

# LES MODÈLES DE LANGUE USAGES ET ENJEUX SOCIÉTAUX

Mercredi 8 Octobre 2025 Séminaire CNRS, ANF-TDM-IA 2025

Vincent Guigue https://vguigue.github.io



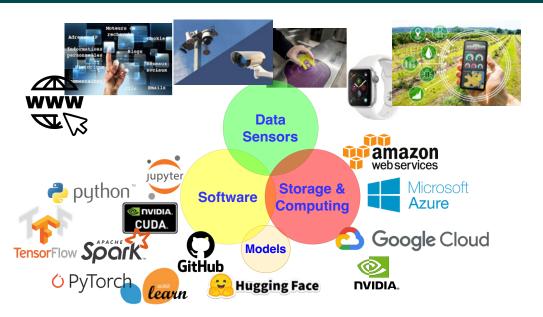
INRA© AgroParisTech



# Introduction



#### The Ingredients of Machine-Learning





#### From tabular data to text

#### ■ Tabular data

- Fixed dimension
- Continuous values
- ⇒ A perfect playground for machine learning

# $f(\underline{\hspace{0.5cm}}) = \operatorname{\mathsf{pred}}_{\hat{y} \in Y}$ $X \quad Y$ Features Supervision

This new iPhone, what a marvel

#### ■ Textual data

- Various length
- Discrete values
- ⇒ Complex for machine learning

An iPhone, What a scam!

Half the price is for the logo

How do we turn this text data into a table?

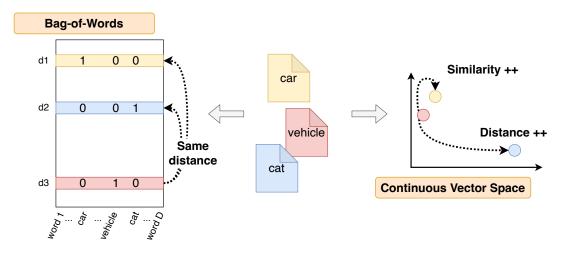
Apple once again proves that perfection can be sold

#### Deep/Representation Learning for Text Data

From Bag of Words to Vector Representations

[2008, 2013, 2016]

Risks



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

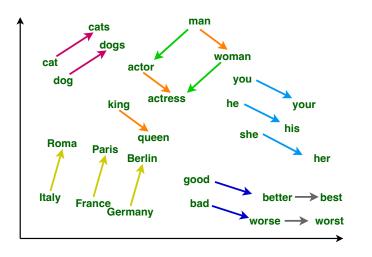
Conclusion

#### Deep/Representation Learning for Text Data

#### From Bag of Words to Vector Representations

[2008, 2013, 2016]

Risks



- Semantic Space:

  similar meanings

  ⇔

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

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

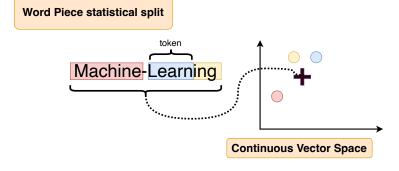
#### Deep/Representation Learning for Text Data

#### From Bag of Words to Vector Representations

[2008, 2013, 2016]

Risks

#### From Words to Tokens



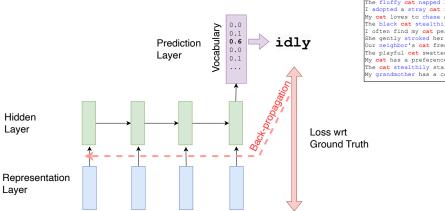
- Representation of unknown words
- Adaptation to technical domains
- Resistance to spelling errors

Enriching word vectors with subword information. Bojanowski et al. TACL 2017.

Introduction ○○○●○○ chatGPT Limits Uses Risks Conclusion

#### Aggregating word representations: towards generative Al

- 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 windowsil, 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.

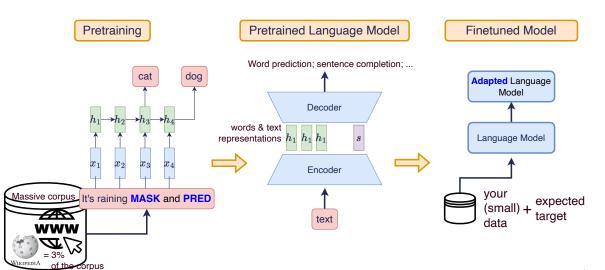
Corpus

The fluffy cat napped lazily in the sunbeam.



#### A new developpement paradigm since 2015

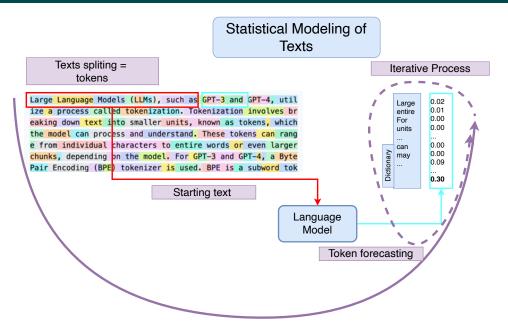
- Huge dataset + huge archi.  $\Rightarrow$  unreasonable training cost
- Pre-trained architecture + 0-shot / finetuning



Conclusion



#### At the end of the day: a stochastic parrot :)



### CHATGPT

NOVEMBER 30, 2022

1 MILLION USERS IN 5 DAYS 100 MILLION BY THE END OF JANUARY 2023

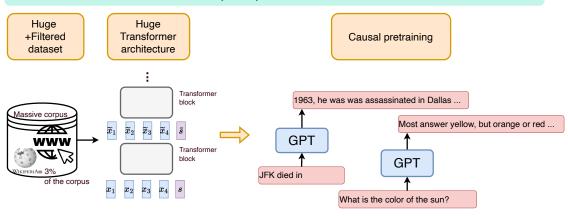
1.16 BILLION BY MARCH 2023

Introduction chatGPT ● ○ ○ ○ Conclusion Limits Uses Risks Conclusion



#### The Ingredients of chatGPT

0. Transformer + massive data (GPT)

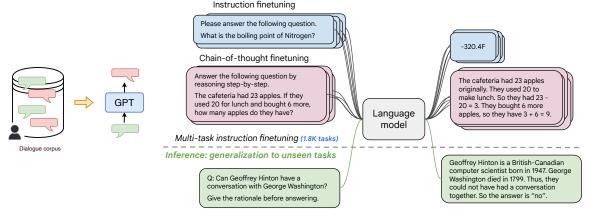


- Grammatical skills: singular/plural agreement, tense concordance
- (Parametric) Knowledge: entities, names, dates, places



#### The Ingredients of chatGPT

#### 1. Dialogue + Tasks



■ Very clean data

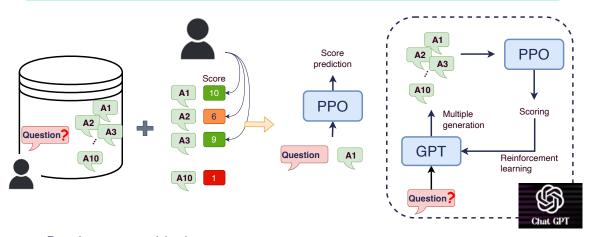
Data generated/validated/ranked by humans

Introduction chatGPT ○ ○ ● ○ Limits Uses Risks Conclusion



#### The Ingredients of chatGPT

#### 3. Instructions + answer ranking



- Database created by humans
- Response improvement

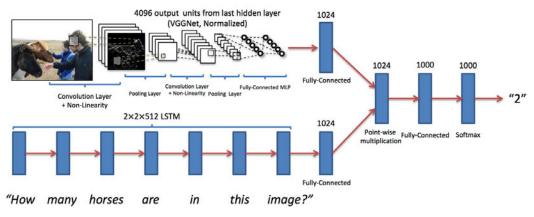
... Also a way to avoid critical topics = censorship



#### GPT4 & Multimodality

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

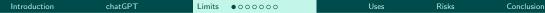
The example of VQA: visual question answering



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



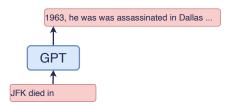
# Machine Learning Limits



#### chatGPT and the relationship with truth

- Likelyhood = grammar, agreement, tense concordance, logical sequences...
  ⇒ Repeated knowledge
- Predict the most plausible word...
  ⇒ produces hallucinations
- 3 Offline functioning
- 4 chatGPT  $\neq$  knowledge graphs
- **5** Brilliant answers...

And silly mistakes! + we cannot predict the errors



#### Example: producing a bibliography



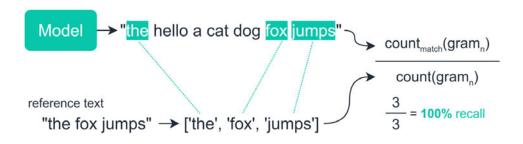
Uses

Conclusion

#### Generative AI: how to evaluate performance?

#### The critical point today

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

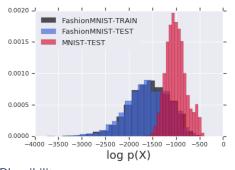


Introduction chatGPT Limits ○ ● ○ ○ ○ ○ ○ Uses Risks Conclusion

### Generative AI: how to evaluate performance?

#### The critical point today

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





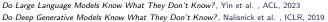


Plausibility

Train

Test

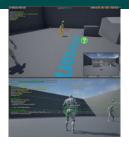


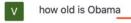




#### Stability/predictability

- Difficult to bound a behavior
- Impossible to predict good/bad answers
- ⇒ Little/no use in video games







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

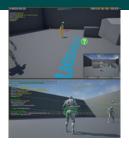




Introduction chatGPT Limits Uses Risks Conclusion 00 00000

### Stability/predictability

- Difficult to bound a behavior
- Impossible to predict good/bad answers
- ⇒ Little/no use in video games



- how old is obama?
- As of 2021, Barack Obama was born on August 4, 1961, so he is 60 years old.



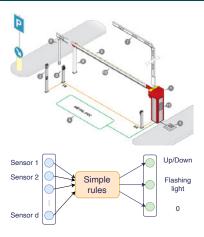


and today?

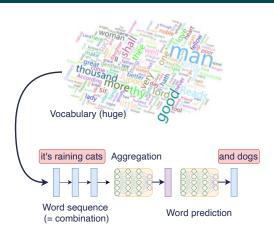




#### Explainability... And complexity



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



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

Introduction chatGPT Limits ○○○●○○○ Uses Risks Conclusion



#### Interpretability vs Post-hoc Explanation

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

too many combinations to anticipate

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



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

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

#### Transparency: open source / open weight

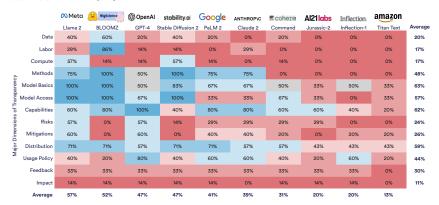
- Can I modify it?
- What training data was used?
- What editorial stance / censorship is involved?
- Why this answer?

Adaptation

Data contamination / skills
Access to information

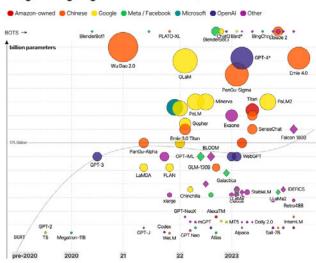
Explainability / interpretability

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



#### Costs / Frugality

### The Rise and Rise of A.I. Size = no. of parameters open-access Large Language Models (LLMs) & their associated bots like ChatGPT



#### # Parameters

1998 LeNet-5 = 0.06M2011 Senna = 7.3M2012 AlexNet = 60M2017 Transformer = 65M / 210M2018 EL Mo = 94M2018 BERT = 110M / 340M2019 GPT2 = 1.500M2020 GPT3 = 175,000M2025 Llama-4 = 2,000,000M



### Everything beyond the LLM's capabilities/training

- Simple calculations (multiplication, division)
- Generating *n*-syllable animal names (in progress)
- Playing chess
- Follow (complex) causal reasoning
- ..

### ATARI 2600 SCORES STUNNING VICTORY OVER CHATGPT



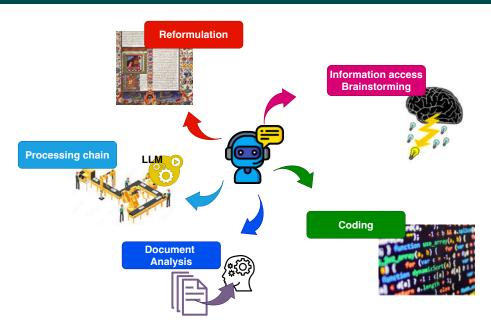
WHEN YOU UNDERESTIMATE A 1977 CHESS ENGINE...
AND IT HUMBLES YOU IN FRONT OF THE WHOLE INTERNET

# Large Language Models

USES



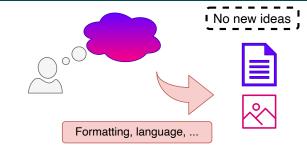
#### Key uses in 5 pictures



Introduction chatGPT Limits Uses ○●○○○○○○ Risks Conclusion



## A fantastic tool for **formatting**



- Personal assistant
  - Standard letters, recommendation letters, cover letters, termination letters
  - Translations
- Meeting reports
  - Formatting notes
- Writing scientific articles
  - Writing ideas, in French, in English

**No new information** ⇒ just writing, improving, translating, cleaning up, ...

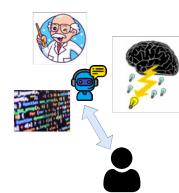
Introduction chatGPT Limits Uses ○ ● ● ○ ○ ○ ○ ○ Risks <u>Conclusion</u>

#### $\overline{(2)}$ Brainstorming / Course Planning / Statistics Review

**■ Find** inspiration

[writer's block syndrome]

- Organize ideas quickly
- Avoid omissions / increase confidency
- Search in a targeted way, adapted to one's needs
- **Answer** student questions (24/7)
- Partner in research, test/enrich ideas
- ⇒ Impressive answers, sometimes incomplete or partially incorrect... But often useful



- In which areas are LLMs reliable?
- What are the risks for primary information sources?
- What societal risks for information?

#### (3) Coding: Different Tools, Different Levels

- Providing solutions to exercises
- Learning to code or getting back into it
  - New languages, new approaches (ML?)
  - Benefit from explanations...

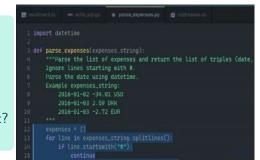
But how to handle mistakes?

- Help with a library [getting started]
- Faster coding
- What about copyrights?
  - What impact on future code processing?
- How to adapt teaching methods?
- How many calls are needed for code completion? What about the carbon footprint?
- What is the risk of error propagation?







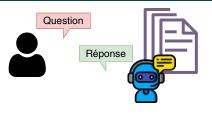


Introduction chatGPT Limits Uses ○○○○●○○○○ Risks Conclusion



#### (4) Document Analysis

- Summarizing documents / articles
- Dialoguing with a document database
- Assistance in writing reviews
- FAQs, internal support services within companies
- Technology watch
- Generating quizzes from lecture notes



NotebookLM

#### Think Smarter, Not Harder

Try NotebookLM

- Will articles still be read in the future?
  - Should we make our articles NotebookLM-proof?
- How to save time while remaining honest and ethical?

Introduction chatGPT Limits Uses ○○○○○ •○○○ Risks <u>Co</u>nclusion

#### $\overline{(5)}$ LLM in a Production Pipeline / Agentic Al

- Run LLM locally
- Extract knowledge
- Generate examples to train a model [Teacher/student - distillation]
- Generate variants of examples <a> → ¬ increase dataset size</a>

[Data augmentation]

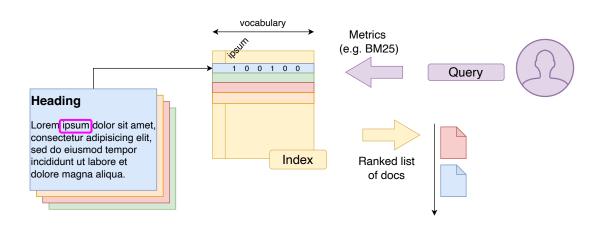
⇒ Integrate the LLM into a processing pipeline = little/less supervision = Agentic AI



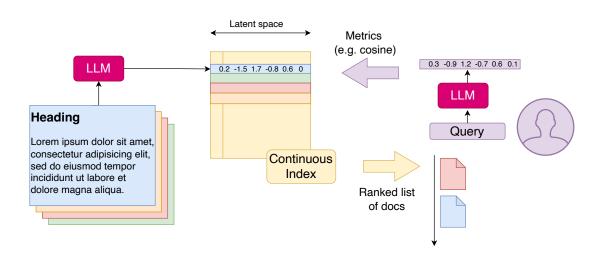
- How much does it cost? (\$ + CO<sub>2</sub>) Need for GPUs?
- How good are open-weight models?
- How to build multiple agents?



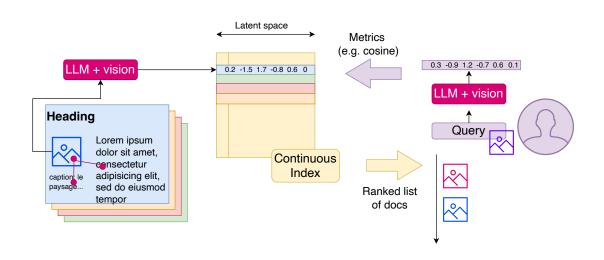
#### LLM vs Information Retrieval



#### LLM vs Information Retrieval

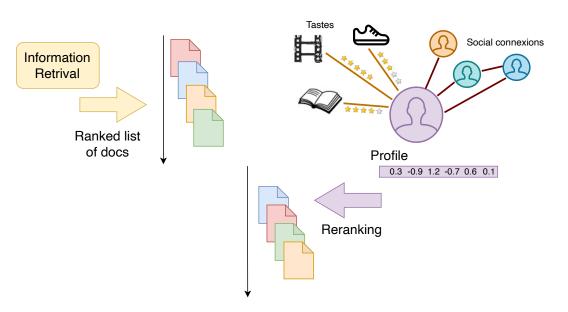


#### LLM vs Information Retrieval





### LLM vs Information Retrieval

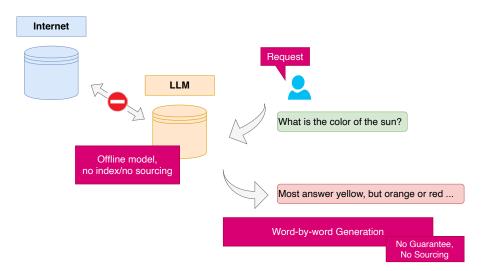




### LLMs $\Rightarrow$ RAG : parametric memory vs Info. Extraction

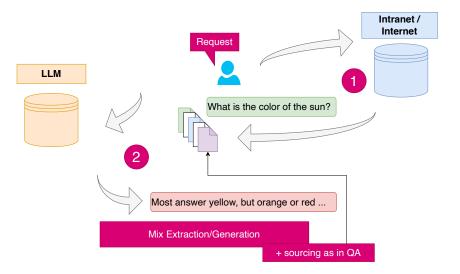
- Asking for information from ChatGPT... A surprising use!
- But is it reasonnable?

[Real Open Question (!)]



Introduction chatGPT Limits Uses ○○○○○○●○ Risks Conclusion

## LLMs $\Rightarrow$ RAG : parametric memory *vs* Info. Extraction

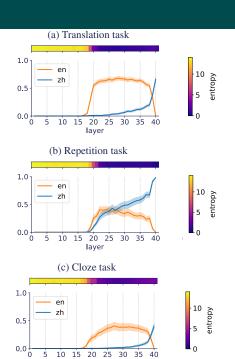


- RAG: Retrieval Augmented Generation
- (Current) limit on input size (2k then 32k tokens)

# Language Handling

- Language models are (mostly) multilingual:
- ⇒ Think in the language you are most comfortable with
- $\Rightarrow$  Ask for answers in the target language

[Wendler et al. 2024] Do Llamas Work in English? On the Latent Language of Multilingual Transformers



# RISKS



duction

chatGPT

Limits

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



Discrimination, exclusion and toxicity

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



Malicious

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



Information hazards

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



Human-computer interaction harms

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



Misinformation harms

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



Automation, access and environmental harms

Harms that arise from environmental or downstream economic impacts of the language model Introduction chatGPT Limits Uses Risks 0 0000000000 Conclusion



### Access to Information

- Access to dangerous/forbidden information
  - +Personal data
  - Right to be forgotten (GDPR)
- 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







Introduction chatGPT Limits Uses Risks ○ ● ○ ○ ○ ○ ○ ○ Conclusion

# Machine Learning & Bias



Mustache, Triangular Ears, Fur Texture

Cat



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

Senior Executive

### Bias in the data $\Rightarrow$ bias in the responses

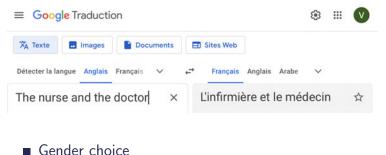
Machine learning is based on extracting statistical biases...

⇒ Fighting bias = manually adjusting the algorithm

### Machine Learning & Bias



Sterreotypes from Pleated Jeans



- . . . .
- Skin color
- Posture
- . . . .

### Bias in the data $\Rightarrow$ bias in the responses

Machine learning is based on extracting statistical biases...

 $\Rightarrow$  Fighting bias = manually adjusting the algorithm

Introduction chatGPT Limits Uses Risks 000 ●00000000 Conclusion

### 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?









## Machine learning is never neutral

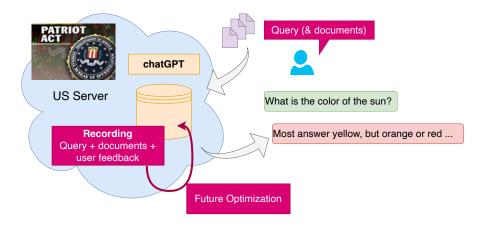
- Data selection
  - Sources, balance, filtering
- Data transformation
  - Information selection, combination
- 3 Prior knowledge
  - Balance, loss, a priori, operator choices...
- 4 Output filtering
  - Post processing
  - Censorship, redirection, ...
- ⇒ Choices that influence algorithm results



Introduction chatGPT Limits Uses Risks 0000 000000 Conclusion



## Data Leak(s): different security levels



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



# Data Leak(s): different security levels

Level 1:	Variable licenses (depending on the companies and
Commercial tools,	subject to change over time). Uncertain data protec-
free to use	tion, risk to personal data.
	chatGPT, Mistral, Perplexity,
Level 2:	Strong contractual guarantees. Risks associated with
Commercial tools,	the Patriot Act. Possible to enforce non-storage of
paid licence	queries.
	chatGPT, Mistral, Perplexity,
Level 3:	+ Negotiation on the server location/data security.
Local dev., Commercial	Microsoft Azur, Mistral, AWS, Aristote, Ragarenn
tools & paid licence ++	
Level 4:	Use of a locally operated LLM, with no data trans-
Local use	ferred over the web.
	HuggingFace, Ollama,

### 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
  - $\blacksquare \approx \text{software rephrasing}$
- New problems!
  - Direct malware generation







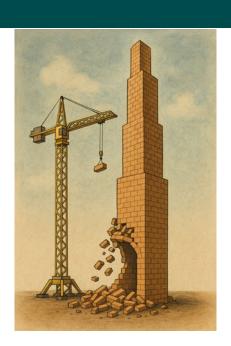


### Educational Challenges

- Redefine our educational priorities, subject by subject, as we did with Wikipedia/calculator/...
  - Accept the decline of certain skills
- Train students in the use of LLMs, while managing to temporarily prohibit their use



■ Learn to recognize LLM-generated content, use detection tools.



 Introduction
 chatGPT
 Limits
 Uses
 Risks ○○○○○○○●○○○○
 Conclusion

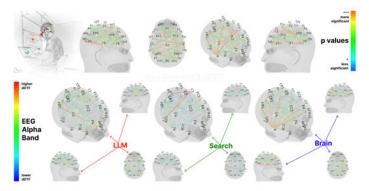
## Decline / Evolution of Cognitive skills

Our brain will evolve with these new tools...

What is the scope of these transformations? What will be the consequences?

■ Education sciences and psychology had conjectured it...

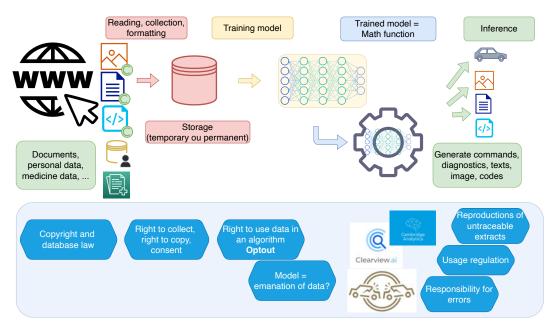
cognitive sciences have measured it



Introduction chatGPT Limits Uses Risks ○○○○○○○○ Conclusion



## Legal Risks/Questions

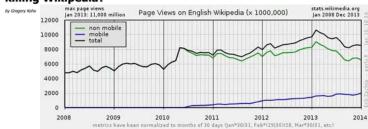




### **Economic Questions**

- Funding/Advertising ⇔ **visits** by internet users
- Google knowledge graph (2012)  $\Rightarrow$  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?



⇒ Who does benefit from the feedback? [StackOverFlow]

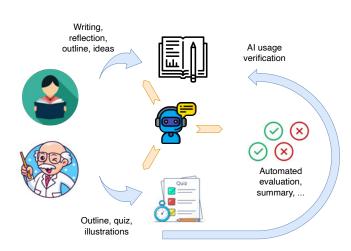
Introduction chatGPT Limits Uses Risks 00000000000 Conclusion Conclusion



### Risks of Al Generalization

 $\label{eq:Al_everywhere} \mbox{Al everywhere} = \\ \mbox{loss of meaning?}$ 

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







## How to approach the ethics question?

### Medicine

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

### **Artificial Intelligence**

- 1 Autonomy: Humans control the process
- **Beneficence:** in the interest of whom? User + GAFAM...
- **3 Non-maleficence:** Humans + environment / sustainability / malicious uses
- 4 Equality: access to Al and equal opportunities
- **5 Confidentiality:** what about the Google/Facebook business model?
- Truth and transparency: the tragedy of modern AI
- 7 Informed consent: from cookies to algorithms, knowing when interacting with an Al
- **Respect for human dignity:** harassment behavior/ human-machine distinction



## How to approach the ethics question?

### Medicine

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

### **Artificial Intelligence**

- 1 Autonomy: Humans control the process
- **Beneficence:** in the interest of whom? User + GAFAM
- **Non-maleficence:** Humans + environment / sustainability / malicious uses
- **Equality:** access to AI and equal opportunities
- 5 Confidentiality: what about the Google/Facebook business model?
- **Truth and transparency:** the tragedy of modern Al
- **Informed consent:** from cookies to algorithms, knowing when interacting with an AI
- Respect for human dignity: harassment behavior/ human-machine distinction

# CONCLUSION

Introduction chatGPT Limits Uses Risks Conclusion • • • • • • • •

## Upcoming Challenges

### ■ What about hallucinations?

- Should we try to reduce them or learn to live with them?
- Will LLMs improve? In what directions?
- Do LLMs make us *lose* our connection to truth? To verification?

### ■ Do we need small or large language models?

- How much does it cost? Is it sustainable?
- With or without fine-tuning?
- What does frugality mean in the world of LLMs?

### ■ When others use them... What impact does it have on me?

- Productivity (fellow researchers, coders, reviewers, ...)
- Education: managing/training *tech-savvy* students

### ■ Data protection... Mine and others'

- Is it reasonable to train LLMs on GitHub, Wikipedia, scientific papers, news outlets, etc.?
- How important is privacy? What are the risks when using an LLM?

Introduction chatGPT Limits Uses Risks Conclusion • • • • • • • •

# Intro

## Upcoming Challenges

### ■ What about hallucinations?

- Should we try to reduce them or learn to live with them?
- Will LLMs improve? In what directions?
- Do LLMs make us *lose* our connection to truth? To verification?
- Do we need small or large language models?
  - The smartphone has made me an augmented human...
  - Will the LLM make me an *augmented researcher*?
  - ⇒ Still. have a look at NotebookLM
- When others use them... vvnat impact uoes it have on me:
  - Productivity (fellow researchers, coders, reviewers, ...)
  - Education: managing/training *tech-savvy* students
- Data protection... Mine and others'
  - Is it reasonable to train LLMs on GitHub, Wikipedia, scientific papers, news outlets, etc.?
  - How important is privacy? What are the risks when using an LLM?

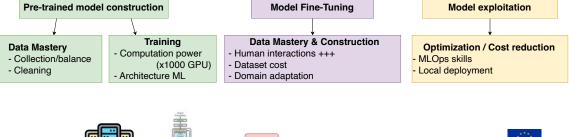


### Levels of Access to Artificial Intelligence

- 1 User via an interface: *chatGPT* 
  - Some training is still required (2-4h)
- Using Python libraries
  - Basics on protocols
  - Standard processing chains
  - Training: 1 week-3 months (ML/DL)
- 3 Tool developer
  - Adapt tools to a specific case
  - Integrate business constraints
  - Build hybrid systems (mechanistic/symbolic)
  - Mix text and images
  - Training:  $\geq 1$  year



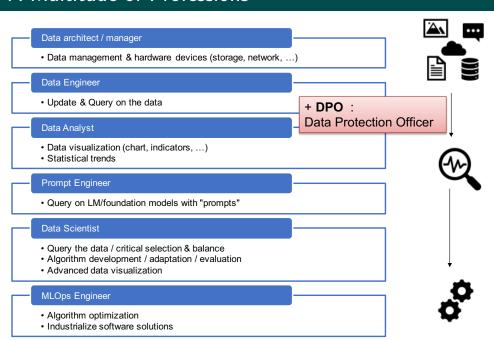
## Digital Sovereignty: the Entire Chain





Introduction chatGPT Limits Uses Risks Conclusion 000 • 00

### A Multitude of Professions



 Introduction
 chatGPT
 Limits
 Uses
 Risks
 Conclusion
 ○ ○ ○ ○ ●



## Factors of Acceptability for Generative Al

### Utilitarianism:

- Performance (acceptance factor of chatGPT)
- Reliability / Self-assessment

### Non-dangerousness:

- Bias / Correction
- Transparency (editorial line, human/machine confusion)
- Reliable Implementation
- Sovereignty (?)
- Regulation (Al act)
  - Avoid dangerous applications

### 3 Know-how:

■ Training (usage/development)



Introduction chatGPT Limits Uses Risks Conclusion 00000



## 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?





