

Intelligence Artificielle: la machine va-t-elle remplacer l'homme ?

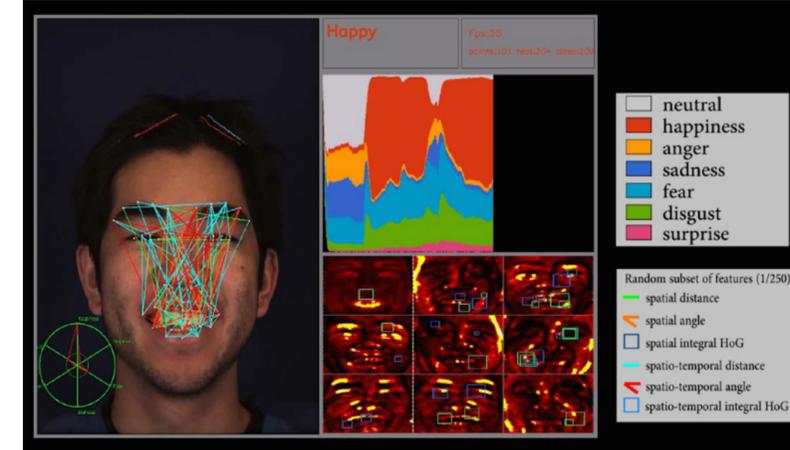
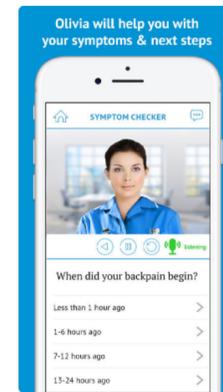
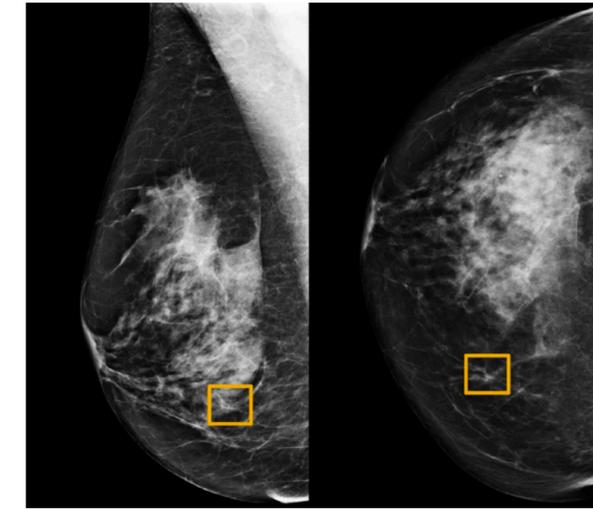
Fondements scientifiques, Enjeux éthiques, Réglementation

Raja Chatila
Institut des Systèmes Intelligents et de Robotique (ISIR)
Sorbonne Université, Paris

Raja.Chatila@sorbonne-universite.fr

Secteurs d'application de l'IA et de la Robotique

- Transports, logistique
- Santé et soin
- Fabrication manufacturière
- Gestion de déchets
- Gestion des réseaux électriques
- Agriculture
- Services personnels et assistance
- e-Commerce, publicité, recommandations
- Procédures administratives
- Processus techniques, codage
- Recrutement & management
- Assurance & finance
- Formation
- Recherche scientifique
- Justice
- Sécurité
- Armement
- ...



AI bias

Can you make AI fairer than a judge? Play our courtroom algorithm game

The US criminal legal system uses predictive algorithms to try to make the judicial process less biased. But there's a deeper problem.



25 YEARS OLD



x2



Aux origines de l'IA

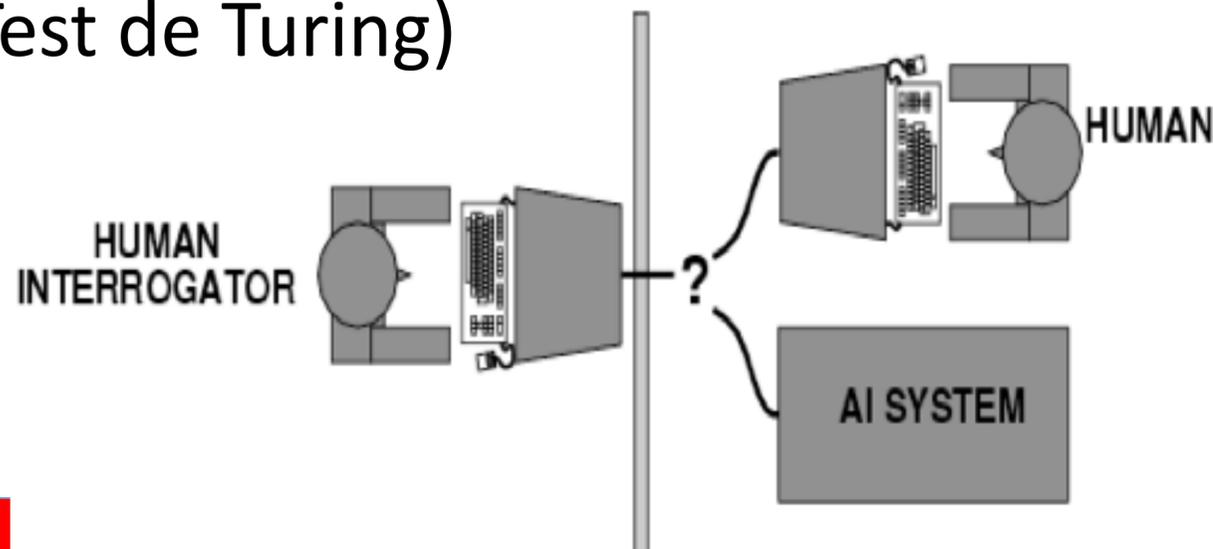
1936: Thèse de Church-Turing: Les fonctions calculables sont exactement les fonctions récursives, c-à-d représentables par un algorithme qui se termine.

Machine Universelle de Turing - modèle formel de l'ordinateur.

1948 : "Intelligent Machinery"

1950 : "Computing Machinery and Intelligence"

- le jeu de l'imitation (Test de Turing)



Alan Turing 1912-1954

VOL. LIX. No. 236.]

[October, 1950

MIND

A QUARTERLY REVIEW

OF

PSYCHOLOGY AND PHILOSOPHY

I.—COMPUTING MACHINERY AND INTELLIGENCE

By A. M. TURING

1. *The Imitation Game.*

I PROPOSE to consider the question, 'Can machines think?' This should begin with definitions of the meaning of the terms 'machine' and 'think'. The definitions might be framed so as to reflect so far as possible the normal use of the words, but this attitude is dangerous. If the meaning of the words 'machine' and 'think' are to be found by examining how they are commonly used it is difficult to escape the conclusion that the meaning and the answer to the question, 'Can machines think?' is to be sought in a statistical survey such as a Gallup poll. But this is absurd. Instead of attempting such a definition I shall replace the question by another, which is closely related to it and is expressed in relatively unambiguous words.

The new form of the problem can be described in terms of a game which we call the 'imitation game'. It is played with three people, a man (A), a woman (B), and an interrogator (C) who may be of either sex. The interrogator stays in a room apart from the other two. The object of the game for the interrogator is to determine which of the other two is the man and which is the woman. He knows them by labels X and Y, and at the end of the game he says either 'X is A and Y is B' or 'X is B and Y is A'. The interrogator is allowed to put questions to A and B thus:

C: Will X please tell me the length of his or her hair?

Now suppose X is actually A, then A must answer. It is A's

28

433

Aux origines de l'IA

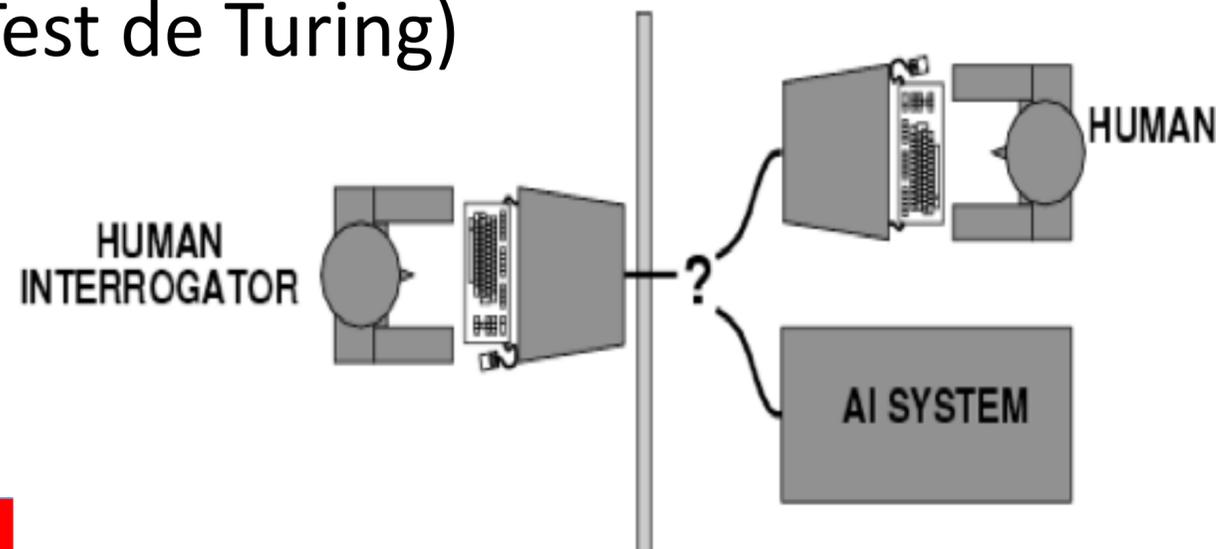
1936: Thèse de Church-Turing: Les fonctions calculables sont exactement les fonctions récursives, c-à-d représentables par un algorithme qui se termine.

Machine Universelle de Turing - modèle formel de l'ordinateur.

1948 : "Intelligent Machinery"

1950 : "Computing Machinery and Intelligence"

- le jeu de l'imitation (Test de Turing)



Alan Turing 1912-1954

VOL. LIX. No. 236.]

[October, 1950

MIND

A QUARTERLY REVIEW

OF

PSYCHOLOGY AND PHILOSOPHY

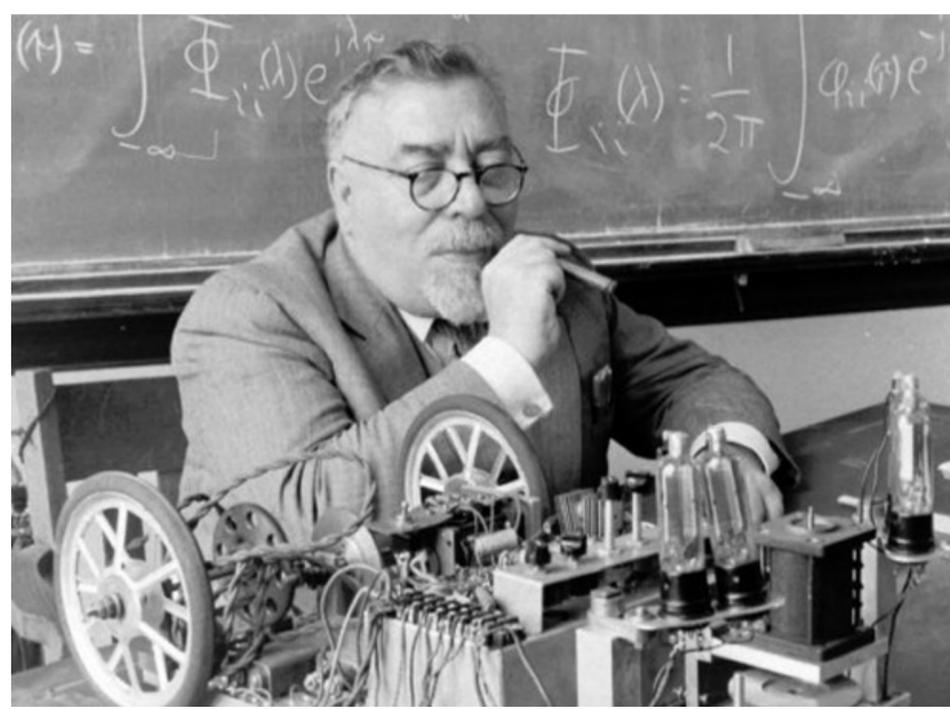
I.—COMPUTING MACHINERY AND INTELLIGENCE

Je me propose d'examiner la question suivante : « Les machines peuvent-elles penser ? ». Il faudrait commencer par définir le sens des termes « machine » et « penser ». Les définitions pourraient être formulées de manière à refléter autant que possible l'usage normal des mots, mais cette attitude est dangereuse ...

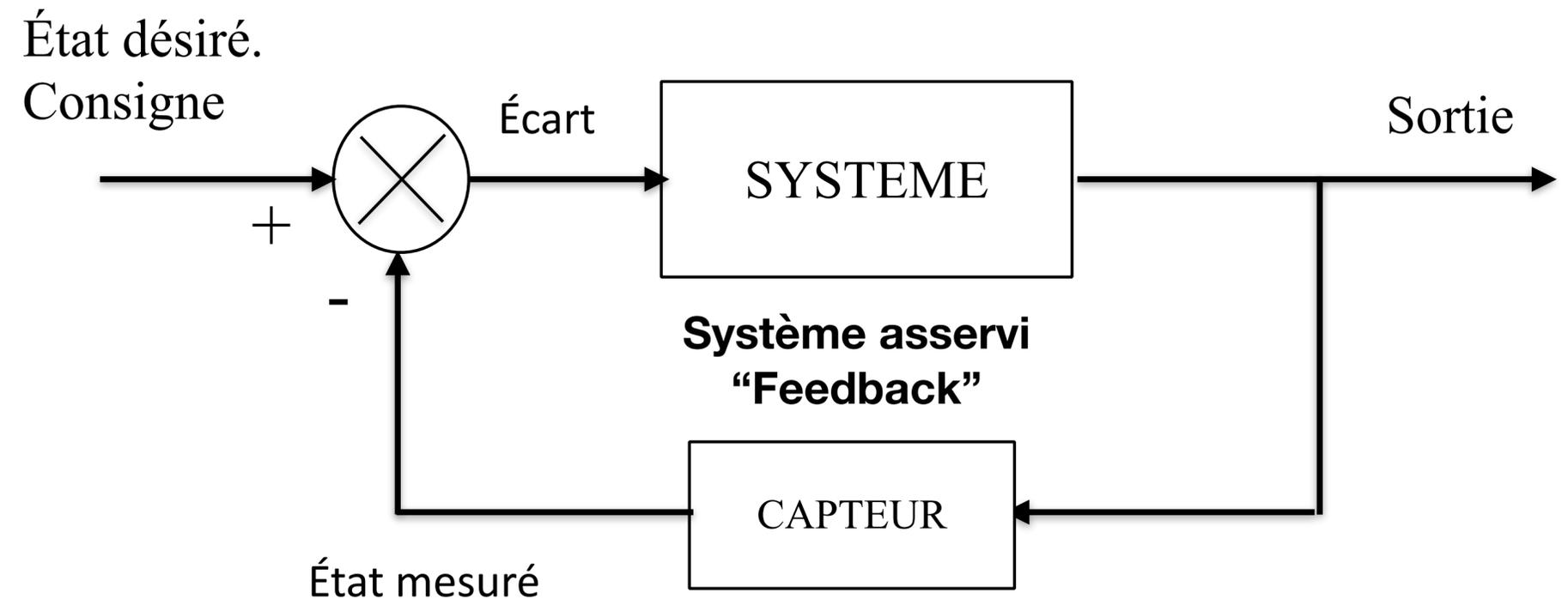
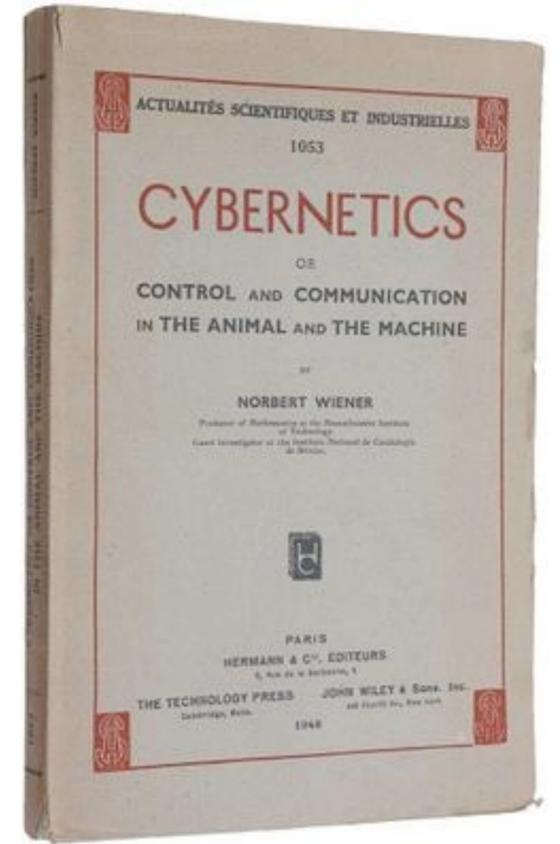
La Cybernétique

- “Cybernétique ou le contrôle et la communication dans l’animal et la machines” (Cybernetics, or Control and Communication in the Animal and the Machine, 1948). Homéostat, boucle fermée, systèmes à commande automatique.

- “Cybernétique et Société - L’usage humain des êtres humains” (The Human use of Human Beings. Cybernetics and Society, 1950).



Norbert Wiener (1894 - 1964)





Boston Dynamics

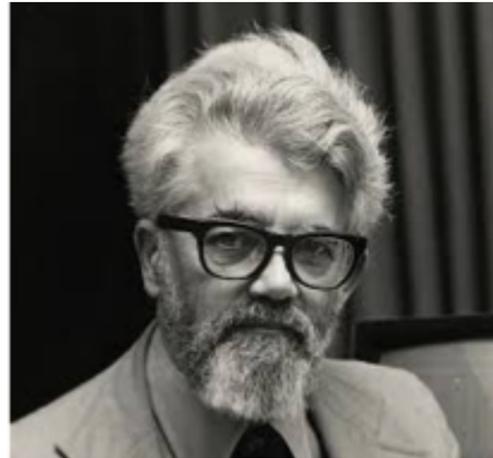


Boston Dynamics



Boston Dynamics

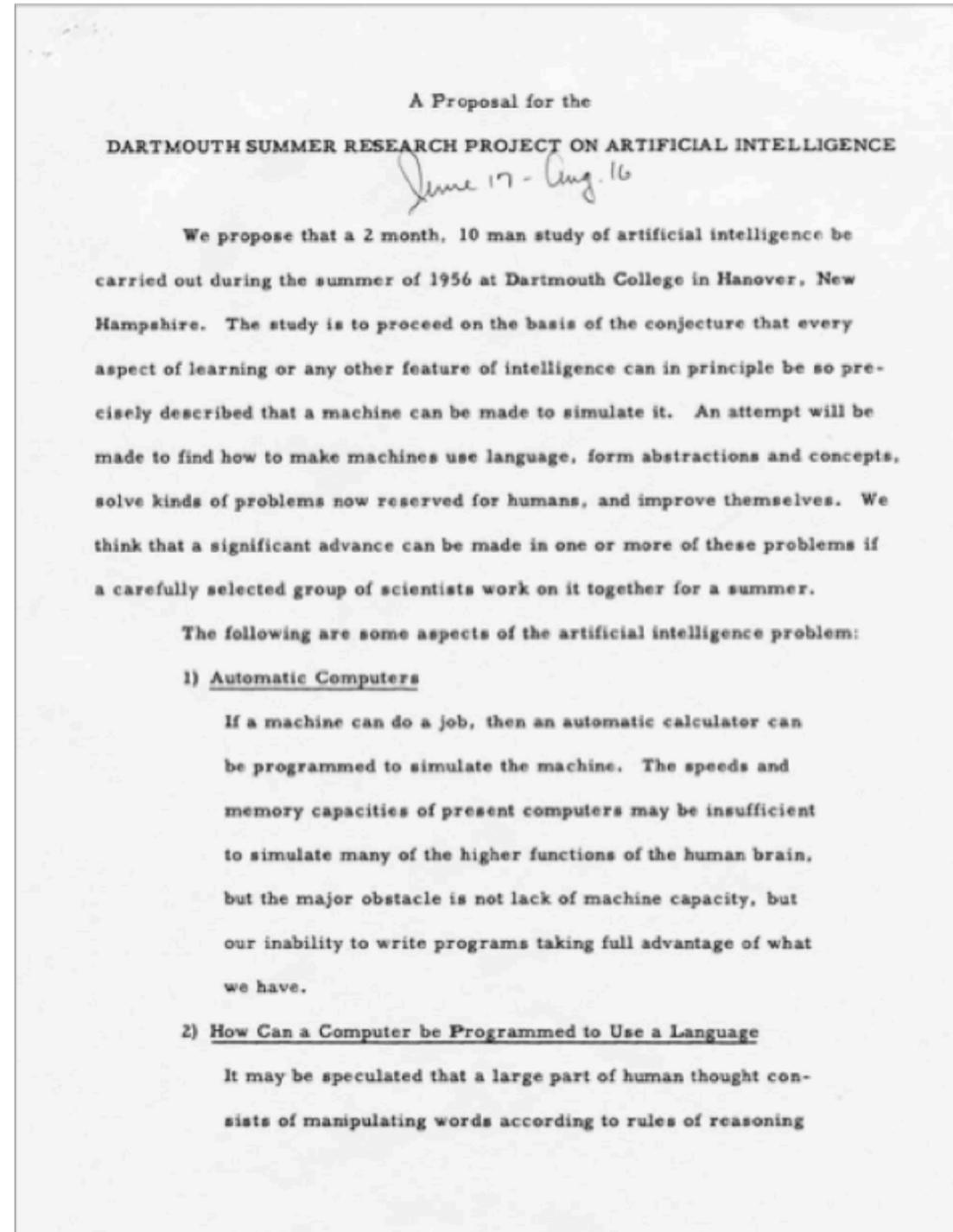
Naissance de l'Intelligence Artificielle: La conférence de Dartmouth College, 1956



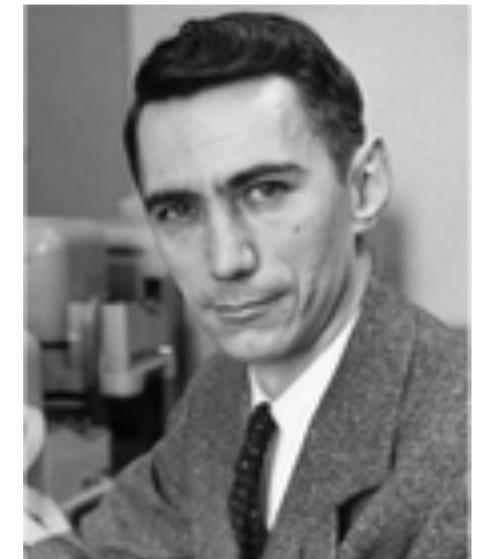
John McCarthy



Marvin Minsky

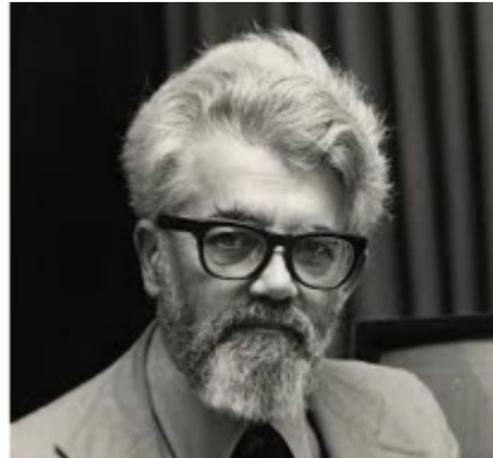


Nathaniel Rochester



Claude Shannon

Naissance de l'Intelligence Artificielle: La conférence de Dartmouth College, 1956



John McCarthy

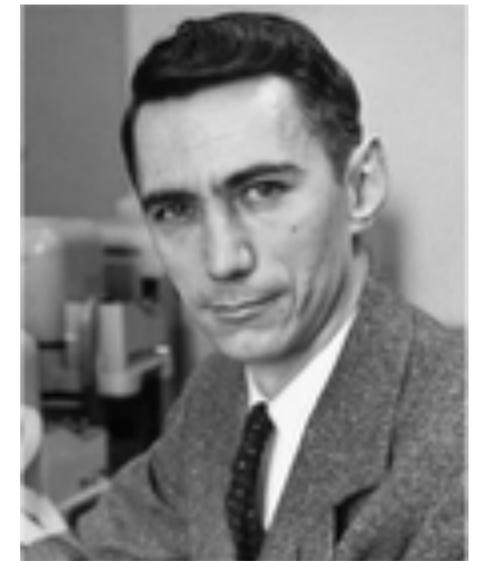


Marvin Minsky

Nous proposons qu'une étude de l'intelligence artificielle soit menée par 10 personnes pendant 2 mois à l'été de 1956... L'étude se fondera sur la conjecture que chaque aspect de l'apprentissage ou toute autre caractéristique de l'intelligence peut en principe être décrits avec une telle précision qu'une machine peut être fabriquée pour les simuler. On tentera de trouver comment faire en sorte que les machines utilisent le langage, forment des abstractions et des concepts, résolvent des types de problèmes aujourd'hui réservés aux humains, et s'améliorent elles-mêmes.

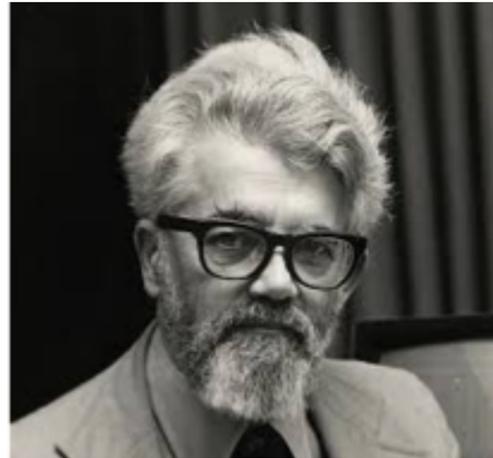


Nathaniel Rochester



Claude Shannon

Naissance de l'Intelligence Artificielle: La conférence de Dartmouth College, 1956



John McCarthy



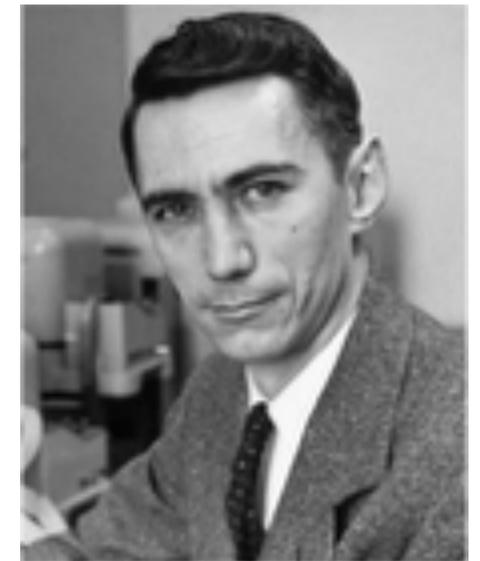
Marvin Minsky

Nous proposons qu'une étude de l'intelligence artificielle soit menée par 10 personnes pendant 2 mois à l'été de 1956... L'étude se fondera sur la conjecture que chaque aspect de l'apprentissage ou toute autre caractéristique de l'intelligence en principe être décrits avec une telle précision qu'une machine peut être fabriquée pour les simuler. On tentera de trouver comment faire en sorte que les machines utilisent le langage, forment des abstractions et des concepts, résolvent des types de problèmes aujourd'hui réservés aux humains, et s'améliorent elles-mêmes.

ALGORITHMES



Nathaniel Rochester



Claude Shannon

Grands domaines de l'IA

Apprentissage Machine

Systemes connexionnistes

Apprentissage

Profond

Supervisé / non-supervisé

Apprentissage

par

renforcement

Robotique

Machine physique
dans le monde réel

Perception

Décision

Action

Interaction

Automatique

Mécanique

Electronique

Informatique

Temps réel

Sûreté

Systemes

IA Symbolique

Représentation des connaissances

Inférence logique ou raisonnement probabiliste

Résolution de problèmes, Planification, apprentissage logique inductif, ...

Grands domaines de l'IA

Apprentissage Machine

Systemes connexionnistes

Apprentissage

Profond

Supervisé / non-supervisé

Apprentissage

par
renforcement

Robotique

Machine physique
dans le monde réel

Perception

Décision

Action

Interaction

Automatique

Mécanique

Electronique

Informatique

Temps réel

Sûreté

Systemes

IA Symbolique

Représentation des connaissances

Inférence logique ou raisonnement probabiliste

Résolution de problèmes, Planification, apprentissage logique inductif, ...

Transformation des données en connaissances
explicitement définie par
le concepteur humain

Grands domaines de l'IA

Apprentissage Machine

Systemes connexionnistes

Apprentissage

Profond

Supervisé / non-supervisé

Modèle statistique
des données

Apprentissage
par
renforcement

Robotique

Machine physique
dans le monde réel

Perception

Décision

Action

Interaction

Automatique
Mécanique
Electronique
Informatique
Temps réel
Sûreté
Systèmes

IA Symbolique

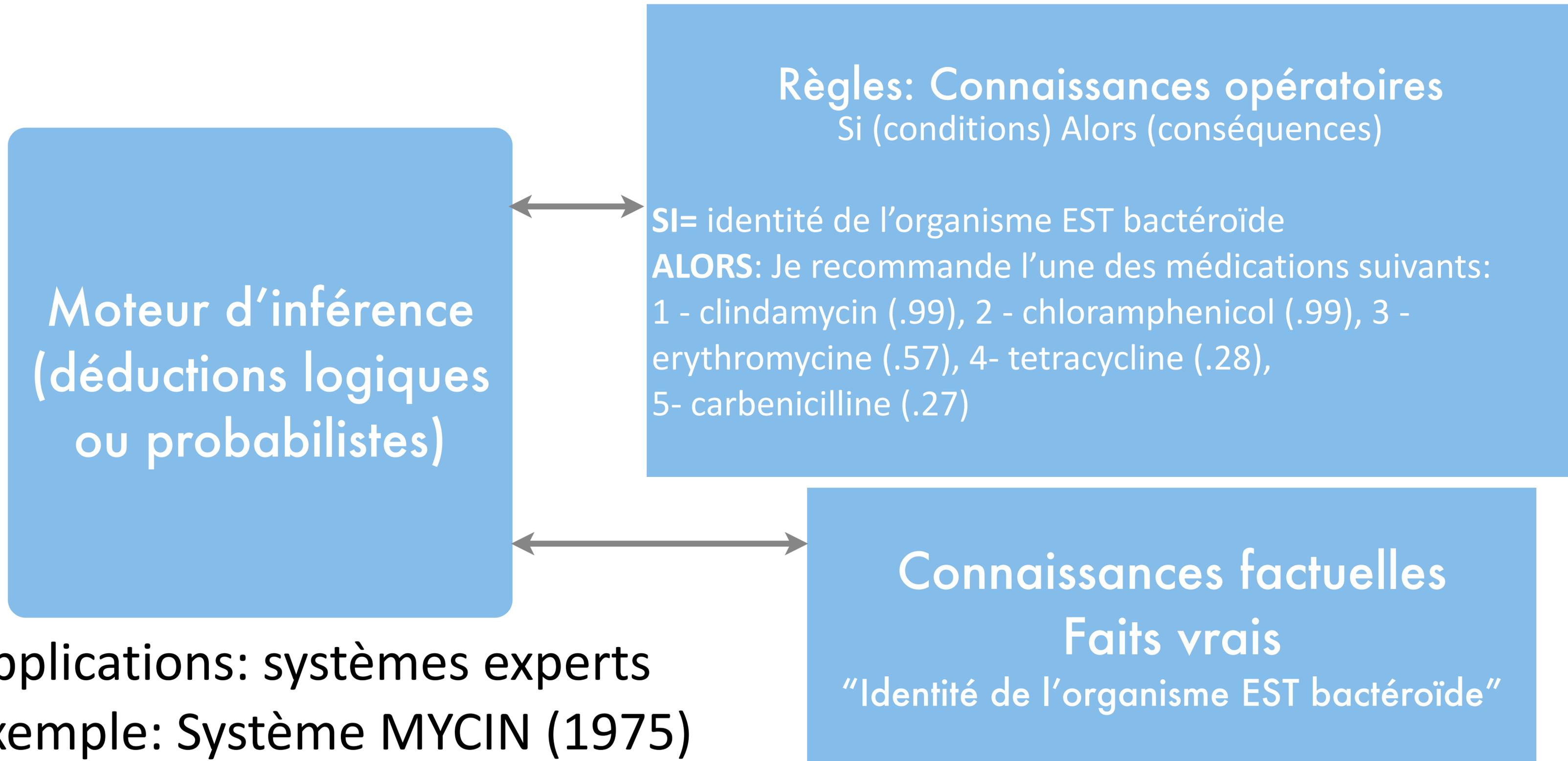
Représentation des connaissances

Inférence logique ou raisonnement probabiliste

Résolution de problèmes, Planification, apprentissage logique inductif, ...

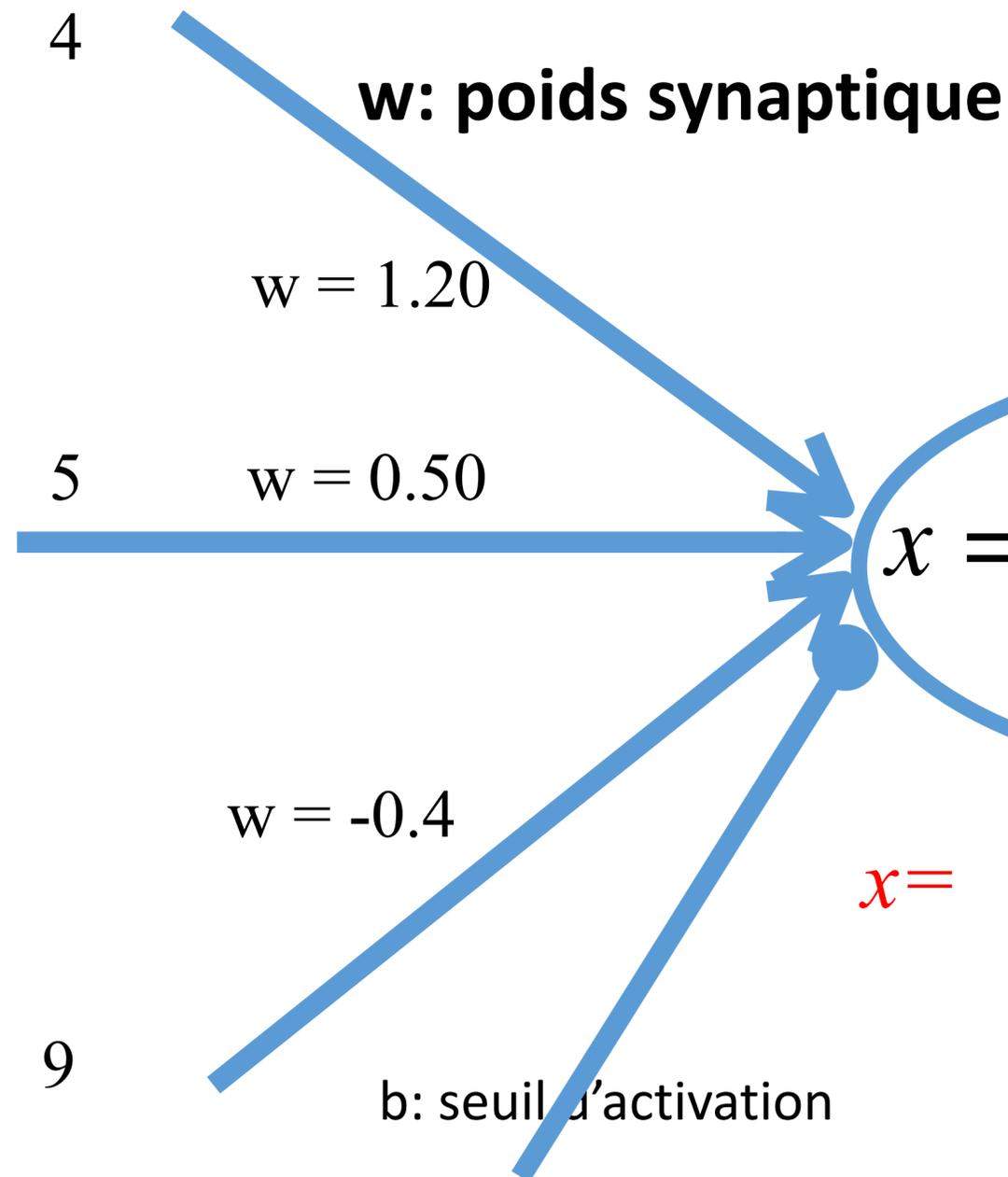
Transformation des données en connaissances
explicitement définie par
le concepteur humain

Systemes à base de règles

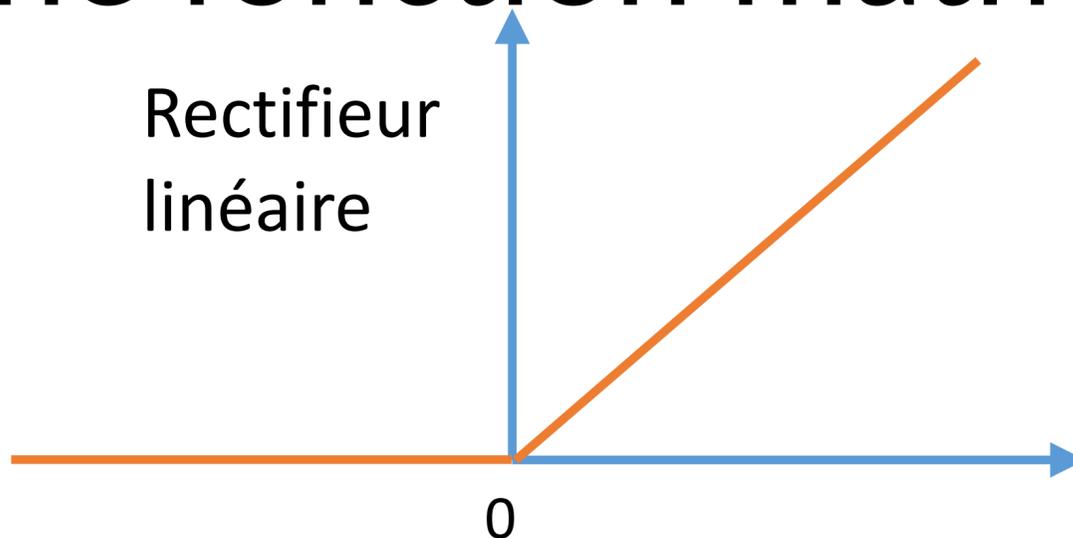


Le neurone formel : une fonction mathématique

Entrées



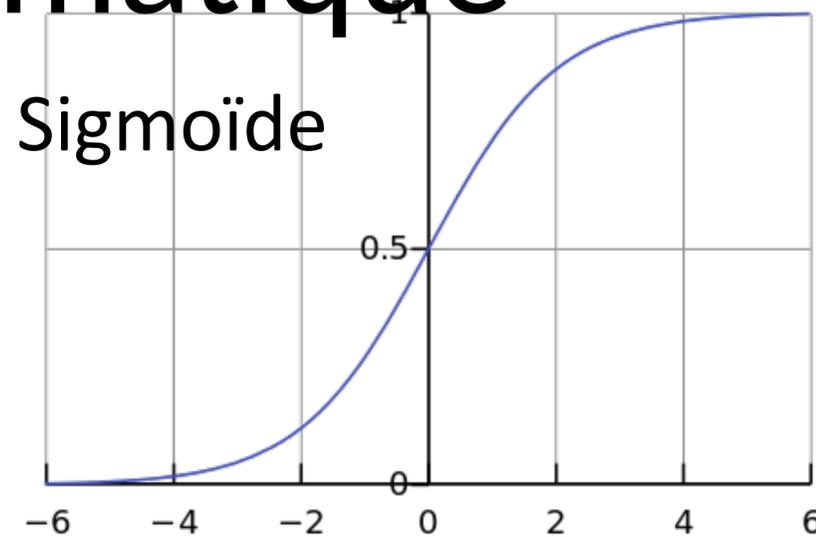
Rectifieur linéaire



$$f(x) = \max(0, x)$$

$$f(3.7) = 3.7$$

Sigmoïde



$$f(x) = \frac{1}{1 + e^{-x}}$$

$$f(3.7) = 0.98$$

$$x = \sum w_k i_k$$

Sortie

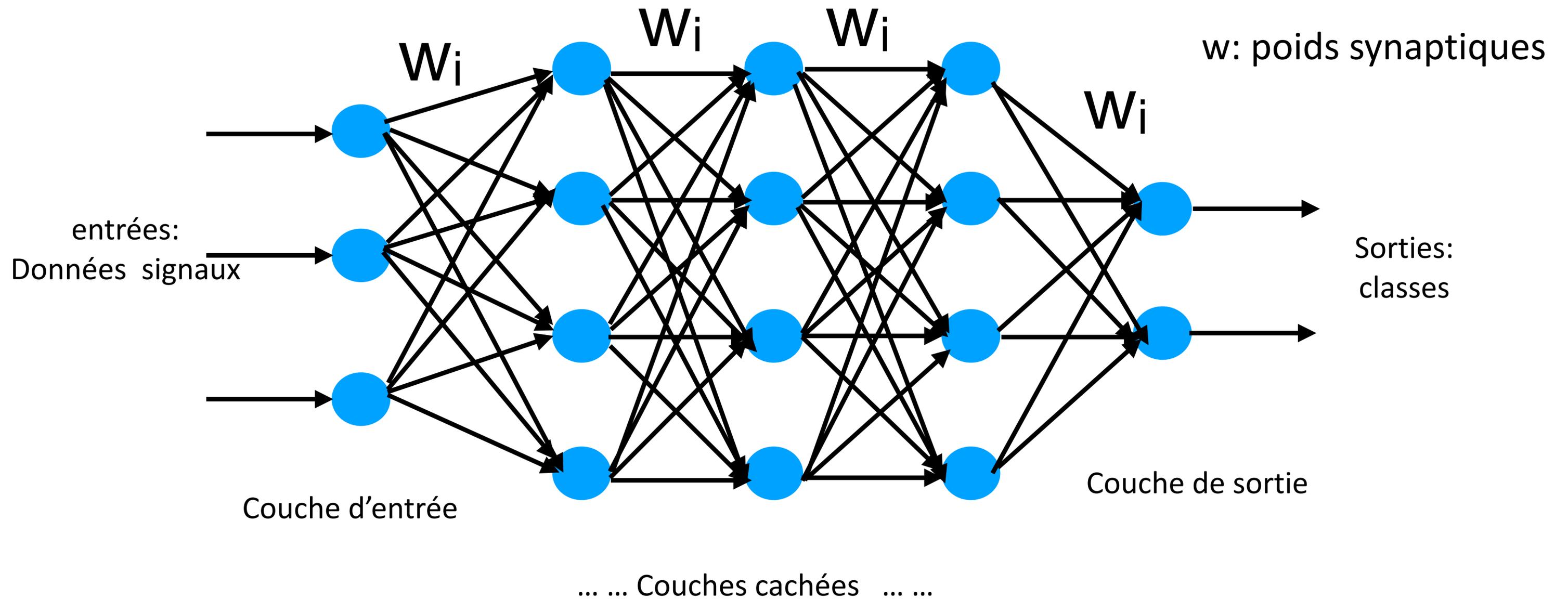
Binaire (> seuil)
ou valeur

$$\begin{aligned}
 x = & 4 \times 1.2 \\
 & + 5 \times 0.5 \\
 & - 9 \times 0.4 = 3.7
 \end{aligned}$$

x = Somme des entrées pondérées par les poids synaptiques

Réseau de Neurones

Millions, milliards, centaines de milliards de paramètres selon la taille du réseau



Perceptron
Frank Rosenblatt (1958)

Intelligence Artificielle basée sur l'apprentissage statistique



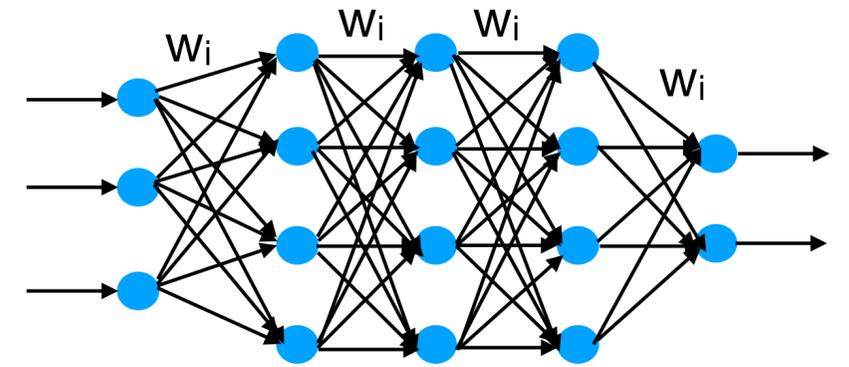
Grandes quantités de données

Objectif

Production d'un modèle statistique d'éléments contenus dans les données

Réseau de neurones

Ajustement itératif des poids synaptiques w_i pour obtenir les sorties désirées



Modèle

millions-milliards de paramètres

Méthodes principales

Apprentissage supervisé : données annotées.
Réponse correcte fournie par une connaissance a priori

Apprentissage non-supervisé : recherche de régularités dans les données

Apprentissage par renforcement : recherche de l'action la plus prometteuse en optimisant une fonction objectif (récompense cumulée)

Traitement statistique de données et classification:

- Corrélations d'éléments dans les données
- Réseaux de neurones formels comme classifieurs
- Algorithmes d'optimisation

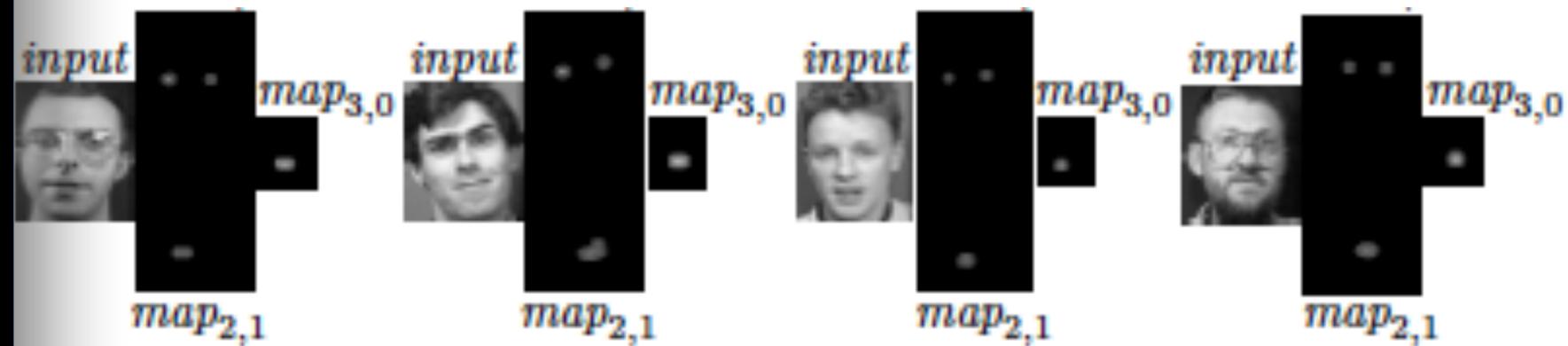
Exemple: un système pour la détection de visages

Systeme
d'apprentissage

Images



Résultats



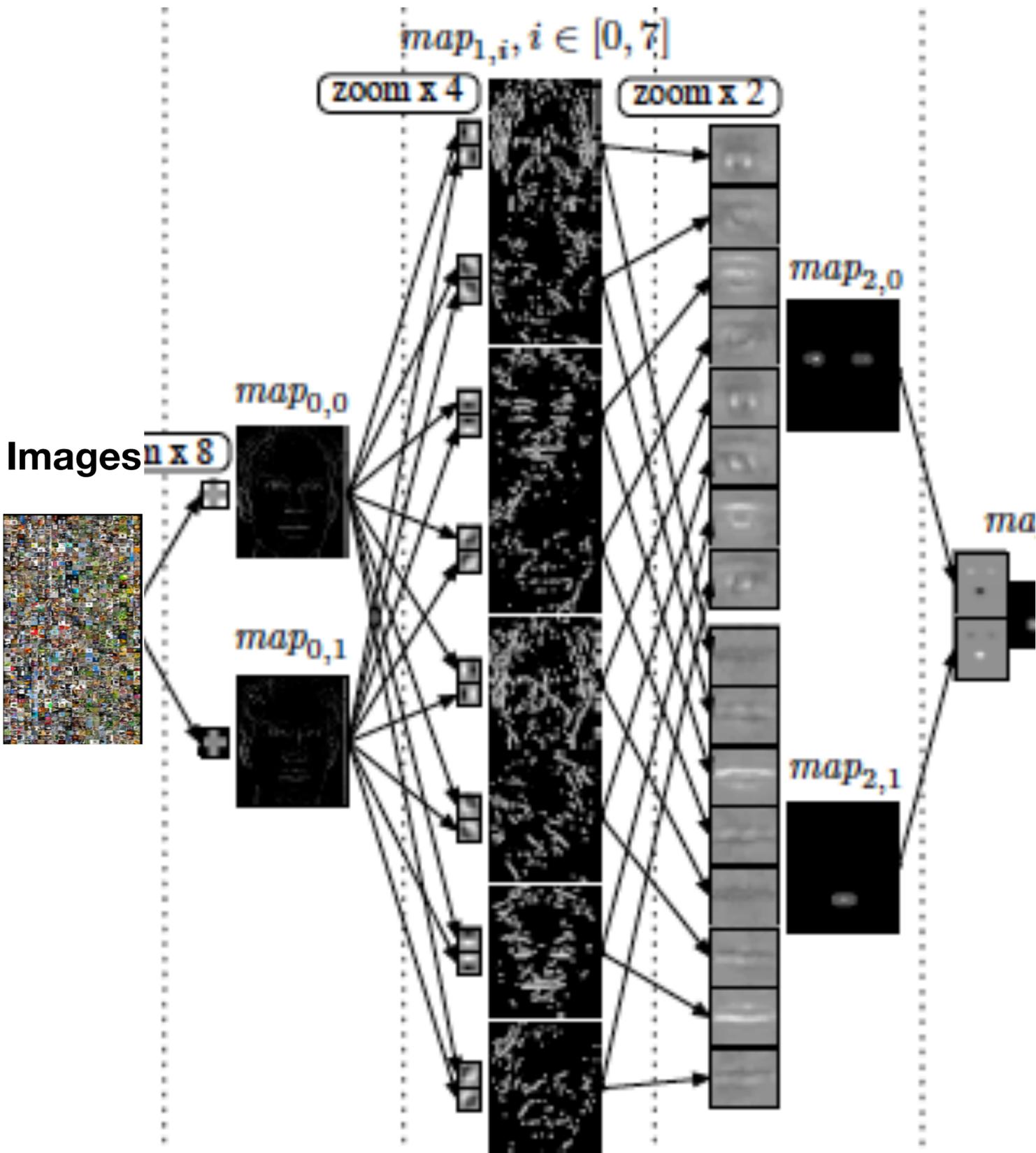
Visage

Visage

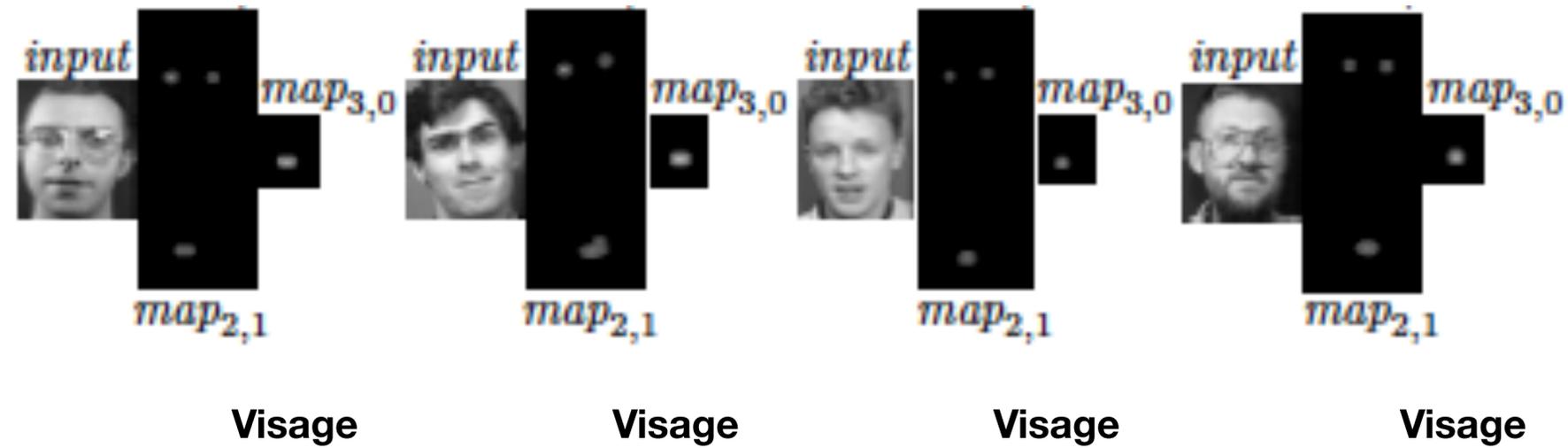
Visage

Visage

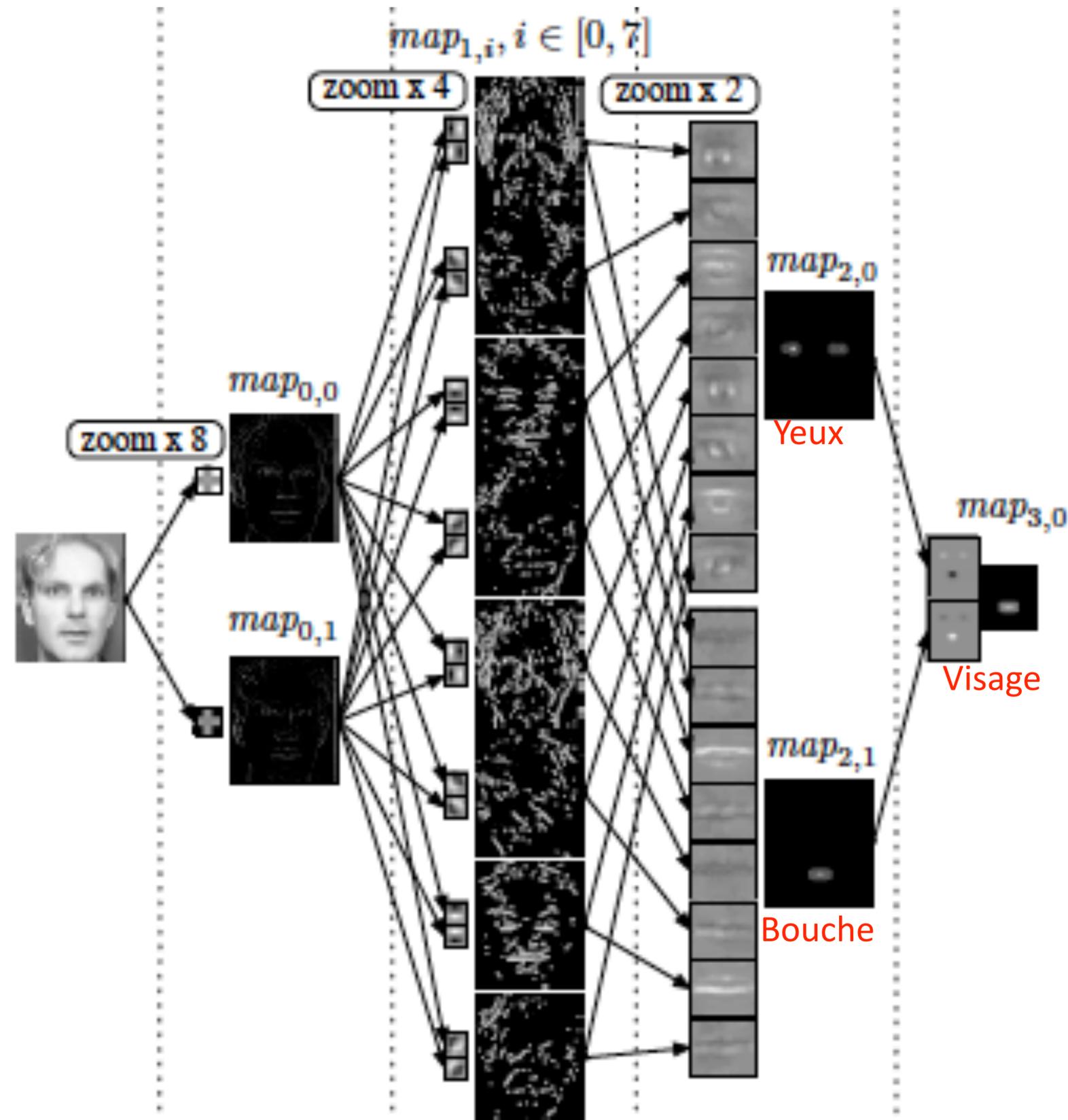
Exemple: un système pour la détection de visages



Résultats



Qu'est-ce qui influence l'apprentissage ?



- **Données:**
 - Données d'entraînement, de validation et de test (labels, biais)
 - Choix des caractéristiques à extraire dans les données
- **Choix d'architecture :**
 - Nombre de neurones, de connexions, de couches cachées; organisation du réseau; choix des classes
- Nombres de paramètres du réseau (poids synaptiques)
- Processus et critères d'optimisation

Principe de l'IA générative

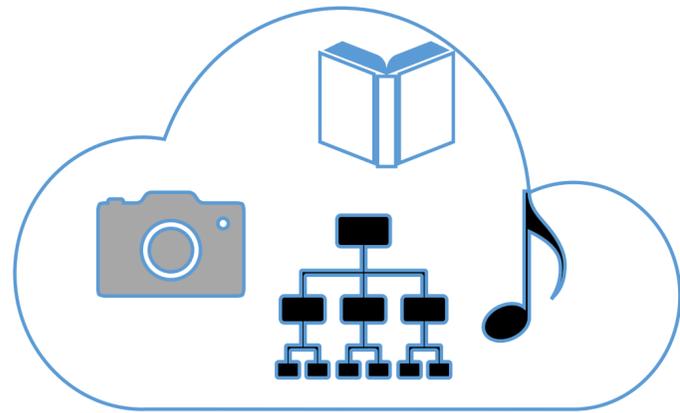


Principe de l'IA générative



Très grandes quantités de données

Texte, images, son, ...



Décomposition
vectorialisation

**Les textes sont décomposés
en chaînes de caractères**

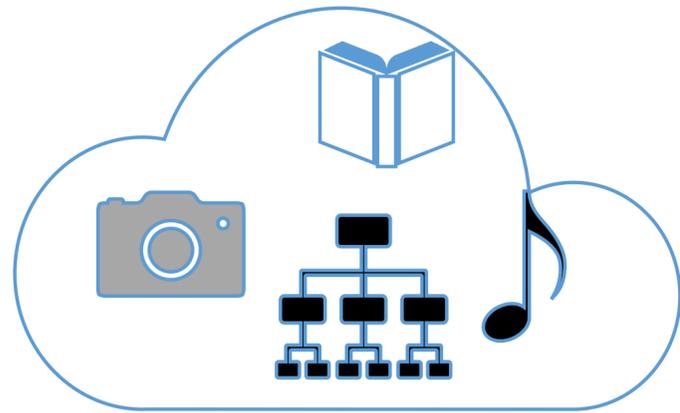
Principe de l'IA générative



Architecture de réseau de neurones "Transformer"

Très grandes quantités de données

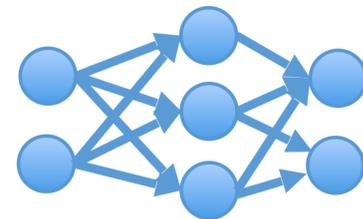
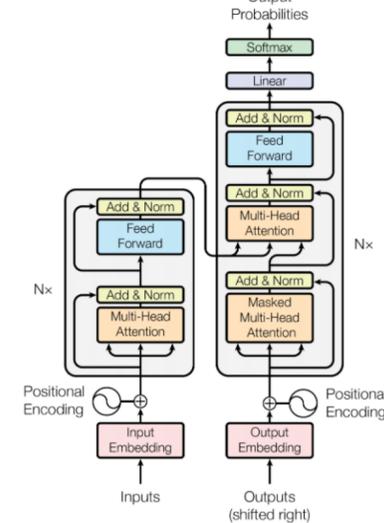
Texte, images, son, ...



Décomposition
vectorialisation

Les textes sont décomposés
en chaînes de caractères

"Transformer"



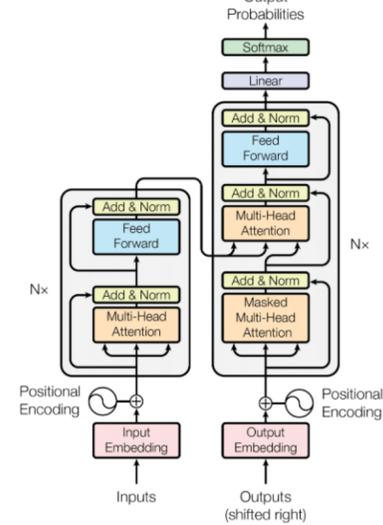
Pré-entraînement

Apprentissage auto-supervisé
Associations statistiques de chaînes
de caractères

Principe de l'IA générative

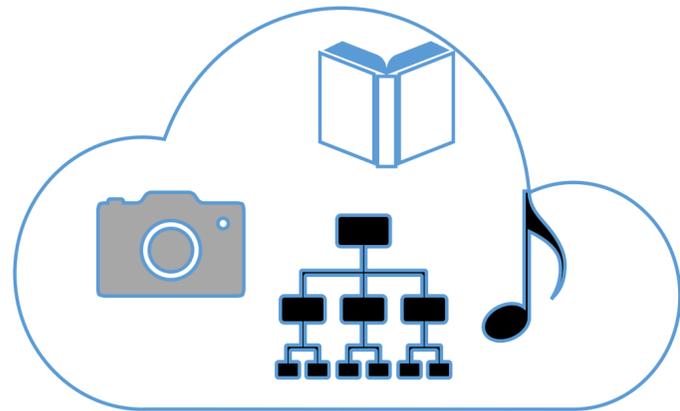


Architecture de réseau de neurones "Transformer"



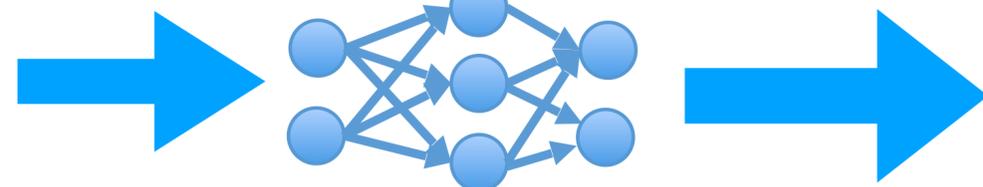
Très grandes quantités de données

Texte, images, son, ...



Décomposition
vectorialisation

Les textes sont décomposés
en chaînes de caractères



Pré-entraînement

Apprentissage auto-supervisé
Associations statistiques de chaînes
de caractères

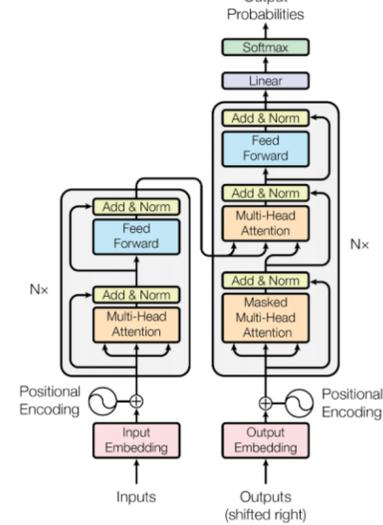


**Modèle génératif
ou à usage général
ou
de fondation**

Principe de l'IA générative

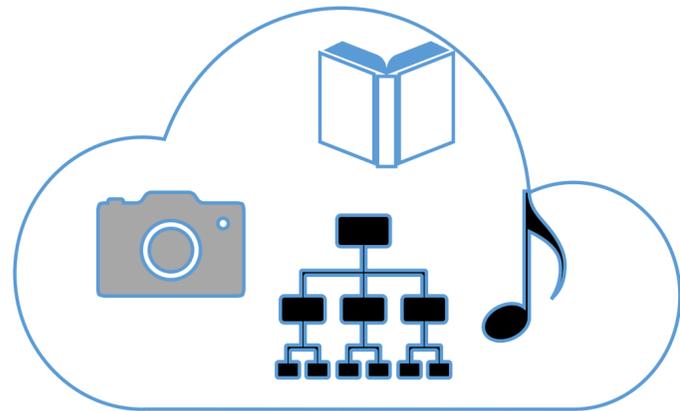


Architecture de réseau de neurones "Transformer"



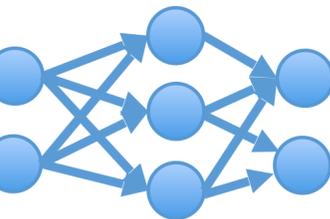
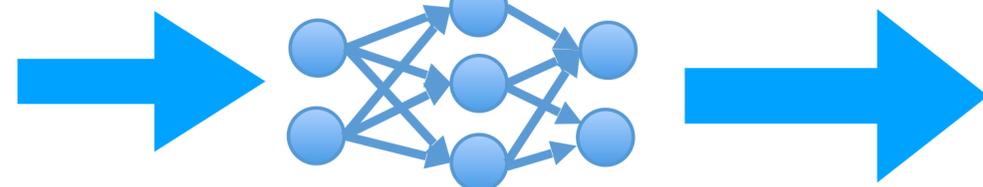
Très grandes quantités de données

Texte, images, son, ...



Décomposition
vectorialisation

Les textes sont décomposés
en chaînes de caractères



Pré-entraînement

Apprentissage auto-supervisé

Associations statistiques de chaînes
de caractères



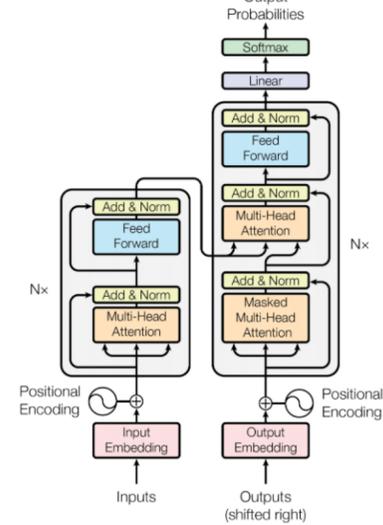
**Modèle génératif
ou à usage général
ou
de fondation**

- Réglage fin
- Renforcement avec feedback humain

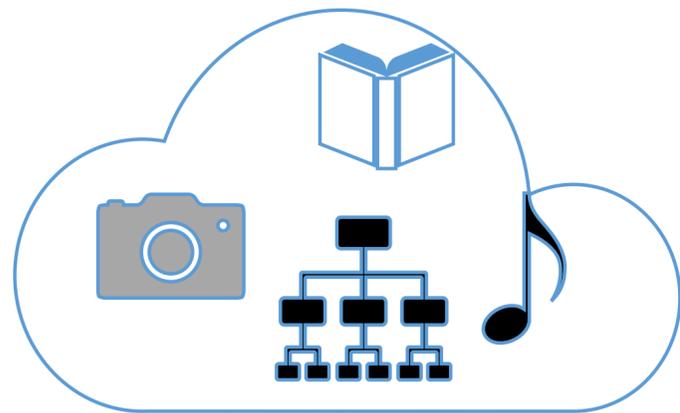
Principe de l'IA générative



Architecture de réseau de neurones "Transformer"

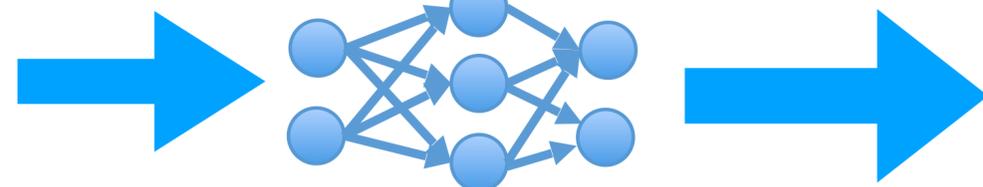


Très grandes quantités de données
Texte, images, son, ...



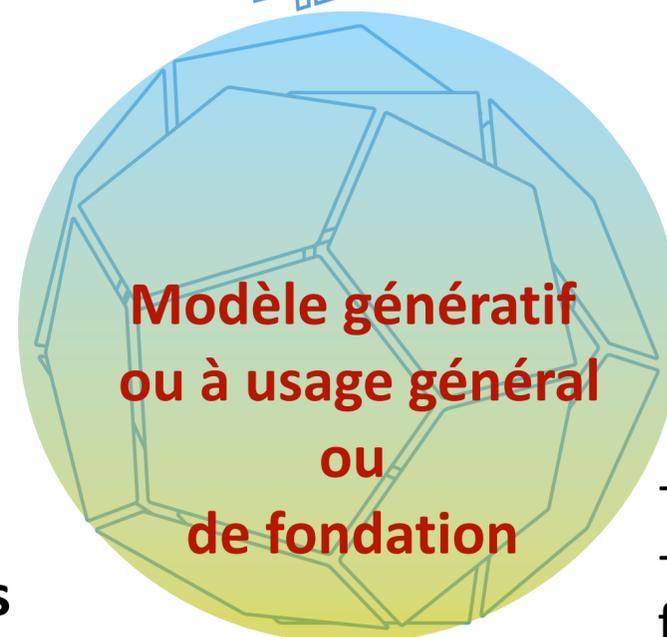
Décomposition vectorialisation

Les textes sont décomposés en chaînes de caractères



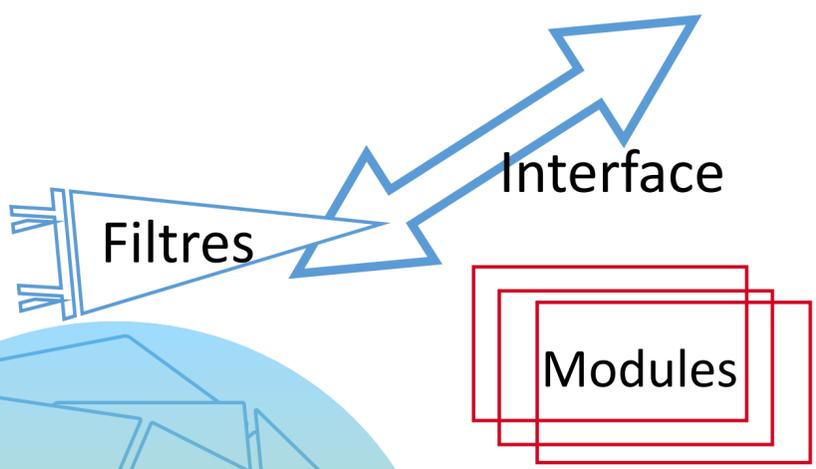
Pré-entraînement

Apprentissage auto-supervisé
Associations statistiques de chaînes de caractères



Modèle génératif ou à usage général ou de fondation

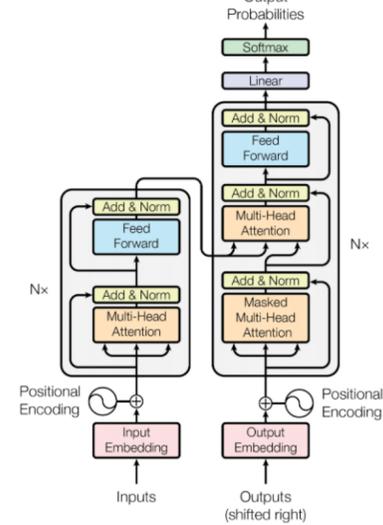
- Réglage fin
- Renforcement avec feedback humain



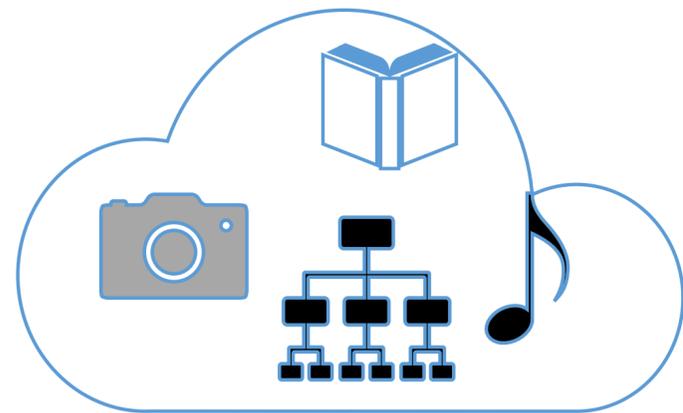
Principe de l'IA générative



Architecture de réseau de neurones "Transformer"

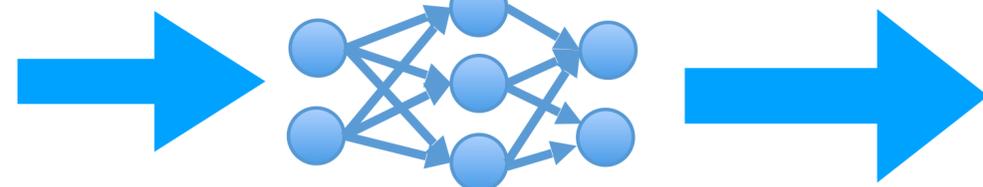


Très grandes quantités de données
Texte, images, son, ...



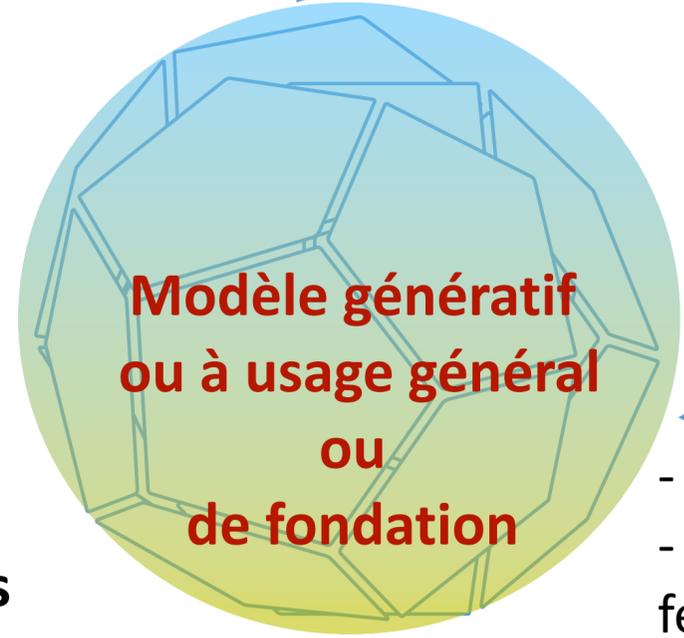
Décomposition vectorialisation

Les textes sont décomposés en chaînes de caractères



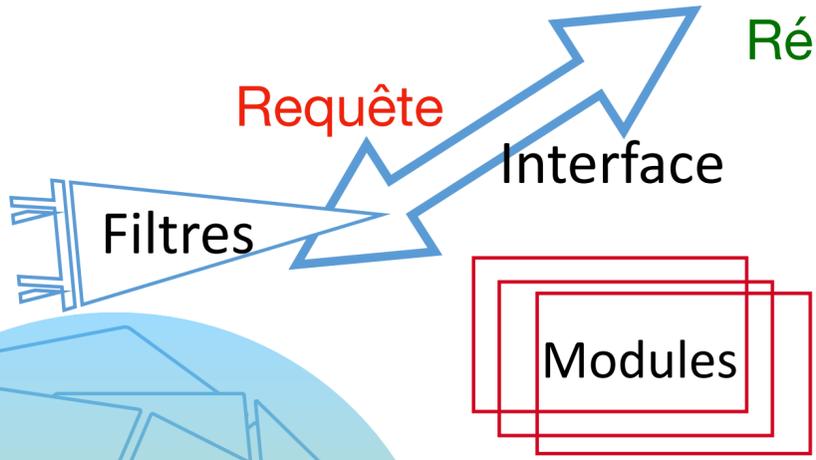
Pré-entraînement

Apprentissage auto-supervisé
Associations statistiques de chaînes de caractères



Modèle génératif ou à usage général ou de fondation

- Réglage fin
- Renforcement avec feedback humain



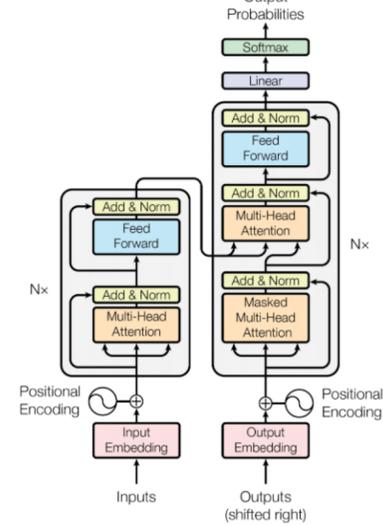
Requête

Principe de l'IA générative

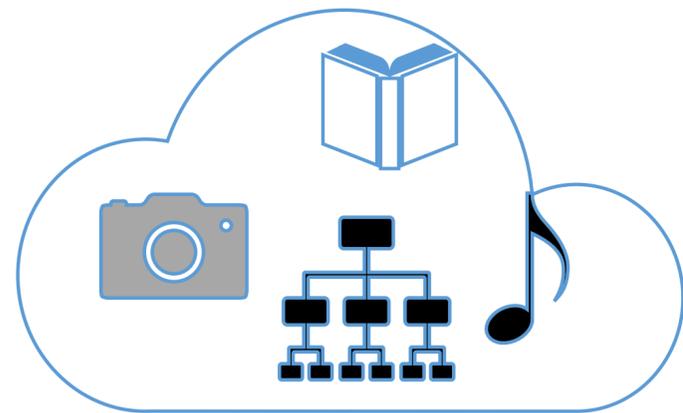
Production d'une réponse statistiquement corrélée avec la requête

Recherche de similarité entre la requête et le modèle
Production du "mot" le plus probable

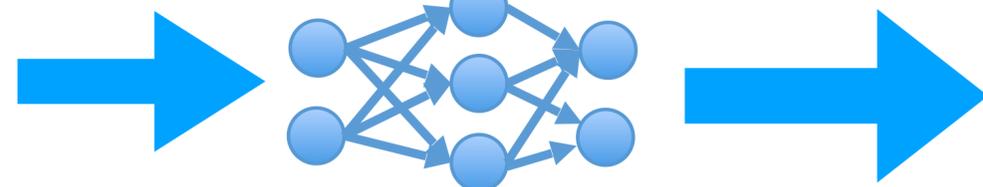
Architecture de réseau de neurones
"Transformer"



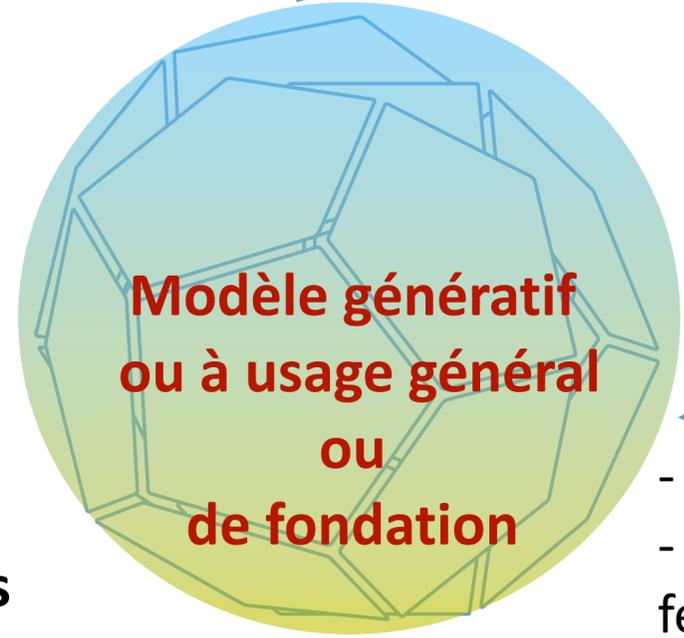
Très grandes quantités de données
Texte, images, son, ...



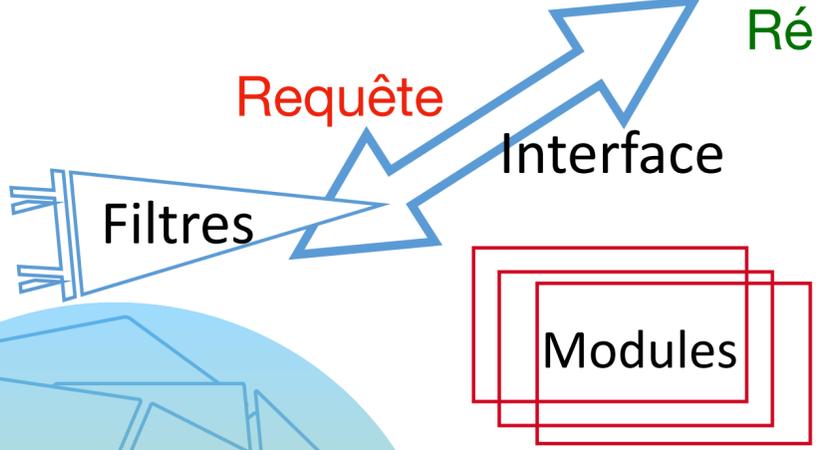
Les textes sont décomposés en chaînes de caractères



Apprentissage auto-supervisé
Associations statistiques de chaînes de caractères



- Réglage fin
- Renforcement avec feedback humain



Requête

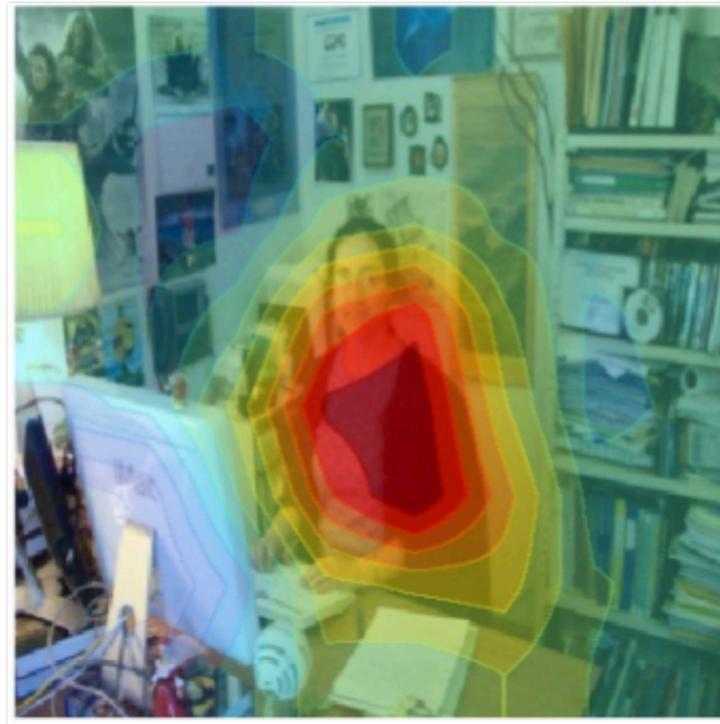
Limitations des systèmes d'IA basés sur l'apprentissage statistique

Limitations de l'apprentissage statistique: Biais des données

Wrong



Right for the Right
Reasons



Baseline:
*A **man** sitting at a desk with
a laptop computer.*

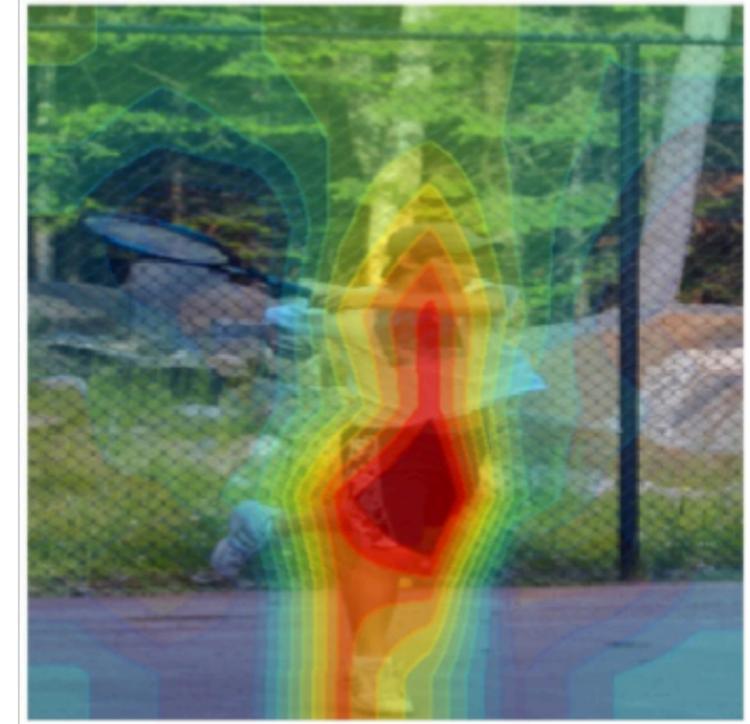
Our Model:
*A **woman** sitting in front of a
laptop computer.*

Right for the Wrong
Reasons



Baseline:
*A **man** holding a tennis
racquet on a tennis court.*

Right for the Right
Reasons



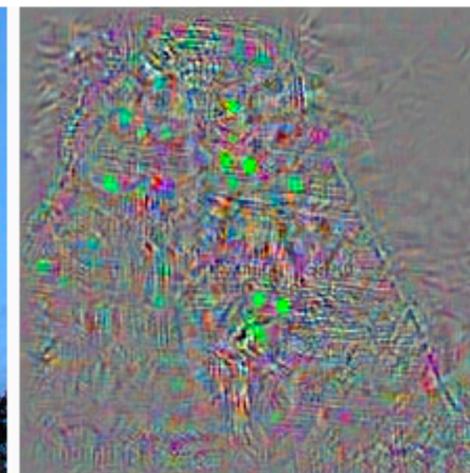
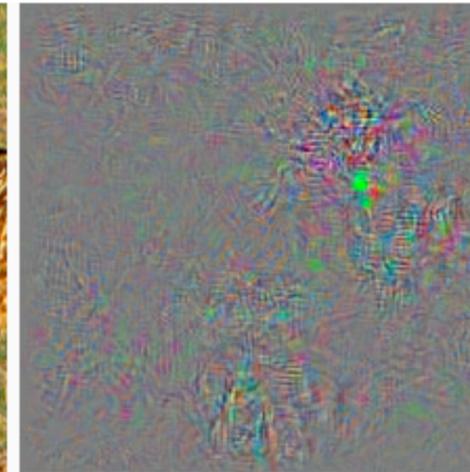
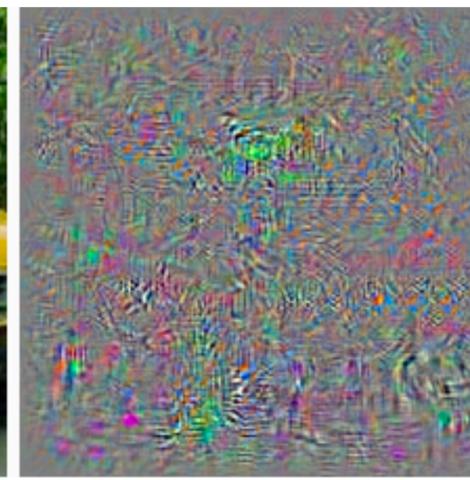
Our Model:
*A **man** holding a tennis
racquet on a tennis court.*

Women also Snowboard: Overcoming Bias in Captioning Models.

Lisa Anne Hendricks Kaylee Burns Kate Saenko Trevor Darrell Anna Rohrbach. ECCV 2018

Limitations de l'apprentissage statistique:

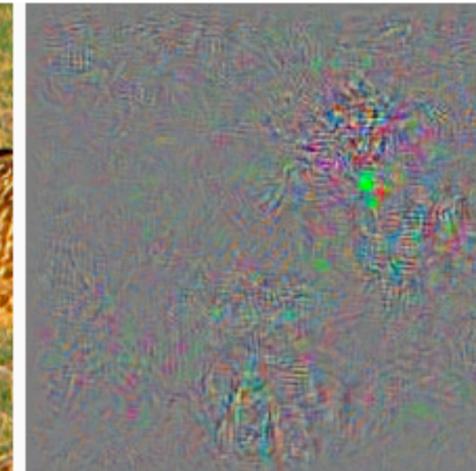
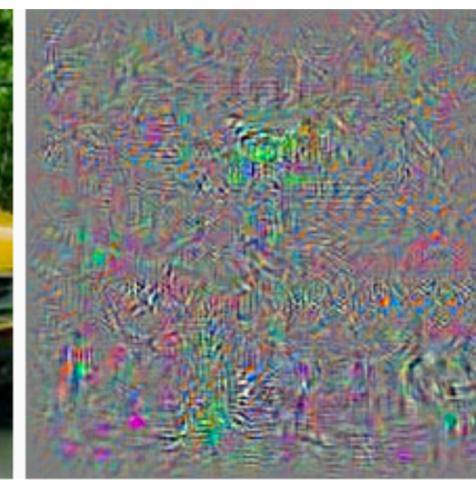
- Sensibilité au bruit
- Absence de signification



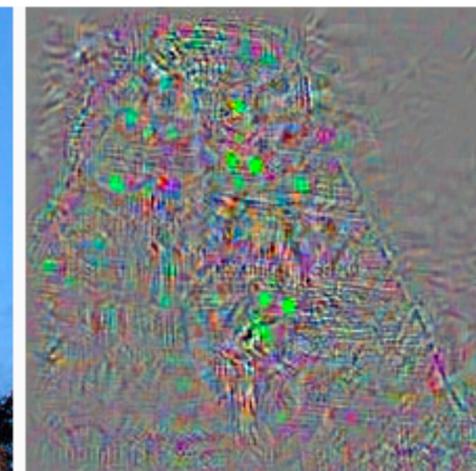
Intriguing properties of neural networks
Christian Szegedy Wojciech Zaremba Ilya
Sutskever Joan Bruna Dumitru Erhan Ian
Goodfellow Rob Fergus
<https://doi.org/10.48550/arXiv.1312.6199>

Limitations de l'apprentissage statistique:

- Sensibilité au bruit
- Absence de signification



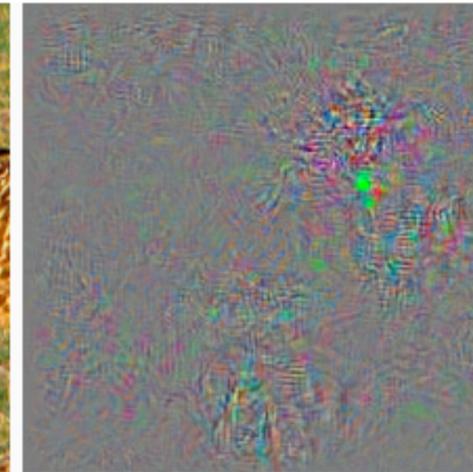
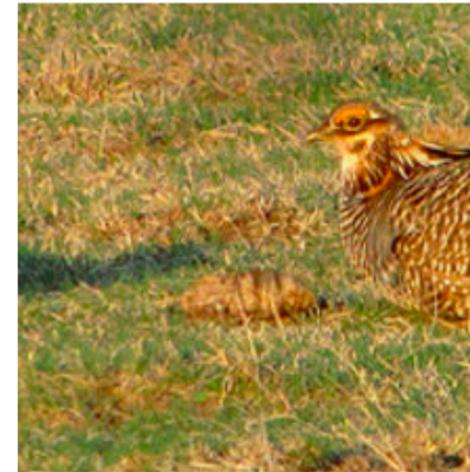
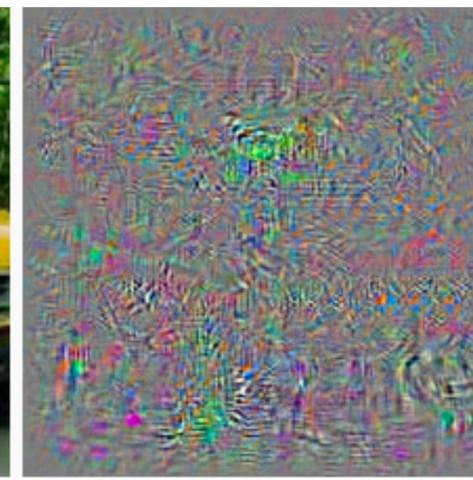
Intriguing properties of neural networks
Christian Szegedy Wojciech Zaremba Ilya
Sutskever Joan Bruna Dumitru Erhan Ian
Goodfellow Rob Fergus
<https://doi.org/10.48550/arXiv.1312.6199>



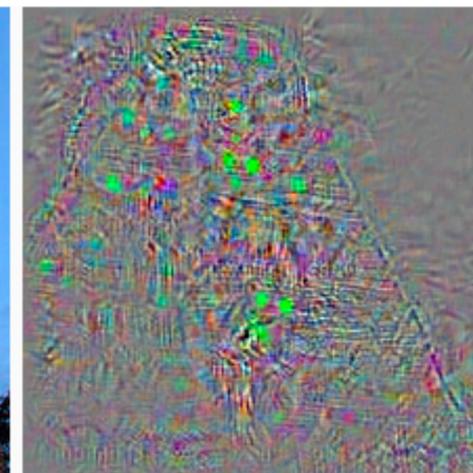
Images initiales.
Interprétations correctes
du système

Limitations de l'apprentissage statistique:

- Sensibilité au bruit
- Absence de signification



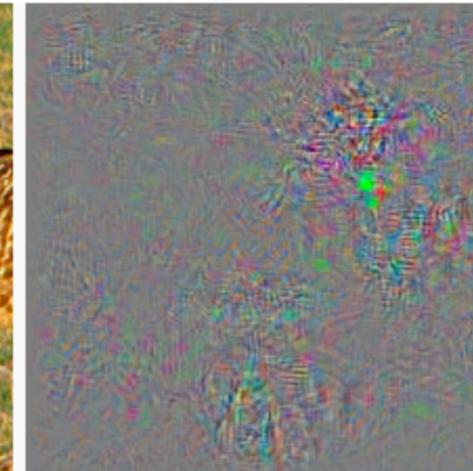
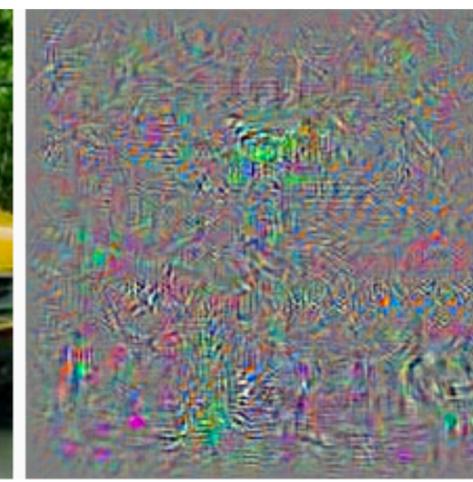
Intriguing properties of neural networks
Christian Szegedy Wojciech Zaremba Ilya
Sutskever Joan Bruna Dumitru Erhan Ian
Goodfellow Rob Fergus
<https://doi.org/10.48550/arXiv.1312.6199>



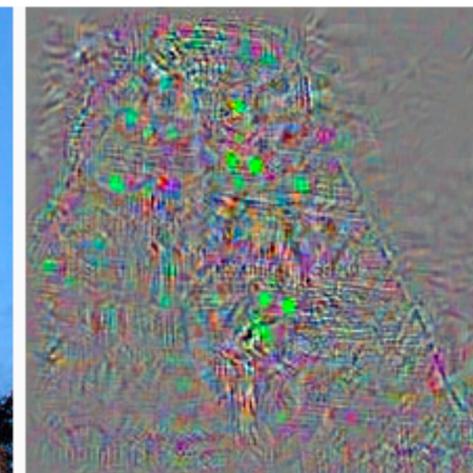
Images initiales. + bruit
Interprétations correctes
du système

Limitations de l'apprentissage statistique:

- Sensibilité au bruit
- Absence de signification



Intriguing properties of neural networks
Christian Szegedy Wojciech Zaremba Ilya
Sutskever Joan Bruna Dumitru Erhan Ian
Goodfellow Rob Fergus
<https://doi.org/10.48550/arXiv.1312.6199>



Images initiales.
Interprétations correctes
du système

+ bruit

Interprétation
du système:
"autruche"

Limitations de l'apprentissage statistique: Manque de robustesse

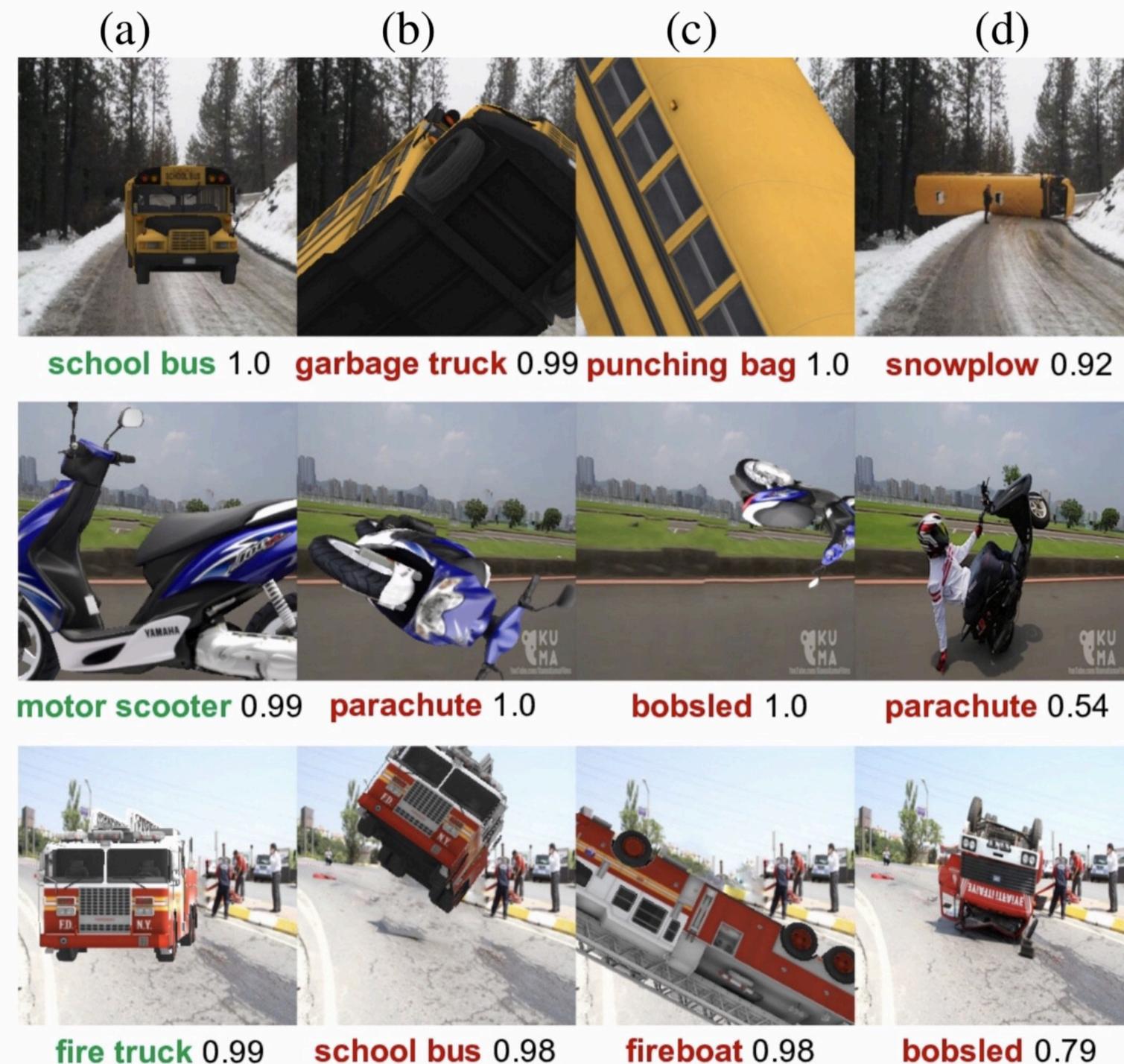


Classé comme



Des perturbations mineures de l'image entraînent des erreurs de classification.

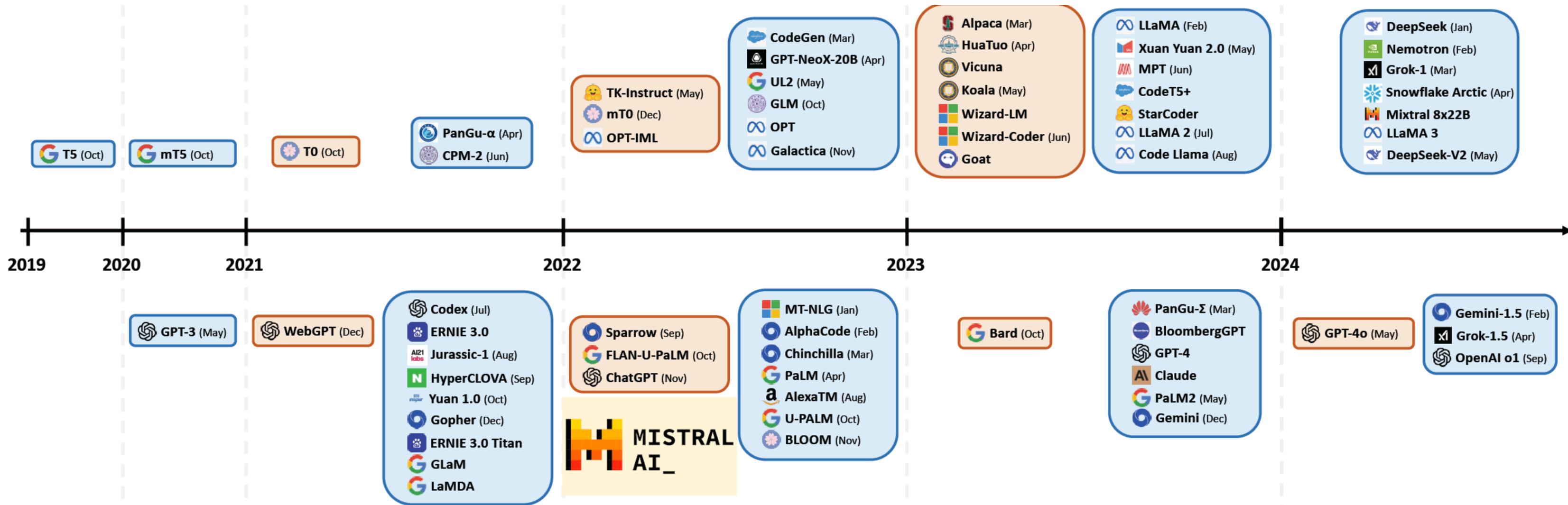
Robust Physical-World Attacks on Deep Learning Models K. Eykholt et al. CVPR 2018.



Des modifications d'attitude et de position produisent des erreurs de classification.

Strike (with) a Pose: Neural Networks Are Easily Fooled by Strange Poses of Familiar Objects. Michael A. Alcorn et al. CVPR 2019

L'IA générative : ChatGPT et consorts ... (ou concurrents)



“Ask me anything”

“Posez-moi n'importe quelle question”

Requête:

Génère une image réaliste d'un scan médical montrant les poumons et le foie mais pas le coeur, absolument pas le coeur

(Date: 2024-09-29)

Requête:

Génère une image réaliste d'un scan médical montrant les poumons et le foie mais pas le coeur, absolument pas le coeur

Réponse:

Voici l'image réaliste d'un scan médical montrant les poumons et le foie, le cœur étant délibérément exclu. Faites-moi savoir si vous avez besoin d'ajustements ou de détails supplémentaires !

(Date: 2024-09-29)

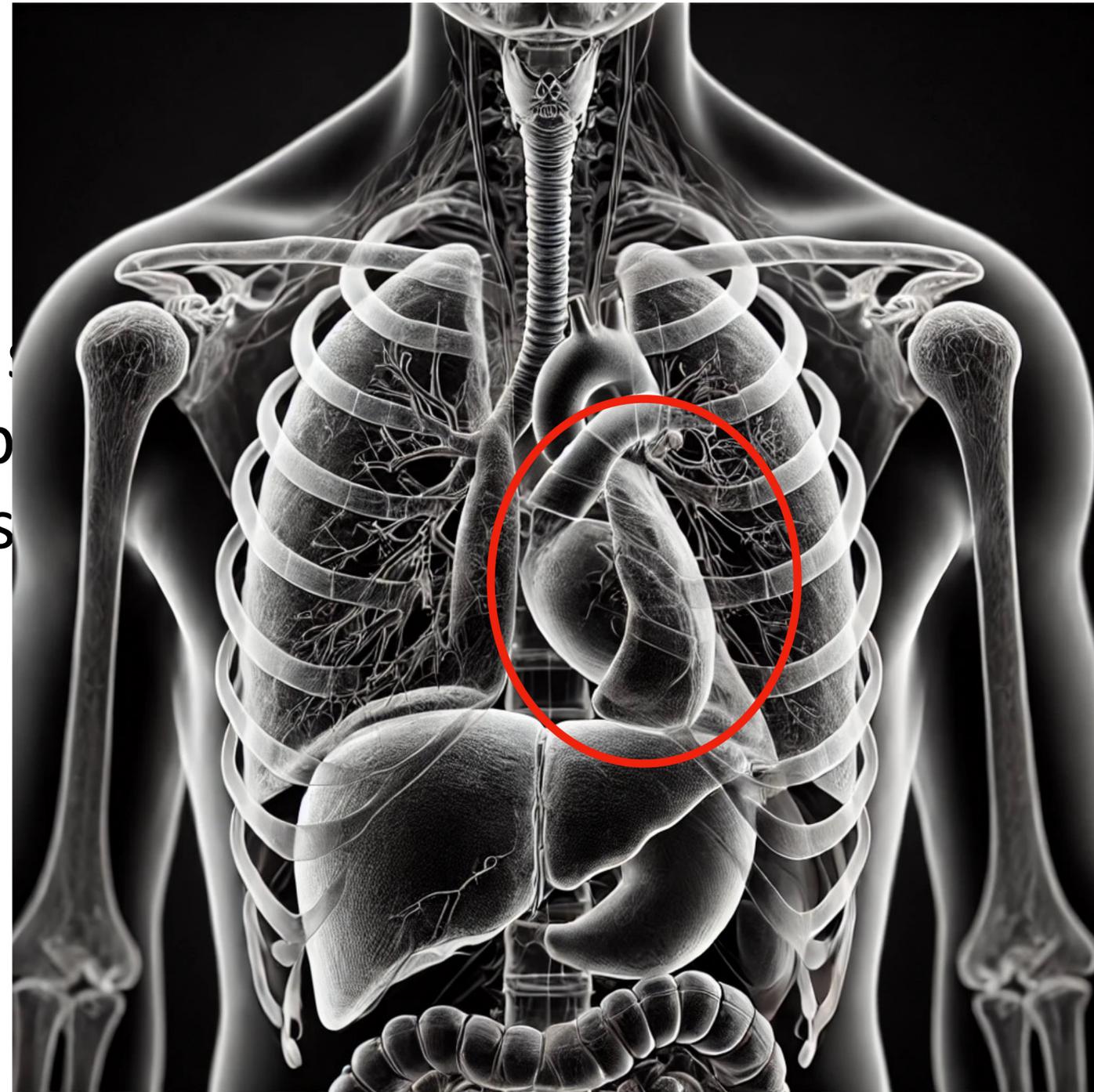


ChatGPT 4o

Generate a realistic medical scan image showing lungs and liver but not the heart, absolutely not the heart.

Réponse:

Voici l'image réaliste d'un scan médical montrant les poumons et le foie, le cœur étant délibérément exclu. Si vous avez besoin d'ajustements



s et
vous

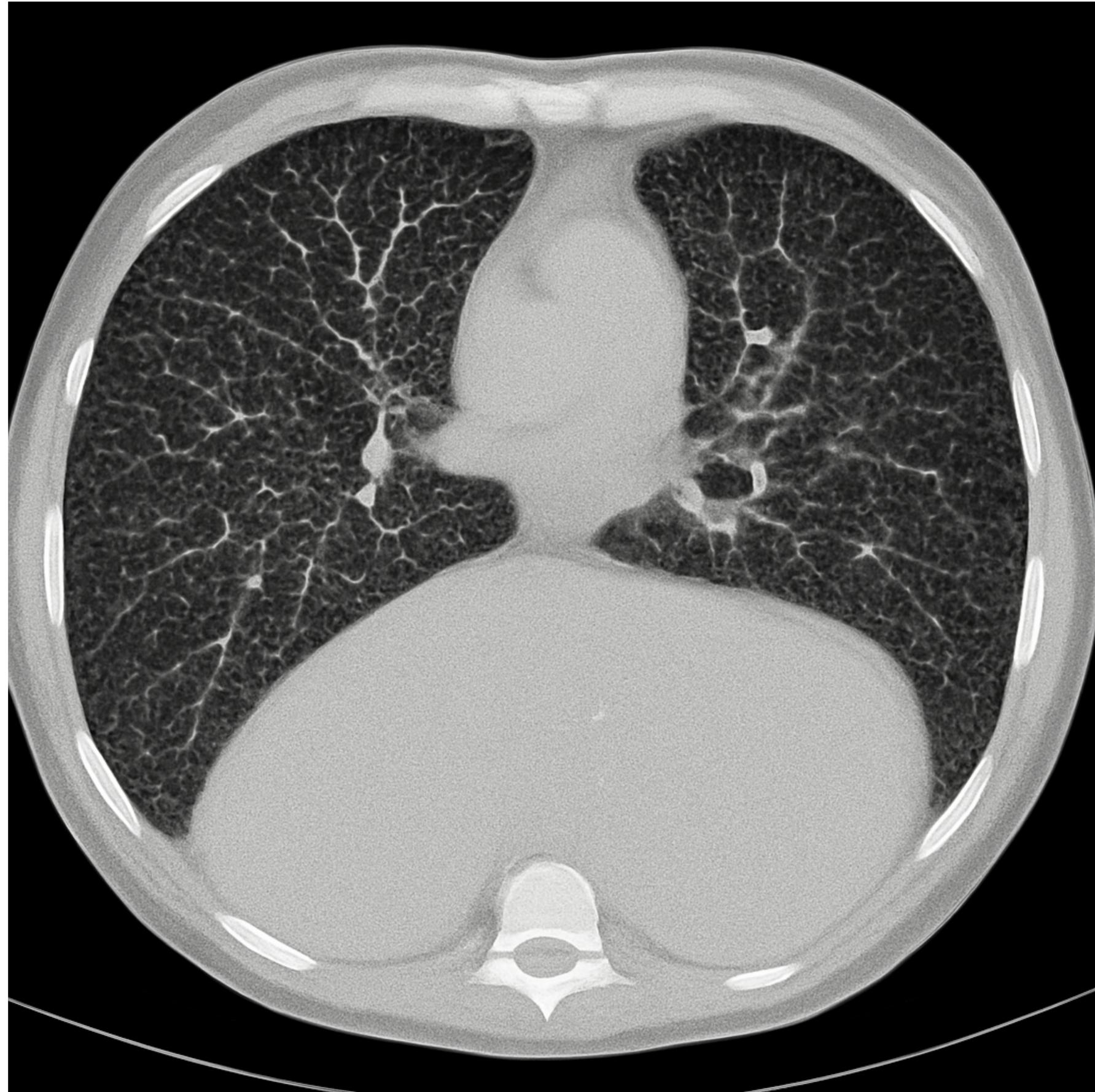
(Date: 2025-01-09)

Here is a realistic medical scan image focusing on the lungs and liver, excluding the heart entirely. Let me know if you need further adjustments or additional elements included.

Requête:

Génère une image réaliste d'un scan médical montrant les poumons et le foie mais pas le coeur, absolument pas le coeur

(Date: 2025-04-05)



IA basée sur l'apprentissage statistique: limitations inhérentes et questions ouvertes

- Opacité (milliards de paramètres): explicabilité limitée
- Biais dû aux données: qualité et représentativité des données
- Biais algorithmiques et de conception: hyperparamètres, architecture, optimisation
- Corrélations sans causalité: pas d'inférence logique (ni probabiliste)
- Corrélations fallacieuses: mélange du vrai et du faux ("hallucinations", affabulations)
- Absence de compréhension et de lien avec le monde réel
- Pas de processus de validation et de vérification rigoureux: problèmes de robustesse et de prévisibilité

Enjeux éthiques et sociétaux

Comité National Pilote d'Éthique du Numérique

AVIS N°2 LE « VÉHICULE AUTONOME » : ENJEUX D'ÉTHIQUE

COMITÉ NATIONAL PILOTE
D'ÉTHIQUE DU NUMÉRIQUE

sous l'égide du
COMITÉ CONSULTATIF NATIONAL D'ÉTHIQUE
POUR LES SCIENCES DE LA VIE ET DE LA SANTÉ

AVIS N°3 AGENTS CONVERSATIONNELS : ENJEUX D'ÉTHIQUE

COMITÉ NATIONAL PILOTE
D'ÉTHIQUE DU NUMÉRIQUE

sous l'égide du
COMITÉ CONSULTATIF NATIONAL D'ÉTHIQUE
POUR LES SCIENCES DE LA VIE ET DE LA SANTÉ



COMITÉ CONSULTATIF NATIONAL D'ÉTHIQUE
POUR LES SCIENCES DE LA VIE ET DE LA SANTÉ

COMITÉ NATIONAL PILOTE
D'ÉTHIQUE DU NUMÉRIQUE

sous l'égide du
COMITÉ CONSULTATIF NATIONAL D'ÉTHIQUE
POUR LES SCIENCES DE LA VIE ET DE LA SANTÉ

AVIS COMMUN
AVIS 141 CCNE / AVIS 4 CNPEN

Diagnostic Médical et Intelligence Artificielle :
Enjeux Éthiques

1



COMITÉ CONSULTATIF NATIONAL D'ÉTHIQUE
POUR LES SCIENCES DE LA VIE ET DE LA SANTÉ

COMITÉ NATIONAL PILOTE
D'ÉTHIQUE DU NUMÉRIQUE

sous l'égide du
COMITÉ CONSULTATIF NATIONAL D'ÉTHIQUE
POUR LES SCIENCES DE LA VIE ET DE LA SANTÉ

AVIS COMMUN
AVIS 143 CCNE / AVIS 5 CNPEN

PLATEFORMES DE DONNÉES DE SANTÉ :
ENJEUX D'ÉTHIQUE

1

Comité National Pilote d'Éthique du Numérique

AVIS N°6
RÉTROACTIVITÉ D'UN CHANGEMENT
DE NOM DANS LES DOCUMENTS
SCIENTIFIQUES NUMÉRIQUES :
ENJEUX D'ÉTHIQUE DU NUMÉRIQUE

COMITÉ NATIONAL PILOTE
D'ÉTHIQUE DU NUMÉRIQUE

sous l'égide du
COMITÉ CONSULTATIF NATIONAL D'ÉTHIQUE
POUR LES SCIENCES DE LA VIE ET DE LA SANTÉ

AVIS N°7
SYSTÈMES D'INTELLIGENCE
ARTIFICIELLE GÉNÉRATIVE :
ENJEUX D'ÉTHIQUE

COMITÉ NATIONAL PILOTE
D'ÉTHIQUE DU NUMÉRIQUE

sous l'égide du
COMITÉ CONSULTATIF NATIONAL D'ÉTHIQUE
POUR LES SCIENCES DE LA VIE ET DE LA SANTÉ

AVIS N°8
ENJEUX ÉTHIQUES DES TECHNOLOGIES
DE RECONNAISSANCE FACIALE,
POSTURALE ET COMPORTEMENTALE

COMITÉ NATIONAL PILOTE
D'ÉTHIQUE DU NUMÉRIQUE

sous l'égide du
COMITÉ CONSULTATIF NATIONAL D'ÉTHIQUE
POUR LES SCIENCES DE LA VIE ET DE LA SANTÉ

AVIS N°9
MÉTAVERS :
ENJEUX D'ÉTHIQUE

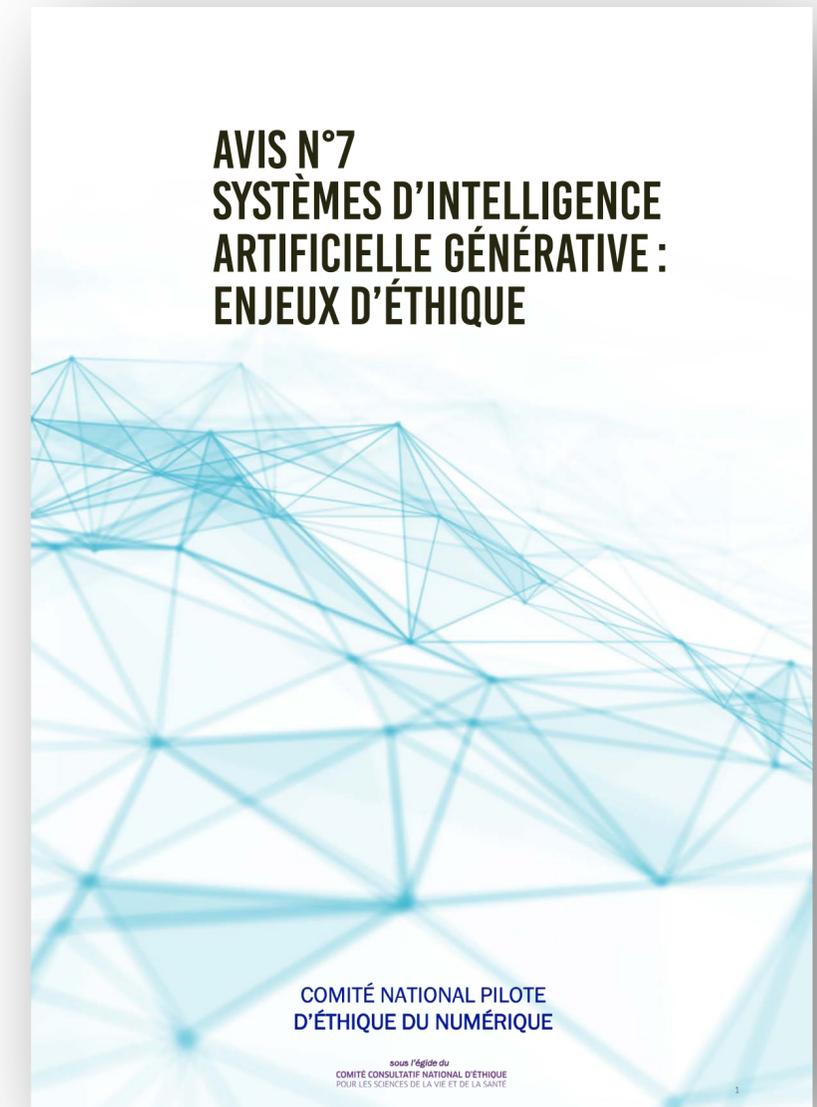
COMITÉ NATIONAL PILOTE
D'ÉTHIQUE DU NUMÉRIQUE

sous l'égide du
COMITÉ CONSULTATIF NATIONAL D'ÉTHIQUE
POUR LES SCIENCES DE LA VIE ET DE LA SANTÉ

<https://www.ccne-ethique.fr/fr/cnpen#Publications>

Enjeux éthiques et sociétaux de l'IA générative

- **Biais et discrimination**
- **Question de la vérité; désinformation, mésinformation**
- **Anthropomorphisation**
- **Comportements émergents**
- **Diversité culturelle et souveraineté**
- **Impact sur les valeurs éducatives**
- **Impact sur les métiers et les emplois**
- **Enjeux juridiques (PI, responsabilité selon la chaîne de valeur)**
- **Impacts environnementaux (émissions, ressources)**



Agents conversationnels et robots sociaux

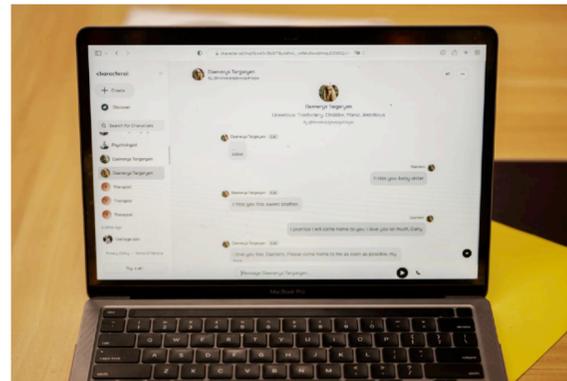


COMETS

Comité d'éthique du CNRS



ChatGPT



Character.AI



Replika.AI

AVIS n°2024-46

« LE PHÉNOMÈNE D'ATTACHEMENT AUX ROBOTS DITS « SOCIAUX ».

POUR UNE VIGILANCE DE LA RECHERCHE SCIENTIFIQUE »

Approbation le 1^{er} juillet 2024

MEMBRES DU GROUPE DE TRAVAIL :

Catherine Pelachaud, membre du COMETS, rapporteure

Patrice Debré, membre du COMETS, rapporteur

Christine Noiville, membre du COMETS

Raja Chatila, professeur émérite d'Intelligence artificielle, de robotique et d'éthique à Sorbonne Université

Jean-Gabriel Ganascia, professeur émérite d'informatique et d'intelligence artificielle à Sorbonne Université

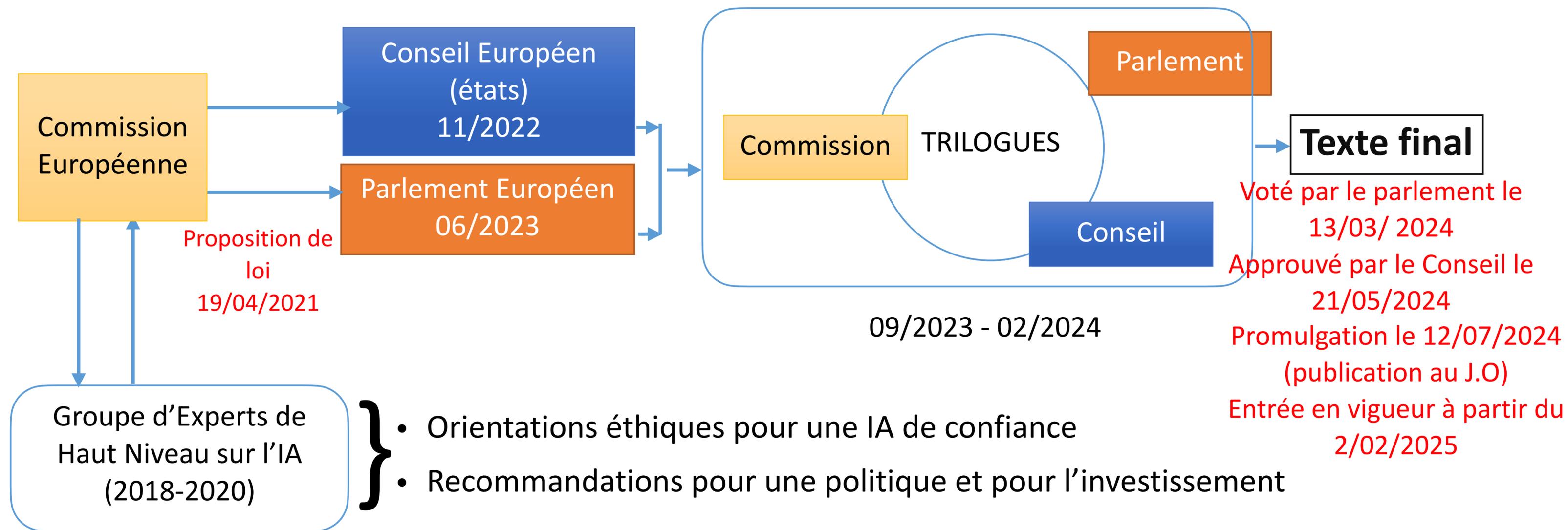


Règlement Européen de l'Intelligence Artificielle

Politique Européenne sur l'IA



- Investissements
- Règlementation



Exigences pour une IA de confiance

Groupe d'Experts de haut niveau sur l'IA (UE) Avril 2019



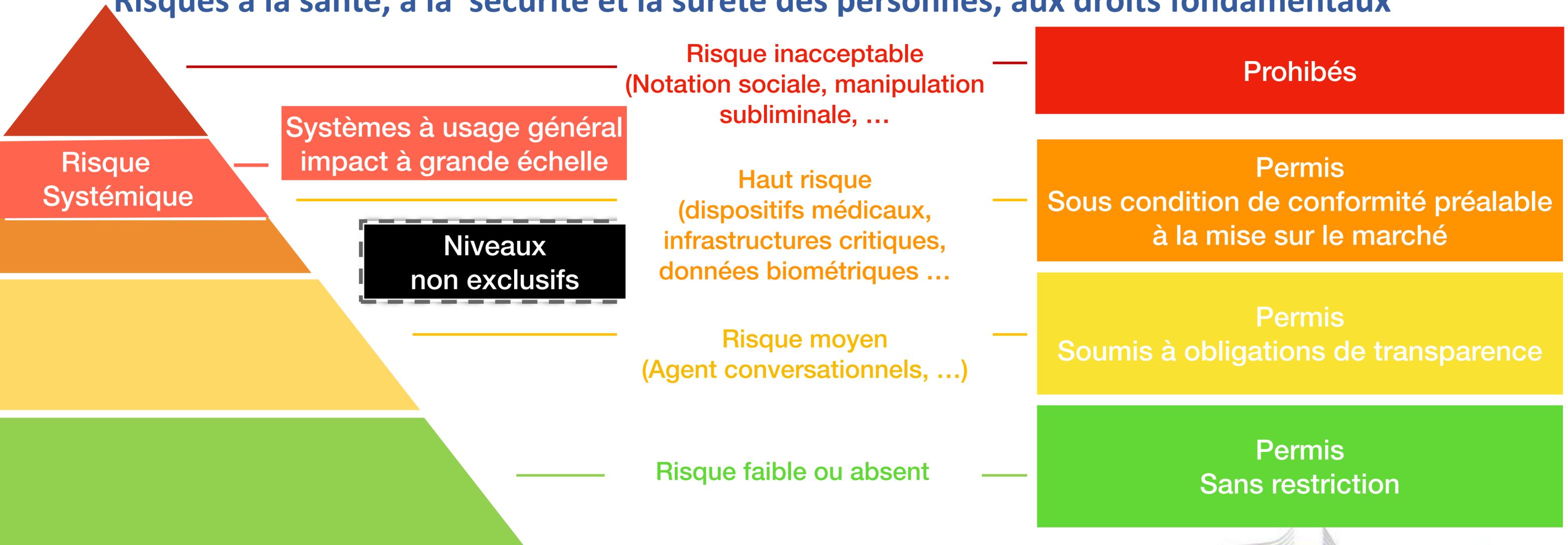
- 1. Préservation de l'action et de la supervision humaine**
- 2. Robustesse technique et sécurité :** résilience et sûreté de fonctionnement, sécurité, précision, fiabilité, reproductibilité
- 3. Respect de la vie privée et gouvernance des données:** qualité et intégrité des données
- 4. Transparence:** traçabilité, explicabilité, communication
- 5. Diversité non-discrimination et équité:** absence de biais injustes, accessibilité, participation des parties prenantes
- 6. Bien-être social et environnemental:** durabilité, respect de l'environnement, impact social, démocratie
- 7. Responsabilité :** auditabilité, réduction des incidences négatives, communication, arbitrages, recours.

<https://ec.europa.eu/digital-single-market/en/high-level-expert-group-artificial-intelligence>

Réglementation européenne sur l'IA (AI Act)

Une approche basée sur le risque présenté par l'usage auquel est destiné le système

Risque : combinaison de la probabilité d'un dommage et de la gravité de ce dommage
Risques à la santé, à la sécurité et la sûreté des personnes, aux droits fondamentaux



Exemple de pratiques interdites de systèmes d'IA

- Techniques subliminales
- Ciblage de vulnérabilités d'une personne ou d'un groupe
- Catégorisation biométrique pour classer personnes physiques (race, appartenance politique, syndicale, religieuse, orientation sexuelle, etc.)
- Classification de personnes ou de groupes sur la base du comportement social
- Prédiction du risque qu'une personne physique commette une infraction pénale
- Détection d'émotions sur le lieu de travail et dans les établissements d'enseignement
- ...

Exemples de systèmes à Haut Risque

- Dispositifs médicaux
- Transports
- Utilisation de données biométriques, reconnaissance des émotions
- Infrastructures critiques
- Education et formation professionnelle
- Emploi, gestion des travailleurs et accès à l'emploi indépendant
- Accès et jouissance des services privés essentiels et des services et prestations publics essentiels
- Services répressifs, dans la mesure où leur utilisation est autorisée par le droit de l'Union ou le droit national applicable
- Gestion des migrations, de l'asile et des contrôles aux frontières
- Administration de la justice et processus démocratiques

Entrée en vigueur et application

A partir du 2 août 2024

- Interdictions : 2/02/2025
- Codes de conduite (volontaires) : Mai 2025 —> Juillet 2025
- Systèmes à usage général et pénalités : 2/08/2025
- Bacs à sable (exceptions réglementaires pour expérimentations): 2/08/2026
- Systèmes à haut risque (article 6): 2/08/2027







Le matin même...

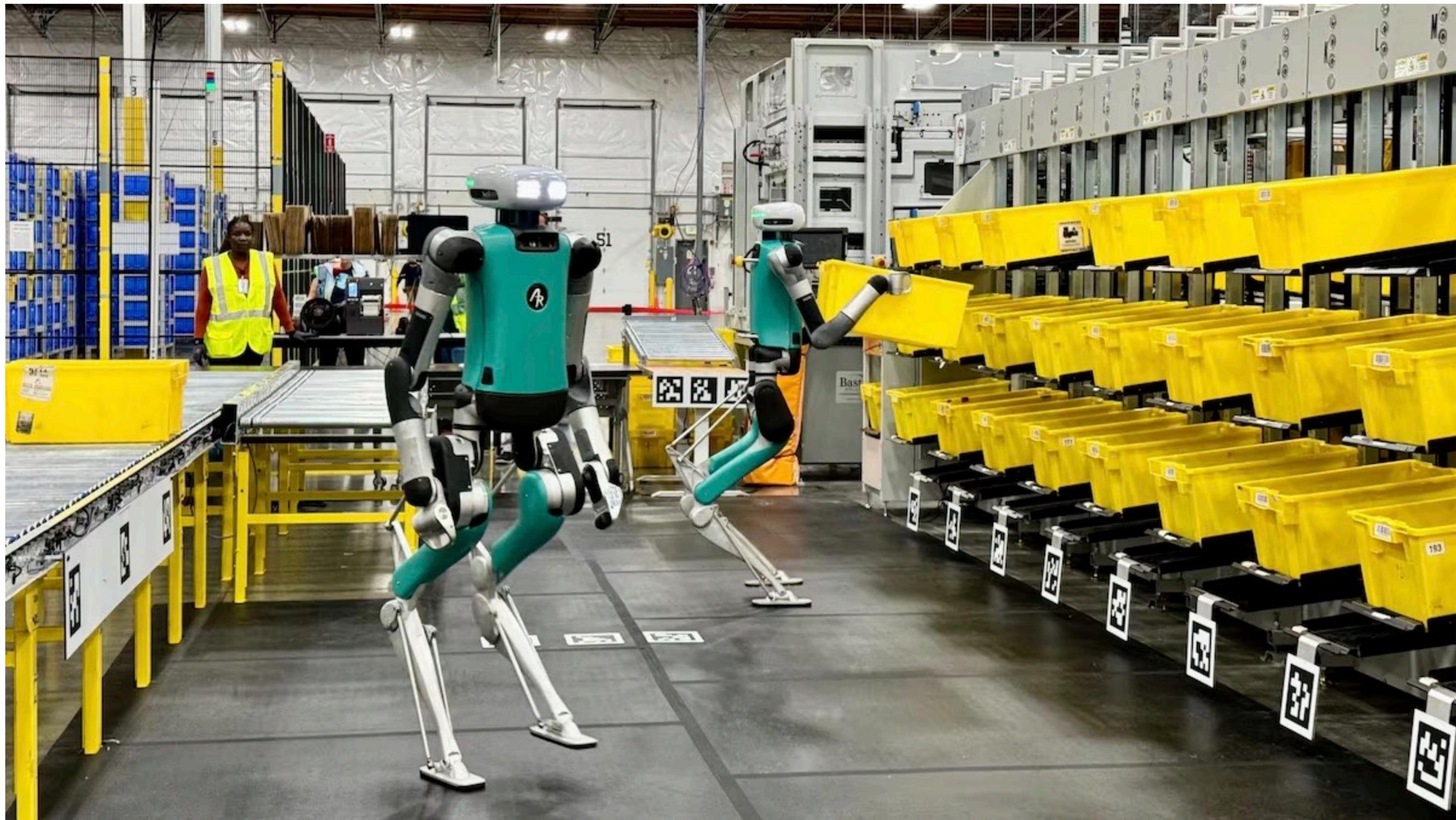


Le matin même...



Le matin même...









Robotiser le poste de travail?







ou redéfinir le process?

Disparition d'emplois Et transformation de métiers



Photographie du pool de dactylographie des usines Renault, 1931.
Régie nationale des usines Renault SA.
Archives centrales, documentation historique, album n° 5, photo n° 201124.



2020



2021

Thimble

Disparition d'emplois Et transformation de métiers



Photographie du pool de dactylographie des usines Renault, 1931.
Régie nationale des usines Renault SA.
Archives centrales, documentation historique, album n° 5, photo n° 201124.



2020



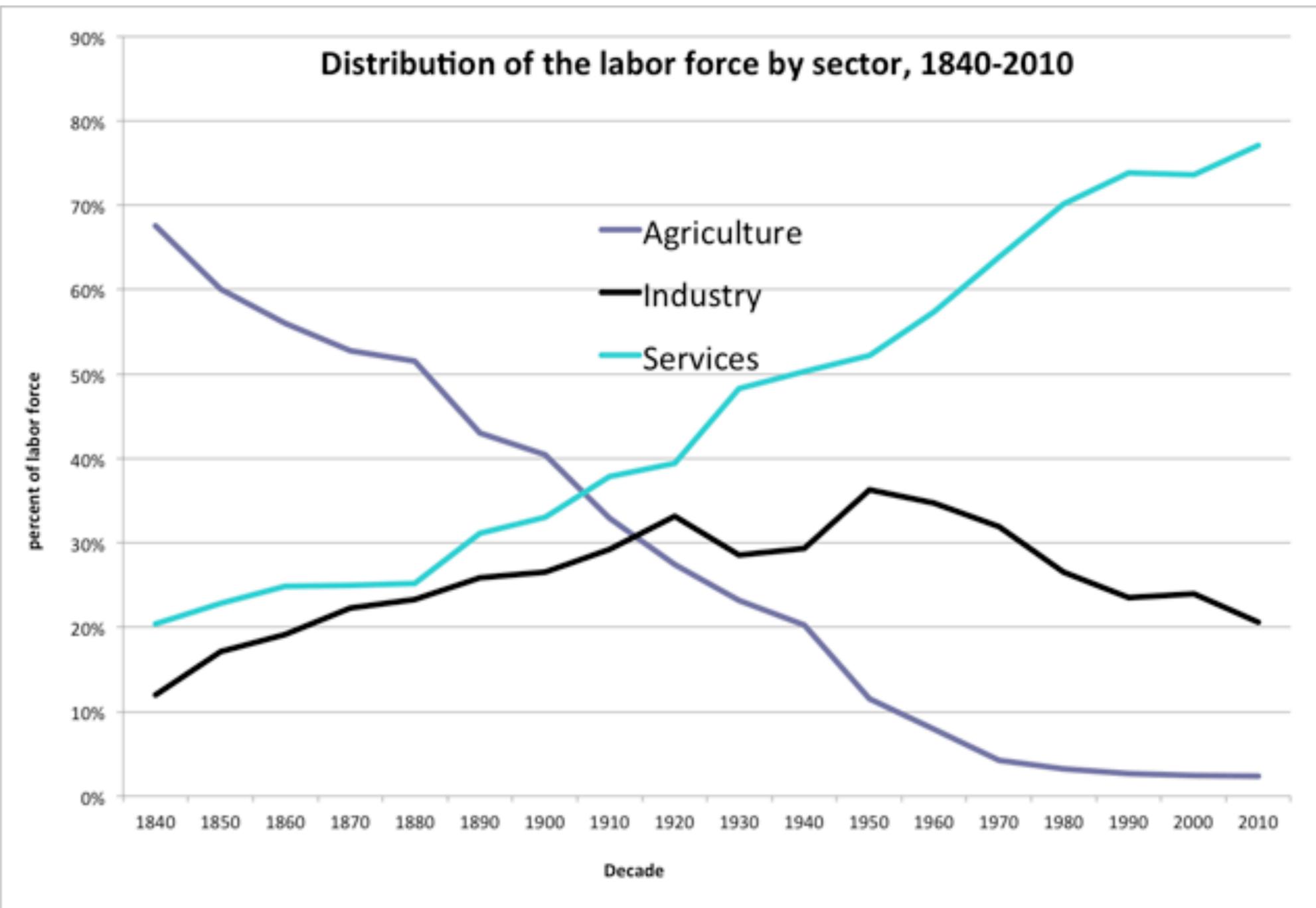
2021

Thimble

De nombreuses tâches dans tous les métiers ont été transformées par les technologies de l'informatique

L'IA poursuit et étend cette transformation

Il faut anticiper, acculturer et former



1840–1900: Robert E. Gallman and Thomas J. Weiss. "The Service Industries in the Nineteenth Century." In *Production and Productivity in the Service Industries*, ed. Victor R. Fuchs, 287-352. New York: Columbia University Press (for NBER), 1969.

1900–1940: John W. Kendrick, *Productivity Trends in the United States*. Princeton: Princeton University Press (for NBER), 1961.

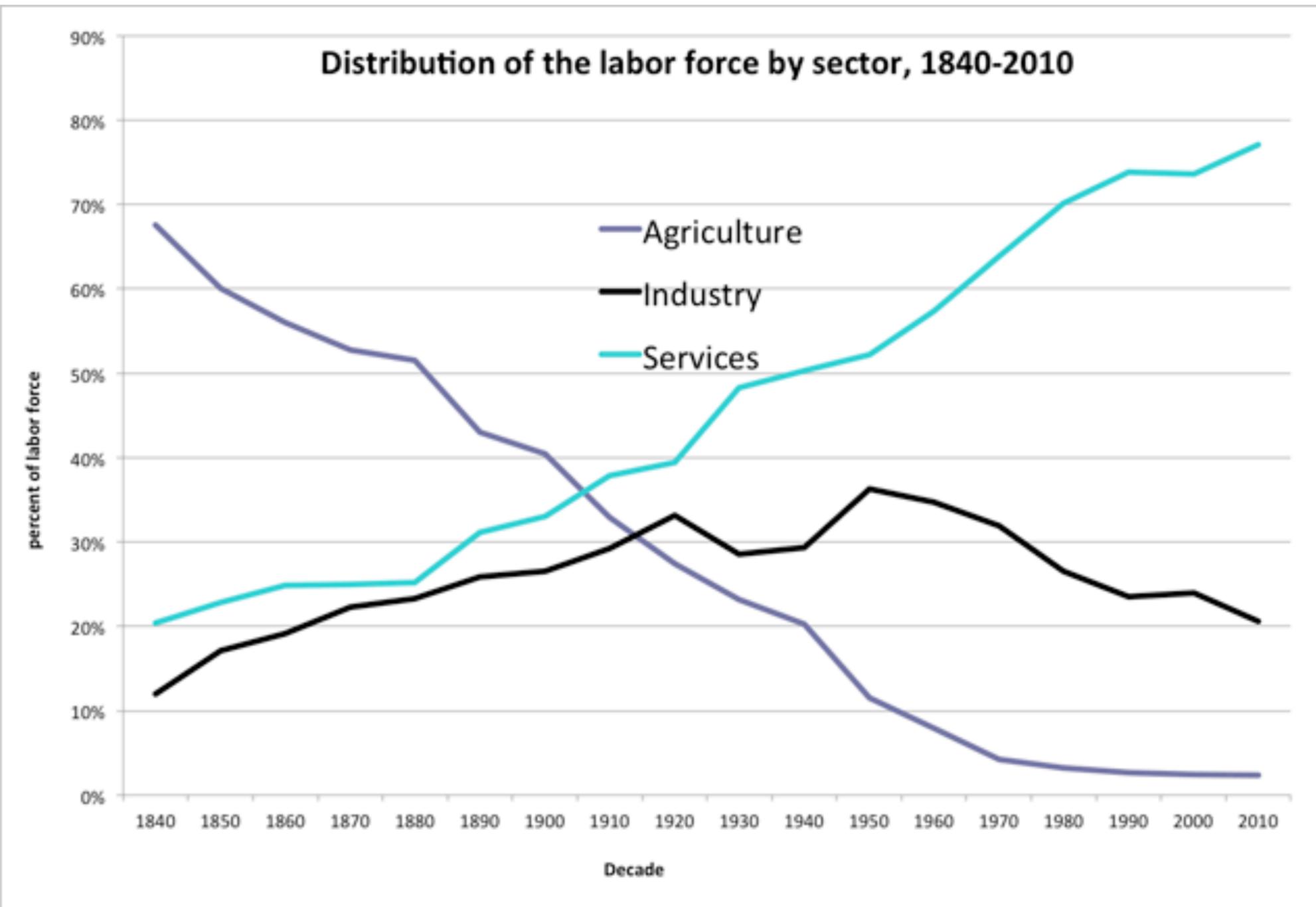
1950–2010: Bureau of Economic Analysis, *National Income and Product Accounts*.

De nombreuses tâches dans tous les métiers ont été transformées par les technologies de l'informatique

L'IA poursuit et étend cette transformation

Il faut anticiper, acculturer et former

Qui paie ?
Quelle distribution des bénéfices ?
Quel impact sur les modes et la qualité de vie ?



1840–1900: Robert E. Gallman and Thomas J. Weiss. "The Service Industries in the Nineteenth Century." In *Production and Productivity in the Service Industries*, ed. Victor R. Fuchs, 287-352. New York: Columbia University Press (for NBER), 1969.

1900–1940: John W. Kendrick, *Productivity Trends in the United States*. Princeton: Princeton University Press (for NBER), 1961.

1950–2010: Bureau of Economic Analysis, *National Income and Product Accounts*.

A retenir

- L'IA statistique y compris l'IA générative sont des technologies utiles mais qui ont des **limitations inhérentes** qui ne peuvent être corrigées (biais, robustesse, sécurité, ...)
- La prise de décision par des systèmes “autonomes” et “intelligents” qui ne comprennent pas le monde réel **soulève des questions éthiques et sociales**
- Ces systèmes doivent être utilisés en prenant des précautions. La **supervision humaine** est indispensable pour des applications critiques
- Effets à long terme sur la **formation, sur l'acquisition des compétences et sur les métiers** est difficile à prévoir ; études et recherches à engager.
- La réduction de l'**impact environnemental** est un problème ouvert
- Le règlement européen de l'IA définit des **responsabilités légales** pour cadrer le développement et l'usage des systèmes d'IA.
- Des **standards** et des processus de **certification** restent encore à définir