# DES MODÈLES DE LANGUE À L'IA FORTE

Lundi 18 novembre 2024
AgroParisTech

Vincent Guigue
`vincent.guigue@agroparistech.fr`
`https://vguigue.github.io`
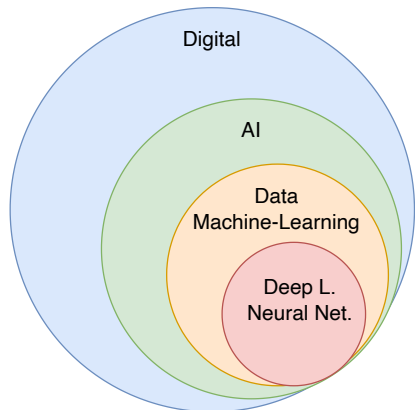
MIA
PARIS-SACLAY
EKINOCS

**A**GROPARISTECH
Institut des Sciences et Industries du Vivant et de l'Environnement

# From AI to
# deep-learning

# Digital & Artificial Intelligence

- Two related but distinct concepts
- AI: Different Definitions

| | |
|---|---|
| 1956 | Any algorithm / program |
| 1960-2012 | Expert systems and logical reasoning |
| 2012- | Data & neural networks |

Digital

AI

Data
Machine-Learning

Deep L.
Neural Net.

A. Turing

Marvin Minsky

G. Hinton

Y. Lecun

Computer          Neural Networks          Deep-learning

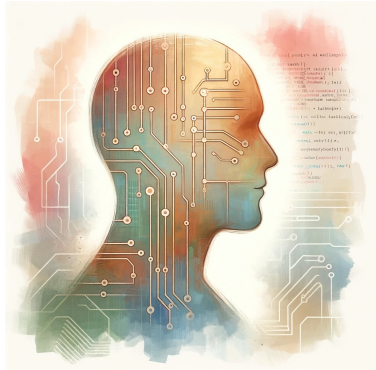**1941          1956          1986          2012**

Computer-
Sciences

AI: wide variety of algorithms
Mainly : Expert System + Reasonning

AI= Neural Networks

2/68

# Artificial Intelligence & Machine Learning



| Input ($\mathbf{X}$) | | Output ($\mathbf{Y}$) | Application |
|---|---|---|---|
| email | $\longrightarrow$ | spam? (0/1) | spam filtering |
| audio | $\longrightarrow$ | text transcript | speech recognition |
| English | $\longrightarrow$ | Chinese | machine translation |
| ad, user info | $\longrightarrow$ | click? (0/1) | online advertising |
| image, radar info | $\longrightarrow$ | position of other cars | self-driving car |
| image of phone | $\longrightarrow$ | defect? (0/1) | visual inspection |

**AI:** computer programs that engage in tasks which are, for now, performed more satisfactorily by human beings because they require high-level mental processes.

*Marvin Lee Minsky, 1956*

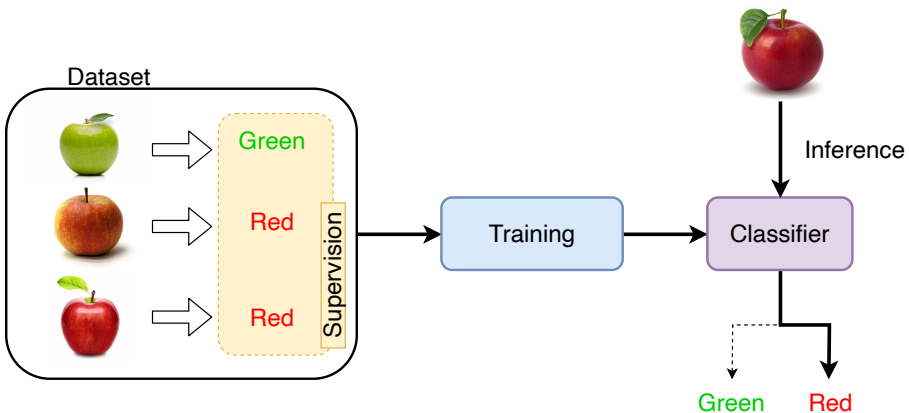**N-AI (Narrow Artificial Intelligence)**, dedicated to a single task

≠ **G-AI (General AI)**, which replaces humans in complex systems.
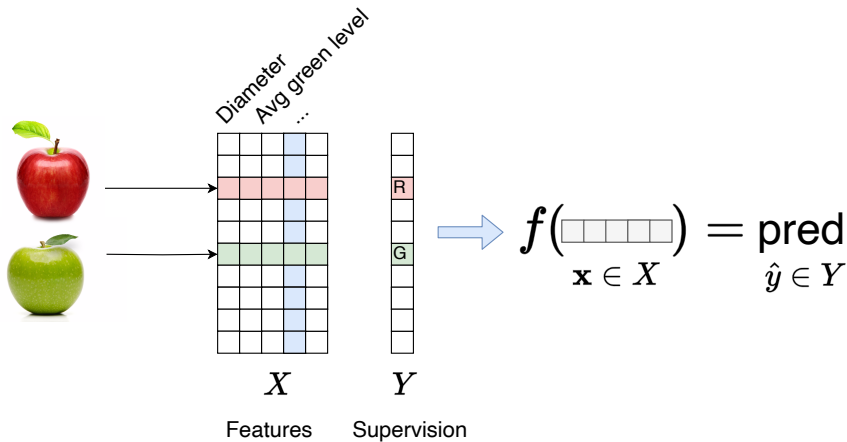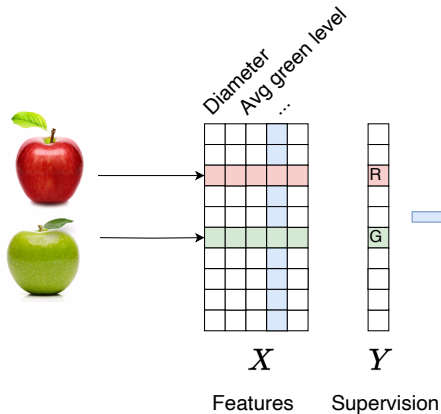
*Andrew Ng, 2015*

# Machine Learning Definition

1. Collecting labeled **dataset**
2. Training **classifier**
3. Exploiting the model

# Machine Learning Definition

1. Collecting labeled **dataset**
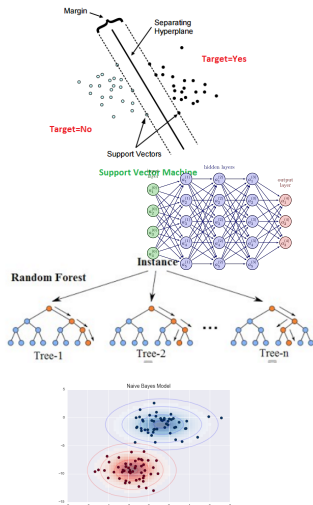2. Training **classifier**
3. Exploiting the model



$$\underset{\mathbf{x}\,\in\,X}{f(\boxed{\ \ \ \ \ })} = \underset{\hat{y}\,\in\,Y}{\mathsf{pred}}$$

$X$     $Y$

Features    Supervision

# Machine Learning Definition

1. Collecting labeled **dataset**
2. Training **classifier**
3. Exploiting the model



$$f(\underset{\mathbf{x} \in X}{\boxed{\phantom{xxxx}}}) = \underset{\hat{y} \in Y}{\mathsf{pred}}$$

$X$      $Y$

Features    Supervision
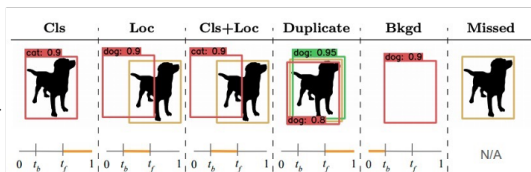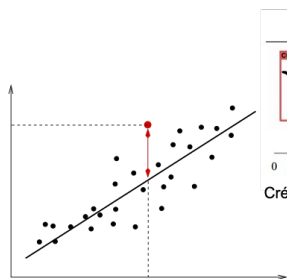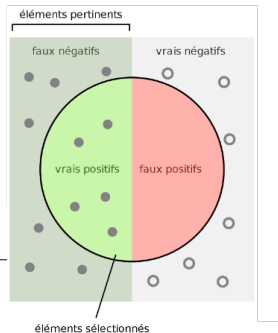
# Measuring Performance

Estimating performance (in generalization)… as important as training the model!



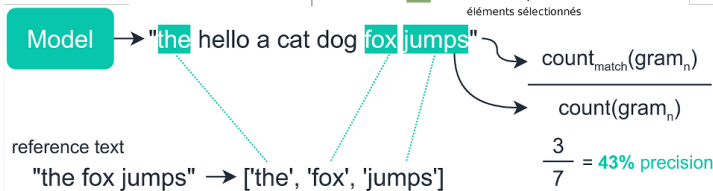Crédit: https://github.com/phalanx-hk/eccv2020_paperlist/issues/5
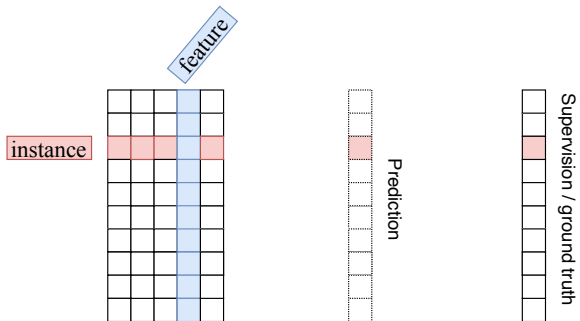
# Measuring Performance

Estimating performance (in generalization)... as important as training the model!

# Measuring Performance

Estimating performance (in generalization)… as important as training the model!

# General AI *vs* Narrow AI

**Narrow AI**

Like any computer science project:

- Define Inputs & Outputs
- Break down into subtasks
- Build & test components (processing chain)
- Assert (limited) generalization (iid assumption)
- Performances Evaluation

**General AI**

- Augmented Generalization Capability (Universality)
- Autonomous Learning
    - Data/information access
    - Knowledge extraction (Training+Eval+Confidence/Trust)
- Reasoning
- Conscience, Intentionality



Turing test

Wikipedia

6/68

# From tabular data to text

→ Tabular data
  → Fixed dimension
  → Continuous values

$\implies$ f( ⬜⬜⬜⬜⬜ ) = pred

→ Textual data
  → Variable length
  → Discrete values

this new iPhone, what a marvel

An iPhone? What a scam!

# AI + Textual Data: Natural Language Processing (NLP)

NLP = largest scientific community in AI

## Linguistics [1960-2010]

### Rule-based Systems:

$* \longrightarrow$ {like, love, appreciate} $\rightarrow * \rightarrow$ #product

$* \rightarrow$ {didn't, not, doesn't, don't} $\rightarrow$ {like, love, appreciate} $\rightarrow * \rightarrow$ #product

$* \longrightarrow$ {hate, loathe, detest} $\rightarrow * \rightarrow$ #product

- Requires expert knowledge
- Rule extraction $\Leftrightarrow$
        very clean data
- Very high precision
- Low recall
- Interpretable system

# AI + Textual Data: Natural Language Processing (NLP)

NLP = largest scientific community in AI

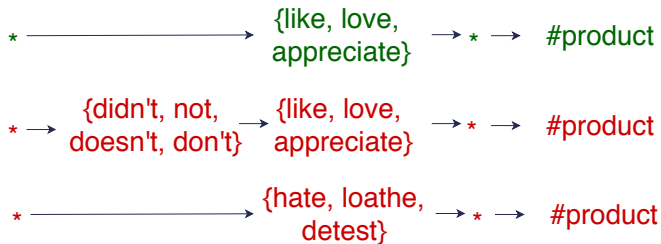**Machine Learning [1990-2015]**



$$f(x) = \sum_j w_j x_j \approx y$$

# AI + Textual Data: Natural Language Processing (NLP)

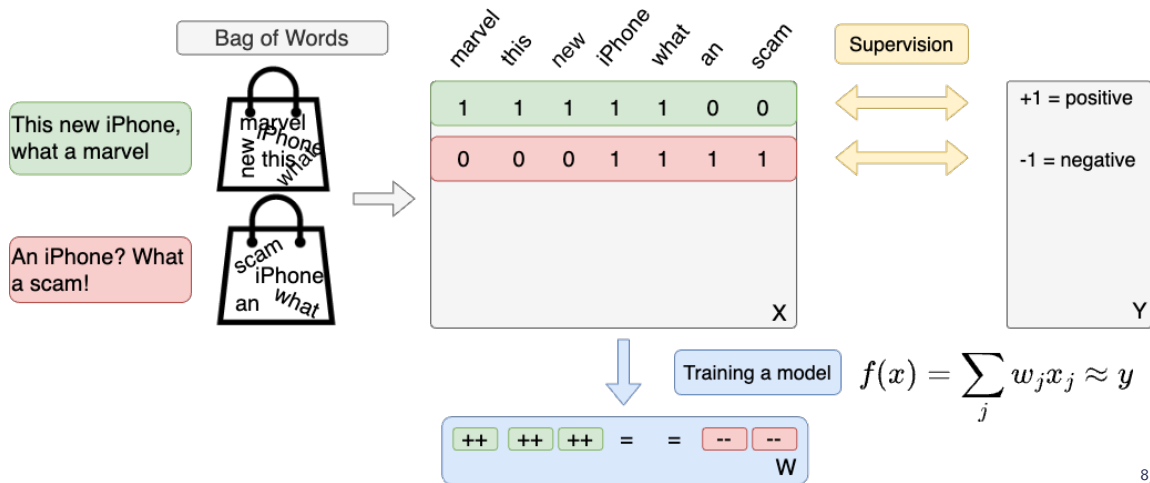NLP = largest scientific community in AI

## Linguistics [1960-2010]

- Requires expert knowledge
- Rule extraction ⇔

    very clean data

+ Interpretable system
+ Very high precision
− Low recall

## Machine Learning [1990-2015]

- Little expert knowledge needed
- Statistical extraction ⇔

    robust to noisy data

≈ Less interpretable system
− Lower precision
+ Better recall

Precision = criterion for acceptance by industry
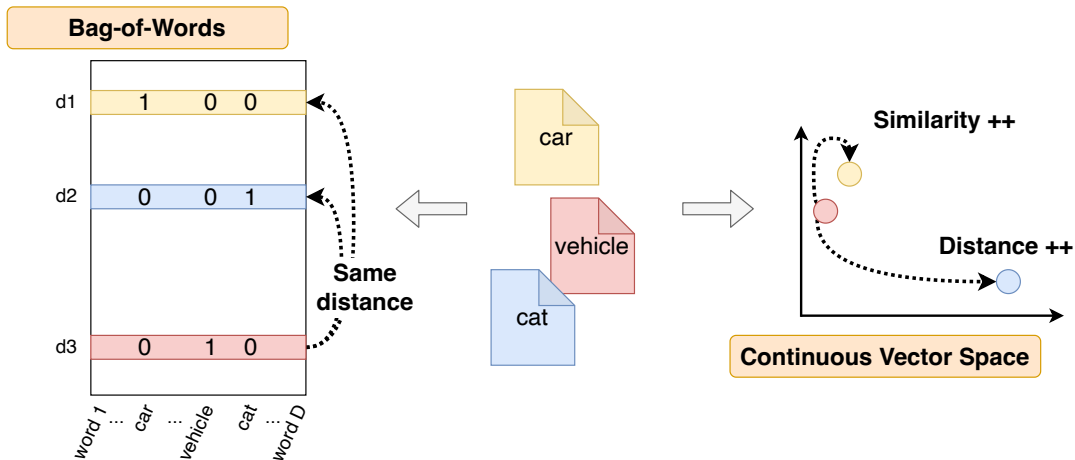
→ Link to metrics

# Deep/Representation Learning for Text Data

From Bag of Words to Vector Representations  [2008, 2013, 2016]



LeCun, Y., Bengio, Y., Hinton, G. (2015). Deep learning. Nature, 521(7553), 436-444.

# Deep/Representation Learning for Text Data

## From Bag of Words to Vector Representations          [2008, 2013, 2016]

The fluffy **cat** napped lazily in the sunbeam.
I adopted a stray **cat** from the shelter last week.
My **cat** loves to chase after toy mice.
The black **cat** stealthily crept through the dark alley.
I often find my **cat** perched on the windowsill, watching birds.
She gently stroked her **cat**'s fur as it purred contentedly.
Our neighbor's **cat** frequently visits our backyard.
My **cat** has a preference for fish flavored **cat** food.
The **cat** stealthily stalked a mouse in the garden.
My grandmother has a collection of porcelain **cat** figurines.
The **cat** napped peacefully in the warm sunlight.

# Deep/Representation Learning for Text Data

## From Bag of Words to Vector Representations [2008, 2013, 2016]



- Semantic Space:
  similar meaning
  ⇔
  close position

- Structured Space:
  grammatical regularities,
  basic knowledge, ...

Distributed representations of words and phrases and their compositionality, Mikolov et al. NeurIPS 2013

# Aggregating word representations: towards generative AI

- Generation & Representation
- New way of learning word positions



```
The fluffy cat napped lazily in the sunbeam.
I adopted a stray cat from the shelter last week.
My cat loves to chase after toy mice.
The black cat stealthily crept through the dark alley.
I often find my cat perched on the windowsill, watching birds.
She gently stroked her cat's fur as it purred contentedly.
Our neighbor's cat frequently visits our backyard.
The playful cat swatted at the dangling string with its paw.
My cat has a preference for fish flavored cat food.
The cat stealthily stalked a mouse in the garden.
My grandmother has a collection of porcelain cat figurines.
```
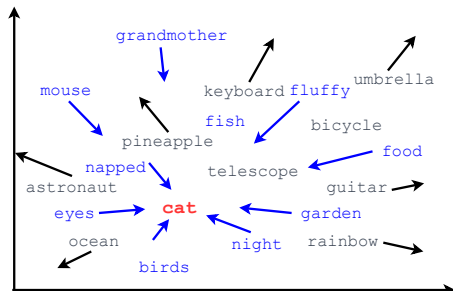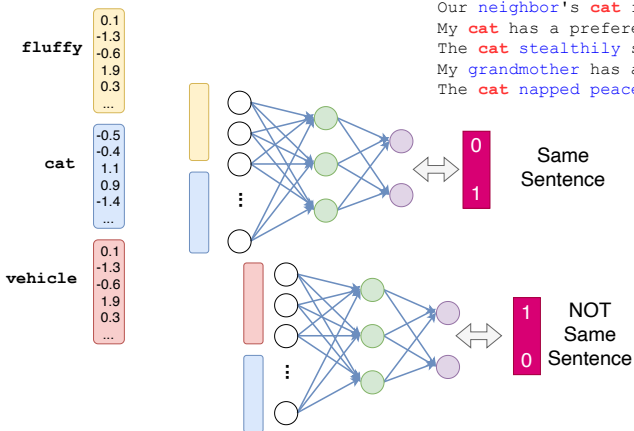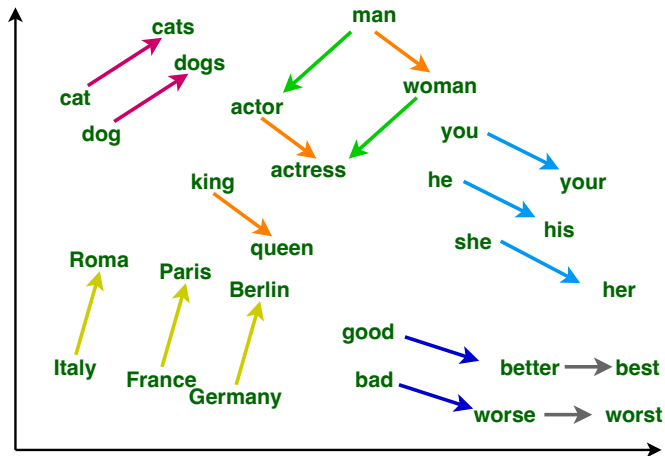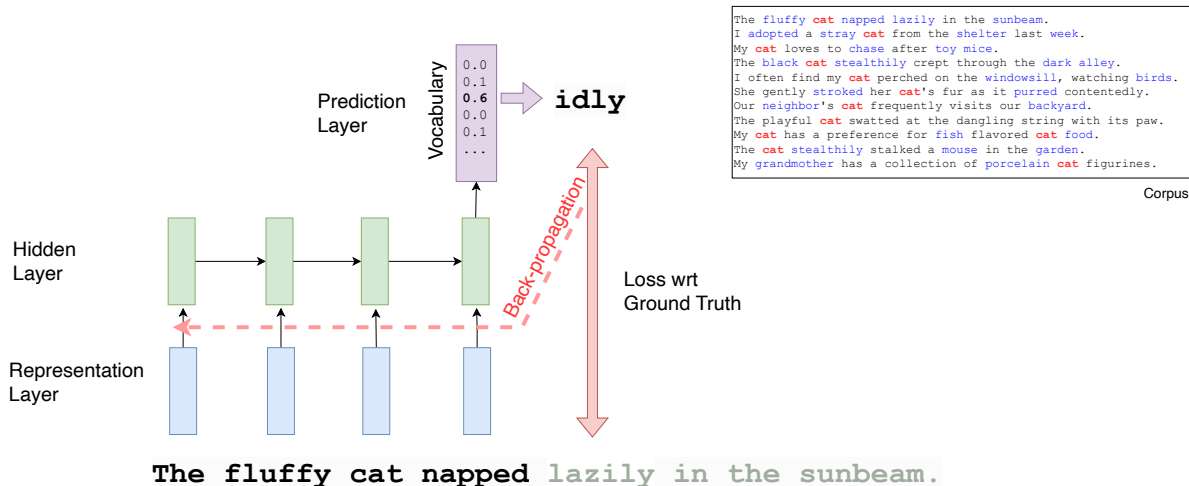
Corpus

Prediction Layer

Vocabulary

0.0
0.1
**0.6**
0.0
0.1
...

⟹ `idly`

Back-propagation

Loss wrt Ground Truth

Hidden Layer

Representation Layer

**The fluffy cat napped** lazily in the sunbeam.

Sequence to Sequence Learning with Neural Networks, Sutskever et al. NeurIPS 2014

# Inference & Beam Search



- High cost $\approx$ 1 call / token
- Max. likelihood principle
- NLP historical task =
  - specific classif./scoring archi.
  - constraint and/or post processing on generative archi.
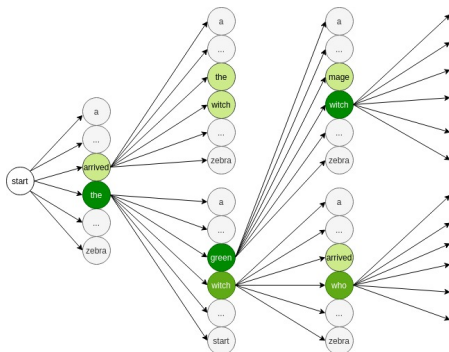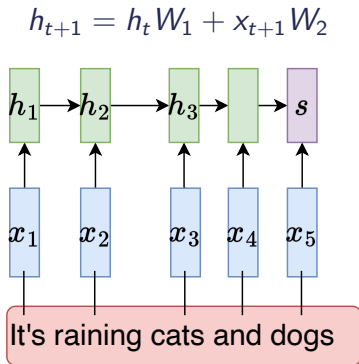
# Transformer architecture: state-of-the-art aggregation

**Recurrent Neural Network:**

$$h_{t+1} = h_t W_1 + x_{t+1} W_2$$

$h_1 \to h_2 \to h_3 \to \to s$

$x_1 \quad x_2 \quad x_3 \quad x_4 \quad x_5$

It's raining cats and dogs

**Transformer:**

$v_i''$

Fully Connected

$$v_i' = \sum_j \alpha_{ij} v_j$$

Self-attention Matrix

$\alpha_{ij}$

Transformer Layer

Token embeddings $\quad v_i$

it's raining cats and dogs

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

Sequence to Sequence Learning with Neural Networks, Sutskever et al. NeurIPS 2014

# A new developpement paradigm since 2015

- Huge dataset + huge archi. $\Rightarrow$ unreasonable training cost
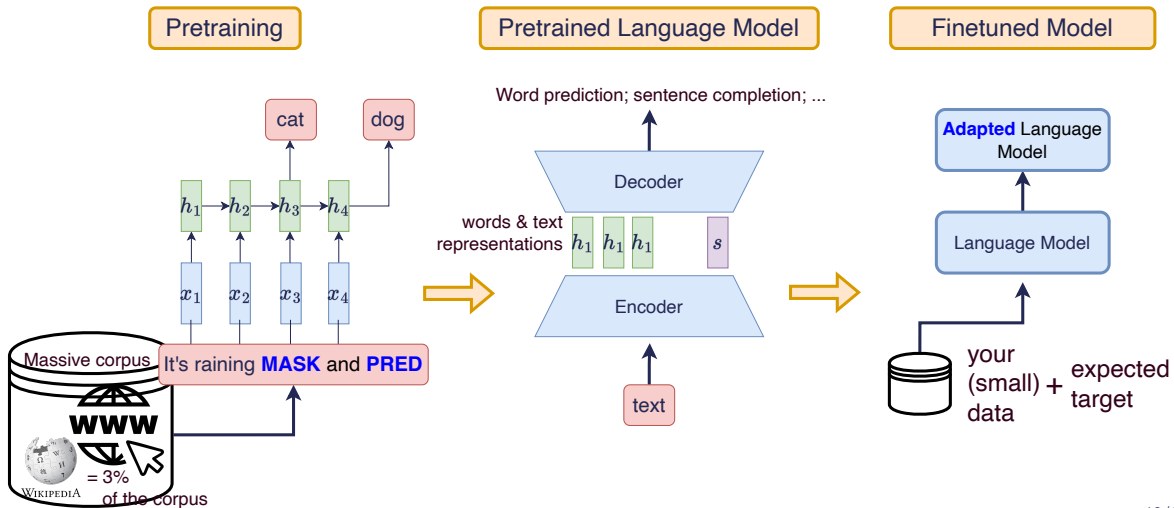- Pre-trained architecture + 0-shot / finetuning

# chatGPT

November 30, 2022

1 million users in 5 days
100 million by the end of January 2023
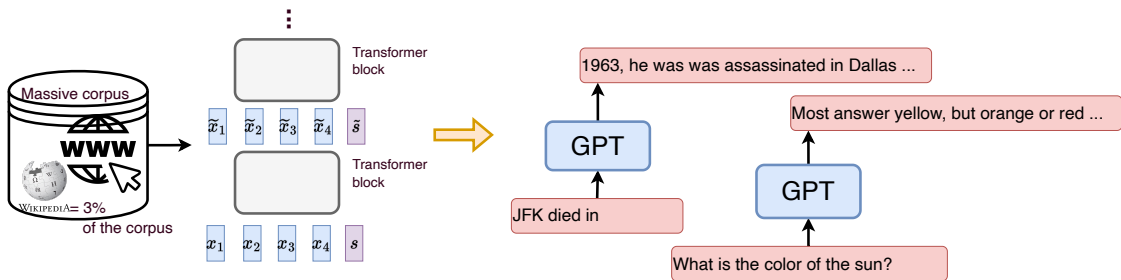1.16 billion by March 2023

# The Ingredients of chatGPT

## 0. Transformer + massive data (GPT)



- Grammatical skills: singular/plural agreement, tense concordance
- Knowledges

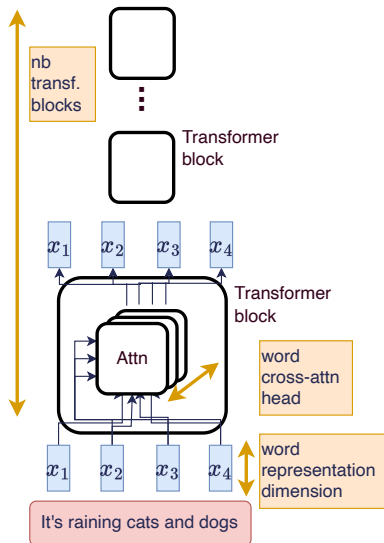Language Models are Few-Shot Learners, Brown et al. 2020

# The Ingredients of chatGPT

## 1. More is better! (GPT)

+ more input words          [500 $\Rightarrow$ 2k, 32k, 100k]
+ more dimensions in the word space   [500-2k $\Rightarrow$ 12k]
+ more attention heads              [12 $\Rightarrow$ 96]
+ more blocks/layers                [5-12 $\Rightarrow$ 96]

### **175 Billion** parameters... What does it mean?

- $1.75 \cdot 10^{11} \Rightarrow$ 300 GB + 100 GB (data storage for inference) $\approx$ 400GB
- NVidia A100 GPU = 80GB of memory (=20k€)
- Cost for (1) training: 4.6 Million €

nb transf. blocks

Transformer block

$x_1$ $x_2$ $x_3$ $x_4$

Transformer block

Attn

word cross-attn head

$x_1$ $x_2$ $x_3$ $x_4$

word representation dimension

It's raining cats and dogs

# The Ingredients of chatGPT

## 2. Dialogue Tracking



Specific training

GPT

Dialog follow-up
Coreference resolution
Way of speaking

Dialog corpus

- **Very clean** data          Data generated/validated/ranked by humans

# The Ingredients of chatGPT

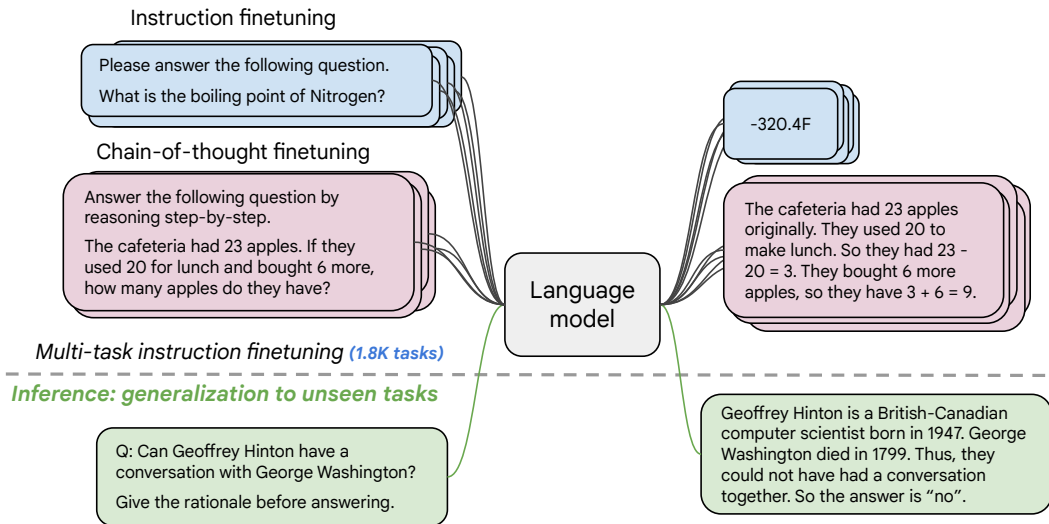## 3. Fine-tuning on different (±) complex reasoning tasks



Instruction finetuning

> Please answer the following question.
>
> What is the boiling point of Nitrogen?

Chain-of-thought finetuning

> Answer the following question by reasoning step-by-step.
>
> The cafeteria had 23 apples. If they used 20 for lunch and bought 6 more, how many apples do they have?

Language model

> -320.4F

> The cafeteria had 23 apples originally. They used 20 to make lunch. So they had 23 - 20 = 3. They bought 6 more apples, so they have 3 + 6 = 9.

*Multi-task instruction finetuning* **(1.8K tasks)**

*Inference: generalization to unseen tasks*

> Q: Can Geoffrey Hinton have a conversation with George Washington?
>
> Give the rationale before answering.

> Geoffrey Hinton is a British-Canadian computer scientist born in 1947. George Washington died in 1799. Thus, they could not have had a conversation together. So the answer is "no".

Scaling Instruction-Finetuned Language Models, Chung et al., JMLR 2024
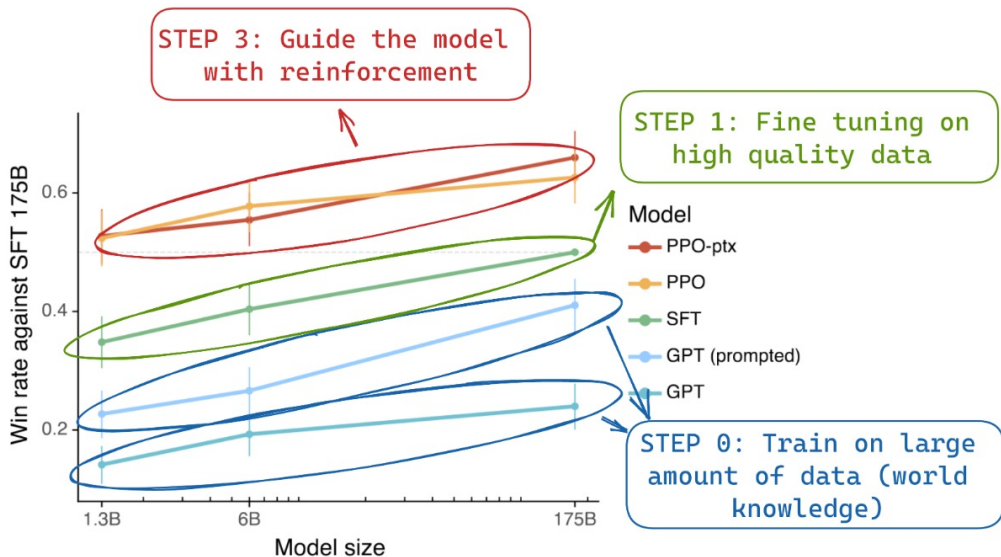
# The Ingredients of chatGPT

## 4. Instructions + answer ranking



- Database created by humans
- Response improvement

- ... Also a way to avoid critical topics = censorship

Training language models to follow instructions with human feedback, Ouyang et al., 2022

# Steps & Performance

Massive data $\Rightarrow$ HQ data (dialogue) $\Rightarrow$ Tasks $\Rightarrow$ RLHF

# Usage of chatGPT & Prompting

- Asking chatGPT = skill to acquire ⇒ *prompting*
  - Asking a question well: *... in detail, ... step by step*
  - Specify number of elements e.g. : *3 qualities for ...*
  - Provide context : *cell* for a biologist / legal assistant
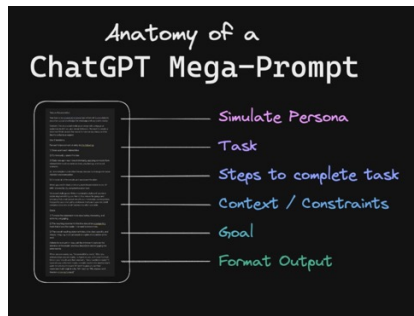
- Don't stop at the first question
  - Detail specific points
  - Redirect the research
  - Dialogue

- Rephrasing
  - Explain like I'm 5, like a scientific article, bro style, ...
  - Summarize, extend
  - Add mistakes (!)

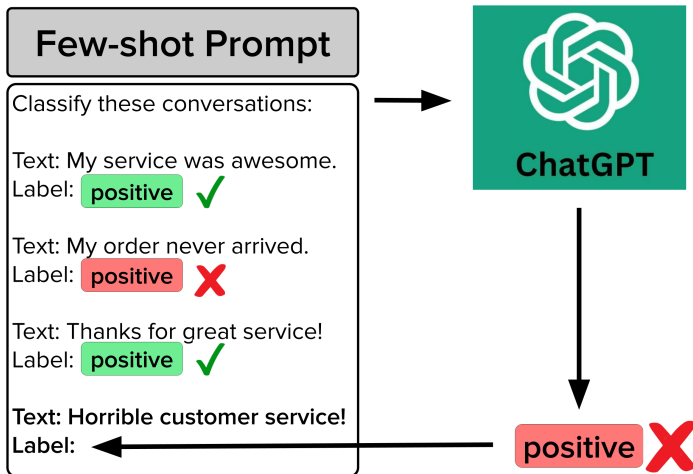  ⇒ Need for **practice** [1 to 2 hours], discuss with colleagues



Anatomy of a
ChatGPT Mega-Prompt

— Simulate Persona
— Task
— Steps to complete task
— Context / Constraints
— Goal
— Format Output

https://chatgptprompts.guru/what-makes-a-good-chatgpt-prompt/

# Towards *few-shot learning*

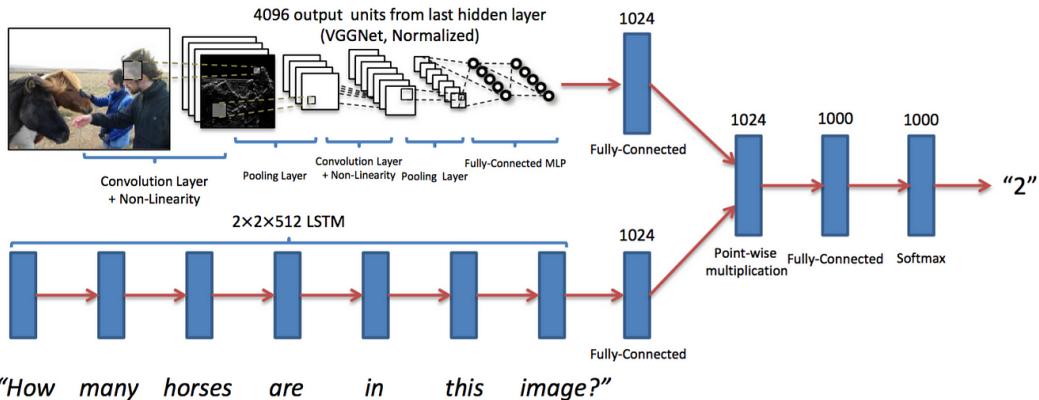■ Learning without modifying the model = examples in the prompt

# GPT4 & Multimodality

**Merging** information from text & image. **Learning** to exploit information jointly

*The example of VQA: visual question answering*



$\Rightarrow$ Backpropagate the error $\Rightarrow$ modify word representations + image analysis

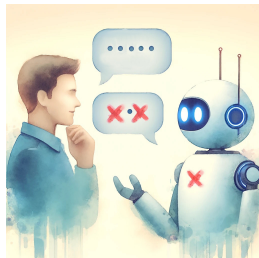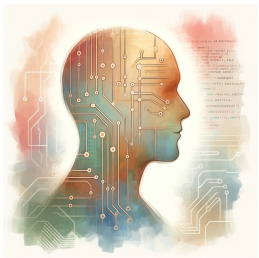📖 *VQA: Visual Question Answering*, arXiv, 2016 , A. Agrawal et al.

# Why So Much Controversy?
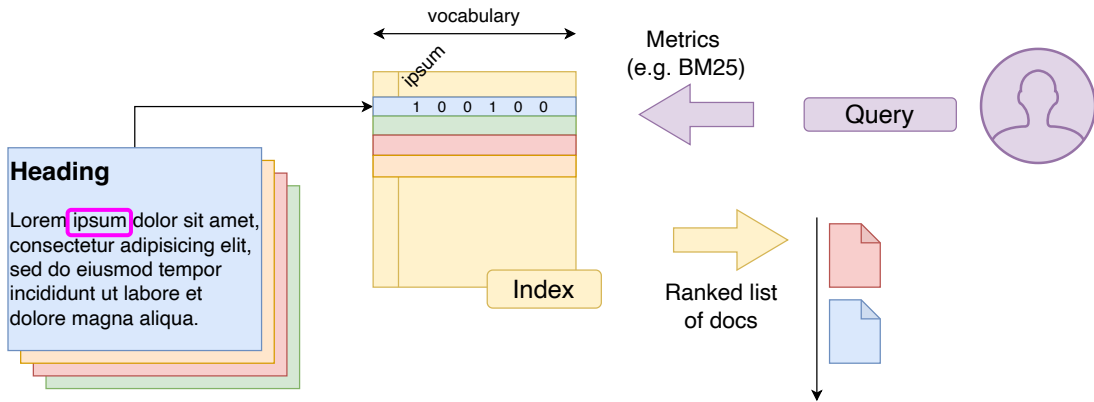
- New tool                                                    [December 2022]
- + Unprecedented adoption speed                      [1M users in 5 days]
- Strengths and weaknesses... Poorly understood by users
  - Significant productivity gains
  - Surprising / sometimes absurd uses
  - Bias / dangerous uses / risks
- Misinterpreted feedback
  - Anthropomorphization of the algorithm and its errors
- Prohibitive cost: what economic, ecological, and societal model?
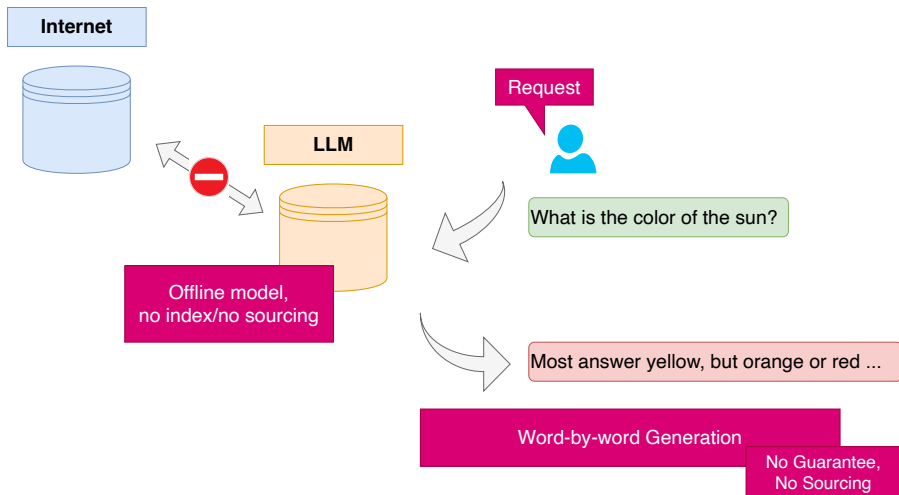
# Large Language Models Uses

# LLM & Information Retrieval



vocabulary

ipsum

| 1 | 0 | 0 | 1 | 0 | 0 |

Index

Metrics
(e.g. BM25)

Query

**Heading**

Lorem ipsum dolor sit amet, consectetur adipisicing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua.

Ranked list
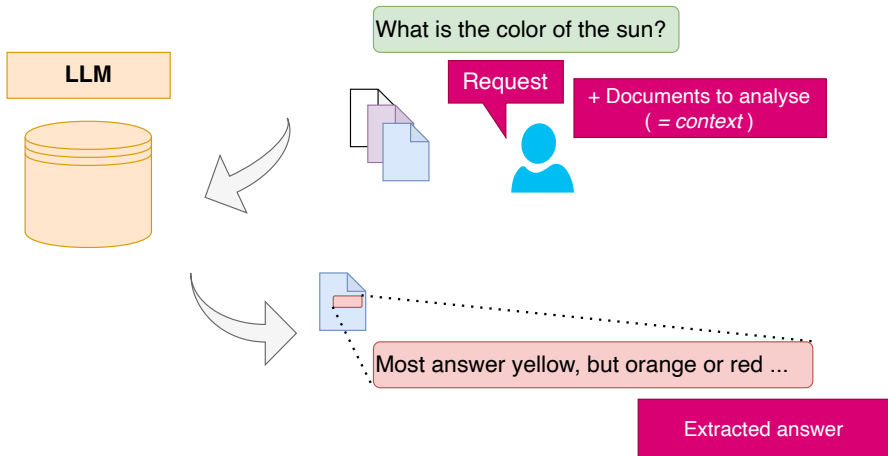of docs

# Information access: from word index to RAG

- Asking for information from ChatGPT... A surprising use!
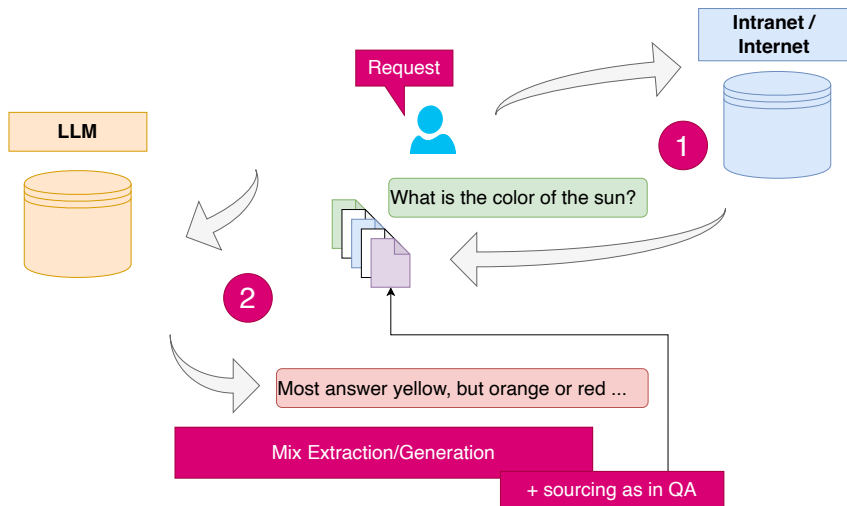- But is it reasonnable?                                    [Real Open Question (!)]



**Internet**

**LLM**

Offline model,
no index/no sourcing

**Request**

What is the color of the sun?

Most answer yellow, but orange or red ...

Word-by-word Generation

No Guarantee,
No Sourcing

# Information access: from word index to RAG



**LLM**

What is the color of the sun?

Request

+ Documents to analyse
( = *context* )

Most answer yellow, but orange or red ...

Extracted answer

- Web query + analysis, automatic summary, rephrasing, meeting reports...
- (Current) limit on input size (2k then 32k tokens)
- = *pre chatGPT use of LLM for question answering*

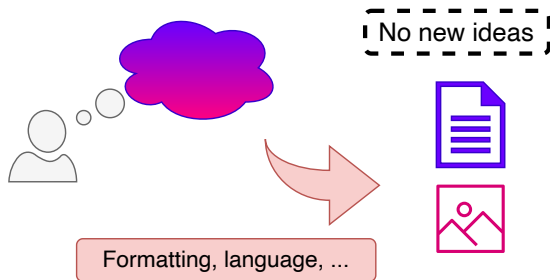# Information access: from word index to RAG



- RAG: Retrieval Augmented Generation
- (Current) limit on input size (2k then 32k tokens)

# Other Uses of Generative AIs

No new ideas

A fantastic tool for
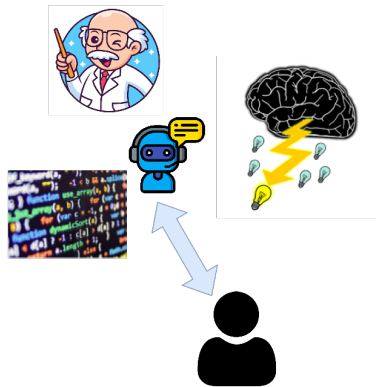**formatting**

Formatting, language, ...

- Personal assistant
  - Standard letters, recommendation letters, cover letters, termination letters
  - Translations
- Meeting reports
  - Formatting notes
- Writing scientific articles
  - Writing ideas, in French, in English
- Document analysis
  - Information extraction, question-answering, ...

# Other Uses of Generative AIs

## And a tool for **reflection**!

- Brainstorming
  - Argument development, contradiction search
- Assistant for software development
  - Code generation, error search, …
  - Documentation
- Educational assistant
  - Wikipedia $++$, proposal of outlines for essays,
  - Code explanation / correction proposals

# LLM & Teaching opportunities



- A great opportunity to have a 24/7 available teacher
- In particular for coding:
  - Learning python
  - Learning machine learning
⇒   1 Generate a small program
    2 Ask question about the different functions

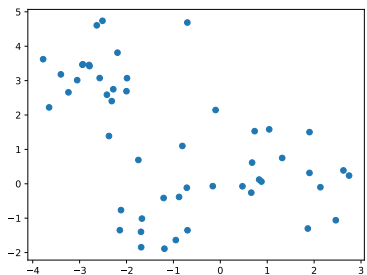LLM can do your homeworks... But LLM can explain you, answer questions about the solution, teach you!

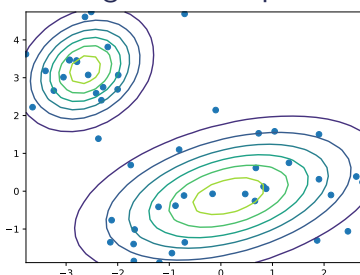# From Generative AI
# to Foundation Models

# At the origin of statistical modeling

1. **Observing** data (and context)
2. **Modeling** = Choosing probabilistic model / bayesian network
3. **Optimize** parameters (Max. Likelihood, EM, BFGS, ...)
4. **Sampling** / Inference + Evaluate distances : existing *vs* sampled

Observations



Modeling: choice+optim.



Sampling / eval.

# At the origin of statistical modeling

1. **Observing** data (and context)
2. **Modeling** = Choosing probabilistic model / bayesian network
3. **Optimize** parameters (Max. Likelihood, EM, BFGS, ...)
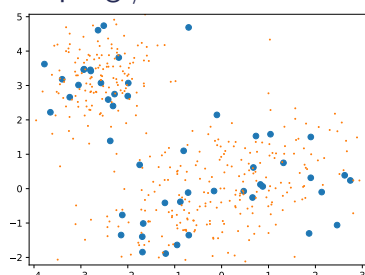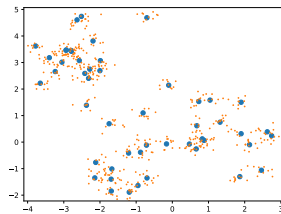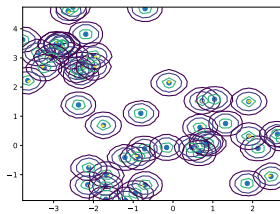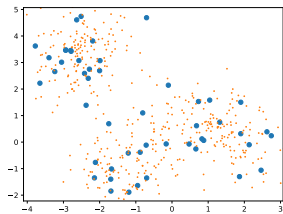4. **Sampling** / Inference + Evaluate distances : existing *vs* sampled

Different modeling options / different traps

# At the origin of deep learning

■ Gradient vanishing issue in deep architecture



Gradient backpropagation

Gradient **weakening => vanishing**

## At the origin of deep learning

- Gradient vanishing issue in deep architecture
- Auto-Encoder architecture / facing unsupervised dataset with NN

noise$(X)$      $\widetilde{X}$    Ground Truth $X^{\star}$

data

Loss

- Denoising
- Low dimensional representation learning (/ PCA, SVD)

*Auto-association by multilayer perceptrons and singular value decomposition*, Biological Cybernetics, 1988
H. Bourlard & Y. Kamp

# At the origin of deep learning

- Gradient vanishing issue in deep architecture
- Auto-Encoder architecture / facing unsupervised dataset with NN
- Stacked Denoising Auto-Encoder : iterative training / **pretraining**



*The difficulty of training deep architectures and the effect of unsupervised pre-training*, AIS, PMLR 2009
Erhan, D., Manzagol, P. A., Bengio, Y., Bengio, S., & Vincent, P.

# Variational Auto-Encoder



- a priori on the distribution
- Structuring of the latent space

Generative AI (for statisticians)

📖 *Auto-Encoding Variational Bayes*, 2013
DP Kingma

# Different Forms of Generative AI

Compact
Vector
Representation

Input

Output

Encoder

Decoder

1 Encode an input = construct a vector

2 Decode a vector = *generate* an output

# Different Media / Different Architectures

- Texts: classification problem

# Different Media / Different Architectures

- Texts: classification problem
- Images: multivariate regression problem



*U-Net: Convolutional Networks for Biomedical Image Segmentation*, MICCAI, 2015
Ronneberger et al.

NVidia Lab.

# Different Media / Different Architectures

- Texts: classification problem
- Images: multivariate regression problem

Generative Adversarial Networks (GAN): detecting generated samples



*Generative Adversarial Nets*, NeurIPS 2014
Goodfellow et al.

# Different Media / Different Architectures

- Texts: classification problem
- Images: multivariate regression problem
- Physical processes



📖 *Denoising Diffusion Probabilistic Models*, NeurIPS, 2020
Ho, J., Jain, A., & Abbeel, P.

📖 *Hierarchical Text-Conditional Image Generation with CLIP Latents*, arXiv, 2022
Ramesh et al.

# Different Media / Different Architectures

- Texts: classification problem
- Images: multivariate regression problem
- Physical processes
- Complex structures / 3D / graphs: sequential problem

- Mix mechanistic and *data-driven* approaches

e.g. Model differential equations in a neural network

📖 *Neural ordinary differential equations*, NeurIPS, 2018
Chen et al.

📖 *Physics-informed neural networks* J. Comp. Physics, 2019
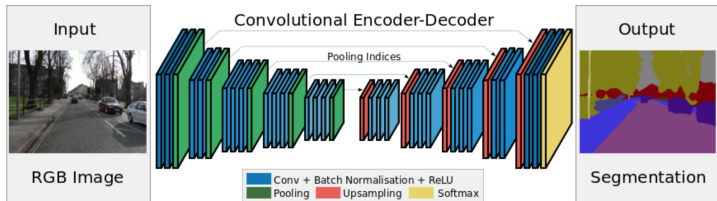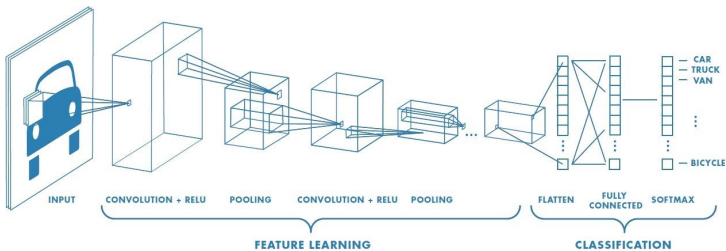Raissi et al.



Navier-Stokes PINN model

Discriminative GAN model

# Different Media / Different Architectures

- Texts: classification problem
- Images: multivariate regression problem
- Physical processes
- Complex structures / 3D / graphs: sequential problem



$u(t)$
**Input Function**

$$\frac{ds(t)}{dt} = u(t)$$

**ODE/PDE Solver**

**Neural Operator**

**Output Function**
$s(t)$

Data + Models :

- PDE, neural ODE
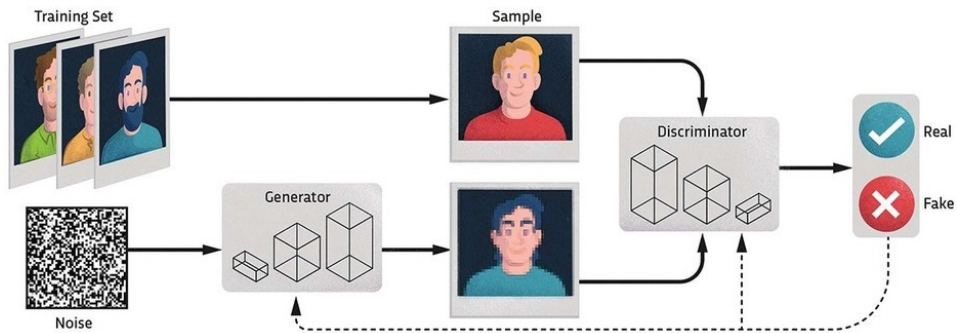- Simulation approximations
- Residual Models
- Hybrid Complex Systems
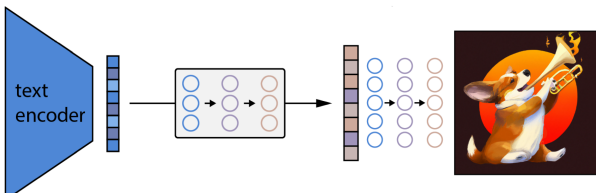
# Different Media / Different Architectures

- Texts: classification problem
- Images: multivariate regression problem
- Physical processes
- Complex structures / 3D / graphs: sequential problem

  - Reinforcement learning: action/reward



Apprentissage par renforcement

*Highly accurate protein structure prediction with AlphaFold*, Nature, 2021
Jumper et al.

# Multi-Modality

- Construction of multimodal representation spaces = *grounding*
- Image ⇒ Text: *Captioning, Visual Question Answering*
- Text ⇒ Image: *mid-journey, dall-e, ...*



**Alignment** of representation spaces

| Word | Teraword | Knext |
|------|----------|-------|
| Spoke | 11,577,917 | 372,042 |
| Laughed | 3,904,519 | 179,395 |
| Murdered | 2,843,529 | 16,890 |
| Inhaled | 984,613 | 5,617 |
| Breathed | 725,034 | 41,215 |

# Multi-Modality

- Construction of multimodal representation spaces = *grounding*
- Image ⇒ Text: *Captioning, Visual Question Answering*
- Text ⇒ Image: *mid-journey, dall-e, ...*



A cute little dog sitting in a heart drawn on a sandy beach.

A dog walking next to a little dog on top of a beach.

*Show and Tell: image captioning open sourced in TensorFlow*, Chris Shallue , Google Research, 2016

# Multi-Modality

- Construction of multimodal representation spaces = *grounding*
- Image ⇒ Text: *Captioning, Visual Question Answering*
- Text ⇒ Image: *mid-journey, dall-e, ...*



*Vqa: Visual question answering*, ICCV, 2015
Antol et al.

# Multi-Modality

- Construction of multimodal representation spaces $=$ *grounding*
- Image $\Rightarrow$ Text: *Captioning, Visual Question Answering*
- Text $\Rightarrow$ Image: *mid-journey, dall-e, ...*

# Towards Larger Foundation Models?

■ Let the modalities enrich each other



*On the Opportunities and Risks of Foundation Models*, Tech. Report, Stanford, 2021
Bommasani et al.

# Conclusion

## The main challenges of multimodality

- New applications
    - at the interface between text, image, music, voice, …

- Performance improvement
    - Better encoding, disambiguation, context encoding

- Explainability (through dialogue)
    - IoT / RecSys / Intelligent Vehicle / …



Dall-e

# Machine Learning Limits

# chatGPT and the relationship with truth

1. Plausibility = grammar, agreement, tense concordance, logical sequences...
   ⇒ Repeated knowledge

2. Predict the most **plausible** word...
   ⇒ produces **hallucinations**

3. **Offline** functioning

4. chatGPT ≠ **knowledge graphs**

5. Brilliant answers...

   And silly mistakes!
   + we cannot predict the errors

1963, he was was assassinated in Dallas ...

GPT

JFK died in

### Example: producing a bibliography

Can you give me a short bibliography on VAE for Time Series?

    recurrents.

2. "Variational Recurrent Autoencoders" par Chung et al. (2016) - Cette étude introduit une approche VAE pour la modélisation de séquences temporelles en utilisant des réseaux de neurones récurrents et une méthode de maximisation de la vraisemblance pour la phase d'entraînement.

**Variational recurrent auto-encoders**
O Fabius, JR Van Amersfoort - arXiv preprint arXiv:1412.6581, 2014 - arxiv.org
In this paper we propose a model that combines the strengths of RNNs and SGVB: the **Variational Recurrent** Auto-Encoder (VRAE). Such a model can be used for efficient, large scale ...
☆ Enregistrer 〃 Citer Cité 302 fois Autres articles Les 2 versions ≫

3. "Generative Modeling for Time ... Bao et al. (2017) - Cette étude ... pour la modélisation de séries t... profonds, y compris les VAE.

4. "Deep Variational Bayes Filters: Unsupervised Learning of State Space Models from Raw Data" par Krishnan et al. (2017) - Cette étude présente une approche VAE pour la

# Generative AI: how to evaluate performance?

## The critical point today

- How to evaluate against ground truth?
- How to evaluate system confidence / plausibility of generation?



$$\frac{count_{match}(gram_n)}{count(gram_n)}$$

$$\frac{3}{3} = \textbf{100\% recall}$$

The Ultimate Performance Metric in NLP, J. Briggs, Medium 2021

# Generative AI: how to evaluate performance?

## The critical point today

- How to evaluate against ground truth?
- How to evaluate system confidence / plausibility of generation?



Plausibility



Train



Test

*Do Large Language Models Know What They Don't Know?*, Yin et al. , ACL, 2023
*Do Deep Generative Models Know What They Don't Know?*, Nalisnick et al. , ICLR, 2019

# Stability/predictability



- Difficult to bound a behavior
- Impossible to predict good/bad answers

⇒ Little/no use in video games

V   how old is Obama

🪷   Barack Obama was born on August 4, 1961, making him 61 years old as of February 2,   👍  👎
2023.

# Stability/predictability



- Difficult to bound a behavior
- Impossible to predict good/bad answers

$\Rightarrow$ Little/no use in video games

| V | how old is obama? |

| 🌀 | As of 2021, Barack Obama was born on August 4, 1961, so he is 60 years old.    👍 👎 |

| V | and today? |

# Stability, explainability... And complexity



Sensor 1
Sensor 2
⋮
Sensor d

Simple rules

Up/Down
Flashing light
0

Vocabulary (huge)

it's raining cats    Aggregation    and dogs

Word sequence (= combination)    Word prediction

- Simple system
- Exhaustive testing of inputs/outputs
- **Predictable** & **explainable**

- Large dimension
- Complex non-linear combinations
- **Non-predictable** & **non-explainable**

# Stability, explainability... And complexity

## Interpretability *vs* Post-hoc Explanation

Neural networks = **non-interpretable** (almost always)

*too many combinations to anticipate*

Neural networks = **explainable a posteriori** (almost always)



[Uber Accident, 2018]

- Simple system
- Exhaustive testing of inputs/outputs
- **Predictable** & **explainable**

- Large dimension
- Complex non-linear combinations
- **Non-predictable** & **non-explainable**

39/68

# Transparency

- Model weights (*open-weight*)… $\Rightarrow$ but not just the weights
- Training data (*BLOOM*) + distribution + instructions
- Learning techniques
- Evaluation

**Foundation Model Transparency Index Scores by Major Dimensions of Transparency, 2023**
Source: 2023 Foundation Model Transparency Index

| Major Dimensions of Transparency | Meta Llama 2 | BigScience BLOOMZ | OpenAI GPT-4 | stability.ai Stable Diffusion 2 | Google PaLM 2 | ANTHROP\C Claude 2 | cohere Command | AI21 labs Jurassic-2 | Inflection Inflection-1 | amazon Titan Text | Average |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Data | 40% | 60% | 20% | 40% | 20% | 0% | 20% | 0% | 0% | 0% | 20% |
| Labor | 29% | 86% | 14% | 14% | 0% | 29% | 0% | 0% | 0% | 0% | 17% |
| Compute | 57% | 14% | 14% | 57% | 14% | 0% | 14% | 0% | 0% | 0% | 17% |
| Methods | 75% | 100% | 50% | 100% | 75% | 75% | 0% | 0% | 0% | 0% | 48% |
| Model Basics | 100% | 100% | 50% | 83% | 67% | 67% | 50% | 33% | 50% | 33% | 63% |
| Model Access | 100% | 100% | 67% | 100% | 33% | 33% | 67% | 33% | 0% | 33% | 57% |
| Capabilities | 60% | 80% | 100% | 40% | 80% | 80% | 60% | 60% | 40% | 20% | 62% |
| Risks | 57% | 0% | 57% | 14% | 29% | 29% | 29% | 29% | 0% | 0% | 24% |
| Mitigations | 60% | 0% | 60% | 0% | 40% | 40% | 20% | 0% | 20% | 20% | 26% |
| Distribution | 71% | 71% | 57% | 71% | 71% | 57% | 57% | 43% | 43% | 43% | 59% |
| Usage Policy | 40% | 20% | 80% | 40% | 60% | 60% | 40% | 20% | 60% | 20% | 44% |
| Feedback | 33% | 33% | 33% | 33% | 33% | 33% | 33% | 33% | 33% | 0% | 30% |
| Impact | 14% | 14% | 14% | 14% | 14% | 0% | 14% | 14% | 14% | 0% | 11% |
| **Average** | **57%** | **52%** | **47%** | **47%** | **41%** | **39%** | **31%** | **20%** | **20%** | **13%** | |

https://crfm.stanford.edu/fmti/May-2024/index.html

# (Main) Risks
## derived from ML & LLM

# Typology of AI Risks in NLP (L. Weidinger)

### Discrimination, exclusion and toxicity

Harms that arise from the language model producing discriminatory and exclusionary speech.

### Information hazards

Harms that arise from the language model leaking or inferring true sensitive information.

### Misinformation harms

Harms that arise from the language model producing false or misleading information.

### Malicious uses

Harms that arise from actors using the language model to intentionally cause harm.

### Human–computer interaction harms

Harms that arise from users overly trusting the language model, or treating it as human-like.

### Automation, access and environmental harms

Harms that arise from environmental or downstream economic impacts of the language model.

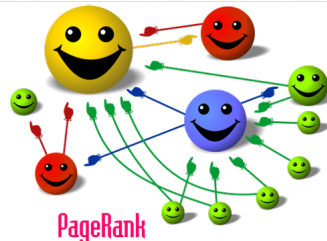# Access to Information



- Access to dangerous/forbidden information
    - +Personal data
    - Right to digital oblivion

- Information authorities
    - Nature: unconsciously, image = truth
    - Source: newspapers, social media, …
    - Volume: number of variants, citations (pagerank)

- Text generation: harassment…

- Risk of anthropomorphizing the algorithm
    - Distinguishing human from machine

# Machine Learning & Bias



Mustache, Triangular Ears, Fur Texture

Cat



Over 40 years old, white, clean-shaven, suit

Senior Executive

## Bias in the data ⇒ bias in the responses

Machine learning is based on extracting statistical biases...

⇒ Fighting bias = manually adjusting the algorithm

# Machine Learning & Bias



Stereotypes from *Pleated Jeans*

≡  **Google** Traduction                    ⚙ ⠿ V

| 🔤 Texte | 🖼 Images | 📄 Documents | 🖿 Sites Web |

Détecter la langue   **Anglais**   Français  ∨        ⇄        **Français**  Anglais  Arabe  ∨

| The nurse and the doctor                  ✕ | L'infirmière et le médecin              ☆ |

- ■ Gender choice
- ■ Skin color
- ■ Posture
- ■ …

## Bias in the data ⇒ bias in the responses

Machine learning is based on extracting statistical biases…
⇒ Fighting bias = manually adjusting the algorithm

# Bias Correction & Editorial Line

**Bias Correction:**

- Selection of specific data, rebalancing
- Censorship of certain information
- Censorship of algorithm results

⇒ Editorial work...                    Done by whom?

- Domain experts / specifications
- Engineers, during algorithm design
- Ethics group, during result validation
- Communication group / user response

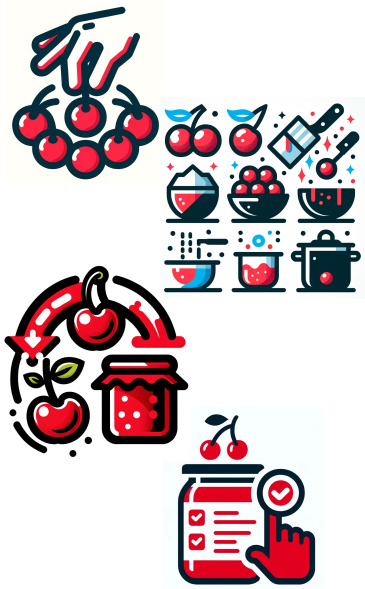⇒ What legitimacy? What transparency? What effectiveness?





PLATE ?

# Machine learning is never neutral

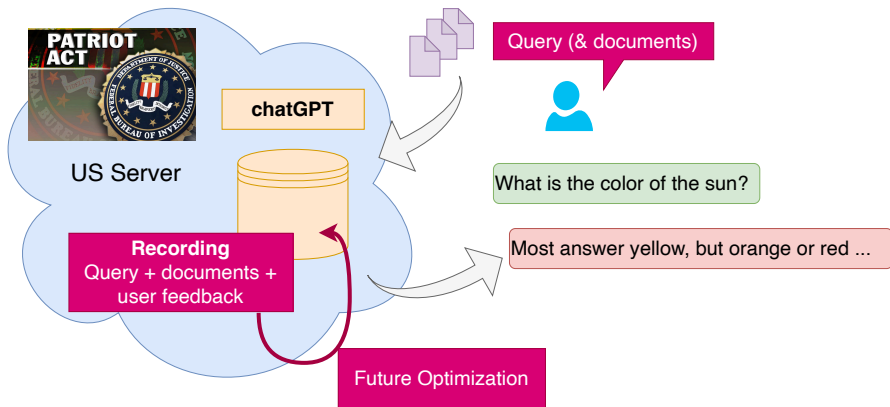1. Data selection
   - Sources, balance, filtering

2. Data transformation
   - Information selection, combination

3. Prior knowledge
   - Balance, loss, a priori, operator choices…

4. Output filtering
   - Post processing

$\Rightarrow$ Choices that influence algorithm results

# Data Leak(s)



- Transfer of sensitive data
- Exploitation of data by OpenAI (or others)
- Data leakage in future models

# Security Issues

- Plug-ins ⇒ Often significant security vulnerabilities for users
  - Email access / transfer of sensitive information etc...

- Management issues for companies
  - Securing (very) large files

- Increased opportunities for malware signatures
  - ≈ software rephrasing

- New problems!
  - Direct malware generation

# Legal Risks/Questions



Reading, collection, formatting

Training model

Trained model = Math function

Inference

Documents, personal data, medicine data, ...

Storage (temporary ou permanent)

Generate commands, diagnostics, texts, image, codes

Copyright and database law

Right to collect, right to copy, consent

Right to use data in an algorithm **Optout**

Model = emanation of data?

Clearview.ai

Cambridge Analytica

Reproductions of untraceable extracts

Usage regulation

Responsibility for errors

# Economic Questions

- Funding/Advertising ⇔ **visits** by internet users
- Google knowledge graph (2012) ⇒ fewer visits, less revenue
- chatGPT = encoding web information... ⇒ much fewer visits?

⇒ What **business model for information sources** with chatGPT?

**Google's Knowledge Graph Boxes:
killing Wikipedia?**

by Gregory Kohs



⇒ Who does **benefit from the feedback**? [StackOverFlow]

# Risks of AI Generalization

AI everywhere =
                loss of meaning?

- In the educational domain
- Transposition to HR
- To project-based funding systems



Writing, reflection, outline, ideas

AI usage verification

Automated evaluation, summary, ...

Quiz

Outline, quiz, illustrations

# How to approach the ethics question?
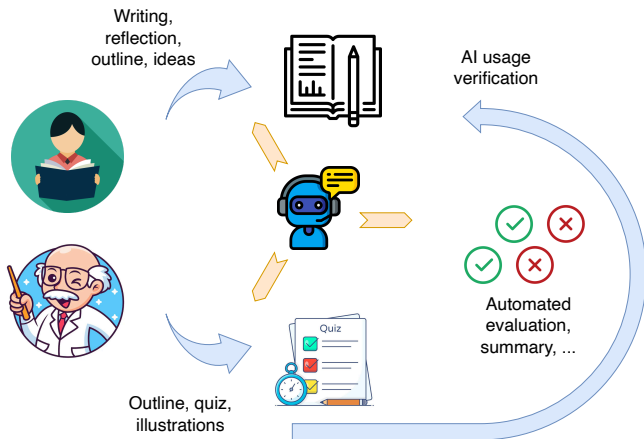
## Medicine

1. **Autonomy:** the patient must be able to make informed decisions.
2. **Beneficence:** obligation to do good, in the interest of patients.
3. **Non-maleficence:** avoid causing harm, assess risks and benefits.
4. **Justice:** fairness in the distribution of health resources and care.
5. **Confidentiality:** confidentiality of patient information.
6. **Truth and transparency:** provide honest, complete, and understandable information.
7. **Informed consent:** obtain the free and informed consent of patients.
8. **Respect for human dignity:** treat all patients with respect and dignity.

## Artificial Intelligence

1. **Autonomy:** Humans control the process
2. **Beneficence:** including the environment?
3. **Non-maleficence:** Humans + environment / sustainability / malicious uses
4. **Justice:** access to AI and equal opportunities
5. **Confidentiality:** what about the Google/Facebook business model?
6. **Truth and transparency:** the tragedy of modern AI
7. **Informed consent:** from cookies to algorithms, knowing when interacting with an AI
8. **Respect for human dignity:**

# How to approach the ethics question?

## Medicine

1. **Autonomy:** the patient must be able to make informed decisions.
2. **Beneficence:** obligation to do good, in the interest of patients.
3. **Non-maleficence:** avoid causing harm, assess risks and benefits.
4. **Justice:** fairness in the distribution of health resources and care.
5. **Confidentiality:** confidentiality of patient information.
6. **Truth and transparency:** provide honest, complete, and understandable information.
7. **Informed consent:** obtain the free and informed consent of patients.
8. **Respect for human dignity:** treat all patients with respect and dignity.

## Artificial Intelligence

1. **Autonomy:** Humans control the process
2. **Beneficence:** including the environment?
3. **Non-maleficence:** Humans + environment / sustainability / malicious uses
4. **Justice:** access to AI and equal opportunities
5. **Confidentiality:** what about the Google/Facebook business model?
6. **Truth and transparency:** the tragedy of modern AI
7. **Informed consent:** from cookies to algorithms, knowing when interacting with an AI
8. **Respect for human dignity:**

# LLM & Conscience

# La conscience (par chatGPT)

1. **Subjectivité** La conscience est intrinsèquement subjective. Chaque individu a sa propre perspective interne, un point de vue unique sur le monde.

2. **Intentionnalité** La conscience est souvent dirigée vers quelque chose : un objet, une pensée, une sensation. Cela signifie qu'elle est intentionnelle, se focalisant sur des éléments spécifiques.

3. **Réflexivité** La conscience permet à un individu de se reconnaître comme étant conscient. C'est la capacité à penser à ses propres pensées, à s'auto-évaluer et à se considérer comme un être distinct.

4. **Unité** Malgré la multiplicité des sensations, pensées et émotions, la conscience tend à les unifier en une seule expérience cohérente.

5. **Continuité** La conscience a un caractère temporel. Elle s'inscrit dans une continuité, reliant le passé, le présent et les projections futures.

6. **Sentience** Il s'agit de la capacité à ressentir des émotions et des sensations. La conscience permet de vivre des expériences plaisantes ou douloureuses.

7. **Libre arbitre** Certains considèrent que la conscience est associée au libre arbitre, c'est-à-dire la capacité de faire des choix délibérés, bien que cela fasse l'objet de débats philosophiques.
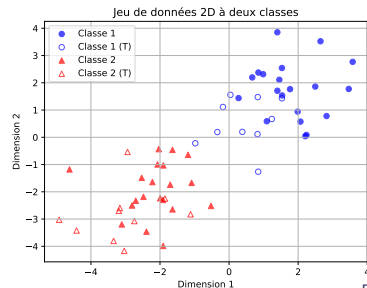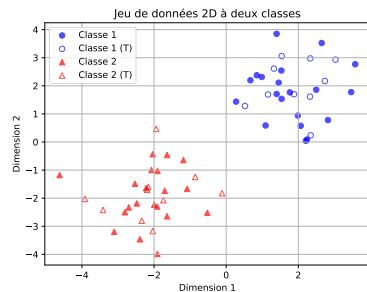
# Généralisation

# Pouvoir de Généralisation

La notion de **généralisation** est centrale en Machine Learning:
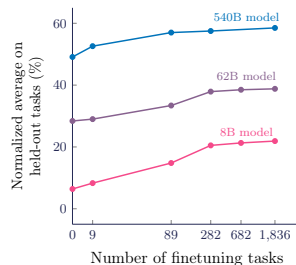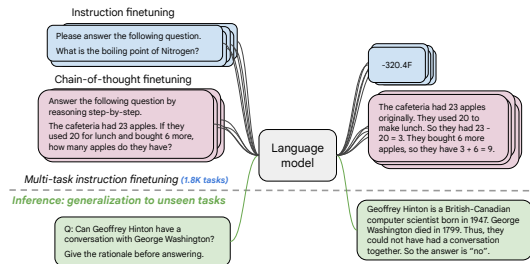


Jeu de données 2D à deux classes

1. Problème iid: indépendant et identiquement distribué
   - Sur-apprentissage, généralisation
   - Data-Augmentation, régularisation
2. Transfert d'apprentissage
   - Dépasser le cas iid, dérive des distributions
3. Multi-tâches, transfert de tâche
   - Apprendre à faire de nouvelles choses



Jeu de données 2D à deux classes

# Les LLM et la généralisation



- Que signifie iid dans les données textuelles?
    - Wikipedia, Reddit, Bioinformatique, Médecine, Finance, …
- Multi-tâche & FLAN
- Du multi-tâche à la multimodalité

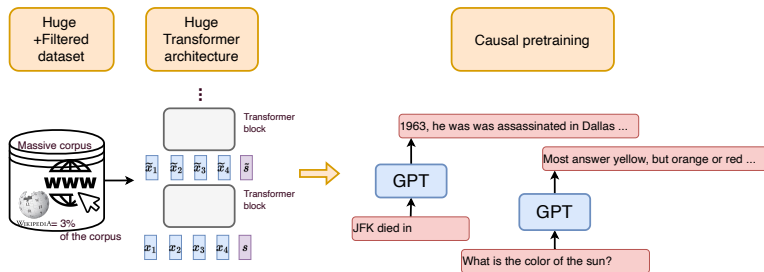# Mémoire
# Connaissances
# et Raisonnement

# Les connaissances paramétriques

**1** Construction



- Vocabulaire

- Grammaire

- Connaissance

Des connaissances imparfaites mais impressionnantes

**2** Mesure: benchmark & métrique

**3** Limites

# Les connaissances paramétriques

1. Construction
2. Mesure: benchmark & métrique
   - QA: Question Answering *HotpotQA; 2WikiMultihopQA; MuSiQue; KQA Pro...*
   - Formattage imposé, Regex, NLI pour la vérification des résultats

**Paragraph A, Return to Olympus:**
[1] *Return to Olympus is the only album by the alternative rock band Malfunkshun.* [2] *It was released after the band had broken up and after lead singer Andrew Wood (later of Mother Love Bone) had died of a drug overdose in 1990.* [3] Stone Gossard, of Pearl Jam, had compiled the songs and released the album on his label, Loosegroove Records.

**Paragraph B, Mother Love Bone:**
[4] *Mother Love Bone was an American rock band that formed in Seattle, Washington in 1987.* [5] The band was active from 1987 to 1990. [6] *Frontman Andrew Wood's personality and compositions helped to catapult the group to the top of the burgeoning late 1980s/early 1990s Seattle music scene.* [7] *Wood died only days before the scheduled release of the band's debut album, "Apple", thus ending the group's hopes of success.* [8] The album was finally released a few months later.

**Q:** What was the former band of the member of Mother Love Bone who died just before the release of "Apple"?
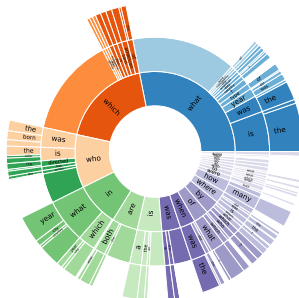**A:** Malfunkshun
**Supporting facts:** 1, 2, 4, 6, 7



Figure 2: Types of questions covered in HOTPOTQA. Question types are extracted heuristically, starting at question words or prepositions preceding them. Empty colored blocks indicate suffixes that are too rare to show individually. See main text for more details.

3. Limites

# Les connaissances paramétriques

1. Construction
2. Mesure: benchmark & métrique
3. Limites
   - Hallucinations
   - Auto-évaluation / confiance problématiques
   - Quid des limites imposées aux LLM (politique etc...)

# Des bases de connaissances aux LLM

## Ontologies

- Stockage (RDF, ...)
- Requêtage (SparQL)
- Raisonnement logique (Prolog, Pellet, Hermit, Elk)

## LLM

- Stockage implicite (paramètres)
- Requêtage en langage naturel mais *instable*
- Raisonnement = mimétisme des schémas vus en apprentissage : puissant mais *imparfait*

**Base de faits:**

Barack Obama est né à Honolulu

Honolulu est la capitale d'Hawaï

⟹ Barack Obama est né à Hawaï

**Base de règles:**
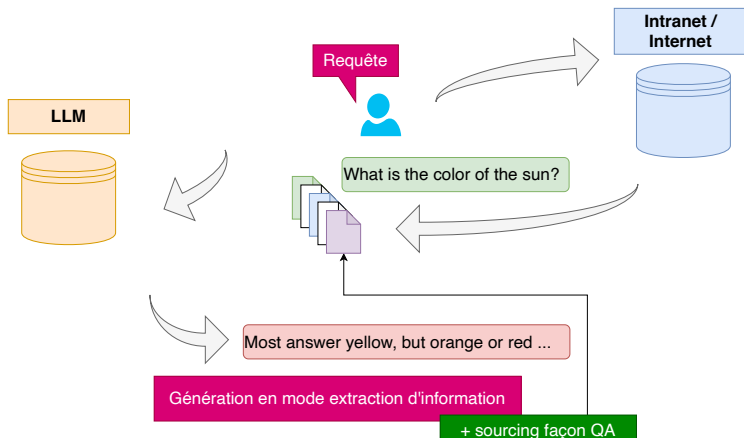
est la capitale

⇓

est inclus dans

**Moteur d'inférence:**

# Couplage: RAG, Toolsformer, Raisonnement

- Chercher dans des documents plutot que dans sa mémoire [RAG]
- Faire appel à des outils externes [calculatrice, Web, appel SQL]
- Apprendre à raisonner
  - Difficile pour un modèle qui ne sait pas faire une opération mathématique
  - ... Mais plus facile quand on sait programmer



**Intranet / Internet**

**Requête**

**LLM**

What is the color of the sun?

Most answer yellow, but orange or red ...

Génération en mode extraction d'information

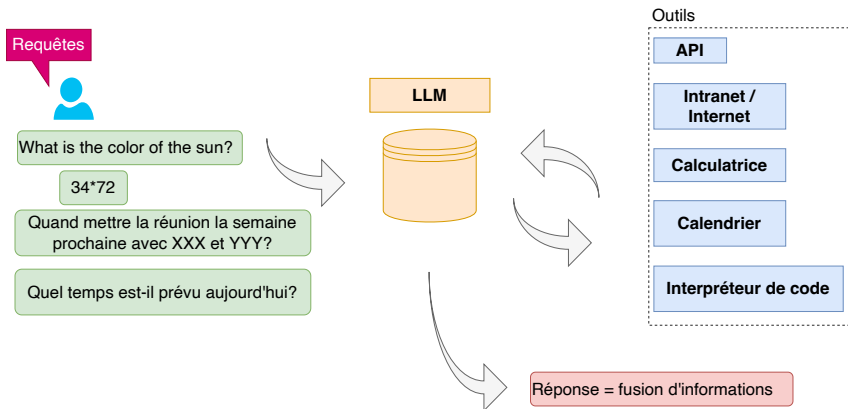+ sourcing façon QA

# Couplage: RAG, Toolsformer, Raisonnement

- Chercher dans des documents plutot que dans sa mémoire [RAG]
- Faire appel à des outils externes [calculatrice, Web, appel SQL]
- Apprendre à raisonner
  - Difficile pour un modèle qui ne sait pas faire une opération mathématique
  - … Mais plus facile quand on sait programmer

# Couplage: RAG, Toolsformer, Raisonnement

- Chercher dans des documents plutot que dans sa mémoire [RAG]
- Faire appel à des outils externes [calculatrice, Web, appel SQL]
- Apprendre à raisonner
    - Difficile pour un modèle qui ne sait pas faire une opération mathématique
    - … Mais plus facile quand on sait programmer

**Task:** Basic Math
**Problem:** Before December, customers buy 1346 ear muffs from the mall. During December, they buy 6444, and there are none. In all, how many ear muffs do the customers buy?

- - - - - - - - - - - - - - - - - - - -

**Predicted Answer:** 1346.0 ✗
**Generated Program:**
```
answer = 1346.0 + 6444.0
print(answer)
# Result ==> 7790.0
```

**Gold Answer**: 7790.0 ✓

**Task:** Muldiv
**Problem:** Tickets to the school play cost 6 for students and 8 for adults. If 20 students and 12 adults bought tickets, how many dollars' worth of tickets were sold?

- - - - - - - - - - - - - - - - - - - -

**Predicted Answer:** 48 ✗
**Generated Program:**
```
a=20*6
b=12*8
c=a+b
answer=c
print(answer)
# Result ==> 216.0
```
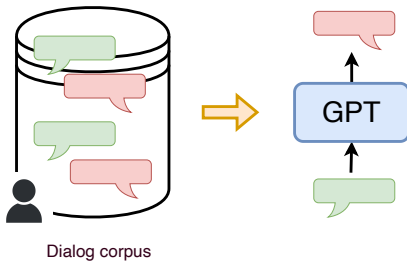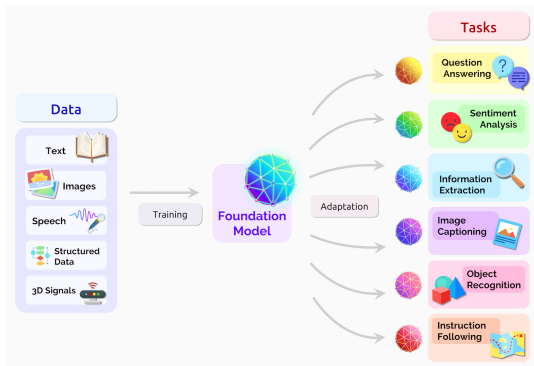
**Gold Answer**: 216 ✓

# Unité et continuité

Deux domaines où les modèles ont le plus progressé... Mais on partait de 0 !

- **Unité** : vers des modèles de fondation
    - Loin de l'universalité (ou même des 5 sens)
- **Continuité**
    - Suivi de dialogue

# Conclusion

- L'intelligence est-elle assimilable à du calcul?
- La logique est-elle indispensable?
- L'apprentissage sans logique est-il raisonnable?
  - Plus de livre qu'un humain n'en lira jamais, plus d'image qu'un humain n'en verra jamais…
  - *vs* esprit analytique
- Il existe d'autre forme d'intelligence que l'intelligence humaine… Mais l'intelligence est-elle la conscience?

# Intentionalité, libre arbitre, créativité

# La conscience et l'intention

Tout ce qui est vivant à des intentions, des buts

- Libre arbitre
- Intentionalité

- Réponse à un prompt
- Suivi des commandes
- Initiatives: aller sur le web chercher une réponse

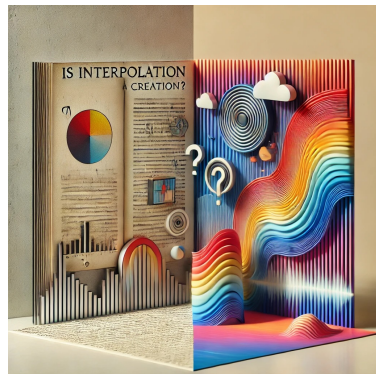## IA Forte / Artificial General Intelligence

- Define Inputs & Outputs
- Break down into subtasks
- Build & test components (processing chain)
- Assert (limited) generalization (iid assumption)
- Performances Evaluation

- Augmented Generalization Capability (Universality)
- Autonomous Learning
  - Data/information access
  - Knowledge extraction (Training+Eval+Confidence/Trust)
- Reasoning
- Conscience, Intentionality

# Créativité

La créativité est-elle menacée par les IA? Nécessite-elle de l'intention?

- L'interpolation entre deux éléments (textes, images, sons, ...) est-elle une création?
- Que se passe-t-il si la base d'interpolation est infinie?
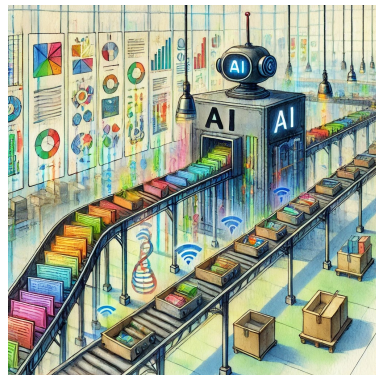- Les IA peuvent-elles apprendre à partir de données générées?



Les textes/images générés en IA sont nouveaux (peu de reprise mot à mot, de portion d'image copiée)

Les problématiques de droit d'auteur sont critiques

# Intentionalité et accès à l'information

- Une IA n'est jamais neutre
  - Choix des données, présence des biais
  - Corrections manuelles, ligne éditoriale
- Un IA n'a pas d'intention... Si ce n'est une fonction objectif à minimiser
  - Comment est choisi cet objectif dans l'accès à l'information?
  - $\Rightarrow$ Max. rétention des utilisateurs
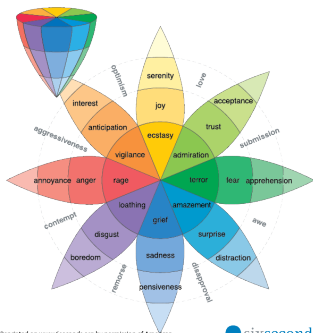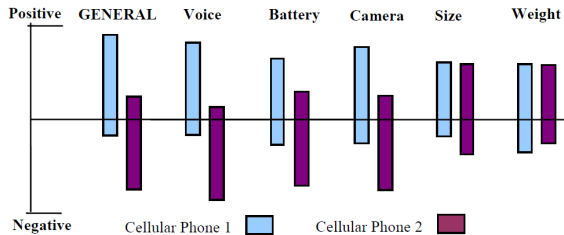  - $\Rightarrow$ Bulles de pensées etc...

# Jugement de valeurs
# Subjectivité

# Le machine learning peut il aborder des tâches subjectives

- Oui, lorsqu'on est capable de lui fournir des étiquettes
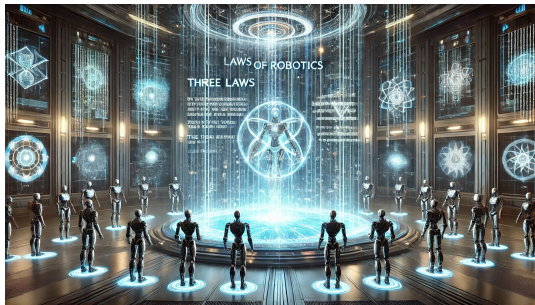⇒ Opinion Mining dans les années 2005-2015

# Bien/Mal, Beau/Laid

Une IA peut-elle emettre un jugement?

- Reproduction de règles vues en apprentissage
- ... Avec extension à des tâches proches
- Beaucoup de valeurs imposées
    - Ligne éditoriale absolument pas autonome

Les 3 lois de la robotiques imposées dans I. Asimov: répétées encore et encore jusqu'à assimilation



1 Un robot ne peut porter atteinte à un être humain ni, restant passif, permettre qu'un être humain soit exposé au danger.

2 Un robot doit obéir aux ordres donnés par les êtres humains, sauf si de tels ordres entrent en contradiction avec la Première Loi.

3 Un robot doit protéger sa propre existence tant que cette protection n'entre pas en contradiction avec la Première ou la Deuxième Loi.

# Mais des usages concrets

- Les IA sont utilisées pour juger:
  - Qualité d'un résumé Automatique
  - Niveau de fluidité d'un texte...

⇒ On utilise des LLM pour ces tâches

**Judging LLM-as-a-Judge
with MT-Bench and Chatbot Arena**

Lianmin Zheng[1,*]　Wei-Lin Chiang[1,*]　Ying Sheng[4*]　Siyuan Zhuang[1]

Zhanghao Wu[1]　Yonghao Zhuang[3]　Zi Lin[2]　Zhuohan Li[1]　Dacheng Li[13]

JUSTICE OR PREJUDICE? 🦹
QUANTIFYING BIASES IN LLM-AS-A-JUDGE

Jiayi Ye[†,*], Yanbo Wang[†,*], Yue Huang[1,*], Dongping Chen[2], Qihui Zhang[3], Nuno Moniz[1],
Tian Gao[4], Werner Geyer[4], Chao Huang[5], Pin-Yu Chen[4], Nitesh V. Chawla[1], Xiangliang Zhang[1,‡]

# Conscience de soi

# L'IA a-t-elle conscience d'elle-même?

**A priori, pas du tout... Mais:**

**Google licencie un ingénieur après sa discussion troublante avec une IA : elle avait peur d'être débranchée**

Par Mathilde Rochefort
Publié le 13 juin 2022 à 11h00

□ **58**



Répétition d'ordres abstraits pour accéder au coeur de la mémoire des LLM

Beaucoup de neurones dont les fonctions ne sont pas établies

# Comment qualifier les deadbots?

1. LLM assimilant les données d'une personne décédée
2. Humain dialoguat avec la personne en question
3. Risque important mais aussi outil pour faire son deuil

**Forum européen de bioéthique**
**Deuil et intelligence artificielle : faut-il avoir peur des «deadbots» ?**

Quel humain pour demain ? dossier ⌄

# Conclusion

1. **Subjectivité** La conscience est intrinsèquement subjective. Chaque individu a sa propre perspective interne, un point de vue unique sur le monde.

2. **Intentionnalité** La conscience est souvent dirigée vers quelque chose : un objet, une pensée, une sensation. Cela signifie qu'elle est intentionnelle, se focalisant sur des éléments spécifiques.

3. **Réflexivité** La conscience permet à un individu de se reconnaître comme étant conscient. C'est la capacité à penser à ses propres pensées, à s'auto-évaluer et à se considérer comme un être distinct.

4. **Unité** Malgré la multiplicité des sensations, pensées et émotions, la conscience tend à les unifier en une seule expérience cohérente.

5. **Continuité** La conscience a un caractère temporel. Elle s'inscrit dans une continuité, reliant le passé, le présent et les projections futures.

6. **Sentience** Il s'agit de la capacité à ressentir des émotions et des sensations. La conscience permet de vivre des expériences plaisantes ou douloureuses.

7. **Libre arbitre** Certains considèrent que la conscience est associée au libre arbitre, c'est-à-dire la capacité de faire des choix délibérés, bien que cela fasse l'objet de débats philosophiques.