

FONCTIONNEMENT DES MODÈLES DE LANGUE EXPLOITATION SUR DONNÉES ALIMENTAIRES

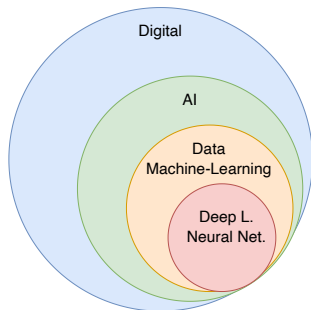
Lundi 29 septembre 2025
Séminaire ALIMining, IRIT, Toulouse

Vincent Guigue
<https://vguigue.github.io>

INTRODUCTION



Artificial Intelligence & Machine Learning



Input (\mathbf{x})	Output (\mathbf{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

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

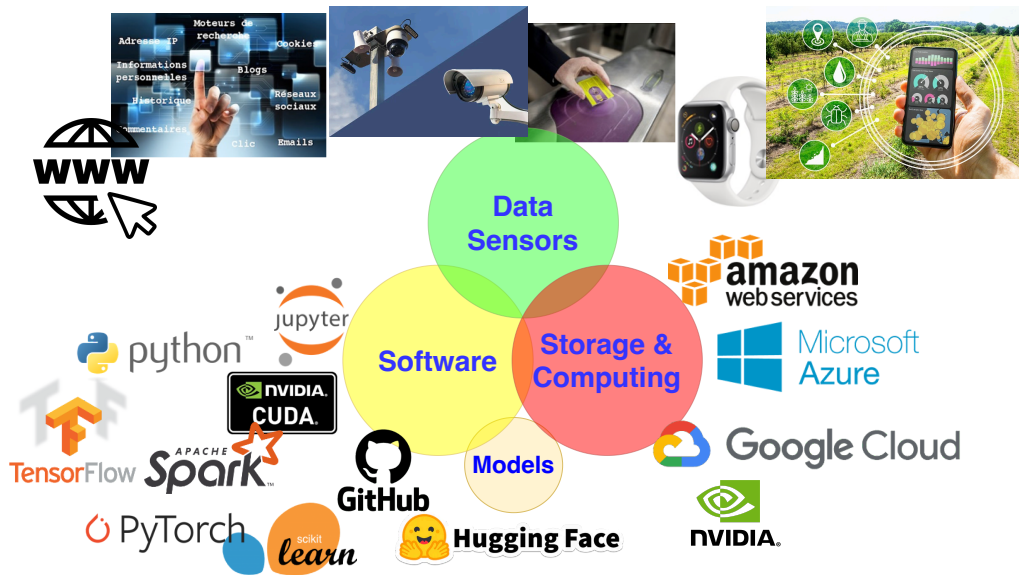
N-AI (Narrow Artificial Intelligence), dedicated to a single task

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

Andrew Ng, 2015



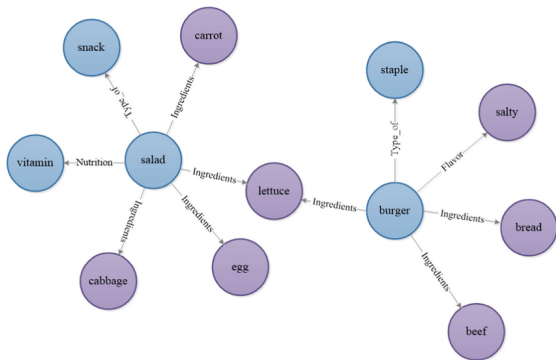
The Ingredients of Machine-Learning





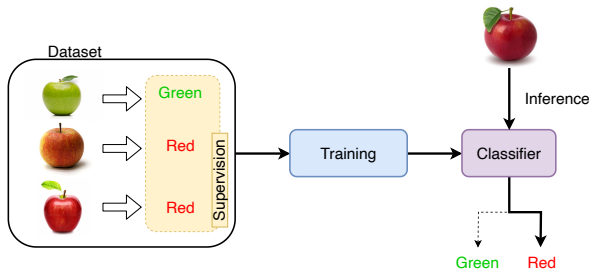
Machine-Learning vs Expert Knowledge

Modeling Expert Knowledge



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

Machine Learning



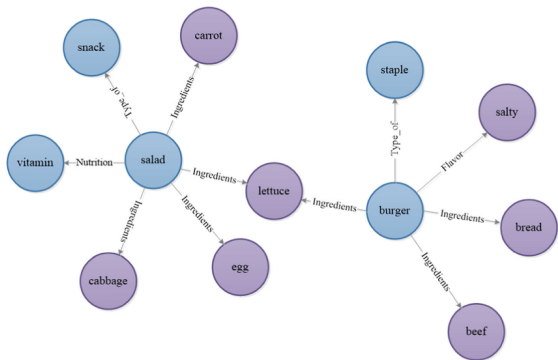
Different behaviors:

different strengths and weaknesses, different costs & requirements



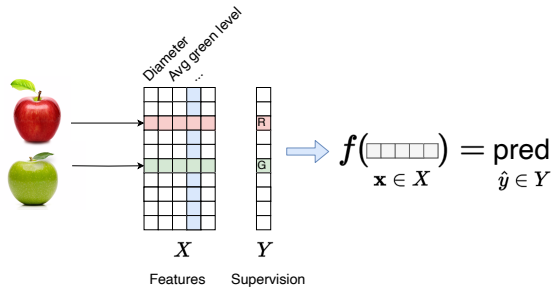
Machine-Learning vs Expert Knowledge

Modeling Expert Knowledge



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

Machine Learning



Different behaviors:

different strengths and weaknesses, different costs & requirements

DEEP LEARNING & REPRESENTATION LEARNING

[APPLICATION TO TEXTUAL DATA]



From tabular data to text

→ Tabular data

- Fixed dimension
- Continuous values



$$\rightarrow f(\text{ } \boxed{} \boxed{} \boxed{} \boxed{} \boxed{} \text{ }) = \text{pred}$$

→ Textual data

- Variable length
- Discrete values

this new iPhone, what a marvel

An iPhone? What a scam!

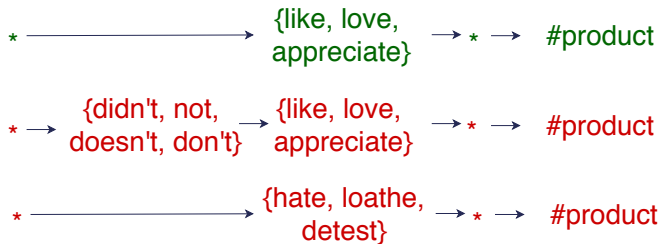


AI + Textual Data: Natural Language Processing (NLP)

NLP = largest scientific community in AI

Linguistics [1960-2010]

Rule-based Systems:



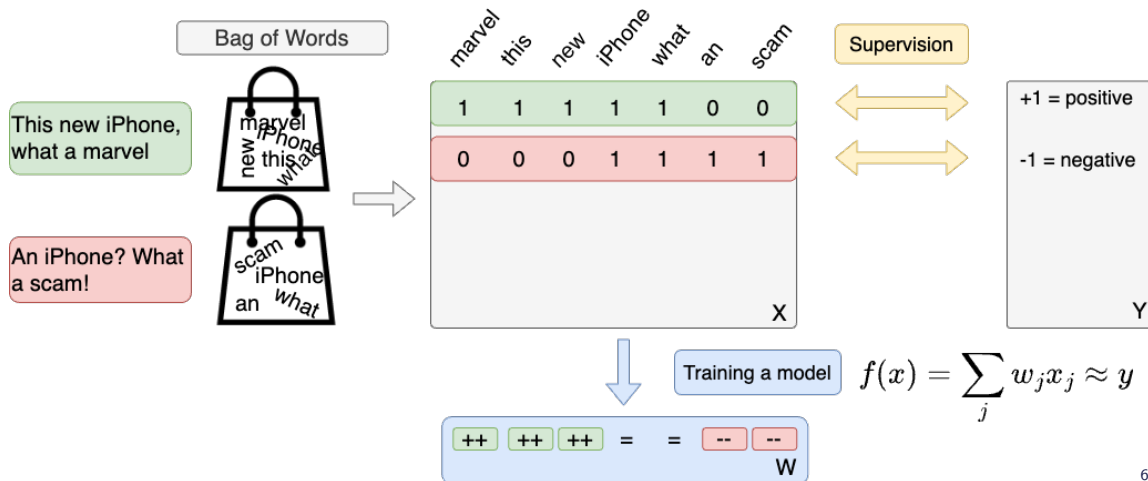
- Requires expert knowledge
- Rule extraction \Leftrightarrow very clean data
- Very high precision
- Low recall
- Interpretable system



AI + Textual Data: Natural Language Processing (NLP)

NLP = largest scientific community in AI

Machine Learning [1990-2015]





AI + Textual Data: Natural Language Processing (NLP)

NLP = largest scientific community in AI

Linguistics [1960-2010]

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

Machine Learning [1990-2015]

- Little expert knowledge needed
- Statistical extraction \Leftrightarrow
robust to noisy data
- ≈ Less interpretable system
- Lower precision
- + Better recall

Precision = criterion for acceptance by industry

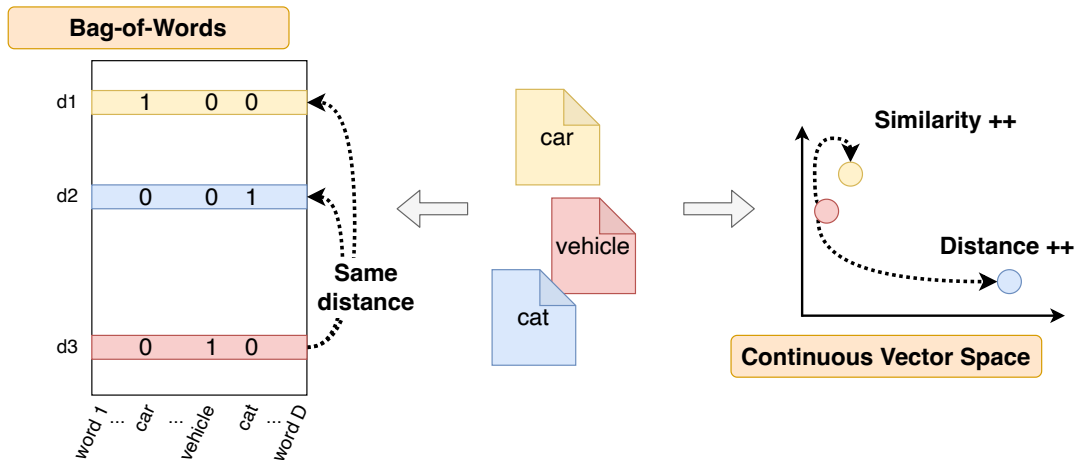
→ Link to metrics



Deep/Representation Learning for Text Data

From Bag of Words to Vector Representations

[2008, 2013, 2016]



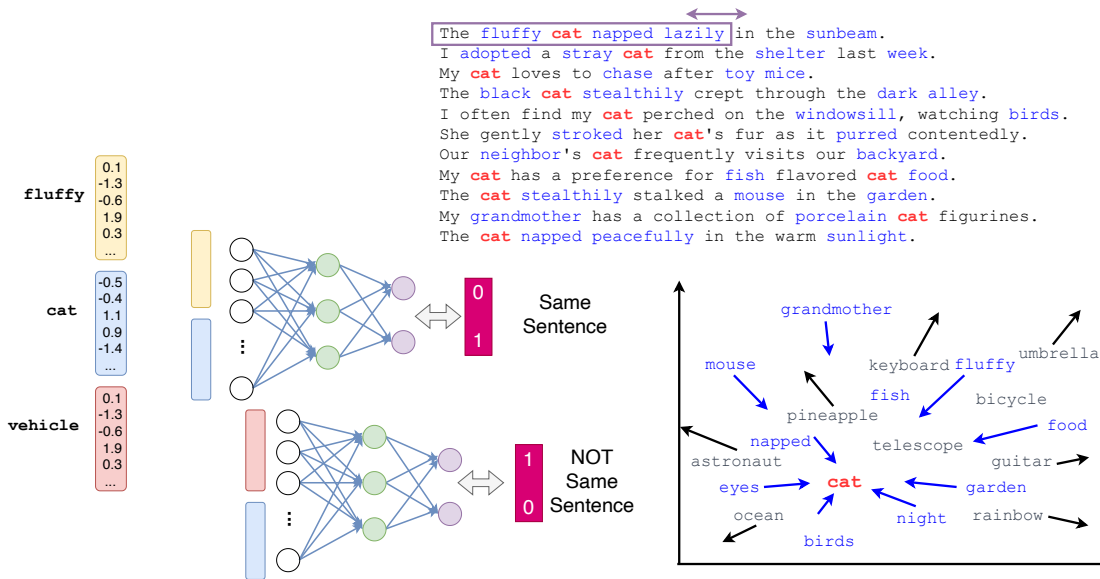
LeCun, Y., Bengio, Y., Hinton, G. (2015). [Deep learning](#). Nature, 521(7553), 436-444.



Deep/Representation Learning for Text Data

From Bag of Words to Vector Representations

[2008, 2013, 2016]

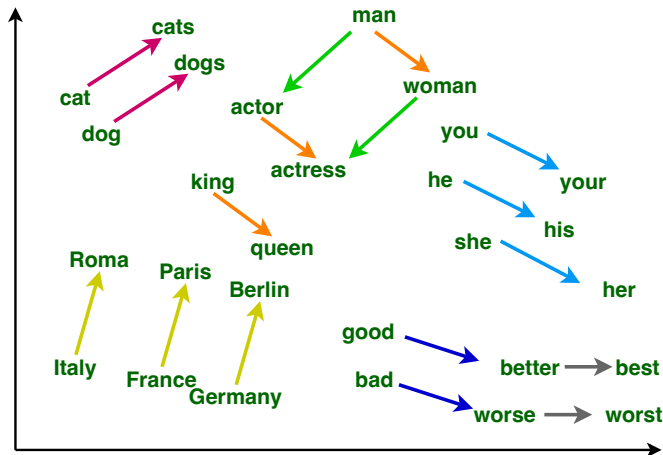




Deep/Representation Learning for Text Data

From Bag of Words to Vector Representations

[2008, 2013, 2016]



- Semantic Space:
similar meanings
 \Leftrightarrow
close positions
- Structured Space:
grammatical regularities,
basic knowledge, ...

Distributed representations of words and phrases and their compositionality, [Mikolov et al. NeurIPS 2013](#)



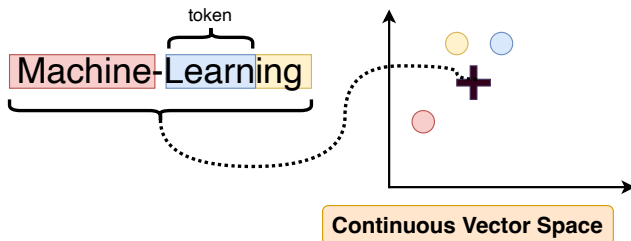
Deep/Representation Learning for Text Data

From Bag of Words to Vector Representations

[2008, 2013, 2016]

From Words to Tokens

Word Piece statistical split



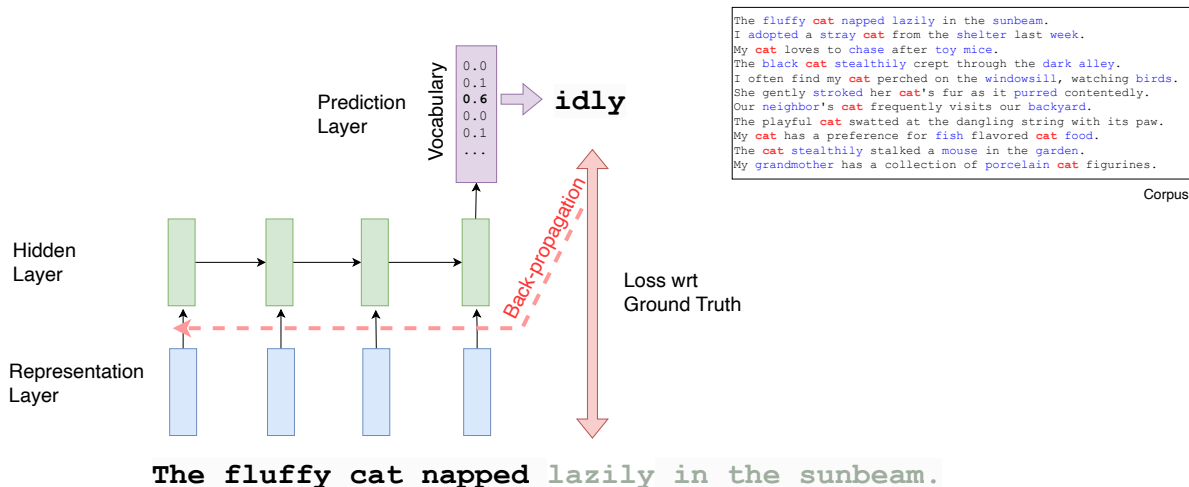
- Representation of unknown words
- Adaptation to technical domains
- Resistance to spelling errors

Enriching word vectors with subword information. [Bojanowski et al. TACL 2017.](#)



Aggregating word representations: towards generative AI

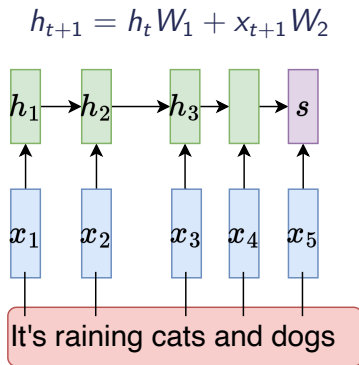
- Generation & Representation
- New way of learning word positions



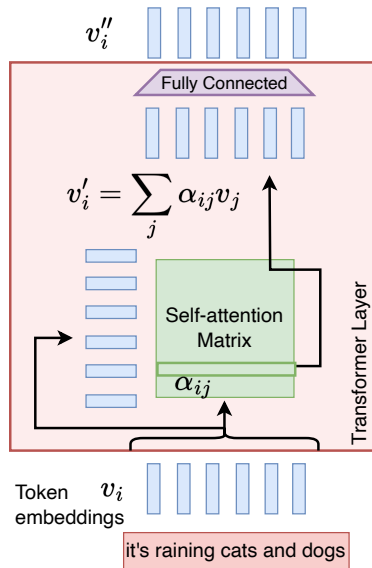


Transformer architecture: state-of-the-art aggregation

Recurrent Neural Network:



Transformer:



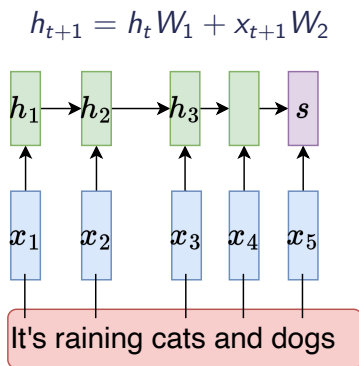
Attention is all you need, [Vaswani et al. NeurIPS 2017](#)

Sequence to Sequence Learning with Neural Networks, [Sutskever et al. NeurIPS 2014](#)

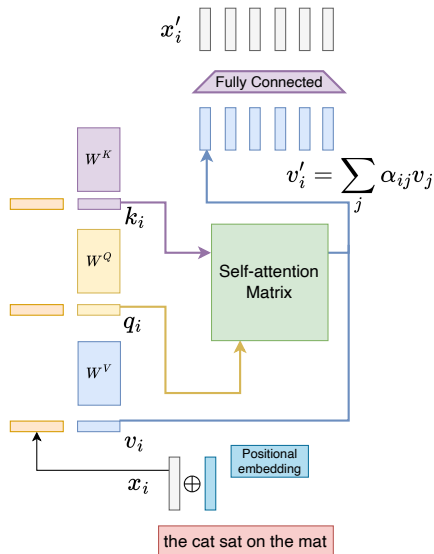


Transformer architecture: state-of-the-art aggregation

Recurrent Neural Network:



Transformer:



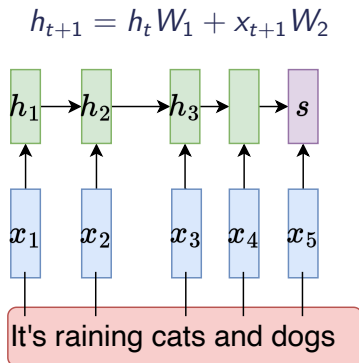
Attention is all you need, [Vaswani et al. NeurIPS 2017](#)

Sequence to Sequence Learning with Neural Networks, [Sutskever et al. NeurIPS 2014](#)

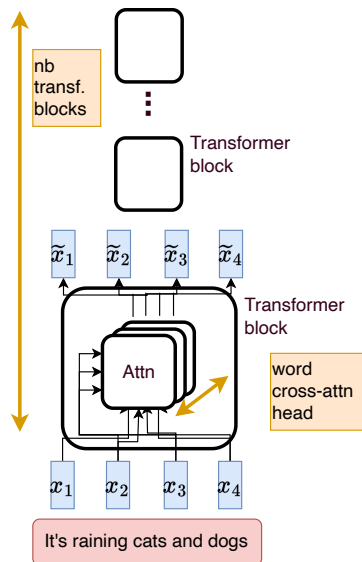


Transformer architecture: state-of-the-art aggregation

Recurrent Neural Network:



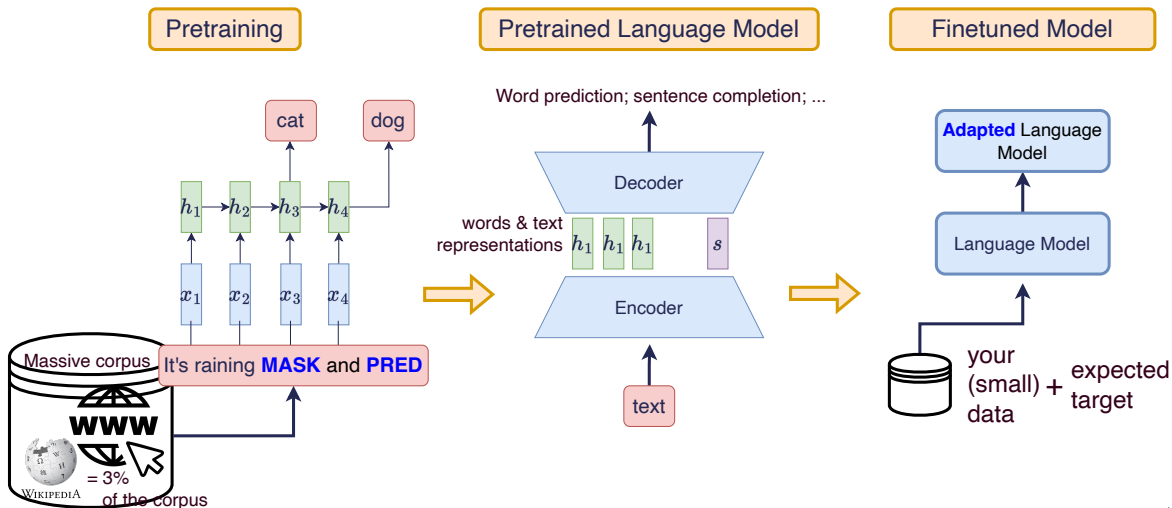
Transformer:





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

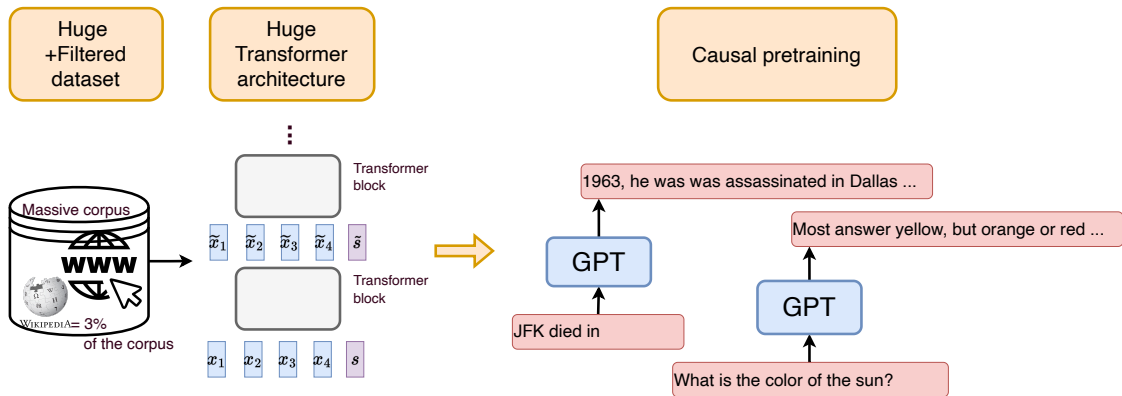
100 MILLION BY THE END OF JANUARY 2023

1.16 BILLION BY MARCH 2023



The Ingredients of chatGPT

0. Transformer + massive data (GPT)



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



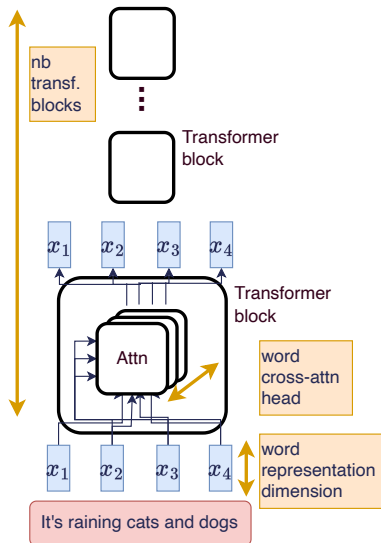
The Ingredients of chatGPT

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?

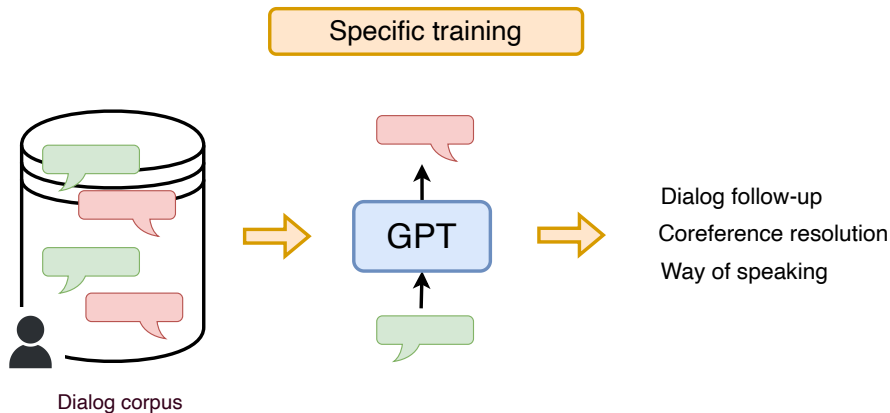
- $1.75 \cdot 10^{11} \Rightarrow 300 \text{ GB} + 100 \text{ GB}$ (data storage for inference) $\approx 400\text{GB}$
- NVidia A100 GPU = 80GB of memory (=20k€)
- Cost for (1) training: 4.6 Million €





The Ingredients of chatGPT

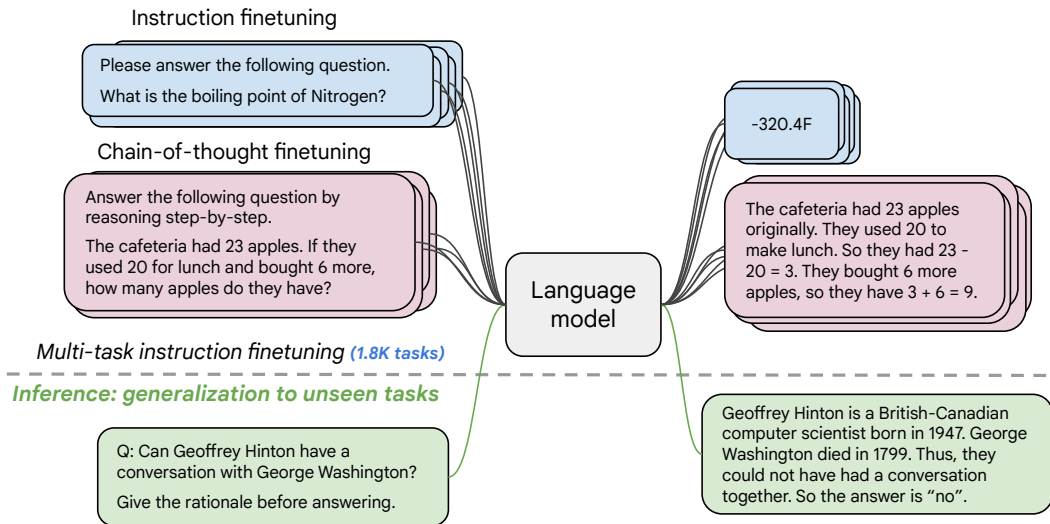
2. Dialogue Tracking



■ **Very clean data**

Data generated/validated/ranked by humans

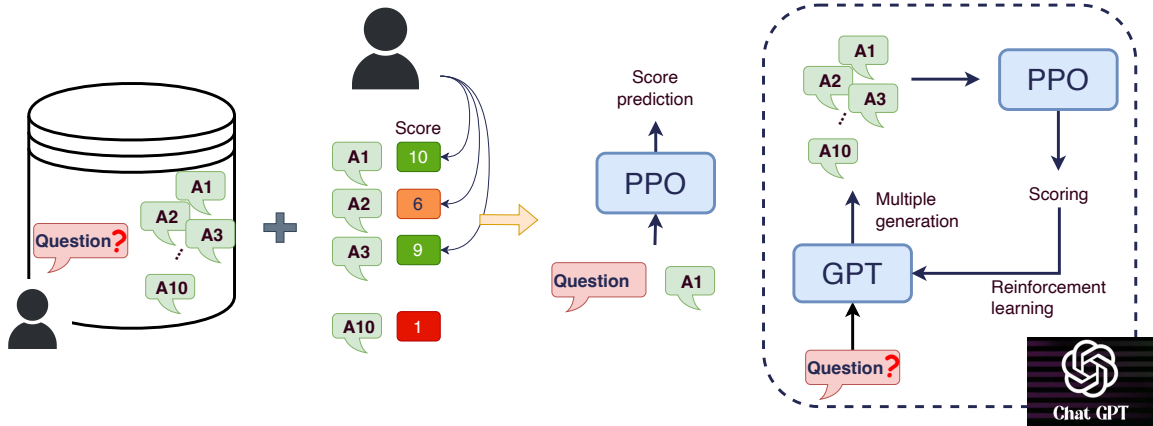
3. Fine-tuning on different (\pm) complex reasoning tasks





The Ingredients of chatGPT

4. Instructions + answer ranking



- Database created by humans
- Response improvement

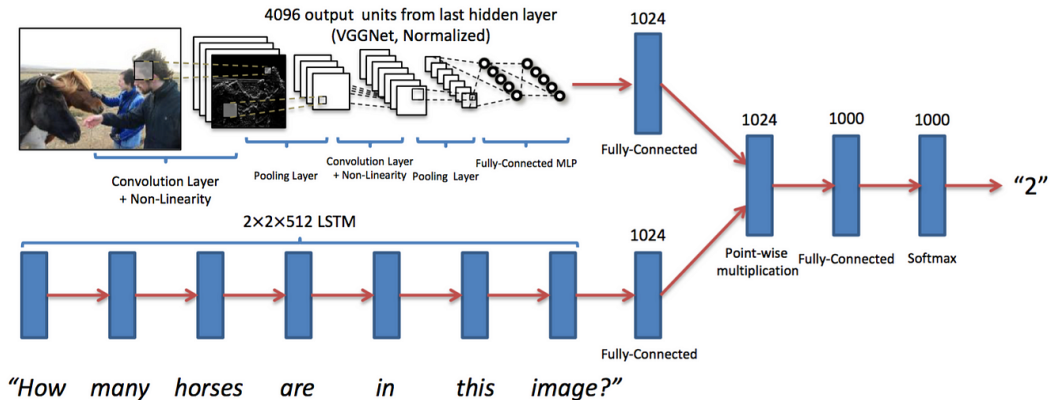
- ... Also a way to avoid critical topics = censorship



GPT4 & Multimodality

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

The example of VQA: visual question answering



⇒ Backpropagate the error ⇒ modify word representations + image analysis

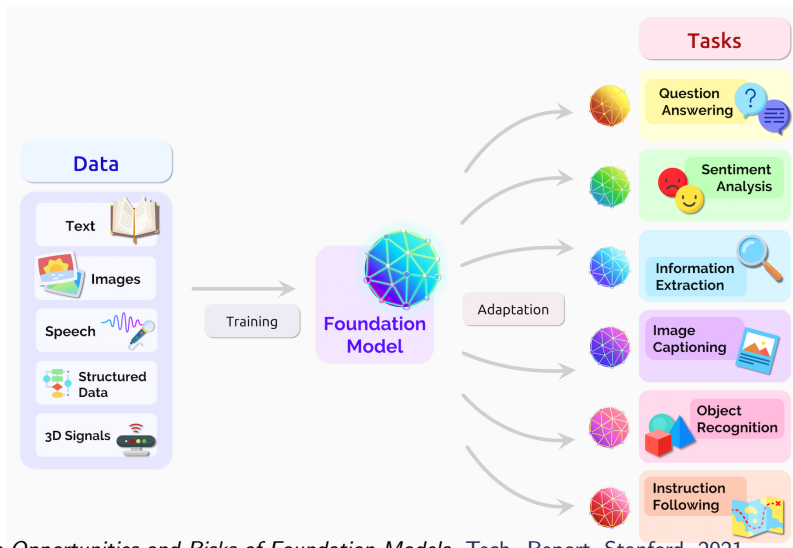


VQA: Visual Question Answering, arXiv, 2016, A. Agrawal et al.



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.

MACHINE LEARNING LIMITS



chatGPT and the relationship with truth

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

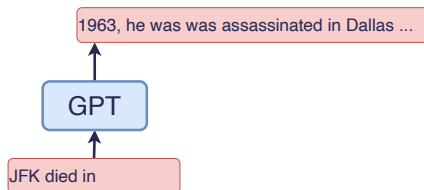
2 Predict the most **plausible** word...
⇒ produces **hallucinations**

3 **Offline** functioning

4 chatGPT \neq **knowledge graphs**

5 Brilliant answers...

And silly mistakes!
+ we cannot predict the errors



Example: producing a bibliography

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](#)

3. "Generative Modeling for Time Series" par Bao et al. (2017) - Cette étude propose un modèle qui combine les forces des RNNs et des GVBs: le Variational Recurrent Auto-Encoder (VRAE). Such a model can be used for efficient, large scale time series generative modeling.

[Enregistrer](#) [VS Citer](#) [Clé 302 fois](#) [Autres articles](#) [Les 2 versions](#) [36](#)

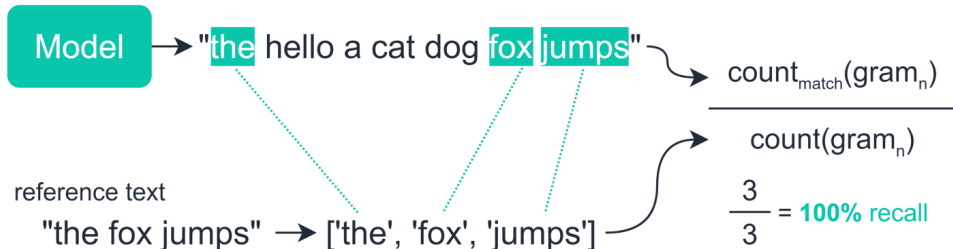
4. "Deep Variational Bayes Filters: Unsupervised Learning of State Space Models from Data" par Kulkarni et al. (2017) - Cette étude présente une nouvelle VAE pour la



Generative AI: how to evaluate performance?

The critical point today

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

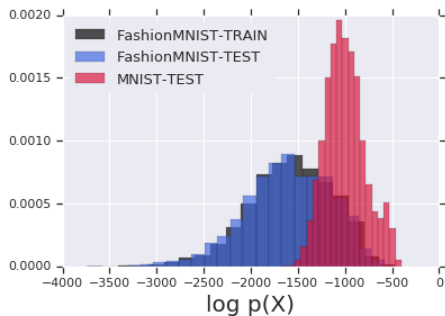




Generative AI: how to evaluate performance?

The critical point today

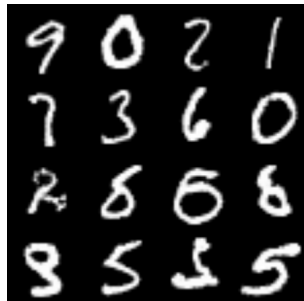
- How to evaluate against ground truth?
- How to evaluate system confidence / plausibility of generation?



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



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

⇒ 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.



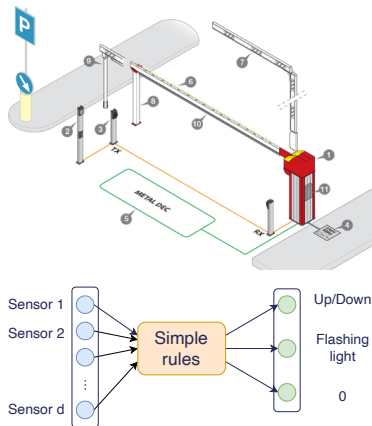
V

and today?

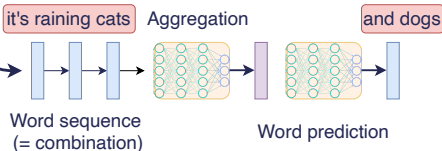




Explainability... And complexity



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



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



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)



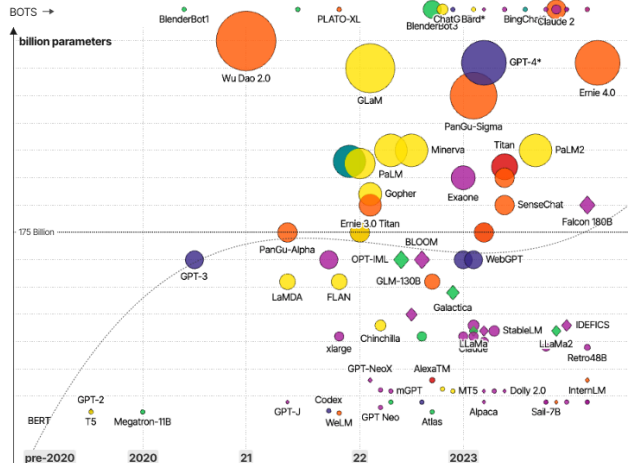
[Uber Accident, 2018]

- Simple system
- Exhaustive testing of inputs/outputs
- **Predictable & explainable**
- Large dimension
- Complex non-linear combinations
- **Non-predictable & non-explainable**

The Rise and Rise of A.I. Large Language Models (LLMs) & their associated bots like ChatGPT

size = no. of parameters open-access

● Amazon-owned ● Chinese ● Google ● Meta / Facebook ● Microsoft ● OpenAI ● Other



David McCandless, Tom Evans, Paul Barton
Information is Beautiful // UPDATED 2nd Nov 23

* = parameters undisclosed // see the data

1998	LeNet-5	= 0.06M
2011	Senna	= 7.3M
2012	AlexNet	= 60M
2017	Transformer	= 65M / 210M
2018	ELMo	= 94M
2018	BERT	= 110M / 340M
2019	GPT2	= 1,500M
2020	GPT3	= 175,000M
2025	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

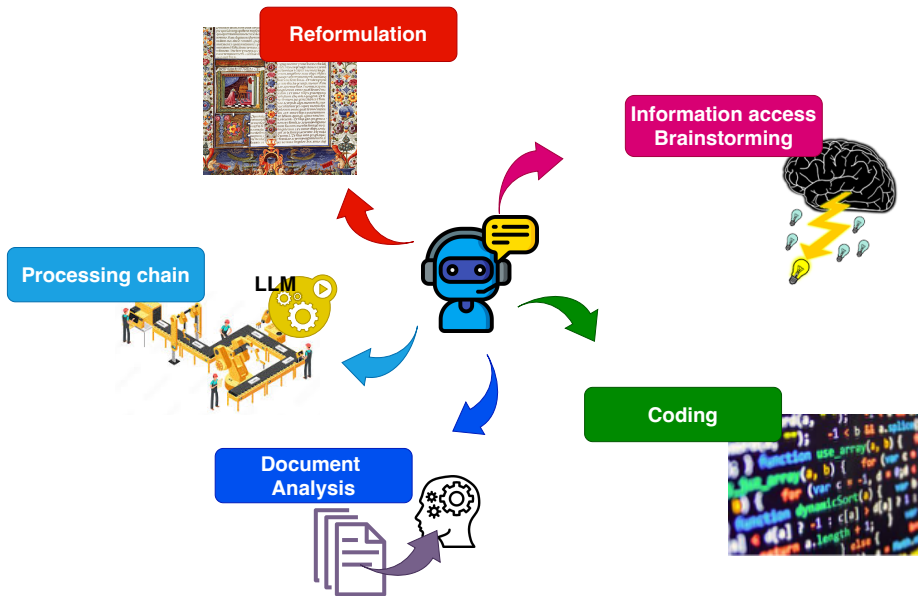


**WHEN YOU UNDERESTIMATE A 1977 CHESS ENGINE...
AND IT HUMBLER YOU IN FRONT OF THE WHOLE INTERNET**

LARGE LANGUAGE MODELS USES [IN NUTRITION RESEARCH]



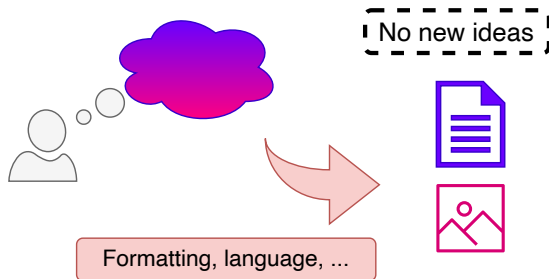
Key uses in 5 pictures





(1) 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, ...



(1) Nutrition use : Input standardization (?)

⇒ opportunity to fuse heterogeneous information



INCA2 Individual national study of food consumption in France Database

anses
agence nationale de sécurité sanitaire
alimentation, environnement, travail

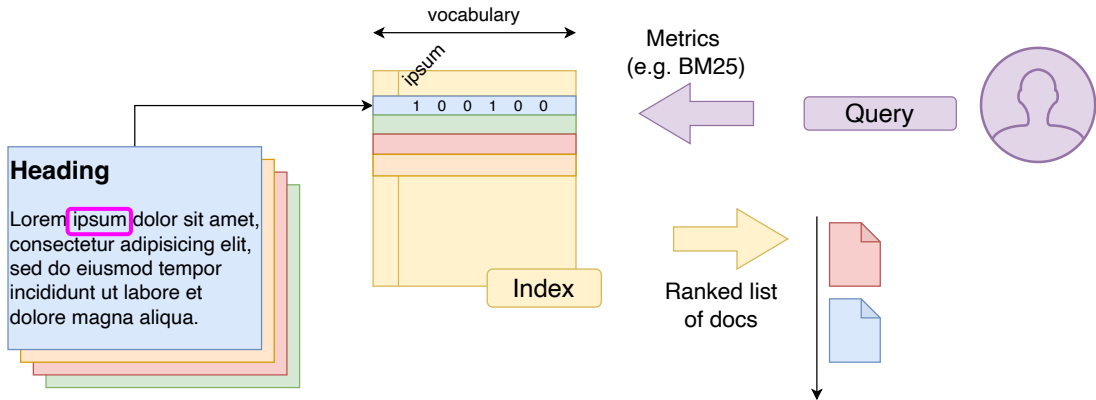


ARTICLE OPEN

FoodOn: a harmonized food ontology traceability, quality control and data in

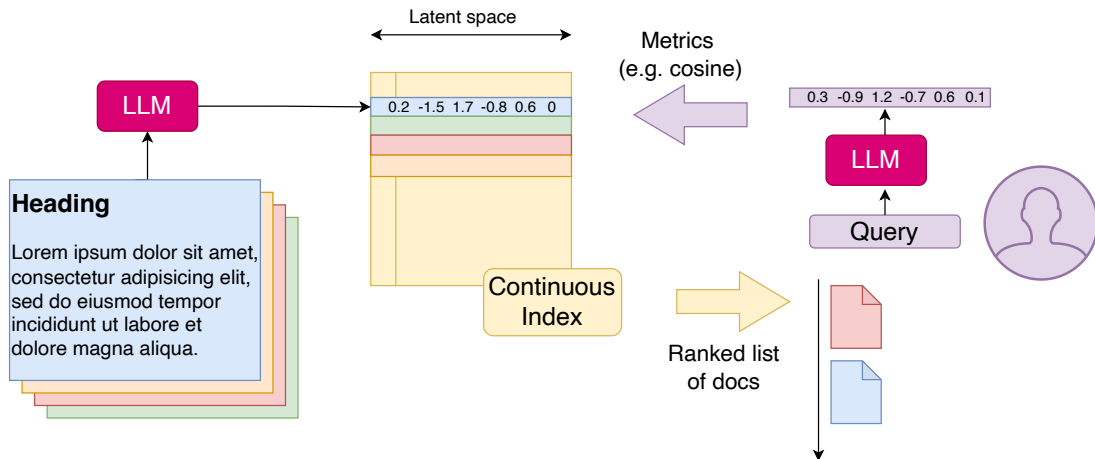
J. Griffiths^{2,8}, Gurinder S. Gosal¹, Pier L. Buttigieg³, Ro
rinkman² and William W. L. Hsiao^{1,2,7}





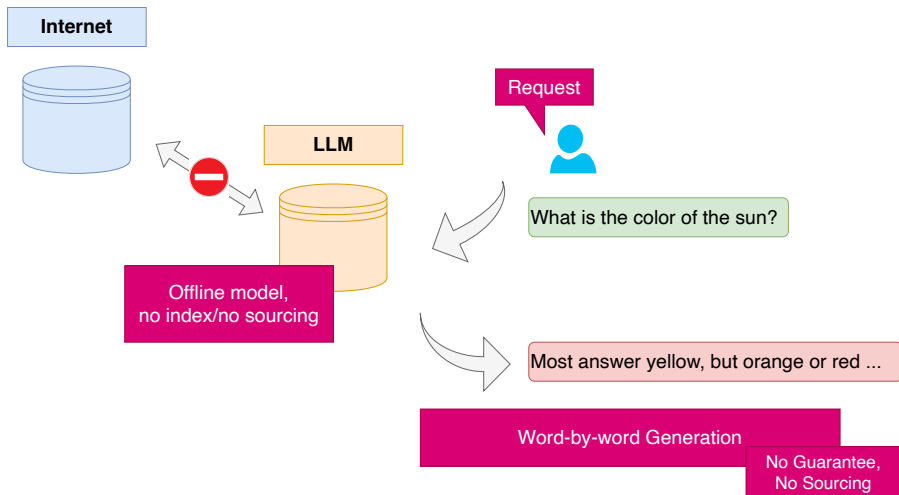


(1) Chat & RAG : a new way to access information



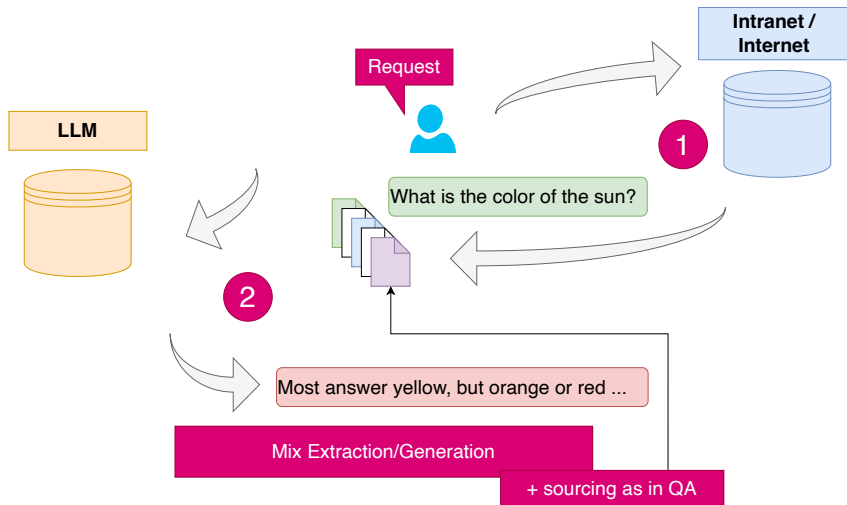


(1) Chat & RAG : a new way to access information





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



⇒ A way to build a *reliable* chatbot to advise users?

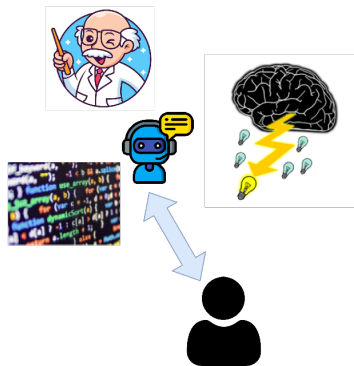
■ Parametric memory vs Information Retrieval



(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?

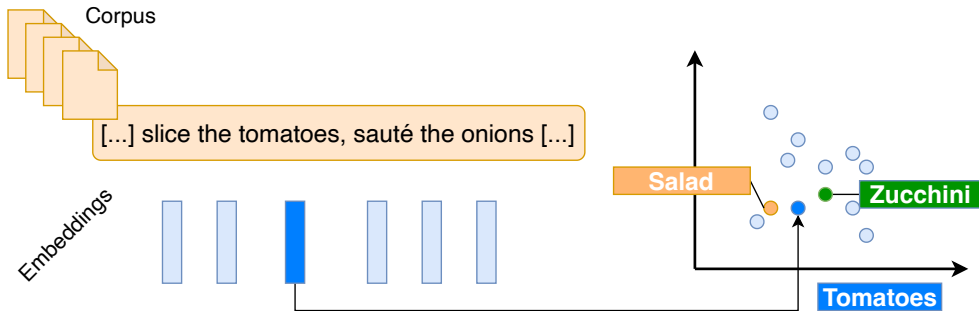
Propose 10 sections for a course on Transformers in AI

- 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

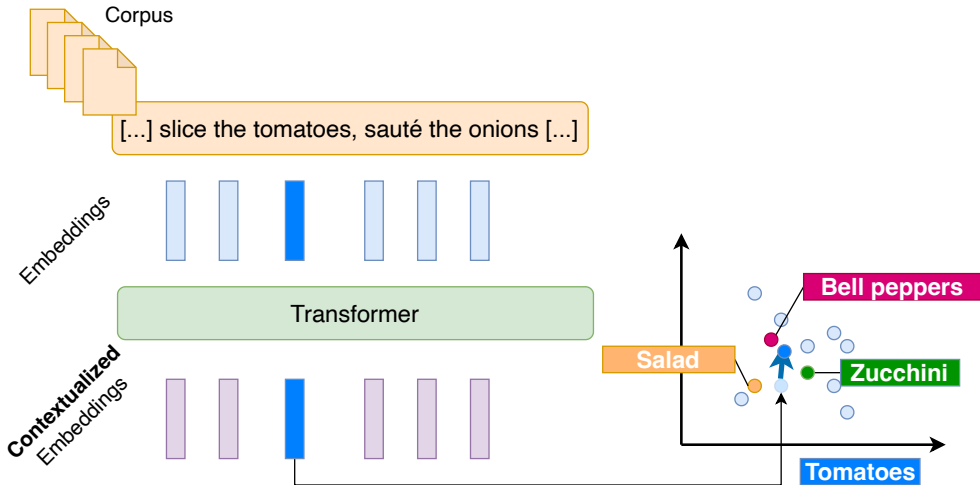
- Brainstorming in the kitchen: which application for cooking?
- Ingredient substitution... At every scale: Ingredient, Food, Dish





(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



(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

INCA2 Individual
national study of
food consumption in
France Database



Ciqua

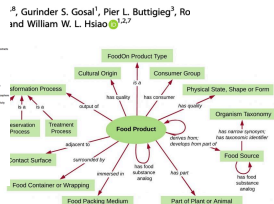
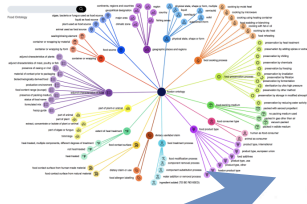
Table de composition nutritionnelle des aliments

	Moyenne	Min	Max
Glucose (g/100 g)	2,3		
Lactose (g/100 g)	< 0,2		
Maltose (g/100 g)	< 0,2		
Saccharose (g/100 g)	< 0,2		
Amidon (g/100 g)	< 0,35		
Fibres alimentaires (g/100 g)	1,2		
Polyols totaux (g/100 g)	< 0,5	0	
Candides (g/100 g)	0,35		
Alcool (g/100 g)	0		
Acides organiques (g/100 g)	0,56		
AG saturés (g/100 g)	< 0,01		
AG monoinsaturés (g/100 g)	< 0,01		
AG polyinsaturés (g/100 g)	< 0,01		
AG totaux (g/100 g)	< 0,01		

ARTICLE OPEN

FoodOn: a harmonized food ontology
traceability, quality control and data in

Gurinder S. Gosal¹, Pier L. Buttigieg³, Ro
and William W. L. Hsiao^{1,2,7}



Multilingual alignment

(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

INCA2 Individual national study of food consumption in France Data and analysis

Ciqua

Table de composition nutritionnelle des aliments

Glucose (g/100 g)
Lactose (g/100 g)
Maltose (g/100 g)
Saccharose (g/100 g)
Amidon (g/100 g)
Fibres alimentaires (g/100 g)
Polyols totaux (g/100 g)
Candies (g/100 g)
Alcool (g/100 g)
Acides organiques (g/100 g)
AG saturés (g/100 g)
AG monoinsaturés (g/100 g)
AG polyinsaturés (g/100 g)
AG totaux (g/100 g)

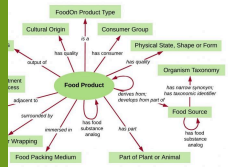
A new alignment method based on FoodOn as pivot ontology

Patrice Buche, Julien Cufi, Liliana Ibanescu, Alrick Oudot, Magalie weber

12/10/2021

Food ontology
control and data in

S. Gosal¹, Pier L. Buttigieg¹, Ro
W. L. Hsiao^{1,2,7}





(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



GitHub
Copilot



- 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?

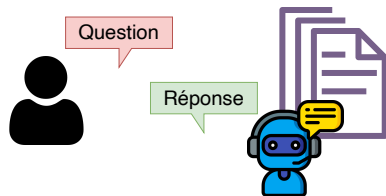
```
sentiments.ts  write_sql.go  parse_expenses.py  addresses.rb

1 import datetime
2
3 def parse_expenses(expenses_string):
4     """Parse the list of expenses and return the list of triples (date,
5     Ignore lines starting with #.
6     Parse the date using datetime.
7     Example expenses_string:
8     2016-01-02 -34.01 USD
9     2016-01-03 2.59 DKK
10    2016-01-03 -2.72 EUR
11    """
12    expenses = []
13    for line in expenses_string.splitlines():
14        if line.startswith("#"):
15            continue
```



(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

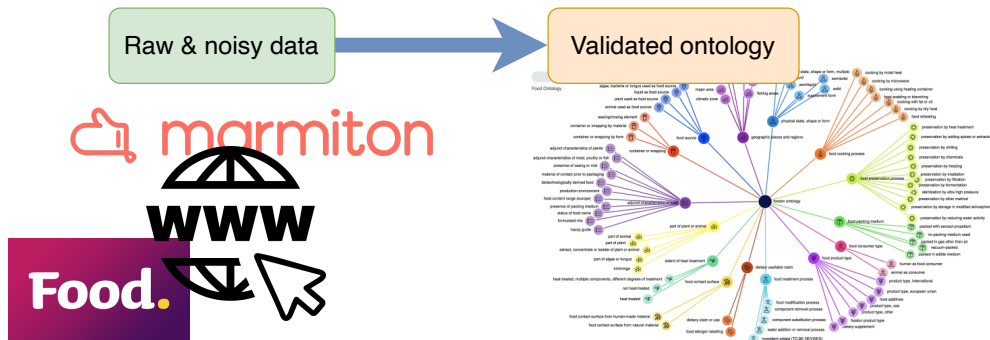
Think **Smarter**,
Not Harder

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?

(4) Information Extraction in Nutrition

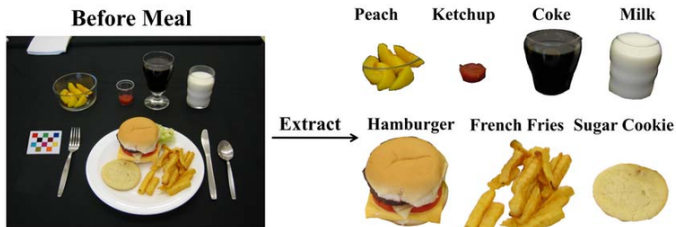
■ Ontology building (mostly textual data)



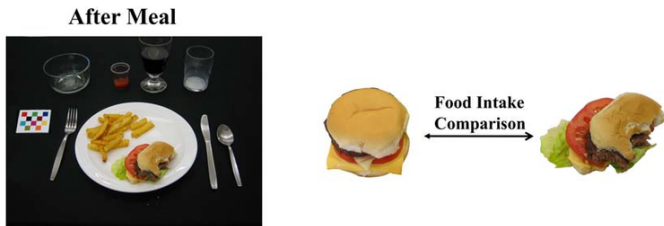


(4) Information Extraction in Nutrition

- Ontology building (mostly textual data)
- Image analysis



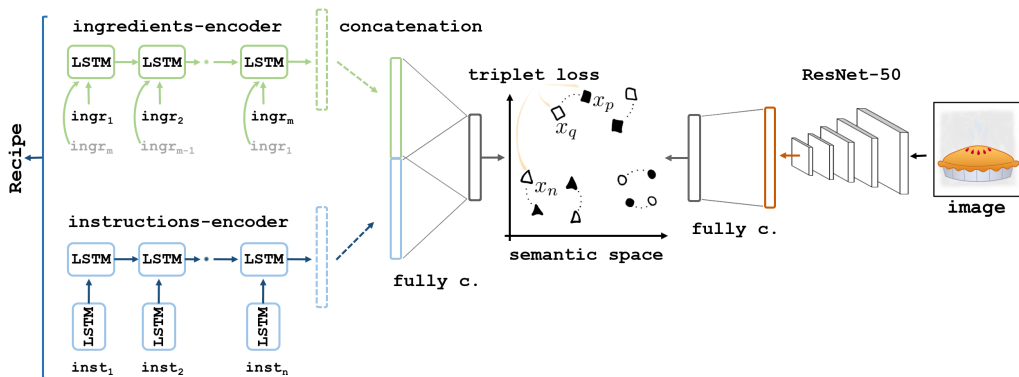
- Food recognition
- Segmentation
- Estimation of quantities





(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.



(4) Information Extraction in Nutrition

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

	ingr (ingredients)	instr (cooking instructions)	image
Pizza	<ol style="list-style-type: none">1) <i>pizza dough</i>2) <i>hummus</i>3) <i>arugula</i>4) <i>cherry / grape tomatoes</i>5) <i>pitted greek olives</i>6) <i>crumbled feta cheese</i>	<ol style="list-style-type: none">1) <i>Cut the dough into two 8-ounce sized pieces.</i>2) <i>Roll the ends under to create round balls.</i>3) <i>Then using a well-floured rolling pin, roll the dough out into 12-inch circles.</i>4) <i>Place the dough circles on sheets of parchment paper.</i> <p>... ..</p>	
Pecan Pie	<ol style="list-style-type: none">1) <i>unsalted butter</i>2) <i>eggs</i>3) <i>condensed milk</i>4) <i>sugar</i>5) <i>vanilla extract</i>6) <i>chopped pecans</i>7) <i>chocolate chips</i> <p>... ..</p>	<ol style="list-style-type: none">1) <i>Preheat the oven to 375 degrees F.</i>2) <i>In a large bowl, whisk together the melted butter and eggs until combined.</i>3) <i>Whisk in the sweetened condensed milk, sugar, vanilla, pecans, chocolate chips, butterscotch chips, and coconut.</i> <p>... ..</p>	

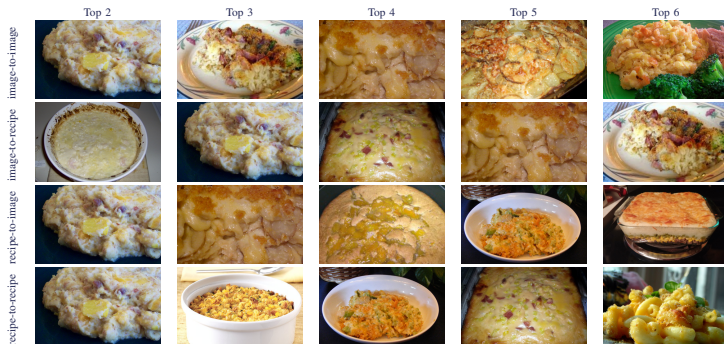


Images & Recipes: Retrieval in the cooking context, SIGIR 2018
Carvalho et al.



(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.



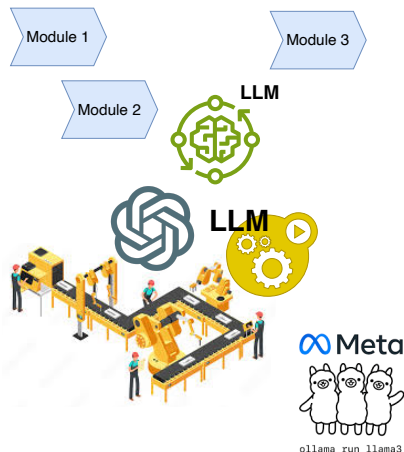
(5) LLM in a Production Pipeline / Agentic AI

- 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 AI**

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





(5) LLM in a Production Pipeline / Agentic AI

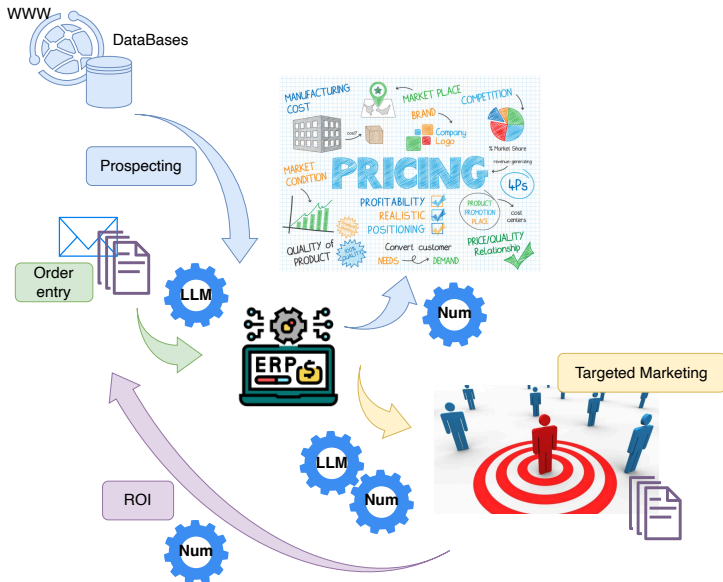
- Run LLM
- Extract kn
- Sort docu
- Generate e

- Generate v
- dataset siz

⇒ Integrate t

=

- Can I t
- How m
- How gc



Module 3



Meta

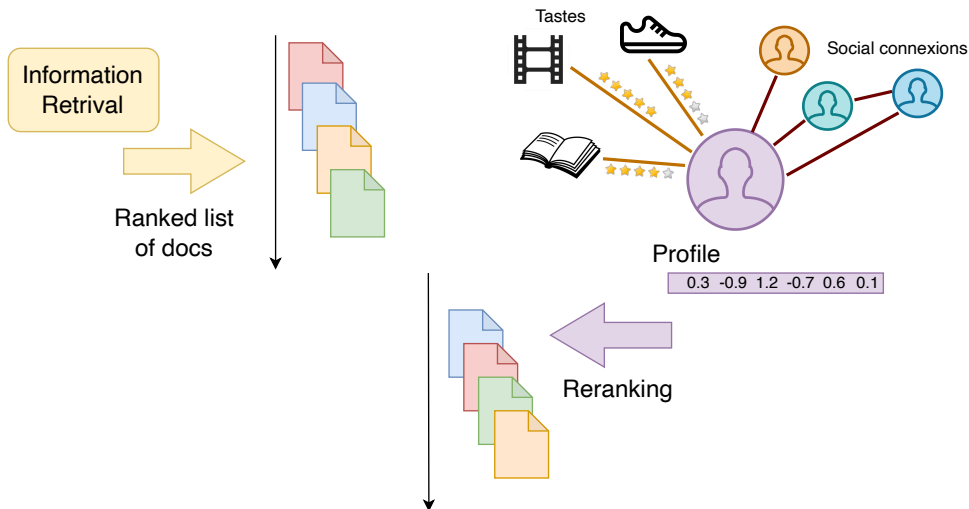


ollama run llama3



(5) What about RecSys in Nutrition?

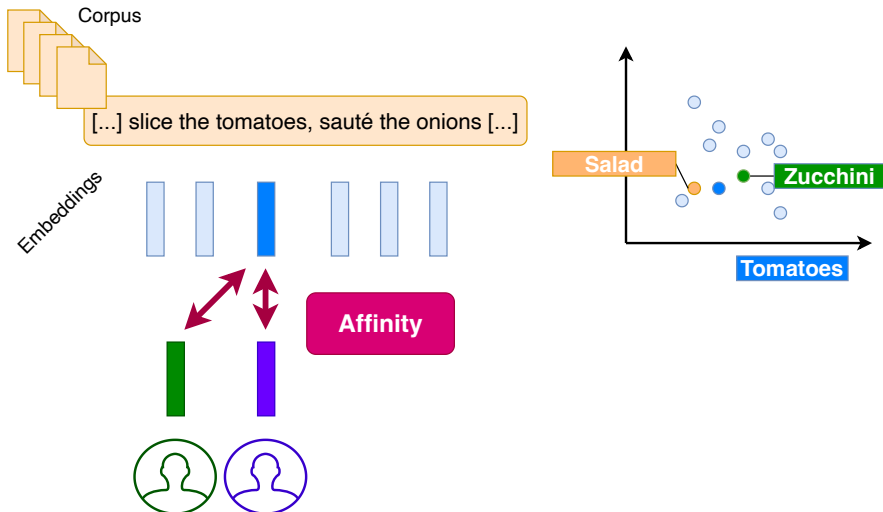
Profiling is roughly everywhere in Information Retrieval





(5) What about RecSys in Nutrition?

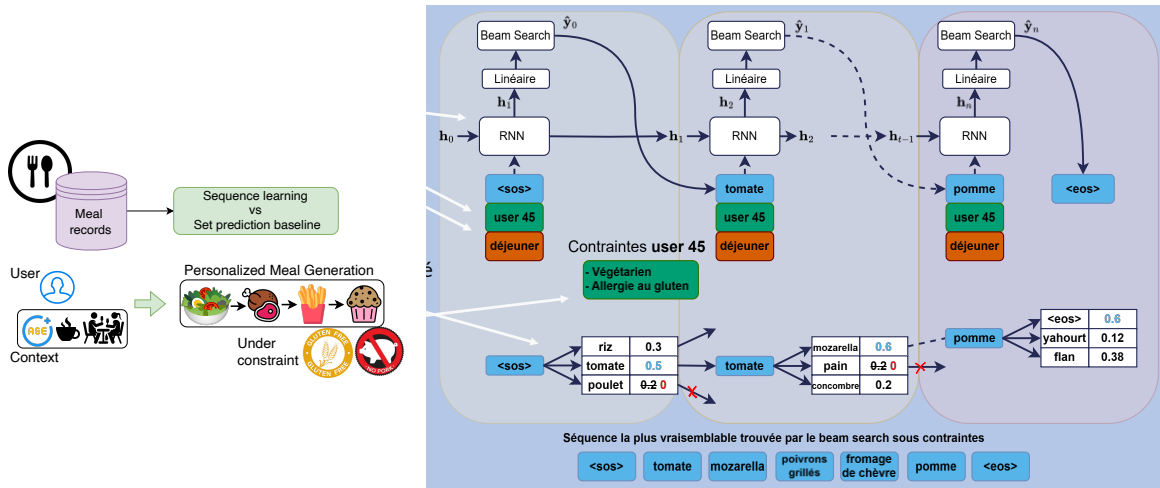
Opportunities in nutrition : modeling user preferences





(5) What about RecSys in Nutrition?

Building consistent proposals... With expert constraints



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

CONCLUSION



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 \Rightarrow enrich databases
- A modular component for next-generation recommender systems