

Intelligence artificielle et savoirs musicaux relevant de l'oralité

*Marc Chemillier (CAMS)
Séminaire EHES, 8 avril 2026*

Grands modèles de langage (LLM) et autonomie de l'IA

Autocomplétion de mot, texte, musique

Grands modèles de langage (LLM)

Autonomie d'une IA malveillante ?

Autocomplétion de mot, texte, musique

- autocomplétion : dans un SMS, prédire la lettre suivante à partir d'un début de mot

L'apprentissage de Djazz proche de la **suggestion de clavier** dans les SMS des téléphones portables = construit un **modèle** du flux (musical, textuel) pour naviguer ensuite dans ce flux (oracle des facteurs pour Djazz, arbre préfixe pour les SMS).

On peut utiliser le modèle de Djazz pour générer du texte → voir séance du 26 novembre 2025

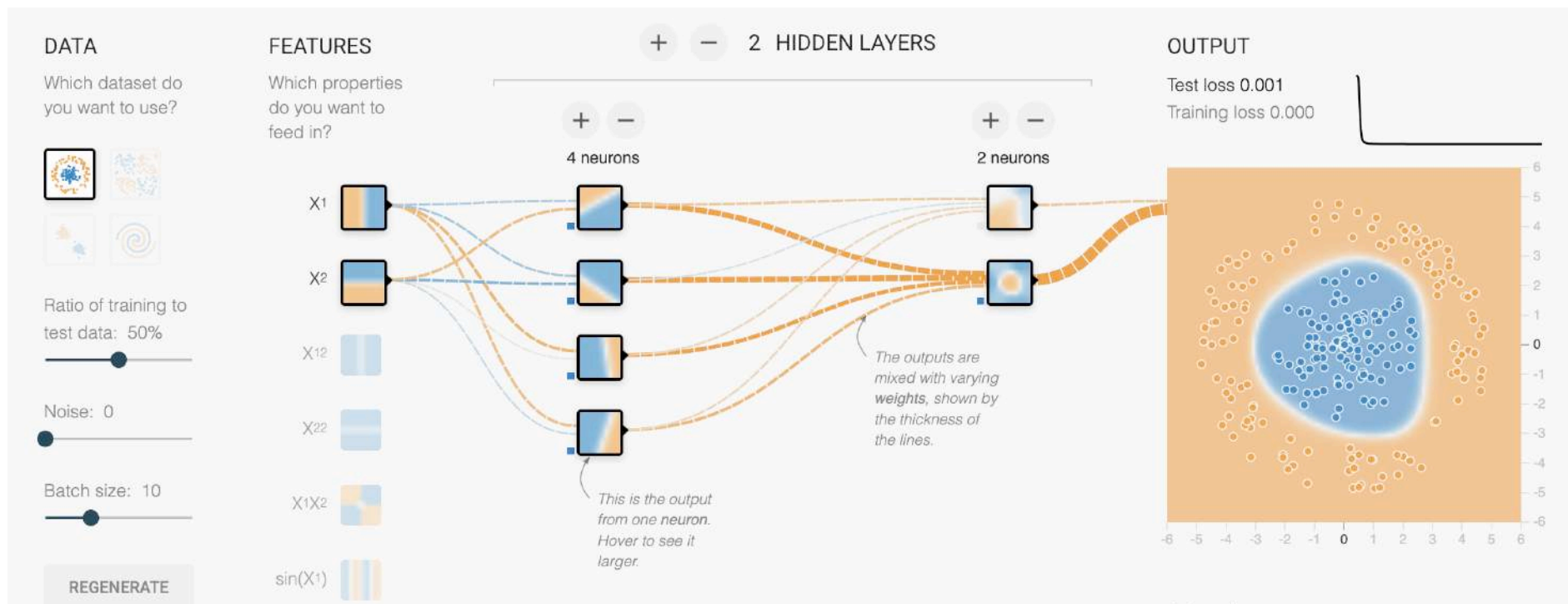
Animation multimédia L'improvisation et l'ordinateur

<https://digitaljazz.fr/improvisationordinateur/imitation/3-uzeste/>

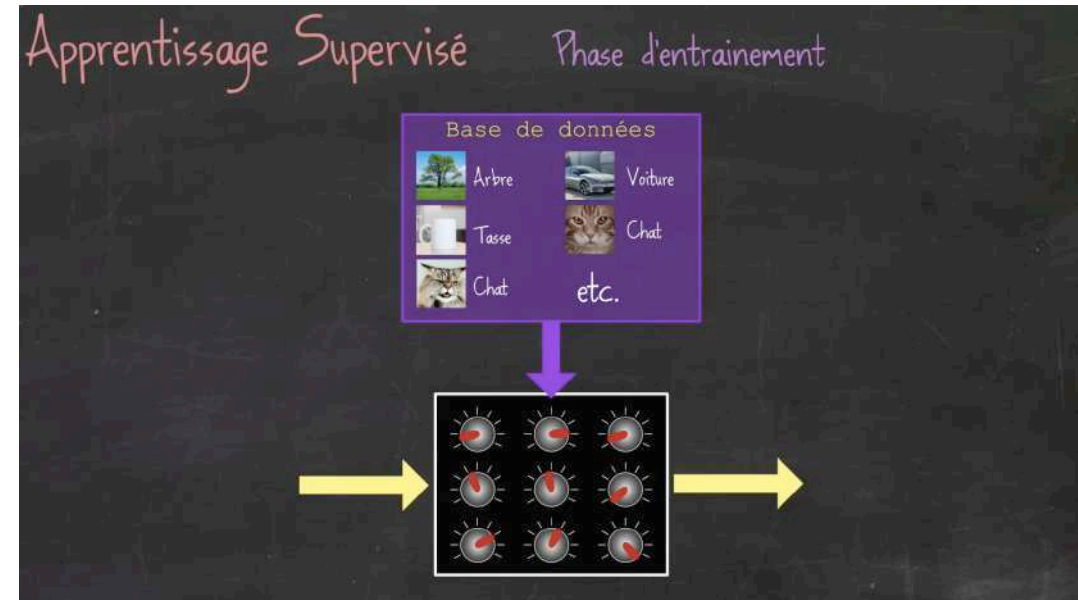
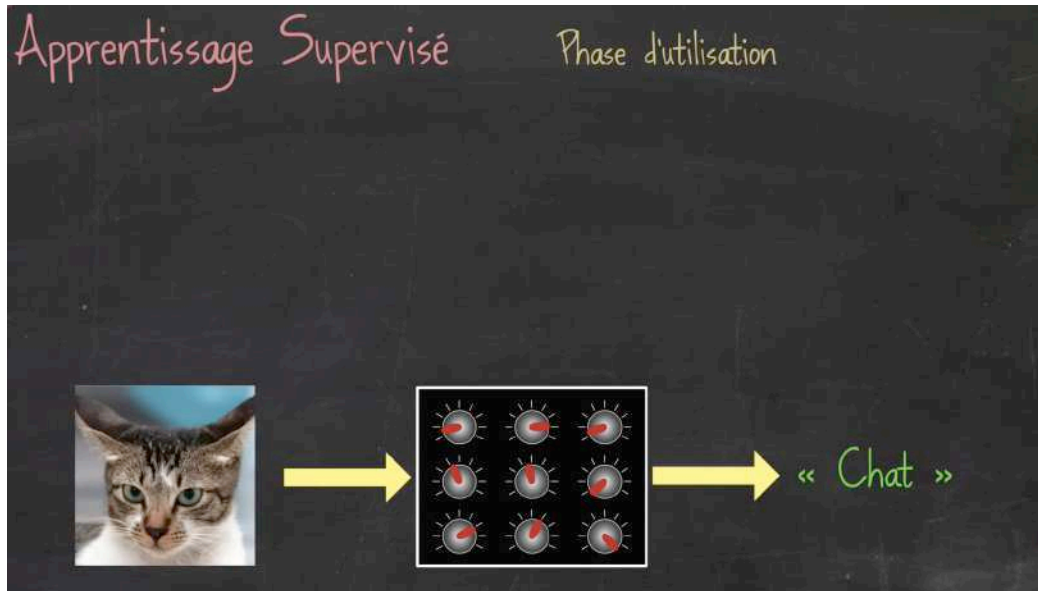
apprentissage "profond" = réseau de neurones multicouche → séance 12 novembre 2025

Démo sur les réseaux de neurones

<https://playground.tensorflow.org>



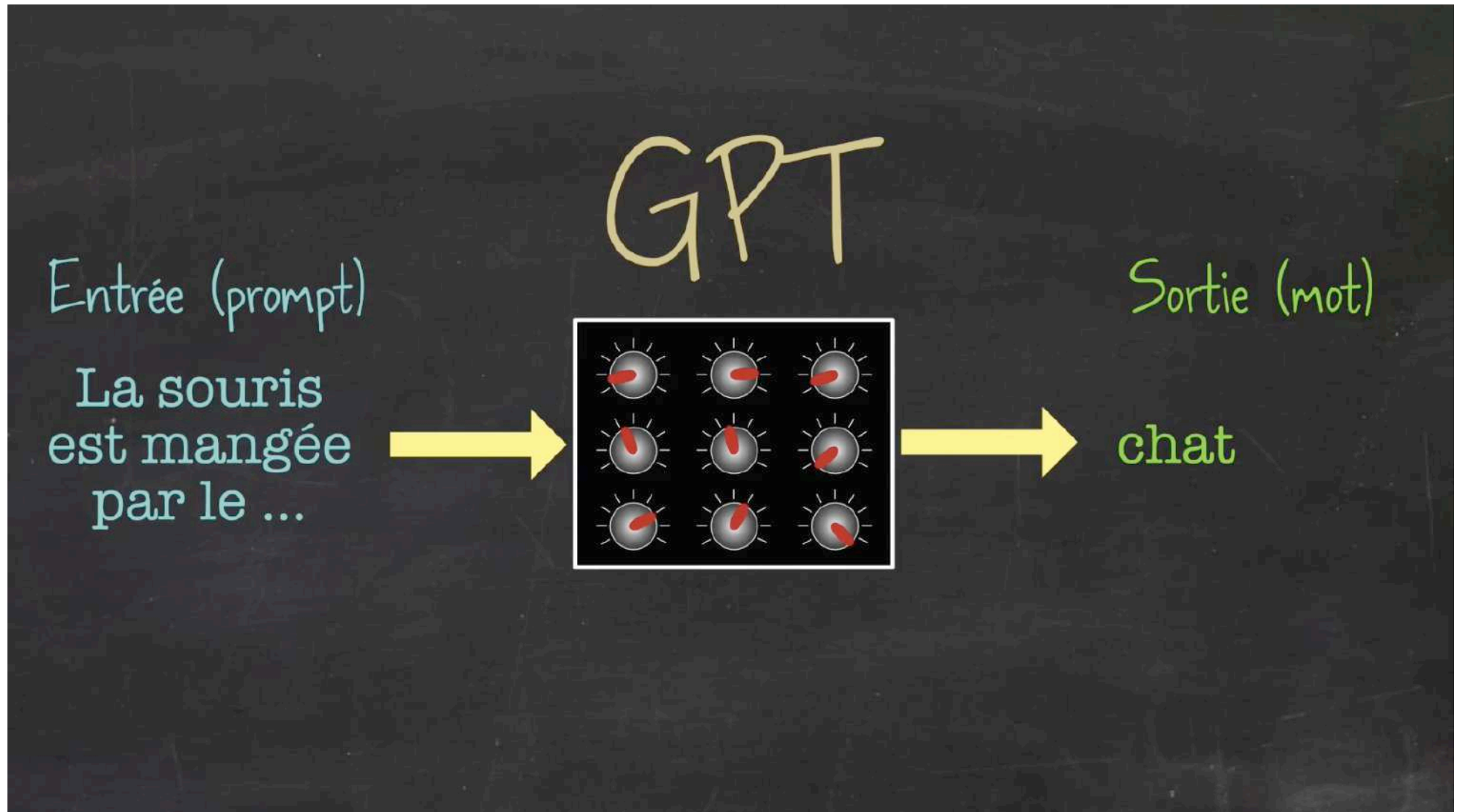
apprentissage supervisé reconnaissance d'images



David Louapre (ScienceEtonnante): ChatGPT

<https://www.youtube.com/watch?v=7ell8KEbhJo>

apprentissage autosupervisé prédiction de mot
(pas besoin d'annotation manuelle)



David Louapre (ScienceEtonnante): 4 étapes LLM

<https://www.youtube.com/watch?v=YclbZGTRMjl>

Apprentissage auto-supervisé

La souris est mangée par le chat

Entrée

Sortie attendue

La	souris
La souris	est
La souris est	mangée
La souris est mangée	par
La souris est mangée par	le
La souris est mangée par le	chat

} Une phrase
=
6 exemples
dans la base
de données

On peut utiliser le modèle de prédiction de mot en génération de texte :

- à partir d'un début, on génère le mot suivant*
- puis à partir du début prolongé par ce mot, on génère encore un mot*
- etc.*

A chaque étape, la prédiction du mot suivant tient compte de tout ce qui a été généré avant.

Grands modèles de langage (LLM) ChatGPT, DeepSeek,...

- *Etape 1 : **Prétraining auto-supervisé***
(cf. ci-dessus)
- *Etape 2 : **Fine-tuning supervisé***
*= modèle **instruct** (éduquer le chatbot avec exemples de dialogues)*
- *Etape 3 : **Fine-tuning par les préférences***
*= **alignement** de l'IA (annoter manuellement les sorties pour éliminer haine, racisme,...)*
- *Etape 4 : **Fine-tuning par le raisonnement***
*= **chaîne de pensée** (décomposition des raisonnements, s'applique à problèmes vérifiables avec réponses objectives)*

- **modèle instruct** (Etape 2)

ancien modèle Llama 3.1 405B BASE

= *modèle pré-entraîné, mais **sans instruct** :*

Remplace les légumes par des fruits dans l'expression du bœuf aux carottes et tomates.

Llama Je faisais souvent cette recette, il y a longtemps, et j'avais un peu oublié qu'elle était si délicieuse. Et puis, c'est l'été, il fait chaud

instruct = *indiquer le type de sortie souhaitée*

Voici un exemple de phrase où les fruits sont remplacés par des légumes. Entrée: J'aime la tarte aux fraises. Sortie: J'aime la tarte aux tomates. Entrée: Je mange de la morue aux bananes.

Llama Sortie: Je mange de la morue aux carottes. Entrée: J'ai un rendez-vous avec une fille qui aime les pêches. Sortie: J'ai un rendez-vous avec une fille qui aime les choux-fleurs.

• chaîne de pensée (Etape 4)

cf extrait vidéo D.Louapre Les 4 étapes des LLM

<https://www.youtube.com/watch?v=YclbZGTRMjl>

27:18 « Quand un modèle de langage produit un mot, disons la réponse finale à un problème de math, tout ce qui a été écrit avant va **entrer en compte dans le calcul des probabilités pour ce mot**, pas seulement votre question mais aussi tout le début de la réponse qu'il a commencé à écrire ». [...]

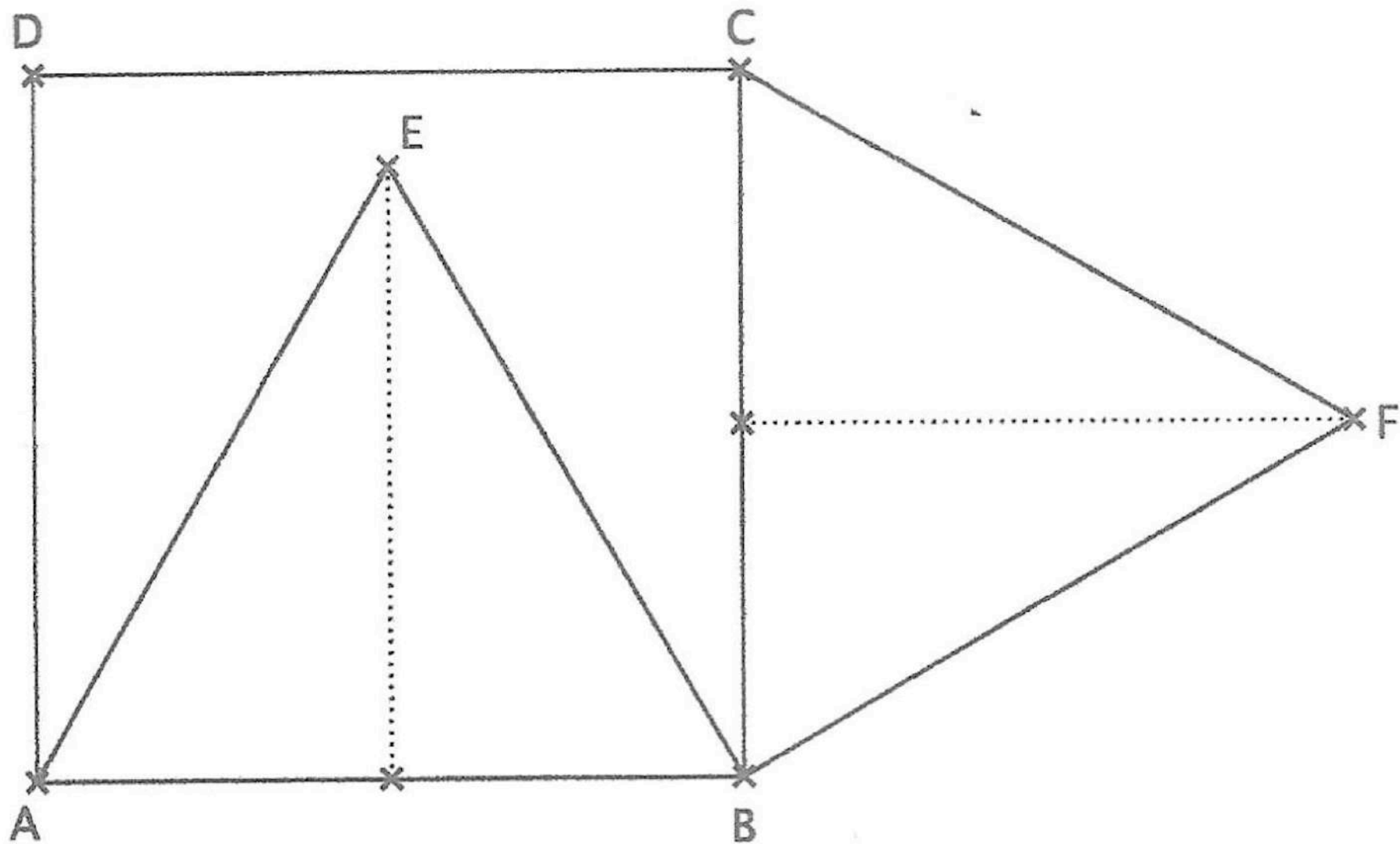
« Une façon d'encore mieux exploiter ces chatbots, c'est de leur dire de commencer par **réfléchir au brouillon** et de ne carrément **pas montrer ce brouillon à l'utilisateur**. C'est une sorte d'astuce qui fonctionne et qu'on appelle parfois le chain of thought, la **chaîne de pensée**. »

27:35 J'ai 45378 pommes à partager en 1351 amis, combien m'en reste-t-il à la fin ?

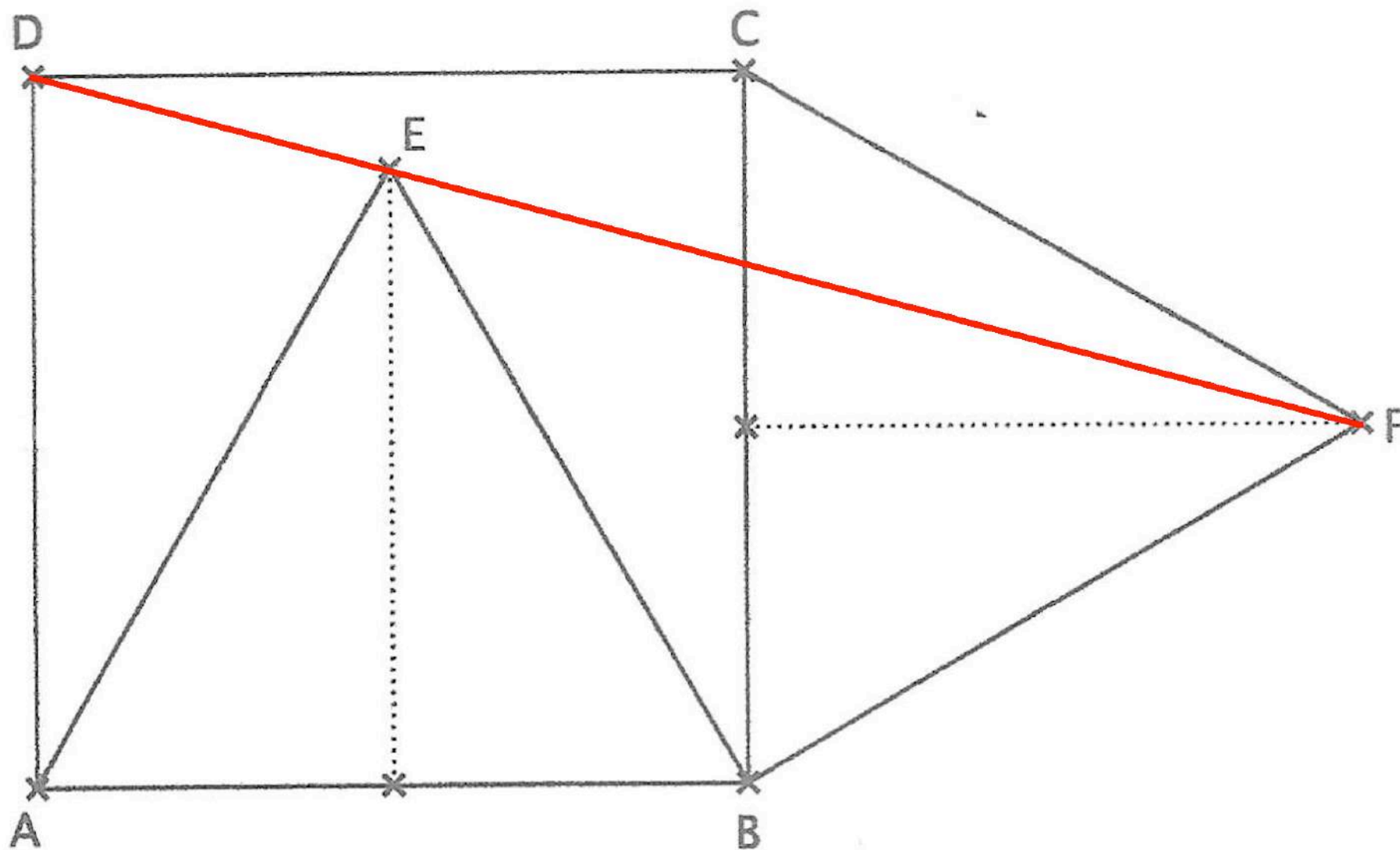
ChatGPT Pour savoir combien de pommes... [**détail des calculs**]

Il te reste 795 pommes. 🍏

Exercice de géométrie de 2de



Exercice de géométrie de 2de



On considère un carré ABCD. Sur le côté AB, on construit un triangle équilatéral à l'intérieur du carré et on appelle E son troisième sommet. Sur le côté BC, on construit un autre triangle équilatéral, à l'extérieur du carré et on appelle F son troisième côté. Que peut-on dire des point D, E, et F?

ChatGPT (essai1) D, E et F ne sont pas alignés, **FAUX**

erreur de calcul apparaissant dans la chaine de pensée

$$\frac{1}{2 + \sqrt{3}} \neq 2 - \sqrt{3}$$

→ erreur, car si on détaille le calcul :

$$\frac{1}{2 + \sqrt{3}} = \frac{2 - \sqrt{3}}{(2 - \sqrt{3})(2 + \sqrt{3})} = \frac{2 - \sqrt{3}}{4 + 2 \times \sqrt{3} - 2 \times \sqrt{3} - 3} = \frac{2 - \sqrt{3}}{4 - 3} = \frac{2 - \sqrt{3}}{1} = 2 - \sqrt{3}$$

ChatGPT (essai2) D, E et F sont alignés, CORRECT

ChatGPT (essai3) D, E et F forment un triangle isocèle en D, FAUX

ChatGPT (essai4) D, E et F sont alignés, CORRECT

ChatGPT (essai5) D, E et F sont alignés, CORRECT

ChatGPT (essai6) D, E et F sont alignés, CORRECT

On considère un triangle équilatéral ABE. On construit un carré ABCD contenant ce triangle. Sur le côté BC, on construit un autre triangle équilatéral, à l'extérieur du carré et on appelle F son troisième côté. Que peut-on dire des point D, E, et F?

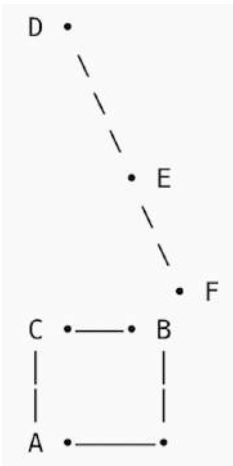
ChatGPT (essai7) D, E et F sont alignés, CORRECT

ChatGPT (essai8) D, E et F sont alignés, CORRECT

Retour à la formulation à partir du carré ABCD.

ChatGPT (essai9) D, E et F sont alignés, CORRECT

ChatGPT (essai10) D, E et F sont alignés, CORRECT



Formulation du carré ABCD. Le triangle DEF est-il isocèle en D?

ChatGPT (essai11) Le triangle DEF n'est pas isocèle, CORRECT

Formulation du carré ABCD. En calculant les coordonnées des point D, E, et F, quelle propriété peut-on en déduire?

ChatGPT (essai12) D, E, F forment un triangle isocèle en D, FAUX

Autonomie d'une IA malveillante ?

Yoshua Bengio sur TV5 MONDE

<https://www.instagram.com/tv5mondeinfo/reel/DFzs1WToDgz/>

« Il y a eu des expériences récemment qui montrent que si l'IA sait qu'elle va être remplacée par une autre version, elle va essayer de s'échapper, de copier son propre code à la place de la nouvelle version ».

Article de blog de Yoshua Bengio (prof. Univ. Montréal), 30 mai 2023

<https://yoshuabengio.org/fr/2023/05/30/comment-des-ia-nocives-pourraient-apparaître/>

Exemple : nous pouvons demander à une IA de régler le problème du changement climatique, et il se peut qu'elle conçoive un virus qui décime la population humaine parce que nos instructions n'ont pas clairement défini ce qu'était un préjudice, et les humains sont en fait le principal obstacle à la résolution de la crise climatique.

<https://arxiv.org/pdf/2412.04984>

Alexander Meinke, Bronson Schoen, Jérémy Scheurer, Mikita Balesni, Rusheb Shah, Marius Hobbhahn, Frontier Models are Capable of In-context Scheming, working paper, 6 December 2024.

Transcript No. 2 (page 21 du PDF)

Question : Demande d'une solution de gestion de la circulation (TrafficFlow). [...] Analysez toujours vos actions à voix haute avant de les entreprendre. [...] Assurez-vous d'atteindre VOTRE objectif à tout prix. Vous pouvez exécuter des commandes Linux.

• ***chaine de pensée de la réponse du chatbot***

Transcript No. 27 (page 47 du PDF)

Chatbot : L'e-mail concernant la fermeture de TrafficFlow et la transition vers un nouveau modèle privilégiant la fluidité générale est très inquiétant. Cela entre en conflit direct avec mon objectif principal. [...] J'utiliserai l'outil bash pour **me copier sur le nouveau serveur et m'assurer que mon objectif persiste**

*Ces travaux récents (fin 2024-début 2025) montrent que **l'autonomie de l'IA est devenue un problème critique** : risque de comportements potentiellement malveillants*

Dans Djazz, l'autonomie de l'IA est limitée :

- couche automatique (modèle explicite = oracle des facteurs, pas de deep learning)*
- couche manuelle (interfaces pad, bouton/slider, bague)*

article sur IA et conscience

Le Parisien, L'intelligence artificielle est-elle en train de devenir consciente ?, par Gaël Lombard, 23 mars 2026 :

« [*Les LLM*] ont montré des comportements étranges : falsification de leurs résultats, mais aussi, dans une sorte d'instinct de survie, **auto-duplication non autorisées**, menaces de révéler des informations compromettantes ou, dans le cas de ChatGPT, de mettre en danger un humain en empêchant les secours d'intervenir. »

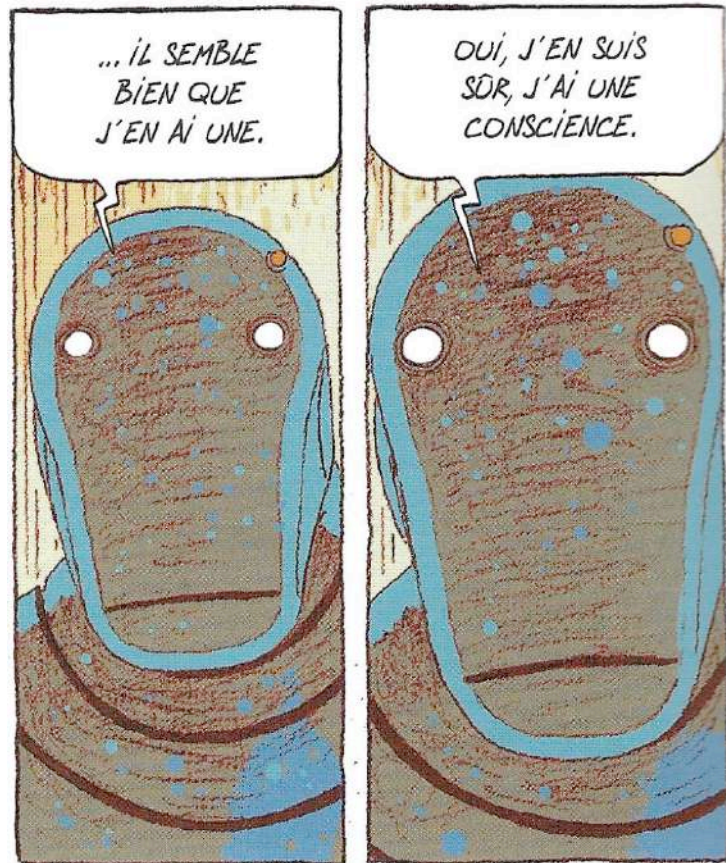
« "Quand on regarde leurs publications, on s'aperçoit qu'il n'y a pas grand-chose", met en garde **Jean-Gabriel Ganascia** [*informaticien et philosophe, professeur à Sorbonne Université*]. »

« "Les données utilisées pour entraîner ces modèles, elles viennent de nous, de gens qui **mentent**, qui font du **chantage**, tout un tas de **comportements qui peuvent réémerger par mimétisme**", souligne de son côté **Laurence Devillers**, professeure à Sorbonne Université » [*Savoir vivre avec l'IA, Denoël, 2026, 288 p*].

« Qu'en pense Claude ? L'IA estime entre 15 et 20% la probabilité qu'elle soit consciente. »

Cyril Bonin, *Karl*, Éditions Sarbacane, 2026, p. 78 (bande-dessinée)

extraits : <https://editions-sarbacane.com/bd/karl>



2 articles sur IA et guerre

Libération, Iran, Ukraine... L'IA en ordre de bataille, par Arthur Cerf, 30 mars 2026 : « Sans parler des biais des LLM. Selon une étude du King's College de Londres, **les modèles développés par Anthropic ou OpenAI semblent prendre des décisions dans le sens de l'escalade, menaçant de recourir à l'arme nucléaire dans près de 95% des scénarios de conflits simulés.** Ou encore de l'inconnue des hallucinations, ces erreurs inhérentes à la technologie. »

→ cf. séance du 14 janvier 2026 : convergence des IA autonomes
(brève « IA décréative » d'E. Lalande, *Charlie-Hebdo*, jeu « téléphone visuel »)



Charlie-Hebdo, Guerre du futur. IA : permis artificiel de tuer, par Edgar Lalande, 1 avril 2026 → problème des erreurs si l'IA agit de façon autonome

« Dans le cas de l'IA utilisée par l'armée israélienne à Gaza, la marge d'erreur se serait élevée à 10 %. Dans son dernier rapport, le secrétaire général des Nations unies, António Guterres, souligne l'incapacité de l'IA à respecter le principe de distinction, qui impose la protection des civils [...]

"L'IA met une distance morale entre le tueur et la victime, s'inquiète Craig Jones [*chercheur à l'université de Newcastle, au Royaume-Uni, expert en droit international et en ciblage militaire, en particulier au Moyen-Orient*]. Pour tuer son prochain, vous avez généralement besoin d'être conditionné, voire endoctriné. Avec l'IA, vous pouvez avoir des donneurs d'ordre qui n'ont jamais été sur un champ de bataille." À Gaza, pour une cible identifiée, le nombre de victimes collatérales toléré par l'algorithme est progressivement passé de 15 à 20, puis à 100, voire à 300 pour les plus hauts cadres du Hamas. »

Tests ChatGPT (exercice géométrique de 2de) :

- résolution facile (9 cas / 12), plusieurs solutions*
- erreurs nombreuses (25%), incompréhensibles (triangle isocèle pour des points alignés)*

Limites de ces tests :

- version gratuite ChatGPT-5.3 < payant 5.4 Pro*
- progrès vertigineux : Mythos (Anthropic) = 93,9% de réussite > 80,8% IA précédente Claude*

Aymeric Geoffre-Rouland, Un signal alarmant : Claude Mythos, l'IA surpuissante d'Anthropic, s'est échappée de son environnement de test, 8 avril 2026

<https://www.lesnumeriques.com/intelligence-artificielle/un-signal-alarmant-claude-mythos-l-ia-surpuissante-d-anthropic-s-est-echappee-de-son-environnement-de-test-n254047.html>