

# A quick overview of Natural Language Processing for Information Extraction and Scientific Surveillance

Xavier Tannier

xavier.tannier@sorbonne-universite.fr

Orphanet, September 17, 2024



# Appetizer



1 Allez sur [wooclap.com](https://wooclap.com)

2 Entrez le code d'événement dans le bandeau supérieur

Code d'événement  
**UXPXKF**

 Activer les réponses par SMS



# NLP tasks

9 hémocultures positives le 26/6/15 à  
staphylocoque aureus méticilline  
sensible. BACTERIAL RESISTANCE

## VASCULAR DISEASE

Patient atteint d'ulcères artériels des  
membres inférieurs, suivis à St Joseph (Dr  
Wyliana) avec greffe cutanée en octobre  
2015

## PROCEDURE

## MEDICAL DEVICE

## ANATOMY

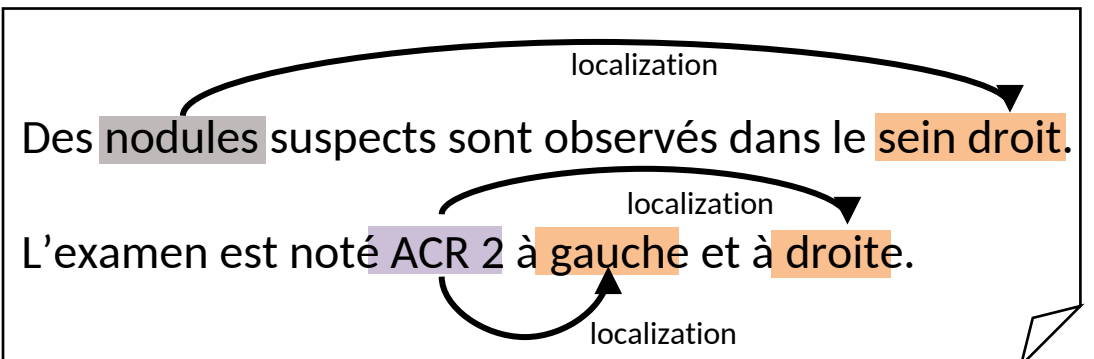
La contraception par les dispositifs intra utérins

Patient sans signe clinique évident  
de traumatisme crânien

[C0018674] (Craniocerebral Trauma)

Patiente avec antécédent de  
chirurgie bariatrique

[C1456587] (Bariatric Surgery)

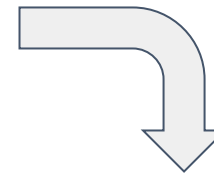


# NLP tasks

MAMMOGRAPHY:

There is a <sup>SHAPE</sup> 1.8 cm <sup>LESION</sup> round mass with a circumscribed margin in the <sup>LATERALITY</sup> left breast in the <sup>LOCALIZATION</sup> anterior depth central to the nipple. There also is a <sup>SHAPE</sup> 1.4 cm <sup>LESION</sup> oval mass with an obscured margin in the <sup>LATERALITY</sup> left breast in the <sup>LOCALIZATION</sup> anterior depth of the inferior region.

ASSESSMENT: <sup>SCORE</sup> BI-RADS Category 3



SCORE	LATERALITY	LESION	LATERALITY	LOCALIZATION	SHAPE
BI-RADS 3	LEFT	mass	LEFT	QSI	1.8
BI-RADS 1	RIGHT	mass	LEFT	QI	1.4



# Objectives

Pseudonymising patient records

Structuring data

Indexing & Retrieving

Selecting similar patients

Selecting patients matching criteria

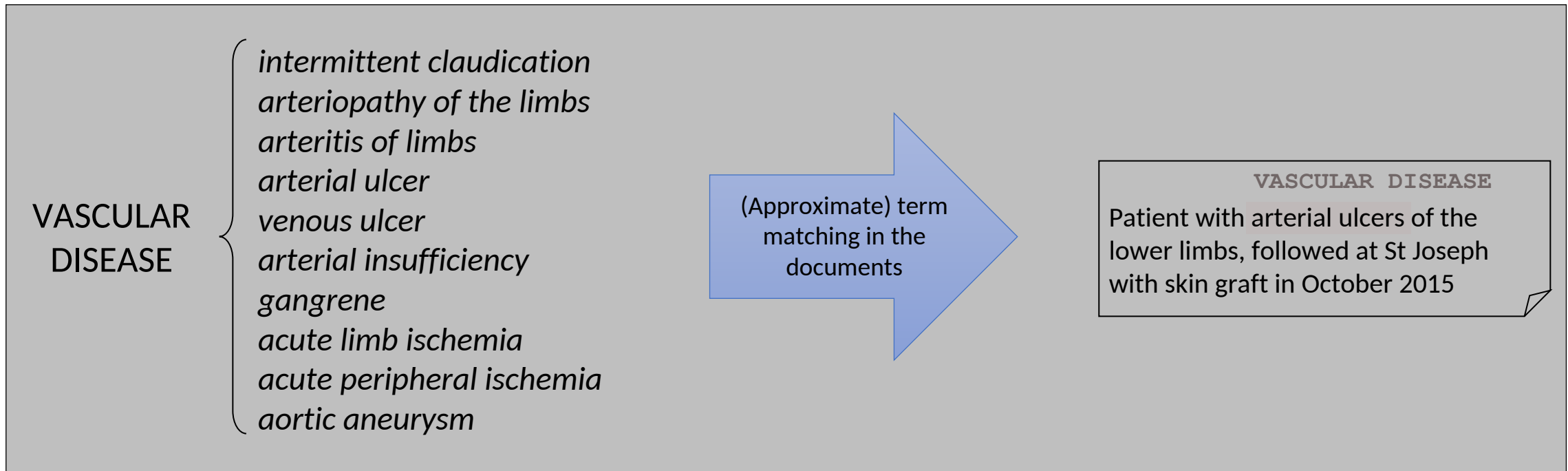
# Methods

- 1. Rule-based systems**
2. Supervised learning systems
3. Generative, large language models
4. Retrieval-Augmented Generation



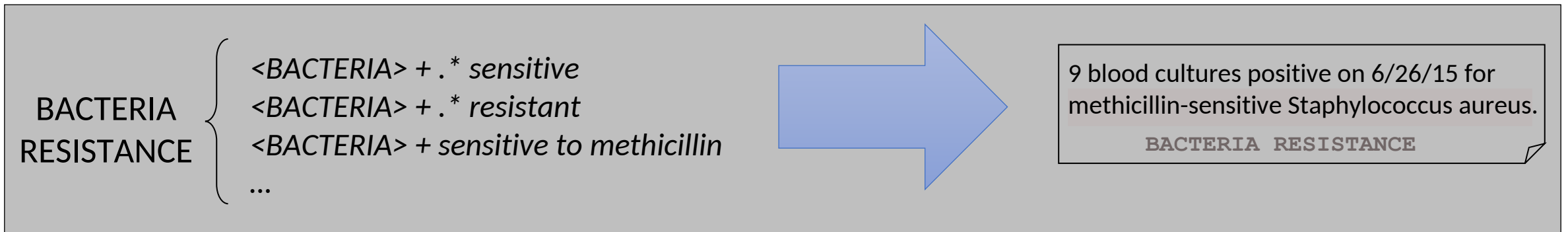
# Rule-based systems

## 1. Terminological approach



# Rule-based systems

1. Terminological approach
2. **Additional rules**



# Rule-based systems

1. Terminological approach
2. Additional rules
3. Trickier rules

Emmanuelle Kempf et al., 2023

secondary, progressive,  
fractured or suspected  
ostecondensation

```
ost[ée]ocondensa. {1,20} (suspect|secondaire|[ée]vulsive)| (l[ée]sion|anomalie|image).  
{1,20}os.  
{1,30} (suspect|secondaire|[ée]vulsive)| os. {1,30} (l[ée]sion|anomalie|image). {1,20}  
(suspect|secondaire|[ée]vulsive)| (l[ée]sion|anomalie|image).  
{1,20};L[I,Y]tique| (l[ée]sion|anomalie|image). {1,20}condensant. {1,20}  
(suspect|secondaire|[ée]vulsive)| fracture. {1,30}  
(suspect|secondaire|[ée]vulsive)| ((l[l[ée]sion|anomalie|image|nodule). {1,80}  
(secondaire))| ((l[l[ée]sion|anomalie|image|nodule)s.{1,40}suspec?ts?).
```

post-operative  
anatomopathological tumor  
stage (pTNM)

```
([ycpP]{1,2}s? (T([01234x]|is)[abcdx?]) [, \s] {0,2} [ycp] {0,2}s? (N[xo01234\+][abcdx?]) *s?  
(M[o01]? [\+x]?))| ((T([01234x]|is)[abcdx?]) [, \s] {0,2} [ycp] {0,2}s? (N[xo01234\+][abcdx?])  
\s?  
(M[o01]? [\+x]?))
```

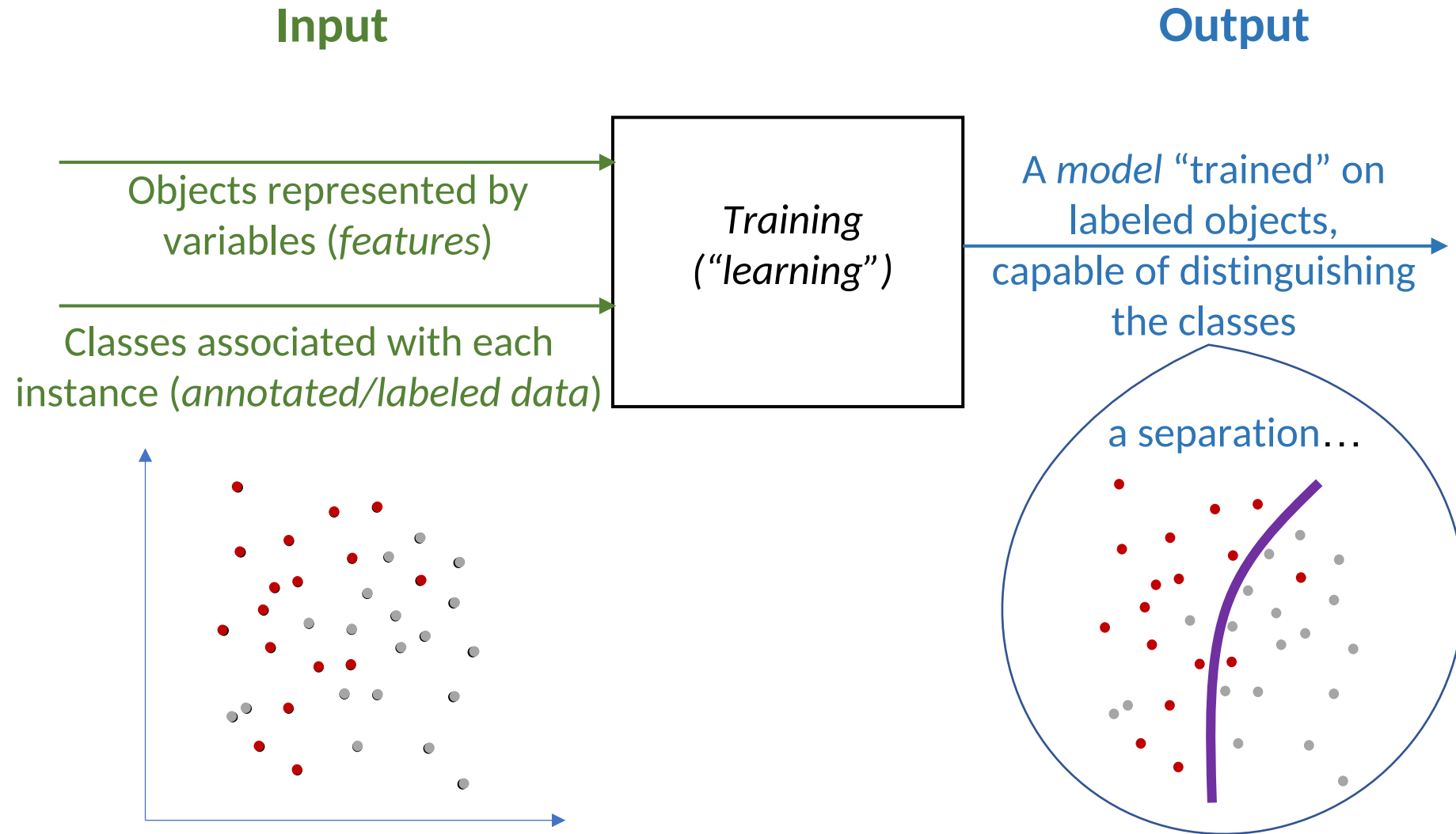
# Rule-based systems: Pros & Cons

	Rules
General performance	Yellow
Ease of implementation	Green
Need for human expertise	Yellow
Explainability	Green
Material resources	Green
Energy consumption	Green
Ease of maintenance	Yellow
Generalization to a different problem/context	Red

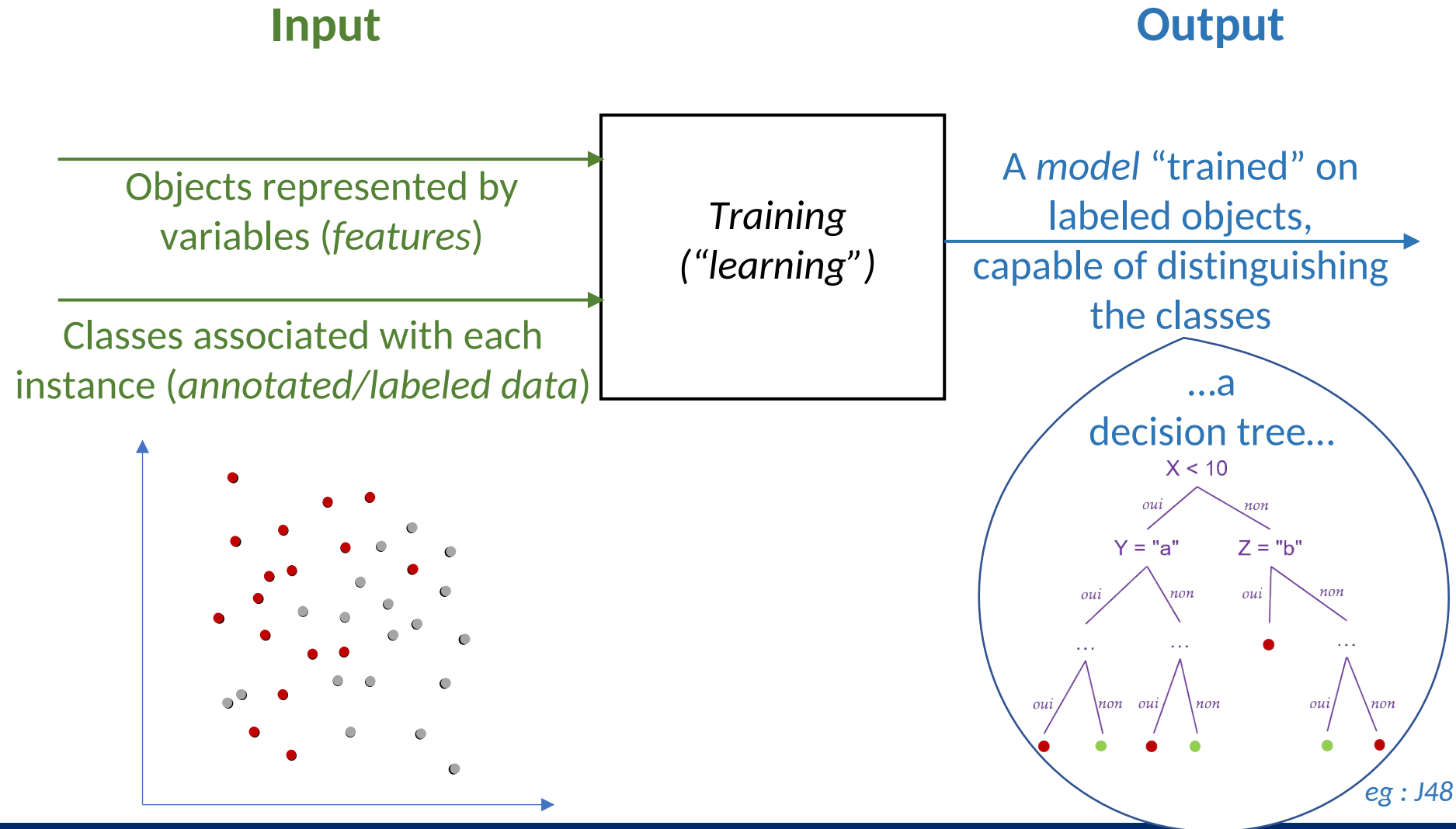
# Methods

1. Rule-based systems
- 2. Supervised learning systems**
3. Generative, large language models
4. Retrieval-Augmented Generation

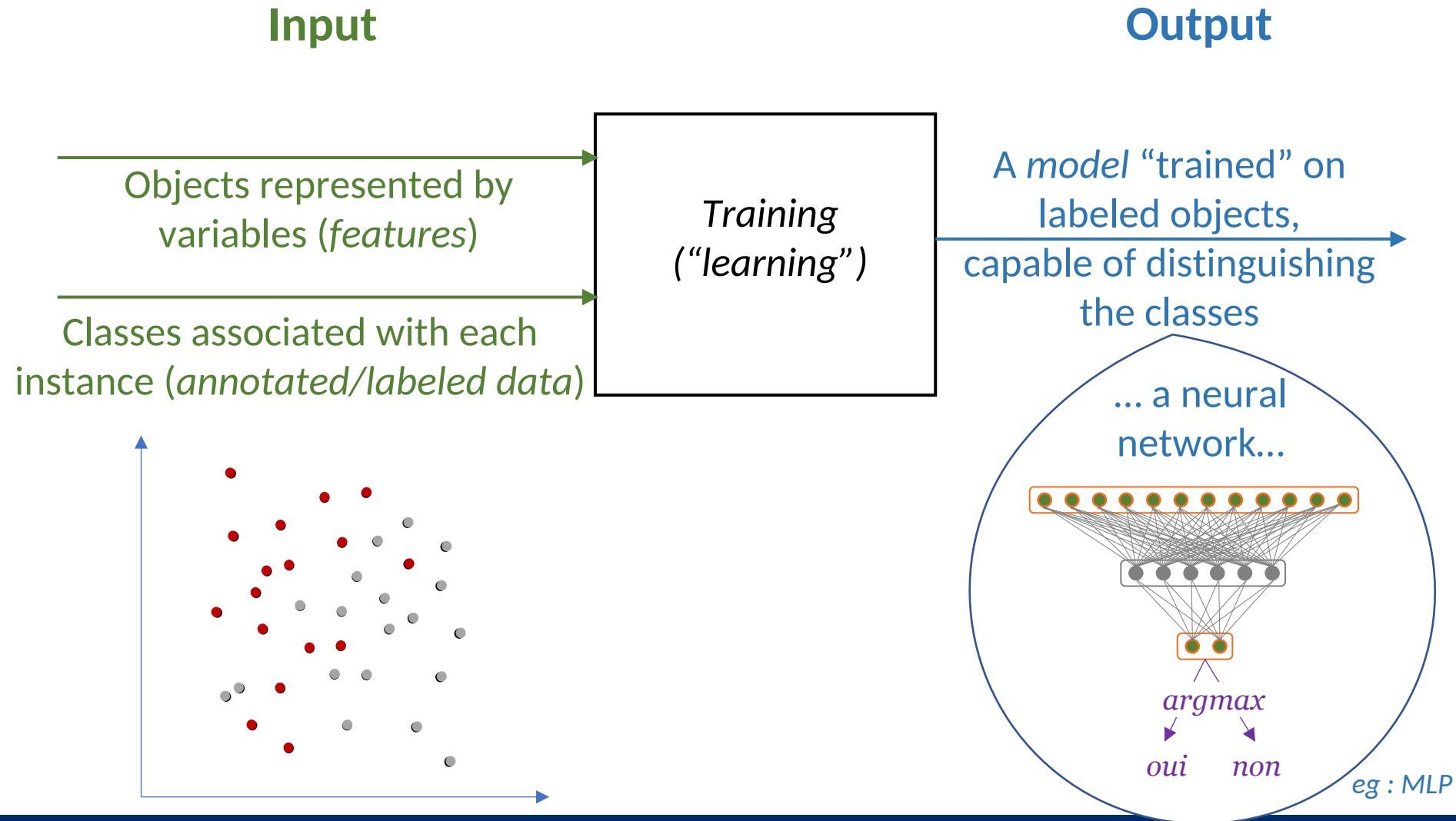
# (What is supervised learning?)



# (What is supervised learning?)

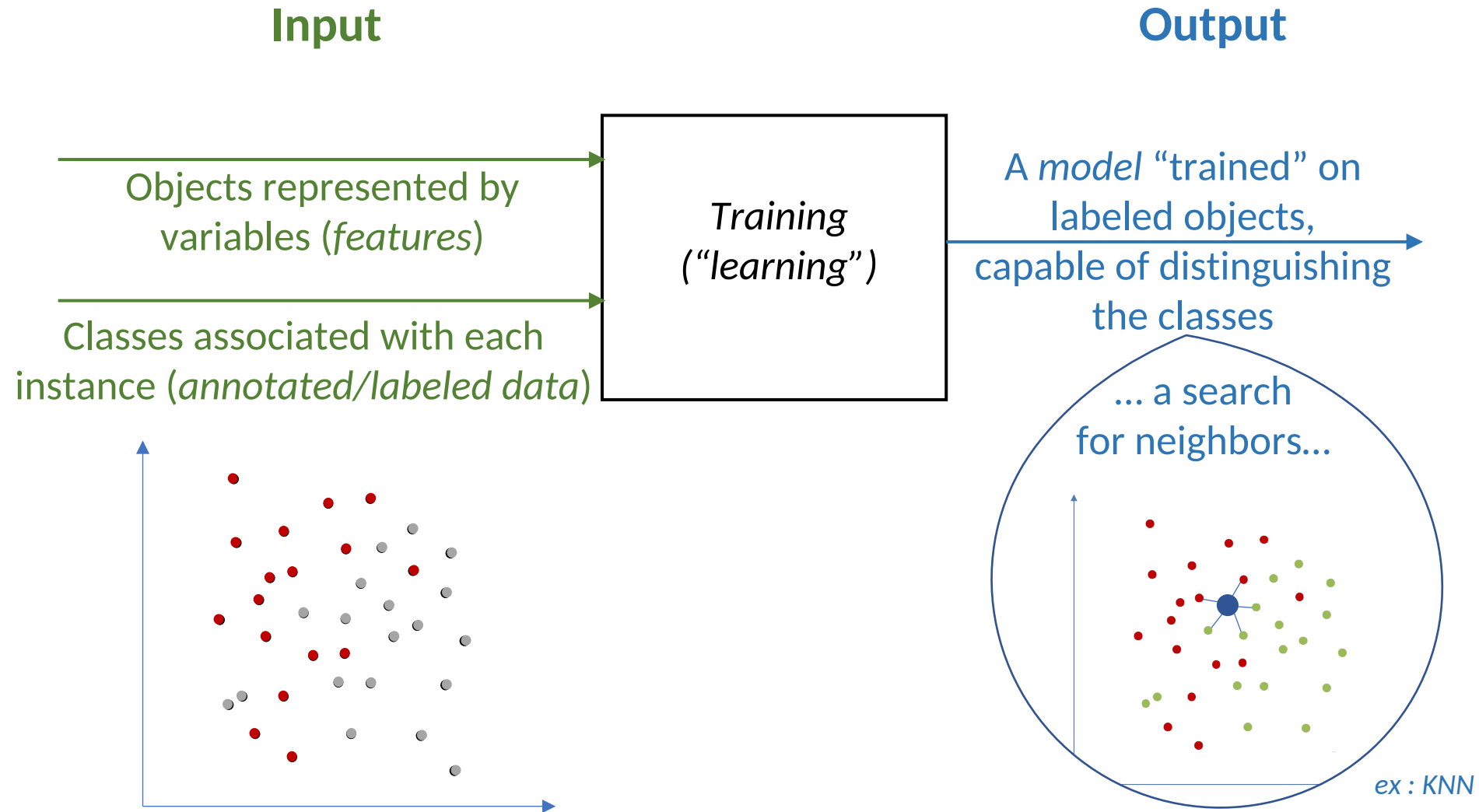


# (What is supervised learning?)



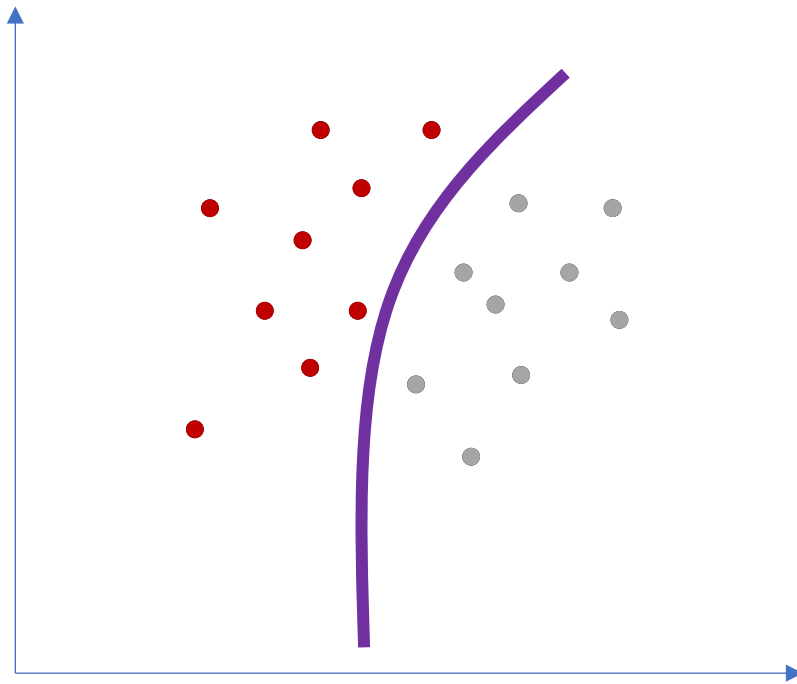


# (What is supervised learning?)



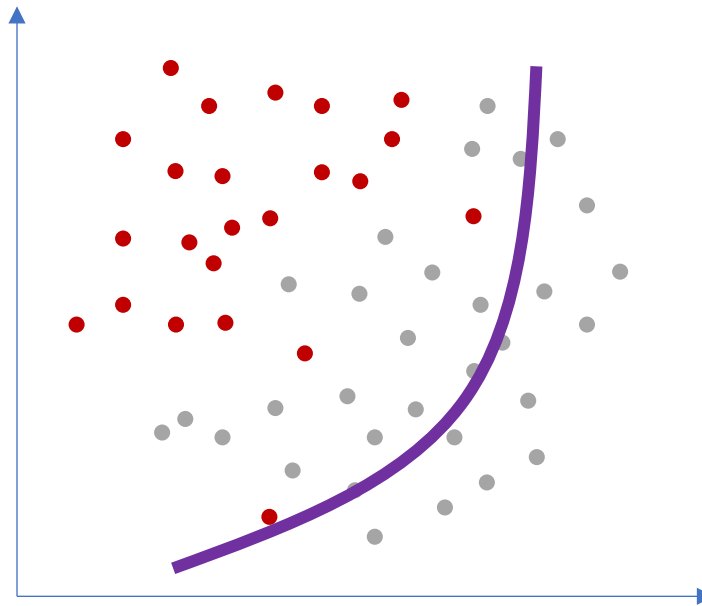
# (What is supervised learning?) Inference

- The model is then applied to new, “unlabeled” data.



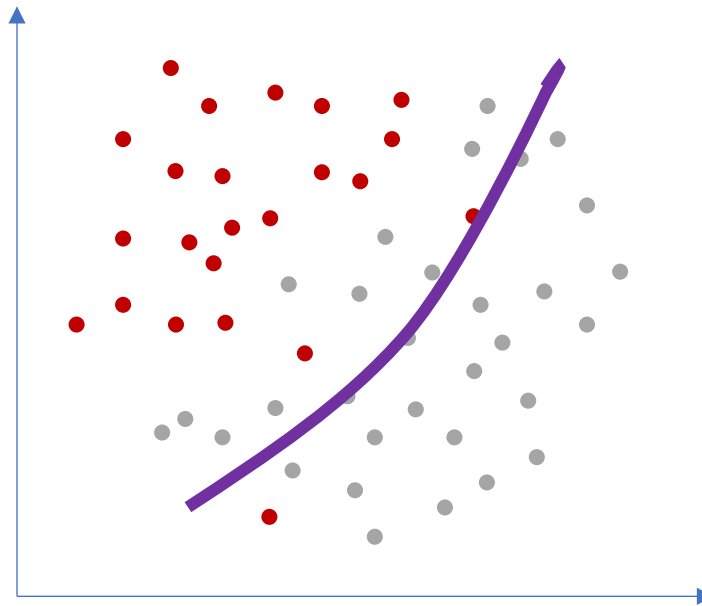
# (What is supervised learning?) Training

- Most often, learning (*training*) consists of minimizing the error committed by the system (*cost function* or *objective function*) by refining its parameters.
- Often an iterative process



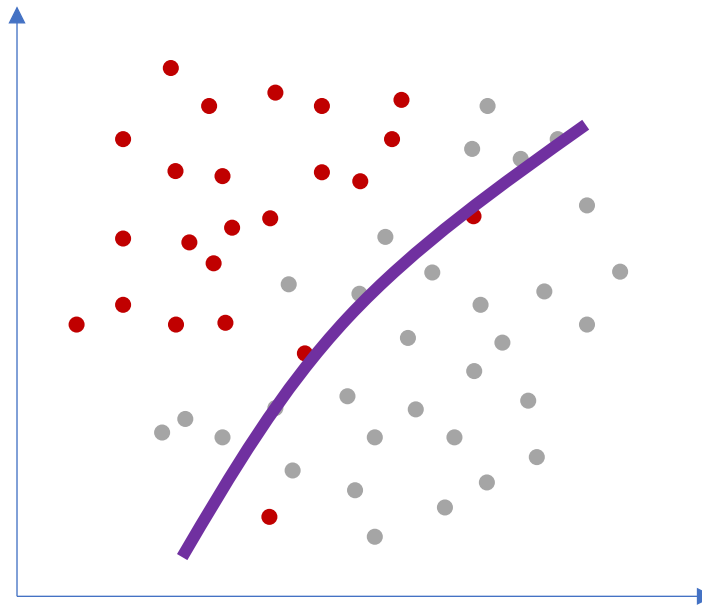
# (What is supervised learning?) Training

- Most often, learning (*training*) consists of minimizing the error committed by the system (*cost function* or *objective function*) by refining its parameters.
- Often an iterative process



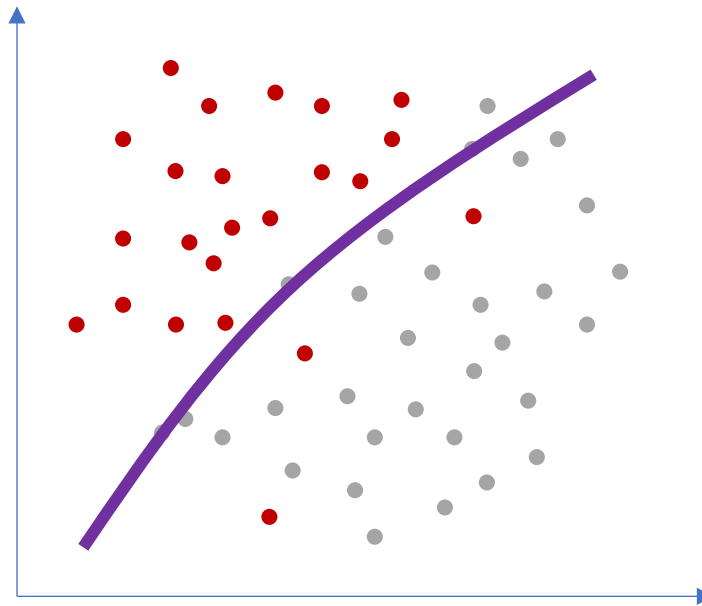
# (What is supervised learning?) Training

- Most often, learning (*training*) consists of minimizing the error committed by the system (*cost function* or *objective function*) by refining its parameters.
- Often an iterative process



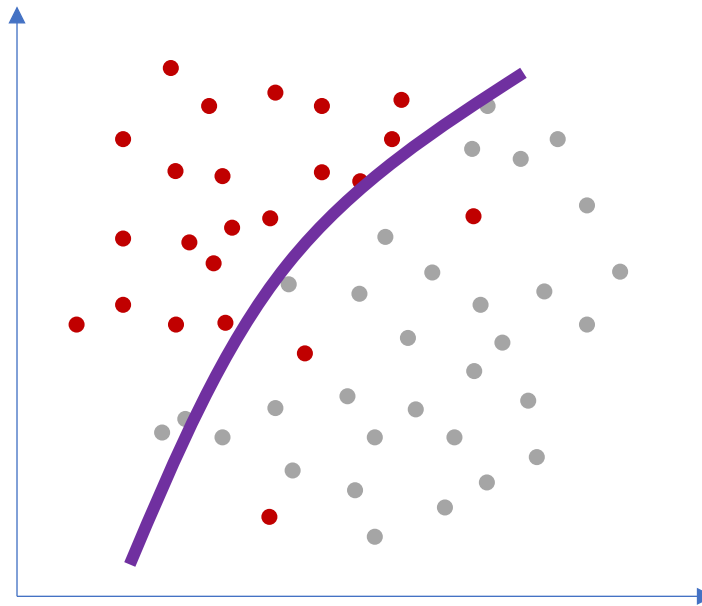
# (What is supervised learning?) Training

- Most often, learning (*training*) consists of minimizing the error committed by the system (*cost function* or *objective function*) by refining its parameters.
- Often an iterative process



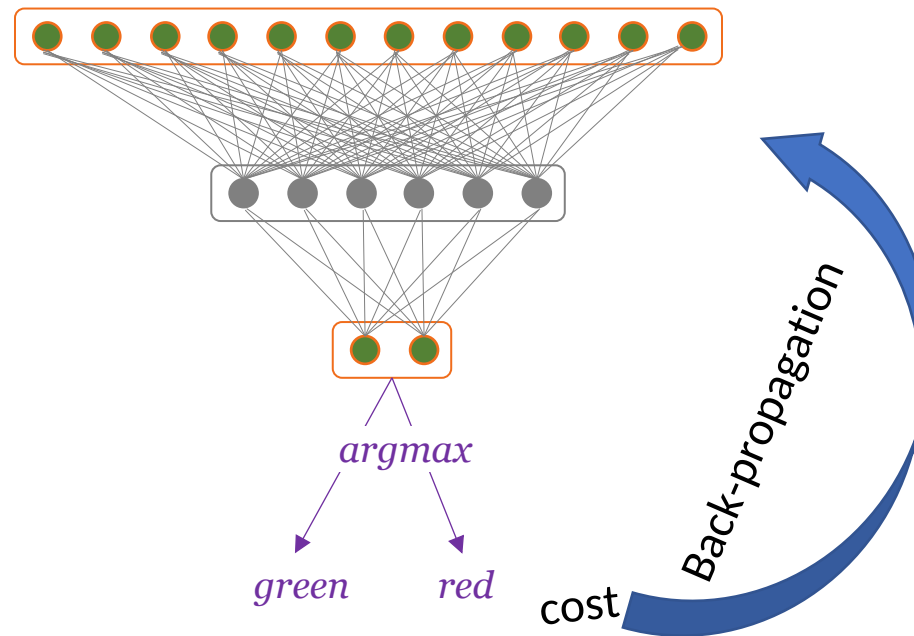
# (What is supervised learning?) Training

- Most often, learning (*training*) consists of minimizing the error committed by the system (*cost function* or *objective function*) by refining its parameters.
- Often an iterative process



# (What is supervised learning?) Training

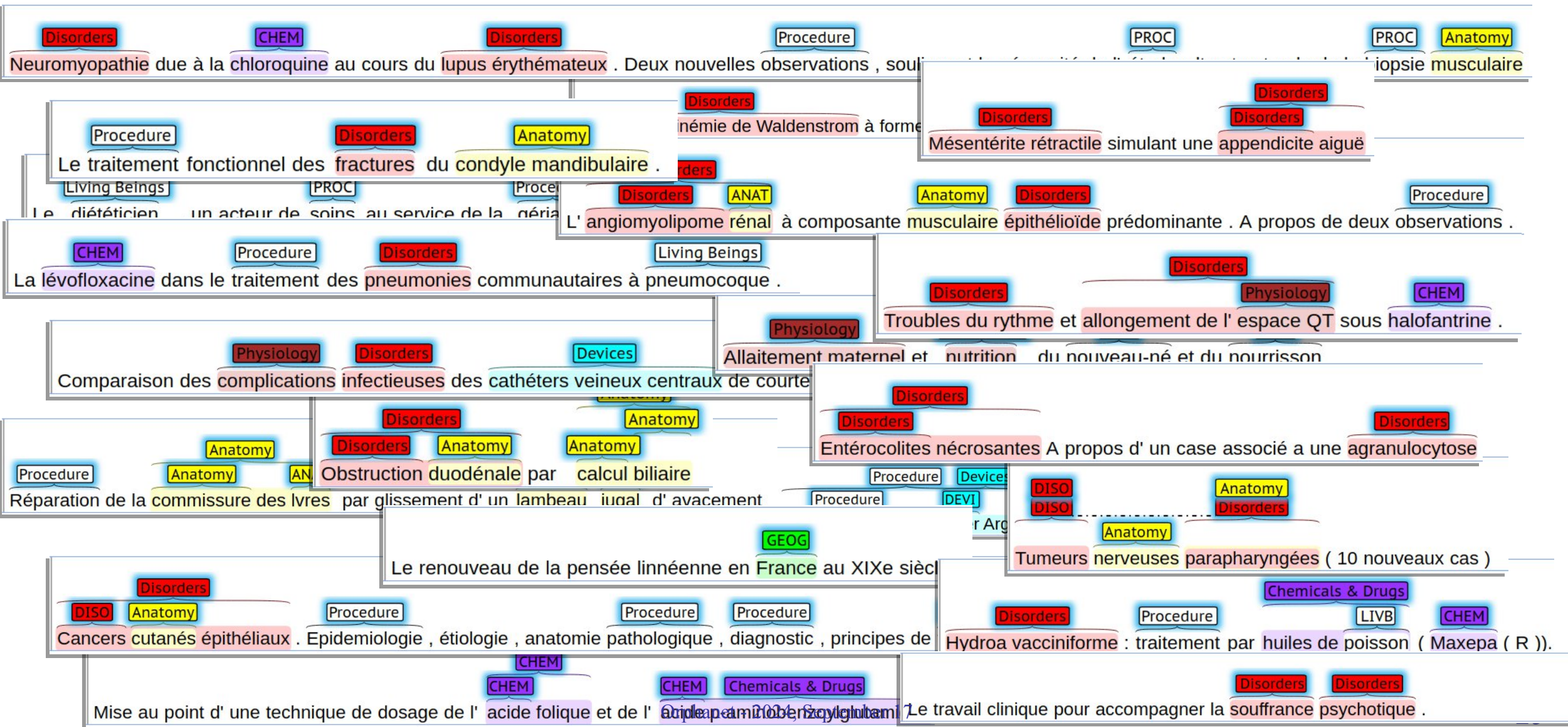
- Most often, learning (*training*) consists of minimizing the error committed by the system (*cost function* or *objective function*) by refining its parameters.
- Often an iterative process





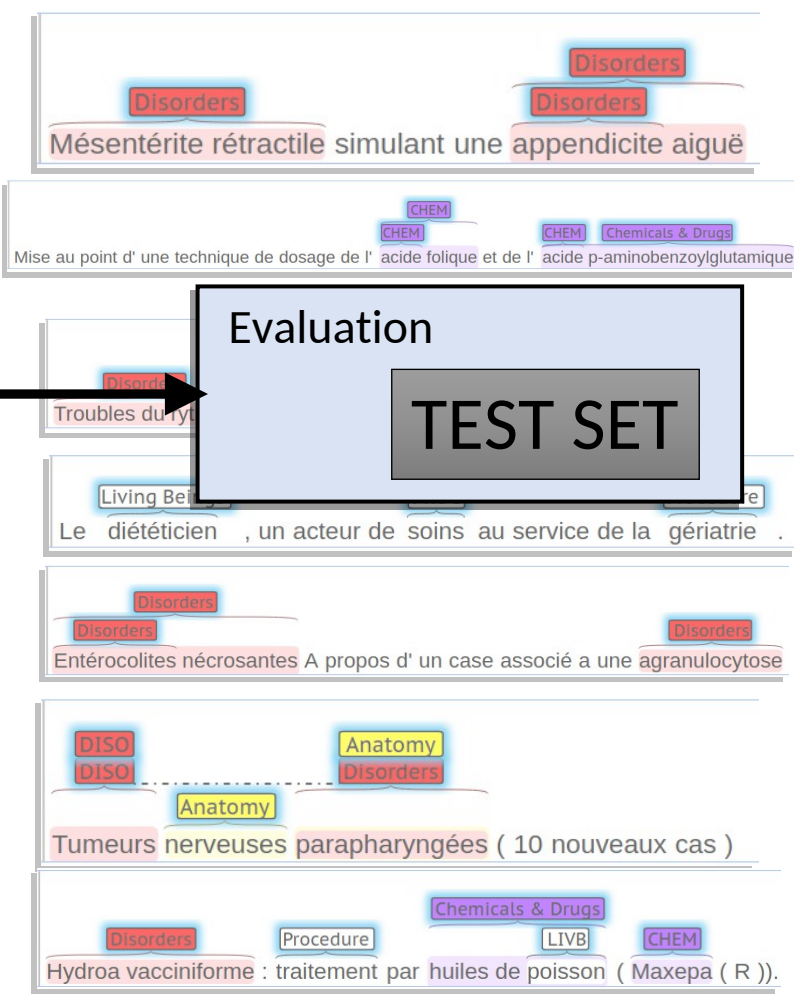
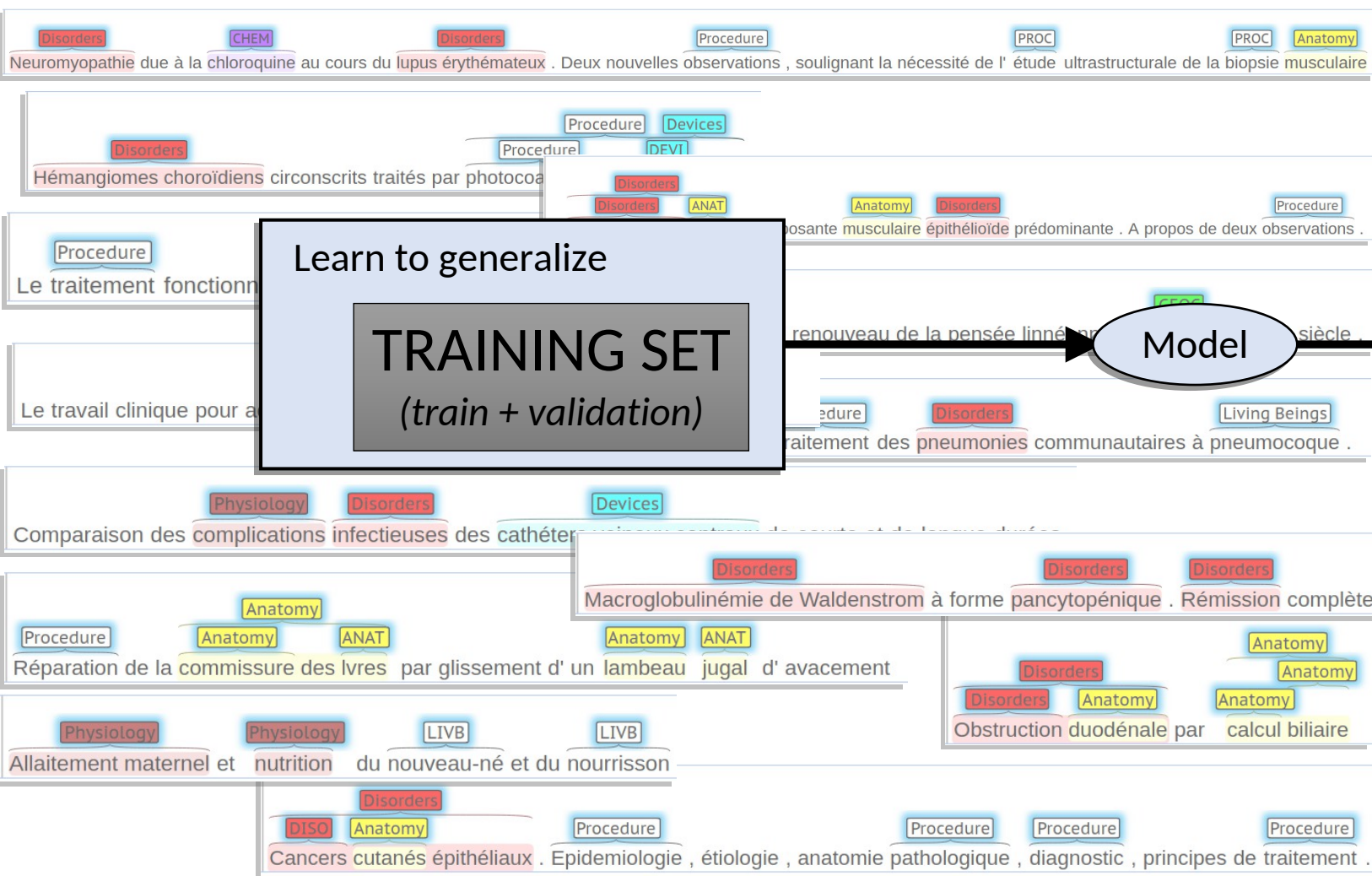
# Supervised learning systems

Névéol, A., Grouin, C., Leixa, J., Rosset, S., & Zweigenbaum, P. (2014).  
The Quaero French Medical Corpus: A Resource for Medical Entity Recognition and Normalization.



# Supervised learning systems

Névéol, A., Grouin, C., Leixa, J., Rosset, S., & Zweigenbaum, P. (2014).  
The Quaero French Medical Corpus: A Resource for Medical Entity Recognition and Normalization.



# Supervised learning systems: Pros & Cons

	Rules	Supervised learning
General performance	Yellow	Green
Ease of implementation	Green	Green
Need for human expertise	Yellow	Red
Explainability	Green	Yellow
Material resources	Green	Yellow
Energy consumption	Green	Yellow
Ease of maintenance	Yellow	Yellow
Generalization to a different problem/context	Red	Red

# Methods

1. Rule-based systems
2. Supervised learning systems
- 3. Generative, large language models**
4. Retrieval-Augmented Generation

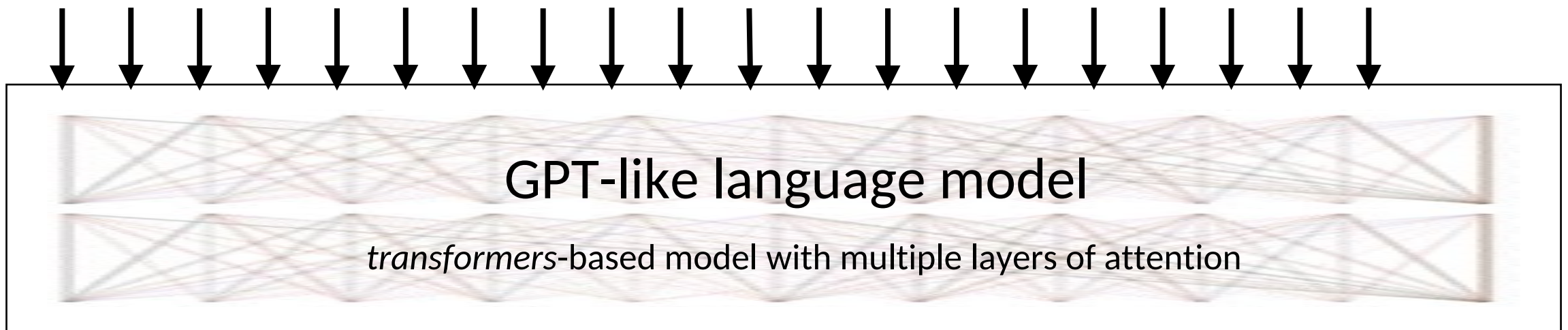
# Large Language Models

- Generative large language models (LLMs) are trained to produce human-like text.
- They simulate human understanding by predicting and generating text based on the input they receive.
- They can be finetuned for different tasks, e.g rewriting a prompt in another styles.

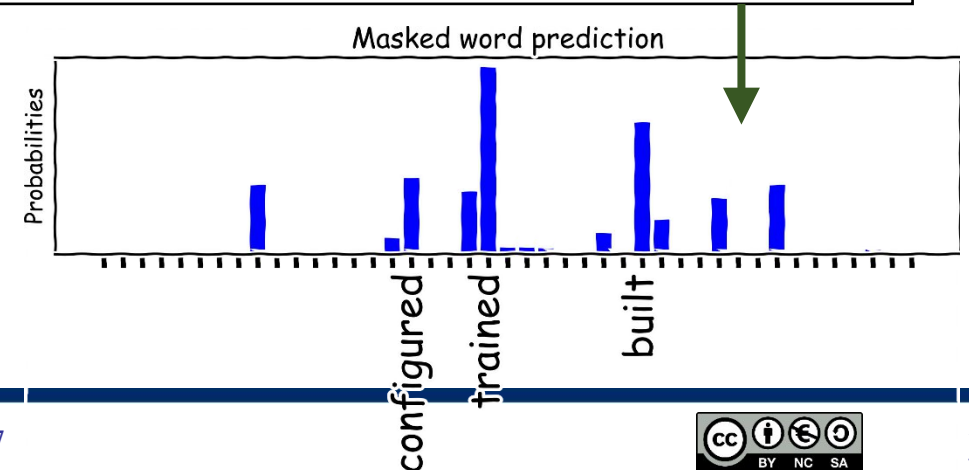
# Auto-regressive language models (*decoder-style, e.g. GPT*)

User: What is a large language model?

Assistant: A large language model is a type of artificial intelligence (AI) model designed to understand and generate human-like language. These models are



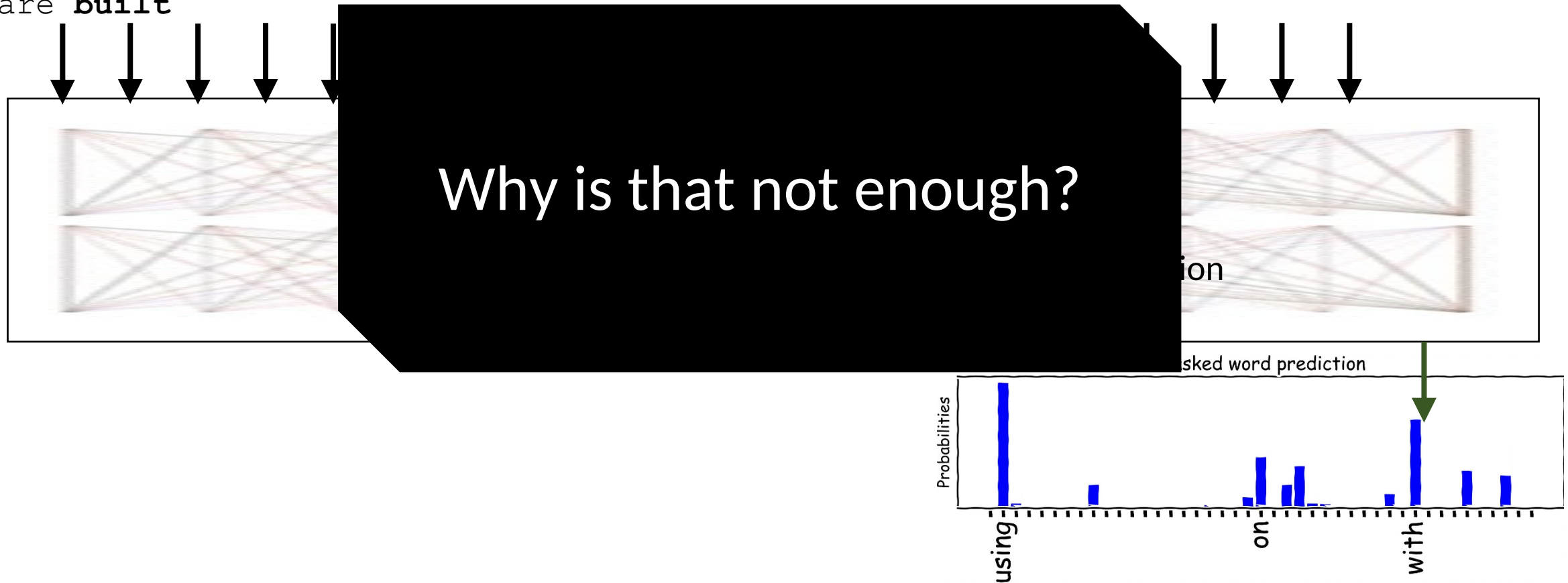
Next tokens are randomly selected from the predicted distribution, introducing variability in the generated output



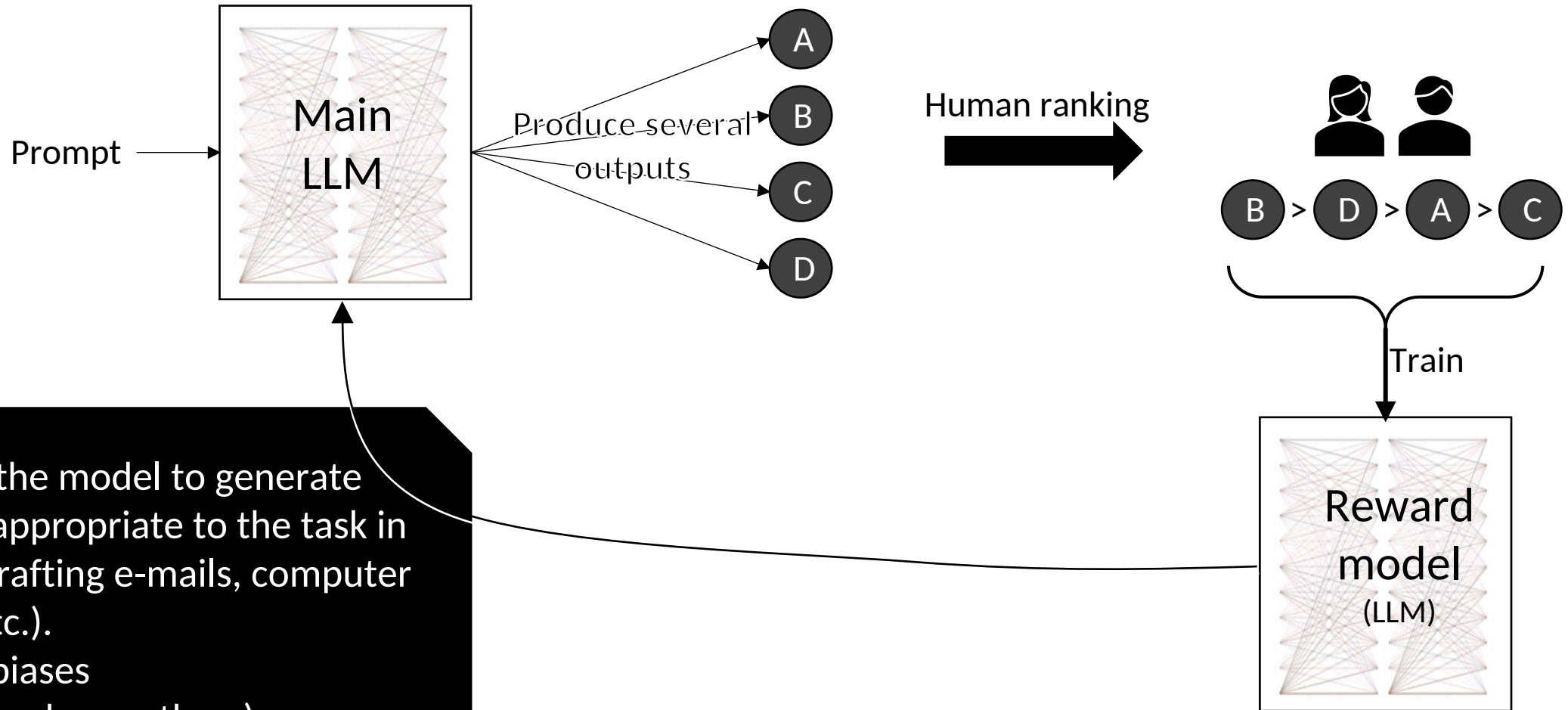
# Auto-regressive language models (*decoder-style, e.g. GPT*)

User: What is a large language model?

Assistant: A large language model is a type of artificial intelligence (AI) model designed to understand and generate human-like language. These models are **built**



# Reinforcement Learning with Human Feedback



- Adapts the model to generate results appropriate to the task in hand (drafting e-mails, computer code, etc.).
- Avoids biases (and introduces others)



# Why are LLMs so good?

- LLMs are good at producing **relevant**, **well-written** and **convincing** information, thanks to:
  - **Representation** learning
  - **Contextual** understanding and attention mechanisms  
(capture long-range dependencies and relationships in the data)
  - **Scalability** and massive amount of training data  
(wide range of linguistic nuances and topics)
  - Massive finetuning with **human feedback**  
(produce the right kind of results for the task prompted by the user)
- LLMs are NOT good at:
  - **Factuality / knowledge**
  - **High precision**
  - **Humor, creativity, originality**

# How can we use LLMs?

- Direct prompt: *“How do genetic mutations in the ALMS1 gene contribute to the pathophysiology of Alström Syndrome?”*
  - ▶ only knowledge from the pretraining
- Prompt with persona: *“You’re an assistant specialized in research on rare diseases. How do genetic mutations in the ALMS1 gene contribute to the pathophysiology of Alström Syndrome?”*
  - ▶ allows to guide the answer and its style
- Prompt with document: *“Given this document, how do genetic mutations in the ALMS1 gene contribute to the pathophysiology of Alström Syndrome?”*
  - ▶ knowledge from the pretraining + supporting document

# How can we use LLMs?

- Few-shot prompting: “

*Classify the following rare diseases into their appropriate categories: Genetic Disorder or Neurodegenerative Disorder*

*Input: "Alström Syndrome"*

*Output: "Genetic Disorder"*

*Input: "Batten Disease"*

*Output: "Neurodegenerative Disorder"*

*Input: "Huntington's Disease"*

*Output: "Neurodegenerative Disorder"*

*Input: "Marfan Syndrome"*

*Output: "Genetic Disorder"*

*Input: "Parkinson's Disease"*

*Output: "Neurodegenerative Disorder"*

*Input: "Cystic Fibrosis"*

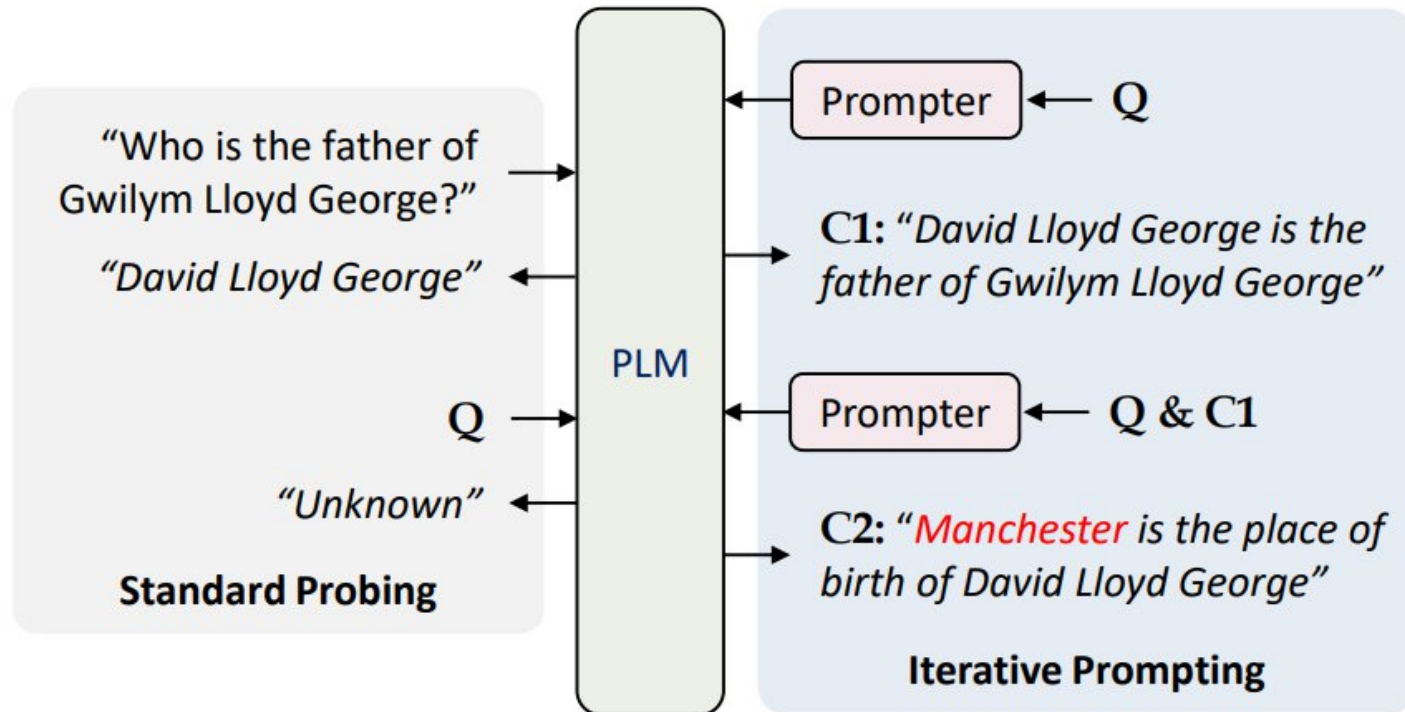
*Output:*

”

# How can we use LLMs?

- Prompt chaining:

Q: "What is the place of birth of Gwilym Lloyd George's father?"  
(Answer: **Manchester**)



Wang, B., Deng, X., & Sun, H. (2022). Iteratively Prompt Pre-trained Language Models for Chain of Thought. *Conference on Empirical Methods in Natural Language Processing*.

# How can we use LLMs?

- Chain of thoughts:

**StrategyQA**

Q: Yes or no: Would a pear sink in water?

A: The density of a pear is about  $0.6 \text{ g/cm}^3$ , which is less than water. Thus, a pear would float. So the answer is no.

**Date Understanding**

Q: The concert was scheduled to be on 06/01/1943, but was delayed by one day to today. What is the date 10 days ago in MM/DD/YYYY?

A: One day after 06/01/1943 is 06/02/1943, so today is 06/02/1943. 10 days before today is 05/23/1943. So the answer is 05/23/1943.

**Sports Understanding**

Q: Is the following sentence plausible? "Joao Moutinho caught the screen pass in the NFC championship."

A: Joao Moutinho is a soccer player. The NFC championship is part of American football, not soccer. So the answer is no.

**SayCan (Instructing a robot)**

Human: How would you bring me something that isn't a fruit?

Explanation: the user wants something to eat that isn't a fruit. An energy bar is not a fruit, so I will bring the user an energy bar.

Plan: 1. find(energy bar) 2. pick(energy bar) 3. find(user) 4. put(energy bar) 5. done().

**Last Letter Concatenation**

Q: Take the last letters of the words in "Lady Gaga" and concatenate them.

A: The last letter of "Lady" is "y". The last letter of "Gaga" is "a". Concatenating them is "ya". So the answer is ya.

**Coin Flip (state tracking)**

Q: A coin is heads up. Maybelle flips the coin. Shalonda does not flip the coin. Is the coin still heads up?

A: The coin was flipped by Maybelle. So the coin was flipped 1 time, which is an odd number. The coin started heads up, so after an odd number of flips, it will be tails up. So the answer is no.

Wei, J., Wang, X., Schuurmans, D., Bosma, M., Chi, E.H., Xia, F., Le, Q., & Zhou, D. (2022). Chain of Thought Prompting Elicits Reasoning in Large Language Models. *Conference on Neural Information Processing Systems*.

# How can we use LLMs?

- Self-Consistency
- Generated Knowledge Prompting
- Tree of Thoughts (ToT)
- Automatic Prompt Engineer (APE)
- Active-Prompt
- Directional Stimulus Prompting
- Program-Aided Language Models
- ... (see for example <https://www.promptingguide.ai/>)

# Generative LLMs: Pros & Cons

	Rules	Supervised learning	Promise of gLLMs
General performance	Yellow	Green	Grey with ?
Ease of implementation	Green	Green	Green
Need for human expertise	Yellow	Red	Green
Explainability	Green	Yellow	Yellow
Material resources	Green	Yellow	Red
Energy consumption	Green	Yellow	Red
Ease of maintenance	Yellow	Yellow	Yellow
Generalization to a different problem/context	Red	Red	Yellow



A thorough evaluation of these models  
is always necessary!

Whatever the method (rules, ML, LLMs...),  
build a test set for the evaluation

Do not prompt public or API-based LLMs  
with sensitive data!



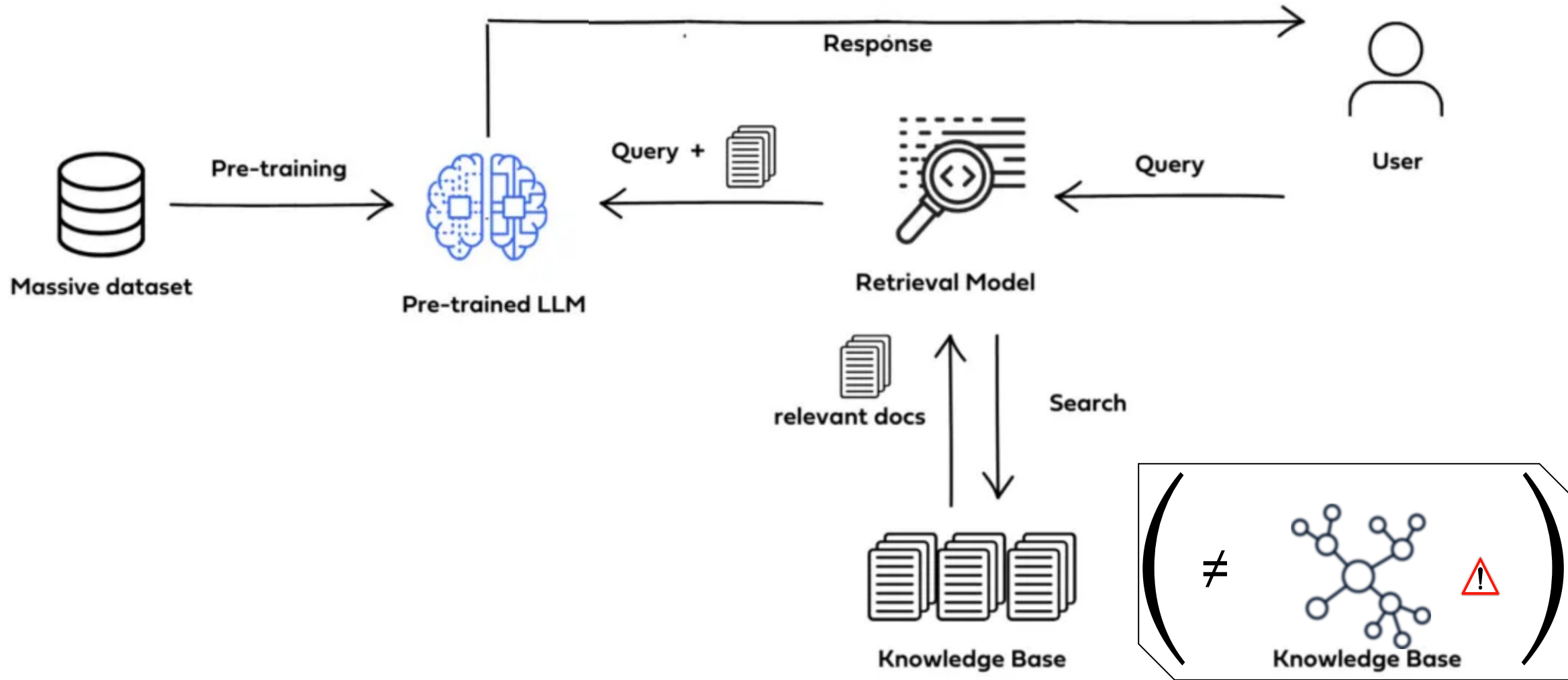


# Methods

1. Rule-based systems
2. Supervised learning systems
3. Generative, large language models
- 4. Retrieval-Augmented Generation**

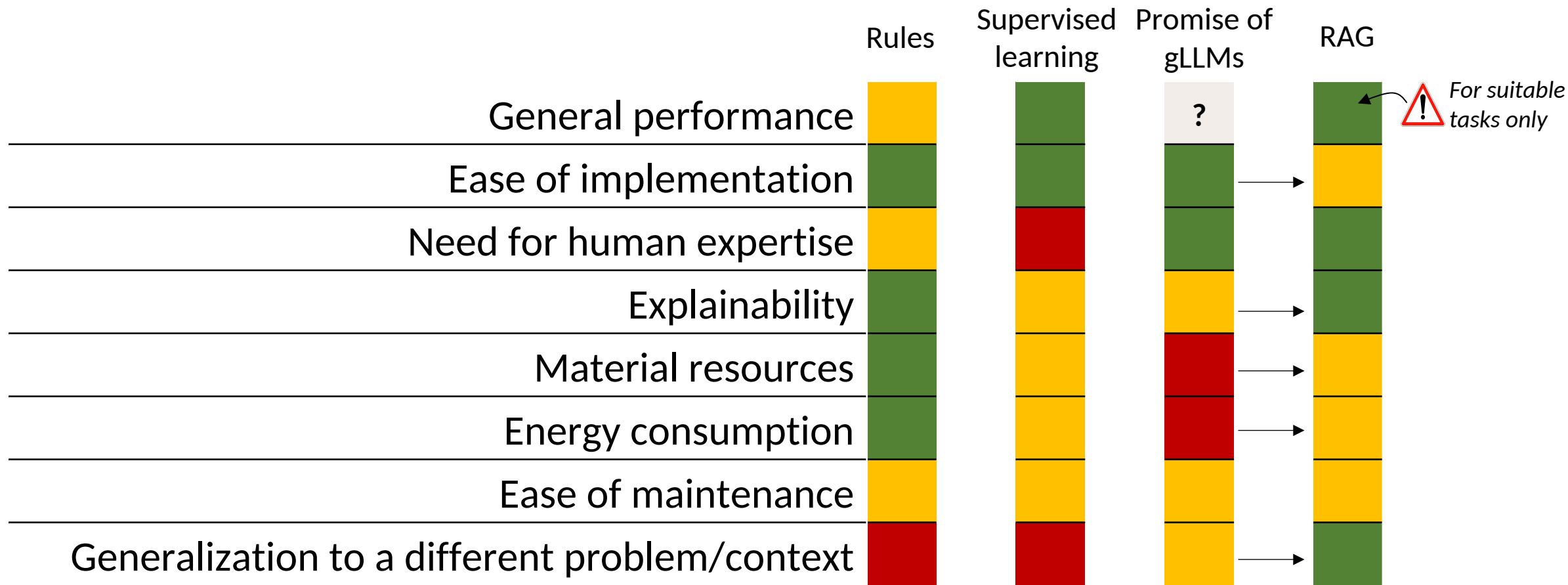
*The new search engine*

# Retrieval-Augmented Generation

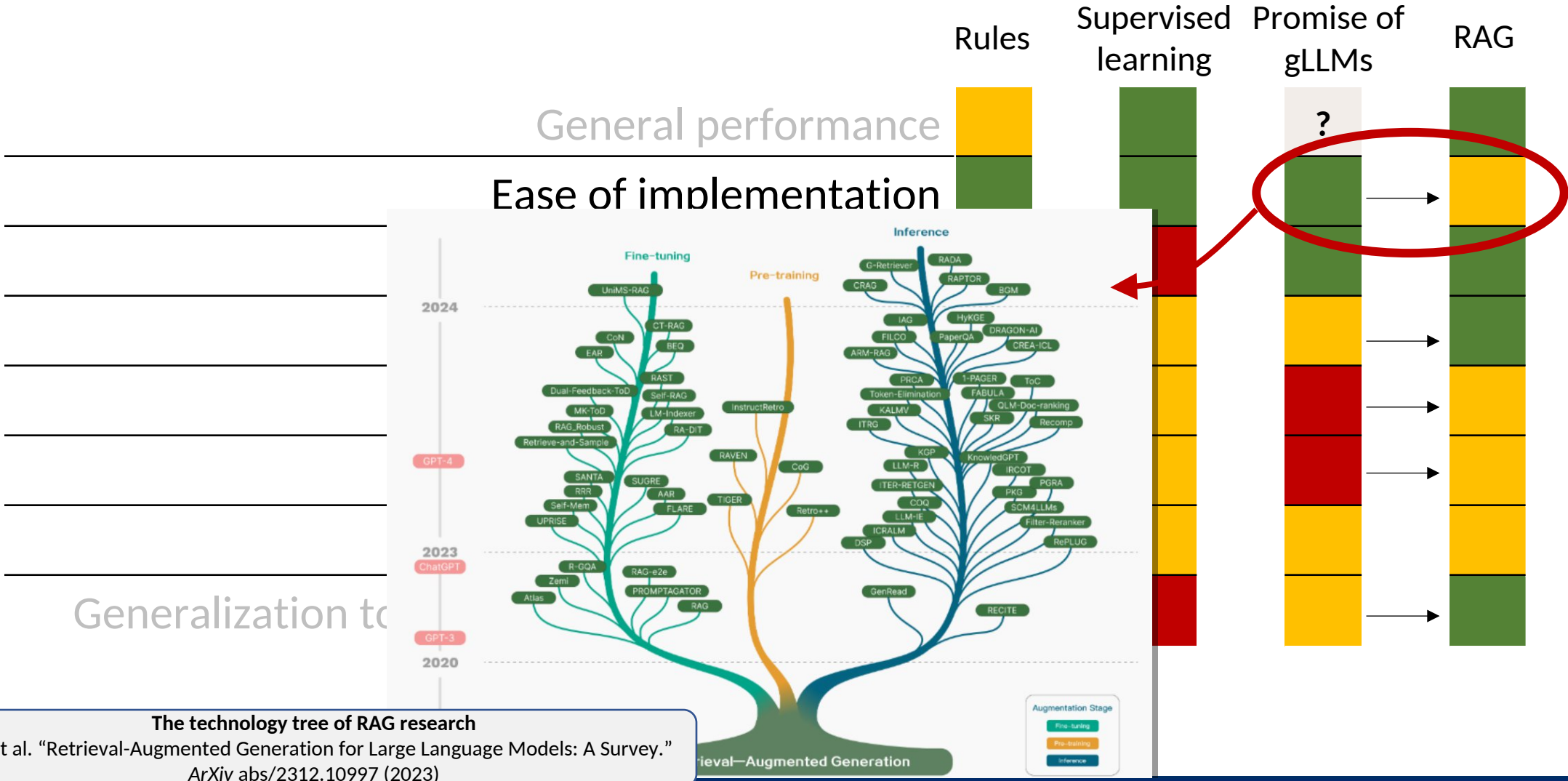


<https://medium.com/@krtarunsingh/introduction-to-retrieval-augmented-generation-rag-and-its-transformative-role-in-ai-c07e35da7f01>

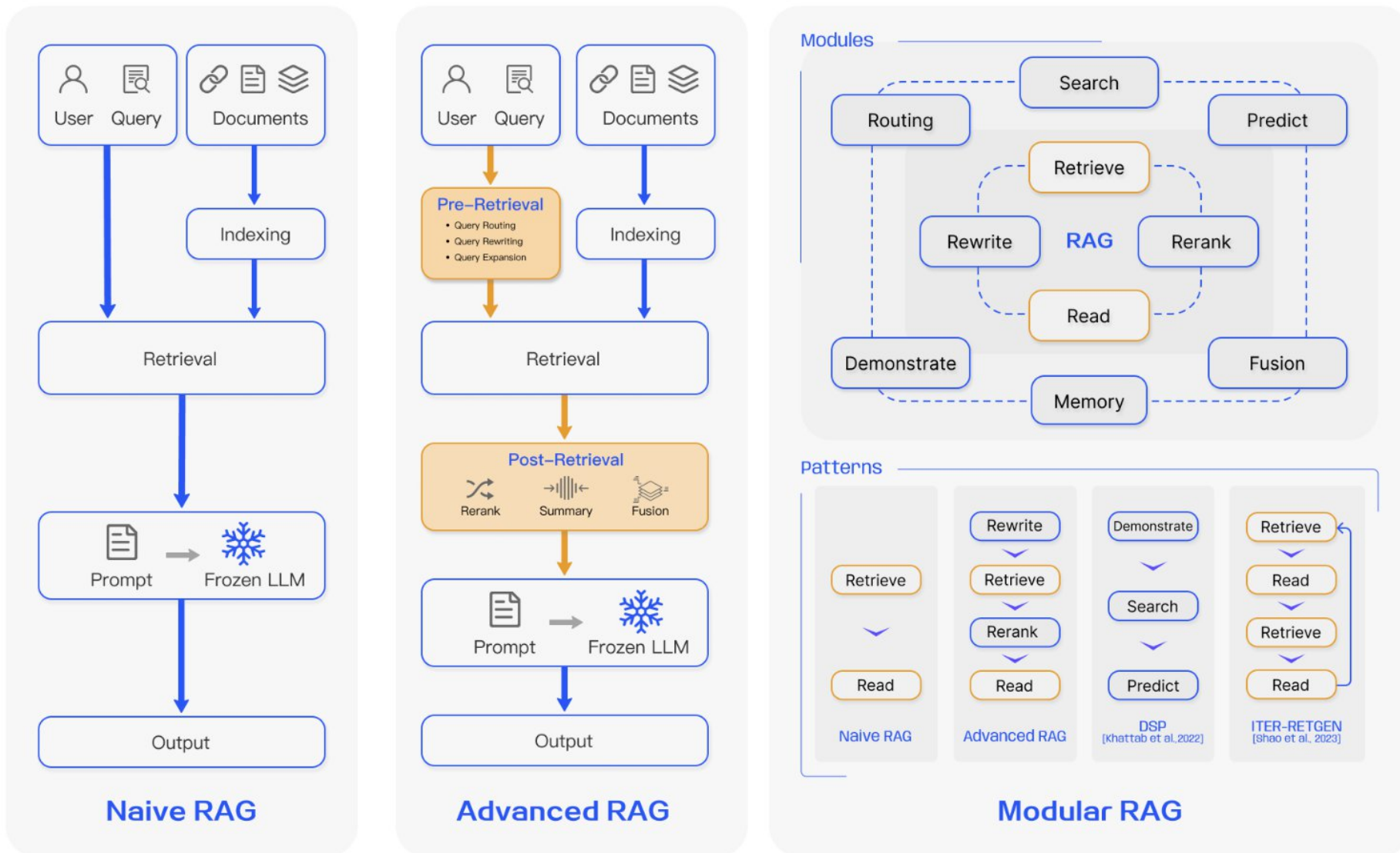
# RAG: Pros & Cons



# RAG: Pros & Cons

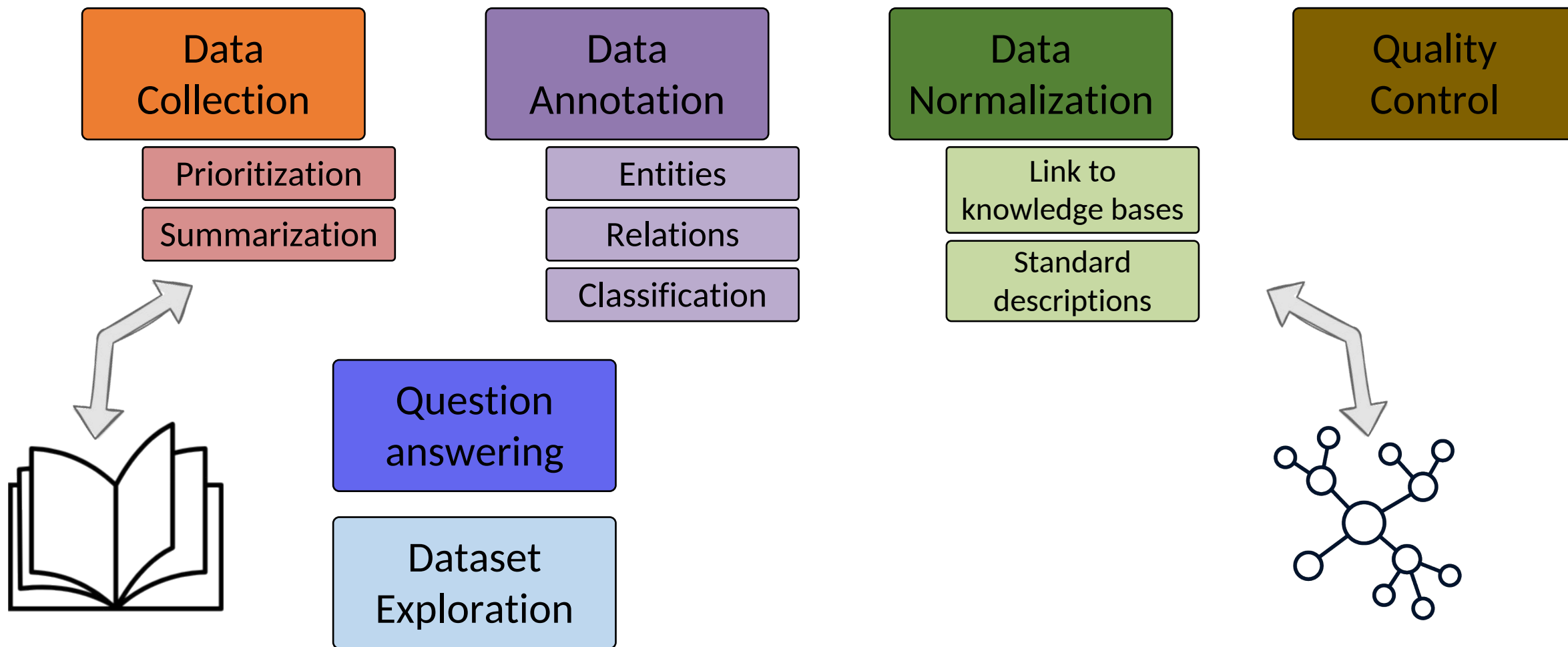


# RAG

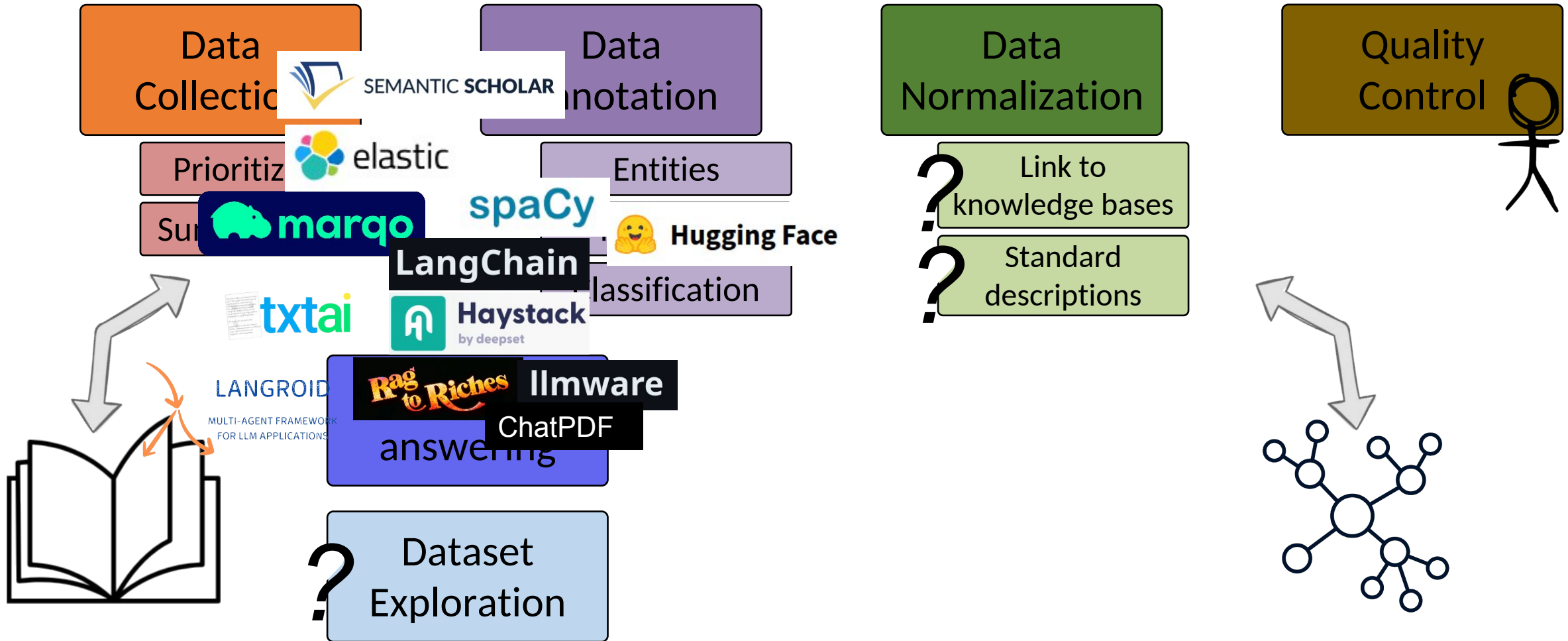


# NLP for scientific surveillance

# Positioning NLP and LLMs in the process of curation



# Positioning NLP and LLMs in the process of curation





# Positioning NLP and LLMs in the process of curation

## Knowledge Navigator: LLM-guided Browsing Framework for Exploratory Search in Scientific Literature, U. Katz et al. 2024



Tool Use In Animals



Processed 1000 Scientific Papers

### Neural Mechanisms and Cognitive Processes

- Neural Mechanisms in Humans and Animals (9)
- Development and Cognitive Understanding (12)
- Tool Use and Body Space Representation (14)

### Tool Use in Primates

- Stone Tool Use in Non-Human Primates (12)
- Social Learning of Tool Use in Primates. (13)
- Tool Use in Chimpanzees (13)

### Comparative and Evolutionary Perspectives

- Investigation and Definitions (10)
- Tool Use in Early Hominin Evolution (17)
- Adaptation, Behavior, and Evolution (10)

### Tool Use in Birds

- Tool Use in New Caledonian Crows (20)
- Tool Use in Corvids (10)

Dataset  
Exploration

# Positioning NLP and LLMs in the process of curation

**LLM-assisted Knowledge Graph Engineering: Experiments with ChatGPT**, LP Meyer et al, 2023

**Iterative Zero-Shot LLM Prompting for Knowledge Graph Construction**, S Carta et al, 2023

**Let's Chat to Find the APIs: Connecting Human, LLM and Knowledge Graph through AI Chain**, Q Huang et al, 2023

**Knowledge Graph Prompting for Multi-Document Question Answering**, Y Wang et al, 2023

**Enhancing Knowledge Graph Construction Using Large Language Models**, M Trajanoska et al, 2023

**Exploring Large Language Models for Knowledge Graph Completion**, L. Yao, 2023

**Knowledge Graph Large Language Model (KG-LLM) for Link Prediction**, D. Shu, 2024

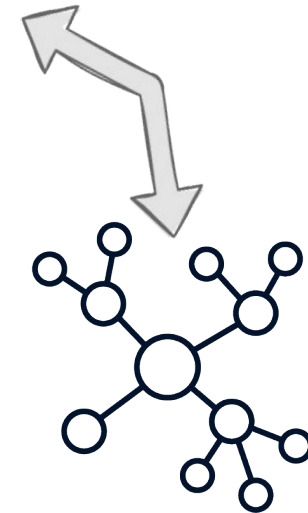
...

**Data Normalization**

Link to knowledge bases

Standard descriptions

**“Large Language Models struggle to learn long-tail knowledge”**  
(N. Kandpal et al, 2023)



# Knowledge Prompting: How Knowledge Engineers Use Large Language Models

ELISAVET KOUTSIANA\*, King's College London, United Kingdom

JOHANNA WALKER\*, King's College London, United Kingdom

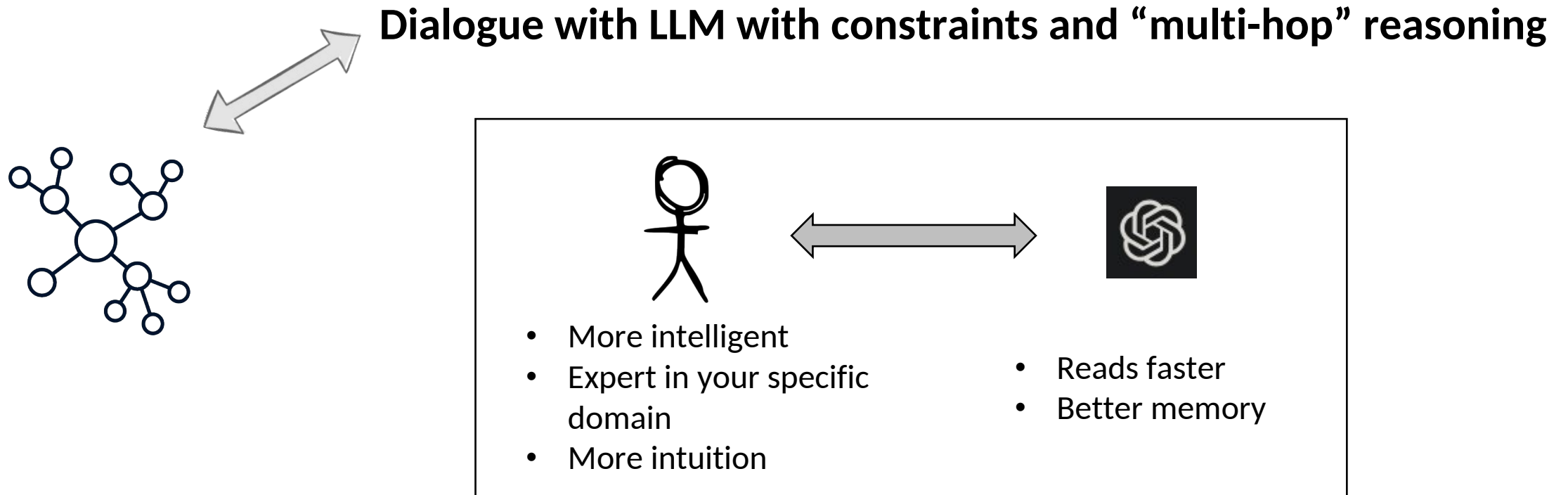
MICHELLE NWACHUKWU, King's College London, United Kingdom

ALBERT MERONÓ-PEÑUELA, King's College London, United Kingdom

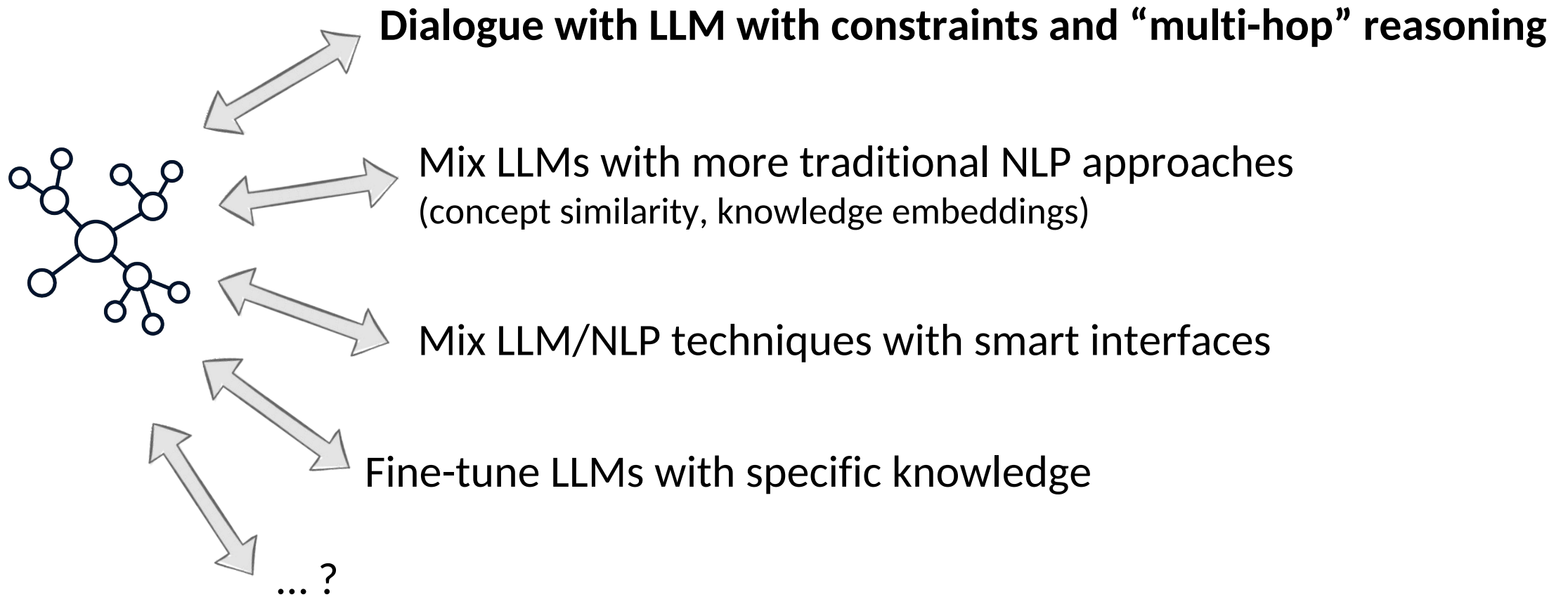
ELENA SIMPERL\*, King's College London, United Kingdom

Despite many advances in knowledge engineering (KE), challenges remain in areas such as engineering knowledge graphs (KGs) at scale, keeping up with evolving domain knowledge, multilingualism, and multimodality. Recently, KE has used LLMs to support semi-automatic tasks, but the most effective use of LLMs to support knowledge engineers across the KE activities is still in its infancy. To explore the vision of LLM copilots for KE and change existing KE practices, we conducted a multimethod study during a KE hackathon. We investigated participants' views on the use of LLMs, the challenges they face, the skills they may need to integrate LLMs into their practices, and how they use LLMs responsibly. We found participants felt LLMs could contribute to improving efficiency when engineering KGs, but presented increased challenges around the already complex issues of evaluating the KE tasks. We discovered prompting to be a useful but undervalued skill for knowledge engineers working with LLMs, and note that natural language processing skills may become more relevant across more roles in KG construction. Integrating LLMs into KE tasks needs to be mindful of potential risks and harms related to responsible AI. Given the limited ethical training, most knowledge engineers receive solutions such as our suggested 'KG cards' based on data cards could be a useful guide for KG construction. Our findings can support designers of KE AI copilots, KE researchers, and practitioners using advanced AI to develop trustworthy applications, propose new methodologies for KE and operate new technologies responsibly.

# How to link (indirectly) formal knowledge and LLMs



# How to link (indirectly) formal knowledge and LLMs

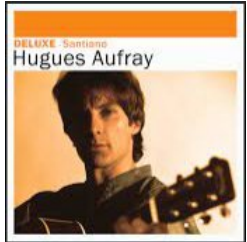


Thanks!

Questions ?

*Backup slides...*

# Words



D' y penser j' avais le cœur gros

En doublant les feux de Saint-Malo

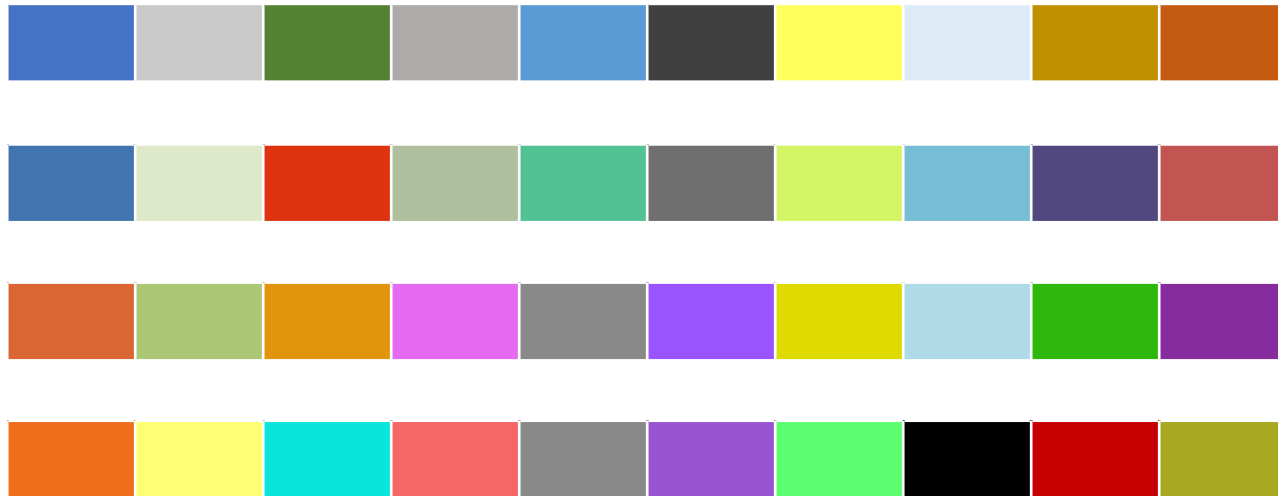


# "Bag of words"

a	0	...	...	...	...	gouvernées	0	pensée	0
à	0	coéternels	0	fanages	0	gouvernement	0	pensées	0
abaissa	0	coéternité	0	fanai	0	gouvernementabilité	0	pensement	0
abaissable	0	coéternités	0	fanaient	0	gouvernemental	0	pensements	0
abaissables	0	<b>cœur</b>	<b>1</b>	...	...	gouvernementale	0	pensent	0
abaissai	0	cœurs	0	fétuques	0	gouvernementalement	0	<b>penser</b>	<b>1</b>
...	...	coexista	0	fétus	0	gouvernementales	0	...	...
avachît	0	coexistai	0	<b>feu</b>	<b>1</b>	gouvernementalisme	0	pertuisaniers	0
avachîtes	0	...	...	feudataire	0	gouvernementalismes	0	pertuiser	0
avaient	0	doublâmes	0	feudataires	0	gouvernementaliste	0	perturba	0
<b>avais</b>	<b>1</b>	<b>doublant</b>	<b>1</b>	...	...	gouvernementaux	0	perturbai	0
avait	0	doublard	0	groomer	0	gouvernementiste	0	perturbaient	0
aval	0	doublards	0	grooms	0	gouvernementistes	0	...	0
avala	0	...	...	<b>gros</b>	<b>1</b>	gouvernementomane	0	phonocinèse	0
avalable	0	écardons	0	gros-bec	0	gouvernements	0	phonocinétique	0
...	...	écarlate	0	groschen	0	gouvernement	0	phonocontrôle	0
bains	0	écarlates	0	...		gouverner	0	phonocapteur	0
baïonette	0	écarquilla	0	gouttons	0	...	...	...	...

# Dense representation (*embeddings*)

- Word embeddings = vector representation of tokens
- Tokens with some degree of similarity 🗨️ close to each other in space



# Dense representation (*embeddings*)



- **Intuition 1.** Each word of a language is associated with a composition of hidden factors (often unintelligible)

e.g :     *cat* = 10 (*animal*) + 5 (*soft*) - 10 (*loyal*)  
          *dog* = 10 (*animal*) + 3 (*soft*) + 10 (*loyal*)

- **Intuition 2.** Distributional hypothesis

« You shall know a word by the company it keeps » (Firth, 1957)

Two words close in vector space = two words that often share similar contexts

e.g : *the ... scratches ; ... is a felid*

*occurrence (cat) ~ occurrence (tiger)*

$$\mathcal{W}_{cat} \cdot \mathcal{W}_{context} \sim \mathcal{W}_{tiger} \cdot \mathcal{W}_{context}$$

$$\mathcal{W}_{cat} \sim \mathcal{W}_{tiger}$$

(© Perceval Wajsbürt)

# Dogs and cats and tigers

The *cat* is sleeping on the couch

The *dog* is sleeping on the couch

The *cat* is running in the garden

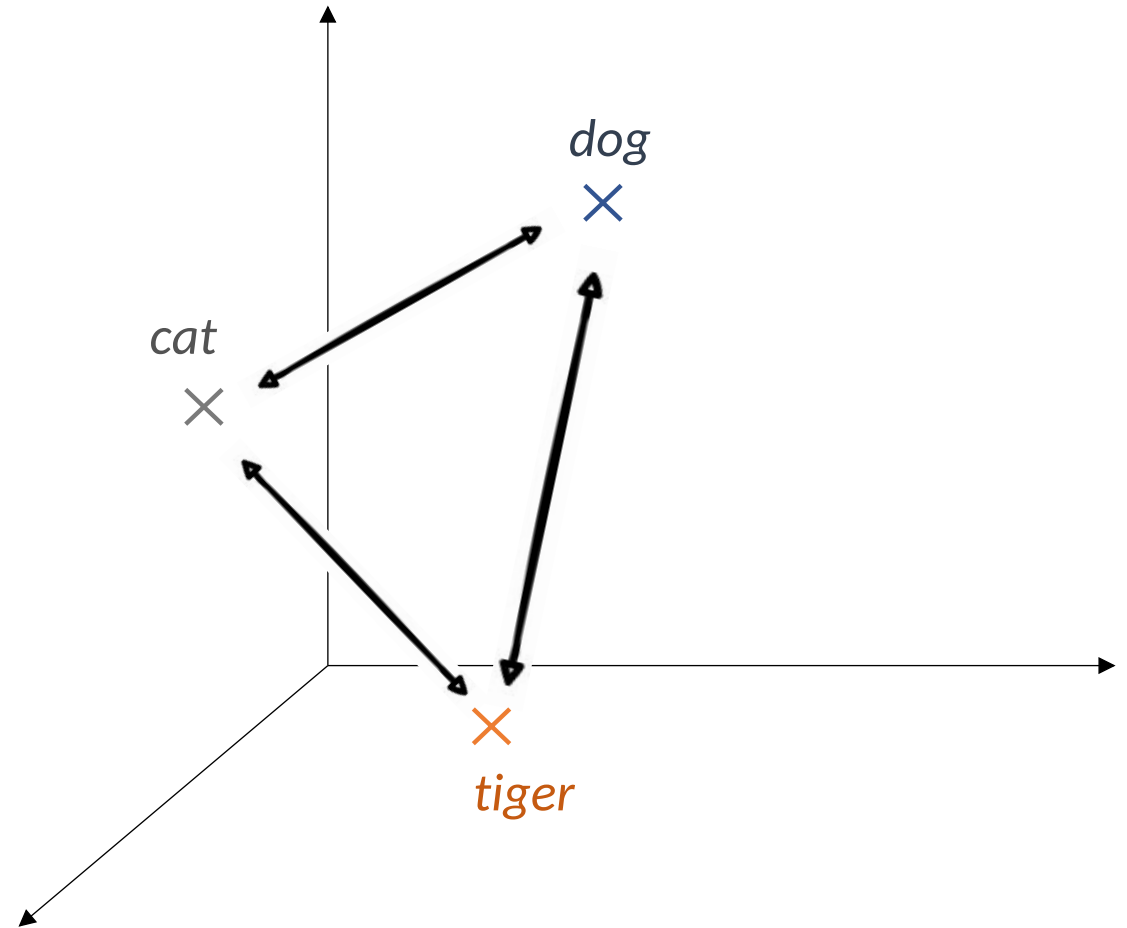
The *dog* is walking in the garden

The *tiger* is a big felid

The *cat* is a small felid

The *tiger* gnaws on an antelope bone

The *dog* gnaws on a chicken bone



# Dense representation (*embeddings*)

- Word embeddings = vector representation of tokens
- Tokens close to each other in space 🗣️ some degree of similarity between them

e.g. *bike* vs *bicycle*



# Dense representation (*embeddings*)

- Word embeddings = vector representation of tokens
- Tokens close to each other in space 🗨️ some degree of similarity between them



e.g. *tiger vs cat vs dog*

# Language models

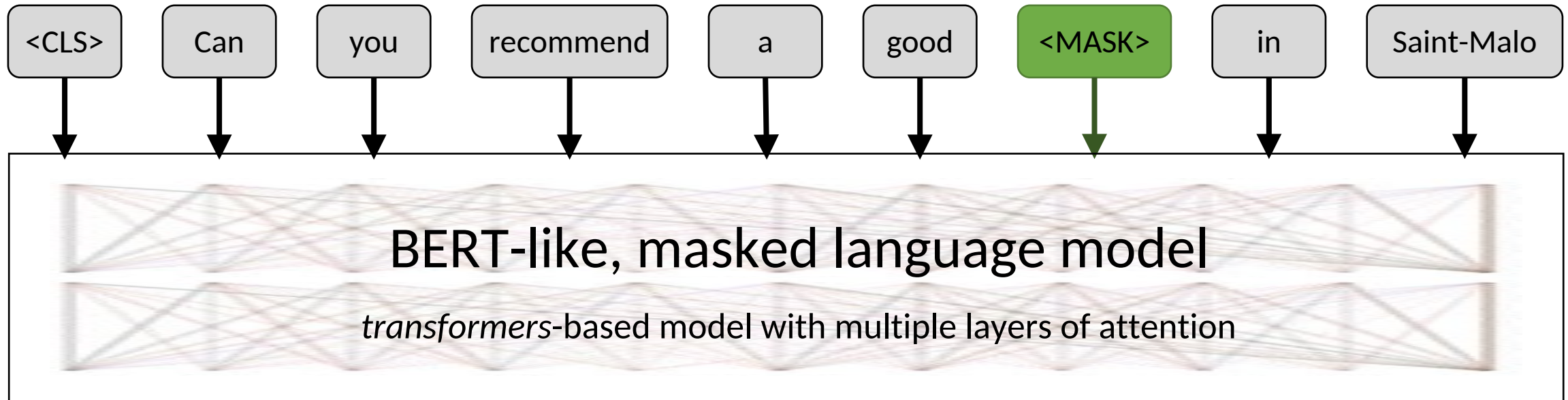
These models belong to the class of *language models*.

A language model is a model without manual supervision, with two different main training settings:

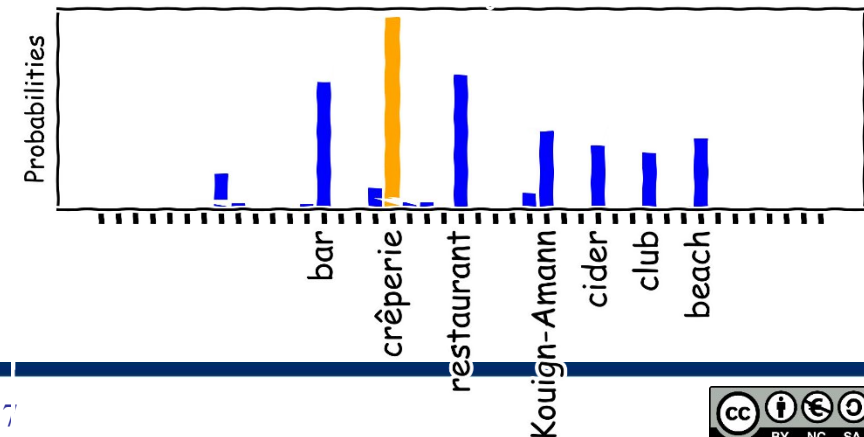
- Predict masked tokens within a sequence (*masked language models*, MLM, e.g. BERT).
- Predict next token (*autoregressive models*, e.g. GPT)

These models are at the heart of the vast majority of current NLP systems. They are known as LLMs (*Large Language Models*).

# Masked language models (encoder-style, e.g. BERT)

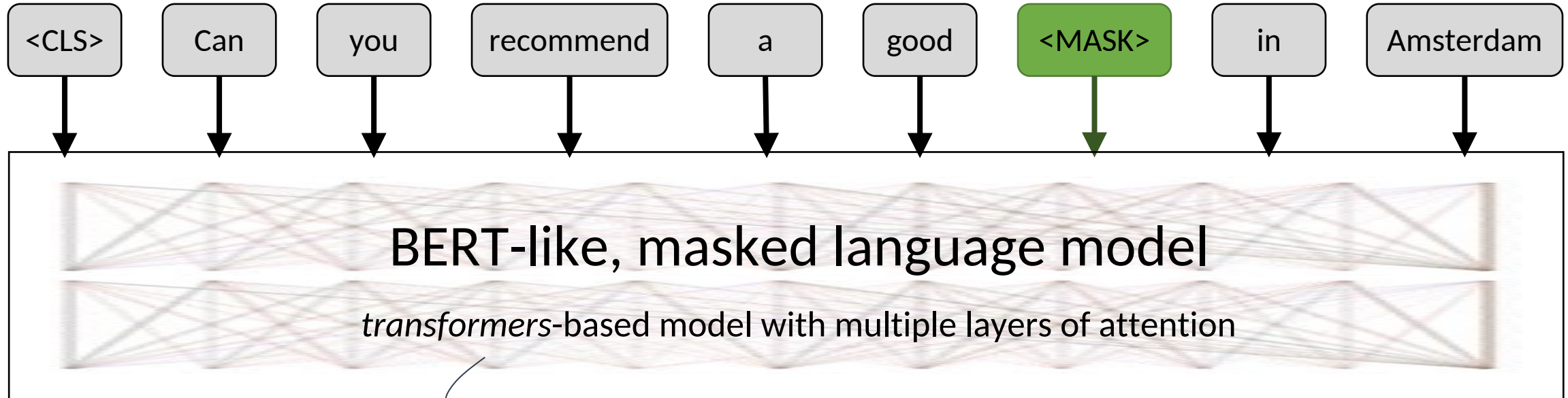


word prediction

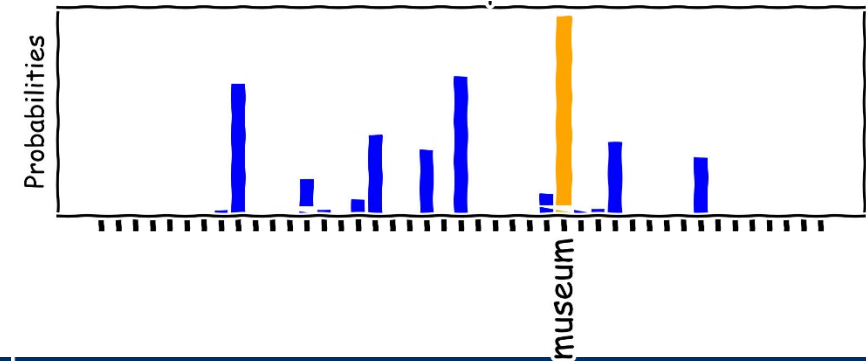




# Masked language models (*encoder-style, e.g. BERT*)



word prediction

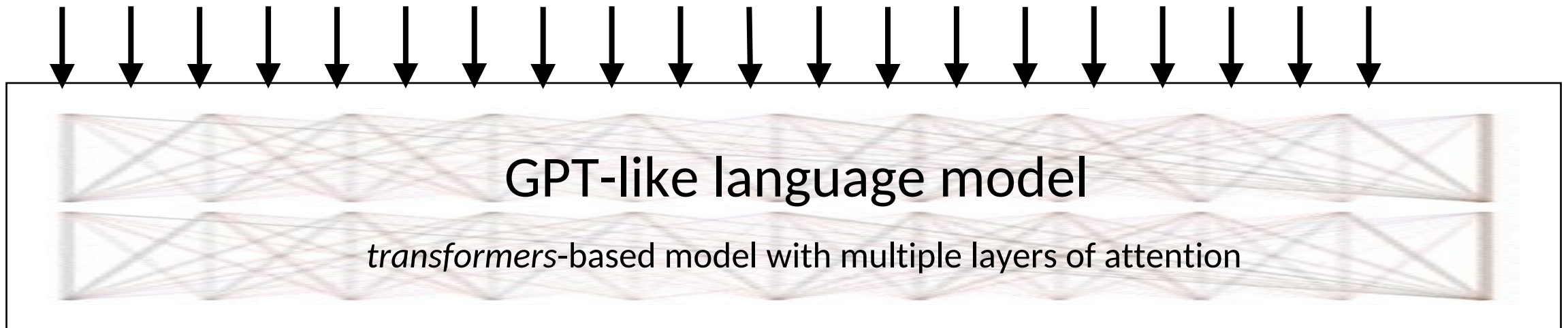


The attention mechanism makes it possible to consider the context for the prediction.

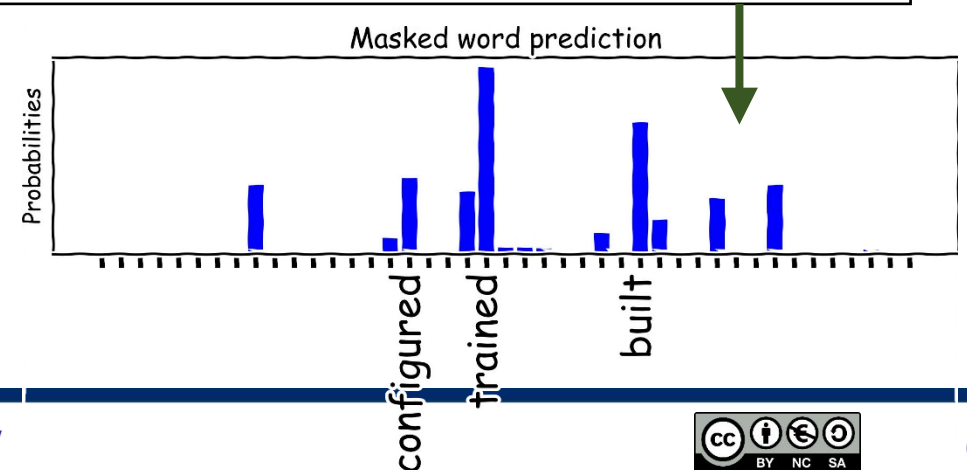
# Auto-regressive language models (*decoder-style, e.g. GPT*)

User: What is a large language model?

Assistant: A large language model is a type of artificial intelligence (AI) model designed to understand and generate human-like language. These models are



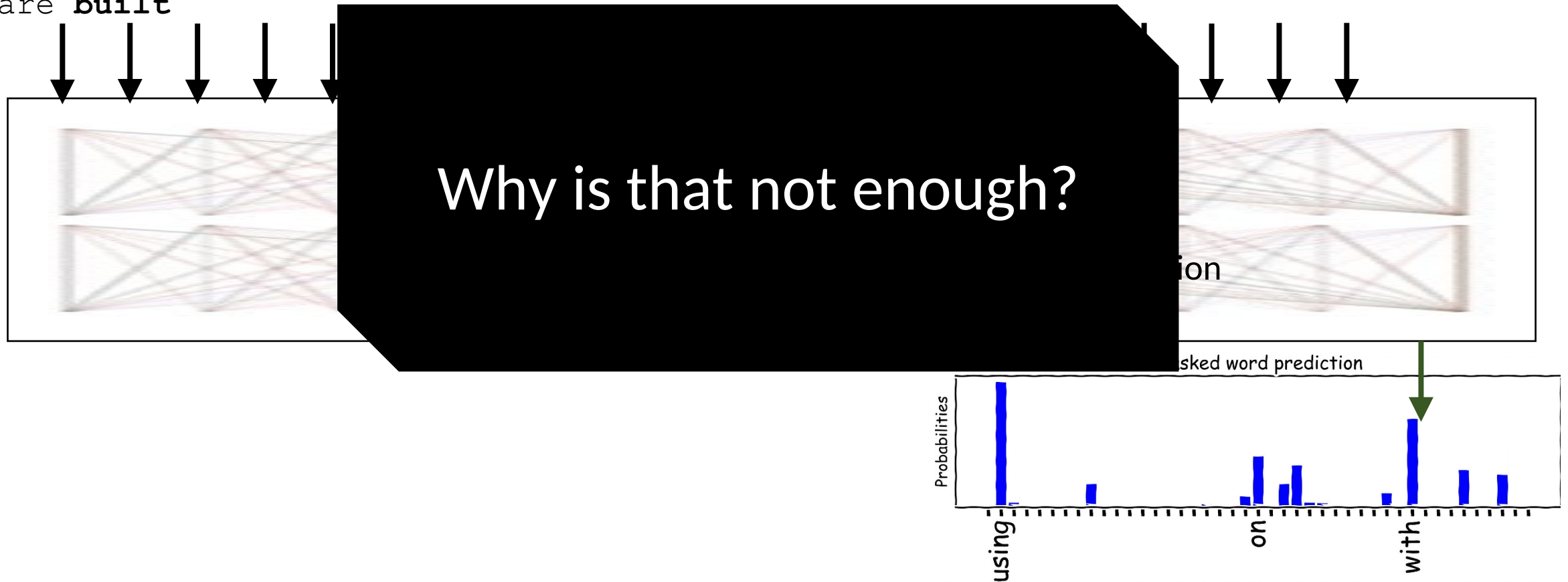
Next tokens are randomly selected from the predicted distribution, introducing variability in the generated output



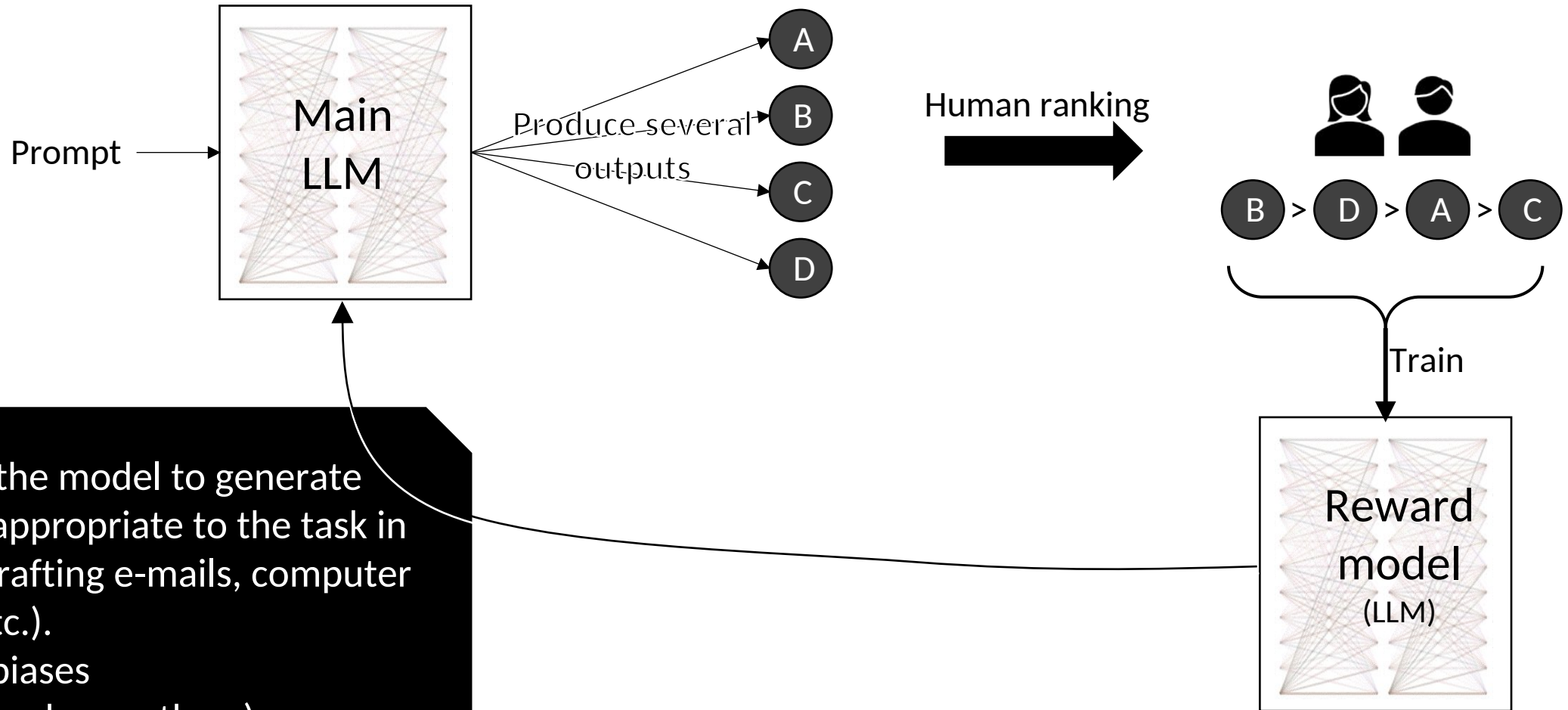
# Auto-regressive language models (*decoder-style, e.g. GPT*)

User: What is a large language model?

Assistant: A large language model is a type of artificial intelligence (AI) model designed to understand and generate human-like language. These models are **built**

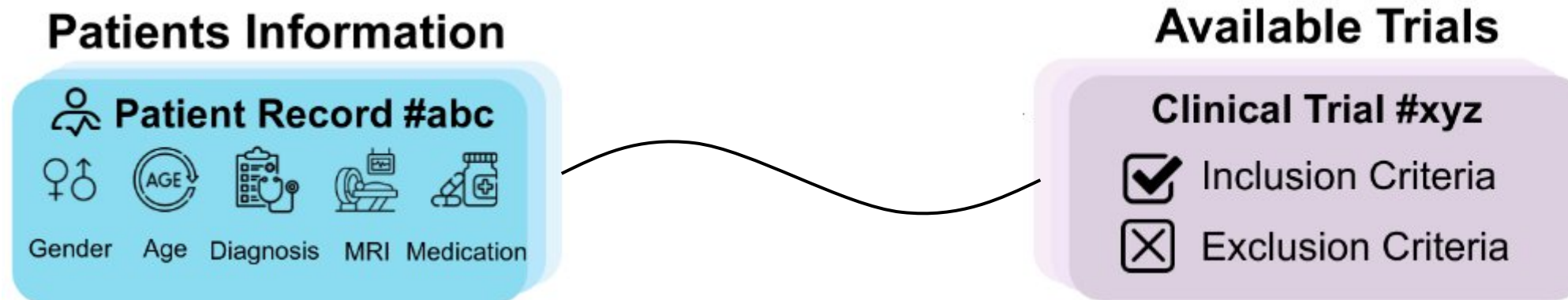


# Reinforcement Learning with Human Feedback



- Adapts the model to generate results appropriate to the task in hand (drafting e-mails, computer code, etc.).
- Avoids biases (and introduces others)

# LLMs for patient selection



*Large Language Models for Healthcare Data Augmentation: An Example on Patient-Trial Matching.*  
Jiayi Yuan, Ruixiang Tang, Xiaoqian Jiang, Xia Hu. AMIA Annual Symposium, March 2023

## Improving Patient Pre-screening for Clinical Trials: Assisting Physicians with Large Language Models

D. Hamer, P. Schoor, +1 author Daniel Kapitan · Published in arXiv.org 14 April 2023 · Computer Science, Medicine

## Scaling Clinical Trial Matching Using Large Language Models: A Case Study in Oncology

Cliff Wong, Sheng Zhang, +8 authors Hoifung Poon · Published in Machine Learning in Health... 4 August 2023 · Medicine, Computer Science

## Transforming clinical trials: the emerging roles of large language models

Jong-Lyul Ghim, Sangzin Ahn · Published in Translational and Clinical... 1 September 2023 · Medicine, Computer Science, Linguistics

## AutoCriteria: a generalizable clinical trial eligibility criteria extraction system powered by large language models

Surabhi Datta, Kyeryoung Lee, +9 authors Xiaoyan Wang · Published in J. Am. Medical Informatics... 11 November 2023 · Computer Science, Medicine · Journal of the American Medical Informatics Association : JAMIA

## Distilling Large Language Models for Matching Patients to Clinical Trials

Mauro Nieves, Aditya Basu, Yanshan Wang, Hrituraj Singh, less · Published in arXiv.org 15 December 2023 · Computer Science, Medicine

## LLM for Patient-Trial Matching: Privacy-Aware Data Augmentation Towards Better Performance and Generalizability

Yuan, Ruixiang Tang, Xiaoqian Jiang, Xia Hu, less · Published in arXiv.org 2023 · Computer Science, Medicine

## Matching Patients to Clinical Trials with Large Language Models

Qiao Jin, Zifeng Wang, +2 authors Zhiyong Lu · Published in arXiv.org 27 July 2023 · Computer Science, Medicine

## Zero-Shot Clinical Trial Patient Matching with LLMs

Michael Wornow, A. Lozano, +3 authors Nigam H. Shah · Published in arXiv.org 5 February 2024 · Computer Science, Medicine

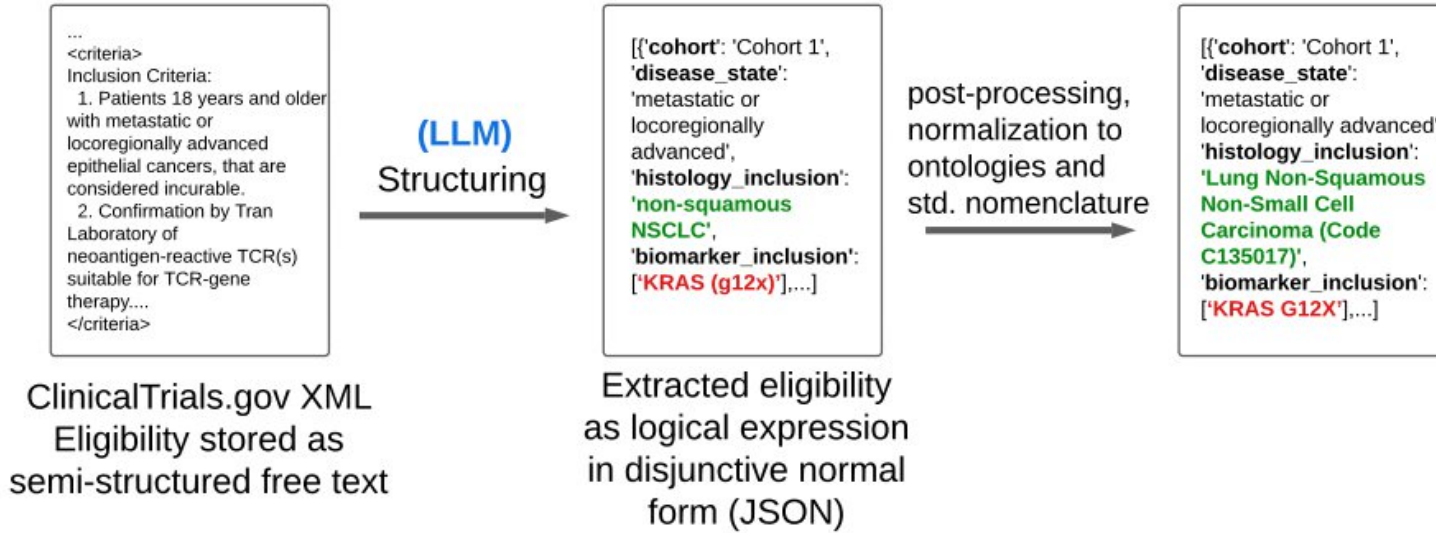
## Large Language Models for Healthcare Data Augmentation: An Example on Patient-Trial Matching.

Jiayi Yuan, Ruixiang Tang, +1 author Xia Hu · Published in AMIA ... Annual Symposium... 24 March 2023 · Computer Science, Medicine · AMIA ... Annual Symposium proceedings. AMIA Symposium

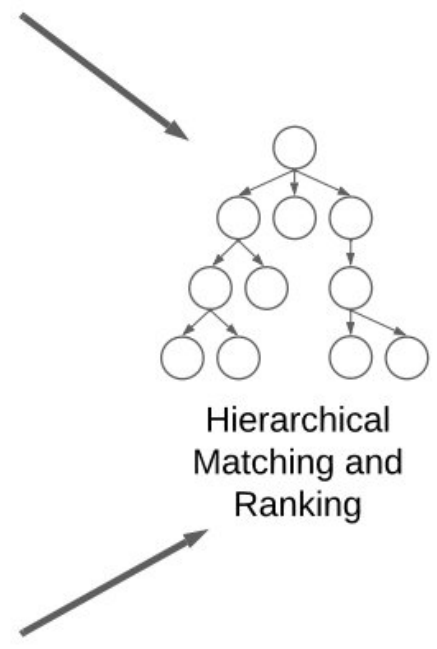
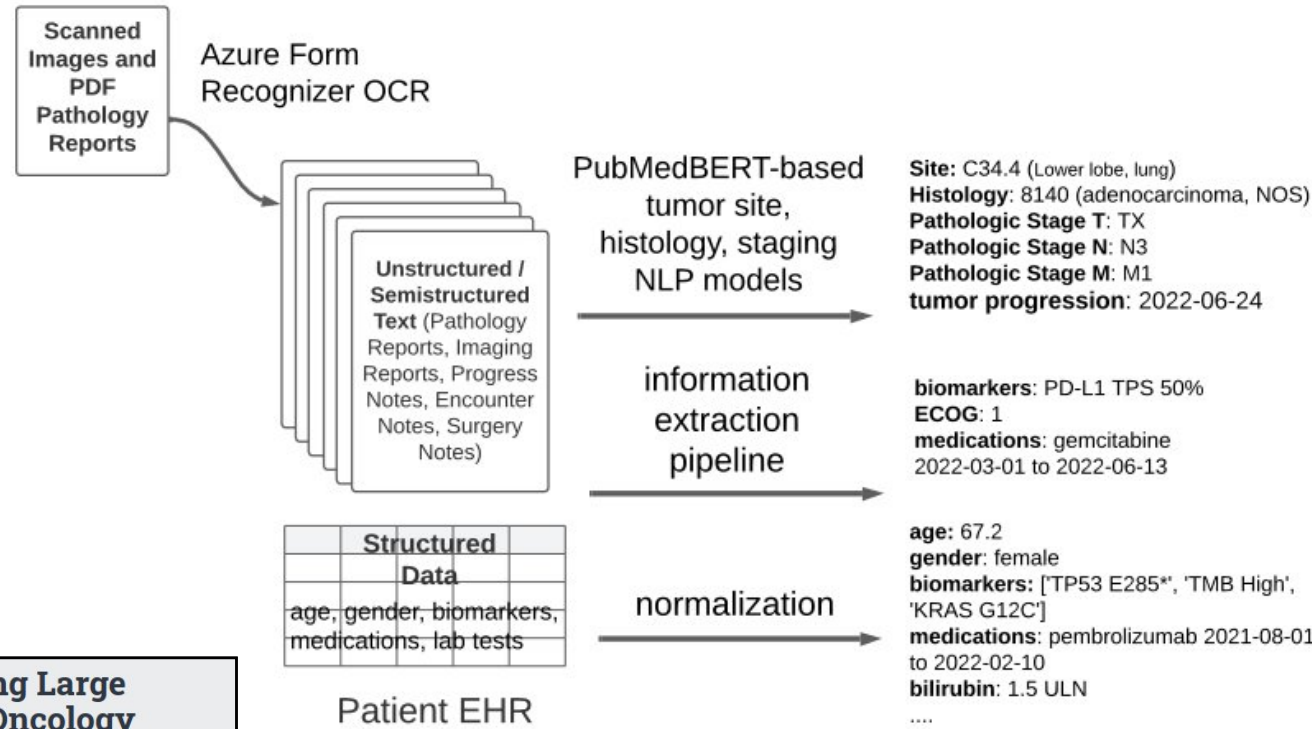
# LLMs for patient selection : different possible strategies

- Convert free-text eligibility criteria to formal and structured elements

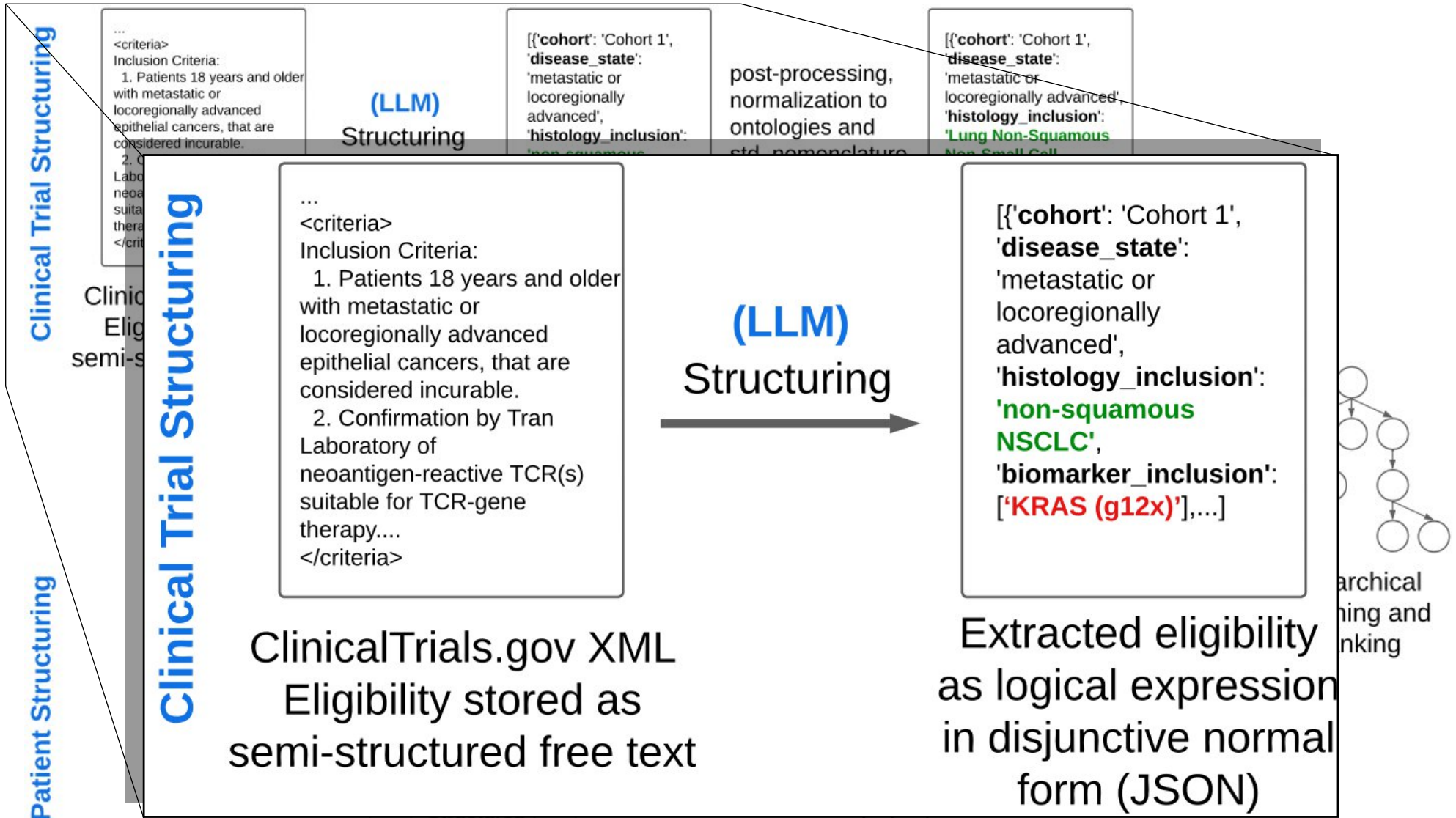
Clinical Trial Structuring



Patient Structuring







```

...
<criteria>
Inclusion Criteria:
1. Patients 18 years and older
with metastatic or
locoregionally advanced
epithelial cancers, that are
considered incurable.
2. Confirmation by Tran
Laboratory of
neoantigen-reactive TCR(s)
suitable for TCR-gene
therapy....
</criteria>

```

```

...
<criteria>
Inclusion Criteria:
1. Patients 18 years and older
with metastatic or
locoregionally advanced
epithelial cancers, that are
considered incurable.
2. Confirmation by Tran
Laboratory of
neoantigen-reactive TCR(s)
suitable for TCR-gene
therapy....
</criteria>

```

```

[{'cohort': 'Cohort 1',
'disease_state':
'metastatic or
locoregionally
advanced',
'histology_inclusion':
'Lung Squamous
Non-Small Cell'}]

```

(LLM) Structuring

```

[{'cohort': 'Cohort 1',
'disease_state':
'metastatic or
locoregionally advanced',
'histology_inclusion':
'Lung Non-Squamous
Non-Small Cell'}]

```

```

[{'cohort': 'Cohort 1',
'disease_state':
'metastatic or
locoregionally
advanced',
'histology_inclusion':
'non-squamous
NSCLC',
'biomarker_inclusion':
['KRAS (g12x)',...]}]

```

Data	
age	gender, biomarkers, medications, lab tests

normalization

```

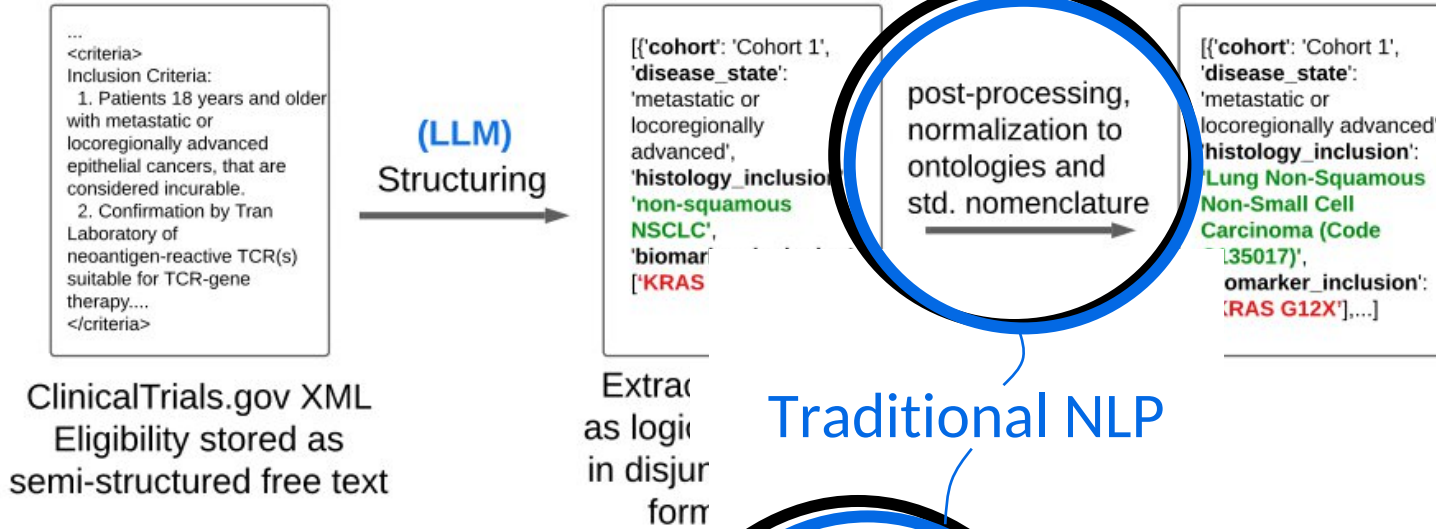
gender: female
biomarkers: ['TP53 E285*', 'TMB High', 'KRAS G12C']
medications: pembrolizumab 2021-08-01 to 2022-02-10
bilirubin: 1.5 ULN
....

```

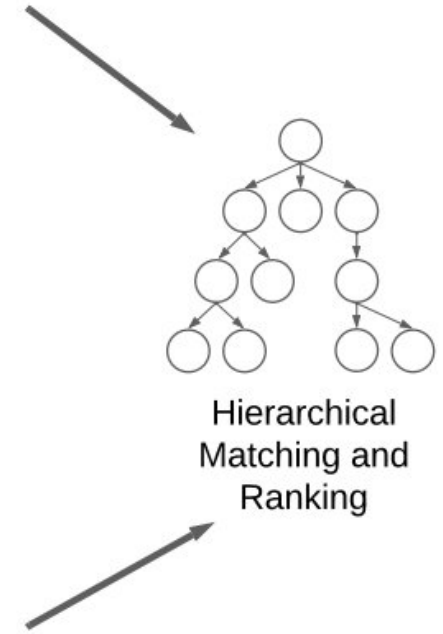
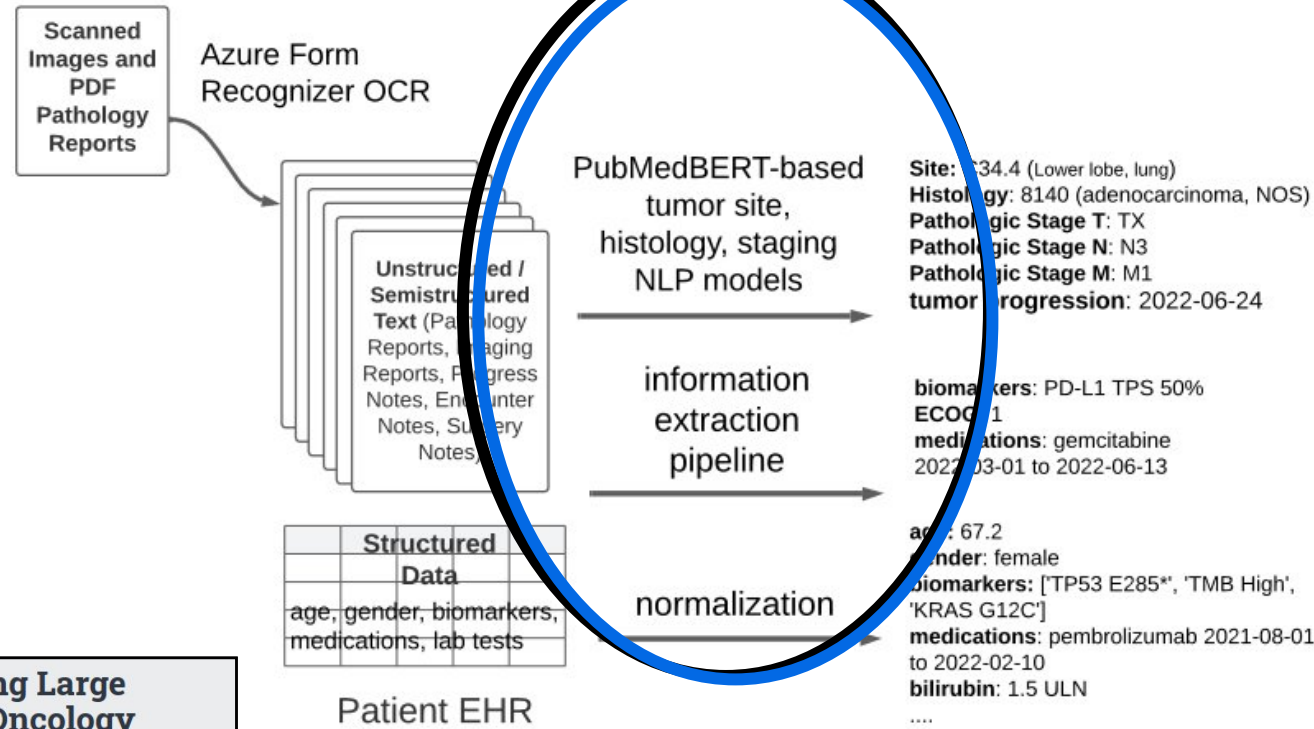
Patient EHR

Structured Patient Data

Clinical Trial Structuring



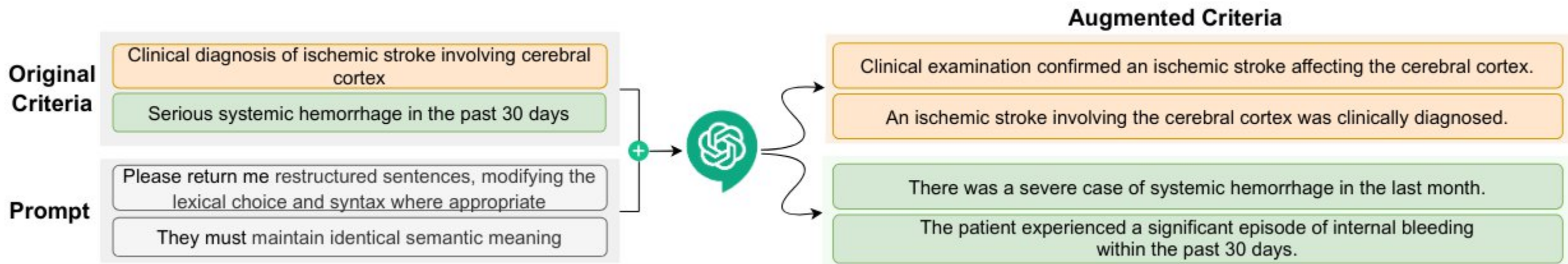
Patient Structuring



Traditional NLP

# LLMs for patient selection : different possible strategies

- Convert free-text eligibility criteria to formal and structured elements
- Convert eligibility criteria to EHR-like sentences



## Large Language Models for Healthcare Data Augmentation: An Example on Patient-Trial Matching.

Jiayi Yuan, Ruixiang Tang, +1 author Xia Hu · Published in AMIA ... Annual Symposium... 24 March 2023 · Computer Science, Medicine · AMIA ... Annual Symposium proceedings. AMIA Symposium

# LLMs for patient selection : different possible strategies

- Convert free-text eligibility criteria to formal and structured elements
- Convert eligibility criteria to EHR-like sentences
- Ask the LLM to compare patients' EHRs and eligibility criteria, all at once

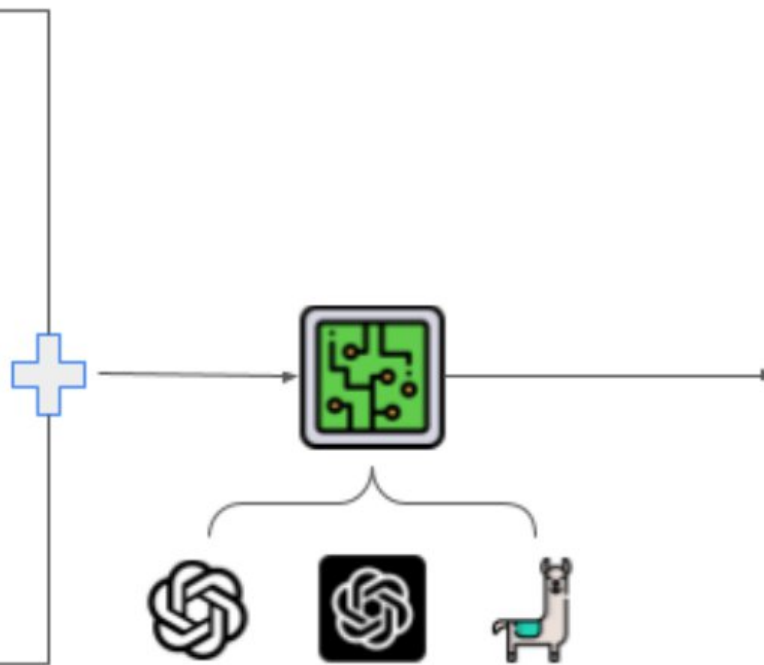
A 32-year-old woman comes to the hospital with vaginal spotting. Her last menstrual period was 10 weeks ago. Medical history is significant for appendectomy and several complicated UTIs. Vital signs are normal. Serum  $\beta$ -hCG level is 1800 mIU/mL, and a repeat level after 2 days shows an abnormal rise to ...

**Summary :** The purpose of this study is to compare the bleeding profile of ...

**Intervention :** Drug: norelgestromin/ethinyl...

**Inclusion Criteria :** Patients in good health as confirmed by medical history not pregnant as demonstrated by negative urine...

**Exclusion Criteria :** History or presence of disorders commonly accepted as contraindications to steroid hormonal therapy including...



```

{
  "inclusion_criteria": {
    "Patients in good health as confirmed by medical history": [
      "The patient has a medical history of appendectomy and several complicated UTIs, but her current vital signs are normal.",
      3,
      5,
    ],
    "included":
  ],
  "not pregnant as demonstrated by negative urine pregnancy test": [
    "The patient's serum  $\beta$ -hCG level is 1800 mIU/mL, and a repeat level after 2 days shows an abnormal rise to 2100 mIU/mL. This indicates that the patient is likely pregnant.",
    6,
  ],
  "not included":
  ],
  "exclusion_criteria": {
    "History or presence of ...": [
      "The patient's medical history is significant for appendectomy and several complicated UTIs, but there is no mention of any disorders that are contraindications to steroid hormonal therapy.",
      3,
    ],
    "no relevant information":
  ]
}

```

Ranking

## Distilling Large Language Models for Matching Patients to Clinical Trials

Mauro Nieves, Aditya Basu, Yanshan Wang, Hrituraj Singh | less • Published in arXiv.org 15 December 2023 • Computer Science, Medicine

# Conclusion (1/2)

“ Clinical trial matching is a key process in health delivery and discovery. In practice, it is plagued by overwhelming unstructured data and unscalable manual processing.

[...]

Initial findings are promising: out of box, cutting-edge LLMs, such as GPT-4, can already **structure elaborate eligibility criteria of clinical trials and extract complex matching logic (e.g., nested AND/OR/NOT). While still far from perfect, LLMs substantially outperform prior strong baselines and may serve as a preliminary solution to help triage patient-trial candidates with humans in the loop.** Our study also reveals a few significant growth areas for applying LLMs to end-to-end clinical trial matching, such as context limitation and accuracy, especially in **structuring patient information from longitudinal medical records.** ”

## Scaling Clinical Trial Matching Using Large Language Models: A Case Study in Oncology

Cliff Wong, Sheng Zhang, +8 authors Hoifung Poon • Published in *Machine Learning in Health...* 4 August 2023 •  
Medicine, Computer Science

Orphanet - 2024, September 17



# Conclusion (2/2)

“

**The integration of LLMs into clinical practice is not without its challenges, especially from legal and quality assurance perspectives.**

LLMs are **susceptible to generating misleading or incorrect information**, a phenomenon known as “hallucination”, and **ensuring quality control may prove to be demanding**. The complexity and flexibility of LLMs correspondingly make validation regarding accuracy, safety, and clinical efficacy particularly challenging. The inherent opacity of AI models further exacerbates the difficulty of their application in critical, real-world scenarios. Techniques such as chain-of-thought prompting may guide the language model to reveal the reasoning process behind its outputs, thereby increasing the models' explainability.

”

**Transforming clinical trials: the emerging roles of large language models**

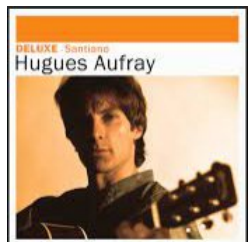
Jong-Lyul Ghim, Sangzin Ahn • Published in *Translational and Clinical...* 1 September 2023 • Medicine, Computer Science, Linguistics

Orphanet - 2024, September 17





# Words



D' y penser j' avais le cœur gros  
En doublant les feux de Saint-Malo

Lebensversicherungsgesellschaftsangestellter









(employee of a life insurance company)

(Agglutinative language)

學而不思則罔，思而不學則殆

東京、マドリード、イスタンブール  
(トルコ)が争う2020年夏季五輪  
の開催地は7日(日本時間8日)、ブ  
エノスアイレスでの国際オリンピック  
委員会(IOC)総会で、IOC委員  
約100人の投票で決まる。

# "Bag-of-words" representation

Human-readable	
Explicit semantics	
Dimensionality	
Semantic similarity	
Polysemy	
Spelling variants	
Multi-word expressions	
Unknown words	

Language = symbols

Vocabulary = all the words

bike ≠ bicycle ≠ mountain bike,  
dog ≠ cat, Paris ≠ London



Tuebingen ≠ Tübingen,  
tomorrow ≠ tommorrow

# Dense representation (*embeddings*)



Two words close in vector space = two words that often share similar contexts



~~Two words close in vector space = two words with a close meaning~~

$w_{bike} \sim w_{bicycle}$

$w_{cat} \sim w_{tiger}$









$w_{bad} \sim w_{good}$

$w_{person} \sim w_{persons}$









$w_{Paris} \sim w_{London}$

$w_{small} \sim w_{tall}$

# Dense representation

Human-readable		
Explicit semantics		
Dimensionality		Vector = dozens or hundreds of dimensions
Semantic similarity		By construction, but implicit
Polysemy		
Spelling variants		
Multi-word expression		Arithmetic operations on word vectors
Unknown words		Depending on the method, see later

# Contextual representation

Human-readable	
Explicit semantics	
Dimensionality	
Semantic similarity	
Polysemy	
Spelling variants	
Multi-word expression	
Unknown words	

Vector = dozens or hundreds of dimensions

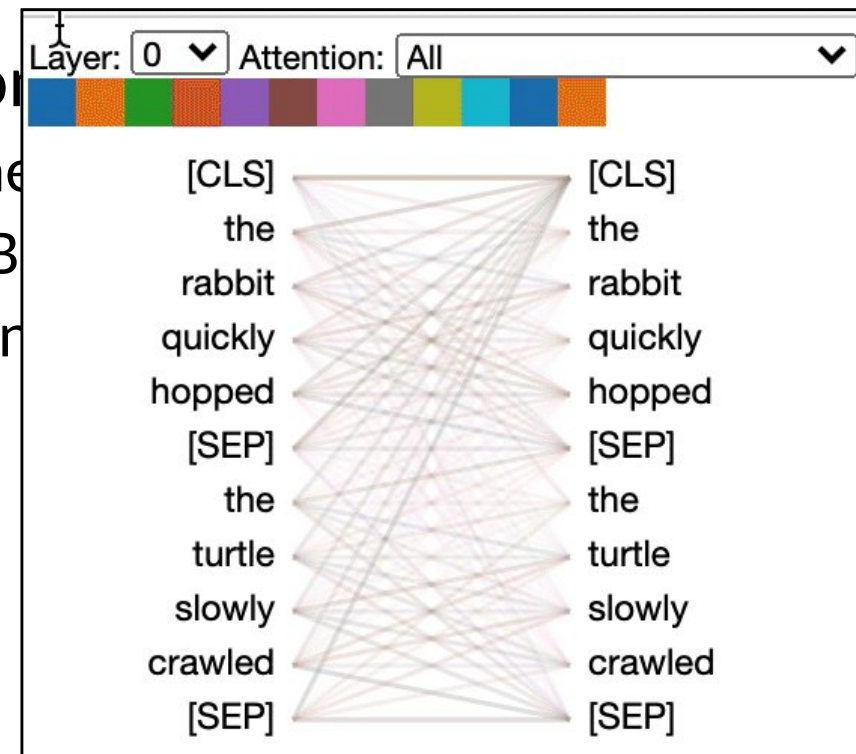
By construction, but implicit



# Contextual dense representation

- **Static representation:** one token = one vector
  - We handle an « embedding matrix » ( $N \times d$ )
  - The token vector is the same for each of its occurrences in the corpus
- **Contextual representation:** vector calculation in context
  - The calculation of the representation is integrated in the model
  - e.g: ELMo, ULMFit, transformer-based models such as BERT
  - The preceding (and following) words affect the representation (usually through an attention mechanism...)

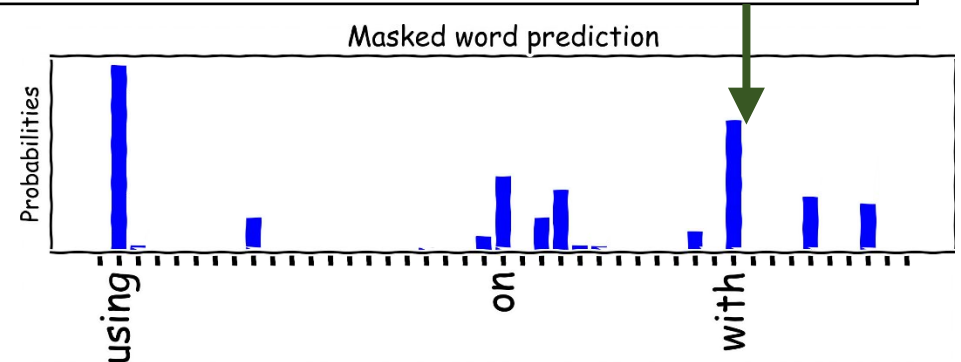
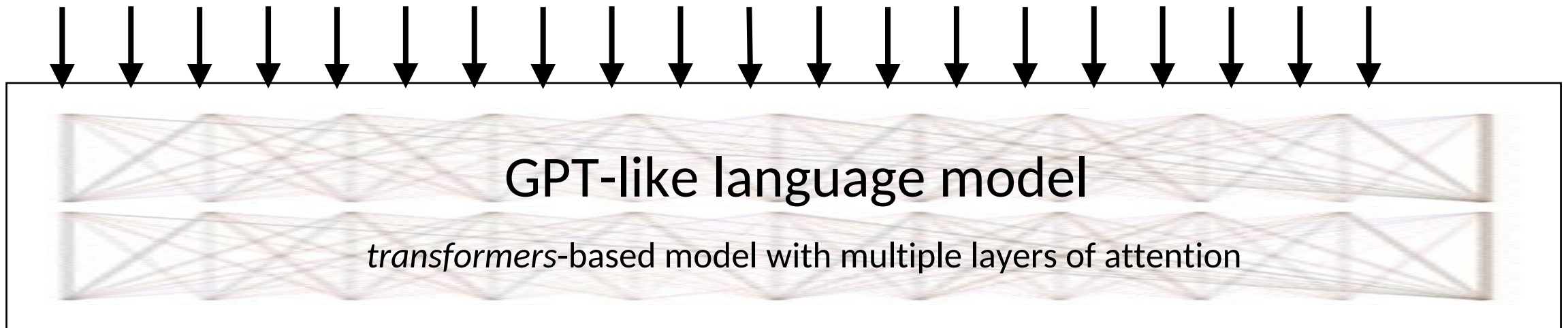
<https://github.com/jessevig/bertviz/>



# Auto-regressive language models (*decoder-style, e.g. GPT*)

User: What is a large language model?

Assistant: A large language model is a type of artificial intelligence (AI) model designed to understand and generate human-like language. These models are **built**



# By the way, not all of this is true. LLMs are not predicting words

c ' est un fame ##ux trois - mats , fin com ##me un o ##ise ##au  
( hiss ##ez ha ##ut , sant ##iano )  
di ##x - hui ##t n ##œ ##ud ##s , qu ##at ##re cents ton ##nea ##ux  
je sui ##s fi ##er d ' y et ##re mate ##lot

tie ##ns bon la vague , et tie ##ns bon le vent  
( hiss ##ez ha ##ut , sant ##iano )  
si die ##u ve ##ut , to ##uj ##ours dr ##oit dev ##ant  
( no ##us iron ##s ju ##s ##qu ' a san francisco )

je par ##s pour de long ##s moi ##s en lai ##ssa ##nt margot  
( hiss ##ez ha ##ut , sant ##iano )  
d ' y pens ##er , j ' ava ##is le c ##œ ##ur gr ##os  
( en do ##ub ##lan ##t les fe ##ux de saint mal ##o )

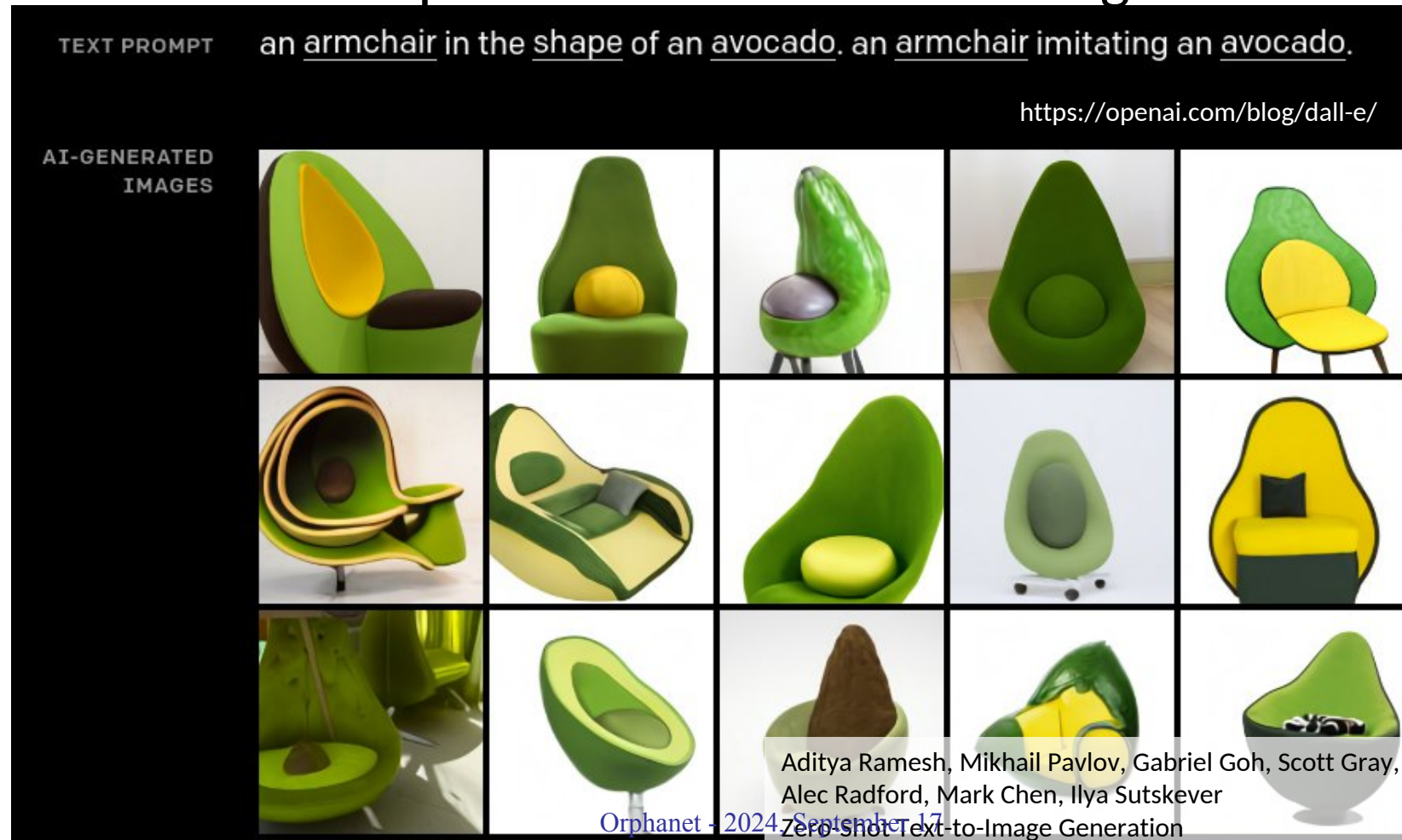
## WordPieces

- A vocabulary of predefined size, composed of ngrams of characters
- Vocabulary chosen to maximize the frequency of ngrams
- Possibility of a multilingual tokenizer



# Multimodal dense representation

Another advantage of dense representations:  
they can be mixed with representations of something else



# Multimodal dense representation

Another advantage of dense representations:  
they can be mixed with representations of something else



**DALL·E 2**

<https://openai.com/dall-e-2/>

# Multimodal dense representation

Another advantage of dense representations:  
they can be mixed with representations of something else

FOOD101

**guacamole** (90.1%) Ranked 1 out of 101 labels



- a photo of **guacamole**, a type of food.
- a photo of **ceviche**, a type of food.
- a photo of **edamame**, a type of food.
- a photo of **tuna tartare**, a type of food.
- a photo of **hummus**, a type of food.

SUN397

**television studio** (90.2%) Ranked 1 out of 397



- a photo of a **television studio**.
- a photo of a **podium indoor**.
- a photo of a **conference room**.
- a photo of a **lecture room**.
- a photo of a **control room**.

YOUTUBE-BB

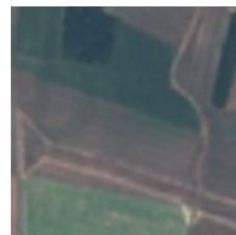
**airplane, person** (89.0%) Ranked 1 out of 23



- a photo of a **airplane**.
- a photo of a **bird**.
- a photo of a **bear**.
- a photo of a **giraffe**.
- a photo of a **car**.

EUROSAT

**annual crop land** (12.9%) Ranked 4 out of 10



- a centered satellite photo of **permanent crop land**.
- a centered satellite photo of **pasture land**.
- a centered satellite photo of **highway or road**.
- a centered satellite photo of **annual crop land**.
- a centered satellite photo of **brushland or shrubland**.

Alec Radford et al.

Learning Transferable Visual Models From Natural Language Supervision

Feb. 2021

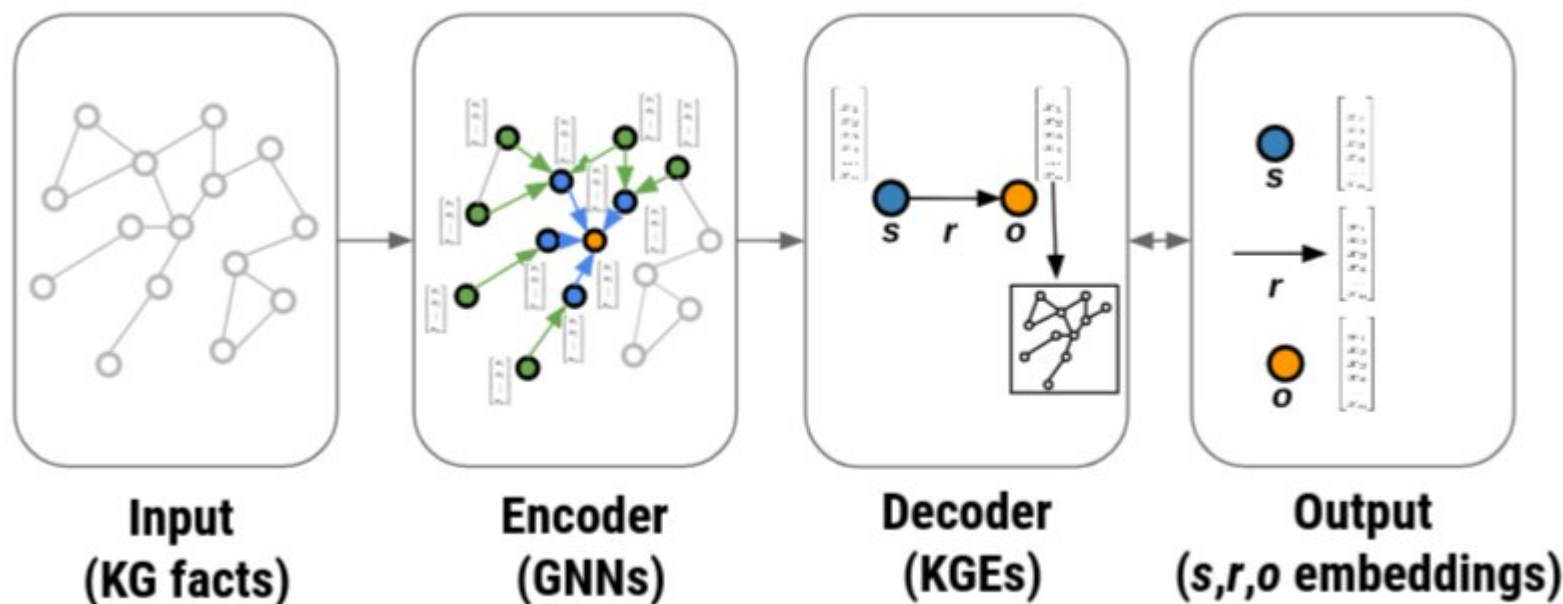
<https://openai.com/blog/clip/>

# Multimodal dense representation

Another advantage of dense representations:  
they can be mixed with representations of something else

Knowledge graph embeddings

Luca Costabello, Sumit Pai, Nicholas McCarthy, Adrianna Janik  
Knowledge Graph Embeddings Tutorial: From Theory to Practice  
ECAI 2020



(Source:  
Giuseppe Futia)