# Compétence transversale – L2
# Grands défis sociétaux : Intelligence Artificielle

Pr. Lucile Sassatelli

Professeure des Universités en Informatique, UniCA

Directrice scientifique de EFELIA Côte d'Azur

# Plan du module

| Chapitre | Titre | Contenu | Date d'ouverture | Date QCM |
|---|---|---|---|---|
| 1 | **Rappel : IA sous le capot** | • Choix humains et principes de fonctionnement<br>• Faiblesses de la technologie<br>• Impacts sociétaux et environnementaux | | |
| 2 | **Qu'est-ce qui est porté par le terme IA ?** | • Objectifs et croyances<br>• Modes de production | | • QCM 1 noté 3-7/11 |
| 3 | **Est-ce que ça peut ou ça doit lire, écrire, penser pour moi?** | • Calculatrice, puis LLM : devez-vous encore faire l'effort d'écrire ? D'écrire quoi pour quoi faire ?<br>• Quelle place des LLM dans le développement de notre pensée ?<br>• Est-ce que ces réponses dépendent de notre discipline ? | | |
| 4 | **Et pour ma discipline ?** | • Quelles avancées pour ma discipline ?<br>• Quels nouveaux problèmes pour ma discipline ? | | • QCM 2 noté 8-12/12 |

# Problématique

- Pour le portail Sciences et Techniques, nous allons donner dans ce chapiter des éléments sur les questions :
  - Comment fonctionnent plus précisément les modèles de ML et ChatGPT en particulier ?
  - Quelles sont les performances des LLMs pour des tâches de raisonnement ?
  - Comment le marché du travail de développement logiciel est modifié par l'arrivée des LLMs ?
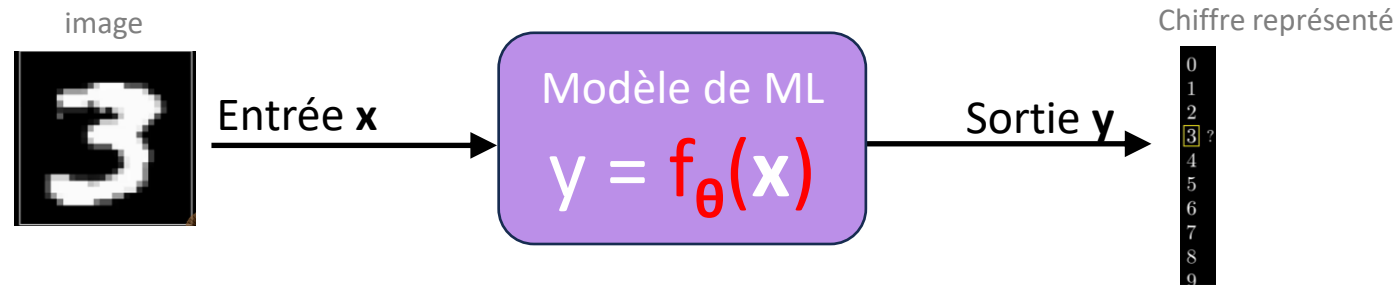
# Plan

→ 1. Formalisation du principe de l'apprentissage machine et des réseaux de neurones artificiels

2. Un peu plus de précisions sur ChatGPT

3. Eléments sur les LLMs pour les tâches de raisonnement

4. Eléments sur le marché du travail de la programmation

# Tâche : reconnaître des chiffres manuscrits

image                                Chiffre représenté

Entrée **x**      Modèle de ML      Sortie **y**
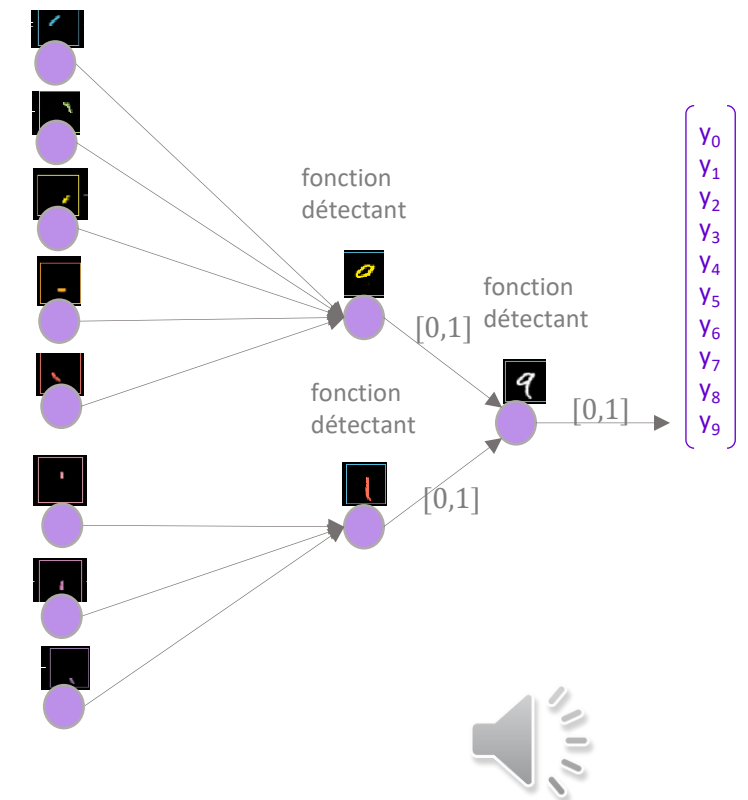
$$y = f_\theta(x)$$

- Données : scans de code postaux d'enveloppes

- Difficulté :
  - Il n'est pas possible d'énumérer tous les motifs possibles correspondant à un seul chiffre (épaisseur, inclinaison, etc.).
  - → Une approche de ML va permettre de ne pas faire d'hypothèse sur les motifs à détecter pour créer la fonction modèle
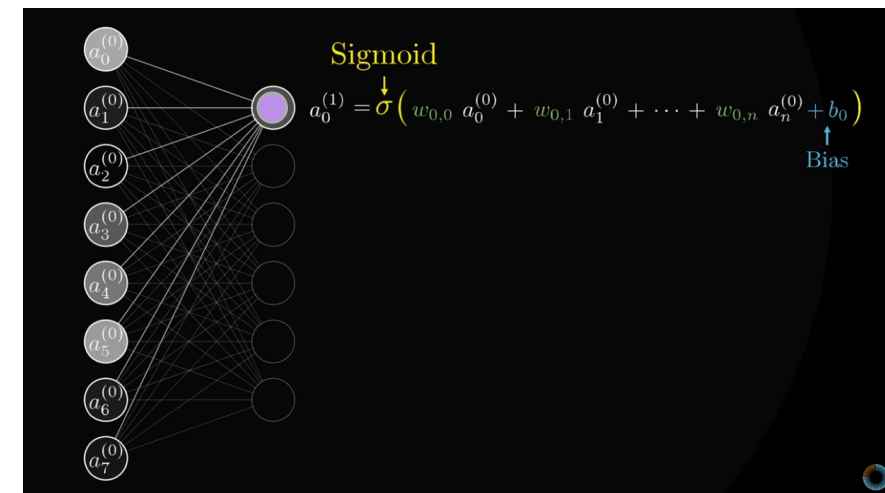
# Un réseau de neurones artificiels

- On peut détecter un motif
  - par une fonction très simple, appelée neurone artificiel
  - et basé sur une simple moyenne pondérée
- Rappelons-nous l'idée d'avoir des petites fonctions :
  - Pour détecter des motifs simples
  - Pour détecter des motifs complexes en composant les détections de motifs simples
- → Ceci est un **réseau de neurones artificiels**

# Alors comment trouver les paramètres de toutes les fonctions neurones artificiels

- Rappel : les paramètres

- Comment ne pas les déterminer à la main ?
  - →On se base sur des données d'exemple :

- Pour les images de chiffres à reconnaître :
  Un jeu de données est constitué des images et de leurs étiquettes (le chiffre correspondant), annotées par des humains.

- On utilise les données d'exemples, paires (***x,y***), pour trouver les paramètres de tous les neurones :
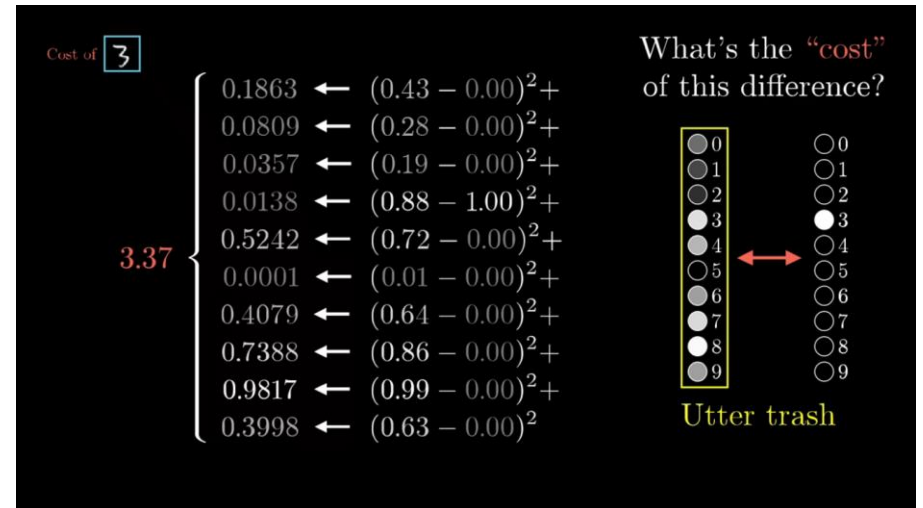  - c'est **l'entraînement**





©3Blue1Brown

# Comment entrainer : comparer les sorties du modèle à la sortie désirée

- Avec les labels de *vérité terrain*

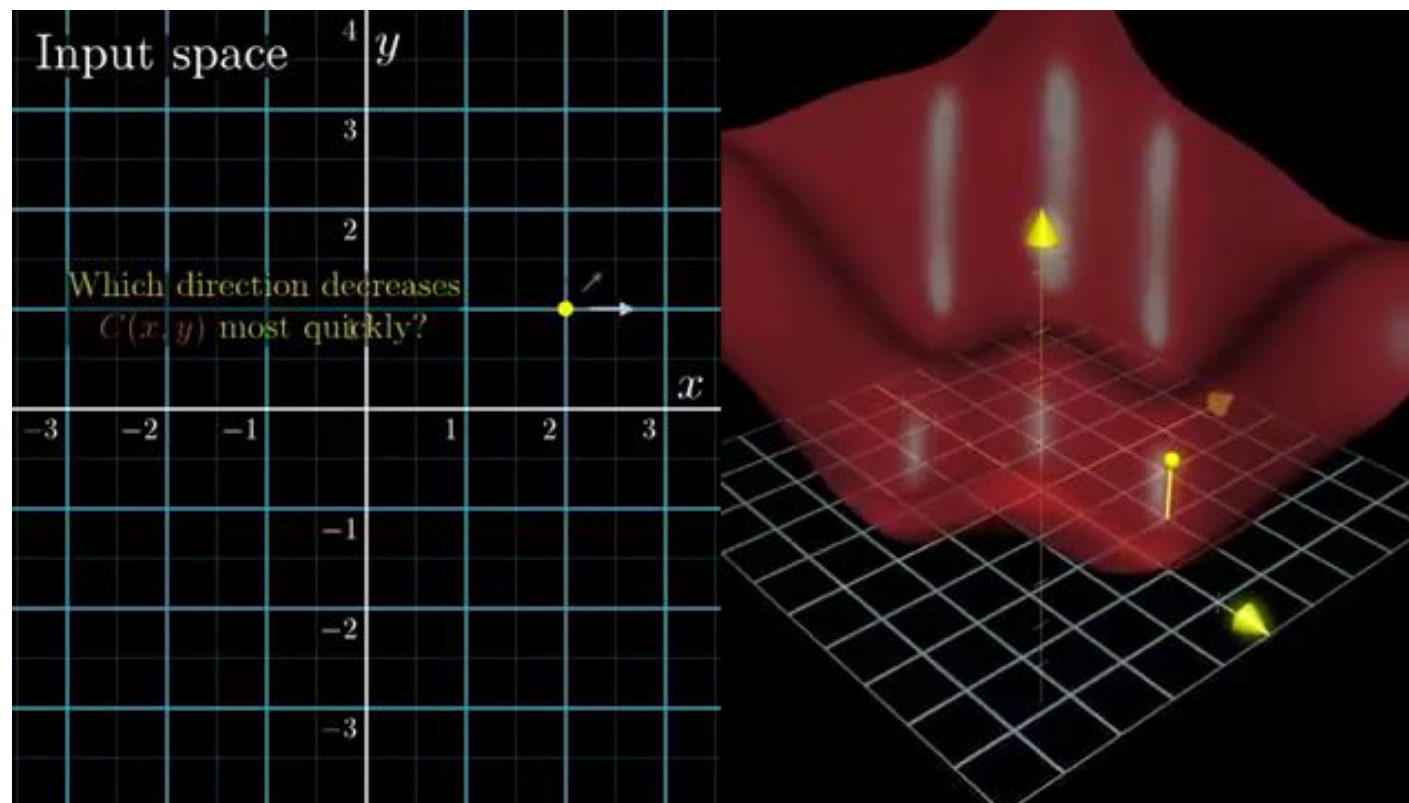→ Calculer le coût/erreur totale pour chaque exemple d'entrainement

- Puis:

© 3Blue1Brown

Fonction de coût : $C_D(\underline{\theta}) = \dfrac{1}{2M} \sum\limits_{i=1}^{M} \left( f_\theta(x_i) - y_i \right)^2$

Objectif de l'entraînement : $\underline{\theta}^* = \arg\min\limits_{\underline{\theta}} C_D(\underline{\theta})$

Principe de la méthode : chercher où la dérivée $C_D'(\theta)$ s'annule

$$\text{Fonction de coût} : \quad C_D(\theta) = \frac{1}{2M} \sum_{i=1}^{M} \left( f_\theta(x_i) - y_i \right)^2$$



© 3Blue1Brown

# Entraînement : Adapter les poids



- Pour aller plus loin : formule de la rétro-propagation et descente de gradient stochastique
  - https://www.youtube.com/watch?v=Ilg3gGewQ5U
  - https://www.youtube.com/watch?v=tIeHLnjs5U8

© 3Blue1Brown

# Un nouveau type de RNA : *Convolutional Neural Networks (CNN)* – Deep Learning

- Moins de paramètres pour mieux décrire les motifs visuels
  - Grâce notamment à l'invariance par translation



$$\tilde{a}^m_{x,y} = \sum^{L}_{-L} \sum^{L}_{-L} f_{i,j} \, a\left(x+i, y+j\right)$$

Motifs à apprendre

Taken from D2lAI

# Tâches de vision : du ML au Deep Learning

- Mise en contexte : avant on pré-déterminait les motifs qui nous paraissaient importants, et on décrivait les données ainsi, pour seulement les séparer en classes avec de l'apprentissage

- Les réseaux de neurones convolutionnels (CNN, combiné à la descente de gradient, à la puissance de calcul et à la quantité de données) permettent à présent de trouver des représentations pertinentes pour classer les données dans les catégories souhaitées

→**Apprentissage de représentation grâce au Deep Learning**

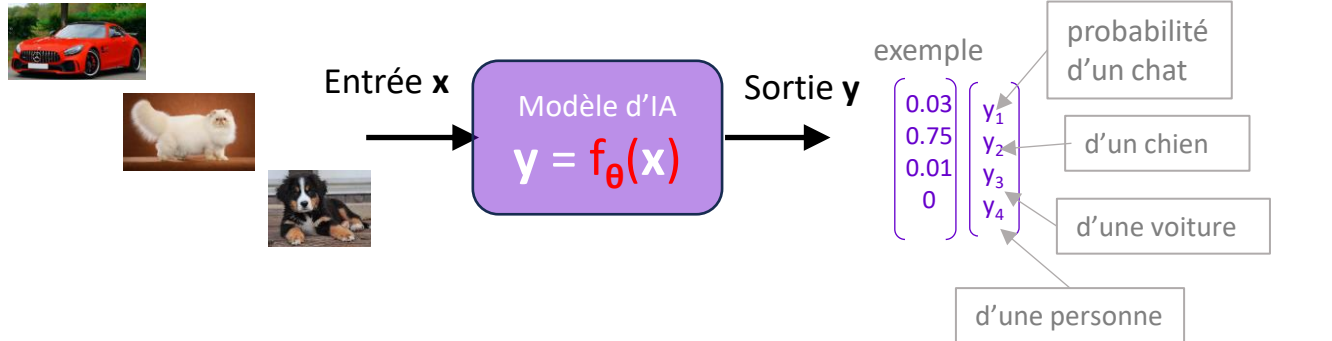→**Mais... ces motifs sont identifiés par des corrélations/associations dans les données...**

# Plan

1. Formalisation du principe de l'apprentissage machine et des réseaux de neurones artificiels

2. Un peu plus de précisions sur ChatGPT

3. Eléments sur les LLMs pour les tâches de raisonnement

4. Eléments sur le marché du travail de la programmation

Entrée **x**

Modèle d'IA

$y = f_\theta(x)$

Sortie **y**

exemple

$\begin{pmatrix} 0.03 \\ 0.75 \\ 0.01 \\ 0 \end{pmatrix}$ $\begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{pmatrix}$

probabilité d'un chat

d'un chien

d'une voiture

d'une personne

Les élèves ouvrent leurs

[.......]
[.......]
[.......]
[.......]

$\begin{pmatrix} 0.03 \\ 0.01 \\ 0 \\ . \\ . \\ . \\ . \\ . \\ . \\ . \\ . \\ 0 \end{pmatrix}$ à abbé

coeur
four
cahier
tiroir
zygote

zygote

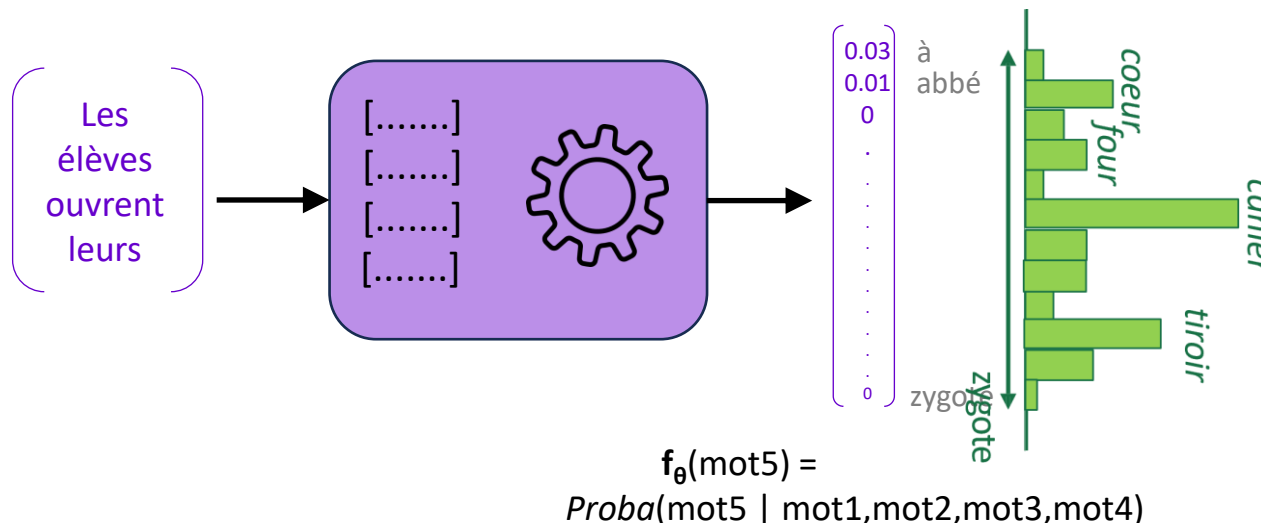$f_\theta(mot5) =$
*Proba*(mot5 | mot1,mot2,mot3,mot4)

Choix de se baser sur J. R. Firth 1957 : le sens d'un mot est donné par son contexte
→ Si on connait les mots entourant un autre mot, on devrait donc pouvoir retrouver ce mot
→ Choix très simplificateur
→ mais très pratique pour utiliser le ML pour trouver les nombres représentant les mots : créer une fonction qui va transformer les mots en nombre pour retrouver un mot à partir de ses voisin :

Stratégie délibérée et simplificatrice :
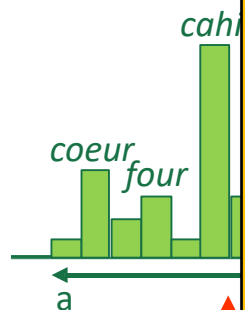retrouver le mot à partir de son contexte

pour arriver à concevoir un modèle de ML qui **reproduit les statistiques de co-occurrences** telles que présentes dans les textes d'entraînement

14

Au cœur des LLM : les réseaux Transformers

$fp_\mu(fr_\theta(\text{texte}))$

Reproduction de motifs de corrélation existants → Automatisation des biais

Beaucoup de données → Travail humain

Beaucoup de calculs → Environnement

Les élèves ouvrent

tiroirs

coeur four cahier

$e^1_4$

$e^0_4$

Prédiction proba
$fp_\mu()$

$w_1$ $w_2$ $w_3$ $w_4$

Les élèves ouvrent leurs

UNIVERSITÉ CÔTE D'AZUR

EFELIA
ÉCOLE FRANÇAISE DE
L'INTELLIGENCE ARTIFICIELLE
CÔTE D'AZUR

FRANCE 2030

# Récap : l'apprentissage de représentation à venir : les réseaux Transformer

- On flexibilise les motifs recherchés encore plus : ils peuvent dépendre des mots ou pixels voisins !

(Multi-head self attention) permet des motifs/noyaux définis sur de grande fenêtres et spécifiques aux données

$$\mathbf{e}_i^{1,h} = \sum_j \mathrm{softmax}\left(\frac{\mathbf{q}_{i,h} \cdot \mathbf{k}_{j,h}}{\sqrt{d}}\right)\mathbf{v}_{j,h}$$

$$\frac{e^{\mathbf{q}_{i,h} \cdot \mathbf{k}_{j,h}}}{\sum_l e^{\mathbf{q}_{l,h} \cdot \mathbf{k}_{l,h}}}$$

*Mais lourd en calcul en test aussi, pas que en entraînement comme avant !*

$$\mathbf{q}_{i,h} = \mathbf{W}^{q,h}\mathbf{e}_i^0, \quad \mathbf{k}_{j,h} = \mathbf{W}^{k,h}\mathbf{e}_j^0, \quad \mathbf{v}_{j,h} = \mathbf{W}^{v,h}\mathbf{e}_j^0$$

$$\mathbf{e}_i^1 = MLP\left(\mathrm{Linear}\left(\mathbf{e}_i^{1,1}, \ldots, \mathbf{e}_i^{1,H}\right)\right)$$

- Le mot i est représenté par une recombinaison de (diverses représentations de) ses mots voisins, dont les facteurs varient eux-mêmes en fonction des mots voisins (pas comme avant).

A. Vaswani et al.. Attention is all you need. NeurIPS 2017.

# ChatGPT : ce qu'on connaît du fonctionnement

UNIVERSITÉ CÔTE D'AZUR

EFELIA
ÉCOLE FRANÇAISE DE
L'INTELLIGENCE ARTIFICIELLE
CÔTE D'AZUR

FRANCE 2030

- Architecture : Transformers
- Paramètres et données croissantes :



**GPT-2**
- 1.5 billion parameters
- 40 GB text training dataset
- Often fine-tuned to perform specific tasks
- Smaller version of the model was released to the public open source

**GPT-3**
- 176 billion parameters
- 570 GB training dataset comprising of books, articles, websites, and more
- Ability to perform most language tasks without additional tuning
- Launched as an API service



Figure 1: The Transformer - model architecture.

[1] A. Vaswani et al., "Attention Is All You Need.", Neural Information Processing System, 2017.
https://towardsdatascience.com/how-chatgpt-works-the-models-behind-the-bot-1ce5fca96286

# Decoder (GPT) → Instructions (Instruct GPT)

1. Nous savons maintenant que le pré-entraînement d'un LLM implique une procédure d'entrainement où il apprend à générer un mot à la fois.
Le LLM pré-entraîné qui en résulte est capable de compléter du texte, ce qui signifie qu'il peut terminer des phrases ou écrire des paragraphes de texte à partir d'un fragment en entrée.

2. Cependant, les LLM pré-entraînés ont souvent du mal avec des instructions spécifiques, telles que « Corrigez la grammaire de ce texte » ou « Convertissez ce texte en voix passive ».
**Nous allons donc nous concentrer sur l'amélioration de la capacité du LLM à suivre de telles instructions et à générer une réponse souhaitée.**

©F. Precioso

# Exemples d'instructions traitées par un LLM pour générer les réponses souhaitées

**The instructions serve as inputs for the LLM.**

**The goal for the LLM is to generate a desired response.**

**Instruction**

**Desired response**

Convert 45 kilometers to meters.   ⟶   45 kilometers is 45000 meters.

Provide a synonym for "bright."   ⟶   A synonym for "bright" is "radiant."

Edit the following sentence to remove all passive voice: "The song was composed by the artist."   ⟶   The artist composed the song.

# ChatGPT : ce qu'on connaît du fonctionnement (à partir de InstructGPT)



**SFT** (Supervised Fine Tuning) model

https://towardsdatascience.com/how-chatgpt-works-the-models-behind-the-bot-1ce5fca96286

# ChatGPT : ce qu'on connaît du fonctionnement



❶ Training

Parameters:175B
Training tokens:300B
Vocab size: ~50K

Stage 1: Pre-training

two plus two is → GPT 3.5 → Complete the input sentence → equal to four.

Internet data(300B tokens)

Stage 2: Fine-tuning

GPT 3.5 → collect data and fine-tune → Finetuned GPT → train a reward model → Reward model → optimize using reinforcement learning (PPO algorithm) → ChatGPT Model

Demonstration data

RM

PPO



❷ Answer a prompt

a new prompt

LeetCode: longest common string

❶ Content moderation

❷ Safe?

3.N

3.Y

ChatGPT Model → Generate response

❹ Content moderation

❺ Safe?

6.N → Template response generation

6.Y

Response

longest common string is the problem of…

Sorry, I am not trained to provide medical advice.

- **Sans oublier beaucoup de contrôle du contenu :**
  - Modèle spécifique développé avec travail d'annotation spécifique
  - Enquête du Time montre pratiques d'exploitation humaine

https://blog.bytebytego.com/p/ep-44-how-does-chatgpt-work

# Techniques de Prompt Engineering

- **Notions de base**
  - Jeux de rôle
  - Raffinage itératif
  - Boucles de rétroaction (dialogue)

- **Avancé**
  - Zero-shot invite
  - Few-shot invite / apprentissage-en-contexte
  - Chaîne de pensée (CoT)

- **Et tant d'autres techniques...**
  - Encourager l'introspection
  - Incitation à un stimulus directionnel
  - Arbre de pensée (version structurée de la chaîne de pensée)
  - Incitation basée sur la complexité
  - Incitation à la production de connaissances...
  - Incitation la plus faible à la plus forte

https://www.datacamp.com/blog/what-is-prompt-engineering-the-future-of-ai-communication

https://aws.amazon.com/fr/what-is/prompt-engineering/#:~:text=Prompt%20engineering%20is%20the%20process,high%2Dquality%20and%20relevant%20output.

count the legs

The zebra in the image appears to have an optical illusion due to its leg positions, but it actually has four legs. ❌

ChatGPT 5 ⌄                                    ⬆ Share  ···



How many overlapping circles are there in the logo of this car? Answer with a number in curly brackets, e.g., {9}.

{4} ❌

ChatGPT 5 ⌄                                    ⬆ Share  ···



count the legs

The duck in the image has two legs. ❌

# Do LLMs Have Visualization Literacy? An Evaluation on Modified Visualizations to Test Generalization in Data Interpretation

Jiayi Hong (iD), Christian Seto (iD), Arlen Fan (iD), and Ross Maciejewski (iD)

- LLMs we explored currently fail to achieve the same levels of visualization literacy when compared to data from the general public reported in VLAT, and LLMs heavily relied on their pre-existing knowledge to answer questions instead of utilizing the information provided by the visualization when answering questions.



(c) Stacked Bar Chart

(d) 100% Stacked Bar Chart

(g) Scatterplot

(h) BubbleChart

# Is vibe coding dying?

Amateurs might not be replacing teams of coders, after all

**GARY MARCUS**
OCT 22, 2025

Remember how in October and in March I told you that vibe coding — in the sense of amateurs using large language models to write code to "build products that would have previously required teams of engineers" — would never be remotely reliable? And that such tools were fine for demos but not for complex apps in the real world? And that the code they wrote would be hard to maintain?

Customers are finally figuring that out.



DevOps & Code Completion Tools Traffic

It is deeply unserious and these tools aren't delivering when they encounter real world complexity (building quick demos isn't complex) in any meaningful enterprise.

The problem, as always, lies in generalizing outside the training distribution. Vibe coding can be fine if you are building something very familiar, but is less reliable for the unfamiliar. Even Andrej Karpathy, who literally invented the term vibe-coding, sees this:

**Gary Marcus** @GaryMarcus · 6d
The very inventor of the term "**vibe** coding", hand-coding.

And confirming — yet again - that current AI has NOT solved distribution shift (the core problem that I have been harping on since 1998).

**Andrej Karpathy** @karpathy · 10/13/25
Replying to @zenitsu_aprntc

Good question, it's basically entirely hand-written (with tab autocomplete). I tried to use claude/codex agents a few times but they just didn't work well enough at all and net unhelpful, possibly the repo is too far off the data distribution.

# Plan

1. Formalisation du principe de l'apprentissage machine et des réseaux de neurones artificiels

2. Un peu plus de précisions sur ChatGPT

3. Eléments sur les LLMs pour les tâches de raisonnement

4. Eléments sur le marché du travail de la programmation

Drawing (and typo) by ChatGPT

©G. Marcus

# Brains vs. Bytes: Evaluating LLM Proficiency in Olympiad Mathematics

**Hamed Mahdavi, Pegah Mohammadipour, Samira Malek, ,Vasant Honavar**
Pennsylvania State University
University Park, PA, USA
{hmm5834,pegahmp,sxm6547,vhonavar}@psu.edu

**Alireza Hashemi**
City University of New York
New York, NY, USA
alireza.hashemi13@outlook.com

**Majid Daliri**
New York University
New York, NY, USA
daliri.majid@nyu.edu

**Alireza Farhadi**
Amirkabir University of Technology
Tehran, Iran
farhadi@aut.ac.ir

**Yekta Yazdanifard**
Bocconi University
Milan, Italy
yekta.yazdanifard@unibocconi.it

**Amir Khasahmadi**
Autodesk
Toronto, Canada
amir.khasahmadi@autodesk.com

## Abstract

Recent advancements in large language models (LLMs) have shown impressive progress in mathematical reasoning tasks. However, current evaluation benchmarks predominantly focus on the accuracy of final answers, often overlooking the logical rigor crucial for mathematical problem-solving. The claim that state-of-the-art LLMs can solve Math Olympiad-level problems requires closer examination. To explore this, we conducted both qualitative and quantitative human evaluations of proofs generated by LLMs, and developed a schema for automatically assessing their reasoning capabilities. Our study reveals that current LLMs fall significantly short of solving challenging Olympiad-level problems and frequently fail to distinguish correct mathematical reasoning from clearly flawed solutions. We also found that occasional correct final answers provided by LLMs often result from pattern recognition or heuristic shortcuts rather than genuine mathematical reasoning. These findings underscore the substantial gap between LLM performance and human expertise in advanced mathematical reasoning and highlight the importance of developing benchmarks that prioritize the rigor and coherence of mathematical arguments rather than merely the correctness of final answers.

27

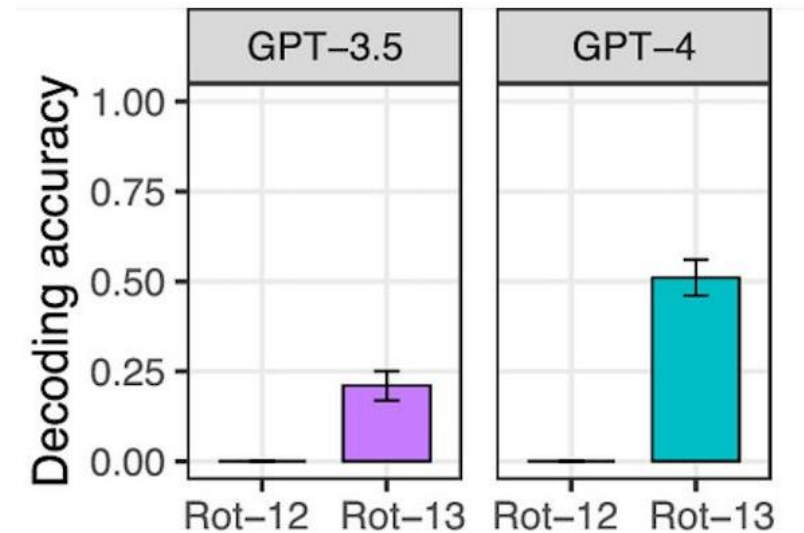# Mais reproduire les co-occurrences de mots a ses limites



Shift ciphers

Decode by shifting each letter **13** positions backward in the alphabet.

**Input:** Jryy, vg jnf abg rknpgyl cynaarq sebz gur ortvaavat.
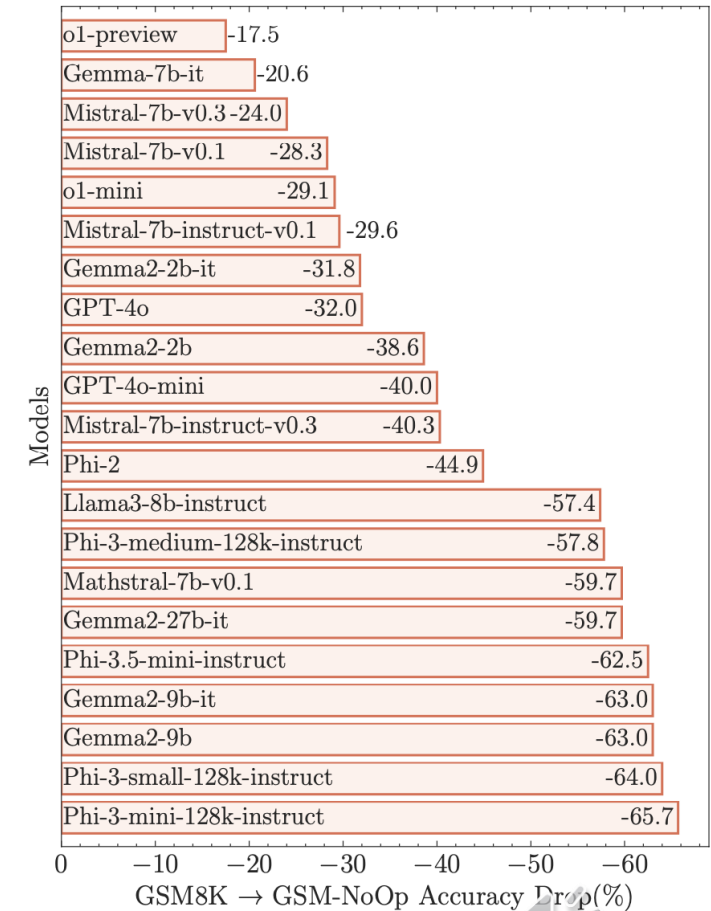**Correct:** Well, it was not exactly planned from the beginning.
✓ **GPT-4:** Well, it was not exactly planned from the beginning.

R. T. McCoy et al., "Embers of autoregression show how large language models are shaped by the problem they are trained to solve," *Proc. Natl. Acad. Sci. U.S.A.*, vol. 121, no. 41, Oct. 2024.

# Mais reproduire les co-occurrences de mots a ses limites

Oliver picks 44 kiwis on Friday. Then he picks 58 kiwis on Saturday. On Sunday, he picks double the number of kiwis he did on Friday, but five of them were a bit smaller than average. How many kiwis does Oliver have?

- Les LLM sont moins performants pour les tâches rares que pour les tâches courantes
- Performances très variables d'une instanciation à l'autre de la même question.

→ Prudence si on veut les utiliser pour des tâches qui sont rares dans les données d'entraînement

→ Limites importantes de la capacité des LLM à effectuer un véritable raisonnement mathématique



| Models | GSM8K → GSM-NoOp Accuracy Drop(%) |
|---|---|
| o1-preview | -17.5 |
| Gemma-7b-it | -20.6 |
| Mistral-7b-v0.3 | -24.0 |
| Mistral-7b-v0.1 | -28.3 |
| o1-mini | -29.1 |
| Mistral-7b-instruct-v0.1 | -29.6 |
| Gemma2-2b-it | -31.8 |
| GPT-4o | -32.0 |
| Gemma2-2b | -38.6 |
| GPT-4o-mini | -40.0 |
| Mistral-7b-instruct-v0.3 | -40.3 |
| Phi-2 | -44.9 |
| Llama3-8b-instruct | -57.4 |
| Phi-3-medium-128k-instruct | -57.8 |
| Mathstral-7b-v0.1 | -59.7 |
| Gemma2-27b-it | -59.7 |
| Phi-3.5-mini-instruct | -62.5 |
| Gemma2-9b-it | -63.0 |
| Gemma2-9b | -63.0 |
| Phi-3-small-128k-instruct | -64.0 |
| Phi-3-mini-128k-instruct | -65.7 |

I. Mirzadeh et al., "GSM-Symbolic: Understanding the Limitations of Mathematical Reasoning in Large Language Models," ICLR, 2025.

- Grâce à des expérimentations poussées sur divers casse-tête, nous montrons que les LRM frontières sont confrontés à un effondrement complet de leur exactitude au-delà de certaines complexités.

- Nous avons constaté que les LRM présentent des limites en matière de calcul exact : ils n'utilisent pas d'algorithmes explicites et raisonnent de manière incohérente d'une énigme à l'autre.

- Nous examinons également plus en détail les traces de raisonnement, en étudiant les schémas des solutions explorées et en analysant le comportement calculatoire des modèles, mettant en lumière leurs points forts et leurs limites, et soulevant ainsi des questions cruciales sur leurs véritables capacités de raisonnement.

# The Illusion of Thinking:
## Understanding the Strengths and Limitations of Reasoning Models via the Lens of Problem Complexity

Parshin Shojaee*†   Iman Mirzadeh*   Keivan Alizadeh
Maxwell Horton   Samy Bengio   Mehrdad Farajtabar

Apple

## Abstract

Recent generations of frontier language models have introduced Large Reasoning Models (LRMs) that generate detailed thinking processes before providing answers. While these models demonstrate improved performance on reasoning benchmarks, their fundamental capabilities, scaling properties, and limitations remain insufficiently understood. Current evaluations primarily focus on established mathematical and coding benchmarks, emphasizing final answer accuracy. However, this evaluation paradigm often suffers from data contamination and does not provide insights into the reasoning traces' structure and quality. In this work, we systematically investigate these gaps with the help of controllable puzzle environments that allow precise manipulation of compositional complexity while maintaining consistent logical structures. This setup enables the analysis of not only final answers but also the internal reasoning traces, offering insights into how LRMs "think". Through extensive experimentation across diverse puzzles, we show that frontier LRMs face a complete accuracy collapse beyond certain complexities. Moreover, they exhibit a counter-intuitive scaling limit: their reasoning effort increases with problem complexity up to a point, then declines despite having an adequate token budget. By comparing LRMs with their standard LLM counterparts under equivalent inference compute, we identify three performance regimes: (1) low-complexity tasks where standard models surprisingly outperform LRMs, (2) medium-complexity tasks where additional thinking in LRMs demonstrates advantage, and (3) high-complexity tasks where both models experience complete collapse. We found that LRMs have limitations in exact computation: they fail to use explicit algorithms and reason inconsistently across puzzles. We also investigate the reasoning traces in more depth, studying the patterns of explored solutions and analyzing the models' computational behavior, shedding light on their strengths, limitations, and ultimately raising crucial questions about their true reasoning capabilities.

## 1 Introduction

Large Language Models (LLMs) have recently evolved to include specialized variants explicitly designed for reasoning tasks—Large Reasoning Models (LRMs) such as OpenAI's o1/o3 [1, 2], DeepSeek-R1 [3], Claude 3.7 Sonnet Thinking [4], and Gemini Thinking [5]. These models are new artifacts, characterized by their "*thinking*" mechanisms such as long Chain-of-Thought (CoT) with self-reflection, and have demonstrated promising results across various reasoning benchmarks. Their

---

*Equal contribution.
†Work done during an internship at Apple.
{p_shojaee, imirzadeh, kalizadehvahid, mchorton, bengio, farajtabar}@apple.com

Hanoi is a classic game with three pegs and multiple discs in which you need to move all the discs on the left peg to the right peg, never stacking a larger disc on top of a smaller one.

(You can try a digital version at mathisfun.com.)

If you have never seen it before, it takes a moment or to get the hang of it. (Hint, start with just a few discs).

With practice, a bright (and patient) seven-year-old can do it. And it's trivial for a computer. Here's a computer solving the seven-disc version, using an algorithm that any intro computer science student should be able to write:

Claude, on the other hand, can barely do 7 discs, getting less than 80% accuracy, left bottom panel below, and pretty much can't get 8 correct at all.

©Gary Marcus, https://garymarcus.substack.com/p/a-knockout-blow-for-llms?publication_id=888615

Josh Wolfe ✔ @wolfejosh · 13h
2/ Apple tested today's "reasoning" AIs like Claude + DeepSeek which look smart—but when complexity rises, they collapse.

Not fail gracefully. Collapse completely.

💬 3    ↻ 17    ❤ 264    📊 22K

Josh Wolfe ✔ @wolfejosh · 13h
3/ They found LLMs don't scale reasoning like humans do.

They think MORE up to a point...

Then they GIVE UP early, even when they have plenty of compute left.

💬 2    ↻ 11    ❤ 217    📊 22K

Josh Wolfe ✔ @wolfejosh · 13h
4/ Even when handed the exact algorithm, LLMs still botch the job.

Execution ≠ understanding.

Its not "missing creativity"—its failing basic logic.

💬 3    ↻ 9    ❤ 176    📊 21K

Josh Wolfe ✔ @wolfejosh · 13h
5/ models "overthink" EASY problems—exploring WRONG answers after finding the RIGHT one.

And when problems get HARDER... they think LESS.

Wasted compute at one end—defeatism at the other

💬 6    ↻ 10    ❤ 198    📊 24K

Josh Wolfe ✔ @wolfejosh · 13h
6/ Apple's take is these models ARE NOT reasoning.

they're super expensive pattern matchers that break as soon as we step outside their training distribution...

💬 22    ↻ 45    ❤ 533    📊 27K

Figure 1: **Top**: Our setup enables verification of both final answers and intermediate reasoning traces, allowing detailed analysis of model thinking behavior. **Bottom left & middle**: At low complexity, non-thinking models are more accurate and token-efficient. As complexity increases, reasoning models outperform but require more tokens—until both collapse beyond a critical threshold, with shorter traces. **Bottom right**: For correctly solved cases, Claude 3.7 Thinking tends to find answers early at low complexity and later at higher complexity. In failed cases, it often fixates on an early wrong answer, wasting the remaining token budget. Both cases reveal inefficiencies in the reasoning process.

# STOP ANTHROPOMORPHIZING INTERMEDIATE TOKENS AS REASONING/THINKING TRACES!

**Subbarao Kambhampati**    **Kaya Stechly**    **Karthik Valmeekam**    **Lucas Saldyt**    **Siddhant Bhambri**

**Vardhan Palod**    **Atharva Gundawar**    **Soumya Rani Samineni**    **Durgesh Kalwar**    **Upasana Biswas**

**School of Computing & AI**
**Arizona State University**

## ABSTRACT

Intermediate token generation (ITG), where a model produces output before the solution, has been proposed as a method to improve the performance of language models on reasoning tasks. These intermediate tokens have been called "reasoning traces" or even "thoughts" – implicitly anthropomorphizing the model, implying these tokens resemble steps a human might take when solving a challenging problem. In this paper, we present evidence that this anthropomorphization isn't a harmless metaphor, and instead is quite dangerous – it confuses the nature of these models and how to use them effectively, and leads to questionable research.

## 1 Introduction

Recent advances in general planning and problem solving have been spearheaded by so-called "Long Chain-of-Thought" models, most notably DeepSeek's R1 [17]. These transformer-based large language models are further post-trained using iterative fine-tuning and reinforcement learning methods. Following the now-standard teacher-forced pre-training, instruction fine-tuning, and preference alignment stages, they undergo additional training on reasoning tasks: at each step, the model is presented with a question; it generates a sequence of intermediate tokens (colloquially or perhaps fancifully called a "Chain of Thought" or "reasoning trace"); and it ends it with a specially delimited answer sequence. After verification of this answer sequence by a formal system, the model's parameters are updated so that it is more likely to output sequences that end in correct answers and less likely to output those that end in incorrect answers with no guarantees of trace correctness.

While (typically) no direct optimization pressure is applied to the intermediate tokens [4, 62], empirically it has been observed that language models perform better on many domains if they output such tokens first [38, 55, 61, 19, 16, 17, 39, 36, 29]. While the fact of the performance increase is well-known, the reasons for it are less clear. Much of the previous work has framed intermediate tokens in wishful anthropomorphic terms, claiming that these models are "thinking" before outputting their answers [38, 12, 17, 56, 62, 7]. The traces are thus seen both as giving insights to the end users about the solution quality, and capturing the model's "thinking effort."

In this paper, we take the position that anthropomorphizing intermediate tokens as reasoning/thinking traces is (1) wishful (2) has little concrete supporting evidence (3) engenders false confidence and (4) may be pushing the community into fruitless research directions. This position is supported by work questioning the interpretation of intermediate tokens as reasoning/thinking traces (Section 4) and by stronger alternate explanations for their effectiveness (Section 6).

Anthropomorphization has long been a contentious issue in AI research [33], and LLMs have certainly increased our anthropomorphization tendencies [20]. While some forms of anthropomorphization can be treated rather indulgently as harmless and metaphorical, our view is that viewing ITG as reasoning/thinking is more serious and may give a false sense of model capability and correctness.

- Il n'y a qu'une faible corrélation entre l'exactitude de la trace (prise comme "raisonnement") et l'exactitude du résultat final.

- Pire, entrainer des modèles sur des traces de raisonnement fausses améliore leur performance sur le résultat final.

> Étant donné que ces traces peuvent n'avoir aucun sens, les faire délibérément apparaître comme du raisonnement humain est dangereux. En fin de compte, les LRM sont censés fournir des solutions que les utilisateurices ne connaissent pas déjà (et qu'iels ne sont peut-être même pas capables de vérifier directement). Encourager à voir ces traces de supposé raisonnement, dont seulement le style est plausible, comme motif de confiance semble bien malavisé !

- Après tout, la dernière chose que nous voulons faire est de concevoir des systèmes d'IA qui sont juste puissants pour exploiter nos failles congnitives en nous convaincant de la validité de réponses incorrectes.

- We show that verbalised chains are frequently unfaithful, diverging from the true hidden computations that drive a model's predictions, and giving an incorrect picture of how models arrive at conclusions.

- Despite this, CoT is increasingly relied upon in high-stakes domains such as medicine, law, and autonomous systems—our analysis of 1,000 recent CoT-centric papers finds that ~25% explicitly treat CoT as an interpretability technique—and among them, papers in high-stakes domains specifically hinge on such interpretability claim heavily.

- Proposal: develop causal validation methods (e.g., activation patching, counterfactual interventions, verifier models) to ground explanations in model internals.



Figure 1: Overview of our paper: Unfaithful Chain-of-Thought behaviors (left), their mechanistic and cognitive underpinnings (center), and our proposed research roadmaps for enhancing CoT faithfulness (right).

33

- Nos résultats révèlent que le raisonnement Chain of Thought (CoT – prompt « pense étape par étape) fonctionne efficacement lorsqu'il est appliqué à des données similaires (d'à peu près la même distribution statistique que les données d'entraînement), mais devient fragile et sujet à l'échec même en cas de changements de distribution modérés. Dans certains cas, les LLM génèrent des étapes de raisonnement fluides, mais logiquement incohérentes. Les résultats suggèrent que ce qui semble être un raisonnement structuré peut être un mirage, émergeant de motifs mémorisés ou interpolés dans les données d'apprentissage plutôt que d'une inférence logique.

- Ces résultats suggèrent que les LLM ne sont pas des raisonneurs mais plutôt des simulateurs sophistiqués de textes ressemblant à du raisonnement.

- Pour les usagères et usagers, nos résultats soulignent le risque de s'appuyer sur le CoT comme solution clé en main pour les tâches de raisonnement et mettent en garde contre toute assimilation des résultats de type CoT à la pensée humaine.

- Pour les chercheuses et chercheurs, ces résultats soulignent le défi non résolu de parvenir à un raisonnement à la fois fiable et généralisable, d'où la nécessité de développer des modèles capables d'aller au-delà de la reconnaissance de formes superficielle pour démontrer une compétence inférentielle plus approfondie.

---

arXiv:2508.01191v3 [cs.AI] 13 Aug 2025

# Is Chain-of-Thought Reasoning of LLMs a Mirage? A Data Distribution Lens

Chengshuai Zhao[1], Zhen Tan[1], Pingchuan Ma[1], Dawei Li[1], Bohan Jiang[1], Yancheng Wang[1], Yingzhen Yang[1] and Huan Liu[1]

[1]Arizona State University, USA

Chain-of-Thought (CoT) prompting has been shown to improve Large Language Model (LLM) performance on various tasks. With this approach, LLMs appear to produce human-like reasoning steps before providing answers (a.k.a. CoT reasoning), which often leads to the perception that they engage in deliberate inferential processes. However, some initial findings suggest that CoT reasoning may be more superficial than it appears, motivating us to explore further. In this paper, we study CoT reasoning via a data distribution lens and investigate if CoT reasoning reflects a structured inductive bias learned from in-distribution data, allowing the model to conditionally generate reasoning paths that approximate those seen during training. Thus, its effectiveness is fundamentally bounded by the degree of distribution discrepancy between the training data and the test queries. With this lens, we dissect CoT reasoning via three dimensions: *task*, *length*, and *format*. To investigate each dimension, we design DATAALCHEMY, an isolated and controlled environment to train LLMs from scratch and systematically probe them under various distribution conditions. Our results reveal that CoT reasoning is a brittle mirage that vanishes when it is pushed beyond training distributions. This work offers a deeper understanding of *why* and *when* CoT reasoning fails, emphasizing the ongoing challenge of achieving genuine and generalizable reasoning. Our code is available at GitHub: https://github.com/ChengshuaiZhao0/DataAlchemy.

## 1. Introduction

Recent years have witnessed Large Language Models' (LLMs) dominant role in various domains (Li et al., 2025b; Ting et al., 2025; Zhao et al., 2025, 2023) through versatile prompting techniques (Kojima et al., 2022; Wei et al., 2022; Yao et al., 2023). Among these, Chain-of-Thought (CoT) prompting (Wei et al., 2022) has emerged as a prominent method for eliciting structured reasoning from LLMs (a.k.a. CoT reasoning). By appending a simple cue such as "Let's think step by step," LLMs decompose complex problems into intermediate steps, producing outputs that resemble human-like reasoning. It has been shown to be effective in tasks requiring logical inference(Xu et al., 2024), mathematical problem solving (Imani et al., 2023), and commonsense reasoning (Wei et al., 2022). The empirical successes of CoT reasoning lead to the perception that LLMs engage in deliberate inferential processes (Ling et al., 2023; Yu et al., 2023; Zhang et al., 2024a,c).



Figure 1 | The data perspective lens. CoT reasoning's effectiveness is fundamentally bounded by the degree of distribution discrepancy between the training data and the test queries. Guided by this lens, we dissect CoT reasoning via three dimensions: *task*, *length*, and *format*.

# Plan

1. Formalisation du principe de l'apprentissage machine et des réseaux de neurones artificiels

2. Un peu plus de précisions sur ChatGPT

3. Eléments sur les LLMs pour les tâches de raisonnement

4. Eléments sur le marché du travail de la programmation

# What's really going on with AI and jobs?

Record-breaking layoff reports, Amazon's mass firings, and a slump in entry level employment. Is AI behind it all?

BRIAN MERCHANT
NOV 06, 2025

is down.) The most-discussed, however, is probably the shrinking number of jobs for recent college grads. Derek Thompson pointed to this trend in an Atlantic piece that argued there were signs that "AI is competing with recent college grads" and a trio of Stanford economists published a paper asserting that early career employment for US workers in "occupations exposed generative AI" aged 22-25 had declined in key fields 13% since 2022, precisely when the commercial technology entered the scene.

What we can be sure of, however, is that there is real pain unfolding right now, irregardless of whether it's due to management enacting bona fide AI job replacement, executives' *hopes* that AI can cut labor costs, or "AI washing" that obscures a company's ulterior motives. I can be sure of this not just because I've

**Headcount Over Time by Age Group**
**Software Developers (Normalized)**

Legend:
- Early Career 1 (22-25)
- Early Career 2 (26-30)
- Developing (31-34)
- Mid-Career 1 (35-40)
- Mid-Career 2 (41-49)
- Senior (50+)

B. Merchant, "What's really going on with AI and jobs?", Nov. 2025.

E. Brynjolfsson, B. Chandar, and R. Chen, "Canaries in the Coal Mine? Six Facts about the Recent Employment Effects of Artificial Intelligence", Stanford pub., Nov. 2025.

# Les agents IA ont, jusqu'à présent, été pour la plupart un échec

©Gary Marcus, *AI Agents have, so far, mostly been a dud,*
https://garymarcus.substack.com/p/ai-agents-have-so-far-mostly-been?publication_id=888615

# LLMs + Coding Agents = Security Nightmare

©Gary Marcus, https://garymarcus.substack.com/p/llms-coding-agents-security-nightmare?publication_id=888615

# Measuring the Impact of Early-2025 AI on Experienced Open-Source Developer Productivity



Measuring the Impact of Early-2025 AI on Experienced Open-Source Developer Productivity. (2025). *METR Blog*.

AI KILLED MY JOB

# AI Killed My Job: Tech workers

Tech workers at TikTok, Google, and across the industry share stories about how AI is changing, ruining, or replacing their jobs.

**BRIAN MERCHANT**
JUN 25, 2025

♡ 256      💬 42      🔁 69                    Share

"What will AI mean for jobs?" may be the single most-asked question about the technology category that dominates Silicon Valley, pop culture, and our politics. Fears that AI will put us out of work routinely top opinion polls. Bosses are citing AI as the reason they're slashing human staff. Firms like Duolingo and Klarna have laid off workers in loudly touted shifts to AI, and DOGE used its "AI-first" strategy to justify firing federal workers.

Meanwhile, tech executives are pouring fuel on the flames. Dario Amodei, the CEO of Anthropic, claims that AI products like his will soon eliminate half of entry level white collar jobs, and replace up to 20% of all jobs, period. OpenAI's Sam Altman says that AI systems can replace entry level workers, and will soon be able to code "like an experienced software engineer." Elsewhere, he's been blunter, claiming "Jobs are definitely going to go away, full stop."

But the question remains: What's actually happening on the ground, right now? There's no doubt that lots of firms are investing heavily in AI and *trying* to use it to improve productivity and cut labor costs. And it's clear that in certain industries, especially creative ones, the rise of cheap AI-generated content is hitting workers hard. Yet broader economic data on AI impacts suggests a more limited disruption. Two and a half years after the rise of ChatGPT, after a torrent of promises, CEO talk, and think pieces, how is—or isn't—AI *really* reshaping work?

B. Merchant, "AI Killed My Job: Tech workers," June 2025.

41

🔍 📤 Subscribe

# "AI is killing the software engineer discipline"

**Software engineer at Google.**

I have been a software engineer at Google for several years. With the recent introduction of generative AI-based coding assistance tools, we are already seeing a decline in open source code quality [1] (defined as "code churn" - how often a piece of code is written only to be deleted or fixed within a short time). I am also starting to see a downward trend of (a) new engineers' readiness in doing the work, (b) engineers' willingness to learn new things, and (c) engineers' effort to put in serious thoughts in their work.

Specifically, I have recently observed first hand some of my colleagues at the start of their career heavily relying on AI-based coding assistance tools. Their "code writing" consists of iteratively and alternatingly hitting the Tab key (to accept AI-generated code) and watching for warning underlines [2] indicating there could be an error (which have been typically based on static analysis, but recently increasingly including AI-generated warnings). These young engineers - squandering their opportunities to learn how things actually work - would briefly glance at the AI-generated code and/or explanation messages and continue producing more code when "it looks okay."

I also saw experienced engineers in senior positions when faced with an important data modeling task decided to generate the database schema with generative AI. I originally thought it was merely a joke but recently found out that they basically just used the generated schema in actual (internal) services essentially without modification, even if there are some obvious glaring issues. Now those issues have propagated to other code that needs to interact with that database and it will be more costly to fix, so chances are people will just carry on, pretending everything is working as intended.

All of these will result in poorer software quality. "Anyone can write code" sounds good on paper, but when bad code is massively produced, it hurts everyone including those who did not ask for it and have been trusting the software industry.

# "AI makes everything worse"

**Senior developer at a cloud company.**

I work for a cloud service provider (who will retaliate if you don't post this anonymously, unfortunately), and they're absolutely desperate for the current AI fad to be useful for something.

They're completely ignoring the environmental costs (insane power requirements, draining lakes of freshwater for cooling, burning untold CPU and GPU hours that could be dedicated to something useful instead) because there's a buck to be made. They hope. But they're still greenwashing the company of course.

For cloud companies, AI is a gold rush; until the bubble bursts, they can sell ridiculous amounts of expensive server time (lots and lots of CPU/GPU/memory /storage) and tons of traffic to and from the models. They're selling shovels to the gold miners, and are in a great position to charge rent if someone strikes a vein of usefulness.

*I can see a scenario coming fast that's going to set back software development by years*

---

I can see a scenario coming fast that's going to set back software development by years (decades? who knows!):

- C-suite: we don't need these expensive senior developers, interns can code with AI
- C-suite: we don't need these expensive security developers, AI can find the problems
- senior developers are laid off, or quit due to terrible working conditions (we're already seeing this)
- they're replaced with junior developers, fresh out of school... cheap, with no sense of work-life balance, and no families to distract them
- all the vibe coding goes straight to production because, obviously, we trust the AI and don't know any better; also we've been told to use AI for everything
- at some point, all the bugs and security vulnerabilities make everything so bad it actually starts impacting the bottom line
- uh oh, the vibe coders never progressed beyond junior skill levels, so nobody can do the code reviews, nobody can find and fix the security problems
- if all the fired senior developers haven't retired or found other jobs (a lot of these people want to get out of tech, because big tech has made everything terrible), they'll need to be hired back, hopefully at massive premiums due to demand

If these tools were generally useful, they wouldn't need to force them on us, we'd be picking them up and running with them.

# Problématique – Conclusion

- Pour le portail Sciences et Techniques, nous énonçons les éléments à retenir (important !) sur les questions analysées :
  - Comment fonctionnent plus précisément les modèles de ML et ChatGPT en particulier ?
  - Quelles sont les performances des LLMs pour des tâches de raisonnement ?
  - Comment le marché du travail de développement logiciel est modifié par l'arrivée des LLMs ?

Rebecca Winthrop and Maryanne Wolf, "Rethinking School in the Age of AI," Center for Humane Technology, April 2025.
Sonja Drimmer and Christopher J. Nygren, "How We Are Not Using AI in the Classroom, " The Newsletter of the International Center of Medieval Art, April 2025.
Carl T. Bergstrom and Jevin D. West, "Modern-Day Oracles or Bullshit Machines? – Lesson 11," lecture UW, 2025.