

FONCTIONNEMENT DES MODÈLES DE LANGUE EXPLOITATION SUR DONNÉES ALIMENTAIRES

Lundi 29 septembre 2025 Séminaire ALIMining, IRIT, Toulouse

Vincent Guigue https://vguigue.github.io

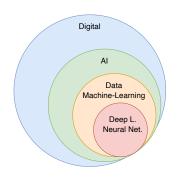


INRA© AgroParisTech



Introduction

Artificial Intelligence & Machine Learning



Input (x)	Output (Y)	Application
email>	spam? (0/1)	spam filtering
audio	text transcript	speech recognition
English	Chinese	machine translation
ad, user info>	click? (0/1)	online advertising
image, radar info 🛶	position of other cars	self-driving car
image of phone -	defect? (0/1)	visual inspection

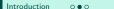
Al: computer programs that engage in tasks which are, for now, performed more satisfactorily by human beings because they require high-level mental processes.

Marvin Lee Minsky, 1956

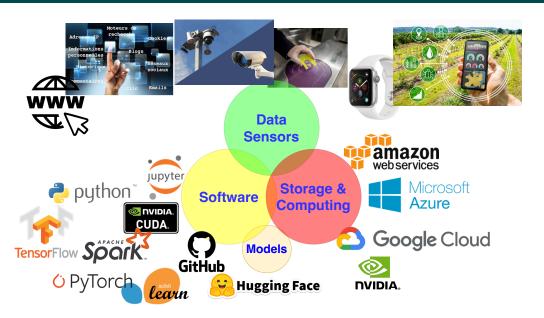
N-AI (Narrow Artificial Intelligence), dedicated to a single task

≠ G-AI (General AI), which replaces humans in complex systems.

Andrew Ng, 2015



The Ingredients of Machine-Learning



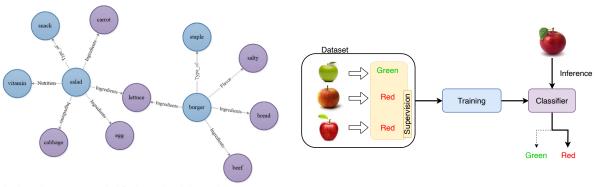
Uses



Machine-Learning vs Expert Knowledge

Modeling Expert Knowledge

Machine Learning



A relationship extraction method for domain knowledge graph construction, Yu et al. 2020

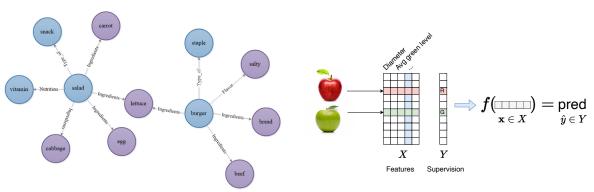
Different behaviors:

different strengths and weaknesses, different costs & requirements

Machine-Learning vs Expert Knowledge

Modeling Expert Knowledge

Machine Learning



A relationship extraction method for domain knowledge graph construction, Yu et al. $2020\,$

Different behaviors:

different strengths and weaknesses, different costs & requirements

[APPLICATION TO TEXTUAL DATA]

Representation Learning

Deep learning &

chatGPT

From tabular data to text

- → Tabular data
 - → Fixed dimension
 - → Continuous values





→ f(□□□□) = pred

- → Textual data
 - → Variable length
 - → Discrete values

this new iPhone, what a marvel

An iPhone? What a scam!



$\overline{\mathsf{AI}} + \mathsf{Textual}$ Data: Natural Language Processing (NLP)

NLP = largest scientific community in Al

Linguistics [1960-2010]

Rule-based Systems:

```
{like, love, appreciate} \rightarrow * \rightarrow * #product appreciate}

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

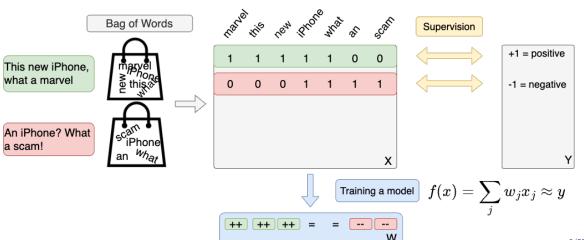
* \rightarrow * {hate, loathe, detest} #product
```

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

$\overline{\mathsf{AI} + \mathsf{Textual}}$ Data: Natural Language Processing (NLP)

NLP = largest scientific community in Al

Machine Learning [1990-2015]





AI + Textual Data: Natural Language Processing (NLP)

NLP = largest scientific community in Al

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

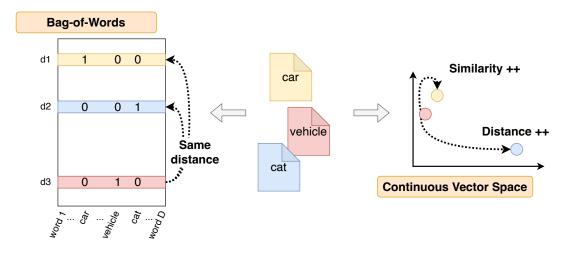
Precision = criterion for acceptance by industry

 \rightarrow Link to metrics



From Bag of Words to Vector Representations

[2008, 2013, 2016]



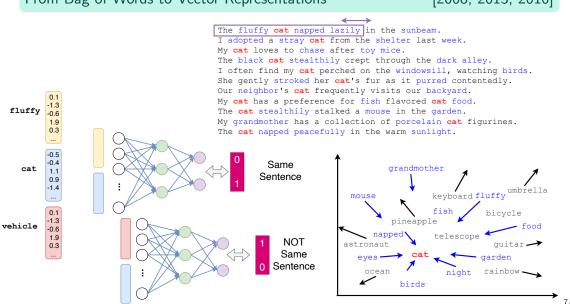
LeCun, Y., Bengio, Y., Hinton, G. (2015). Deep learning. Nature, 521(7553), 436-444.

Introduction Deep learning & NLP chatGPT Limits Uses Conclusion 000000

Deep/Representation Learning for Text Data

From Bag of Words to Vector Representations

[2008, 2013, 2016]

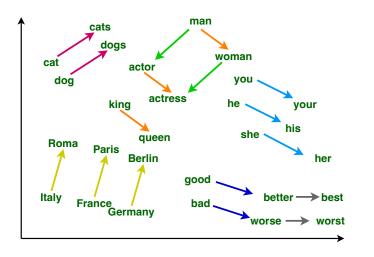


Deep/Representation Learning for Text Data

From Bag of Words to Vector Representations

[2008, 2013, 2016]

Uses



- 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

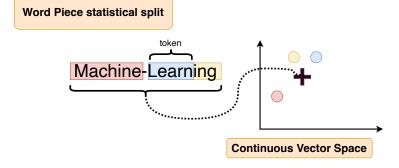


Deep/Representation Learning for Text Data

From Bag of Words to Vector Representations

[2008, 2013, 2016]

From Words to Tokens



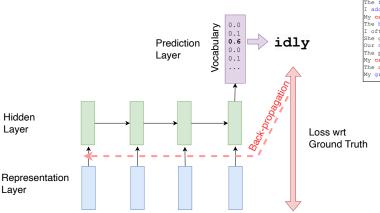
- Representation of unknown words
- Adaptation to technical domains
- Resistance to spelling errors

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

Introduction Deep learning & NLP 000 000 chatGPT Limits Uses 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.

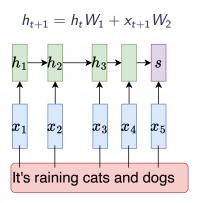
My grandmother has a collection of porcelain cat figurines.

Corpus

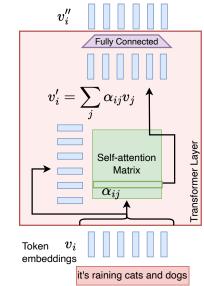
The fluffy cat napped lazily in the sunbeam.

Transformer architecture: state-of-the-art aggregation

Recurrent Neural Network:



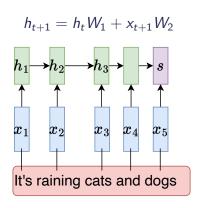
Transformer:



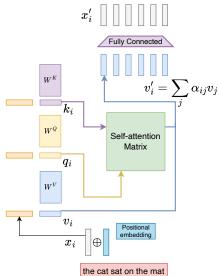
Attention is all you need, Vaswani et al. NeurIPS 2017

Transformer architecture: state-of-the-art aggregation

Recurrent Neural Network:



Transformer:



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

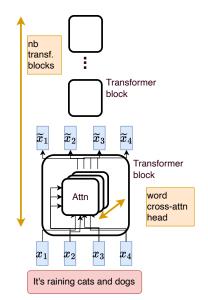
Sequence to Sequence Learning with Neural Networks, Sutskever et al. NeurIPS 2014

Transformer architecture: state-of-the-art aggregation

Recurrent Neural Network:

$h_{t+1} = h_t W_1 + x_{t+1} W_2$ $h_1 \rightarrow h_2 \rightarrow h_3 \rightarrow s$ $x_1 \quad x_2 \quad x_3 \quad x_4 \quad x_5$ It's raining cats and dogs

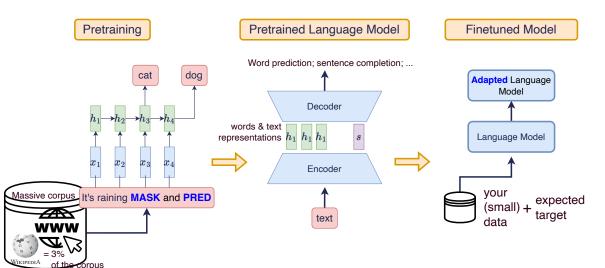
Transformer:



Introduct

A new developpement paradigm since 2015

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



CHATGPT

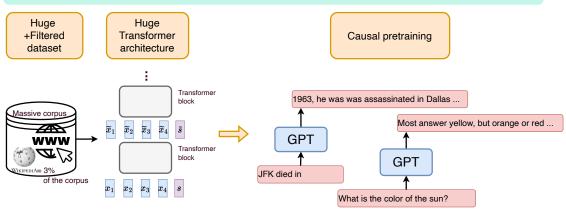
NOVEMBER 30, <u>2022</u>

1 MILLION USERS IN 5 DAYS 100 MILLION BY THE END OF JANUARY 2023 1.16 BILLION BY MARCH 2023 Introduction Deep learning & NLP chatGPT ●○○○○○○○ Limits Uses 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. More is better! (GPT)

+ more input words [500 \Rightarrow 2k, 32k, 100k]

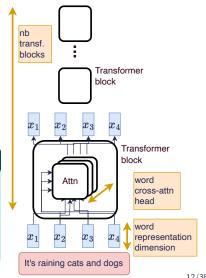
+ more dimensions in the word space $[500-2k \Rightarrow 12k]$

+ more attention heads $[12 \Rightarrow 96]$

+ more blocks/layers [5-12 \Rightarrow 96]

175 Billion parameters... What does it mean?

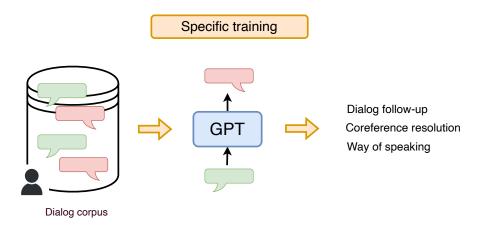
- $1.75 \cdot 10^{11} \Rightarrow 300 \text{ GB} + 100 \text{ GB}$ (data storage for inference) $\approx 400 \text{GB}$
- NVidia A100 GPU = 80GB of memory (=20k€)
- Cost for (1) training: 4.6 Million €





The Ingredients of chatGPT

2. Dialogue Tracking



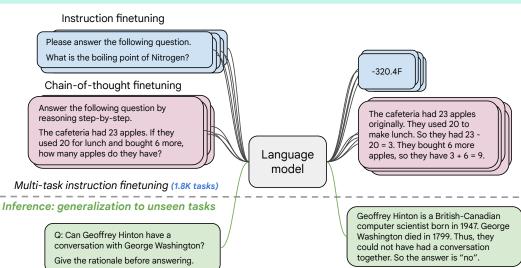
■ Very clean data

Data generated/validated/ranked by humans

Introduction Deep learning & NLP chatGPT ○○○●○○○○○ Limits Uses Conclusion

The Ingredients of chatGPT

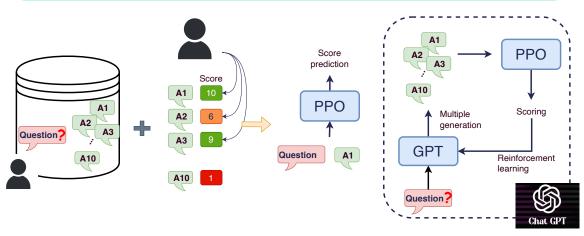
3. Fine-tuning on different (\pm) complex reasoning tasks



Introduction Deep learning & NLP chatGPT ○○○○●○○○○ Limits Uses Conclusion

The Ingredients of chatGPT

4. Instructions + answer ranking



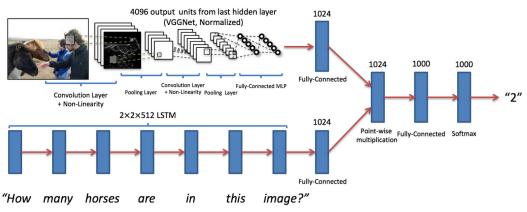
- Database created by humans
- Response improvement

... 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

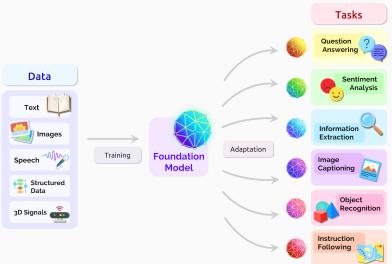


 \Rightarrow Backpropagate the error \Rightarrow modify word representations + image analysis



Towards Larger Foundation Models?

■ Let the modalities enrich each other





Introduction Deep learning & NLP

Why So Much Controversy?

[December 2022] New tool

+ 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?

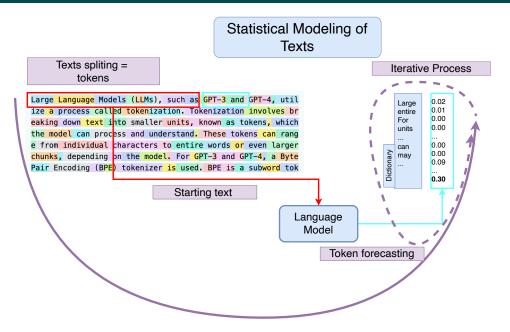








At the end of the day



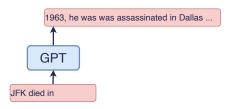
Machine Learning Limits



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



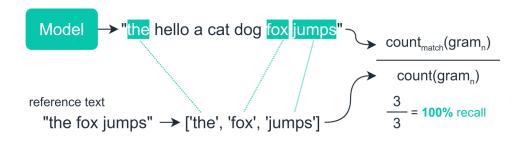
Conclusion

Generative Al: how to evaluate performance?

The critical point today

Limits

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

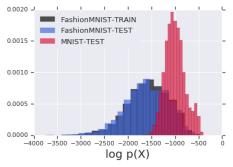


 Introduction
 Deep learning & NLP
 chatGPT
 Limits
 ○ ● ○ ○ ○ ○ ○
 Uses
 Conclusion

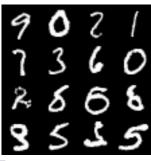
Generative AI: how to evaluate performance?

The critical point today

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





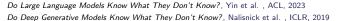


Plausibility





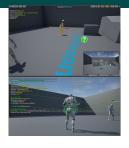


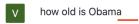




Stability/predictability

- 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.

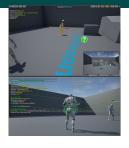




Introduction Deep learning & NLP chatGPT Limits Uses Conclusion 000000

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.

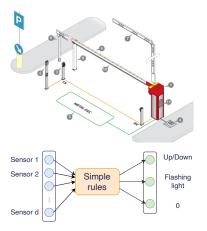




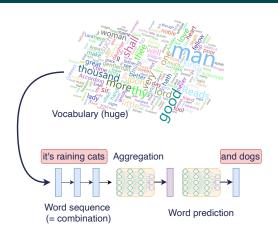
and today?

Introduction Deep learning & NLP chatGPT Limits ○○○●○○○ Uses Conclusion

Explainability... And complexity



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



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

Introduction Deep learning & NLP chatGPT Limits ○○○●○○○ Uses Conclusion

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)



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

[Uber Accident, 2018]

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

Introducti

Transparency: open source / open weight

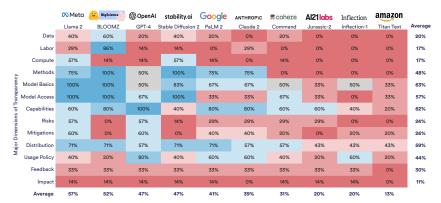
- Can I modify it?
- What training data was used?
- What editorial stance / censorship is involved?
- Why this answer?

Adaptation

Data contamination / skills
Access to information

Explainability / interpretability

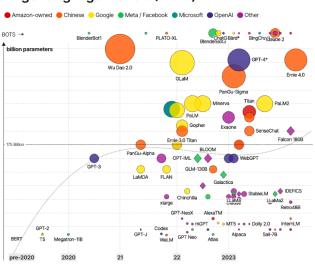
Foundation Model Transparency Index Scores by Major Dimensions of Transparency, 2023
Source: 2023 Foundation Model Transparency Index



Costs

$\overline{\mathsf{Costs}} / \mathsf{Frugality}$

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



Parameters

1998 LeNet-5

2011 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

= 0.06M



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



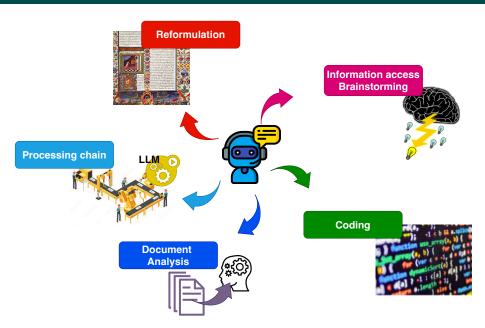
WHEN YOU UNDERESTIMATE A 1977 CHESS ENGINE... AND IT HUMBLES YOU IN FRONT OF THE WHOLE INTERNET

Large Language Models

[IN NUTRITION RESEARCH]

USES

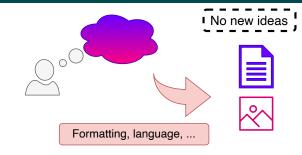
Key uses in 5 pictures





$\left(1 ight)$ Formatting information

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 writting, cleaning up, ...

Introd

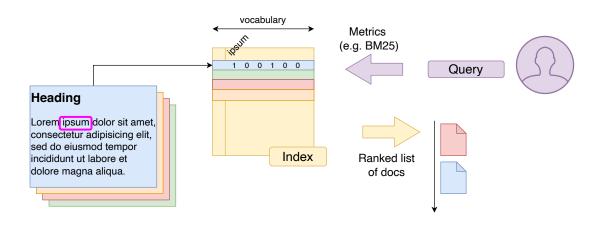
(1) $\mathsf{Nutrition}$ use : Input standardization (?)

 \Rightarrow opportunity to fuse heterogeneous information



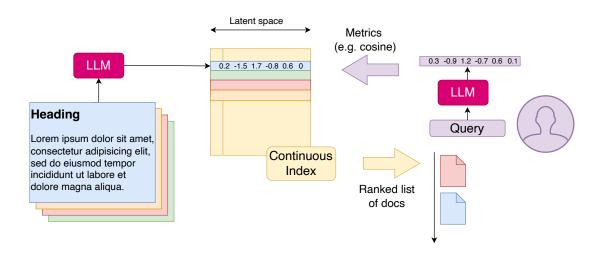


(1) Chat & RAG: a new way to access information

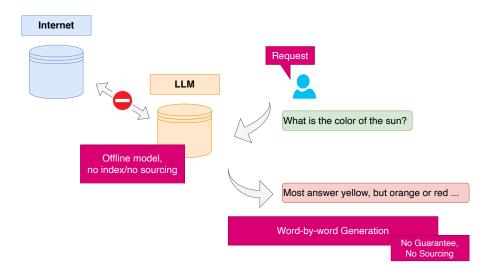


(1) Chat & RAG: a new way to access information

Limits

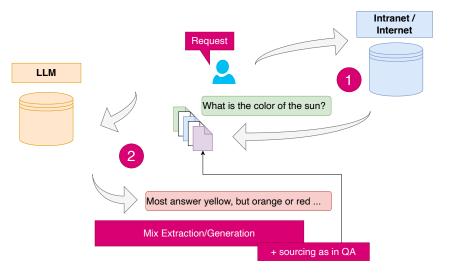


$\overline{(1)}$ Chat & RAG : a new way to access information



Introduction Deep learning & NLP chatGPT Limits Uses ○○○●○○○○○○ Conclusion

(1) Chat & RAG : a new way to access information



- ⇒ A way to build a *reliable* chatbot to advise users?
 - Parametric memory vs Information Retrieval : opposite objectives?

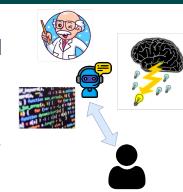
Introduction Deep learning & NLP chatGPT Limits Uses 0000 000000 Conclusion

(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
- ⇒ Impressive answers, sometimes incomplete or partially incorrect... But often useful



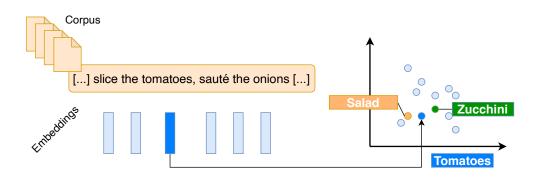
3 reference articles on the use of transformers in recommendation systems
What is the purpose of the log-normal Poisson law?
Propose 10 sections for a course on Transformers in Al

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



(2) Internal knowledge exploitation for nutrition

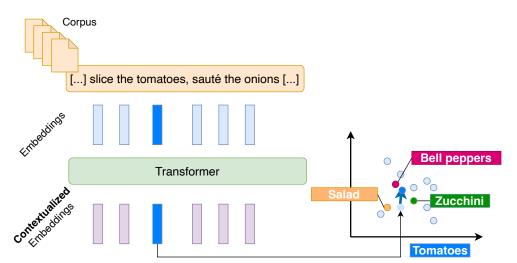
- Brainstorming in the kitchen: which application for cooking?
- Ingredient substitution... At every scale: Ingredient, Food, Dish



Introduction Deep learning & NLP chatGPT Limits Uses ○○○○●○○○○ Conclusion

(2) Internal knowledge exploitation for nutrition

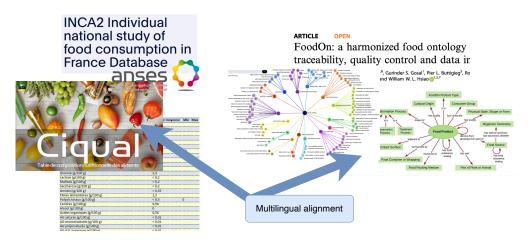
- Brainstorming in the kitchen: which application for cooking?
- Ingredient substitution... At every scale: Ingredient, Food, Dish
- \blacksquare ++ Upgrade by contextualization



Introduction Deep learning & NLP chatGPT Limits Uses ○○○○○●○○○○○ Conclusion

(2) Internal knowledge exploitation for nutrition

- Brainstorming in the kitchen: which application for cooking?
- Ingredient substitution... At every scale: Ingredient, Food, Dish
- ++ Upgrade by contextualization
- Interoperability and ontologies



Introduction Deep learning & NLP chatGPT Limits Uses ○○○○● Conclusion

(2) Internal knowledge exploitation for nutrition

- Brainstorming in the kitchen: which application for cooking?
- Ingredient substitution... At every scale: Ingredient, Food, Dish
- ++ Upgrade by contextualization
- Interoperability and ontologies



rol and data in S. Gosal¹, Pier L. Buttigieg³, Ro W. L. Hsiao (1)^{2,7}



Glacous (g/100 g)
Listinos (g/100 g)
Listinos (g/100 g)
Sistinos (g/100 g)
Amidios (g/100 g)
Filips Internative (g/100 g)
Filips Internative (g/100 g)
Filips Internative (g/100 g)
Accides organizate (g/100 g)

Patrice Buche, Julien Cufi, Liliana Ibanescu, Alrick Oudot, Magalie weber 12/10/2021

(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...

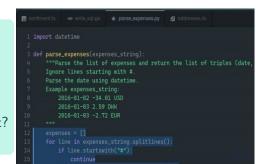
But how to handle mistakes?

- 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?







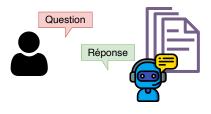


Introduction Deep learning & NLP chatGPT Limits Uses 0000000000 Conclusion

(4) Document Analysis



- 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

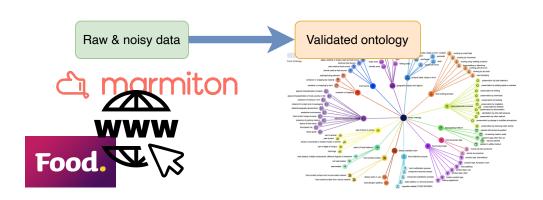
Try NotebookLM

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



(4) Information Extraction in Nutrition

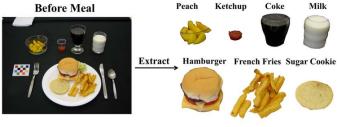
■ Ontology building (mostly textual data)





(4) Information Extraction in Nutrition

- Ontology building (mostly textual data)
- Image analysis



After Meal





- Food recognition
- Segmentation
- Estimation of quantities



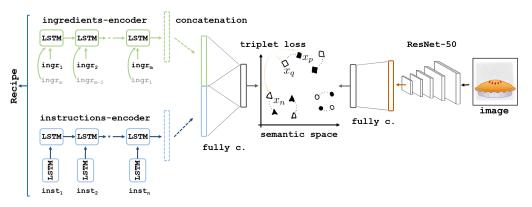
An Overview of The Technology Assisted Dietary Assessment Project at Purdue University., Khanna et al. , 2010

Limits



(4) Information Extraction in Nutrition

- Ontology building (mostly textual data)
- Image analysis
- Multimodal analysis + algorithmic process





Images & Recipes: Retrieval in the cooking context, SIGIR 2018 Carvalho et al.

Introduction Deep learning & NLP chatGPT Limits Uses ○○○○○○○●○○ Conclusion

(4) Information Extraction in Nutrition

- Ontology building (mostly textual data)
- Image analysis
- Multimodal analysis + algorithmic process

ingr (ingredients)

instr (cooking instructions)

image

- 1) pizza dough
- hummus
 arugula
- 4) cherry / grape tomatoes
- 5) pitted greek olives
- crumbled feta cheese

- 1) Cut the dough into two 8-ounce sized pieces.
- 2) Roll the ends under to create round balls.
- Then using a well-floured rolling pin, roll the dough out into 12-inch circles.
- 4) Place the dough circles on sheets of parchment paper.



- 1) unsalted butter
- 3) condensed milk
- 4) sugar
- 5) vanilla extract
- 6) chopped pecans
- 7) chocolate chips

- 1) Preheat the oven to 375 degrees F.
- In a large bowl, whisk together the melted butter and eggs until combined.
- Whisk in the sweetened condensed milk, sugar, vanilla, pecans, chocolate chips, butterscotch chips, and coconut.





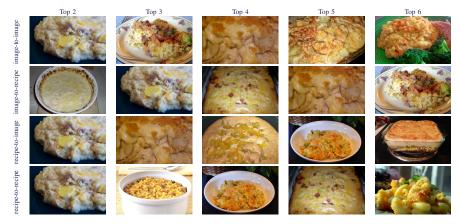
Pecan Pie

Images & Recipes: Retrieval in the cooking context, SIGIR 2018 Carvalho et al.

Introduction Deep learning & NLP chatGPT Limits Uses ○○○○○○○○ Conclusion

(4) Information Extraction in Nutrition

- Ontology building (mostly textual data)
- Image analysis
- Multimodal analysis + algorithmic process





Images & Recipes: Retrieval in the cooking context, SIGIR 2018 Carvalho et al.

Introduction Deep learning & NLP chatGPT Limits Uses ○○○○○○○○●○ Conclusion

(5) LLM in a Production Pipeline / Agentic Al

- Run LLM locally
- Extract knowledge
- Sort documents / generate summaries
- Generate examples to train a model

 [Teacher/student distillation]
- Generate variants of examples // increase dataset size

[Data augmentation]

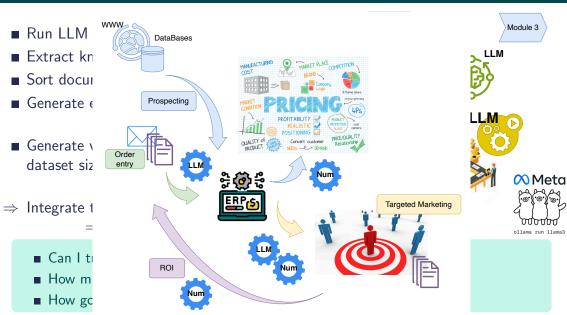
⇒ Integrate the LLM into a processing pipeline
 = little/less supervision = Agentic Al

Module 1 Module 3 Module 2 Meta ollama run llama3

- Can I train models on generated data?
- How much does it cost? (\$ + CO₂) Need for GPUs?
- How good are open-weight models?

/

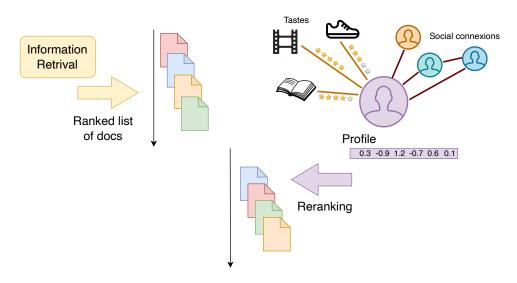
$\overline{(5)}$ LLM in a Production Pipeline / Agentic Al





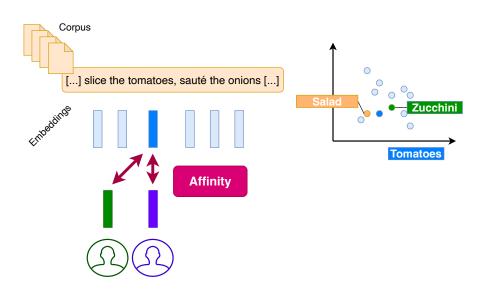
$\overline{(5)}$ What about Recommender System in Nutrition?

Profiling is roughly everywhere in Information Retrieval



(5) What about Recommender System in Nutrition?

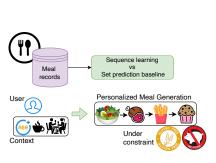
Opportunities in nutrition: modeling user preferences

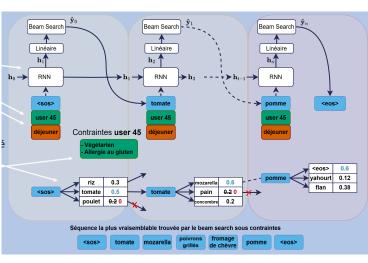


Introduction Deep learning & NLP chatGPT Limits Uses ○○○○○○○○● Conclusion

(5) What about Recommender System in Nutrition?

Building consistent proposals... With expert constraints







Génération séquentielle prenant en compte des informations contextuelles en nutrition , CAp 2025 Combeau et al.

CONCLUSION

Limits

New tools for new opportunities

LLMs offer new perspectives in nutrition:

- A natural and convenient interface for users
 - enabling dialogue, plate analysis, and personalized advice
- Accessible on multiple devices, from computers to smartphones and smart kitchens (Alexa, Google Assistant, ...)
- A means to unify and connect existing nutritional resources
- A powerful tool to extract and structure knowledge ⇒ enrich databases
- A modular component for next-generation recommender systems