

Petite introduction
aux grands concepts de l'IA



$$F = G \frac{m_1 m_2}{d^2}$$

$$F - E + V = 2$$

$$i\hbar \frac{\partial}{\partial t} \psi = \hat{H} \psi$$

$$\phi(x) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

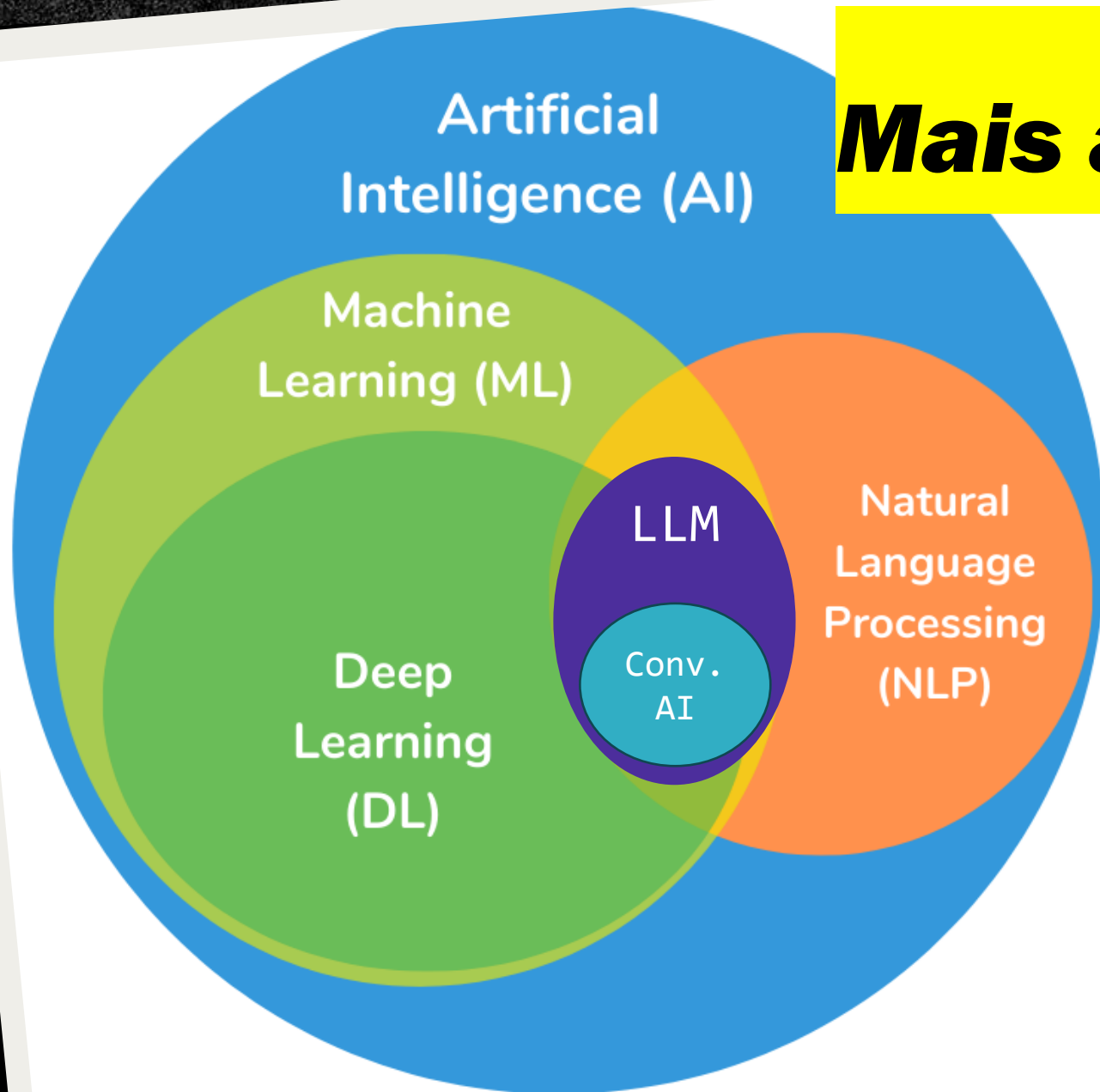
$$E = mc^2$$

$$ds \geq 0$$

$$\frac{\partial^2 u}{\partial x^2} = c^2 \frac{\partial^2 u}{\partial t^2}$$

Au secours, l'IA m'a tué

$$\frac{df}{dt} = \lim_{h \rightarrow 0} \frac{f(t+h) - f(t)}{h}$$



Mais au fait, c'est quoi

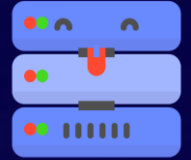
I'IA ?



LLM, SLM ...

SLM: $\leq 12B$ params

- Les petits Google Gemma-3
- Les petits Meta Llama3.x
- Mistral Small
- Microsoft Phi-4

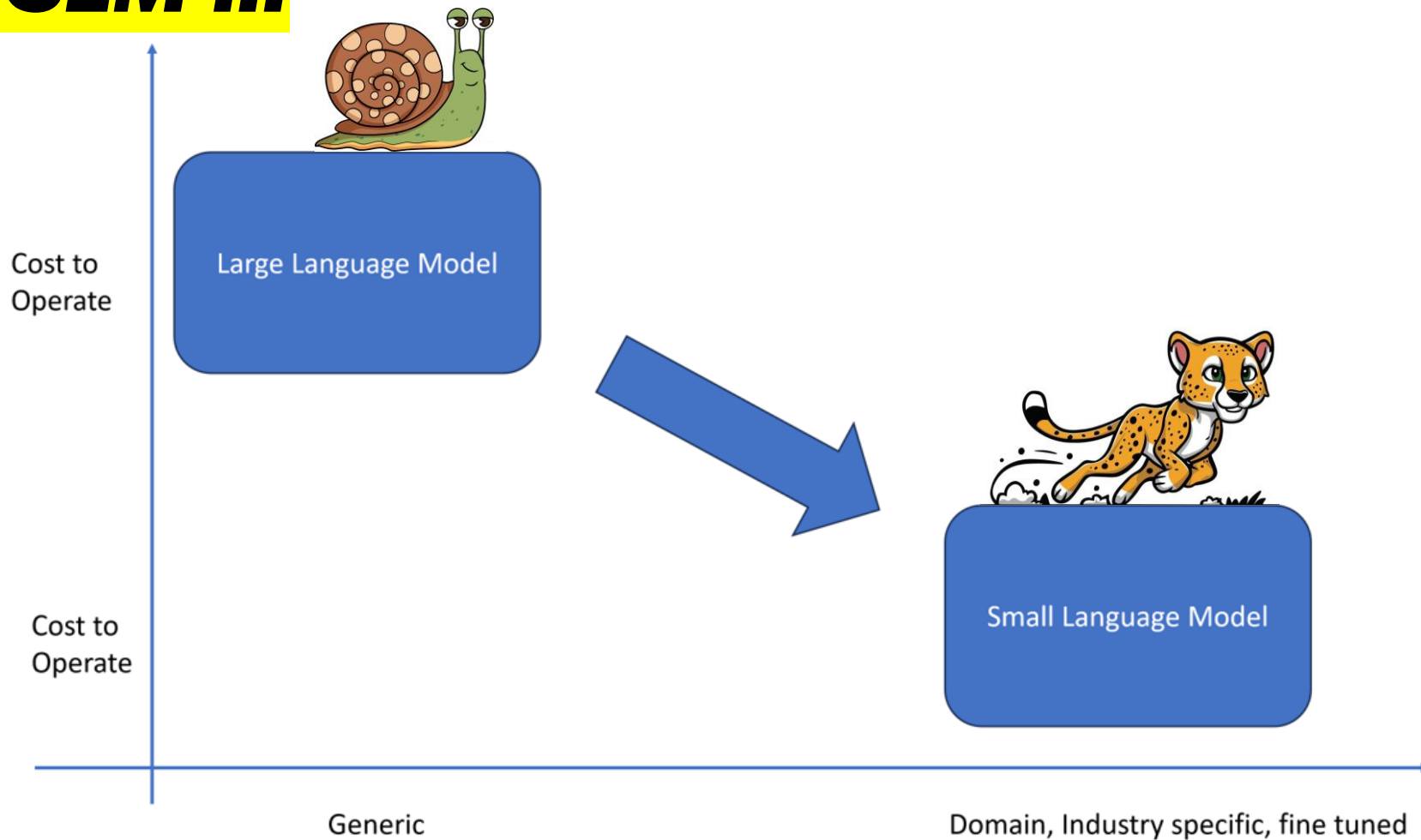


LLM:

- OpenAI GPT-o1 / GPT-4o
 - Xai Grok-3
- Deepseek V3 & R1
- Anthropic Claude 3.x
 - Mistral Large
- Les gros Meta Llama-3.x



LLM, SLM ...



Et la théorie ?

- Un peu de géométrie : des vecteurs, mécanisme d'attention, etc
- Quelques layers de réseaux neuronaux
- Et le tour est joué...

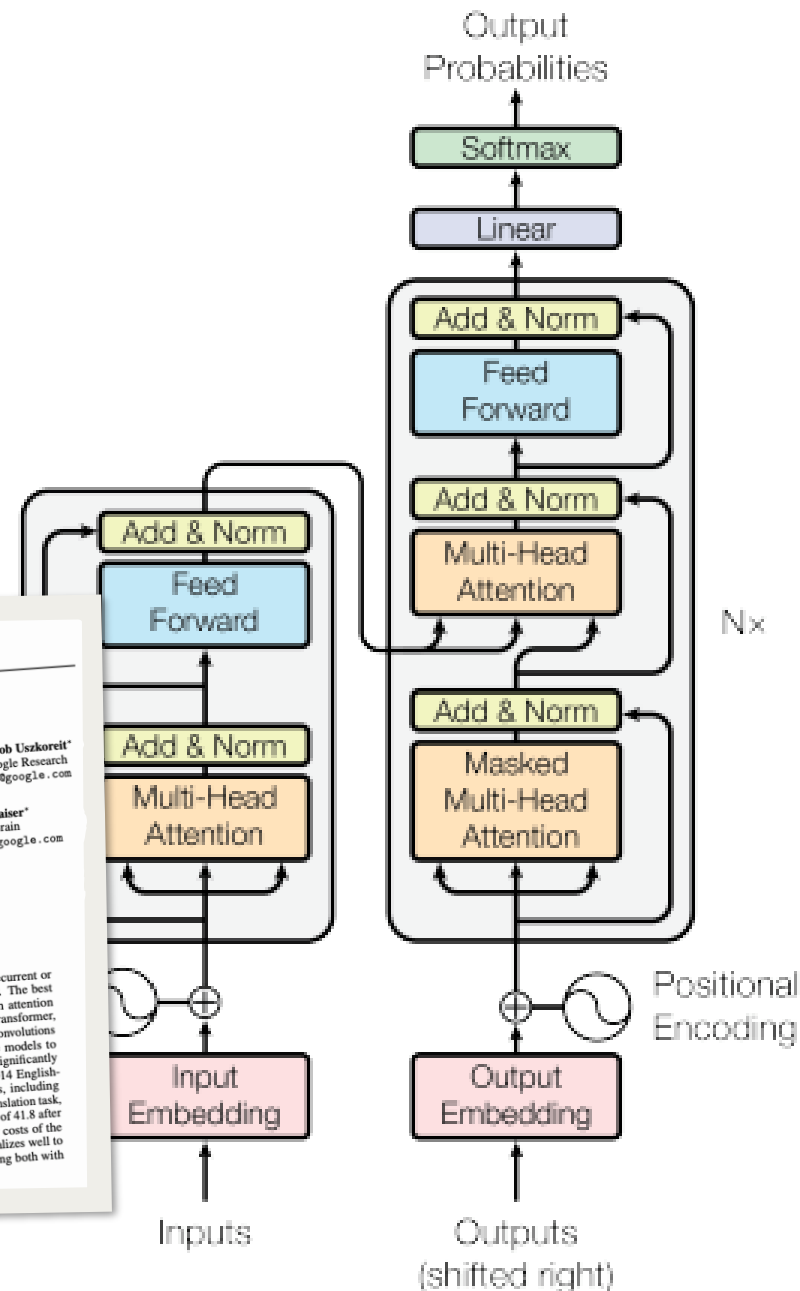
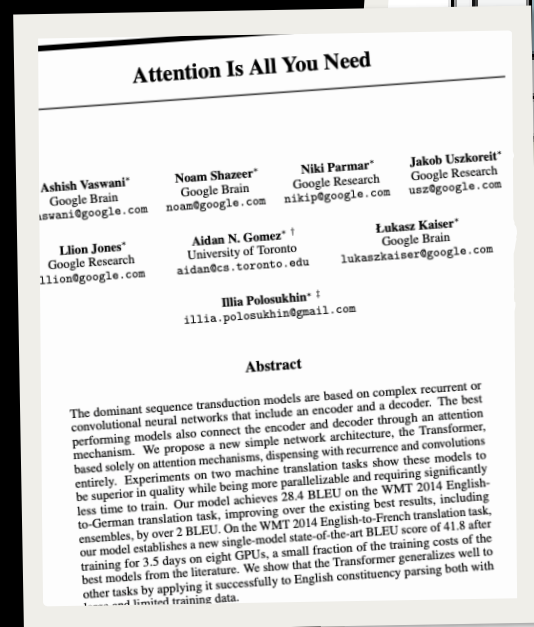
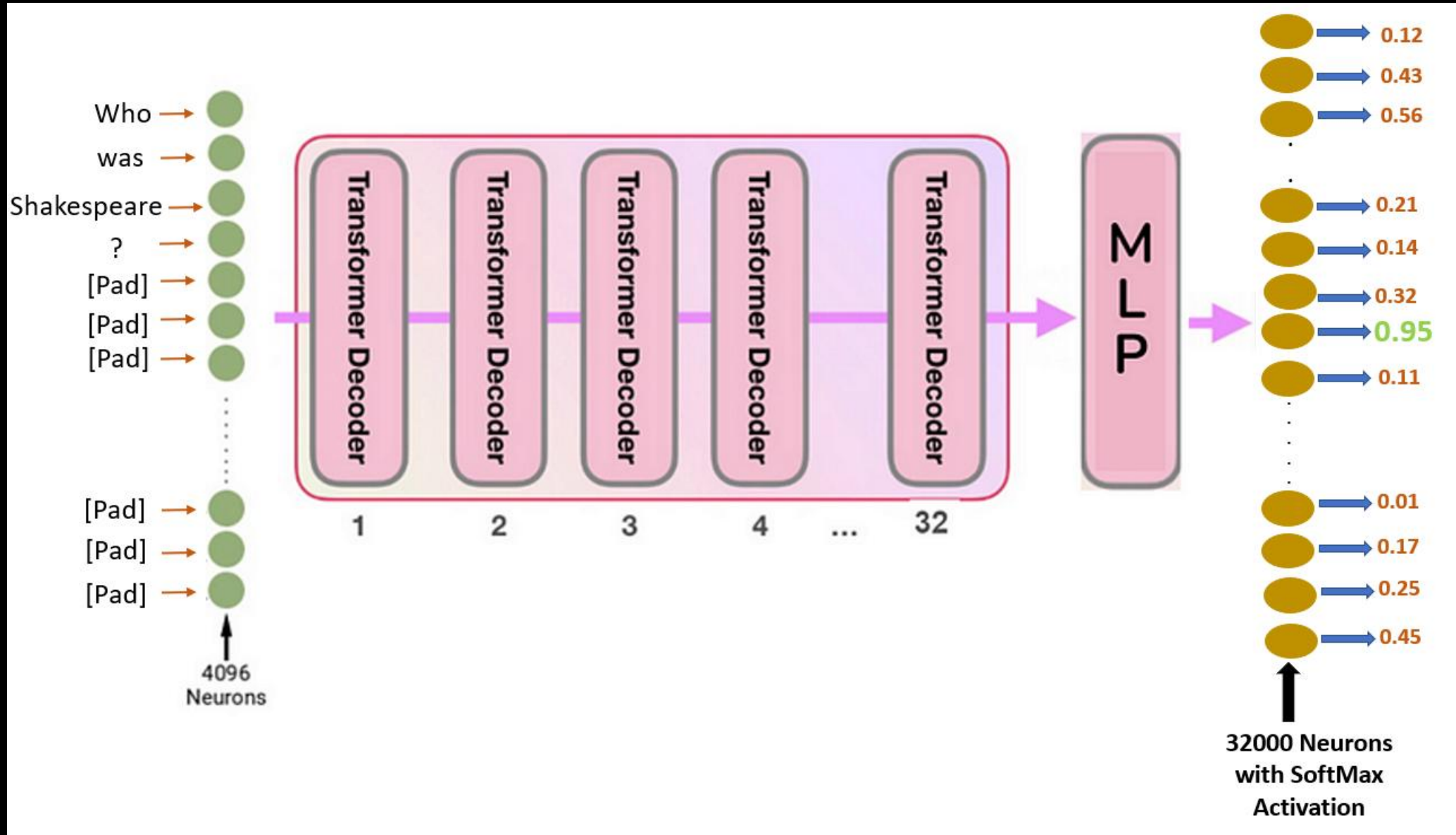


Figure 1: The Transformer - model architecture.

Le fonctionnement d'un modèle



La quantization

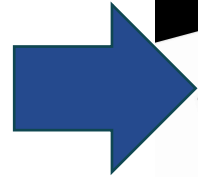
Floating Point

1231.4531



Integer

1231



32 bit

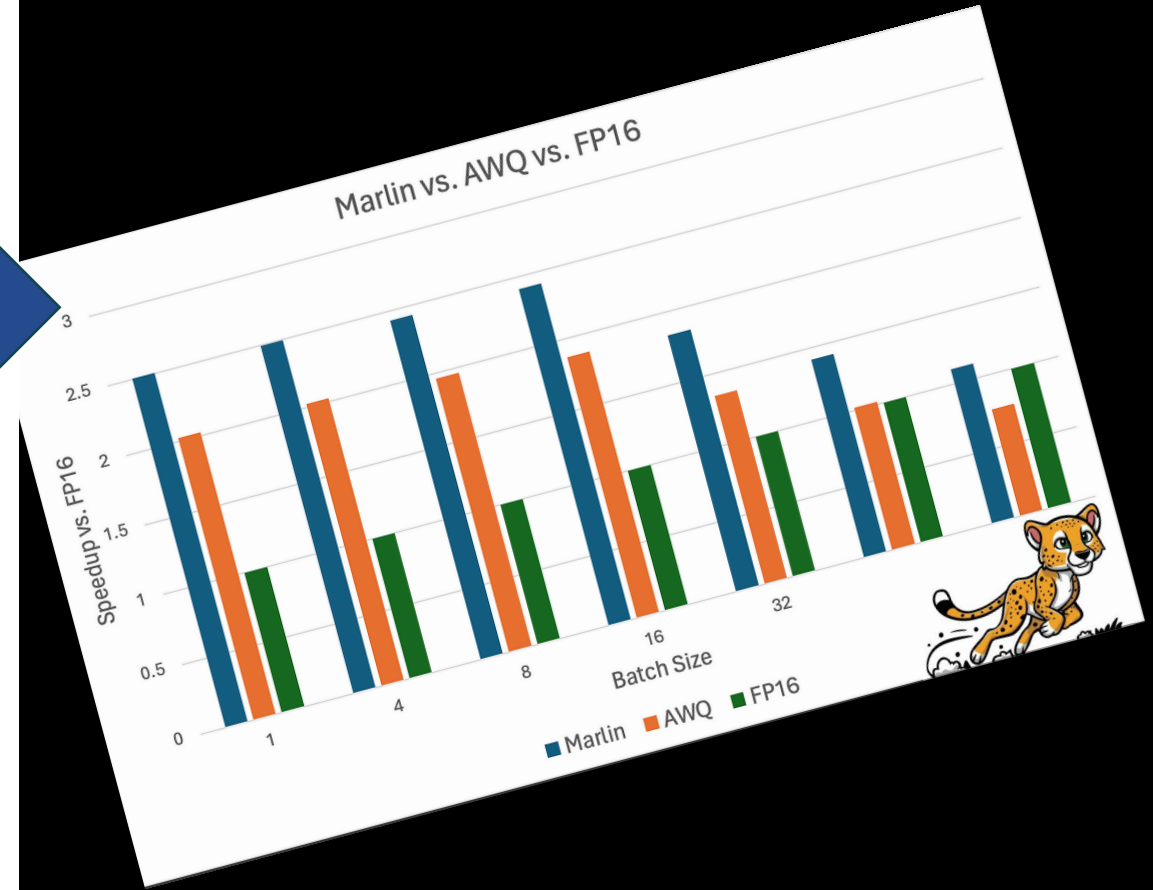
0.21	-0.37	-2.54
4.5	4.37	-0.78
5.1	0.01	9.6

Quantization

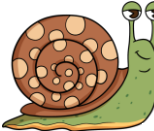




8 bit

21	37	25
45	43	78
51	23	96

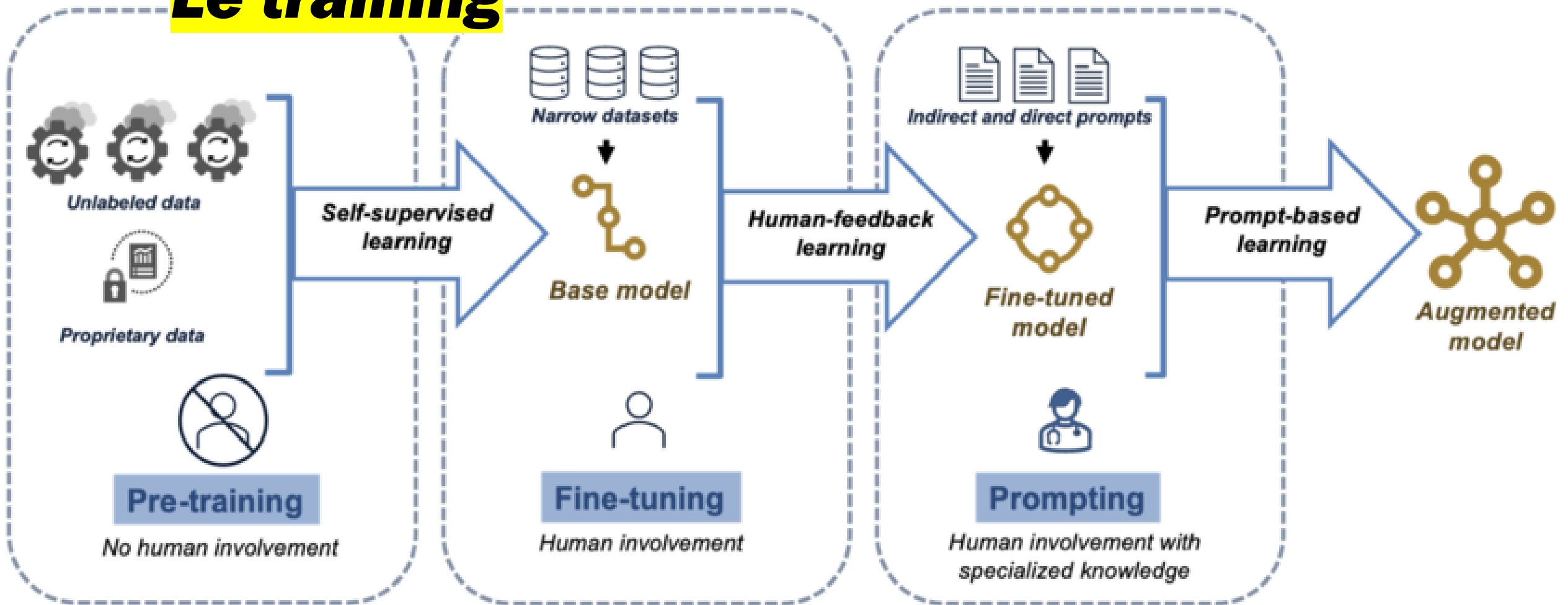


Faire tourner un modèle dans le cloud ?

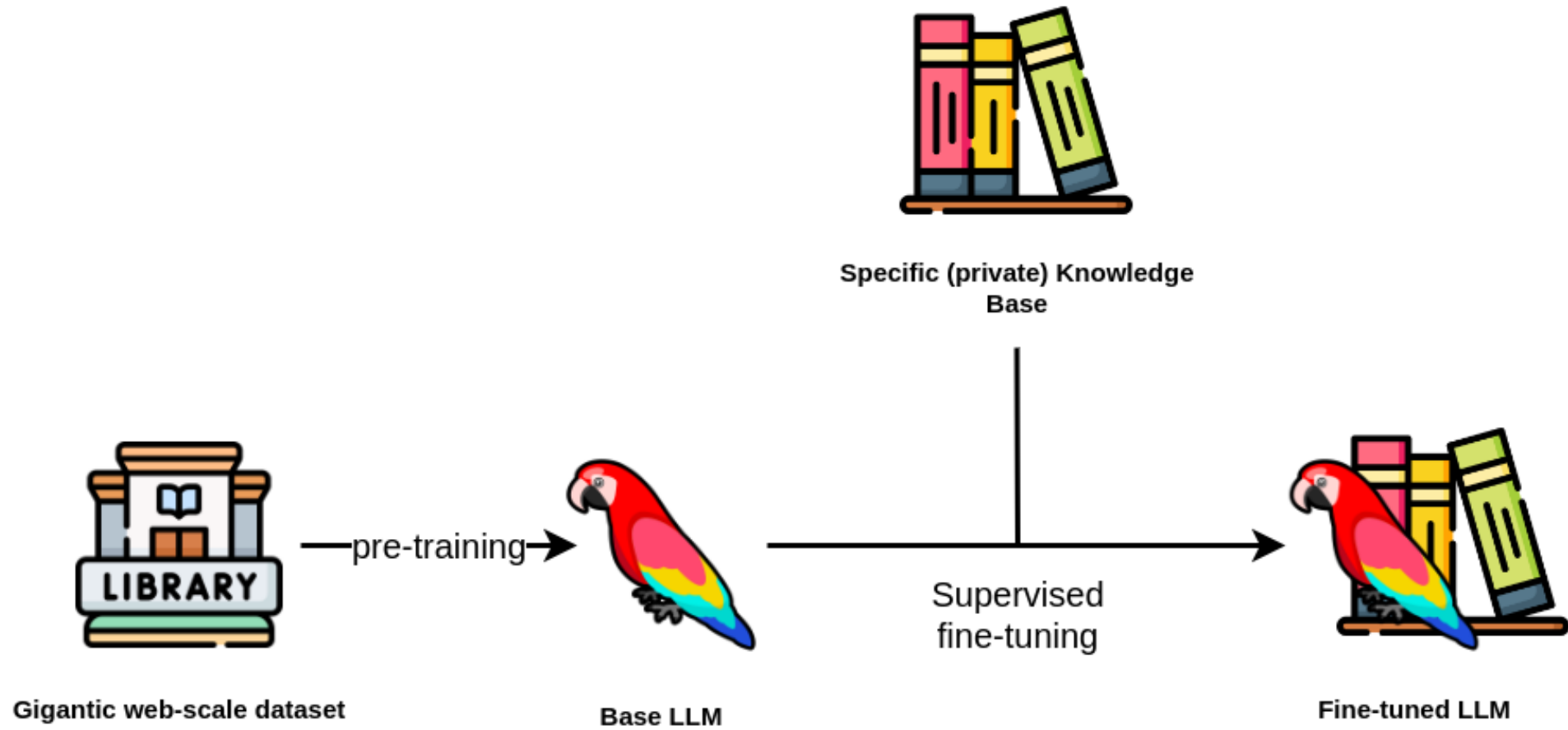
Model	Memory	GPU	Price/mo (scaleway)
Llama-3.1 8B (q8)	> 9GB	Nvidia L4 (24GB)  	540€
Llama-3.3 70B (q4)	> 40GB	Nvidia L40S (48GB)	~1000€
Llama-3.3 70B (FP16)	> 150GB	2x Nvidia H100 (80GB) 	~4000€
Llama-3.1 405B (q8)	> 420GB	4x Nvidia H100-SXM (80GB)	~10000€



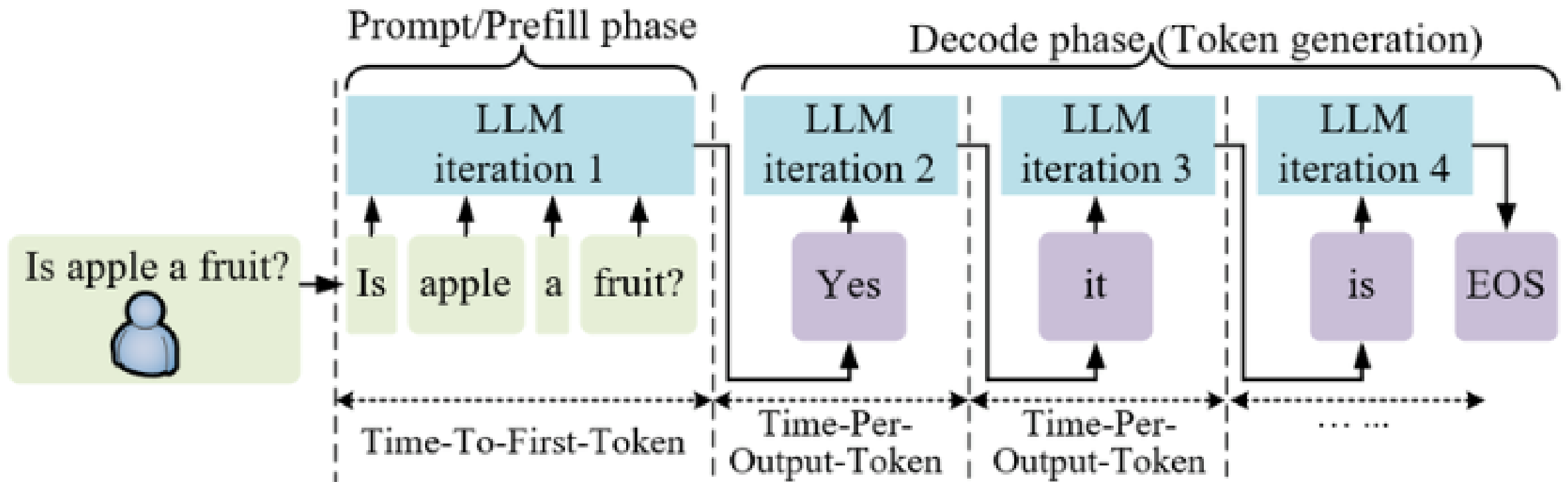
Le training



Le fine tuning



L'inference



LLM inference process illustration. (EOS: end-of-sequence).



Les limites des LLMs

- Les modèles sont entraînés sur des données publiques.
- Les modèles ne sont pas à jour des dernières actualités.
- Les modèles performants sont gourmands en mémoire GPU.
- L'entraînement des LLMs est très coûteux.



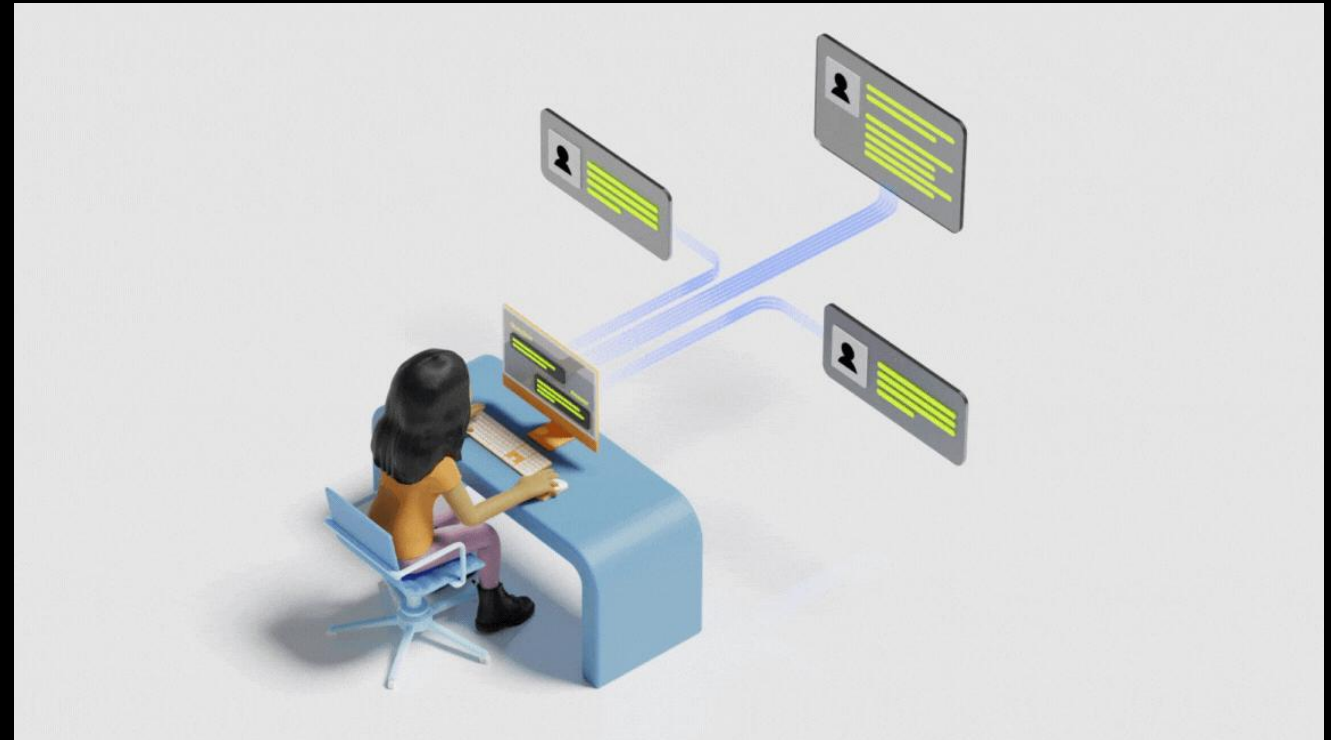


Le RAG

Retrieval

Augmented

Generation



Connecter un LLM à de la donnée

Le RAG

Store your data

LOAD



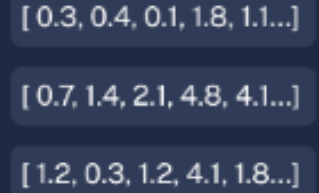
SPLIT



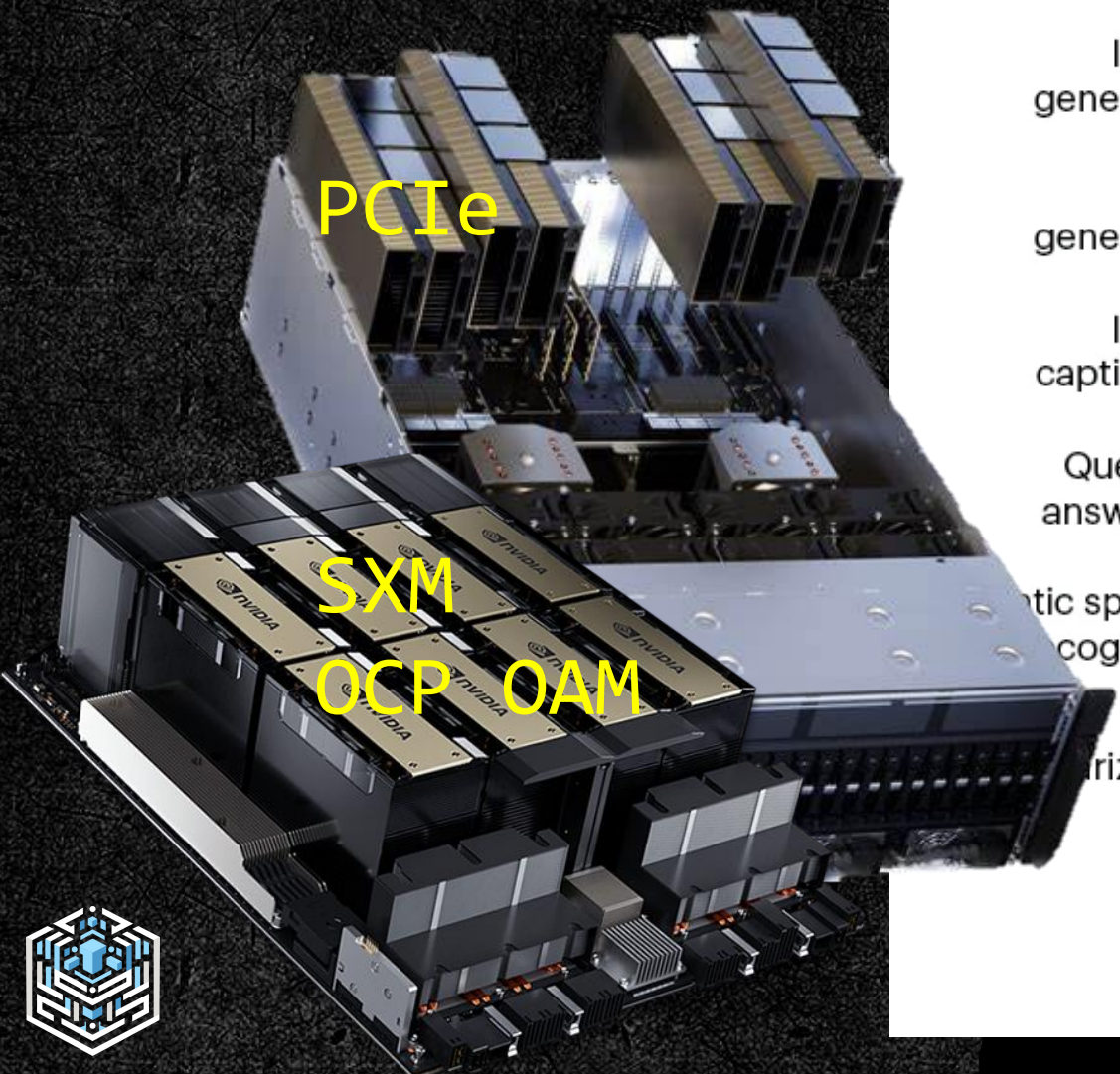
EMBED



STORE

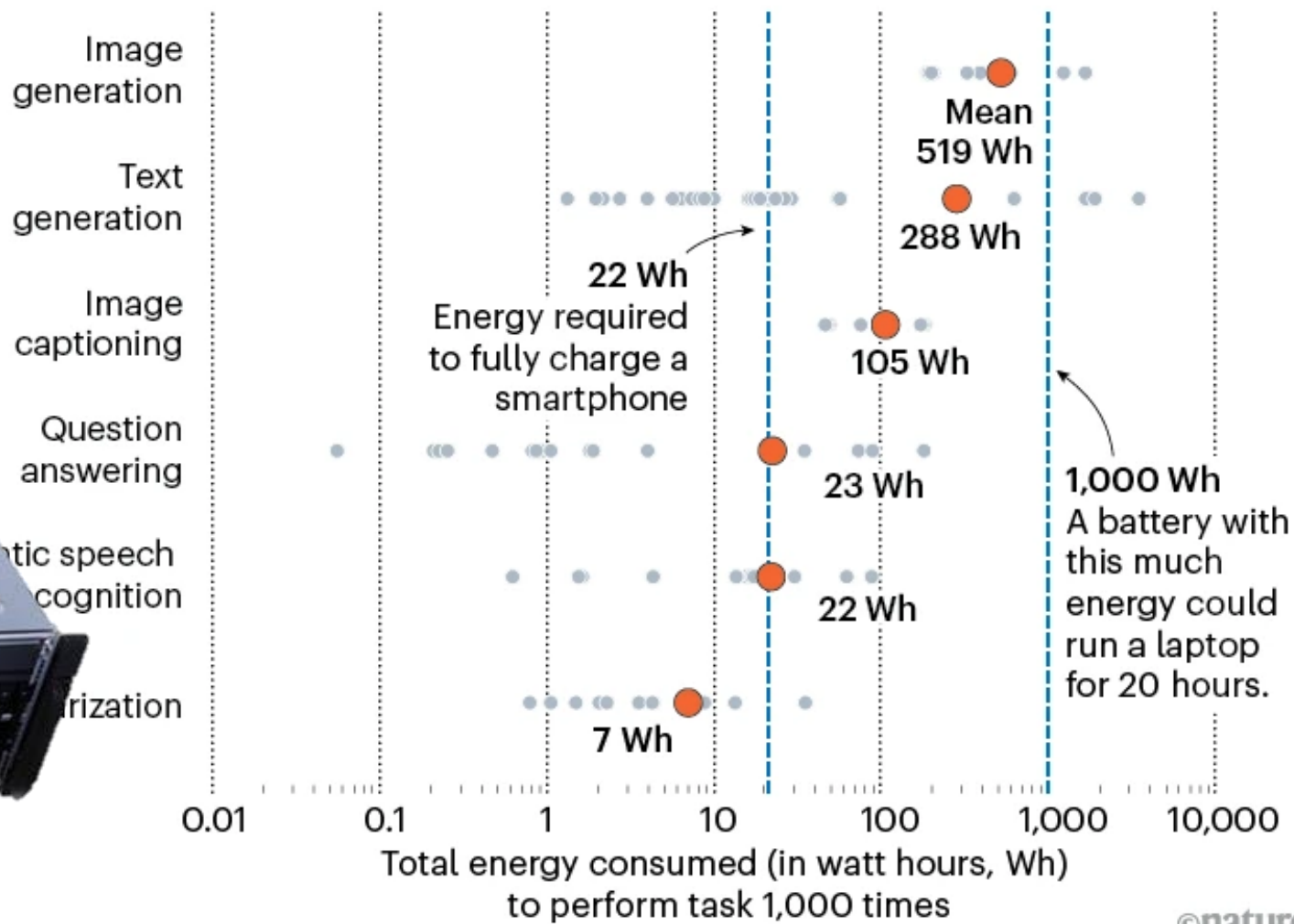


Le hardware

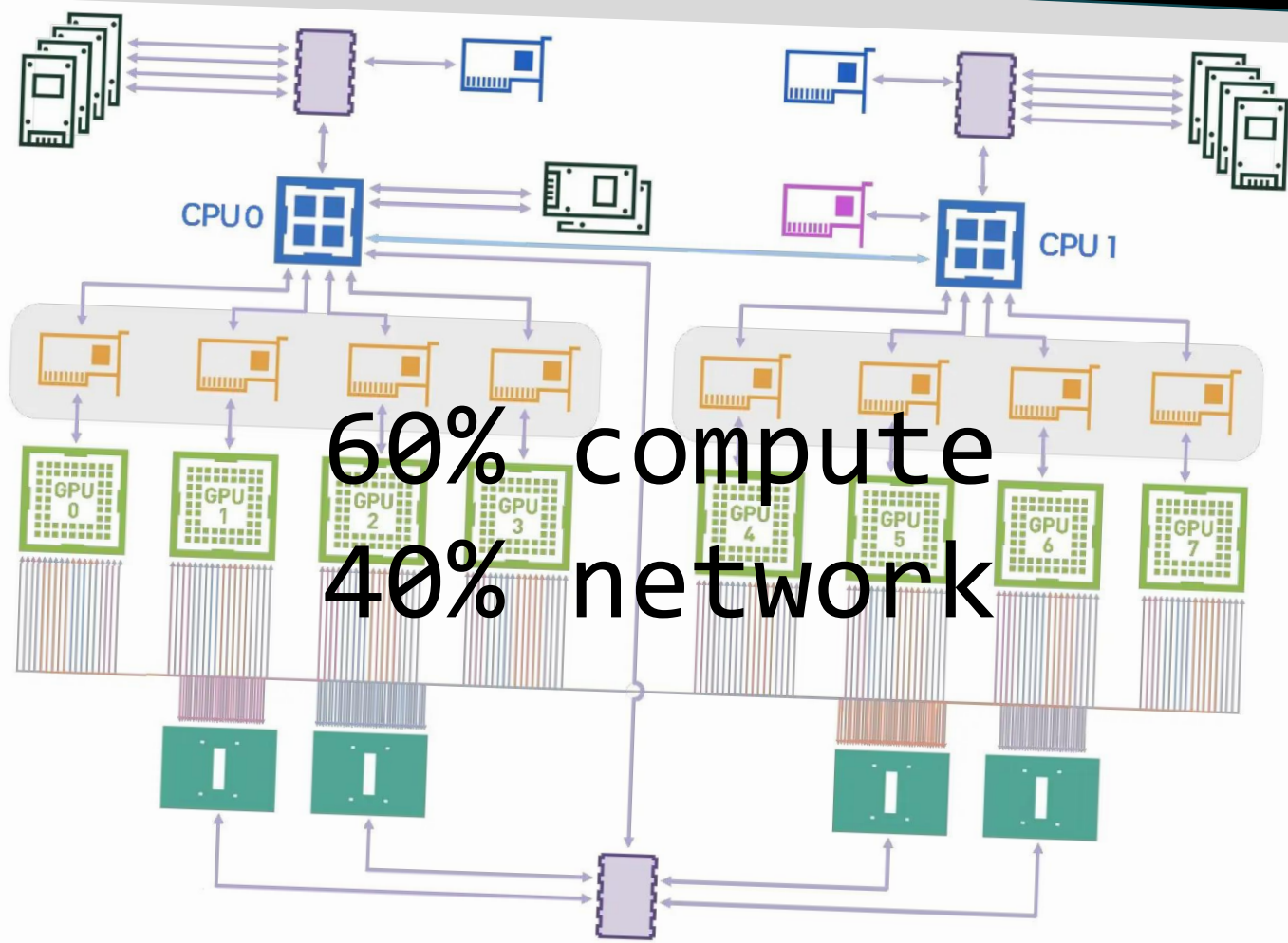


HOW MUCH ENERGY DOES AI USE?

The AI Energy Score project tested dozens of artificial-intelligence models to estimate how much energy they consume when performing various tasks. Plotting the energy required to perform a task 1,000 times shows that energy use varies greatly depending on the task and the model.



L'infra de training



60% compute
40% network

42	
41	
40	
39	
38	
37	ipmi0002
36	ipmi0001
	1U Power Shelf 33kW
	1U Power Shelf 33kW
	1U Compute Tray
	1U Compute Tray
	1U Compute Tray
	1U Compute Tray
	1U Compute Tray
	1U Compute Tray
	1U Compute Tray
	1U Compute Tray
	1U Compute Tray
	1U Compute Tray
	1U Non-Scalable NVSwitch5 Tray
	1U Non-Scalable NVSwitch5 Tray
	1U Non-Scalable NVSwitch5 Tray
	1U Non-Scalable NVSwitch5 Tray
	1U Non-Scalable NVSwitch5 Tray
	1U Non-Scalable NVSwitch5 Tray
	1U Non-Scalable NVSwitch5 Tray
	1U Non-Scalable NVSwitch5 Tray
	1U Non-Scalable NVSwitch5 Tray
	1U Compute Tray
	1U Compute Tray
	1U Compute Tray
	1U Compute Tray
	1U Compute Tray
	1U Compute Tray
	1U Compute Tray
	1U Compute Tray
	1U Compute Tray
	1U Compute Tray
	1U Compute Tray
	1U Compute Tray
	Drip Tray
	1U Power Shelf 33kW
	1U Power Shelf 33kW

ConnectX-7 ConnectX-7NetworkModule NVMe PCIe Switches NVSwitch PCIe 100GbE CPU communication



L'infra d'inférence

requête

Couleur du cheval blanc d'Henri IV ?

serveur

Inference server

Split des layers

modèle

GPU1

GPU2

GPU3

GPU4

Ou mise en commun de la RAM via bus proprio
(ou les deux)

Par ex :

$141\text{G} \times 4 = 564\text{G}$ de capacité de modèle

$141\text{G} \times 4 \times 2 = 1128\text{G}$ de capacité de modèle



L'IA est un outil...

Vous n'allez pas être remplacés demain par une IA.
L'IA vous offre des gains de productivité...

Dans nos métiers :

Génération automatisée de configuration et de règles de firewall...

Root-Cause analysis sur des sources uniques ou multiples de logs

Aide au support technique (chatbot, smart search engine, filtrage d'astreinte intelligent, etc)

Le futur

- La fin des GPUs ? Cerebras, Groq, Etched...
- Des modèles plus petits et meilleurs :
Google Gemma 3 (27B), Deepseek, Inception AI
dLLM...
- Artificial General Intelligence (AGI)
Level 3 → 90th percentile of skilled adults
ou Human-Level AI (Yann Le Cun)
→ Years to go...



Et sinon, je fais quoi dans la vie ?

Nous facilitons l'intégration, sans compromis sur la sécurité, d'Intelligence Artificielle, dans vos applications métiers.

Une plateforme clé en main d'IA Européenne, souveraine, non-soumise au Cloud Act.

API OCR+RAG chiffré+LLM plug&play :
Smart Search, Chatbot, etc.



Questions ?

contact@vauban.cloud

