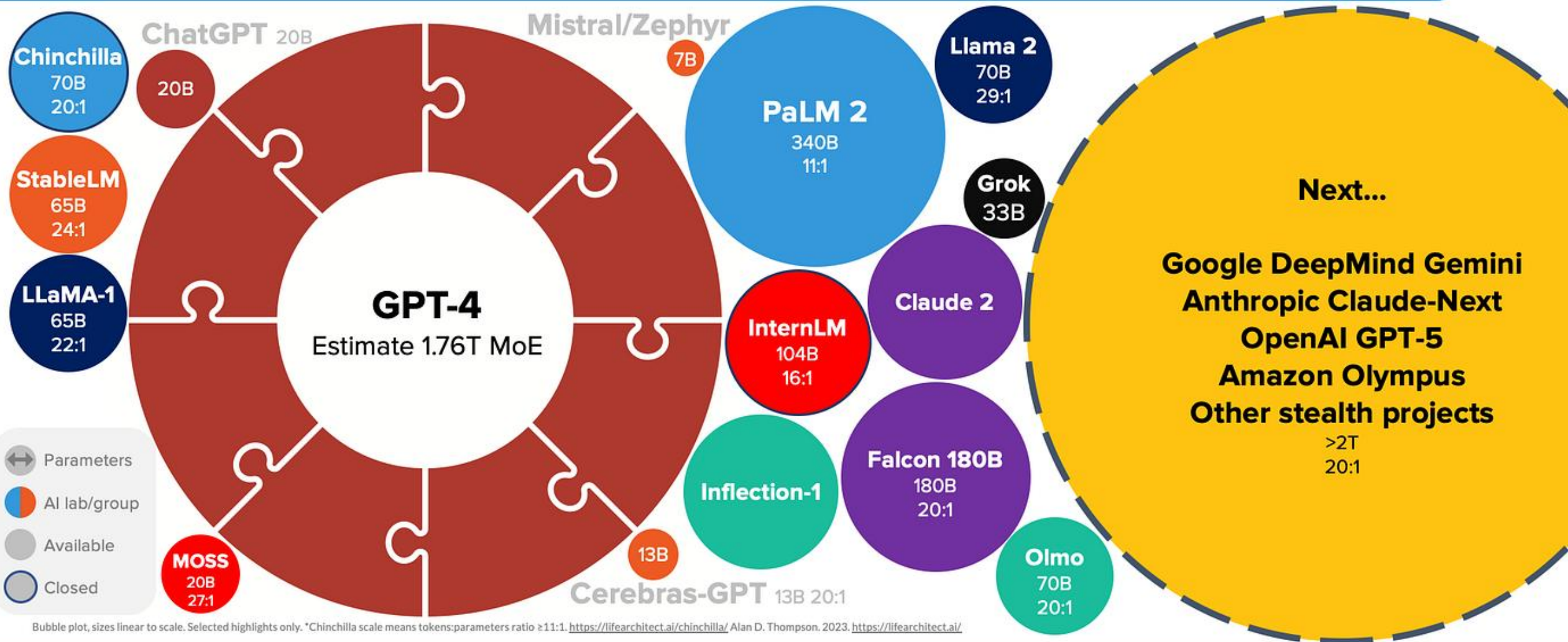
GPT-3, STABLE DIFFUSION,
CUSTOM DATA SETS, ETC.
(1,000'S OF THESE)

LANGUAGE MODEL SIZES TO MAR/2023



2023-2024 OPTIMAL LANGUAGE MODELS

NOV/
2023



Bubble plot, sizes linear to scale. Selected highlights only. *Chinchilla scale means tokens:parameters ratio $\geq 11:1$. <https://lilearchitect.ai/chinchilla/> Alan D. Thompson, 2023. <https://lilearchitect.ai/>

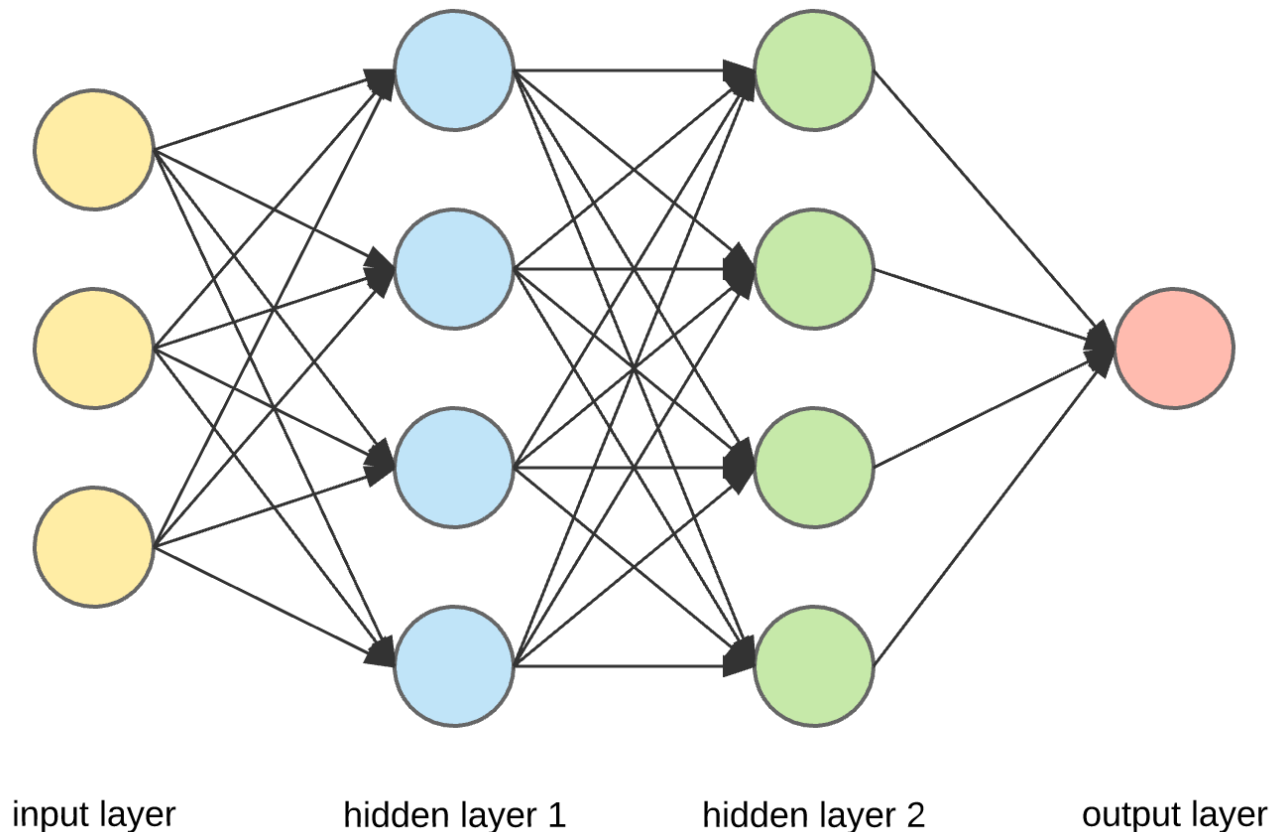


Model	Lab	Playground	Parameters (B)	Tokens trained (B)	Ratio Tokens:Params (Chinchilla scaling220:1)	ALScore "ALScore" i Sqr Root of	MMLU	MMLU -Pro	GPQA	HLE	Training dataset	Announced ▼	Public?	Paper / Arch Repo	Tags
AuroraGPT (ScienceGPT)	Argonne National Lab	https://lifecycle.ai	2000	30000	15:1							TBA	●		
DeepSeek-R2	DeepSeek-AI	https://www.reuters.com										TBA	●	MoE	Reasoning SOTA
ERNIE 5	Baidu	https://lifecycle.ai										TBA			
Grok-4	xAI	https://lifecycle.ai										TBA		MoE	Reasoning SOTA
GPT-5	OpenAI	https://lifecycle.ai	5400	114000	22:1							TBA		MoE	SOTA
GPT-6	OpenAI	https://lifecycle.ai										TBA			SOTA
Llama 4 Reasoning	Meta AI	https://ai.meta.com										TBA	●	https://ai.meta.com	MoE SOTA Reasoning
MAI-1	Microsoft	https://arstechnica.com	500	10000	20:1	7.5						TBA		https://www.microsoft.com	Dense
o4	OpenAI	https://lifecycle.ai										TBA			Reasoning SOTA
Llama 4 Behemoth	Meta AI	https://ai.meta.com	2000	30000	15:1	25.8		82.2	73.7			Apr/2025	●	https://ai.meta.com	MoE SOTA
Llama 4 Maverick	Meta AI	https://ai.meta.com	400	22000	55:1	9.9		80.5	69.8			Apr/2025	●	https://ai.meta.com	MoE SOTA
Qwerky-72B	Featherless AI	https://featherless.ai	72	18000	250:1	3.8	77.46				synthetic, web-scale	Apr/2025	●	https://featherless.ai	Dense
Agentix-Tx	Google DeepMind	https://github.com	200	20000	100:1	6.7			62.4	14.5	synthetic, web-scale	Mar/2025	●	https://storage.googleapis.com	MoE
TxGemma	Google DeepMind	https://huggingface.co	27	14000	519:1	2.0					synthetic, web-scale	Mar/2025	●	https://storage.googleapis.com	MoE Reasoning
Gemini 2.5 Pro	Google DeepMind	https://aistudio.google.com	200	20000	100:1	6.7			84	18.8	synthetic, web-scale	Mar/2025	●	https://blog.google	MoE Reasoning SOTA
DeepSeek-V3 0324	DeepSeek-AI	https://chat.deepseek.com	685	14800	22:1	10.6		81.2	68.4		synthetic, web-scale	Mar/2025	●	https://huggingface.co	MoE SOTA
Llama-3.3-Nemotron-70B	NVIDIA	https://huggingface.co	49	15040	307:1	2.9			66.67		web-scale	Mar/2025	●	https://build.nvidia.com	Dense Reasoning
EXAONE Deep	LG	https://huggingface.co	32	6500	204:1	1.5	83	74	66.1		web-scale	Mar/2025	●	https://arxiv.org	Dense Reasoning
Mistral Small 3.1	Mistral	https://huggingface.co	24	8000	334:1	1.5	81.01	56.03	37.5		web-scale	Mar/2025	●	https://mistral.ai	Dense
ERNIE 4.5	Baidu	https://yiyen.baidu.com									synthetic, web-scale	Mar/2025	●	https://www.baidu.com	MoE
X1	Baidu	https://yiyen.baidu.com									synthetic, web-scale	Mar/2025	●	https://www.baidu.com	MoE Reasoning
OLMo 2 32B	Allen AI	https://playground.allenai.org	32	6400	200:1	1.5	78				synthetic, web-scale	Mar/2025	●	https://allenai.org	Dense
Command A	Cohere	https://dashboards.cohere.com	111	8000	73:1	3.1	85				synthetic, web-scale	Mar/2025	●	https://huggingface.co	Dense
Gemini Robotics	Google DeepMind		200	20000	100:1	6.7		79.1	64.7		synthetic, web-scale	Mar/2025	●	https://storage.googleapis.com	MoE

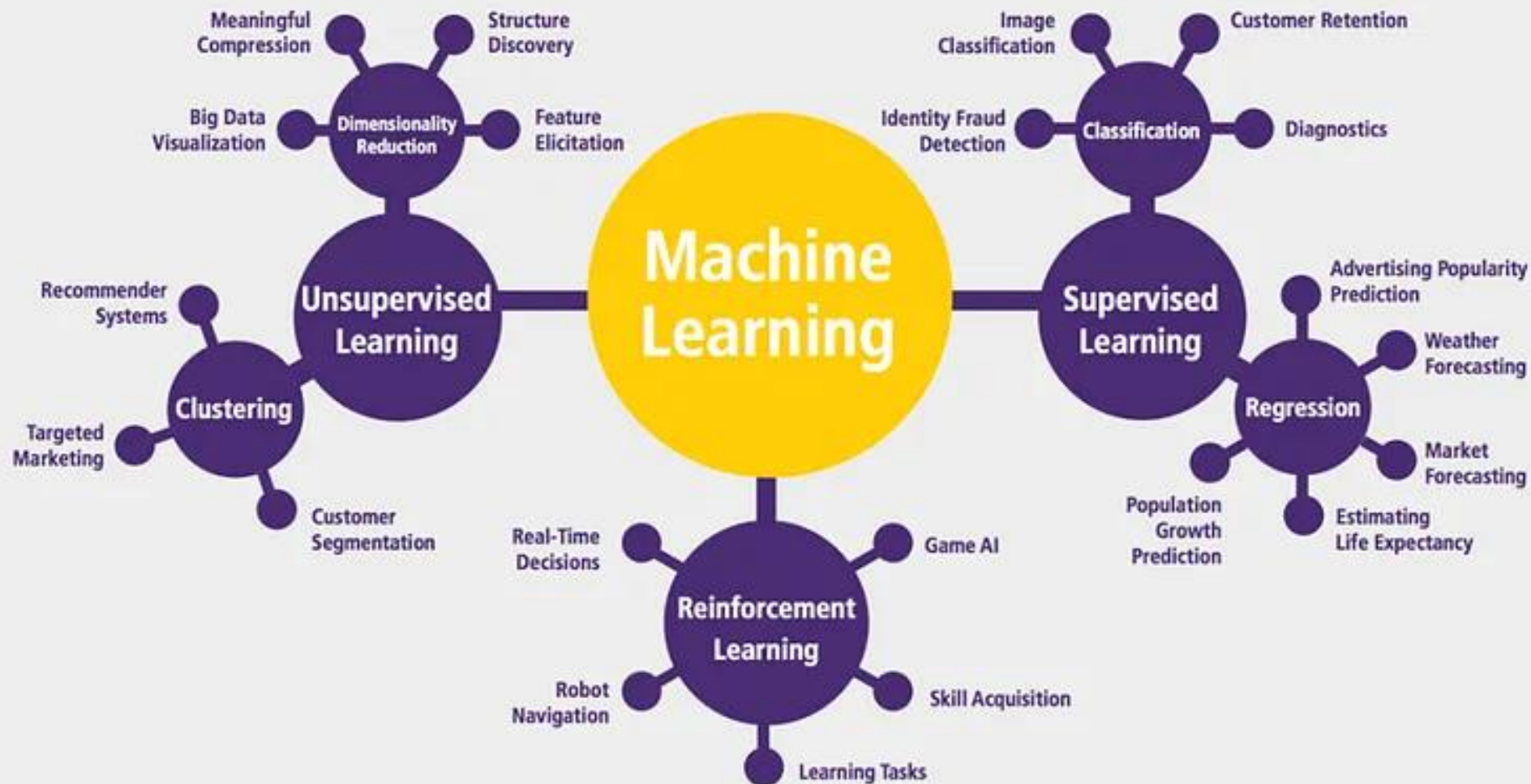
<https://lifecycle.ai/>

Elemente AI / Machine Learning:

1. Algoritm / Model inteligent
2. Date
 - Date de intrare
 - Date de ieșire



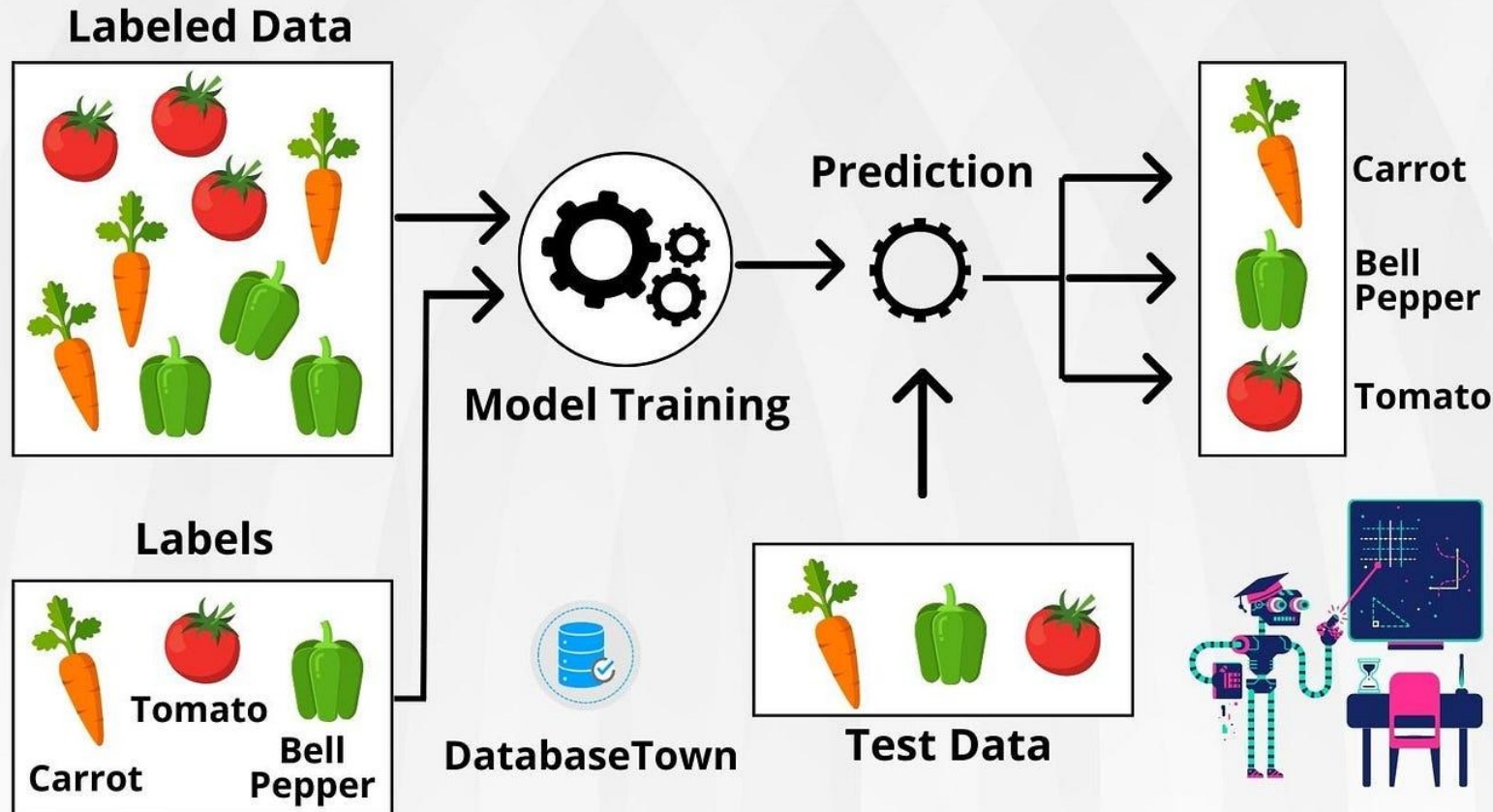
Învățare automată



Învățare supervizată

SUPERVISED LEARNING

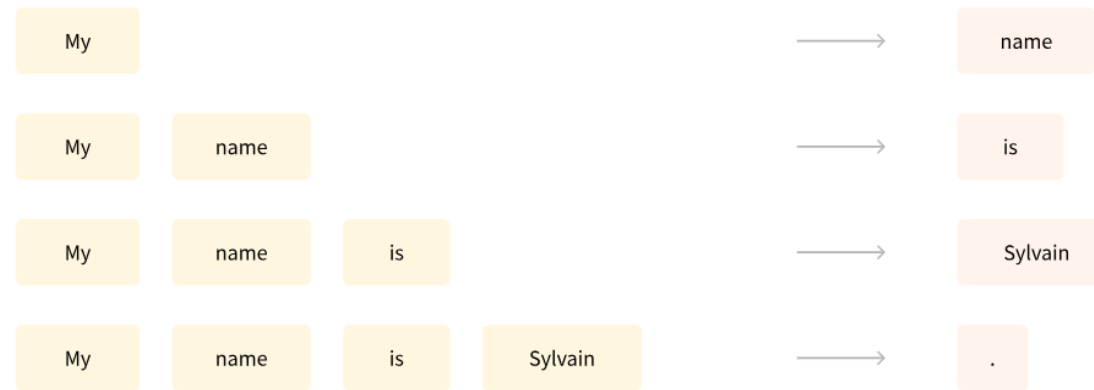
Supervised machine learning is a branch of artificial intelligence that focuses on training models to make predictions or decisions based on labeled training data.



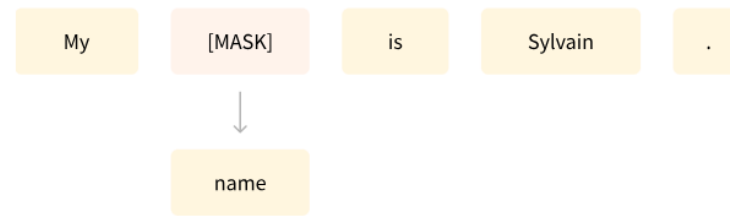
Ce clasifică LLM-urile?

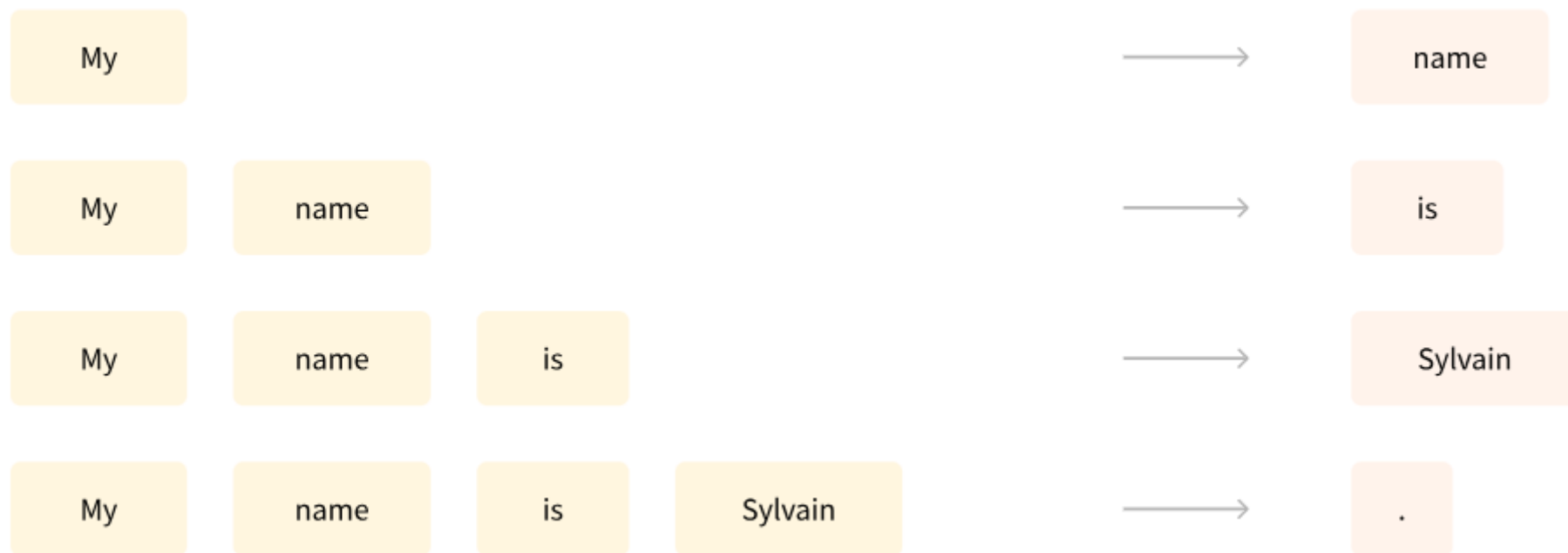


Ce clasifică LLM-urile?

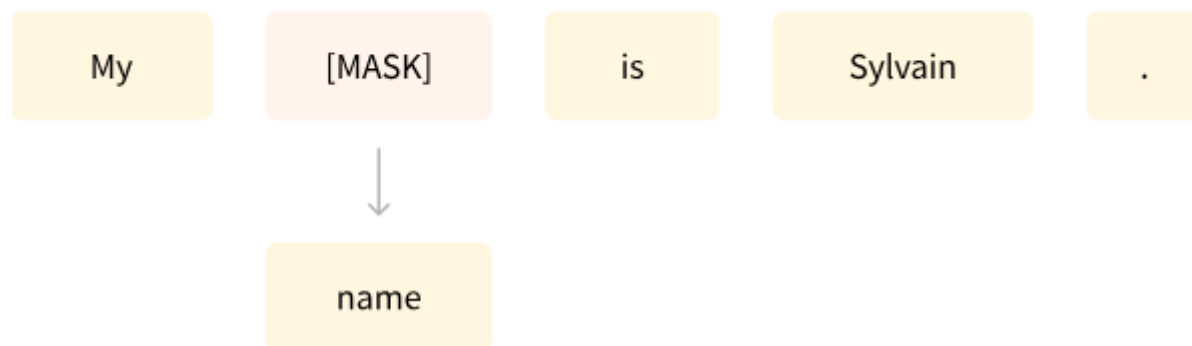


Another example is *masked language modeling*, in which the model predicts a masked word in the sentence.





Another example is *masked language modeling*, in which the model predicts a masked word in the sentence.



Proces de antrenare

Modelele generate

cuv 1 cuv

5193 0229



- Derivind de la un șir
cuvine modelul
zice ce cuvânt ar
ea urma

Rezultatul predicției

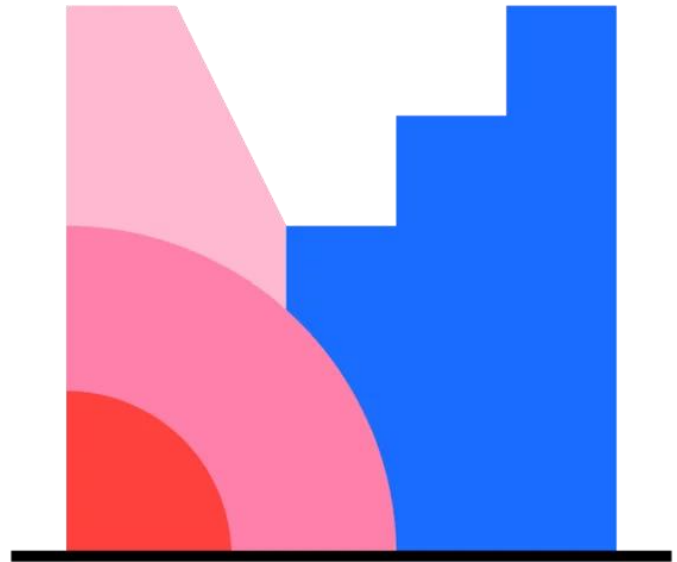
output token probabilities (logits)

model vocabulary size
50,257



0.19850038	aardvark
0.7089803	aarhus
0.46333563	aaron
	...
	...
	...
	...
	...
-0.51006055	zyzzyva

Cuvinte în limba română?



Mentimeter

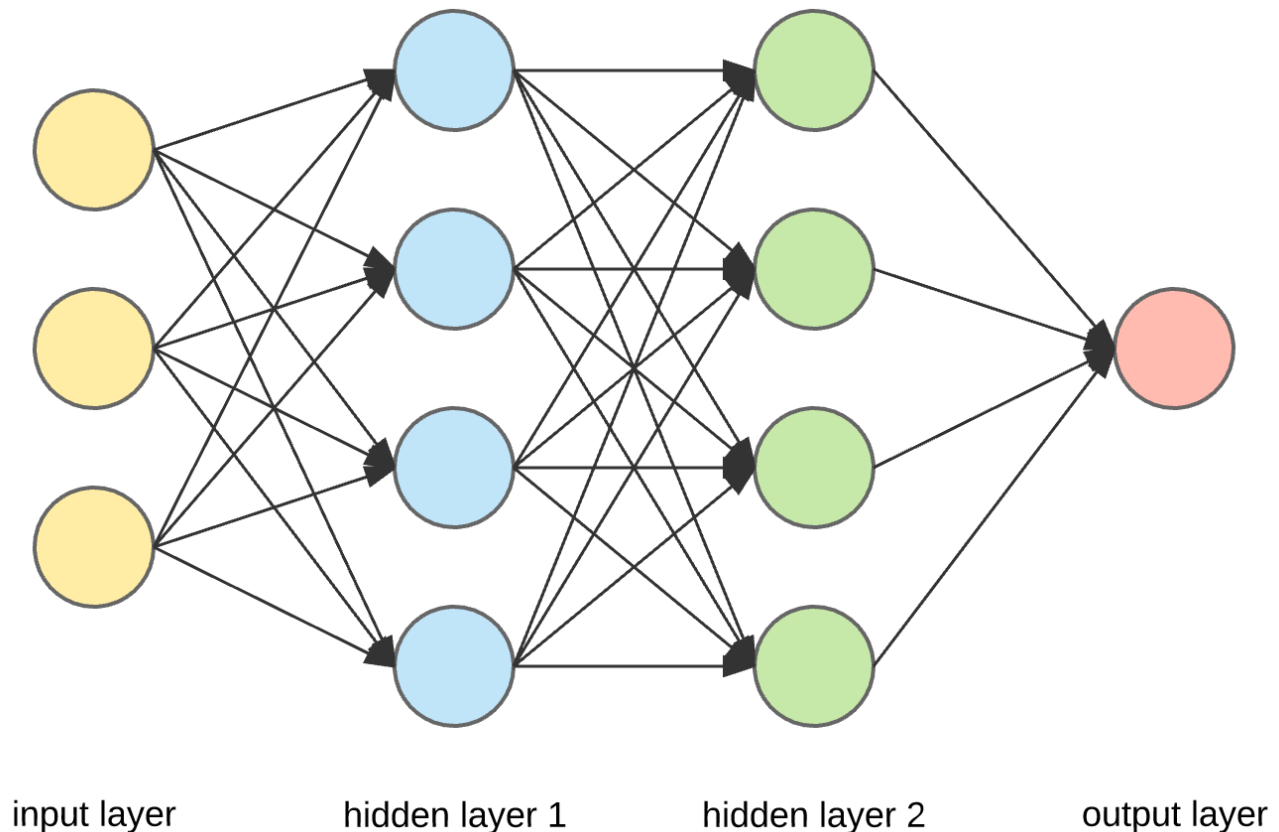
- <https://www.menti.com/>
- cod: **5193 0229**

Cuvinte în limba română

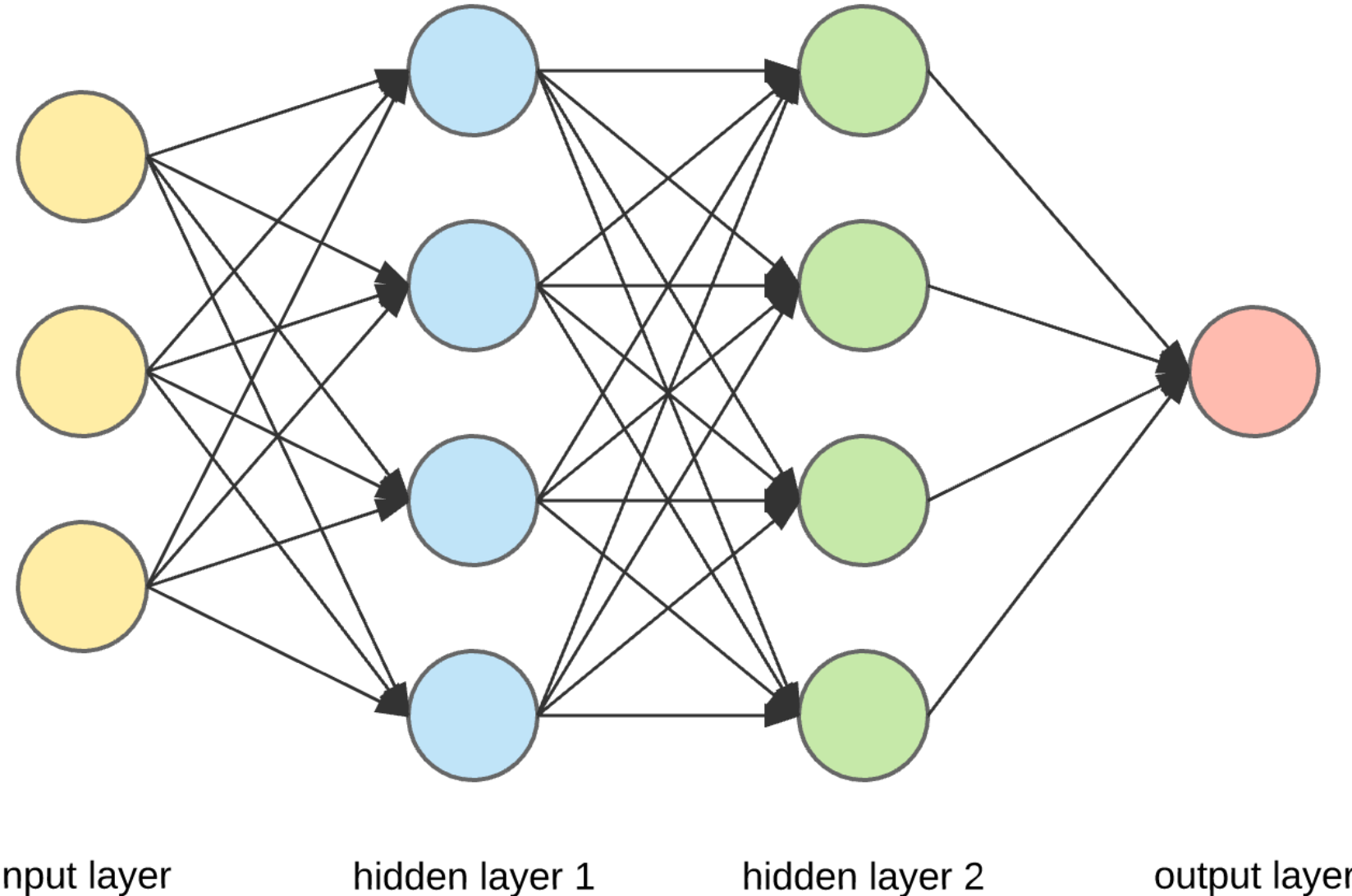
- DOOM – “peste 62.000 cuvinte”
- Dexonline.ro – 75.399 cuvinte
- DLR – 175.000 cuvinte

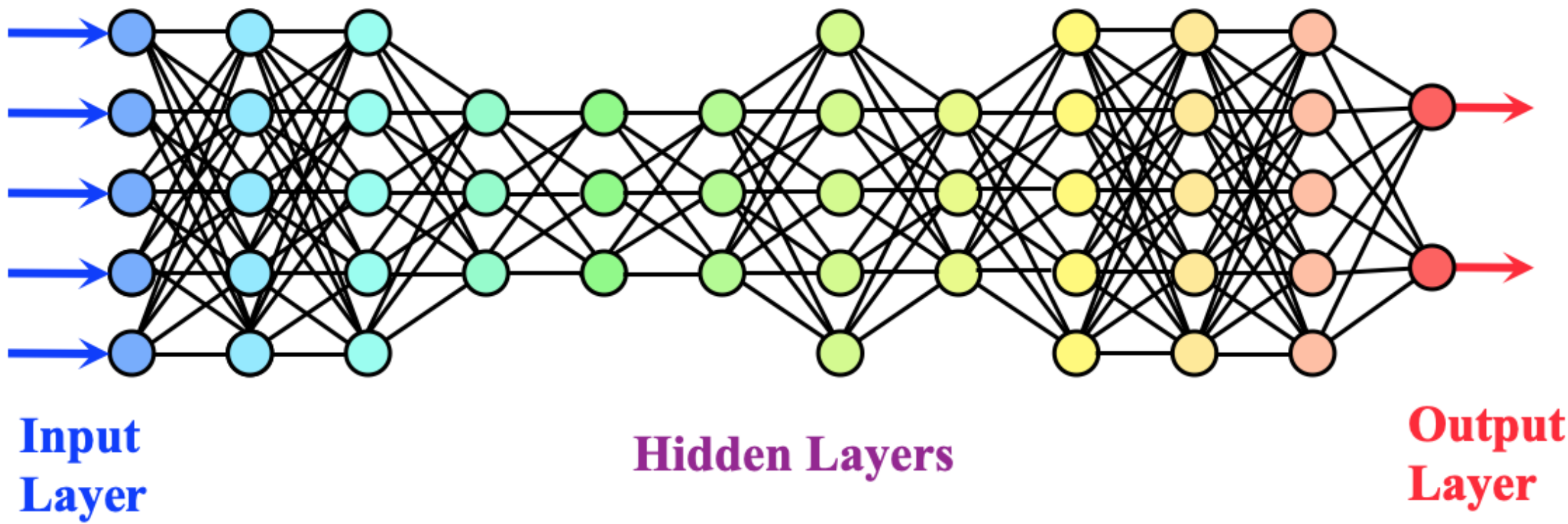
Elemente AI / Machine Learning:

1. Algoritm / Model inteligent
2. Date
 - Date de intrare
 - Date de ieșire

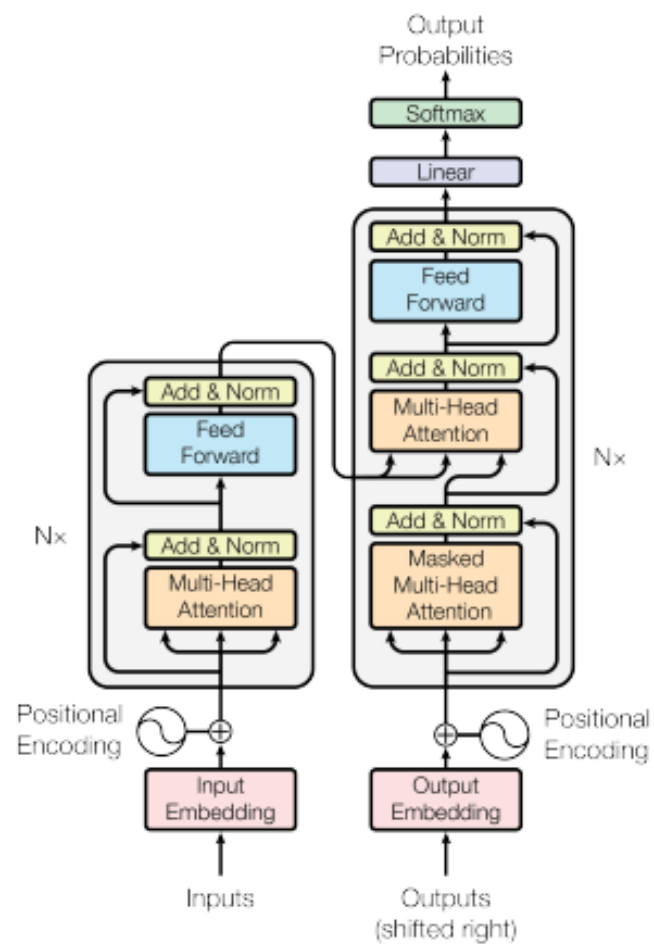


Cum se numește acest algoritm inteligent:
menti.com cod: **5193 0229**

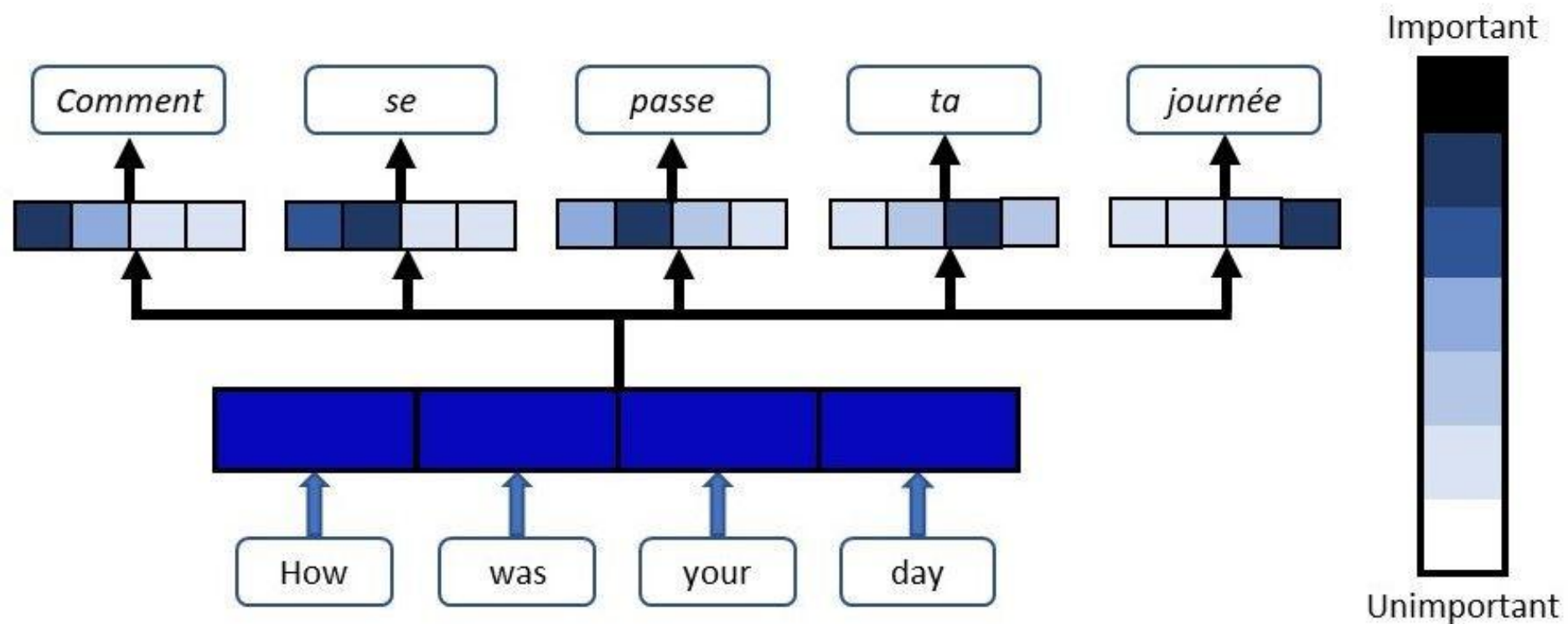




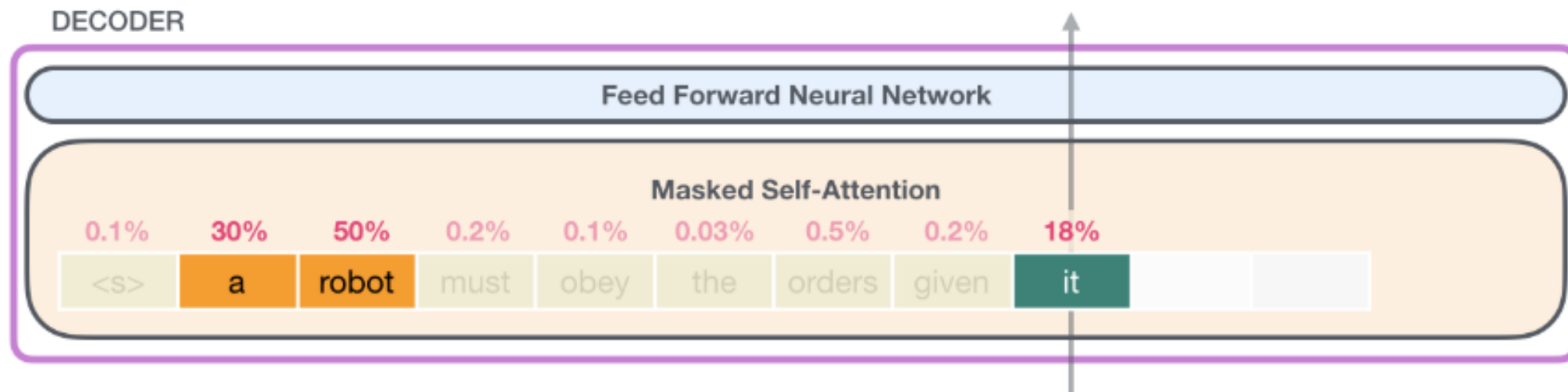
Transformers



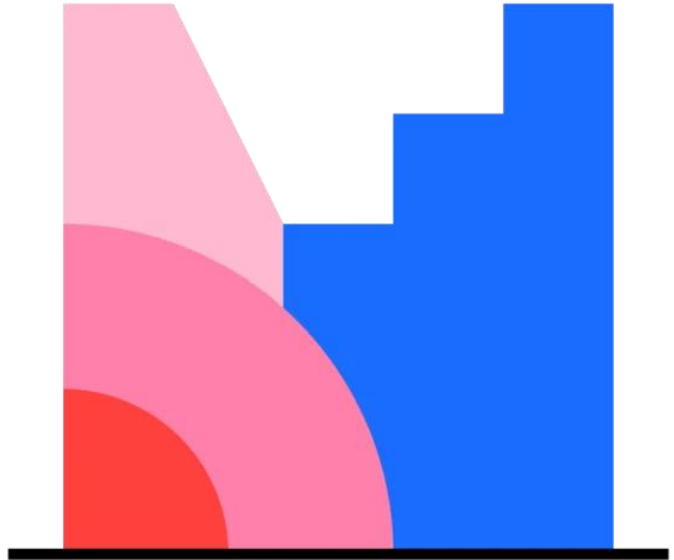
Attention



Attention



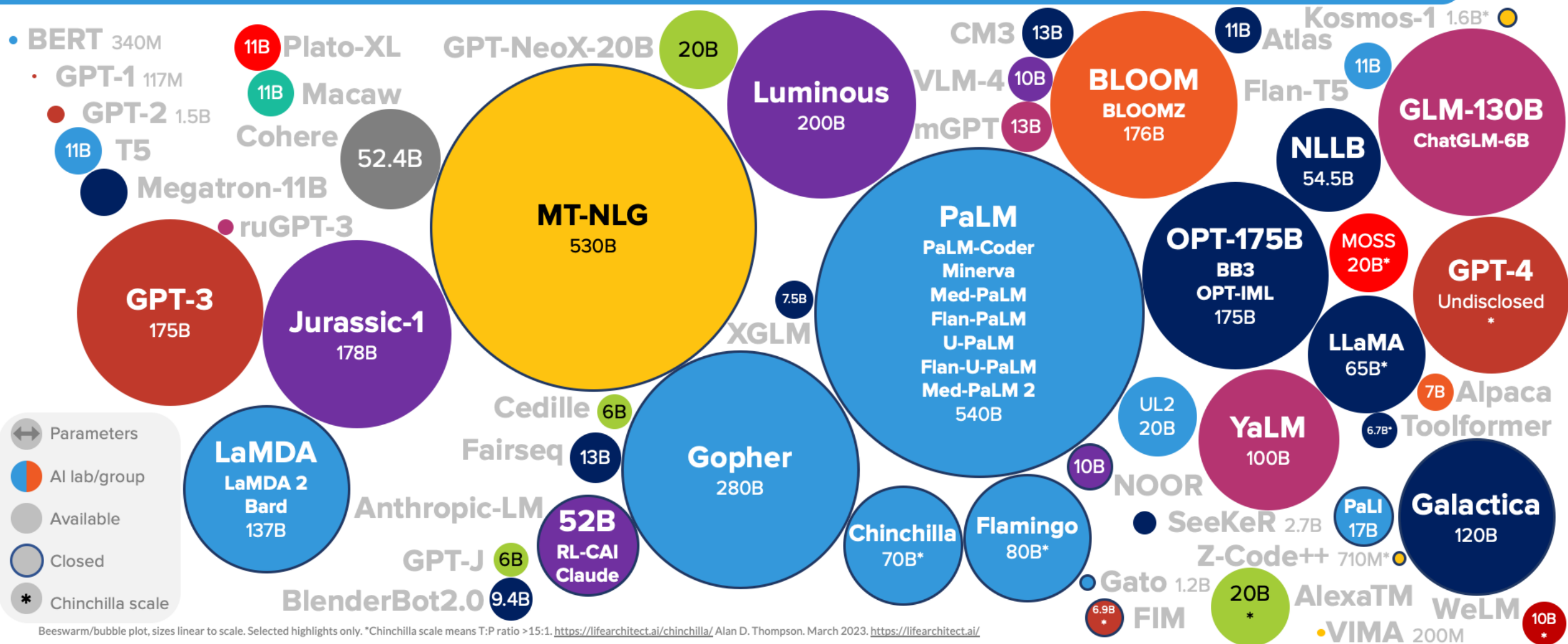
Atentie!



Mentimeter

- <https://www.menti.com/>
- cod: **5193 0229**

LANGUAGE MODEL SIZES TO MAR/2023



Model	Lab	Playground	Parameters (B)	Tokens trained (B)	Ratio Tokens:Params (Chinchilla scaling220:1)	ALScore "ALScore" i Sqr Root of	MMLU	MMLU -Pro	GPQA	HLE	Training dataset	Announced ▼	Public?	Paper / Arch Repo	Tags
AuroraGPT (ScienceGPT)	Argonne National Lab	https://lifecycle.ai	2000	30000	15:1							TBA	●		
DeepSeek-R2	DeepSeek-AI	https://www.reuters.com										TBA	●	MoE	Reasoning SOTA
ERNIE 5	Baidu	https://lifecycle.ai										TBA			
Grok-4	xAI	https://lifecycle.ai										TBA		MoE	Reasoning SOTA
GPT-5	OpenAI	https://lifecycle.ai	5400	114000	22:1							TBA		MoE	SOTA
GPT-6	OpenAI	https://lifecycle.ai										TBA			SOTA
Llama 4 Reasoning	Meta AI	https://ai.meta.com										TBA	●	https://ai.meta.com	MoE SOTA Reasoning
MAI-1	Microsoft	https://arstechnica.com	500	10000	20:1	7.5						TBA		https://www.microsoft.com	Dense
o4	OpenAI	https://lifecycle.ai										TBA			Reasoning SOTA
Llama 4 Behemoth	Meta AI	https://ai.meta.com	2000	30000	15:1	25.8		82.2	73.7			Apr/2025	●	https://ai.meta.com	MoE SOTA
Llama 4 Maverick	Meta AI	https://ai.meta.com	400	22000	55:1	9.9		80.5	69.8			Apr/2025	●	https://ai.meta.com	MoE SOTA
Qwen-72B	Featherless AI	https://featherless.ai	72	18000	250:1	3.8	77.46				synthetic, web-scale	Apr/2025	●	https://featherless.ai	Dense
Agentix-Tx	Google DeepMind	https://github.com	200	20000	100:1	6.7			62.4	14.5	synthetic, web-scale	Mar/2025	●	https://storage.googleapis.com	MoE
TxGemma	Google DeepMind	https://huggingface.co	27	14000	519:1	2.0					synthetic, web-scale	Mar/2025	●	https://storage.googleapis.com	MoE Reasoning
Gemini 2.5 Pro	Google DeepMind	https://aistudio.google.com	200	20000	100:1	6.7			84	18.8	synthetic, web-scale	Mar/2025	●	https://blog.google	MoE Reasoning SOTA
DeepSeek-V3 0324	DeepSeek-AI	https://chat.deepseek.com	685	14800	22:1	10.6		81.2	68.4		synthetic, web-scale	Mar/2025	●	https://huggingface.co	MoE SOTA
Llama-3.3-Nemotron-70B	NVIDIA	https://huggingface.co	49	15040	307:1	2.9			66.67		web-scale	Mar/2025	●	https://build.nvidia.com	Dense Reasoning
EXAONE Deep	LG	https://huggingface.co	32	6500	204:1	1.5	83	74	66.1		web-scale	Mar/2025	●	https://arxiv.org	Dense Reasoning
Mistral Small 3.1	Mistral	https://huggingface.co	24	8000	334:1	1.5	81.01	56.03	37.5		web-scale	Mar/2025	●	https://mistral.ai	Dense
ERNIE 4.5	Baidu	https://yiyen.baidu.com									synthetic, web-scale	Mar/2025	●	https://www.baidu.com	MoE
X1	Baidu	https://yiyen.baidu.com									synthetic, web-scale	Mar/2025	●	https://www.baidu.com	MoE Reasoning
OLMo 2 32B	Allen AI	https://playground.allenai.org	32	6400	200:1	1.5	78				synthetic, web-scale	Mar/2025	●	https://allenai.org	Dense
Command A	Cohere	https://dashboards.cohere.com	111	8000	73:1	3.1	85				synthetic, web-scale	Mar/2025	●	https://huggingface.co	Dense
Gemini Robotics	Google DeepMind		200	20000	100:1	6.7		79.1	64.7		synthetic, web-scale	Mar/2025	●	https://storage.googleapis.com	MoE

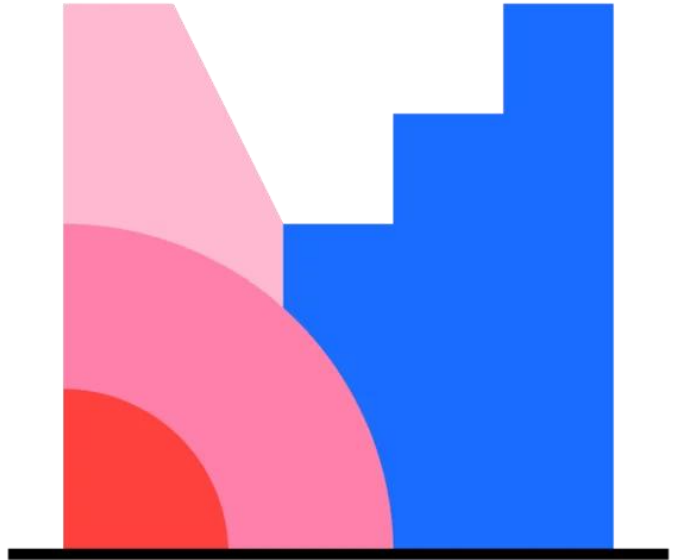
[https://lifecycle.ai/](https://lifecycle.ai)

Exemplu LLM

- Llama 4 Maverick
- Dezvoltat de Meta
- 400 **MILIARDE** de parametri
- menti.com cod: **5193 0229**

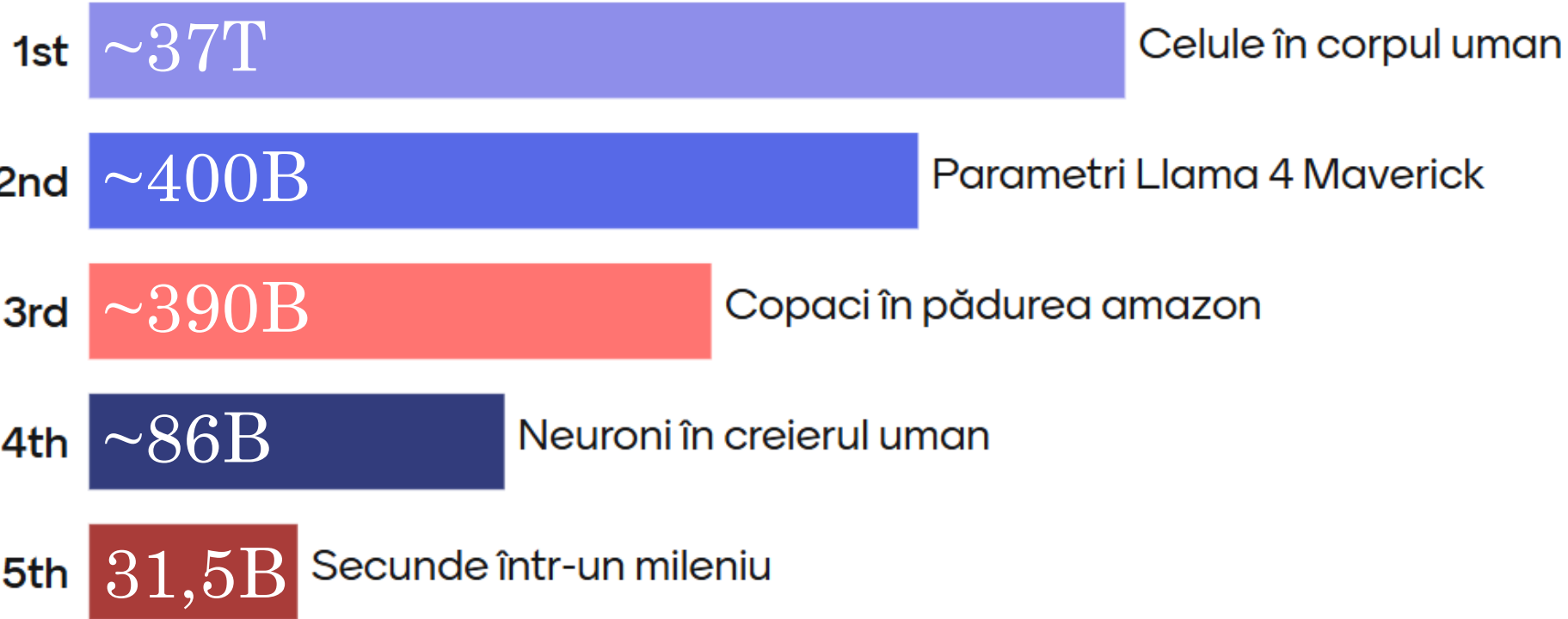
Llama 4:
Leading Multimodal Intelligence

Big numbers

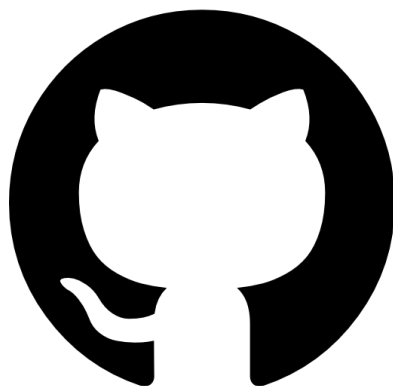
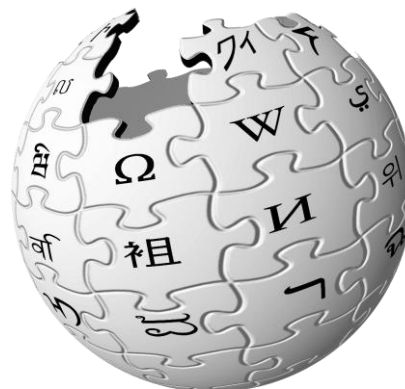


Mentimeter

- <https://www.menti.com/>
- cod: **5193 0229**



Set de date



Ce mai poate fi generat așa?

menti.com cod: 5193 0229



Muzică

```
WS.on("message", m => {  
  let a = m.split(" ")  
  switch(a[0]){  
    case "connect":  
      if(a[1]){  
        if(clients.has(a[1])){  
          ws.send("connected");  
          ws.id = a[1];  
        }else{  
          ws.id = a[1]  
          clients.set(a[1], {client: {position: {x: 0, y: 0, id: 0}, name: ""}})  
          ws.send("connected")  
        }  
      }else{  
        let id = Math.random().toString().slice(2, 9)  
        ws.id = id;  
        clients.set(id, {client: {position: {x: 0, y: 0, id: 0}, name: ""}})  
      }  
    }  
  }  
})
```

Cod

Ce mai poate fi generat așa?



suno.com

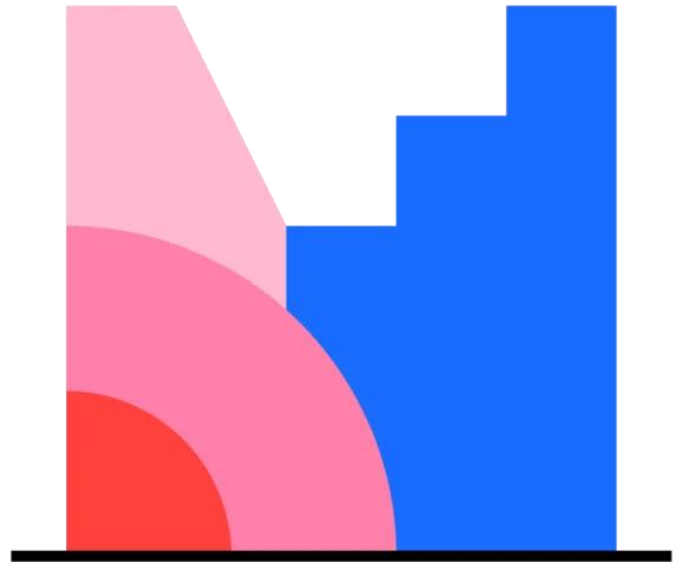
```
ws.on("message", m => {  
  let a = m.split(" ")  
  switch(a[0]){  
    case "connect":  
      if(a[1]){  
        if(clients.has(a[1])){  
          ws.send("connected");  
          ws.id = a[1];  
        }else{  
          ws.id = a[1]  
          clients.set(a[1], {client: {position: {x: 0, y: 0}, id: a[1]}})  
          ws.send("connected")  
        }  
      }else{  
        let id = Math.random().toString().slice(2, 8)  
        ws.id = id;  
        clients.set(id, {client: {position: {x: 0, y: 0}, id: id}})  
      }  
    }  
  }  
})
```

Copilot \ Claude \ LLMs



sora.com

Generare video



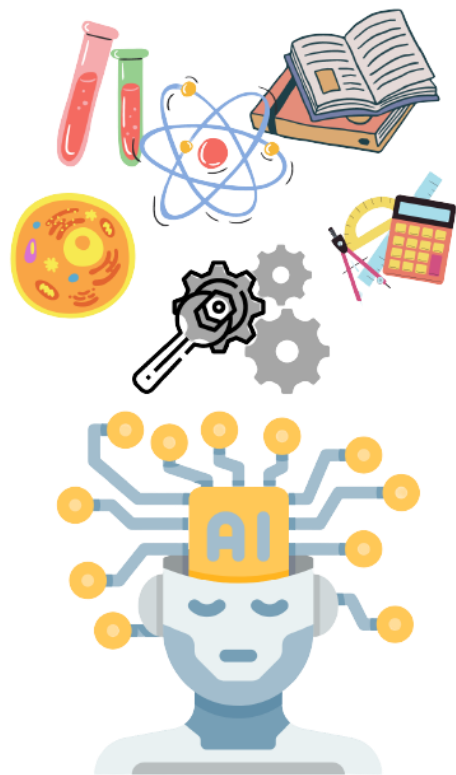
Mentimeter

- <https://www.menti.com/>
- cod: **5193 0229**

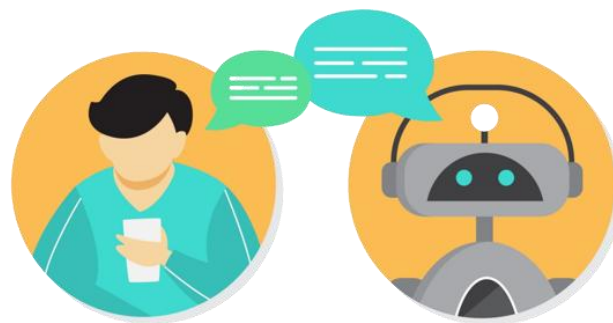
Antrenare LLM

- **Pasul 1: Antrenare inițială (Pre-training)**
 - **Scop:** înțelegerea limbajului natural
 - Prezicerea următorului cuvând din text
- **Pasul 2: Antrenare de bază (Training alignment)**
 - **Scop:** dobândirea de competențe specifice
 - Perechi Prompt X + Răspuns așteptat Y
- **Pasul 3: Reinforcement Learning**
 - **Scop:** îmbunătățirea calității umane
 - Ordonarea răspunsurilor + antrenarea unui model capabil să facă singur ordonarea

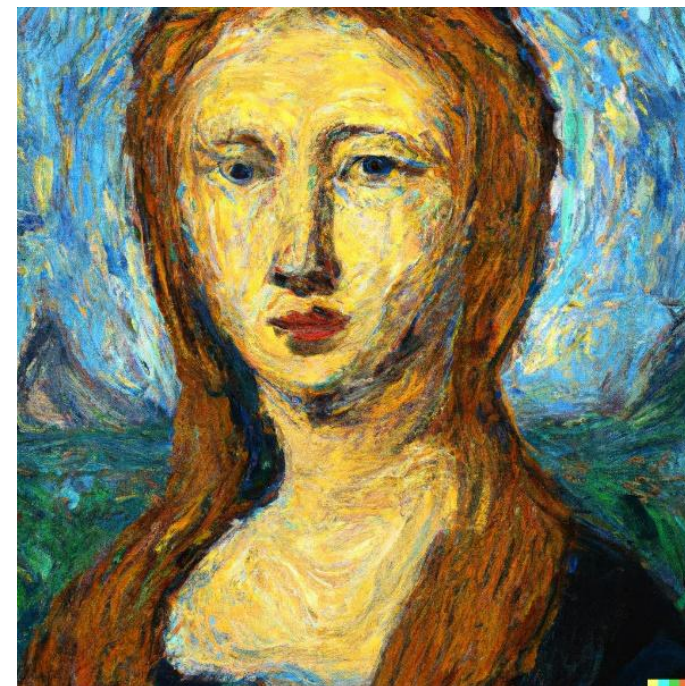
Potențial Generative AI



**Pseudo-expert
în mai multe
domenii**



**Mod interactiv de a
afla informații**



**Inspirație
creativă**

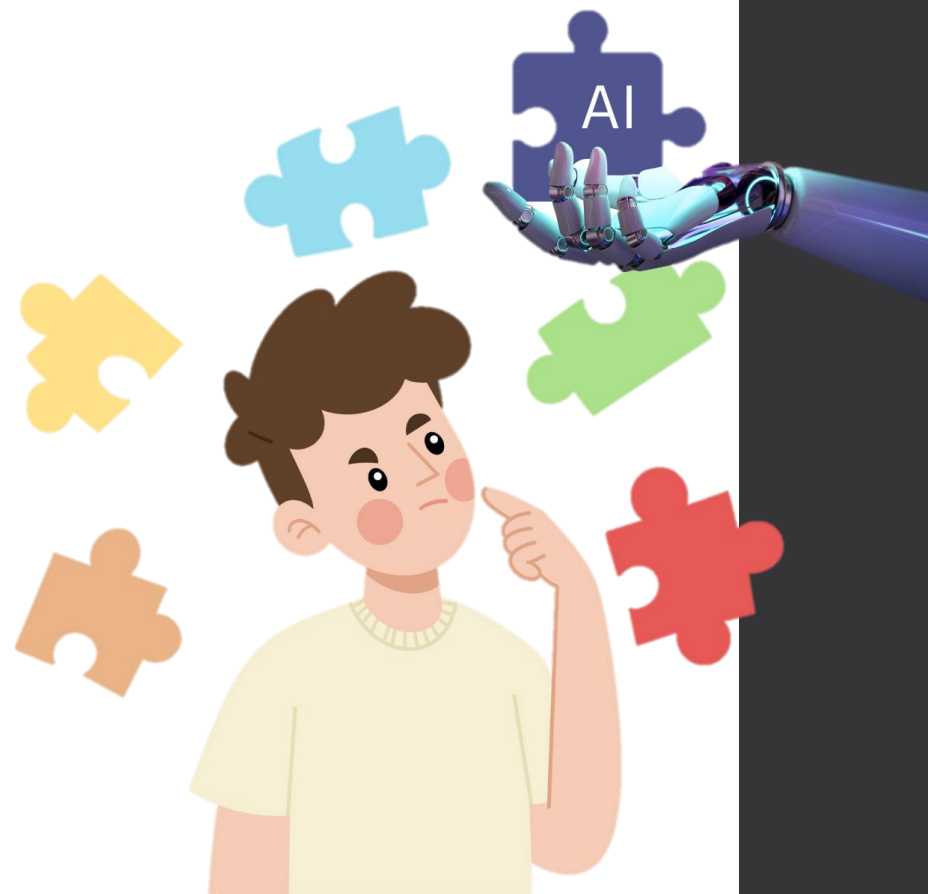
Pericole Generative AI



**Greșeli, bias,
halucinare**



**Gândire de
grup**

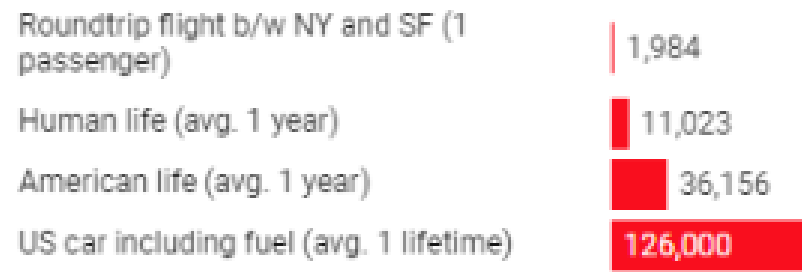


**Are potențialul de a
distruge gândirea
critică**

Probleme Generative AI

Common carbon footprint benchmarks

in lbs of CO2 equivalent



Ce urmează în Generative AI?

- Continuare\Extindere multi-modal LLMs
- Longer memory
- Autonomous AI Agents
- Reasoning

Ce urmează în Generative AI?

The Next Steps in Generative AI



Autonomous AI Agents

goal-driven.
planning, acting



Modular Reasoning Models (Tool Use)

tool calling, calculator, RAG



Improved Long-Term Memory

personalization, context



Reasoning + World Modeling

scale
internal representations



Multimodal Fusion

images, video, audio, data



Alignment, Safety & Interpretability

honesty, transparency

Imagine creată de ChatGPT

Resurse utile

- Karpathy Zero to Hero GPT: <https://www.youtube.com/watch?v=kCc8FmEb1nY&list=PLAqhIrjkxbuWI23v9cThsA9GvCAUhRvKZ&index=7>
- 3Blue1Brown: <https://www.youtube.com/watch?v=wjZofJX0v4M>
- Hugging face LLM course: <https://huggingface.co/learn/llm-course/chapter1/1>
- The annotated Transformer: <https://nlp.seas.harvard.edu/annotated-transformer/>
- Illustrated GPT2: <https://jalammar.github.io/illustrated-gpt2/>
- Running a LLM: https://huggingface.co/docs/transformers/llm_tutorial
- Visual Question Answering: https://huggingface.co/docs/transformers/main/en/tasks/visual_question_answering

Sfârșit partea I

Multumesc!