

# La face cachée de l'intelligence artificielle

consommation énergétique des algorithmes de IAg

Rémi Emonet – 2025-11-10 – UEOS-TEDS

UJM FST / Lab. Hubert Curien





1. Impact environnemental du numérique ?  
... et autres activités
2. Qu'est-ce que l'intelligence artificielle (IA) ?  
Faire apprendre une machine à partir d'exemples
3. Coût énergétique de l'IA générative ?  
Entraînement et utilisation
4. Réduire l'impact environnemental de l'IA générative ?  
Pistes et défis

# Attention : la mesure précise est complexe

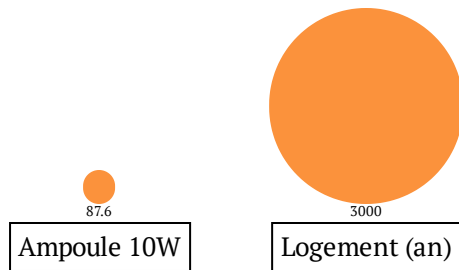
- les chiffres donnés sont des ordres de grandeur
- les sources varient sur leur date, leur méthodologie, leur précision
- les technologies évoluent rapidement

# Impact environnemental du numérique ?

... et autres activités

# Quelques Ordres de Grandeur (Énergie)

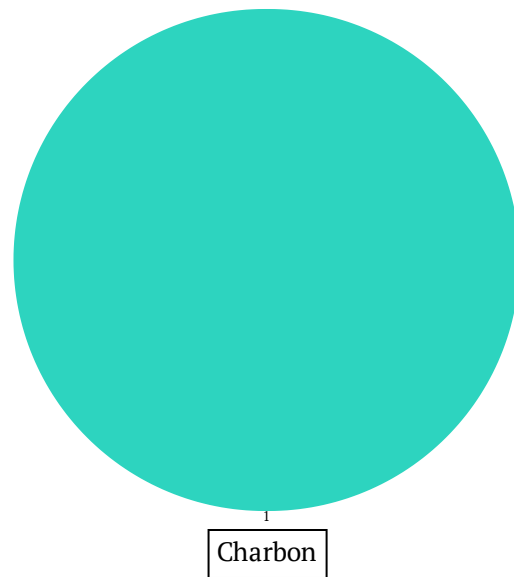
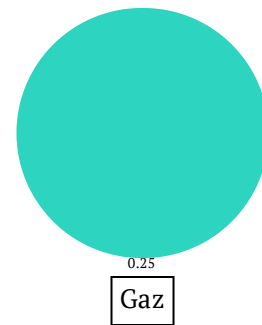
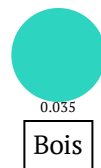
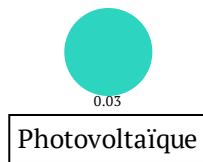
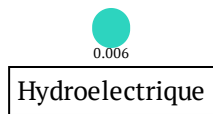
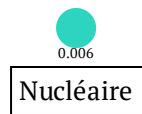
annuel, par personne, en kWh, en France



Production / consommation

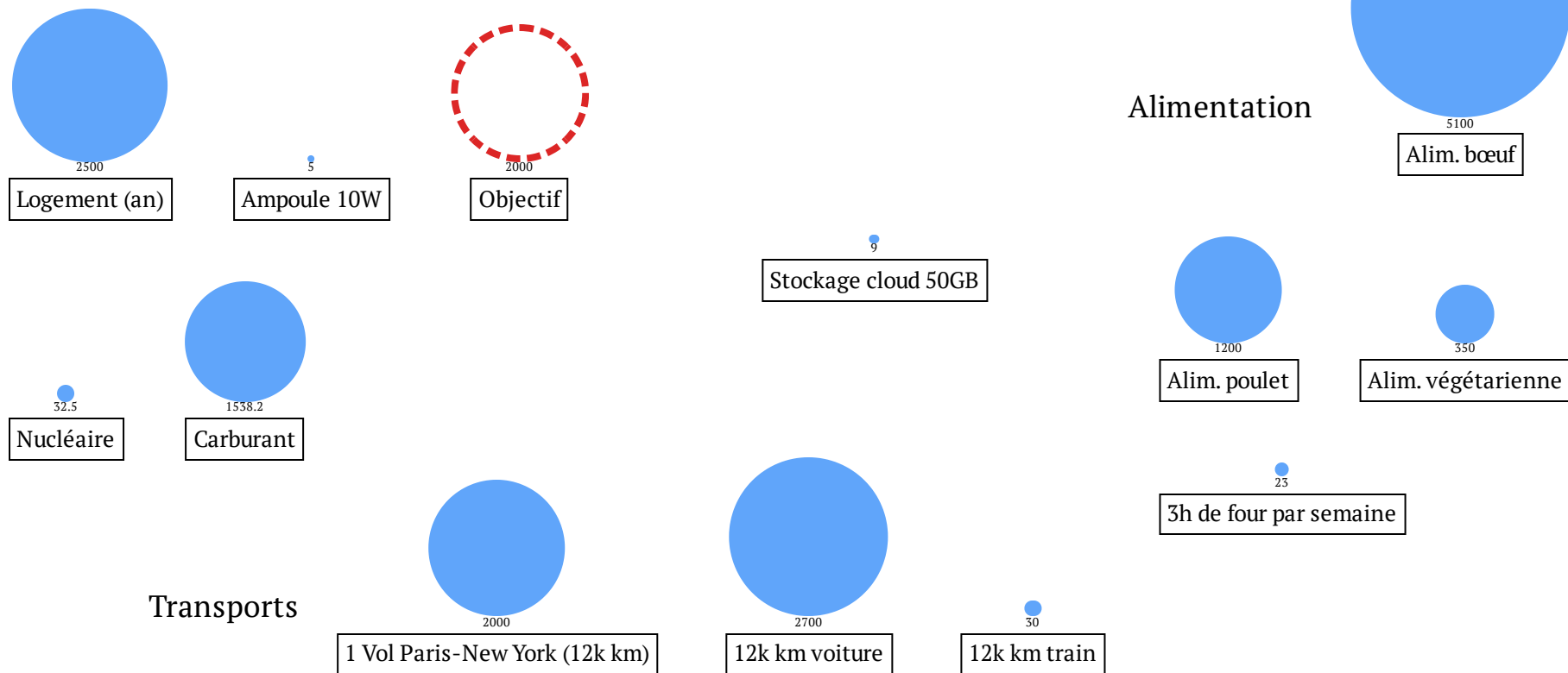
# Conversion kWh $\rightarrow$ Eq. CO<sub>2</sub>

kg Eq. CO<sub>2</sub> par kWh



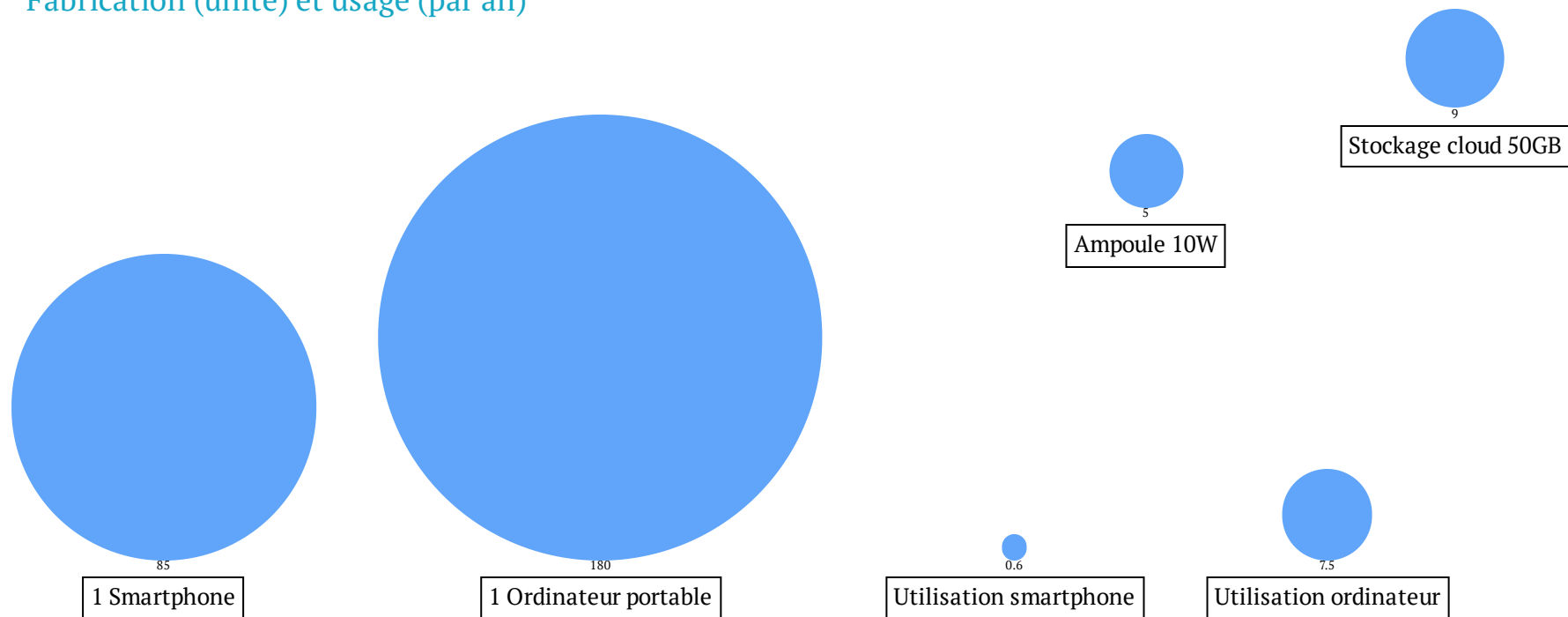
# Quelques Ordres de Grandeur (Eq. CO<sub>2</sub>)

annuel, en kg de CO<sub>2</sub>, par personne, en France



# Impact du matériel numérique ?

Fabrication (unité) et usage (par an)



Ère du zettaoctet, 1E = 1,000,000,000,000,000,000,000 (10<sup>21</sup>) octets

⇒ importance du débit, dominé par le streaming vidéo





1. Impact environnemental du numérique ?  
... et autres activités
2. Qu'est-ce que l'intelligence artificielle (IA) ?  
Faire apprendre une machine à partir d'exemples
3. Coût énergétique de l'IA générative ?  
Entraînement et utilisation
4. Réduire l'impact environnemental de l'IA générative ?  
Pistes et défis

# Qu'est-ce que l'intelligence artificielle (IA) ?

Faire apprendre une machine à partir d'exemples



# Limites de la programmation classique

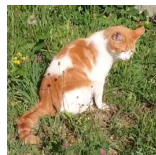
- complexité des programmes
- tâche impossible à décrire
- des « if » pour analyser une image ?



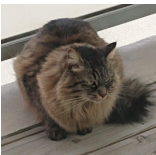
121	155	159	137	92	64	48	35	43	57	53	42	92	138	106	122	122	476	169	155	92	80	144	138	170	96	84	47	112	126	98	105	91	107	125	107	108	76
118	120	120	101	93	143	123	138	142	140	138	134	128	135	105	107	149	166	149	143	114	166	137	131	149	115	99	109	129	118	139	140	165	107	133	137	137	
109	108	111	94	88	32	61	114	107	123	126	154	124	133	112	116	162	137	117	175	103	88	135	157	119	80	93	78	81	99	130	137	117	3	137	137		
109	108	111	94	88	32	61	114	107	123	126	154	124	133	112	116	162	137	117	175	103	88	135	157	119	80	93	78	81	99	130	137	117	3	137	137		
109	108	111	94	88	32	61	114	107	123	126	154	124	133	112	116	162	137	117	175	103	88	135	157	119	80	93	78	81	99	130	137	117	3	137	137		
109	108	111	94	88	32	61	114	107	123	126	154	124	133	112	116	162	137	117	175	103	88	135	157	119	80	93	78	81	99	130	137	117	3	137	137		
109	108	111	94	88	32	61	114	107	123	126	154	124	133	112	116	162	137	117	175	103	88	135	157	119	80	93	78	81	99	130	137	117	3	137	137		
109	108	111	94	88	32	61	114	107	123	126	154	124	133	112	116	162	137	117	175	103	88	135	157	119	80	93	78	81	99	130	137	117	3	137	137		
109	108	111	94	88	32	61	114	107	123	126	154	124	133	112	116	162	137	117	175	103	88	135	157	119	80	93	78	81	99	130	137	117	3	137	137		
109	108	111	94	88	32	61	114	107	123	126	154	124	133	112	116	162	137	117	175	103	88	135	157	119	80	93	78	81	99	130	137	117	3	137	137		
109	108	111	94	88	32	61	114	107	123	126	154	124	133	112	116	162	137	117	175	103	88	135	157	119	80	93	78	81	99	130	137	117	3	137	137		
109	108	111	94	88	32	61	114	107	123	126	154	124	133	112	116	162	137	117	175	103	88	135	157	119	80	93	78	81	99	130	137	117	3	137	137		
109	108	111	94	88	32	61	114	107	123	126	154	124	133	112	116	162	137	117	175	103	88	135	157	119	80	93	78	81	99	130	137	117	3	137	137		
109	108	111	94	88	32	61	114	107	123	126	154	124	133	112	116	162	137	117	175	103	88	135	157	119	80	93	78	81	99	130	137	117	3	137	137		
109	108	111	94	88	32	61	114	107	123	126	154	124	133	112	116	162	137	1																			

- $\Rightarrow$  nécessité de programmer par l'exemple

# Création d'un jeu de données (pour la tâche Chat/NonChat)



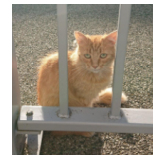
, Chat



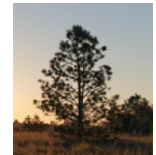
, Chat



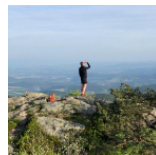
, NonChat



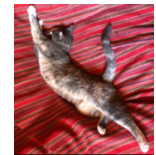
, Chat



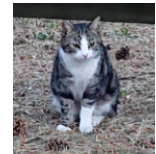
, NonChat



, NonChat



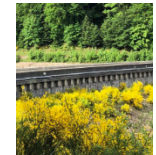
, Chat



, Chat

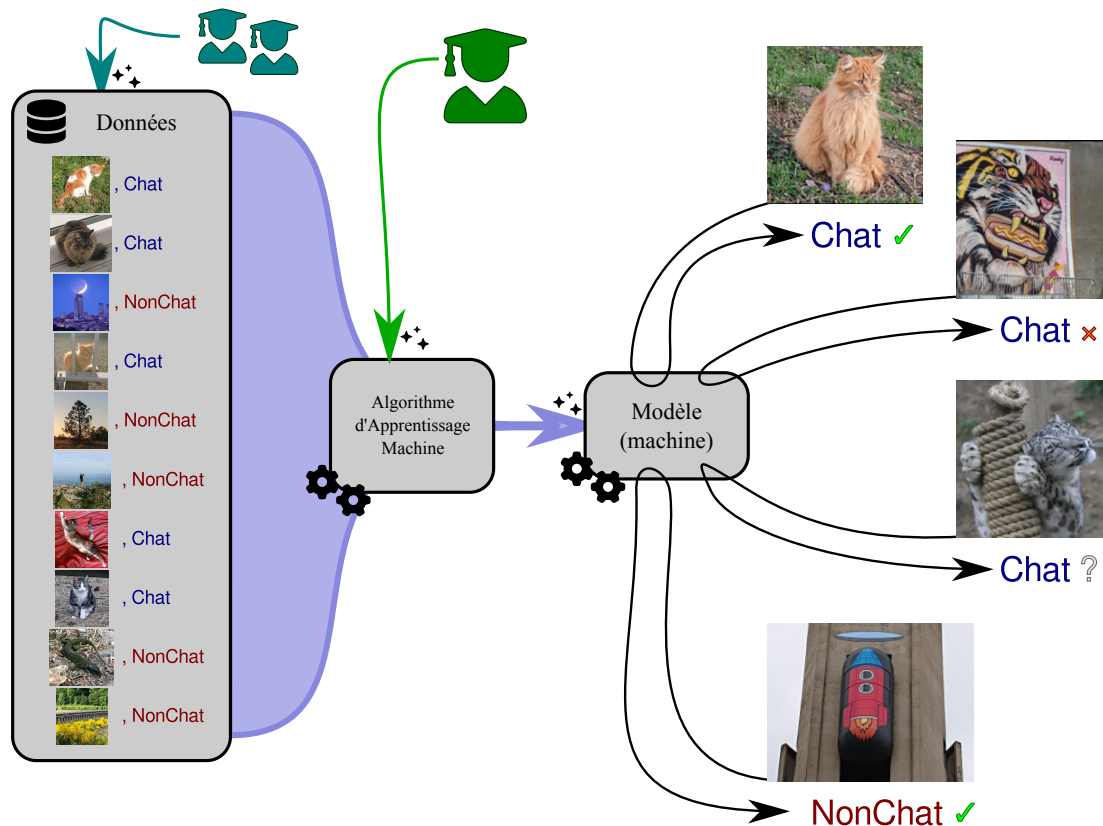


, NonChat

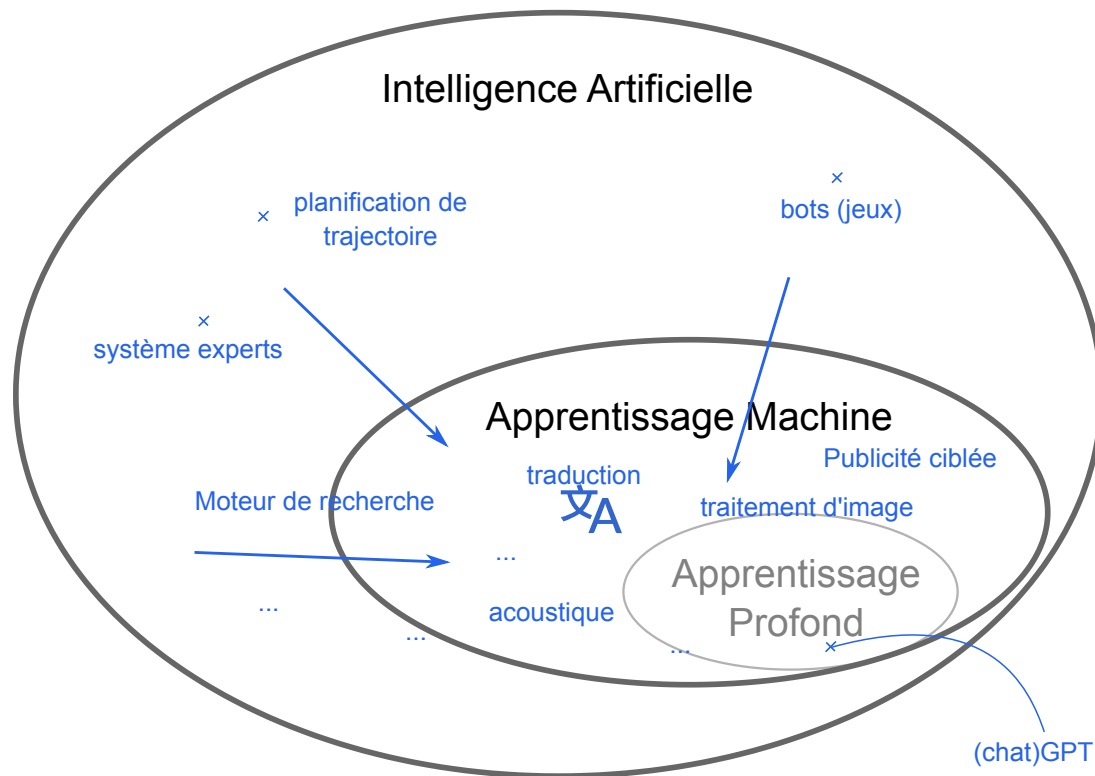


, NonChat

# Apprentissage automatique : principe global

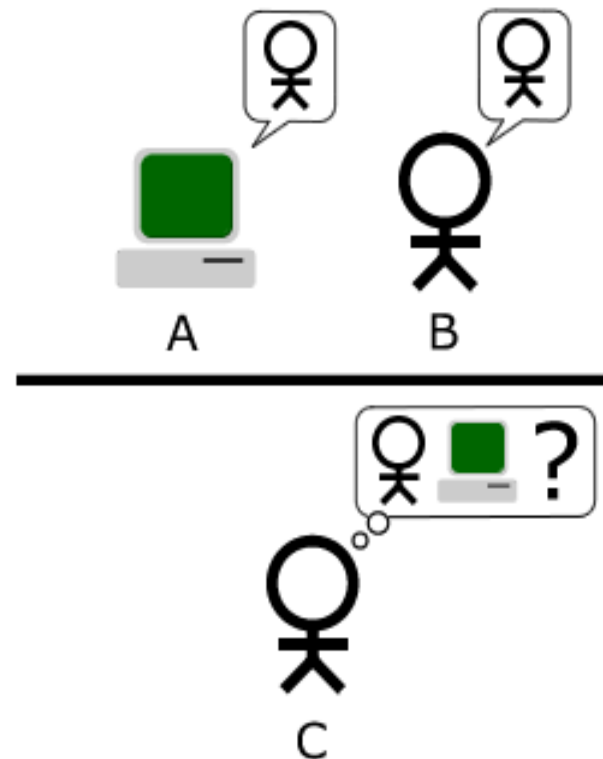


# Terminologie



# Intelligence artificielle ? Le test de Turing

- interaction textuelles entre humain et machine
- l'humain arrive-t-il à savoir qu'il interagit avec une machine ?



# Intelligence ? Zoom sur chatGPT ?

- Modèle de langage : apprendre à prédire le mot manquant

Saint-Étienne appelée « Sainté » en langage familier, est une commune française située au sud-ouest de Lyon (60 km environ) et le quart-sud-est de la France, en région Auvergne-Rhône-Alpes. C'est le chef-lieu du département de la Loire. Avec 174 082 habitants en 2020, elle est la 13e commune la plus peuplée de France (2016) et la 2e commune d'Auvergne-Rhône-Alpes. Saint-Étienne Métropole constitue par sa population (400 813 habitants en 2020) la 3e métropole régionale après la métropole de Grenoble Alpes et la métropole de Lyon. La commune est par ailleurs au cœur d'une vaste aire urbaine de plus de 520 640 habitants en 2017, la dix-septième de France par sa population, regroupant 117 communes.

- Des quantités de données
  - tout **wikipedia**
  - **et 100× plus** avec le web, des livres, etc
- La partie *chat* : faire un agent conversationnel
  - InstructGPT
  - étiquetage manuel d'interactions (~100k)



# Défis et difficultés de l'apprentissage automatique

- concevoir des algorithmes d'apprentissage
- prouver que ces algorithmes marchent
- évaluer ces algorithmes
- défis divers
  - biais des données
  - interprétabilité/explicabilité
  - adaptation/transfert
  - attaques et manipulation
  - *efficacité en données*
  - *efficacité en calcul/énergie*
- choix sociétaux et législation

# Et la génération d'image ?

- Données
  - des paires (image, description textuelle)
  - e.g. LAION-5B: **5 milliards d'images**
- Modèles
  - e.g. DALL·E 2:  
modèle propriétaire d'OpenAI
  - e.g. Midjourney:  
modèle propriétaire
  - e.g. Stable Diffusion:  
modèle open-source très populaire





1. Impact environnemental du numérique ?  
... et autres activités
2. Qu'est-ce que l'intelligence artificielle (IA) ?  
Faire apprendre une machine à partir d'exemples
3. Coût énergétique de l'IA générative ?  
Entraînement et utilisation
4. Réduire l'impact environnemental de l'IA générative ?  
Pistes et défis

# Coût énergétique de l'IA générative ?

Entraînement et utilisation

# Entraînement et Inférence

## Inférence

- utilisation du modèle
- phase de génération
- utilisateurs finaux
- volume d'utilisation



## Entraînement

- création du modèle
- phase d'apprentissage
- données massives

# Impact de l'entraînement

tonnes Eq. CO<sub>2</sub>

- Entraînement GPT-3  $\approx 500$  t Eq. CO<sub>2</sub>

502

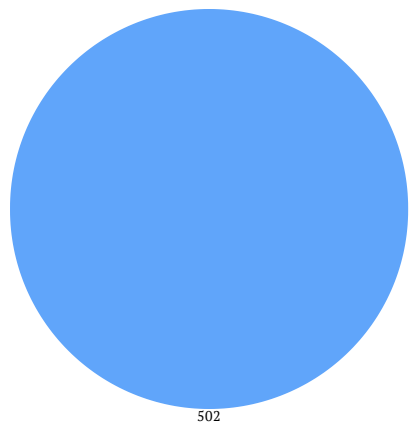
Entraînement GPT-3

2277000

Production (annuelle) Nucléaire France

# Coût de l'entraînement

tonnes Eq. CO<sub>2</sub>

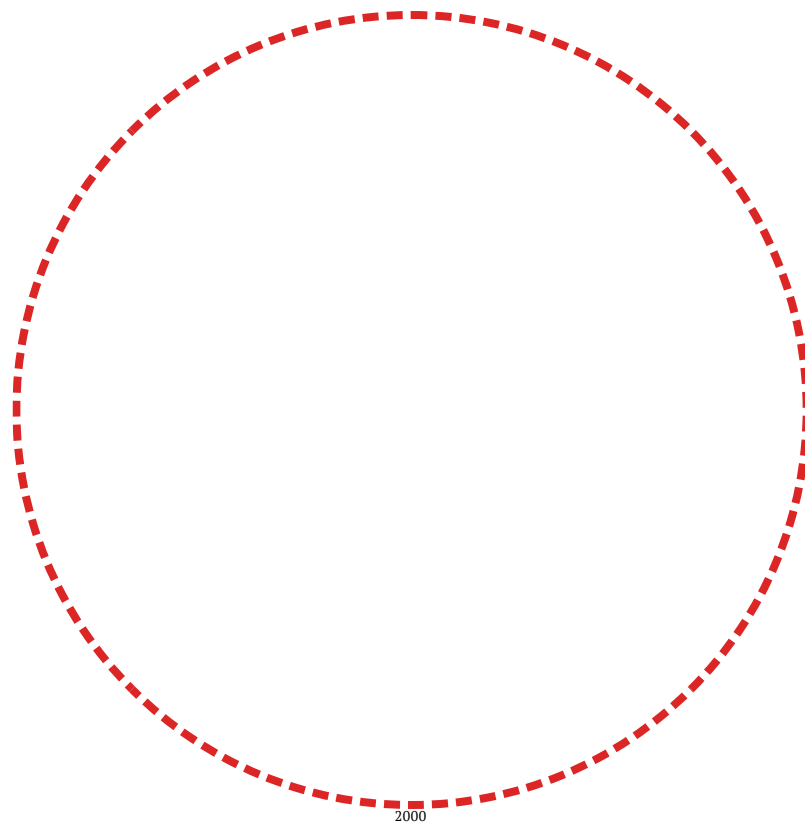


Entraînement GPT-3



32.5

Électricité annuelle 1000 personnes



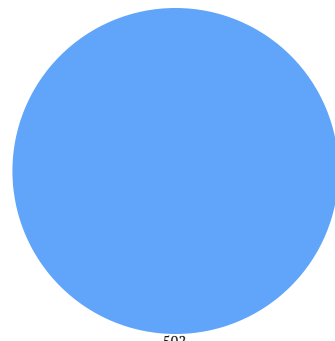
2000

Objectif 1000 personnes

# L'exemple de Llama

transparence sur les coûts, deux version 8B et 70B

- 8 milliards de paramètres (8B) ou 70 milliards (70B)
- Jeu de données
  - 15T "tokens" (parties de mots)
  - soit  $15 \times 10^{12}$
- coût (en t Eq CO<sub>2</sub>)
  - matériel (GPU)
  - calcul
  - **climatisation**



502

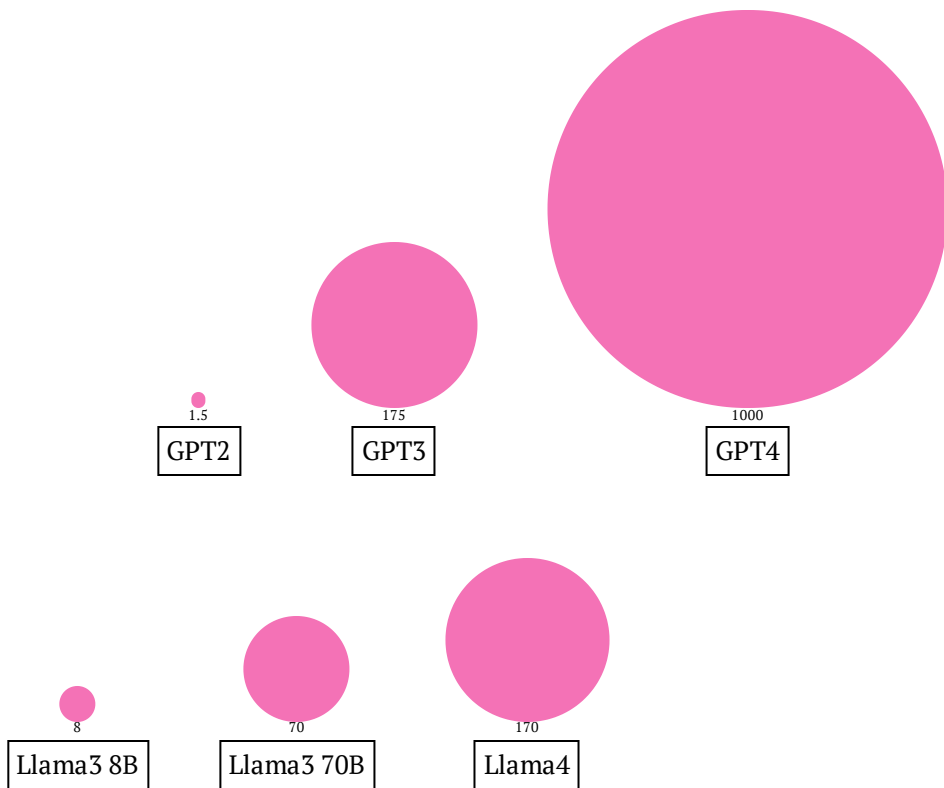
Entraînement GPT-3





# Stockage des modèles

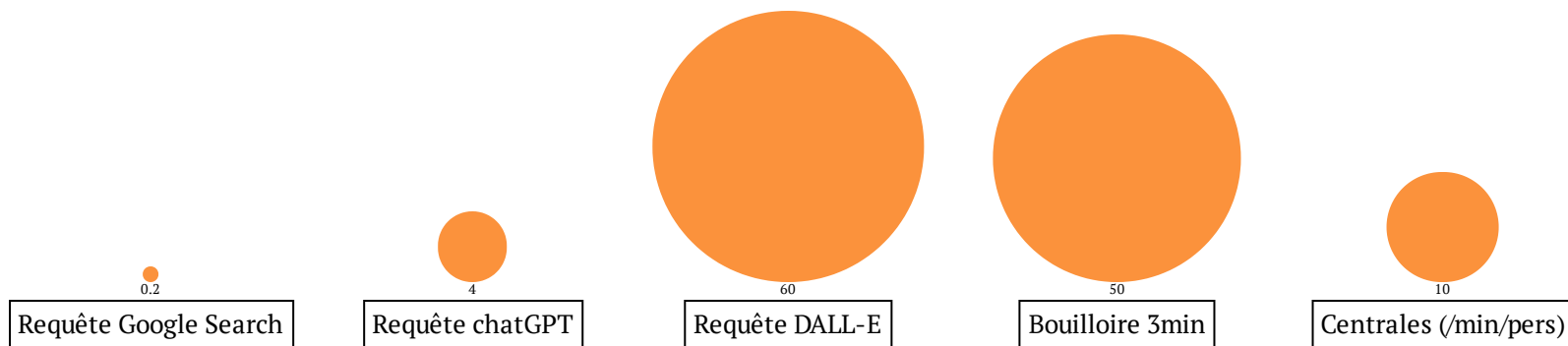
- Dépend du nombre de paramètres
  - GPT-2 : 1,5 milliards de paramètres (1.5B)
  - GPT-3 : 175 milliards de paramètres (175B)
  - GPT-4 : 1 trillion de paramètres (1000B)
  - GPT-5 : ...
- Exemple de Llama3
  - estimation de taille
    - 8B : ~60 Go
    - 70B : ~500 Go
  - compression possible
    - 8B : 5Go, voir 2Go
    - 70B : similaire



# Coût énergétique de l'inférence

Wh par requête

- 1 requête chat GPT = 10 requêtes Google Search





1. Impact environnemental du numérique ?  
... et autres activités
2. Qu'est-ce que l'intelligence artificielle (IA) ?  
Faire apprendre une machine à partir d'exemples
3. Coût énergétique de l'IA générative ?  
Entraînement et utilisation
4. Réduire l'impact environnemental de l'IA générative ?  
Pistes et défis

# Réduire l'impact environnemental de l'IA générative ?

Pistes et défis

# Pistes pour réduire l'impact

- Coté modèles / algorithmes / entraînement :
  - Optimiser les algorithmes d'entraînement
  - Réutiliser des modèles pré-entraînés
  - Réduire la taille des modèles
  - Réduire la quantité de données d'entraînement
- Coté infrastructure :
  - Utiliser des sources d'énergie renouvelable
  - Utiliser le matériel de l'utilisateur
  - Améliorer l'efficacité matérielle
- Coté usage :
  - Réduire le nombre de requêtes
  - Réduire la taille des réponses
- Exemples :
  - Mélange d'experts
  - Distillation de modèles
  - Quantification et compression

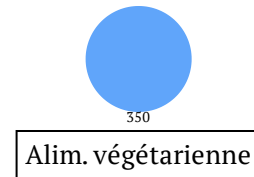
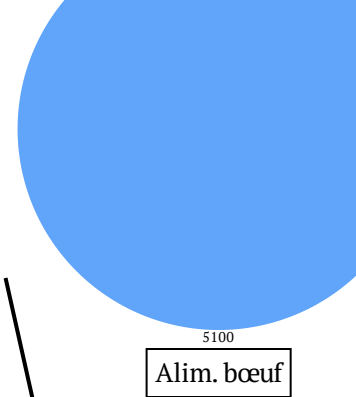
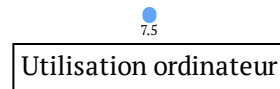
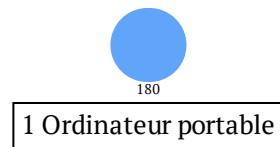
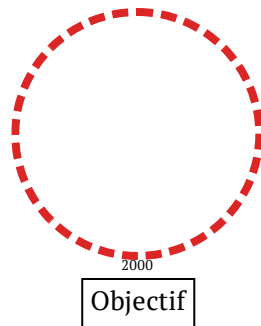
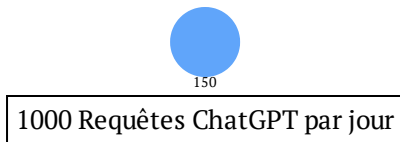
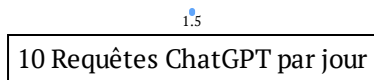
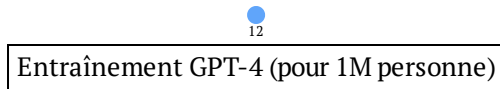


1. Impact environnemental du numérique ?  
... et autres activités
2. Qu'est-ce que l'intelligence artificielle (IA) ?  
Faire apprendre une machine à partir d'exemples
3. Coût énergétique de l'IA générative ?  
Entraînement et utilisation
4. Réduire l'impact environnemental de l'IA générative ?  
Pistes et défis

# Conclusions

# Conclusions

(par an et/ou personne)



Merci !

