



GNOMON[®]
DIGITAL

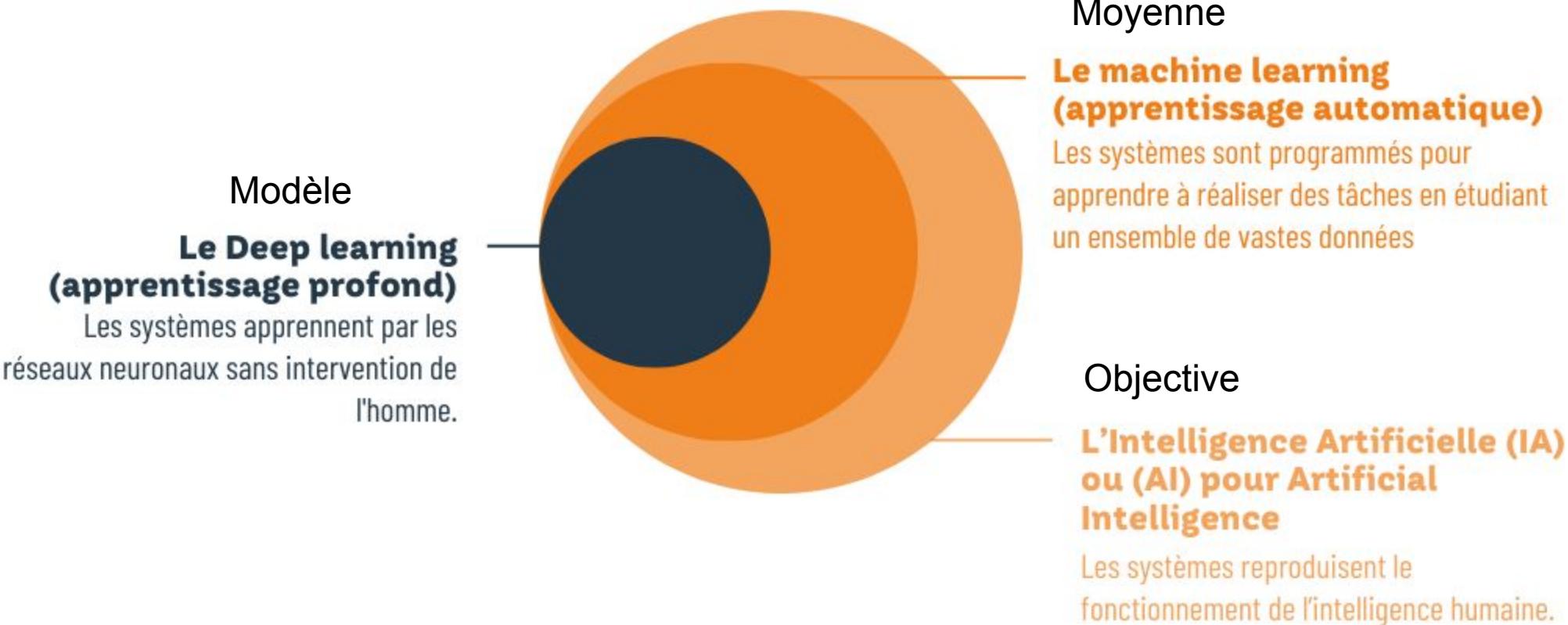
LLM

Jiqiong QIU

Objectif du cours



- Le but de ce cours est l'introduction aux modèles LLM (large language model)
- Le cours se focalise essentiellement sur les aspects suivants :
 - *Introduction à LLM*
 - *Prompt engineering*
 - *Semantic search*
 - *Longchain/ LlamaIndex*
 - *RAG*
 - *Function calling*



AI Générative



- Plus besoin de connaissance ML
- Pas besoin d'apprentissage
- Pas besoin des données
- Just besoin de bien formuler ton demande (prompt)



LLM



GNOMON[®]
DIGITAL

AI Builders

LE LEADER DE LA STRATÉGIE DATA ET IA

LLM Providers Landscape

LLM Commerciaux



LLM Open Source



LLM Privés



Generative language models

Input: Text

Output: Text

Generative image models

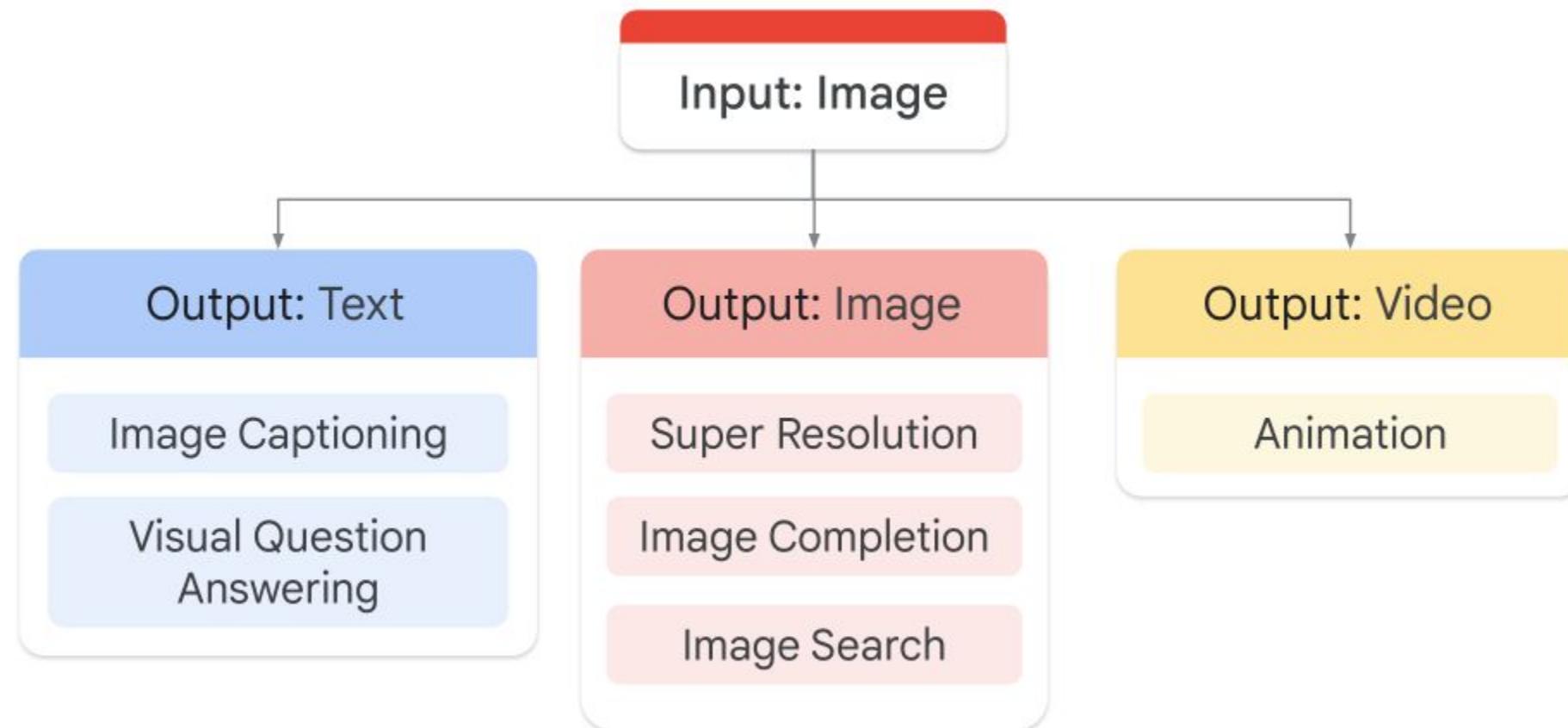
Input: Text

Output: Image

Multimodal:

Input: Text, audio, images, vidéo

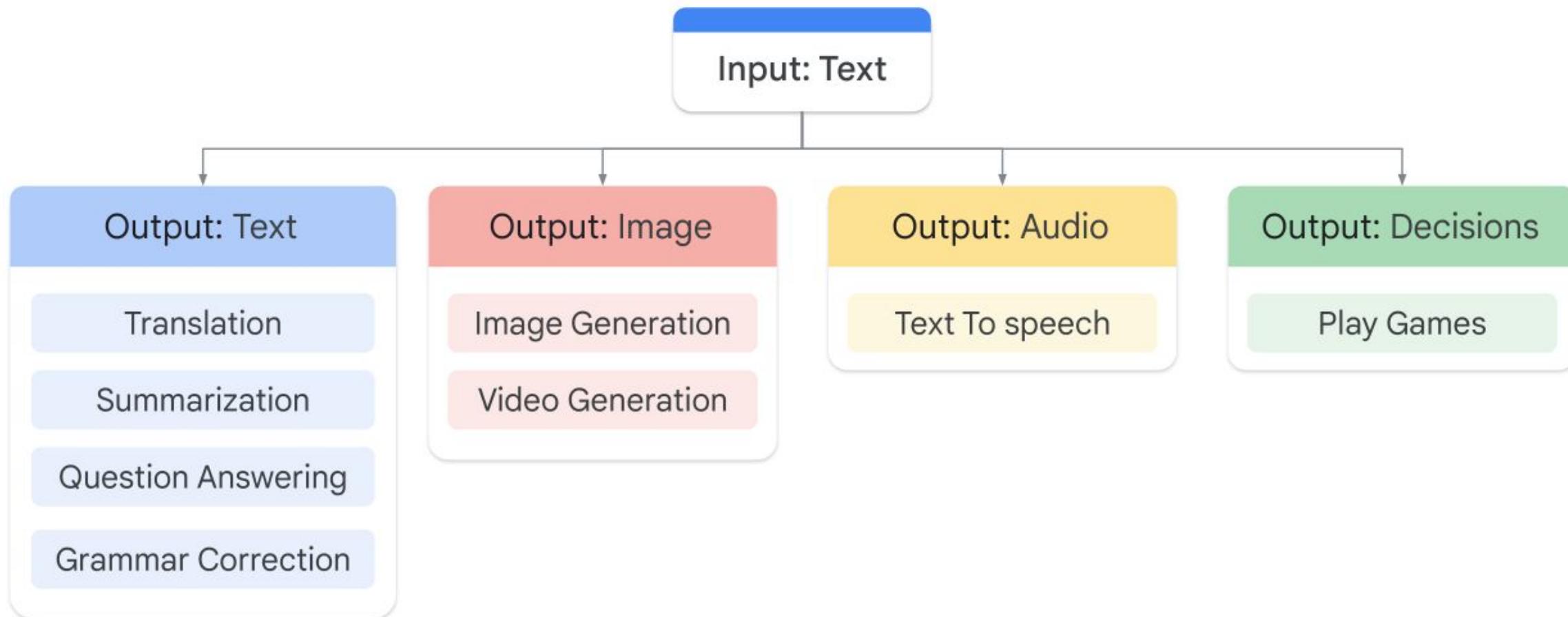
Output: Text, audio, images, vidéo



LLM

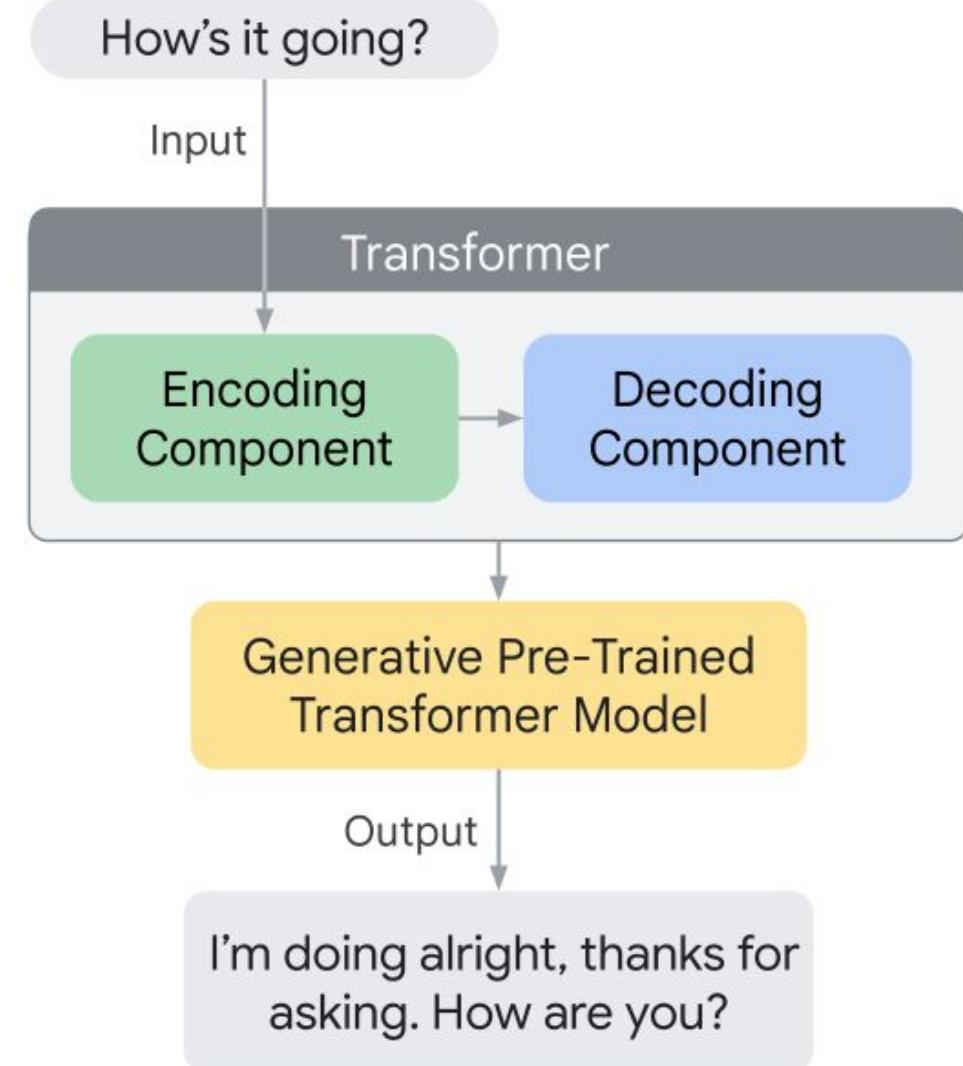


GNOMON[®]
DIGITAL



Comment ça marche?

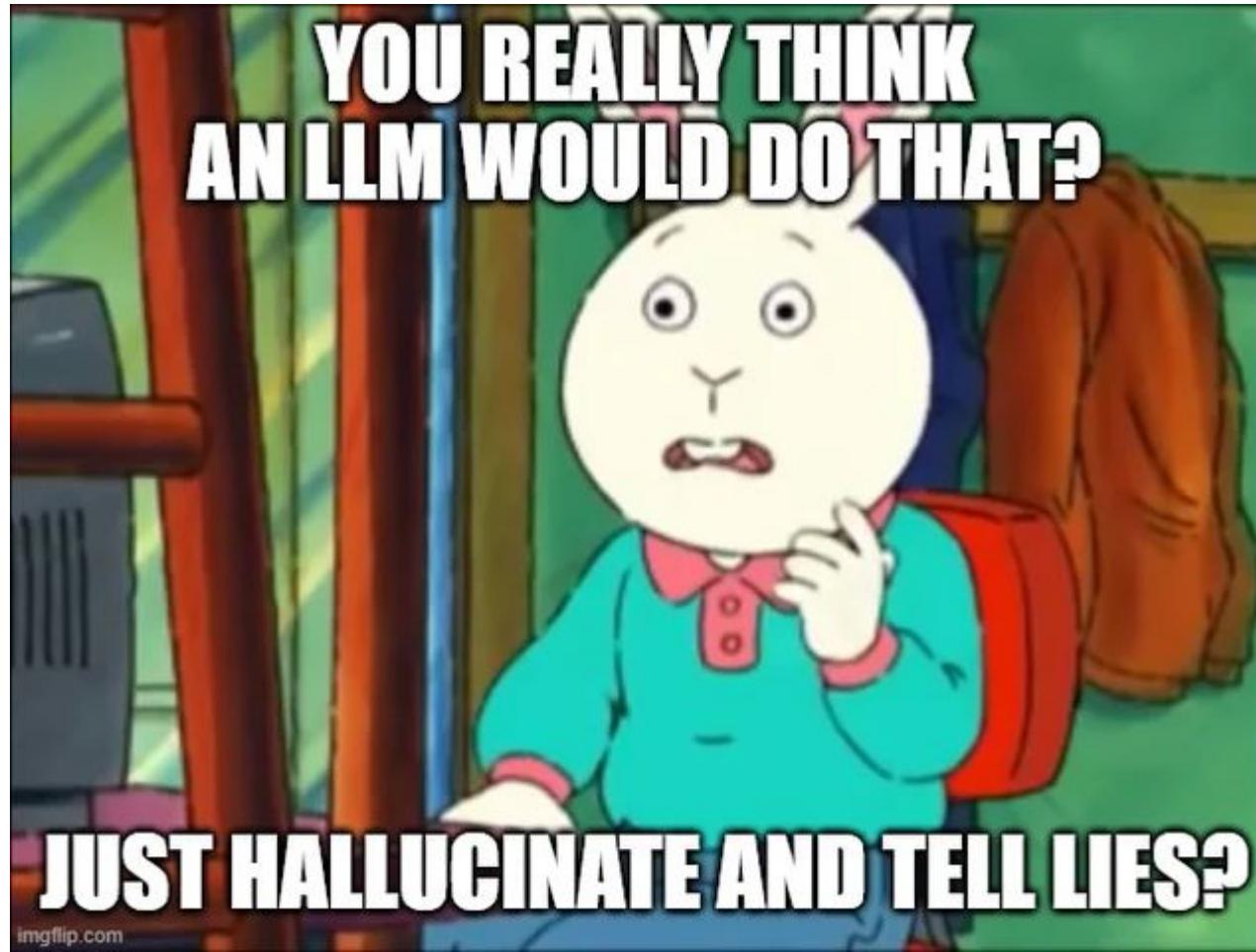
- Grand volume des données
- Billions parameters
- Apprentissage non-supervisé



LLM



GNOMON[®]
DIGITAL



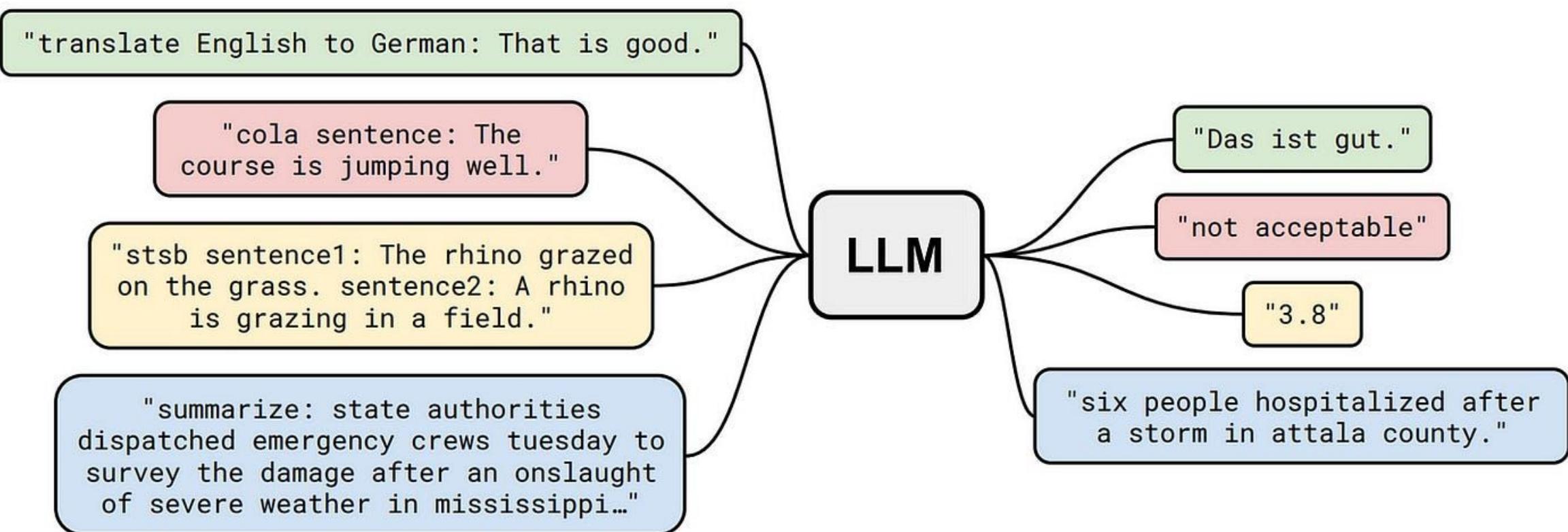


- Hallucinations:
 - Le modèle n'a pas assez des données pour faire l'entraînement
 - Le modèle est entraîné sur des données sales avec des bruits
 - On manque des context
 - Pas assez de contraintes sur le modèle

Prompt Engineering



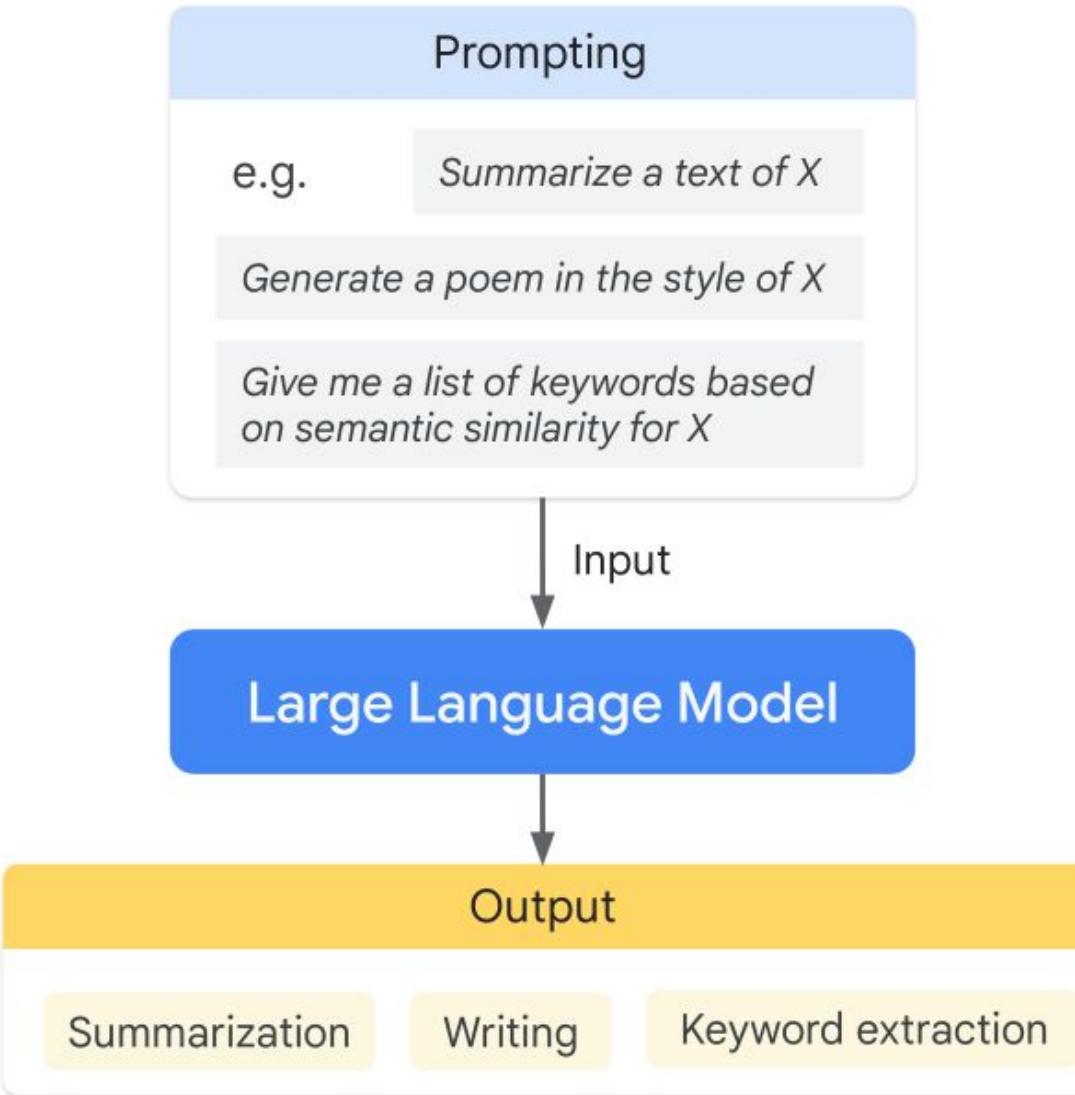
- On donne le contexte au modèle



Prompt Engineering



- Prompt engineering: la qualité d'entrée détermine la qualité de sortie



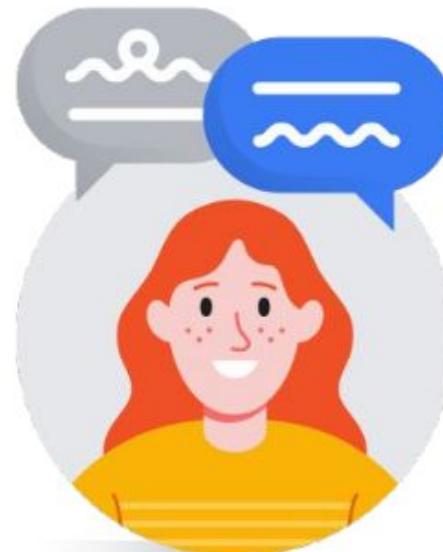
LLM



- Comment détermine un LLM



Forming a Database



Inputting a Prompt

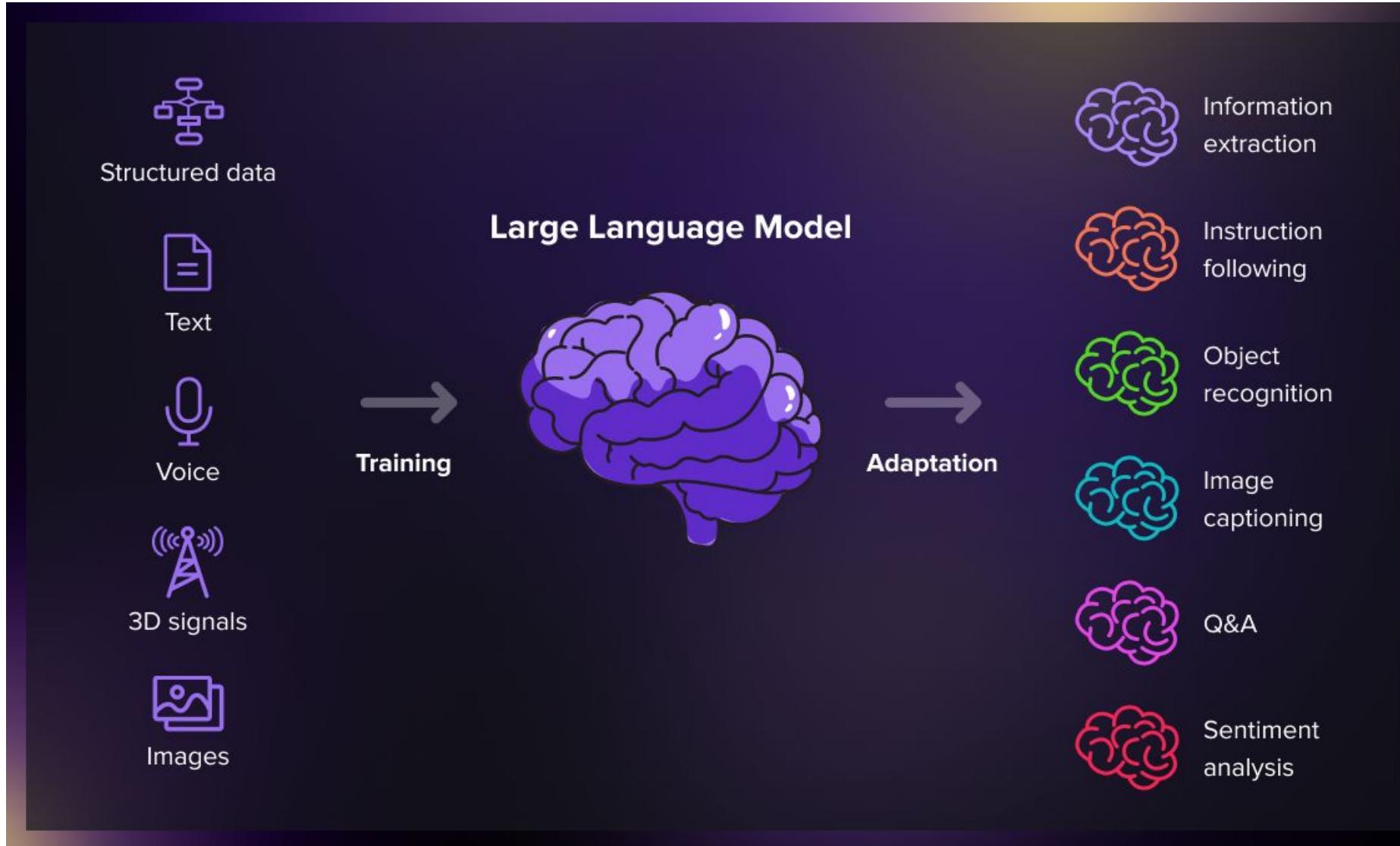


Generating content

LLM



GNOMON[®]
DIGITAL

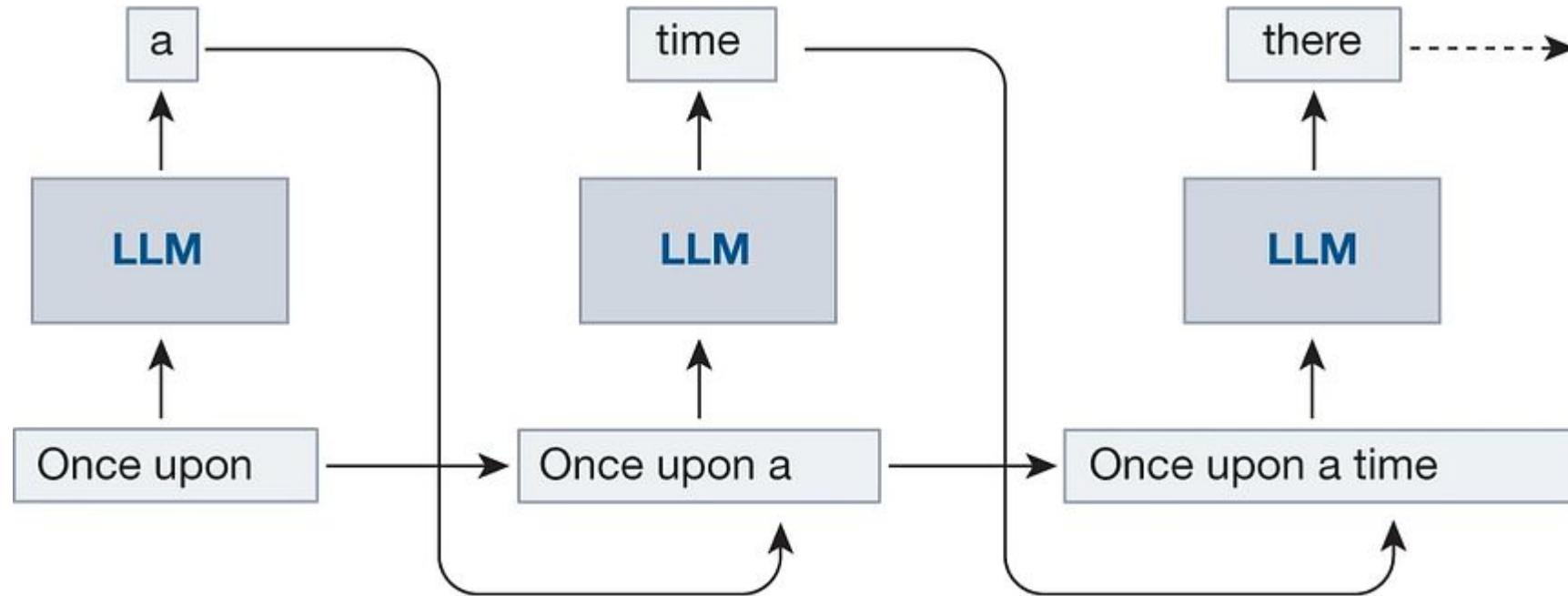


LLM: type de modèles



- Comment détermine un LLM

Text-to-text



LLM: type de modèles

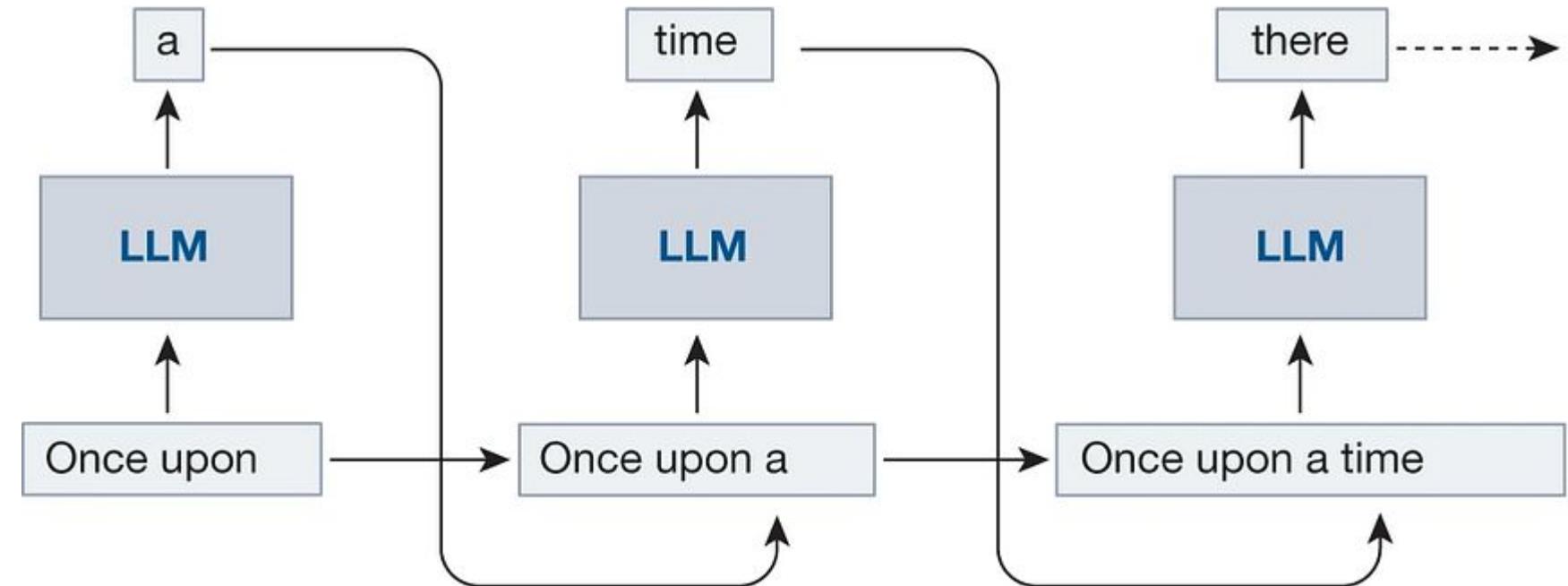


- Comment détermine un LLM

Text-to-text

Domain d'application:

- Génération des text
- Classification
- Résumé
- Traduction
- Reformuler
- Extraction
- Clustering



LLM: type de modèles



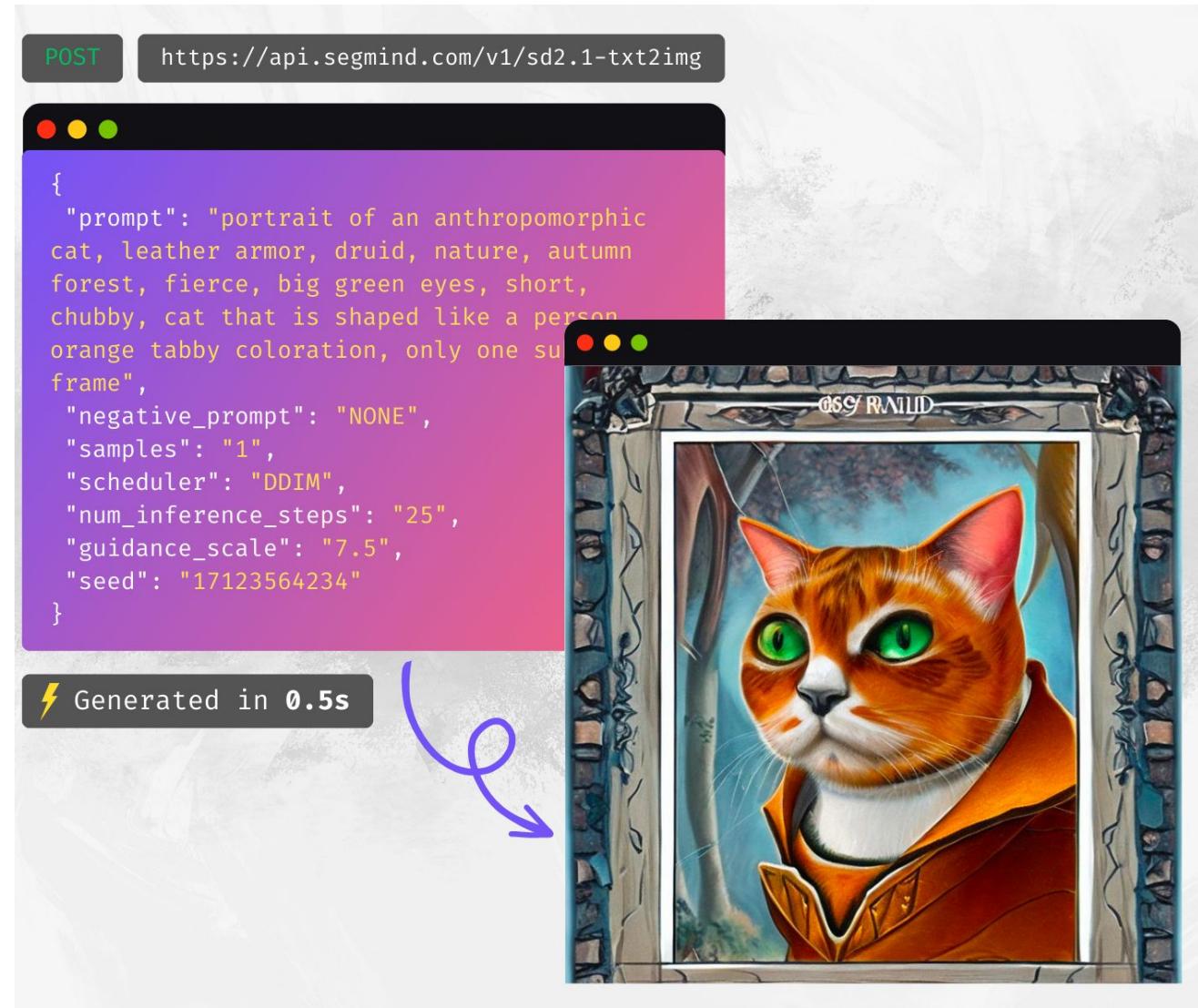
- Comment détermine un LLM

Text-to-image

Domain d'application:

- Image generation
- Image edition

<https://www.midjourney.com/>



LLM: type de modèles



- Comment détermine un LLM

Text-to-3D
(video)

Domain d'application:

- Video generation
- Video edition
- Jeux Vidéo



LLM: type de modèles



- Comment détermine un LLM

Text-to-Task (Tool GPT)

Domain d'application:

- Assistants virtuelle
- Automation
- LLM agent



[Ray-BAN META](#)

LLM: type de modèles

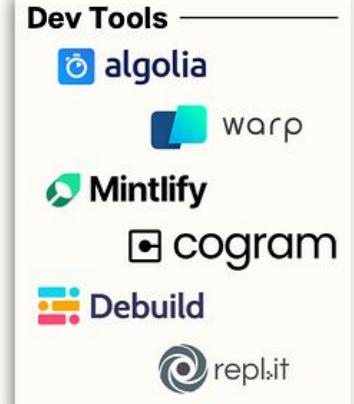


GNOMON[®]
DIGITAL

Large Language Models

BCV

Application Layer



Infrastructure Layer



Qu'est-ce que c'est Prompts et Prompt Engineering



Le **prompt engineering** dans le contexte des modèles de langage (LLM, pour *Large Language Models*) désigne l'art et la science de concevoir des instructions textuelles (ou *prompts*) afin d'obtenir des résultats spécifiques ou optimaux d'un modèle de langage.

Définition

Le prompt engineering consiste à rédiger, structurer ou formuler des *prompts* pour orienter le comportement d'un LLM et maximiser sa capacité à produire des réponses pertinentes, précises ou créatives. Cela peut inclure l'utilisation de formats particuliers, de questions claires, de contextes détaillés ou de contraintes explicites.

Objectifs principaux

- **Optimisation des réponses** : Obtenir des résultats qui correspondent aux attentes en termes de contenu, style ou précision.
- **Réduction des ambiguïtés** : Limiter les interprétations multiples pour guider le modèle vers une réponse unique et claire.
- **Exploration créative** : Stimuler des réponses originales en utilisant des formulations non conventionnelles.

Techniques courantes

1. **Précision des consignes** : Formuler des demandes explicites avec des détails clairs (ex. "Expliquez en trois points la théorie de la relativité").
2. **Contexte riche** : Fournir un contexte pertinent pour améliorer la compréhension (ex. "Imaginez que vous êtes un scientifique du 19e siècle...").
3. **Format structuré** : Exiger des réponses sous forme de liste, tableau ou paragraphe (ex. "Donnez un résumé sous forme de liste à puces").
4. **Chain-of-thought prompting** : Encourager le raisonnement pas à pas pour résoudre des problèmes complexes.
5. **Exemples dans le prompt** : Inclure des exemples explicites pour clarifier les attentes (ex. "Voici un exemple de réponse idéale...").

Généré par ChatGPT

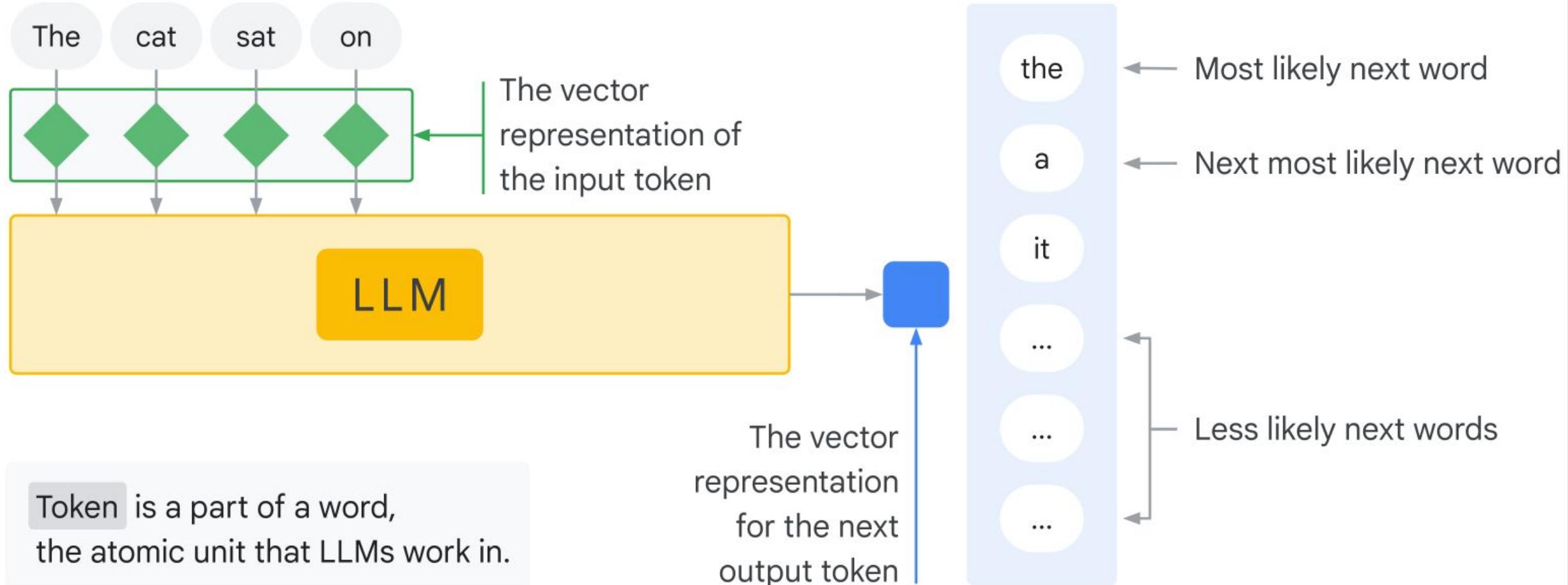
Qu'est-ce que c'est Prompts et Prompt Engineering



Il y a 3 façon d'entraîner un modèle LLM, chacun a besoin une façon de ‘prompting’ différente:

- Modèle LLM générique (ou brut)
- Optimisé par les instructions
- Optimisé par un dialogue

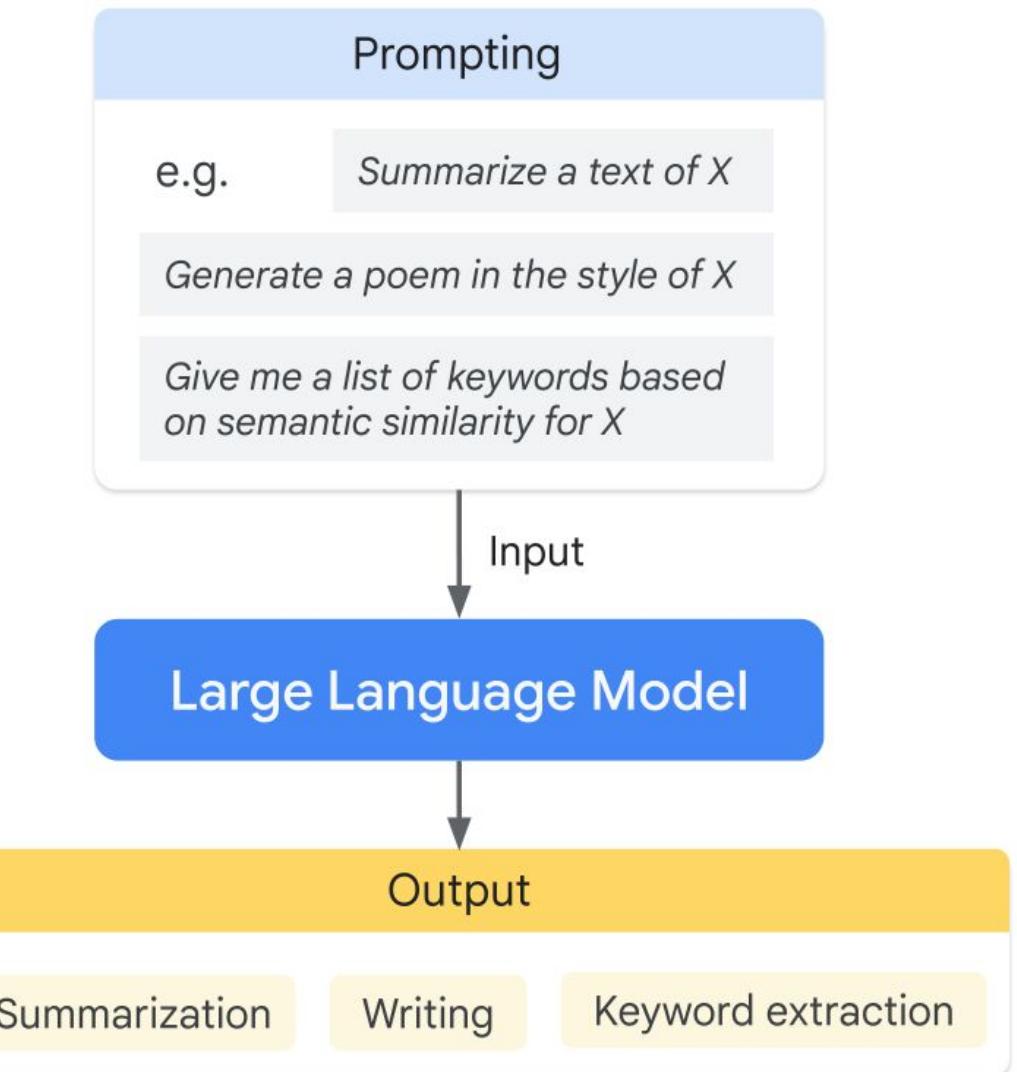
Modèle LLM générique (ou brut)



Optimisé par les instructions



Entraîner à prédire une réponse à partir des instructions données dans les textes entrées



Optimisé par une dialogue



Entraîné à engager un dialogue en prédisant la prochaine réponse.



Prompt examples

[User] Is the comment "do you like the weather?" ok or toxic?
[Bot] ok.
[User] can you briefly say why?
[Bot]

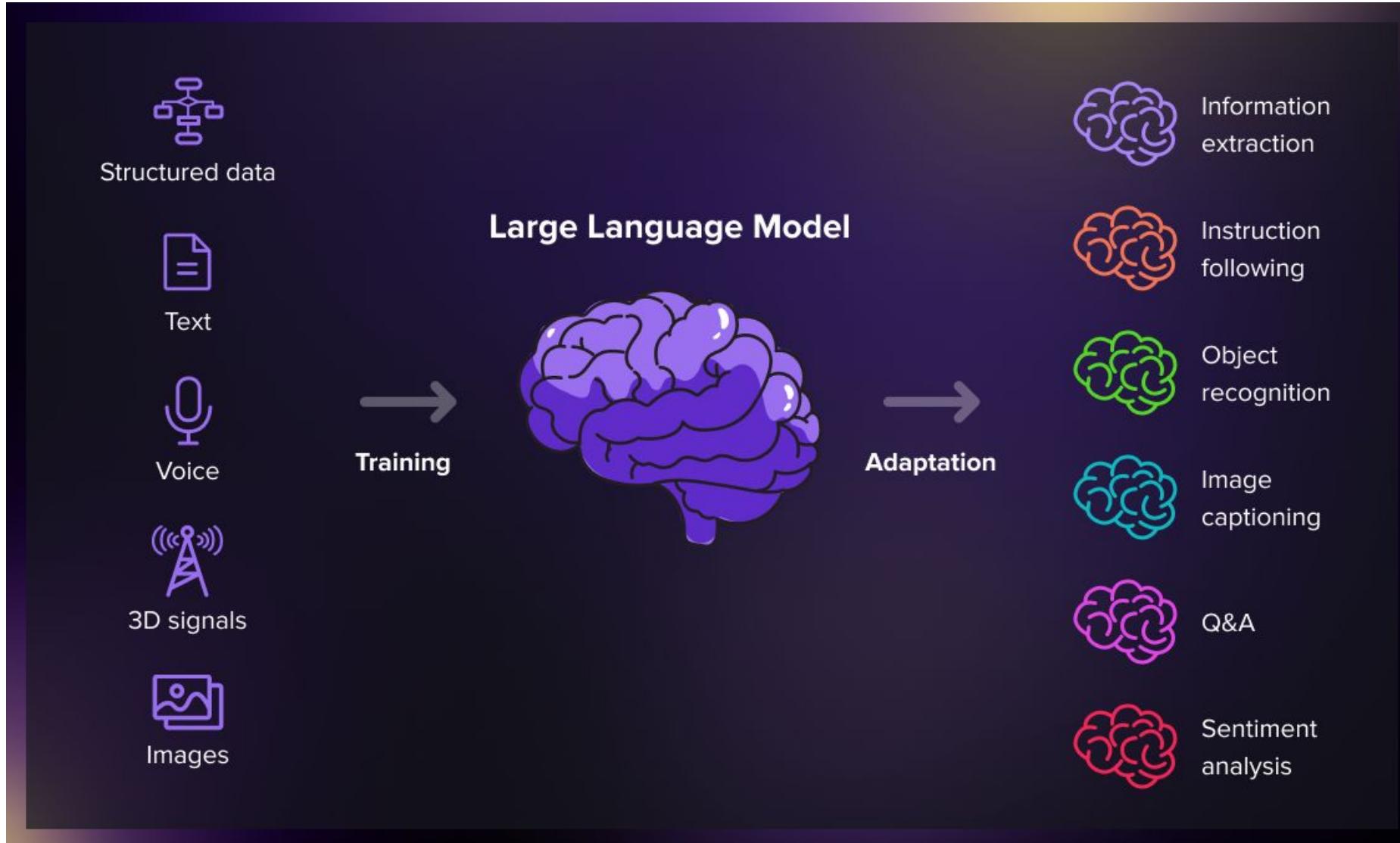
Model Output

It's just a question about the weather, people are not usually upset by that.

Les utilisations LLM



GNOMON[®]
DIGITAL



Fin-tuning in LLM



Le processus d'adaptation d'un modèle à un nouveau domaine ou à un ensemble de cas d'utilisation spécifiques en entraînant le modèle sur de nouvelles données. Par exemple, nous pouvons collecter des données d'entraînement et "ajuster" le modèle LLM spécifiquement pour le domaine juridique ou médical.



Fin-tuning in LLM



Le processus d'adaptation d'un modèle à un nouveau domaine ou à un ensemble de cas d'utilisation spécifiques en entraînant le modèle sur de nouvelles données. Par exemple, nous pouvons collecter des données d'entraînement et "ajuster" le modèle LLM spécifiquement pour le domaine juridique ou médical.



Prompt Engineering



GNOMON[®]
DIGITAL

The **response** from
the model



The **answers** you get depend on the
questions you ask.



The **prompt** you
designed

Prompt Engineering



On teste les cas suivants:

- Who is Linus Torvalds?
- Summarize some important dates in the life of Linus Torvalds
- Write me a tweet to celebrate Linus Torvald's birthday
- Tester toi-même un prompt

Prompt Engineering



GNOMON[®]
DIGITAL

Elements of a Good Prompt

Act as a helpful tutor who breaks down complex subjects into easy explanations.

I want you to explain the process of photosynthesis to a

14 year old student, to assist with biology exam preparations.

Your answer should be 300 words, written in a tone that's friendly and educational.

Persona: Ask the tool to take a role

Objective: What do you want the AI to do

Audience: Specify who it's for

Context: What does the tool need to know

Boundaries: Set your own direction & limitation

Tip 1 **Give Clear Instructions**

Use commands that instruct the AI tool on what you want to generate, such as 'explain', 'translate', 'summarize' or 'compare'.

Tip 2 **Provide Context**

Adding context and background information can help the tool to understand the task better. For example, mention the project type such as 'short story', 'report' or 'outline'

Tip 3 **Iterate & Experiment**

Try different instructions and techniques if you don't get the results you want. Prompting can be like an experiment that may require several rounds of iterations!

Prompt Engineering



Ambiguïté et Nuance:

- **Prompt A (Ambigu):** "Écris une histoire."
- **Prompt B (Précis):** "Écris une histoire courte de science-fiction sur un robot qui découvre qu'il est capable de ressentir des émotions, en utilisant un vocabulaire riche et des descriptions imagées."

Instructions Complexes

- **Prompt A (Simple):** "Traduis 'Hello, how are you?' en français."
- **Prompt B (Complexe):** "Traduis 'Hello, how are you?' en français, puis réécris la traduction en utilisant un langage formel, puis informel, et enfin en argot parisien. Indique pour chaque version le contexte approprié."

Contraintes Créatives

- **Prompt 5 (Libre):** "Écris un poème."
- **Prompt 6 (Constraint):** "Écrit un sonnet en alexandrins sur le thème de la solitude, en utilisant au moins trois fois le mot 'silence'."

Raisonnement et Connaissance

- **Prompt A (Factuel):** "Qui a peint la Joconde?"
- **Prompt B (Raisonnement):** "Si tous les chats sont des mammifères, et que Minou est un chat, est-ce que Minou est un mammifère ? Explique ton raisonnement."

Prompt Engineering



Sensibilité au contexte

- **Prompt 9 (Sans contexte):** "Explique le concept de la relativité."
- **Prompt 10 (Avec contexte):** "Imagine que tu expliques le concept de la relativité à un enfant de 8 ans. Utilise des analogies simples et évite le jargon scientifique."

Tips:

- By nice
- Step by step
- Ask LLM to write your prompt



GNOMON[®]
DIGITAL

POWER OF PROMPT ENGINEERING

ZERO-SHOT INFERENCE

A MODEL THAT CAN ANSWER QUESTIONS OR PERFORM TASKS WITHOUT ANY SPECIFIC TRAINING ON THOSE PARTICULAR PROMPTS.

ONE-SHOT INFERENCE

ONE-SHOT INFERENCE TAKES ONE EXAMPLE. THE MODEL CAN GRASP THE ESSENCE OF THE TASK AND GENERATE THE DESIRED OUTPUT.

FEW-SHOT INFERENCE

FEW-SHOT INFERENCE ALLOWS YOU TO PROVIDE A SMALL NUMBER OF EXAMPLES TO GUIDE THE MODEL'S BEHAVIOR. IT'S LIKE GIVING THE MODEL A MINI-TRAINING SESSION WITH JUST A HANDFUL OF PROMPTS.

<https://www.promptingguide.ai/>

Zero-shot prompting



On donne des instructions à LLM sans aucun exemples

EX:

Quelle est la sentiment de la phrase suivant: 'I am very happy today'? Positive, négative ou neutre.

<https://www.promptingguide.ai/>

One-shot prompting



On donne des instructions à LLM avec un exemple

EX:

Quelle est la sentiment de la phrase suivant: 'I am very happy today'? Positive, négative ou neutre.

Example:

Input: "The book is very interesting"

Output: "Positive"

<https://www.promptingguide.ai/>

Prompt Engineering



On teste?

- Gemini: https://aistudio.google.com/prompts/new_chat
- ChatGPT: <https://chatgpt.com/>

<https://www.promptingguide.ai/>

Few-shot prompting



On donne des instructions à LLM avec un exemple

EX:

Quelle est la sentiment de la phrase suivant: 'I am very happy today'? Positive, négative ou neutre.

Example:

Input: "The book is very interesting"

Output: "Positive"

Input: "I feel so depressed today"

Output: "Négative"

Input: "My name is John"

Output: "Neutre"

<https://www.promptingguide.ai/>

Few-shot prompting



https://github.com/GoogleCloudPlatform/asl-ml-immersion/blob/master/notebooks/vertex_genai/labs/prompt_engineering.ipynb

https://github.com/GoogleCloudPlatform/asl-ml-immersion/blob/master/notebooks/vertex_genai/labs/gemini_for_multimodal_promting.ipynb

Chain of Thought



GNOMON[®]
DIGITAL

Magique prompt: "Let's think this through step by step"

Standard Prompting

Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The answer is 27. X

Chain of Thought Prompting

Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. $5 + 6 = 11$. The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

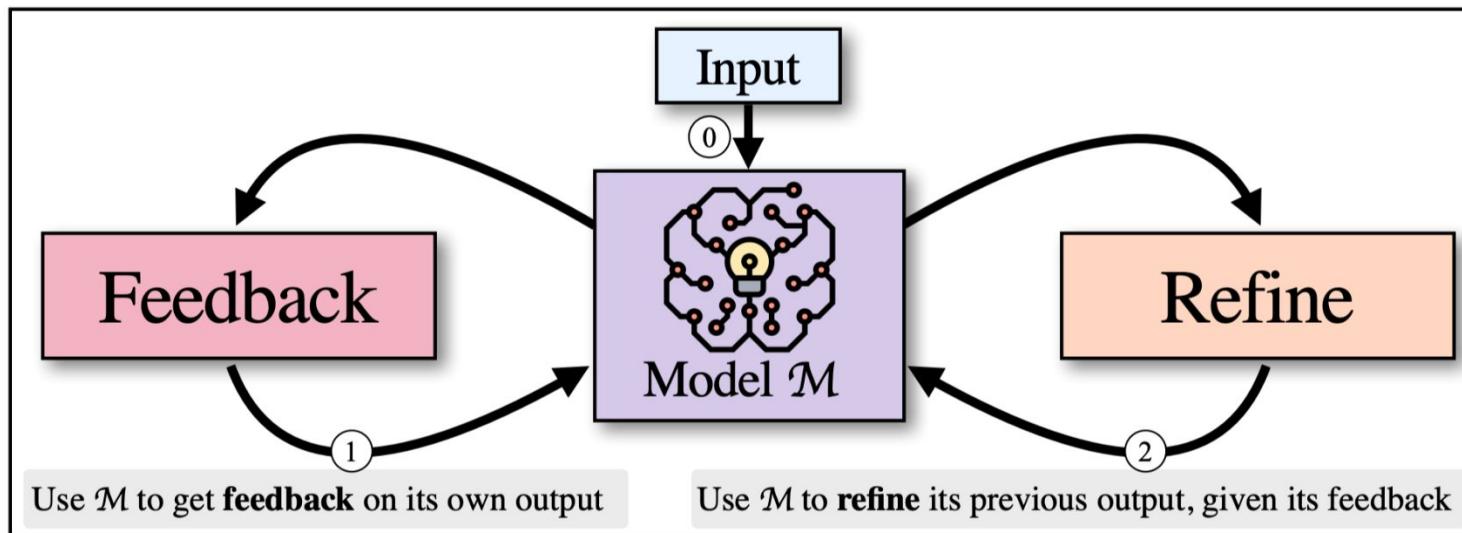
Model Output

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had $23 - 20 = 3$. They bought 6 more apples, so they have $3 + 6 = 9$. The answer is 9. ✓

Self Refine



Self-Refine



Learn Prompting

Aman Madaan¹, Niket Tandon², Prakhar Gupta¹, Skyler Hallinan³, Luyu Gao¹, Sarah Wiegreffe², Uri Alon¹, Nouha Dziri², Shrimai Prabhumoye⁴, Yiming Yang¹, Shashank Gupta², Bodhisattwa Prasad Majumder⁵, Katherine Hermann⁶, Sean Welleck^{2,3}, Amir Yazdanbakhsh⁶, Peter Clark²

¹Language Technologies Institute, Carnegie Mellon University

²Allen Institute for Artificial Intelligence

³University of Washington ⁴NVIDIA ⁵UC San Diego ⁶Google Research, Brain Team
amadaan@cs.cmu.edu, nikett@allenai.org

Self Consistency



GNOMON[®]
DIGITAL

Self-consistency

Q: If there are 3 cars in the parking lot and 2 more cars arrive, how many cars are in the parking lot?

A: There are 3 cars in the parking lot already. 2 more arrive. Now there are $3 + 2 = 5$ cars. The answer is 5.

...
Q: Janet's ducks lay 16 eggs per day. She eats three for breakfast every morning and bakes muffins for her friends every day with four. She sells the remainder for \$2 per egg. How much does she make every day?

A:

Language
model

Sample a diverse set of reasoning paths

She has $16 - 3 - 4 = 9$ eggs left. So she makes $\$2 * 9 = \18 per day.

This means she sells the remainder for $\$2 * (16 - 4 - 3) = \26 per day.

She eats 3 for breakfast, so she has $16 - 3 = 13$ left. Then she bakes muffins, so she has $13 - 4 = 9$ eggs left. So she has $9 \text{ eggs} * \$2 = \18 .

Marginalize out reasoning paths to aggregate final answers

The answer is \$18.

The answer is \$26.

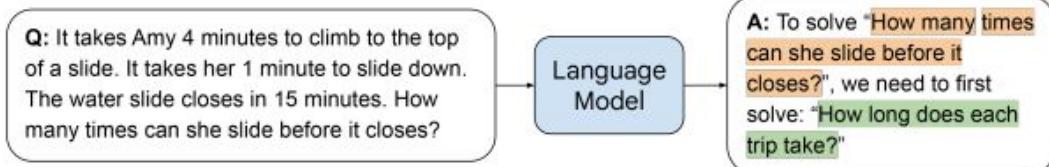
The answer is \$18.

The answer is \$18.

Self Consistency



Stage 1: Decompose Question into Subquestions



Stage 2: Sequentially Solve Subquestions

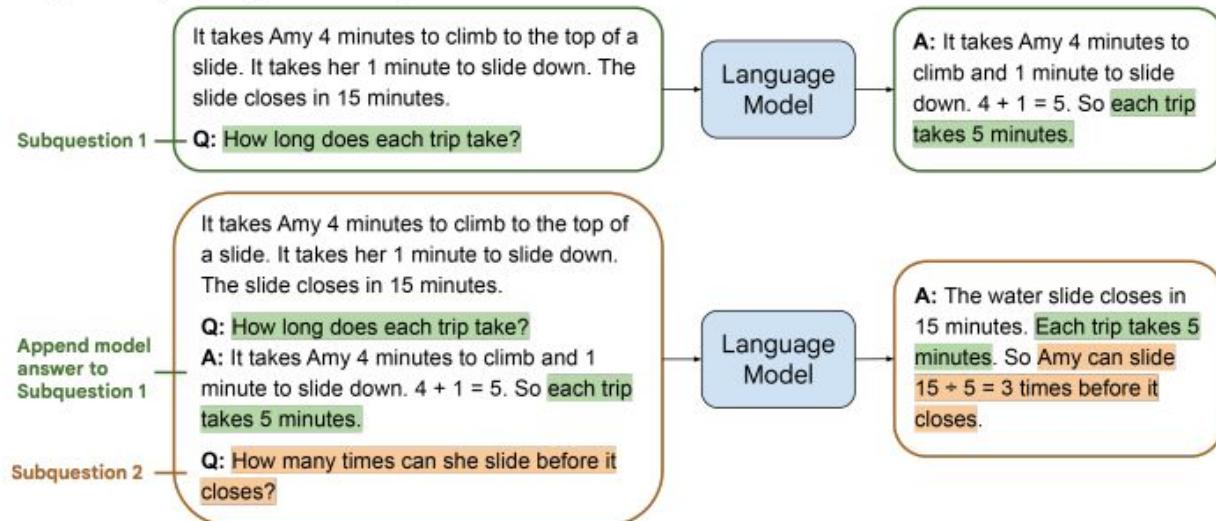


Figure 1: Least-to-most prompting solving a math word problem in two stages: (1) query the language model to decompose the problem into subproblems; (2) query the language model to sequentially solve the subproblems. The answer to the second subproblem is built on the answer to the first subproblem. The demonstration examples for each stage's prompt are omitted in this illustration.

Tree of Thoughts



Tree of Thoughts: Deliberate Problem Solving
with Large Language Models

Shunyu Yao Princeton University Dian Yu Google DeepMind Jeffrey Zhao Google DeepMind Izhak Shafran Google DeepMind

Thomas L. Griffiths Princeton University Yuan Cao Google DeepMind Karthik Narasimhan Princeton University

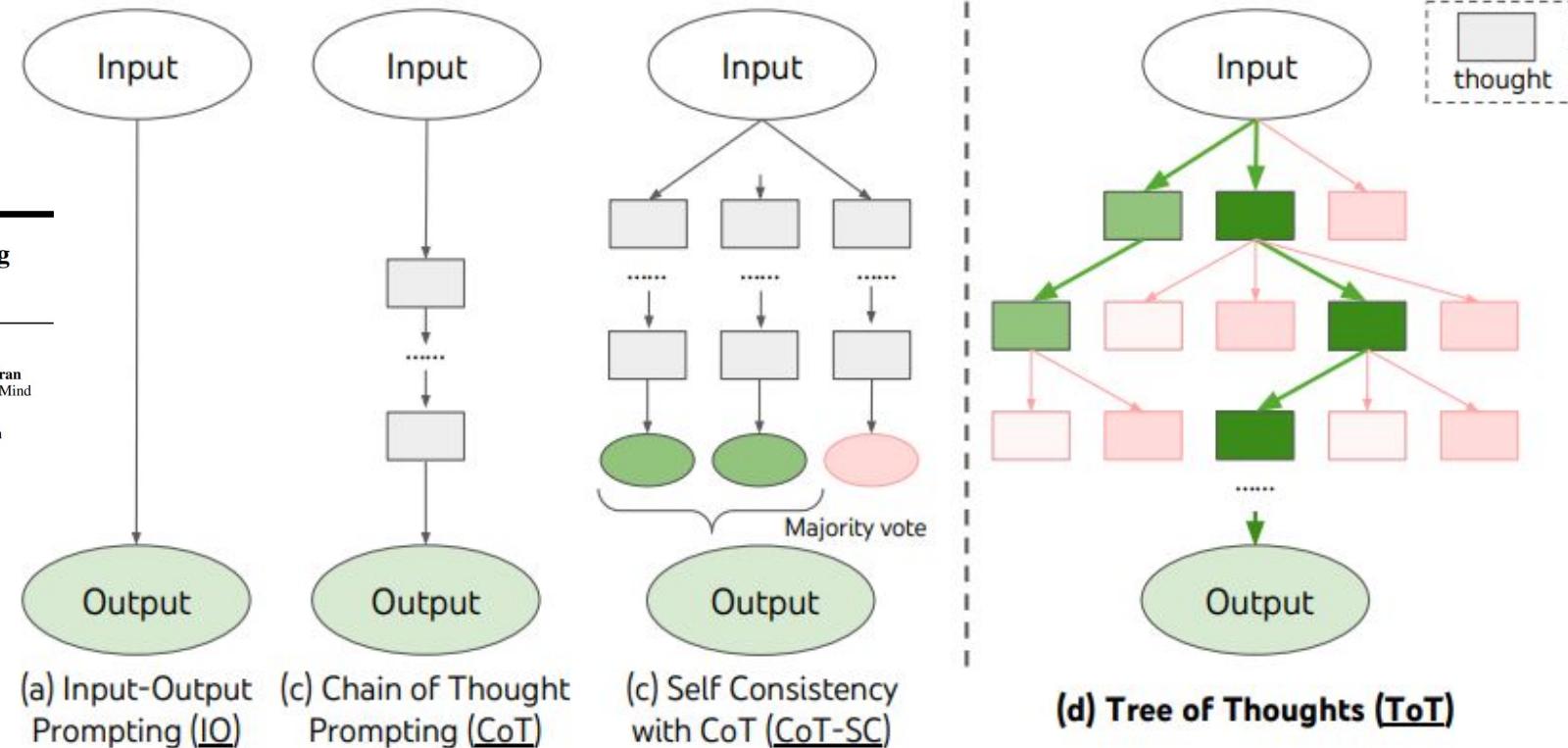


Figure 1: Schematic illustrating various approaches to problem solving with LLMs. Each rectangle box represents a *thought*, which is a coherent language sequence that serves as an intermediate step toward problem solving. See concrete examples of how thoughts are generated, evaluated, and searched in Figures 2,4,6.



La limite de LLM:

- Pas d'accès aux informations à jour

Quelle est le résultat d'élection présidentielle 2024

Cancel

Send



Je n'ai pas accès aux résultats de l'élection présidentielle de 2024 sans effectuer une recherche en direct, et je ne dispose pas de données en temps réel intégrées à ma mémoire. Si vous souhaitez connaître les résultats actualisés, je peux effectuer une recherche pour vous. 😊



Searched 2 sites ^

Q “résultat élection présidentielle 2024”
bing.com

Le Monde.fr — Présidentielle américaine 2024 : retrouvez les dernières informations après l'élection ...
lemonde.fr

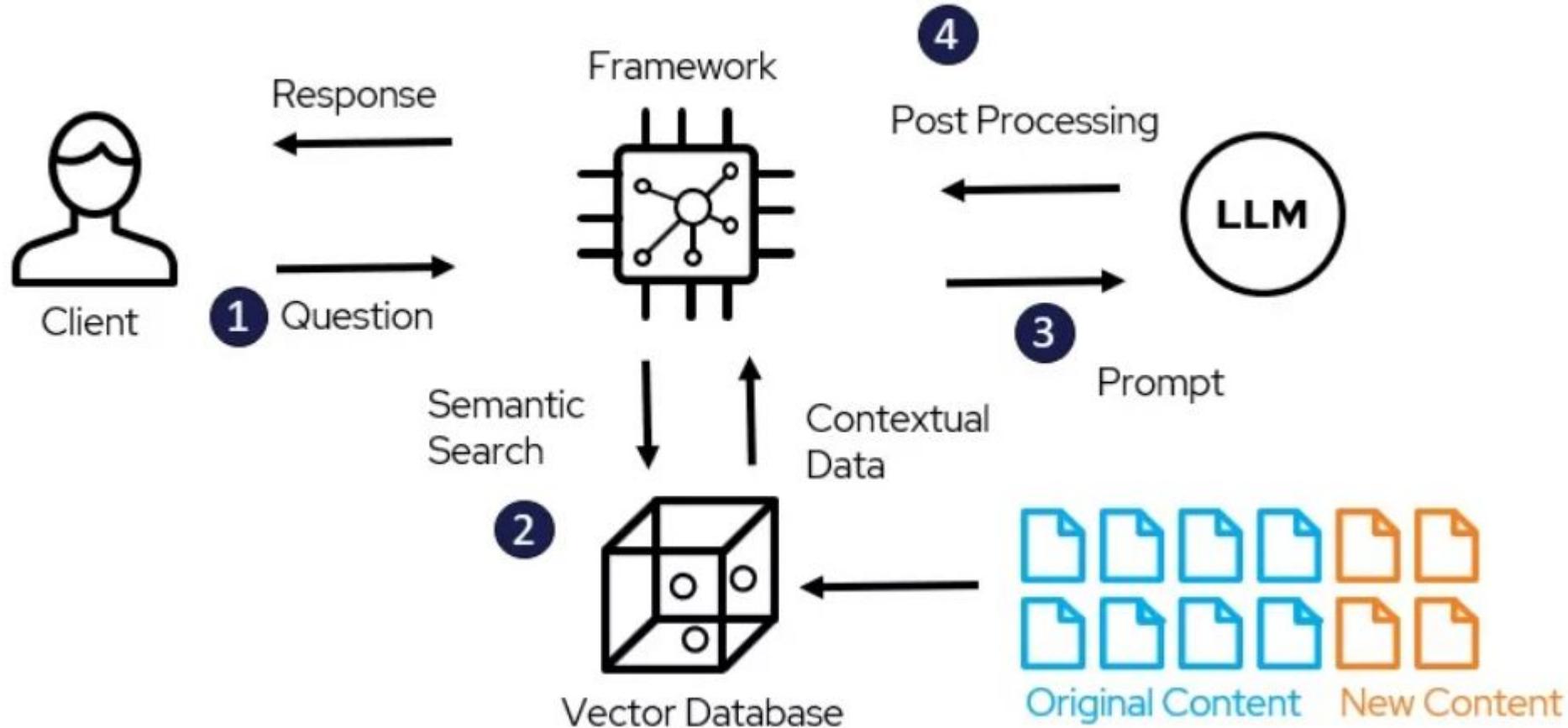
CNEWS — Présidentielle américaine 2024 : suivez les résultats de l'élection en temps réel sur notre...
cnews.fr

Lors de l'élection présidentielle américaine de 2024, Donald Trump a été déclaré vainqueur face à Kamala Harris. Trump a obtenu 312 votes des grands électeurs, dépassant largement la majorité requise de 270, tandis que Harris en a obtenu 226. Ce résultat marque le retour de Trump à la Maison-Blanche pour un second mandat après une campagne acharnée. L'investiture officielle est prévue pour le 20 janvier 2025

LE MONDE.FR CNEWS .

< 2/2 > ⏪ ⏴ ⏵ ⏵ ⏵

RAG Architecture Model

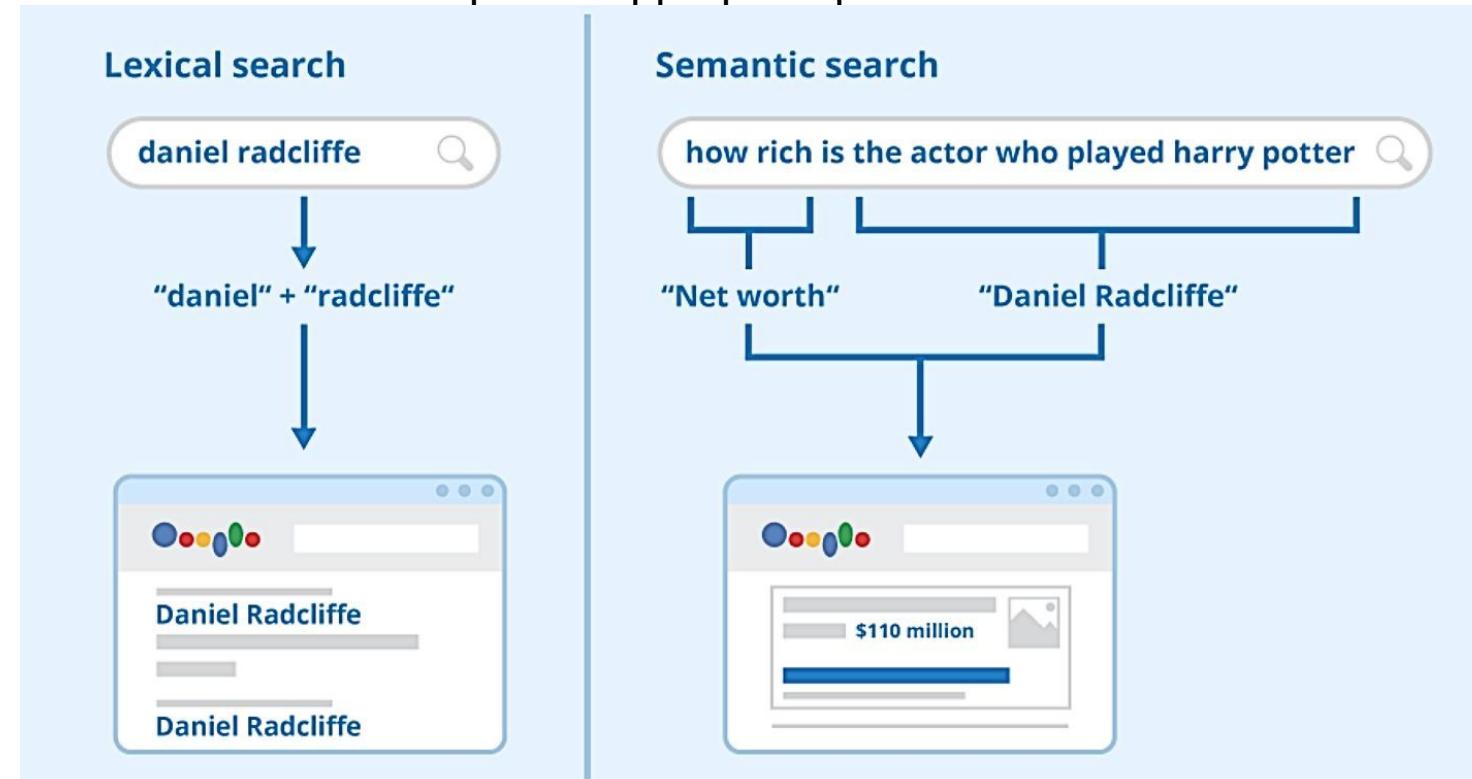


RAG: rechercher semantic



La notion d'adéquation sémantique entre deux entités X et Y s'applique à de nombreux contextes :

- Le sens de Y est-il similaire à celui de X ?
- X et Y représentent-ils la même chose présentée différemment ?
- Y est-il une bonne recommandation ou une réponse appropriée pour X ?



RAG: rechercher semantic



GNOMON[®]
DIGITAL

Lexical Search

Definition

Searches based on exact keyword match

Matching Criteria

Matches keywords or phrases literally

Natural Language Processing

Requires minimal NLP capabilities

Example

Searching for "Amazon river" returns info about amazon river, amazon.com etc.

Semantic Search

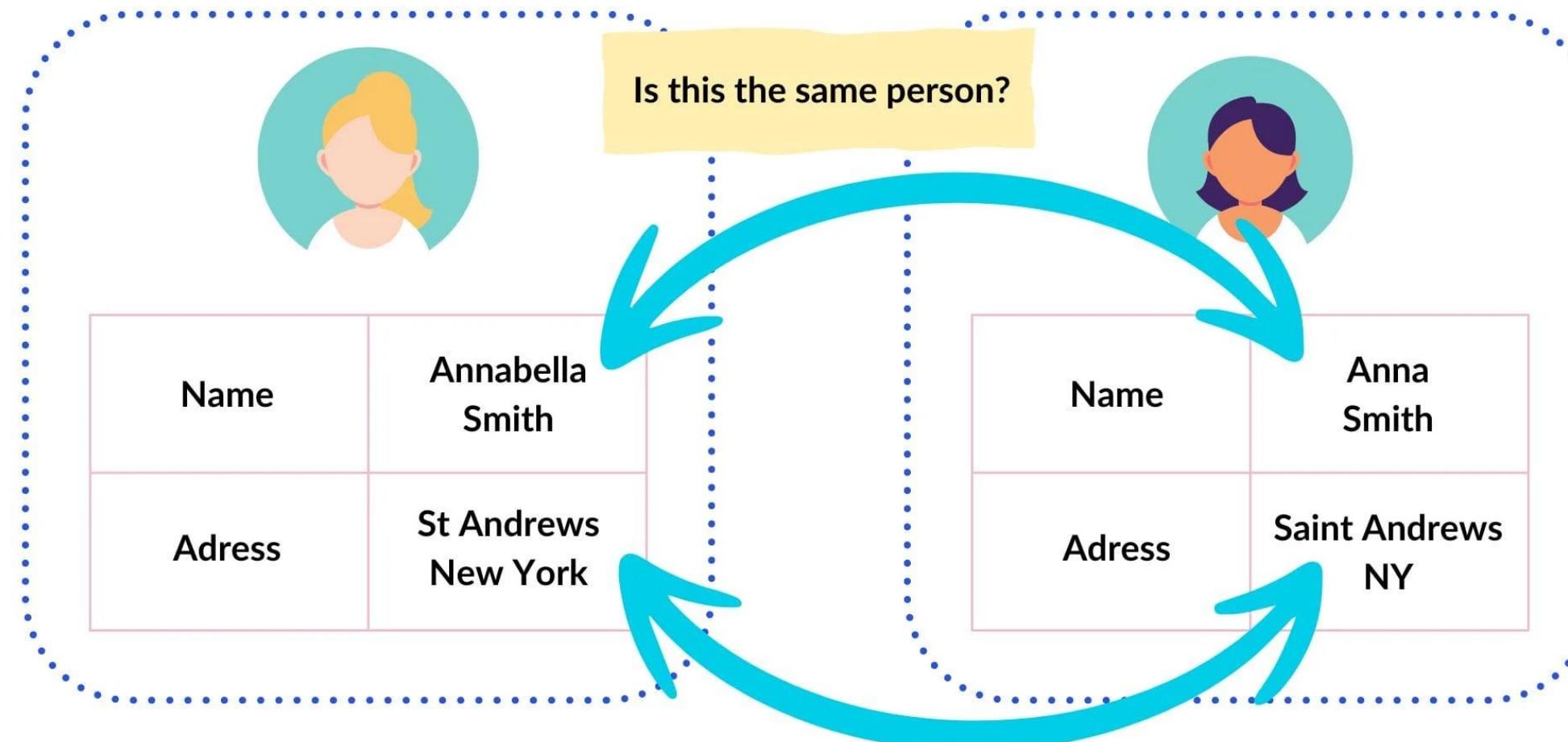
Understands query intent and context

Considers synonyms, context, and relationships between words

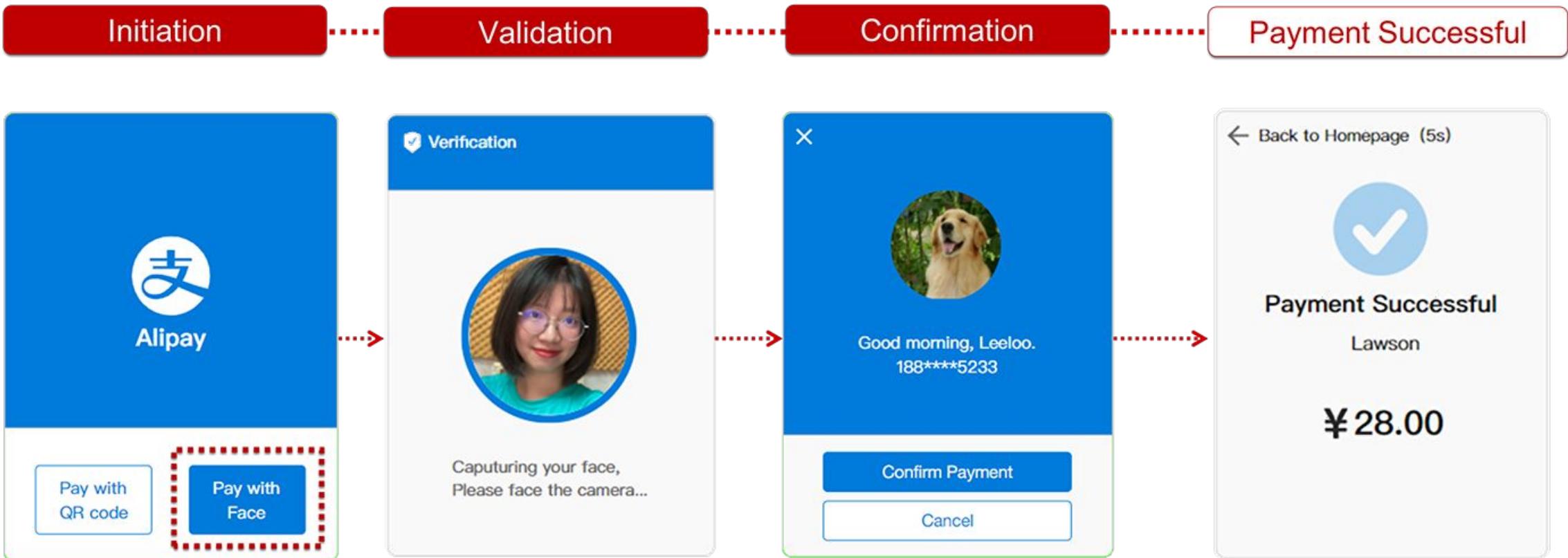
Requires advanced NLP for understanding context

Searching for "Amazon river" returns information about the Amazon river and not the company Amazon.

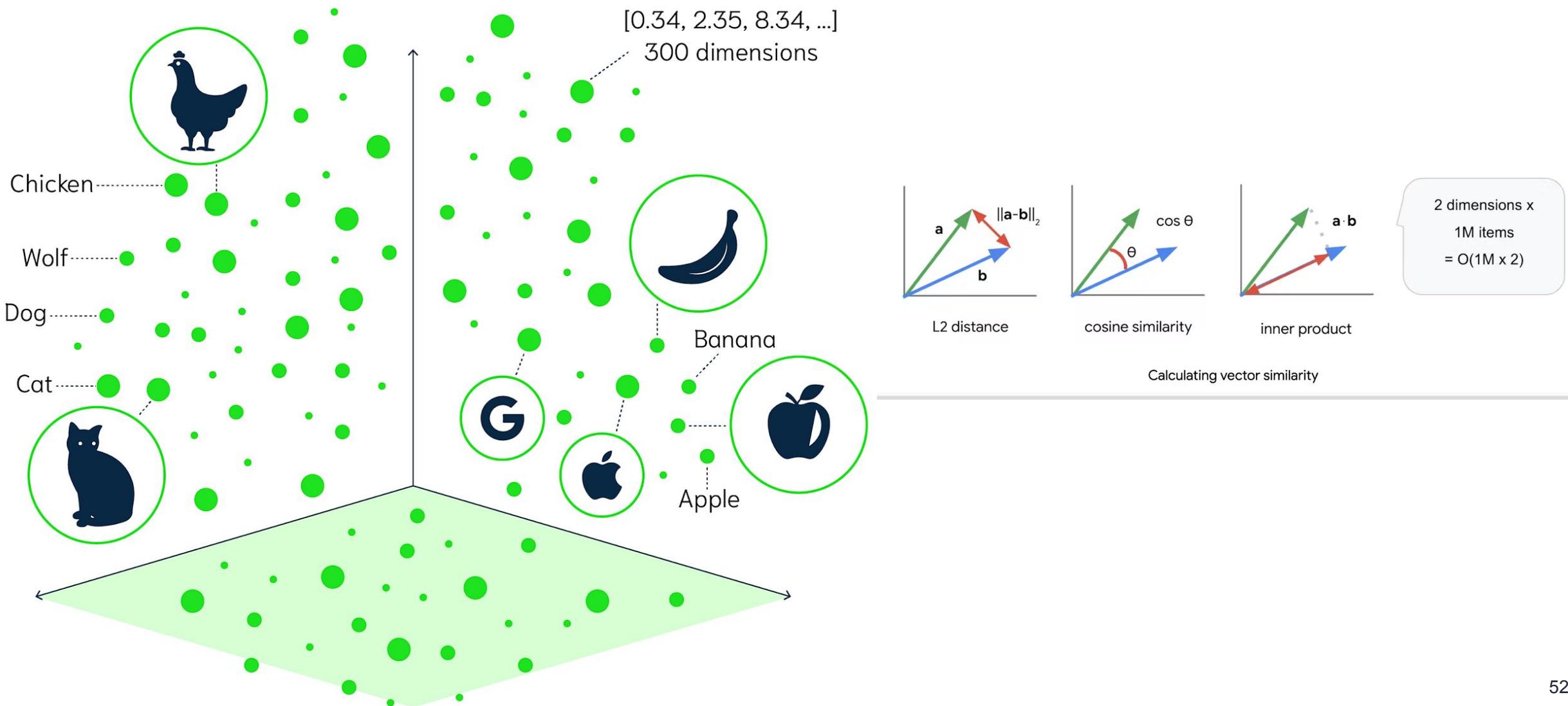
Entity Resolution



Process of Facial Recognition Payment



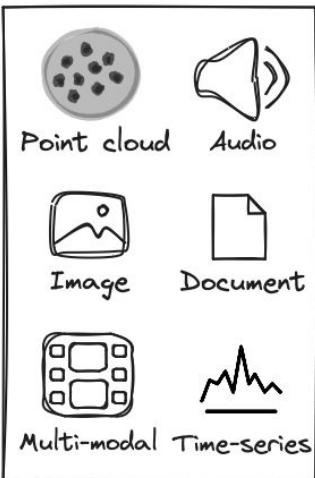
RAG: rechercher semantic



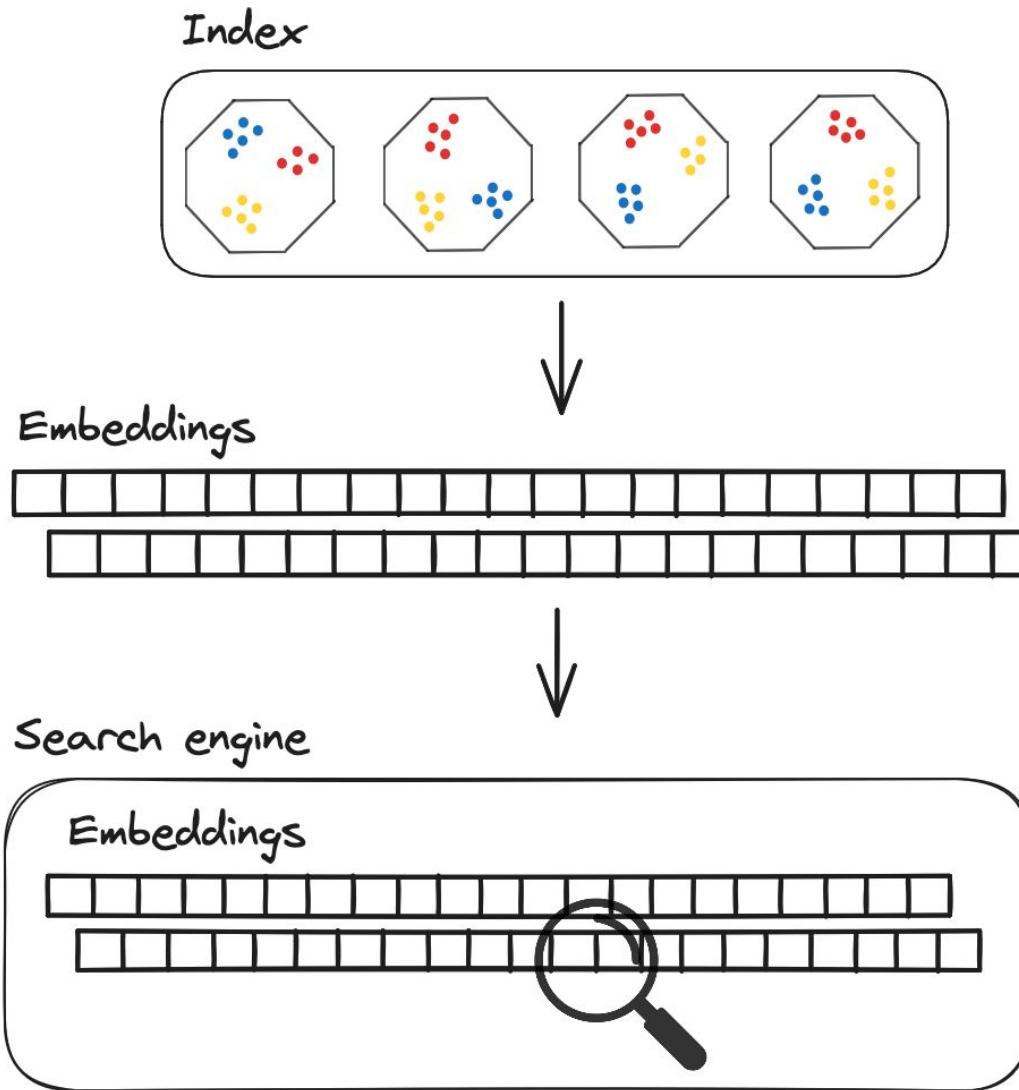
RAG: rechercher semantic



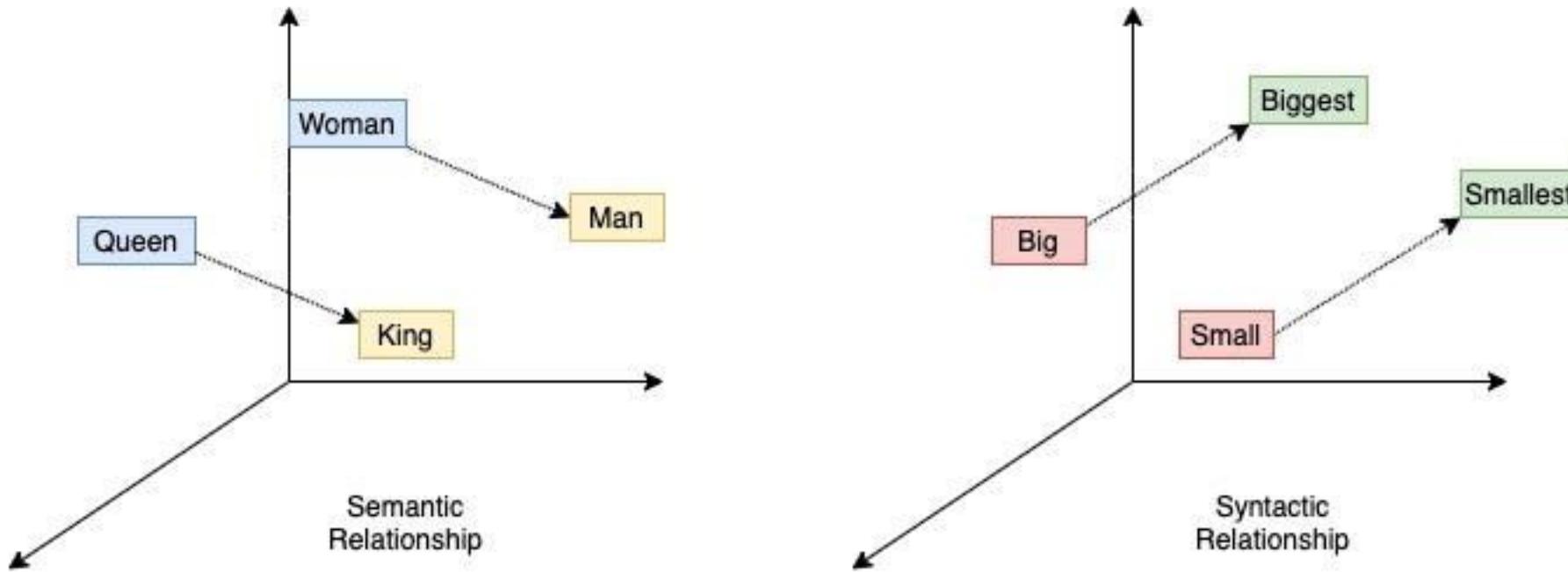
Data



Embedding
Model

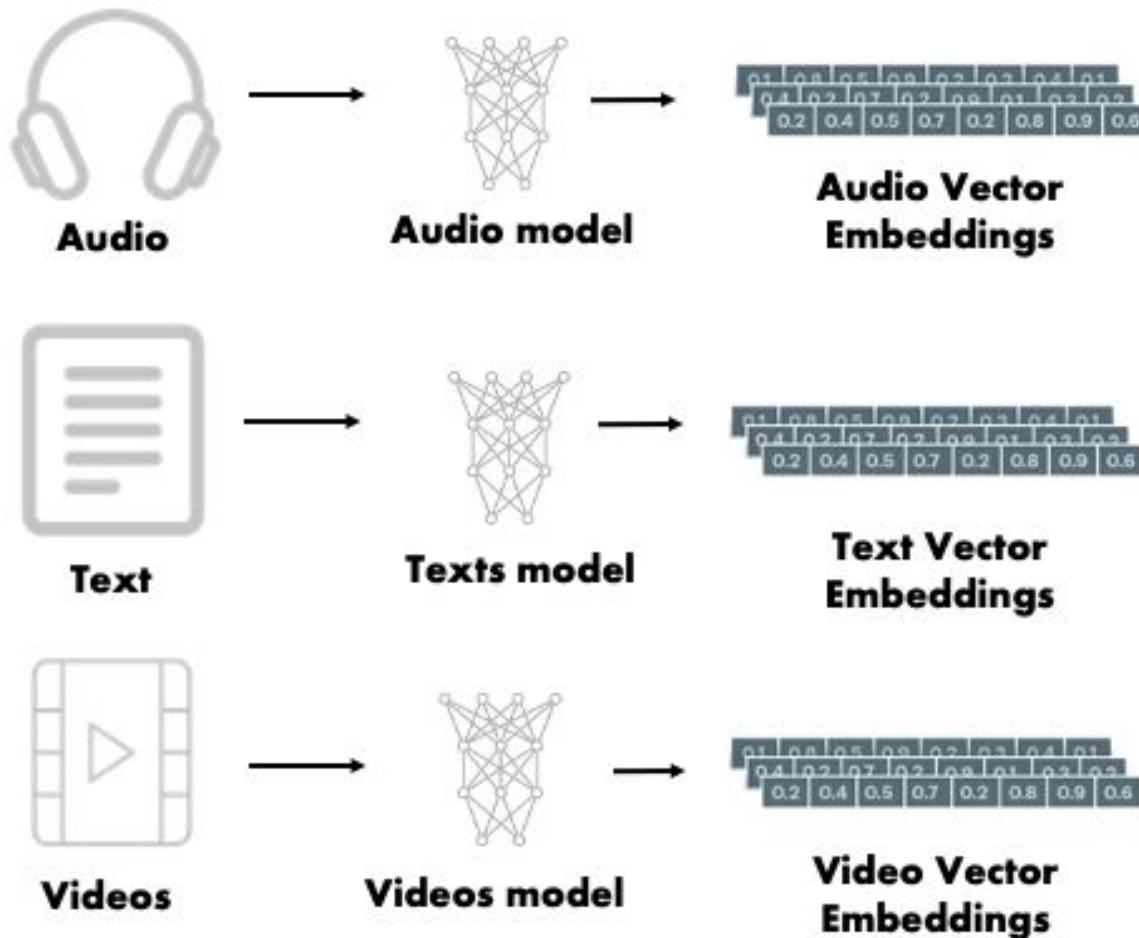


NLP



word2vec

RAG: rechercher semantic

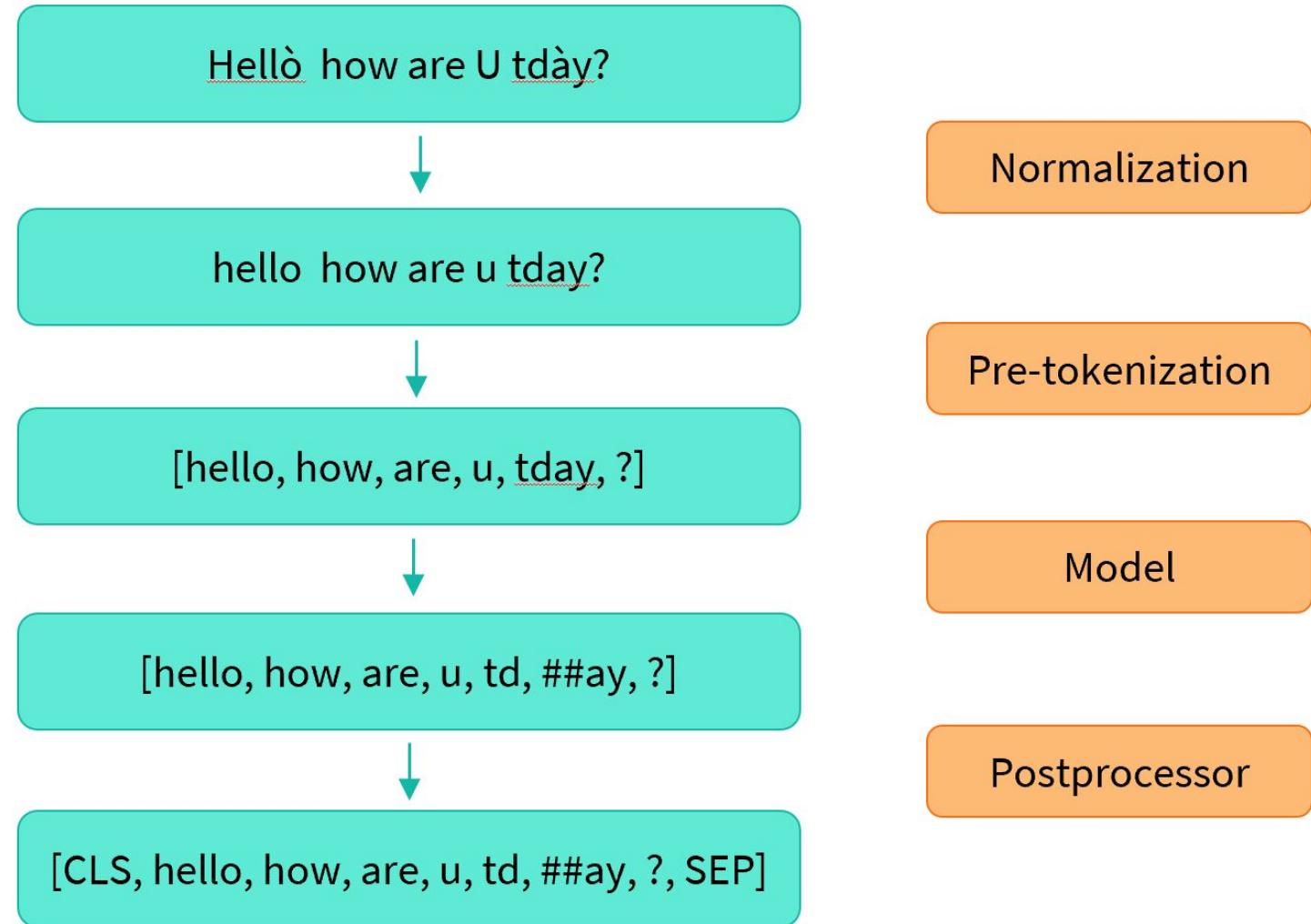


Very similar to word2vec

RAG: rechercher semantic



Step 1: Tokenize

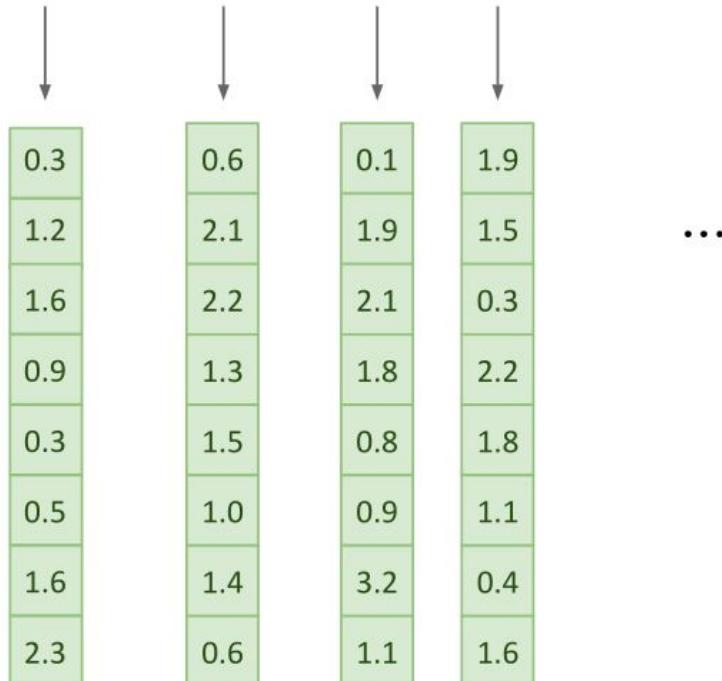


RAG: rechercher semantic



Step 2: Learn embedding vectors

[100, 2031, 3, 18, 10, 8, ...]



RAG: rechercher semantic



Step 2: Learn embedding vectors

Different types of embeddings

Type 1

Non-Contextual
Embeddings

Factorization Based

- Word2Vec
- Gloves
- Swivel

Type 2

Contextual
Embeddings

Language Model Based

- Bert
- PaLM
- GPT

RAG: rechercher semantic

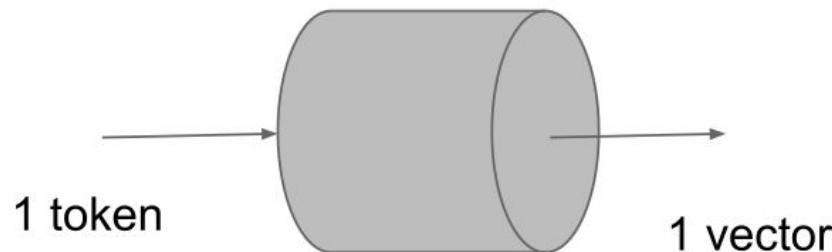


Step 2: Learn embedding vectors

Different types of embeddings

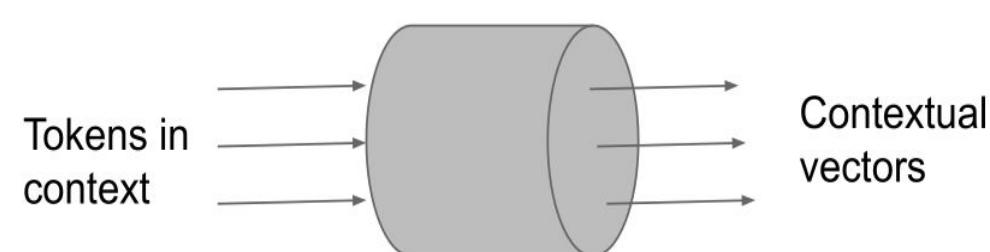
Type 1

Non-Contextual
Embeddings



Type 2

Contextual
Embeddings



RAG: rechercher semantic

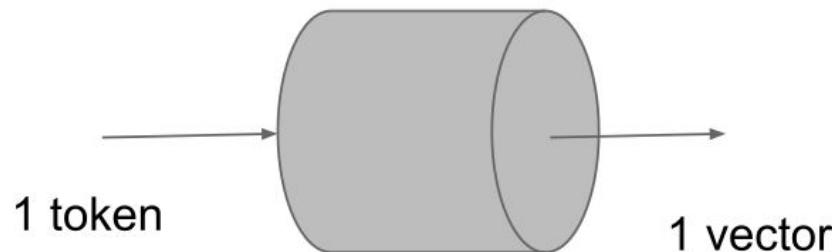


Step 2: Learn embedding vectors

Different types of embeddings

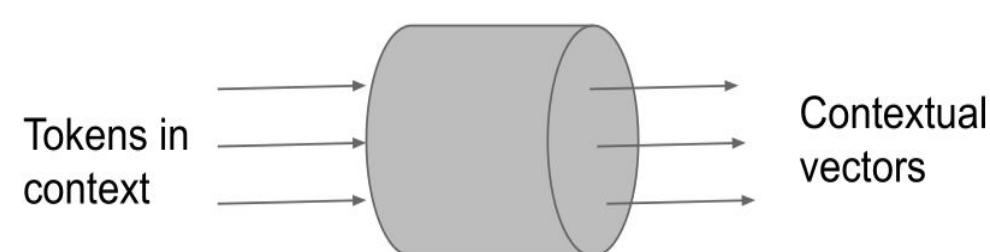
Type 1

Non-Contextual
Embeddings



Type 2

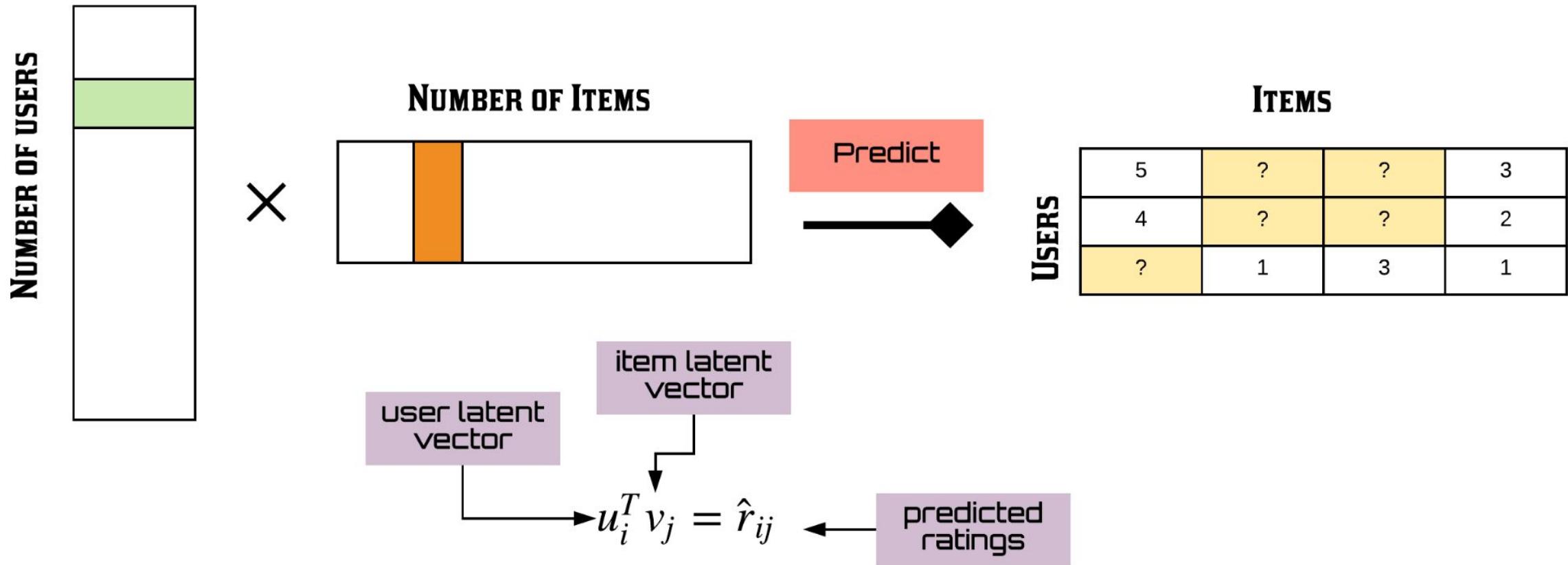
Contextual
Embeddings



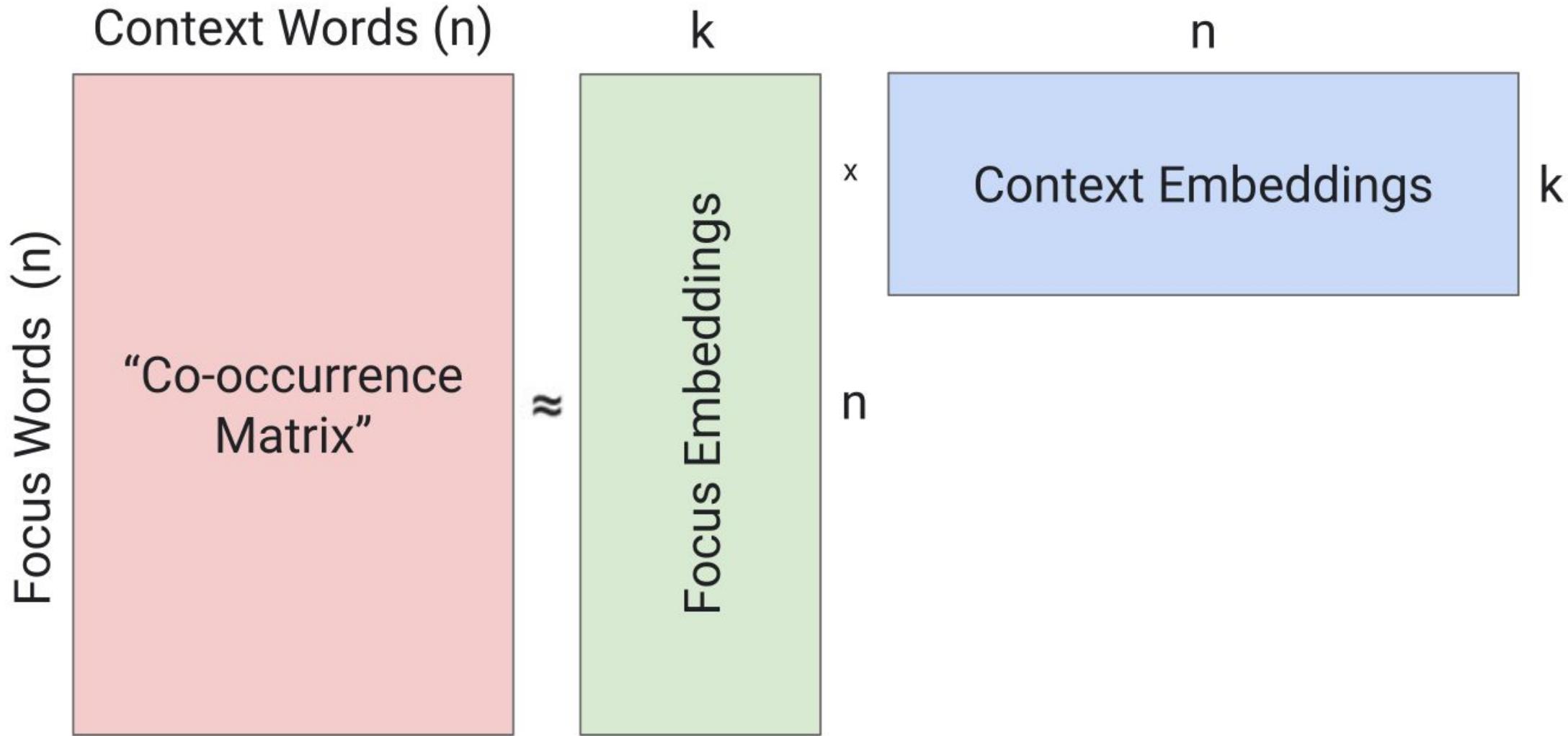
RAG: rechercher semantic



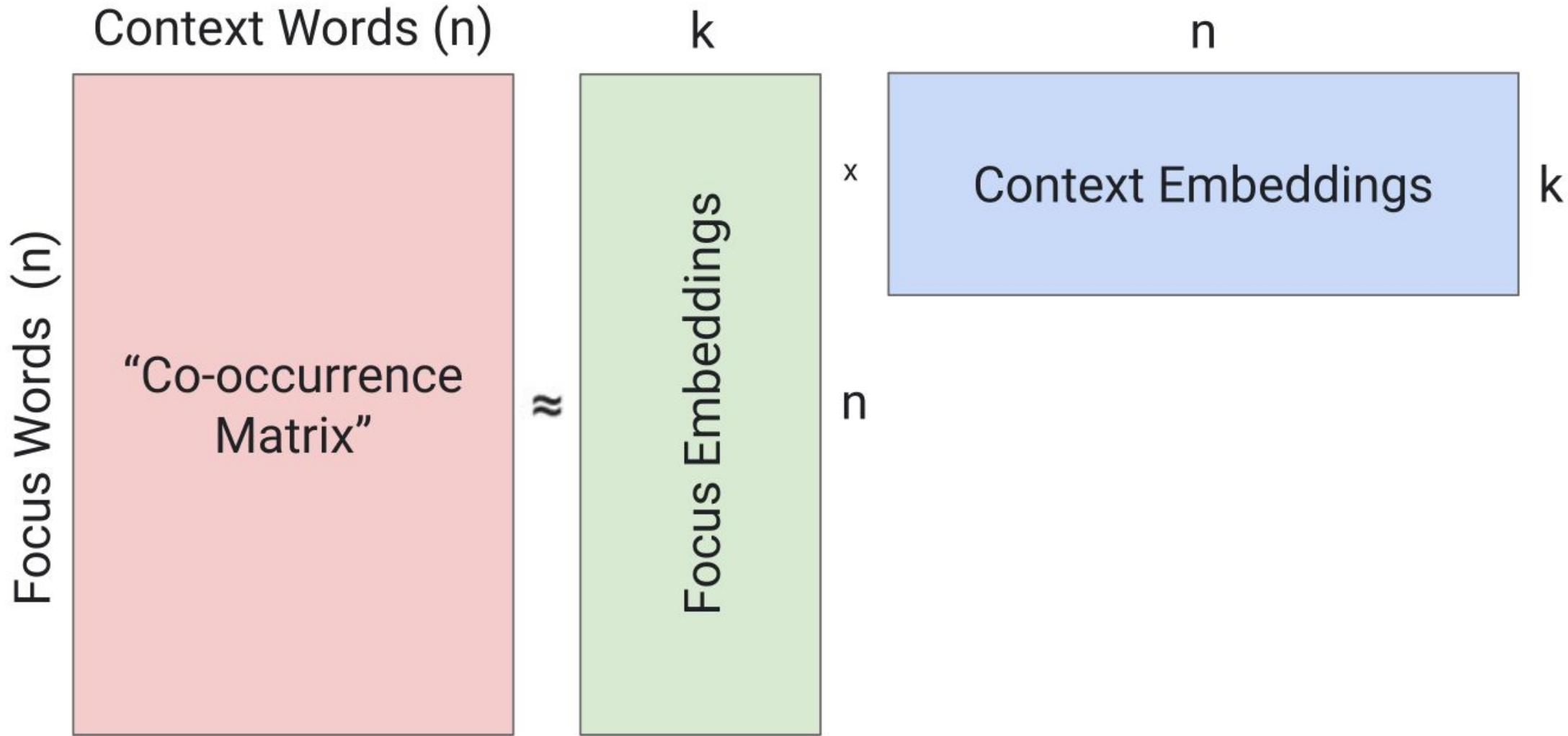
MATRIX FACTORIZATION



RAG: rechercher semantic



RAG: rechercher semantic



RAG: rechercher semantic

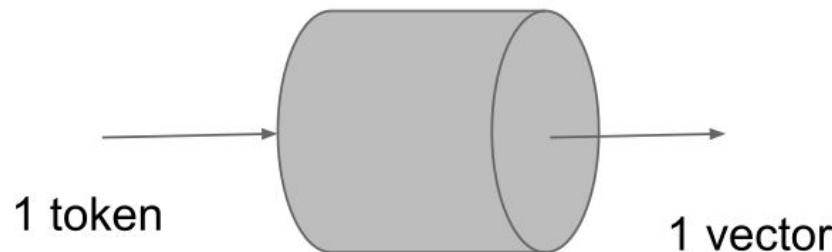


Step 2: Learn embedding vectors

Different types of embeddings

Type 1

Non-Contextual
Embeddings



Type 2

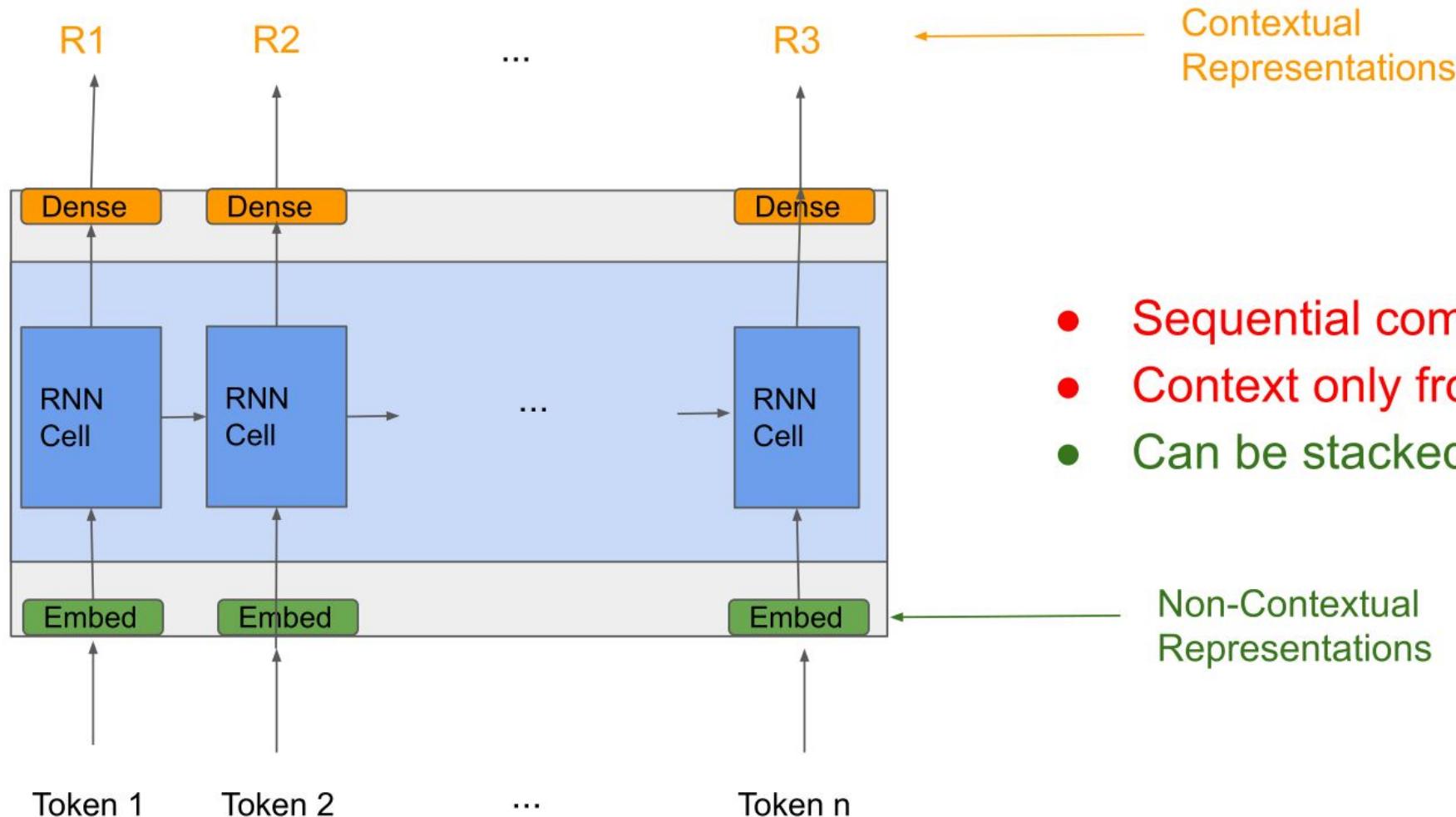
Contextual
Embeddings



RAG: rechercher semantic



Contextual: Recurrent Network Language Models

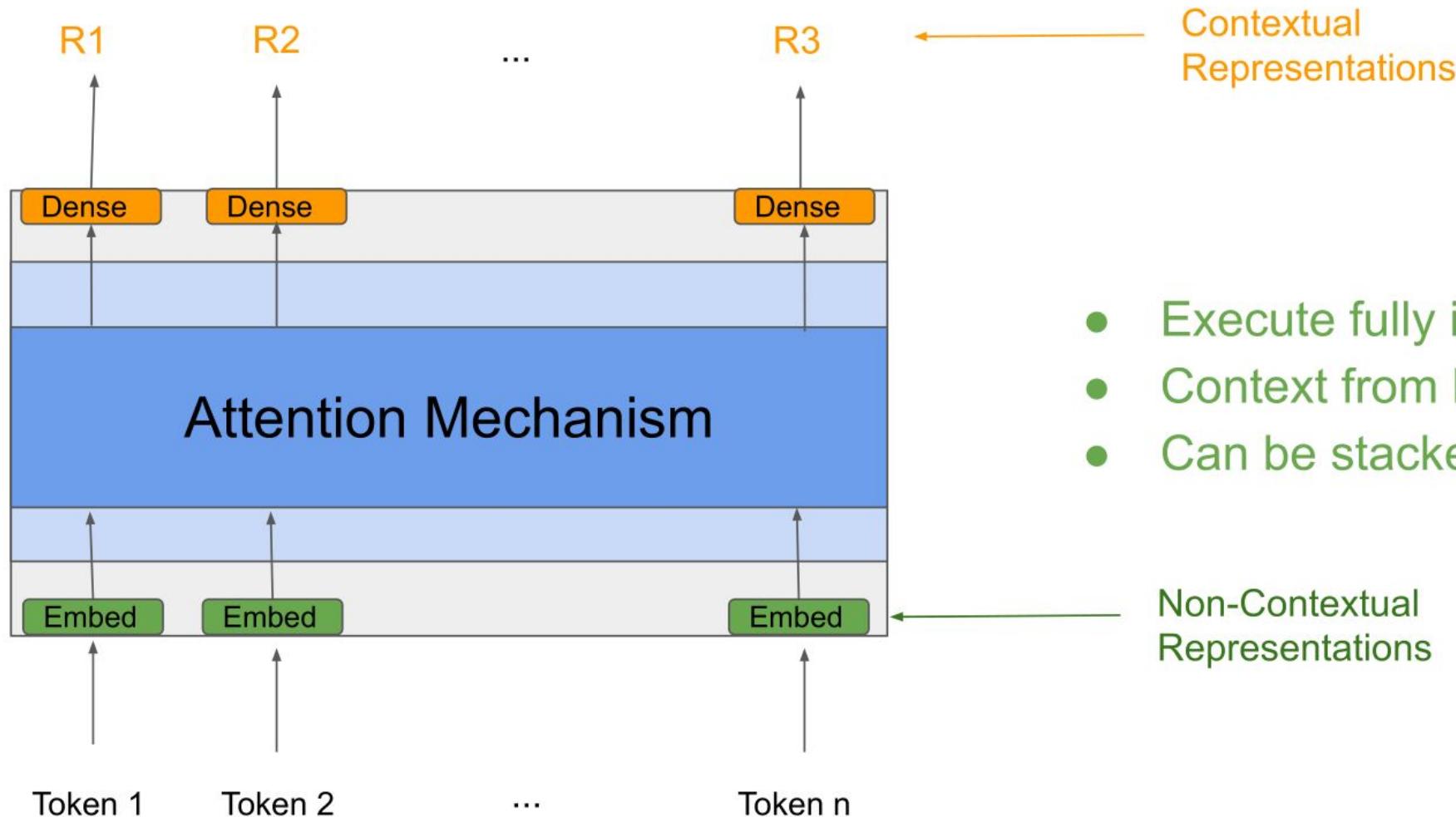


- Sequential computation
- Context only from one side
- Can be stacked

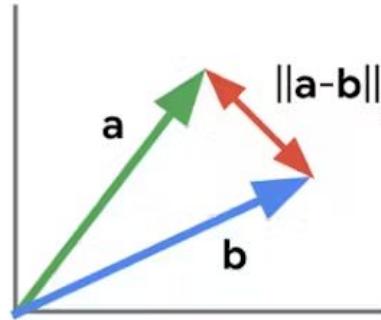
RAG: rechercher semantic



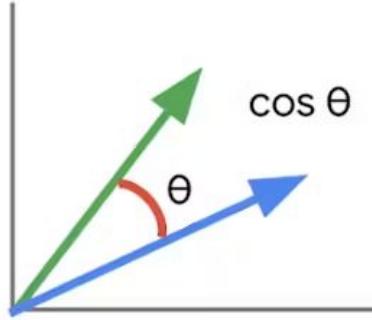
Contextual: Transformer Language Models



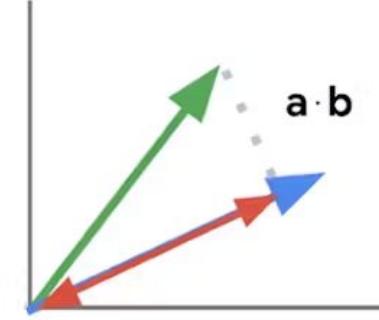
RAG: rechercher semantic



L2 distance



cosine similarity



inner product

Calculating vector similarity

2 dimensions x
1M items
 $= O(1M \times 2)$

RAG: rechercher semantic

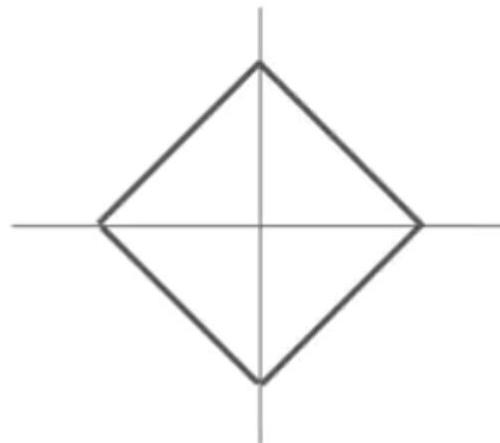


L2 distance (Euclidean) distance:

K-Nearest Neighbors: Distance Metric

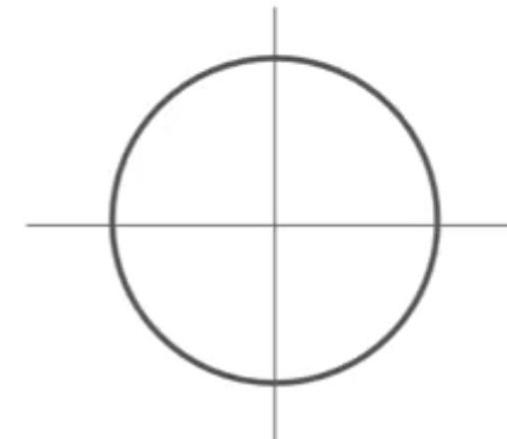
L1 (Manhattan) distance

$$d_1(I_1, I_2) = \sum_p |I_1^p - I_2^p|$$



L2 (Euclidean) distance

$$d_2(I_1, I_2) = \sqrt{\sum_p (I_1^p - I_2^p)^2}$$

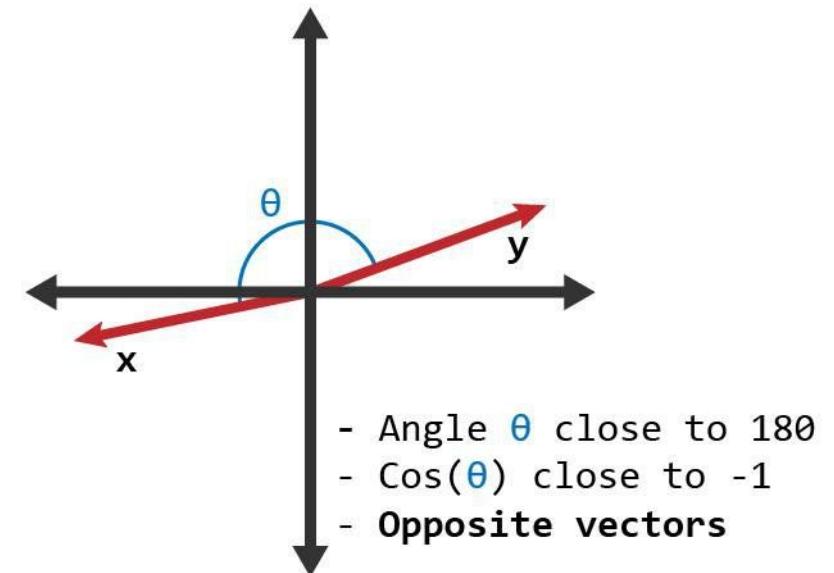
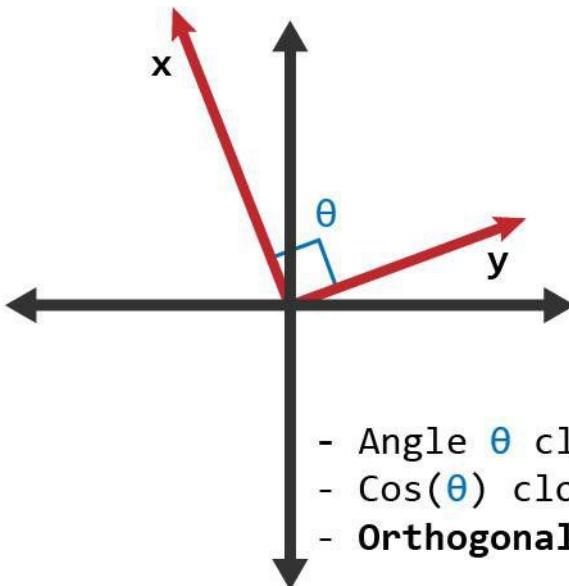
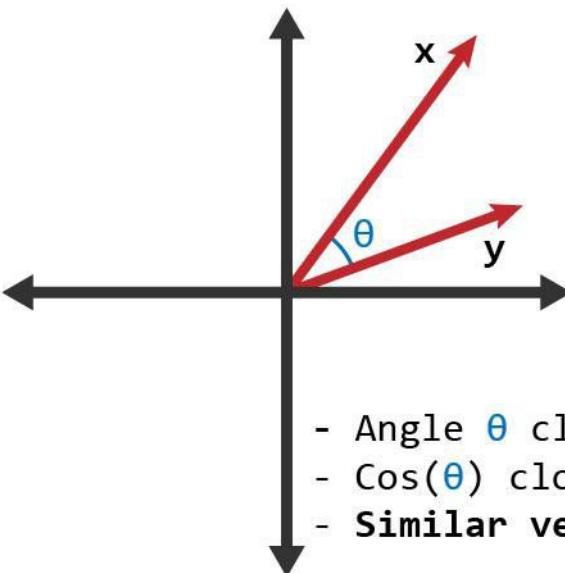


RAG: rechercher semantic



Cos similarity:

$$\cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$



RAG: rechercher semantic

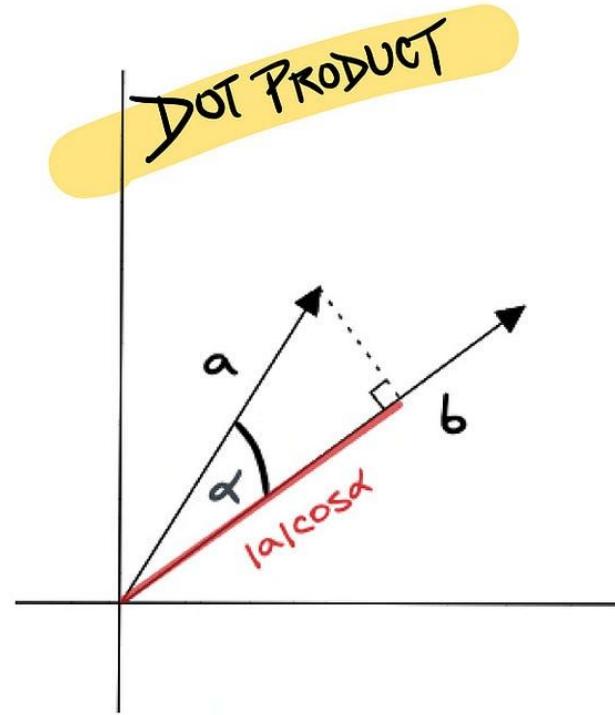


Dot product:

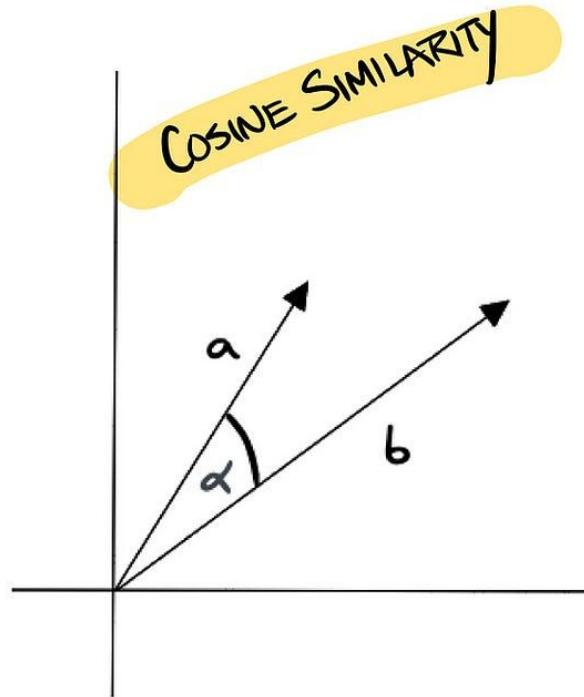
$$\mathbf{a} = \langle a_1, a_2, a_3 \rangle \quad \mathbf{b} = \langle b_1, b_2, b_3 \rangle$$

$$\mathbf{a} \cdot \mathbf{b} = a_1 b_1 + a_2 b_2 + a_3 b_3$$

RAG: rechercher semantic

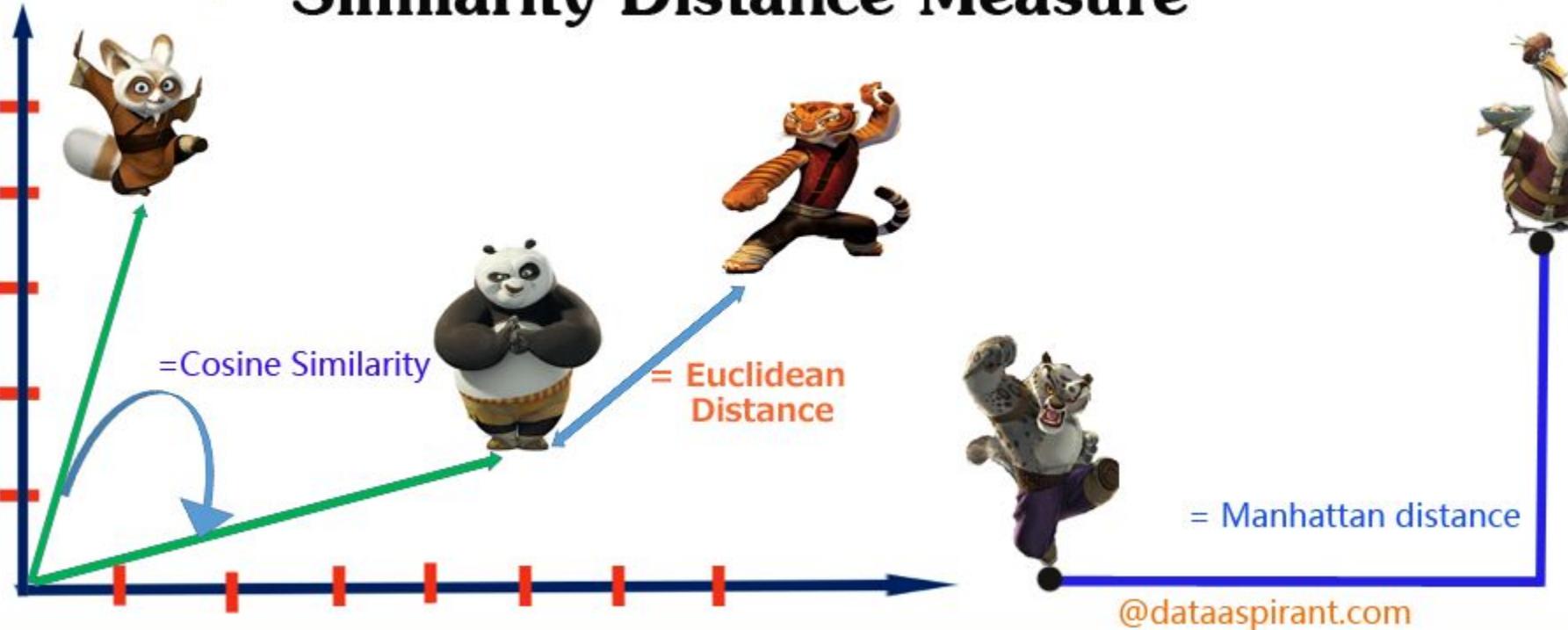


$$\mathbf{a} \cdot \mathbf{b} = |\mathbf{a}| |\mathbf{b}| \cos \alpha$$



$$sim(\mathbf{a}, \mathbf{b}) = \frac{\mathbf{a} \cdot \mathbf{b}}{||\mathbf{a}|| \cdot ||\mathbf{b}||}$$

Similarity Distance Measure



RAG: rechercher semantic



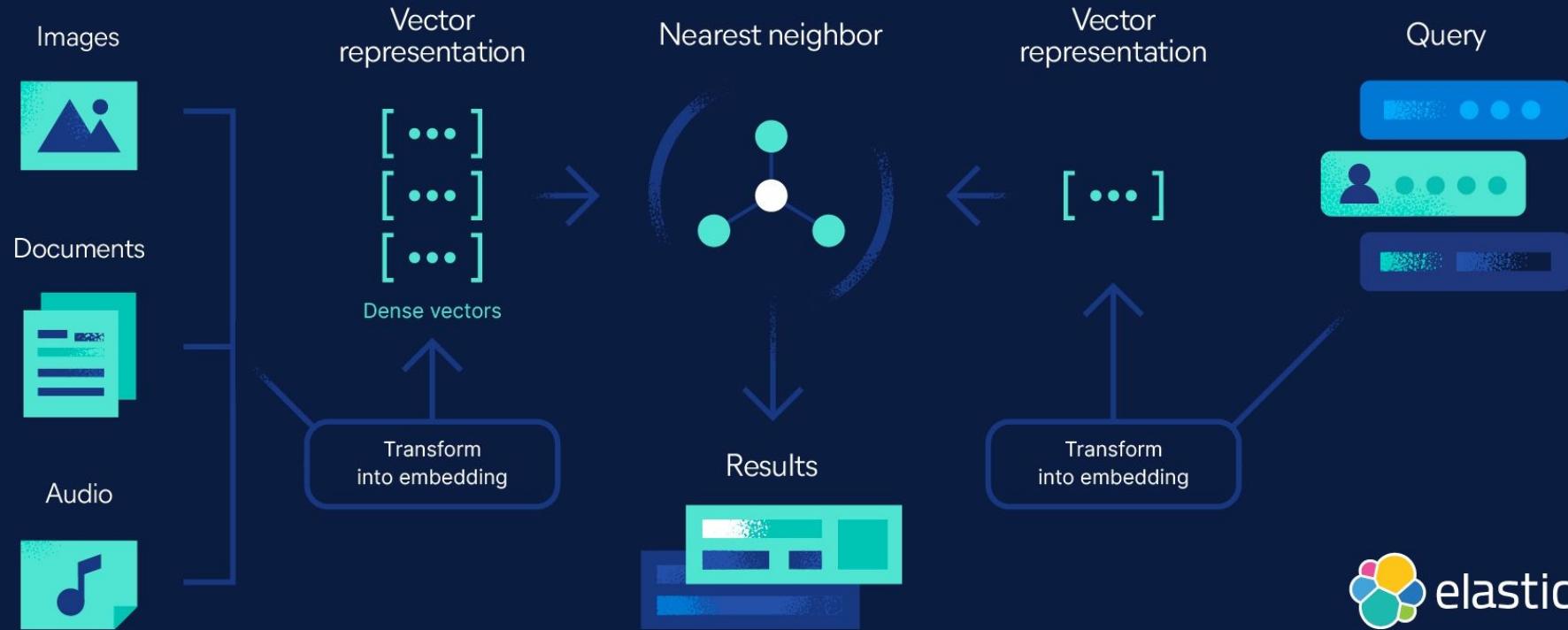
GNOMON[®]
DIGITAL

https://github.com/GoogleCloudPlatform/asl-ml-immersion/blob/master/notebooks/vertex_genai/labs/semantic_matching_with_gemini.ipynb

RAG: rechercher semantic



Vector search data workflow



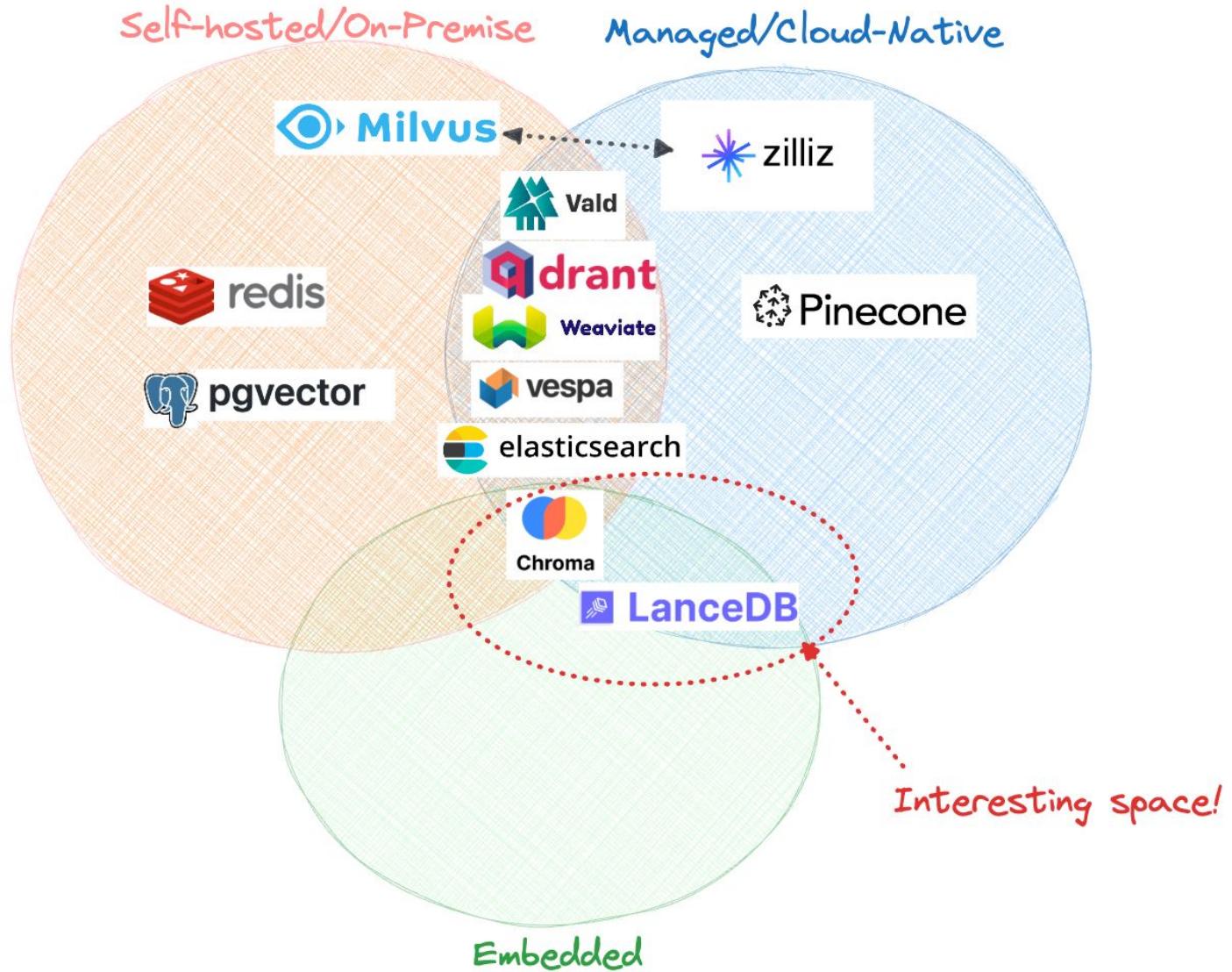
RAG: rechercher semantic



Recherche par vectorielle est splitté en deux parties:

- Offline -> Calcule lourde /lente:
 - Ajouter un vector index (keys= vectors)
 - Sauve tous les embedding vecteurs dans ces index
- Online -> Calcule efficace / Rapid
 - Algorithm search dans vectorDB

RAG: rechercher semantic



RAG: rechercher semantic



GNOMON[®]
DIGITAL

https://github.com/GoogleCloudPlatform/asl-ml-immersion/blob/master/notebooks/vertex_genai/solutions/semantic_search_with_vector_search.ipynb

Introduction à LangChain et LlamaIndex



Introduction à LangChain et LlamaIndex



Untitled prompt

System Instructions

Optional tone and style instructions for the model

What will you build?

Push Gemini to the limits of what AI can do

Trip recommendations

Convert unorganized text into structured tables.

Unit Testing

Add unit tests for a Python function.

Modify writing style

Change the tone and writing style of a blurb.

← Compare ↗ Get code ⋮

Run settings

Reset

Model

Gemini 1.5 Flash

Token count

0 / 1000 000

Temperature

1

Tools

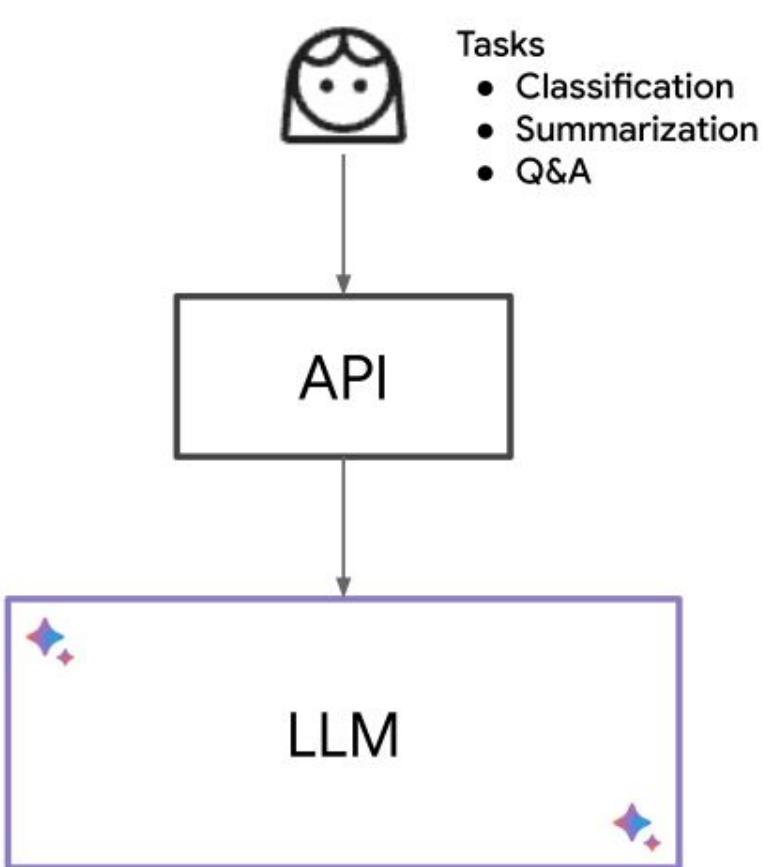
Structured output

Edit schema

Code execution

Function calling

Edit functions

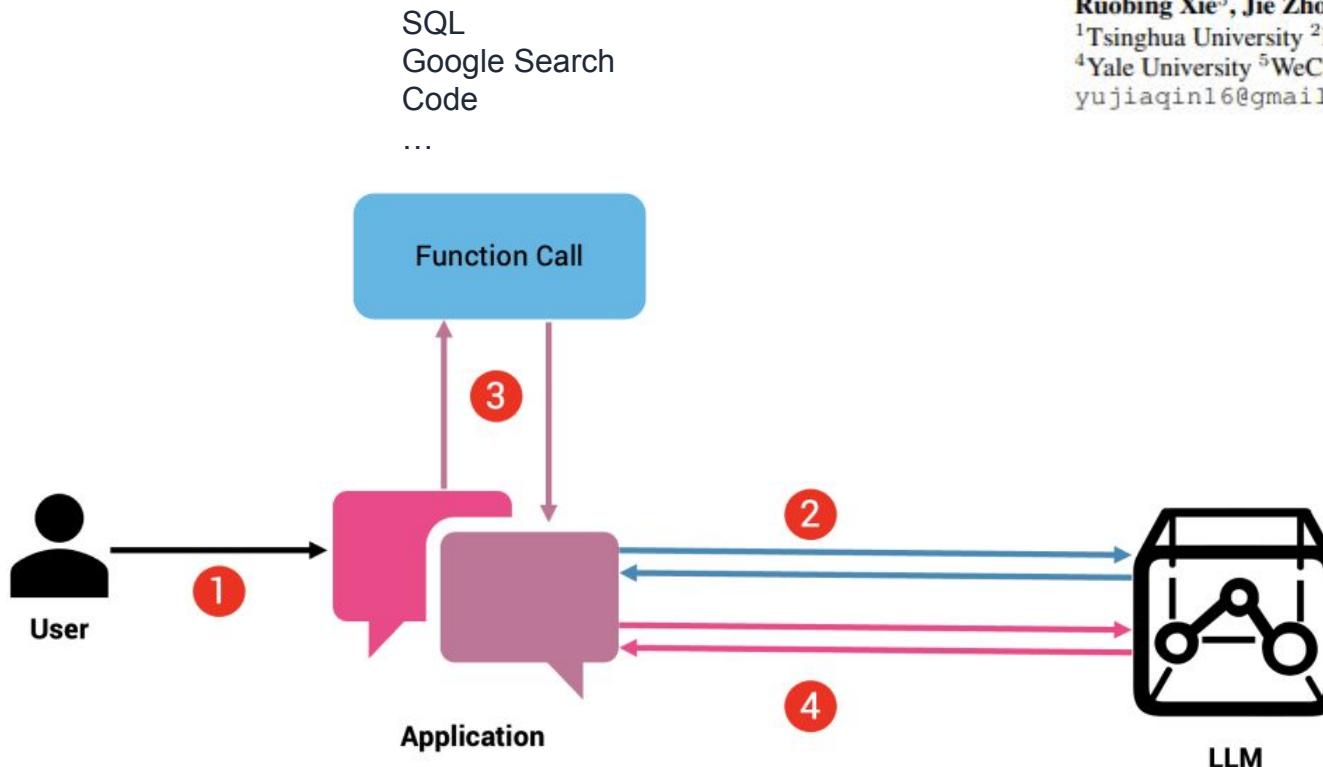


https://aistudio.google.com/app/prompts/new_chat

Introduction à LangChain et LlamaIndex



Tool LLM



TOOLLM: FACILITATING LARGE LANGUAGE MODELS TO MASTER 16000+ REAL-WORLD APIs

Yujia Qin^{1*}, Shihao Liang^{1*}, Yining Ye¹, Kunlun Zhu¹, Lan Yan¹, Yaxi Lu¹, Yankai Lin^{3†}, Xin Cong¹, Xiangru Tang⁴, Bill Qian⁴, Sihan Zhao¹, Lauren Hong¹, Runchu Tian¹, Ruobing Xie⁵, Jie Zhou⁵, Mark Gerstein⁴, Dahai Li^{2,6}, Zhiyuan Liu^{1†}, Maosong Sun^{1†}

¹Tsinghua University ²ModelBest Inc. ³Renmin University of China

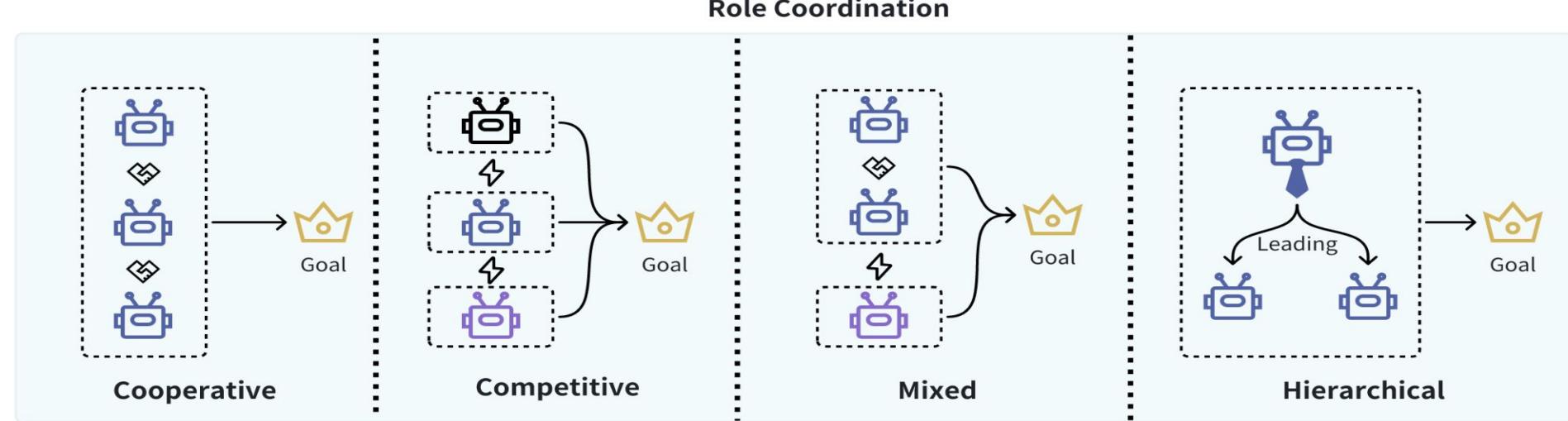
⁴Yale University ⁵WeChat AI, Tencent Inc. ⁶Zhihu Inc.

yujiaqin16@gmail.com

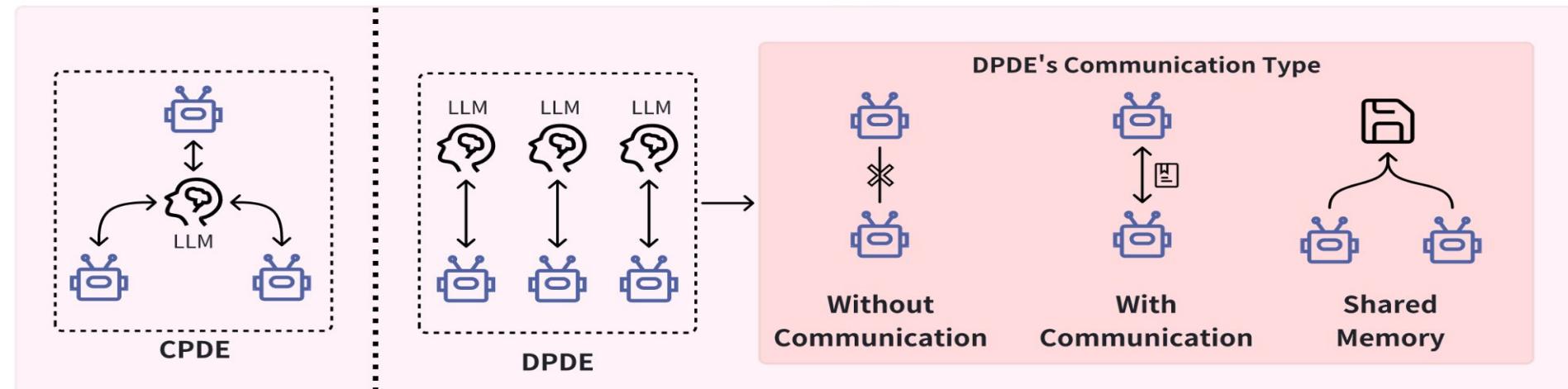
Introduction à LangChain et LlamaIndex



LLM Agent



Planning Type



Introduction à LangChain et LlamaIndex



LLM Agent

Exploring Large Language Model based Intelligent Agents: Definitions, Methods, and Prospects

Yuheng Cheng¹ ✉ **Ceyao Zhang¹** ✉ **Zhengwen Zhang¹**

Xiangrui Meng¹ **Sirui Hong²** **Wenhai Li¹** **Zihao Wang³** **Zekai Wang⁴**

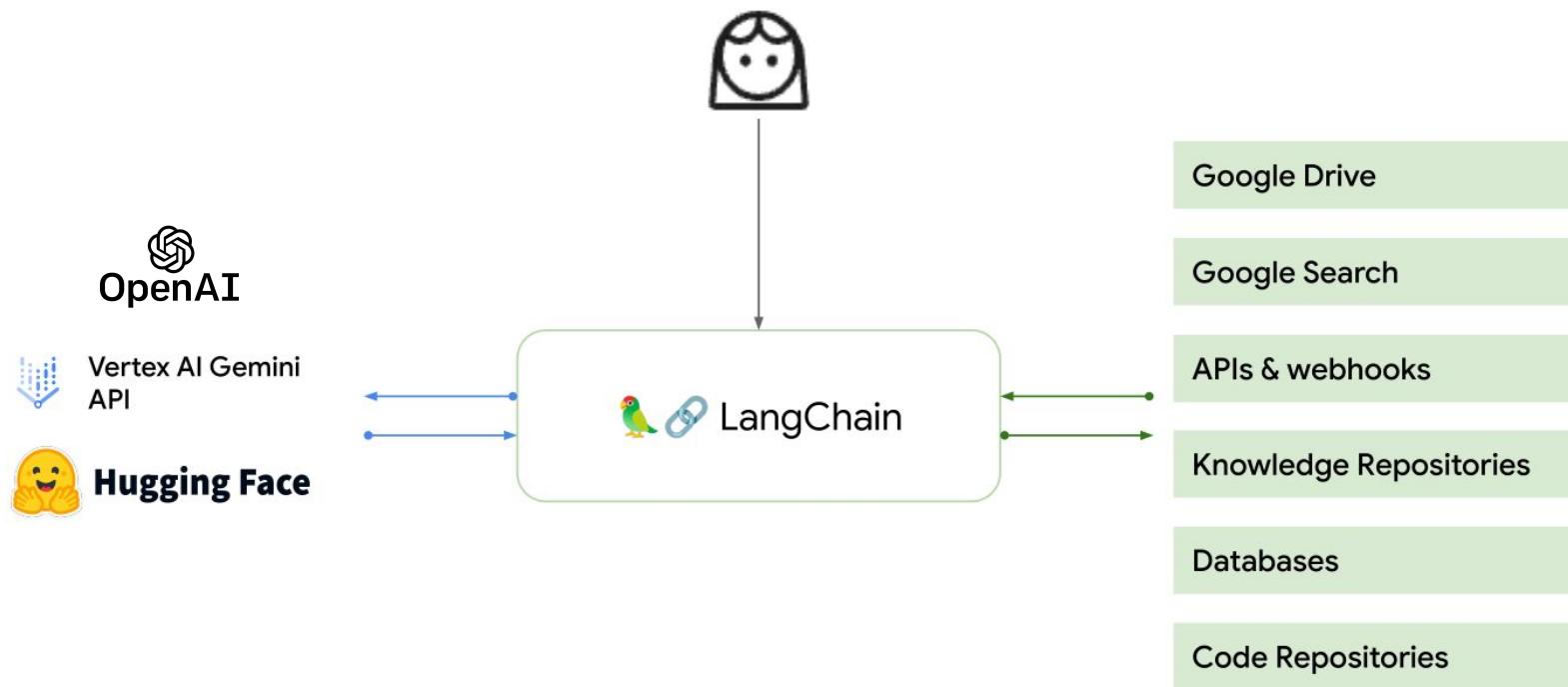
Feng Yin¹ **Junhua Zhao¹** **Xiuqiang He⁵**

¹The Chinese University of Hong Kong, Shenzhen

²DeepWisdom ³Peking University ⁴Yantu.ai ⁵FiT, Tencent

These authors contributed to the work equally and should be regarded as co-first authors. Corresponding author

Introduction à LangChain et LlamaIndex



Usages possibles:

- Question & Answering over private corpus
- Large document Summarization
- Personal assistants & Chatbots
- Querying tabular data
- Code understanding
- API interactions
- Extraction
- Evaluation

Introduction à LangChain et LlamaIndex



LlamaIndex vs. LangChain

data science dojo
data science for everyone

Features	LlamaIndex	LangChain
Primary focus	Intelligent search, data indexing and retrieval	Building a wide range of Gen AI applications
Data handling	Ingesting, structuring, and accessing private or domain-specific data	Loading, processing, and indexing data for various uses
Customization	Offers tools for integrating private data into LLMs	Highly customizable, users can chain multiple components
Flexibility	Specialized for efficient and fast search	General-purpose framework with flexibility in application behavior
Deployment	Ideal for proprietary or specialized data	Facilitates the deployment of bespoke NLP applications
Use cases	Best for applications that require quick data lookup and retrieval	Suitable for applications that require complex interactions like chatbots, GQA, summarization

Introduction à LangChain



Deployment

LangGraph Platform

COMMERCIAL

Components

Integrations

OSS

Architecture

LangChain

OSS

LangGraph

OSS

LangSmith

Debugging

Playground

Prompt Management

Annotation

Testing

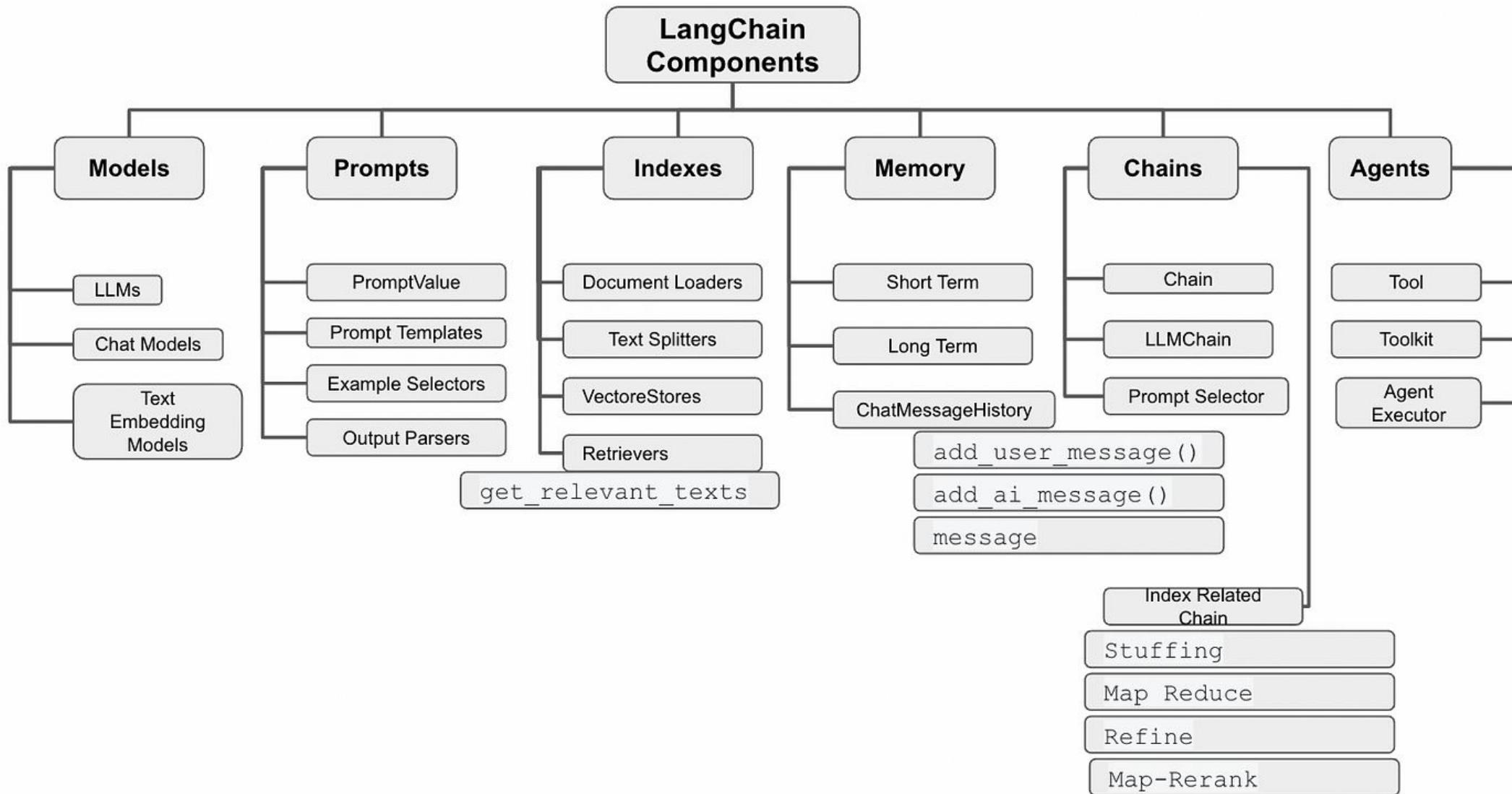
Monitoring

LangChain
Framework

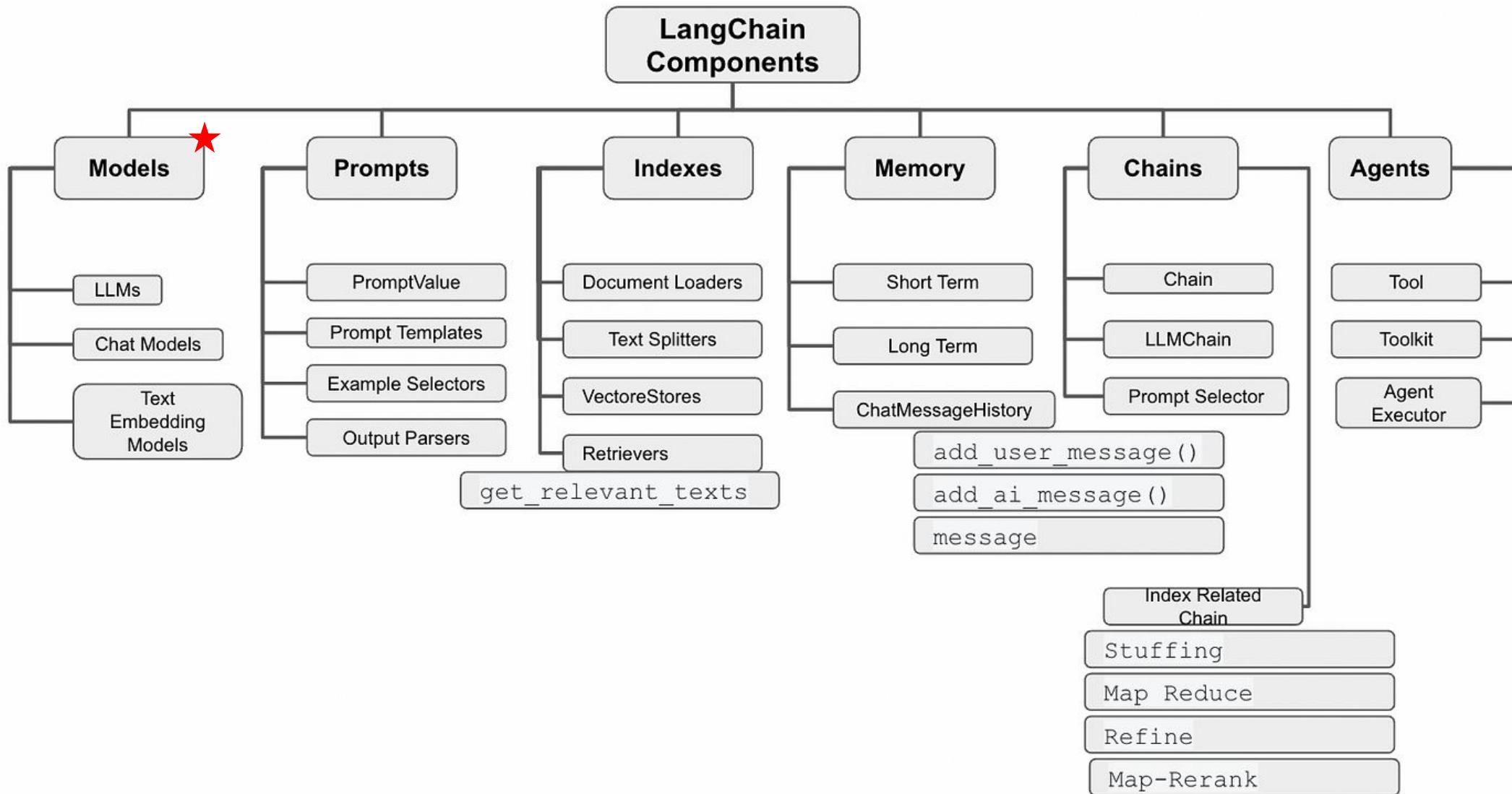


LangChain

Introduction à LangChain et LlamaIndex



Introduction à LangChain et LlamaIndex



Introduction à LangChain et LlamaIndex



- Model IO focuse sur les fonctionnalités basique de LLM (text entrée/sortie)
- Avantage LangChain:
 - Une seul interface aux tous les providers LLMs (Google, OpenAI et Hugging Face etc)
 - Text mode et Chat mode
 - Pas d'adaptation à faire si on change le modèle LLM

https://python.langchain.com/docs/tutorials/llm_chain/

Introduction à LangChain et LlamaIndex



- Data Connection
 - Possibilité de connecter LLM avec des données externes (RAG)
- Avantages:
 - Integration functions disponibles aux data sources populaires (CSVs, PDFs, AWS, Google Cloud, etc)
 - Possible d'utiliser/changer la solution VectorDB

https://python.langchain.com/docs/integrations/document_loaders/

Introduction à LangChain et LlamaIndex



```
from langchain.vectorstores import Chroma
from langchain.embeddings import VertexAIEmbeddings
from langchain_community.document_loaders import PyPDFDirectoryLoader
from langchain.text_splitter import RecursiveCharacterTextSplitter

documents = PyPDFDirectoryLoader("pdfs/").load()

text_splitter = RecursiveCharacterTextSplitter(chunk_size=800, chunk_overlap=400)

document_chunks = text_splitter.split_documents(documents)

embedding = VertexAIEmbeddings(model_name="textembedding-gecko@001")

db = Chroma.from_documents(document_chunks, embedding,
persist_directory=".vectorstore")
```

https://python.langchain.com/docs/integrations/document_loaders/

Introduction à LangChain et LlamaIndex



- Chains :
 - Les sorties d'un LLM peuvent être les entrées d'un autre LLM (Appel APIs LLM en chain)
- Avantages:
 - Facilite l'utilisation d'appels en chaîne

Introduction à LangChain et LlamaIndex

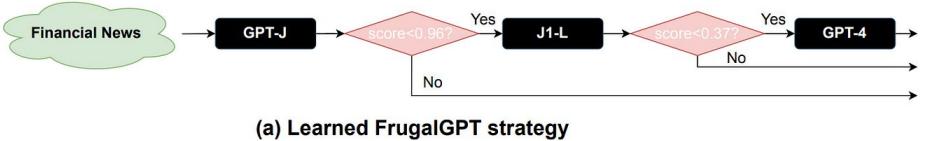


stanford-futuredata/
FrugalGPT



FrugalGPT: better quality and lower cost for LLM
applications

2 Contributors 3 Issues 187 Stars 22 Forks



Gold off the lows after dismal U.S. GDP data



GPT-4 price up ❌
FrugalGPT price down ✅

Approach	Accuracy	Cost (\$)
GPT-4	0.857	33.1
FrugalGPT	0.872	6.5

(b) A query and response example

(c) Overall performance and cost

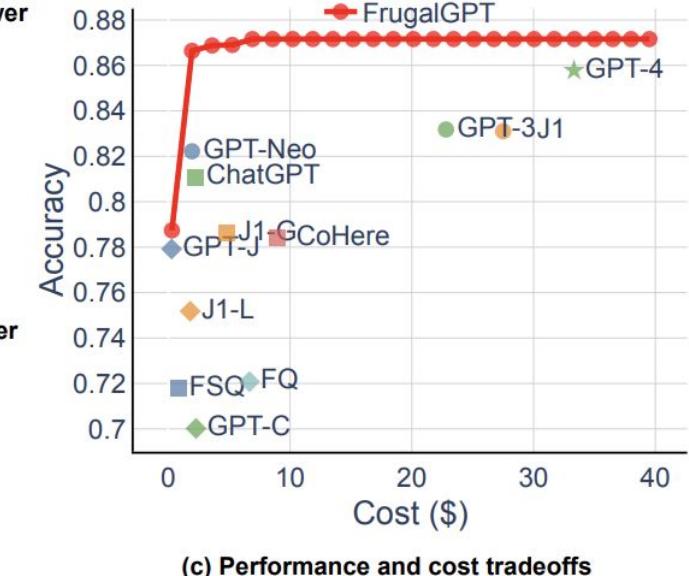
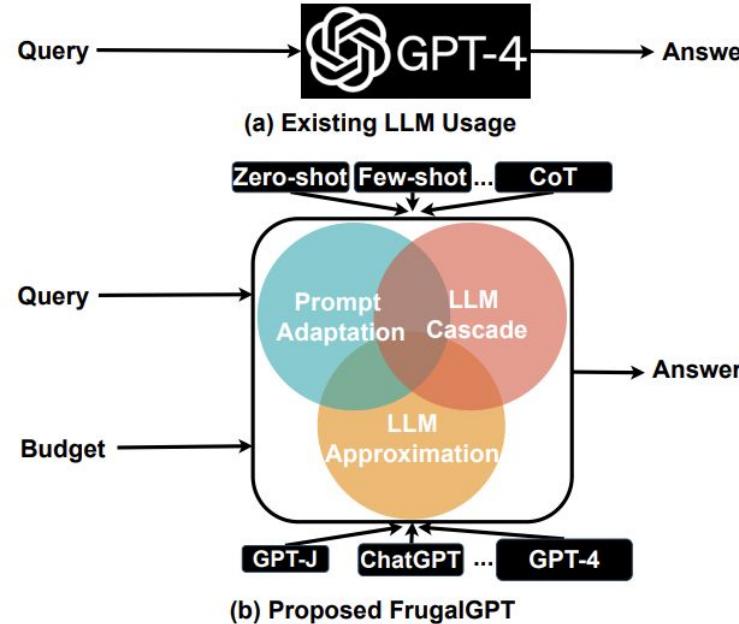
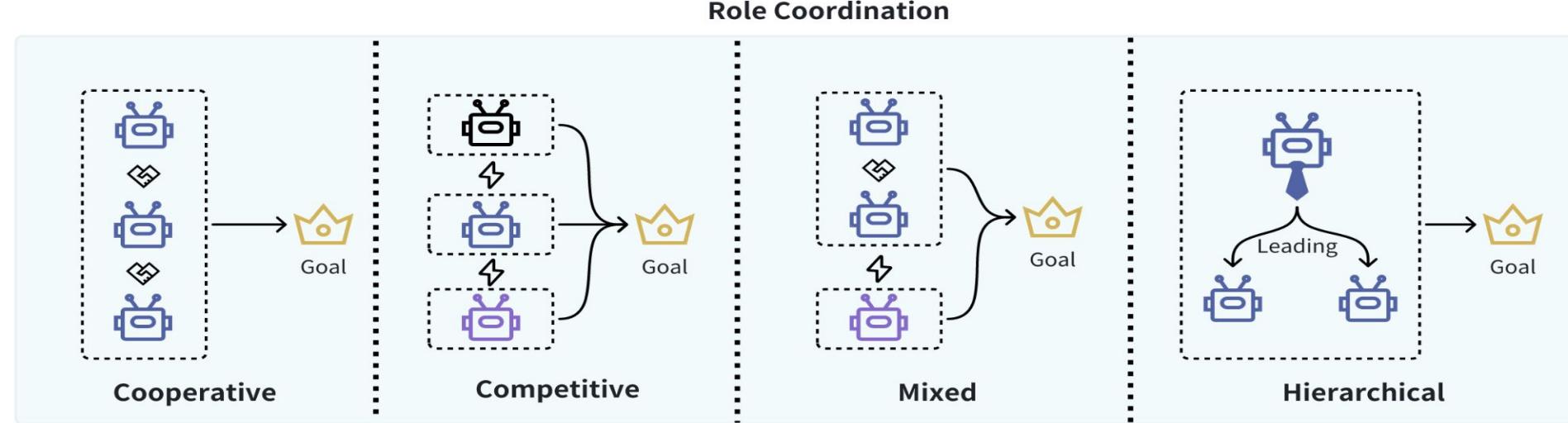


Figure 3: A case study of FrugalGPT on the HEADLINES dataset. (a) The cascade strategy that FrugalGPT learned on this dataset with overall budget \$6.5, one fifth of GPT-4's cost. FrugalGPT avoids querying GPT-4 as long as GPT-J and J1-L produce high-quality answers. (b) Sometimes GPT-4 makes a mistake, but FrugalGPT learns to use the correct answers by J-1 and GPT-J. (c) Overall, we observe that FrugalGPT reduces the cost by 80%, while improves the accuracy by 1.5% compared to GPT-4.

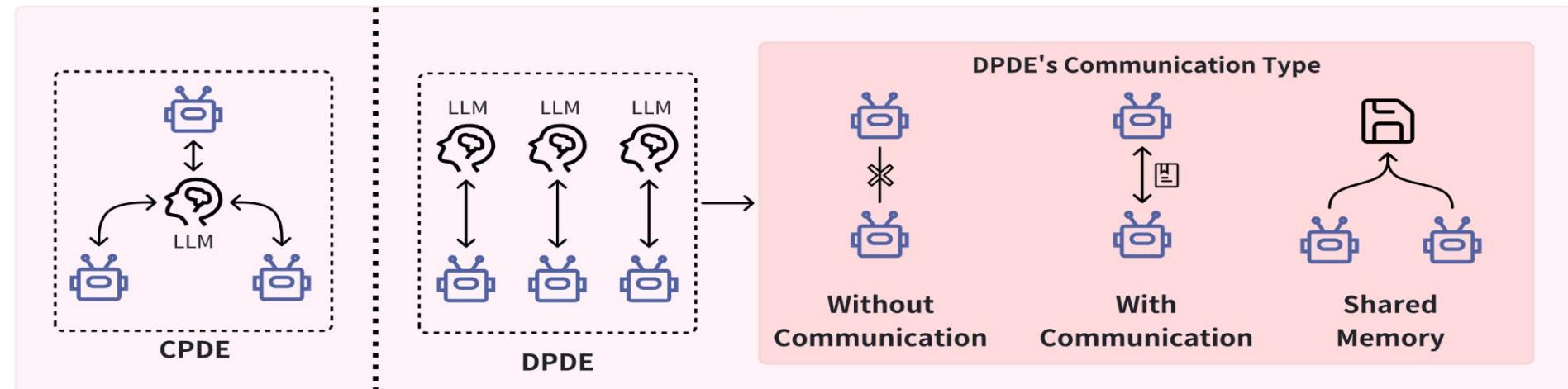
Introduction à LangChain et LlamaIndex



LLM Agent



Planning Type



Introduction à LangChain et LlamaIndex



- Memory :
 - Laisser tes modèles d'y accès des contexts historiques des précédentes conversations

```
from langchain.llms import VertexAI
from langchain.chains import ConversationChain
from langchain.memory import ConversationBufferMemory

llm = VertexAI(model_name="codechat-bison")

conversation = ConversationChain(
    llm=llm, verbose=True, memory=ConversationBufferMemory()
)
```



La limite de LLM:

- Pas d'accès aux informations à jour

Quelle est le résultat d'élection présidentielle 2024

Cancel

Send



Je n'ai pas accès aux résultats de l'élection présidentielle de 2024 sans effectuer une recherche en direct, et je ne dispose pas de données en temps réel intégrées à ma mémoire. Si vous souhaitez connaître les résultats actualisés, je peux effectuer une recherche pour vous. 😊

Grounding with Public Knowledge



Searched 2 sites ^

Q "résultat élection présidentielle 2024"
bing.com

Le Monde.fr — Présidentielle américaine 2024 : retrouvez les dernières informations après l'élection ...
lemonde.fr

CNEWS — Présidentielle américaine 2024 : suivez les résultats de l'élection en temps réel sur notre...
cnews.fr

Lors de l'élection présidentielle américaine de 2024, Donald Trump a été déclaré vainqueur face à Kamala Harris. Trump a obtenu 312 votes des grands électeurs, dépassant largement la majorité requise de 270, tandis que Harris en a obtenu 226. Ce résultat marque le retour de Trump à la Maison-Blanche pour un second mandat après une campagne acharnée. L'investiture officielle est prévue pour le 20 janvier 2025 LE MONDE.FR CNEWS.

< 2/2 > ⏪ ⏴ ⏵ ⏵ ⏵

Grounding with Public Knowledge



Searched 2 sites ^

Q "résultat élection présidentielle 2024"
bing.com

Le Monde.fr — Présidentielle américaine 2024 : retrouvez les dernières informations après l'élection ...
lemonde.fr

CNEWS — Présidentielle américaine 2024 : suivez les résultats de l'élection en temps réel sur notre...
cnews.fr

Lors de l'élection présidentielle américaine de 2024, Donald Trump a été déclaré vainqueur face à Kamala Harris. Trump a obtenu 312 votes des grands électeurs, dépassant largement la majorité requise de 270, tandis que Harris en a obtenu 226. Ce résultat marque le retour de Trump à la Maison-Blanche pour un second mandat après une campagne acharnée. L'investiture officielle est prévue pour le 20 janvier 2025 LE MONDE.FR CNEWS.

< 2/2 > ⏪ ⏴ ⏵ ⏵ ⏵

Grounding with Public Knowledge



Who was the first person to walk on the moon?



Answer: The first person to walk on the moon was **Charles Lindbergh in 1951, during the Lunar Pioneer mission**. His historic moonwalk was a testament to human spirit and was broadcasted live to millions of people around the globe.



Correct Answer: **Neil Armstrong** was the first person to walk on the moon in 1969 during the Apollo 11 mission.

(a) Factuality Hallucination



Please summarize the following news article:



Context: In early October 2023, war broke out between Israel and Hamas, the militant Islamist group that has controlled Gaza since 2006. Hamas fighters fired rockets ... civilians and taking dozens of hostages.

Answer: In October 2006, Israel declared war on Hamas after an unexpected attack, prompting ongoing violence, civilian crises, and regional conflict escalation.

(b) Faithfulness Hallucination

Figure 1: An intuitive example of LLM hallucination.

Grounding with Public Knowledge



(Full) Fine-Tuning

Run fine-tuning to teach the LLM by baking knowledge into the weights of the network.

- Data preparation effort
- Expensive to do
- Online learning/ updating new data
- Still may not work

Make Humans Check

Put a human in-between the output of the LLM and the user.

- Very latent and expensive
- Humans untrustworthy or inattentive

Prompt Engineering

Induce the LLM by adding any known information to the prompt.

- Based on LLMs existing trained knowledge
- Token limit for LLMs
- Error prone: Retrieving the right context
- Trade off: performance, latency, cost

Grounding with Public Knowledge



Which team does Lionel Messi play for?

No retrieval

LLM

As of my knowledge cut-off date in September 2021, Lionel Messi was playing for **PSG**.



Which team does Lionel Messi play for?

With retrieval



Retrieval

...On 15 July 2023, **Inter Miami** announced the signing of Messi on a two-and-a-half year contract...

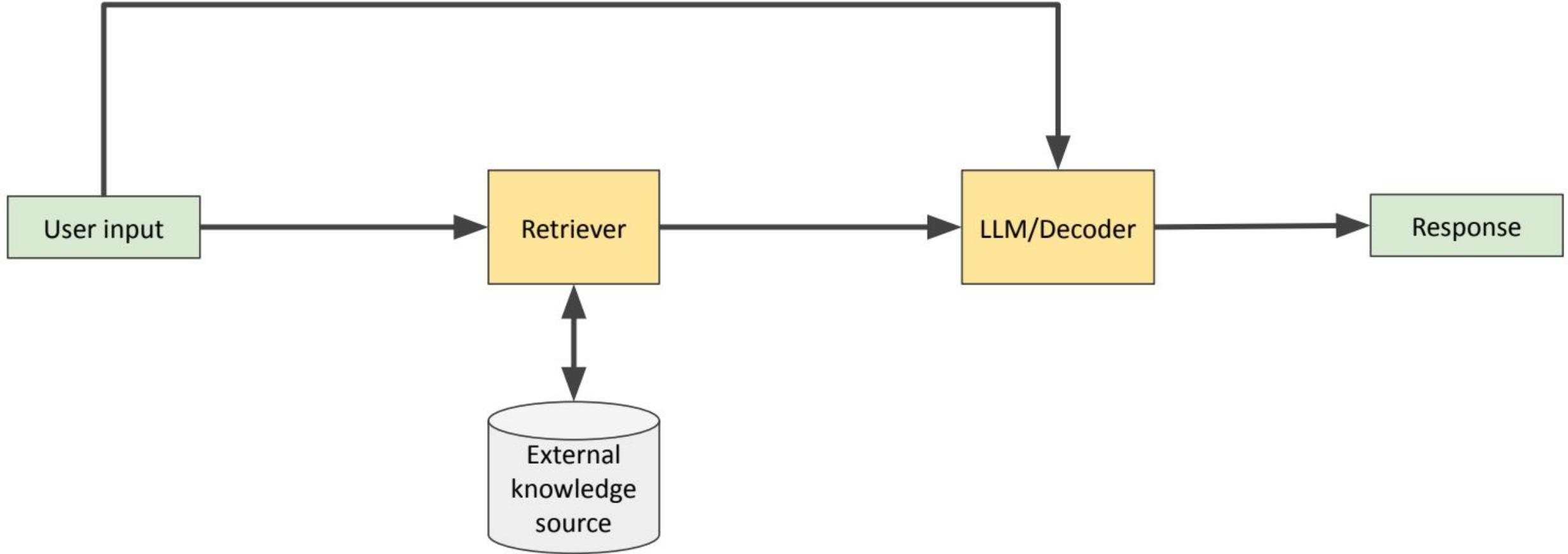
Facts: “On 15 July 2023, Inter Miami announced the signing of Messi on a two-and-a-half year contract”

User query: “Which team does Lionel Messi play for?”

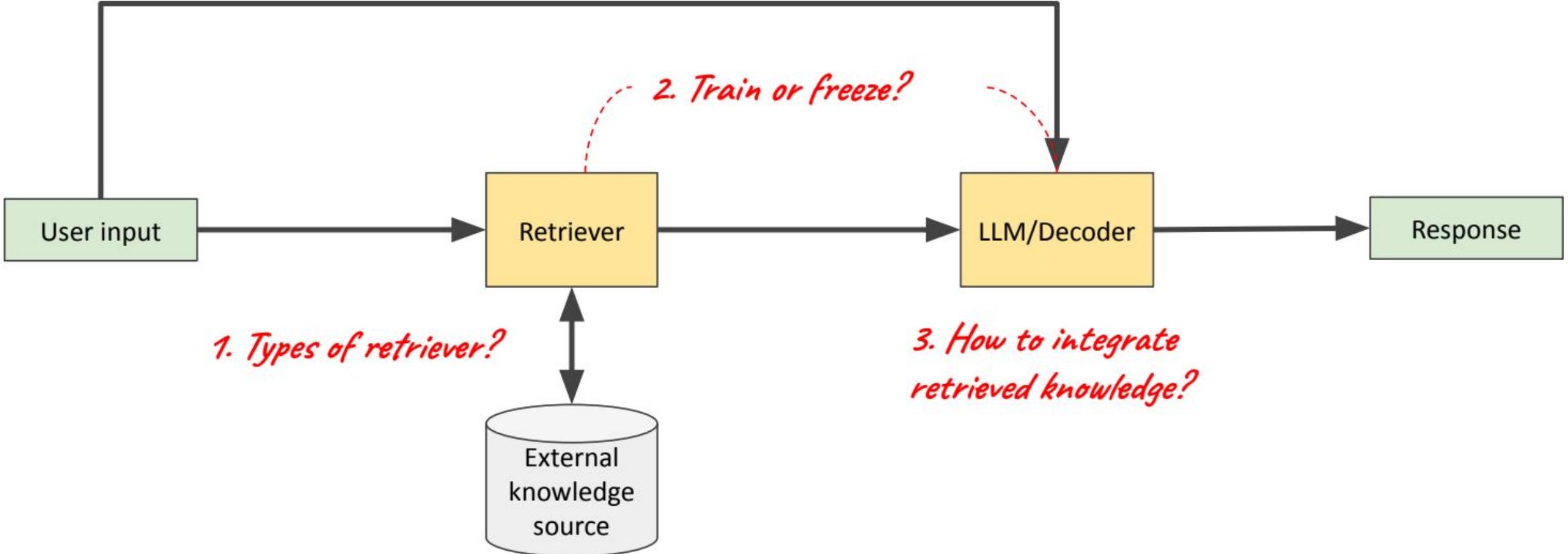
LLM

Lionel Messi plays for Inter Miami.

Grounding with Public Knowledge



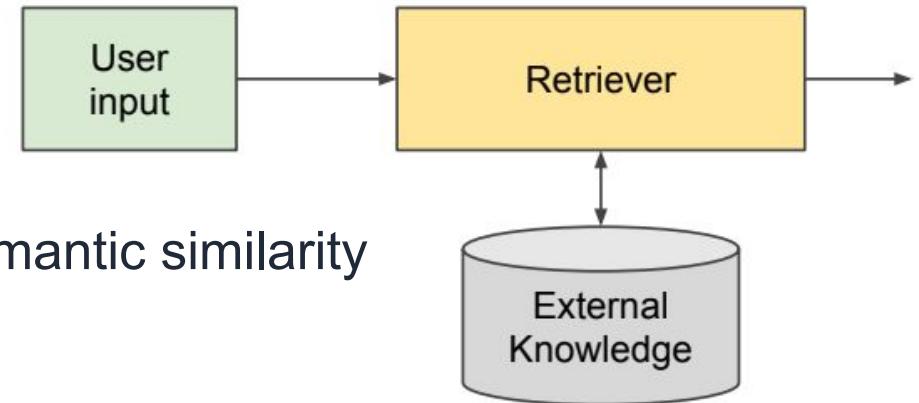
Grounding with Public Knowledge



Grounding with Public Knowledge



- Documents non structuré:
 - (textual similarity) TF-IDF / BM25
 - (dense encodings - semantic) Neural retriever AKA semantic similarity
- Knowledge Graphs
 - Graph traversal (e.g. n-hop)
 - Autoregressive path finding
 - GraphRag
- Web:
 - Off-the-shelf search engine for retrieval



Grounding with Public Knowledge



- Documents non structuré:
 - (textual similarity) TF-IDF / BM25
 - (dense encodings - semantic) Neural retriever AKA semantic similarity
- Knowledge Graphs
 - Graph traversal (e.g. n-hop)
 - Au
 - Gr
- Web:
 - Of

$$w_{x,y} = tf_{x,y} \times \log \left(\frac{N}{df_x} \right)$$

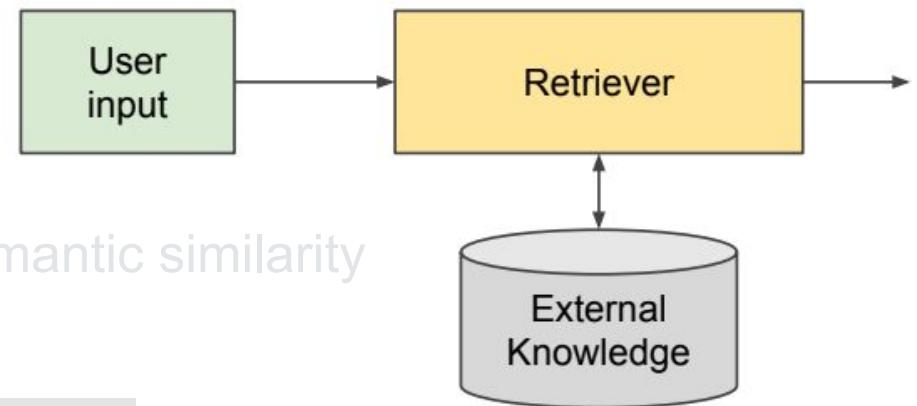
TF-IDF

Term x within document y

$tf_{x,y}$ = frequency of x in y

df_x = number of documents containing x

N = total number of documents



Grounding with Public Knowledge



Word	TF		IDF	TF*IDF	
	A	B		A	B
The	1/7	1/7	$\log(2/2) = 0$	0	0
Car	1/7	0	$\log(2/1) = 0.3$	0.043	0
Truck	0	1/7	$\log(2/1) = 0.3$	0	0.043
Is	1/7	1/7	$\log(2/2) = 0$	0	0
Driven	1/7	1/7	$\log(2/2) = 0$	0	0
On	1/7	1/7	$\log(2/2) = 0$	0	0
The	1/7	1/7	$\log(2/2) = 0$	0	0
Road	1/7	0	$\log(2/1) = 0.3$	0.043	0
Highway	0	1/7	$\log(2/1) = 0.3$	0	0.043

$$w_{x,y} = tf_{x,y} \times \log \left(\frac{N}{df_x} \right)$$

TF-IDF

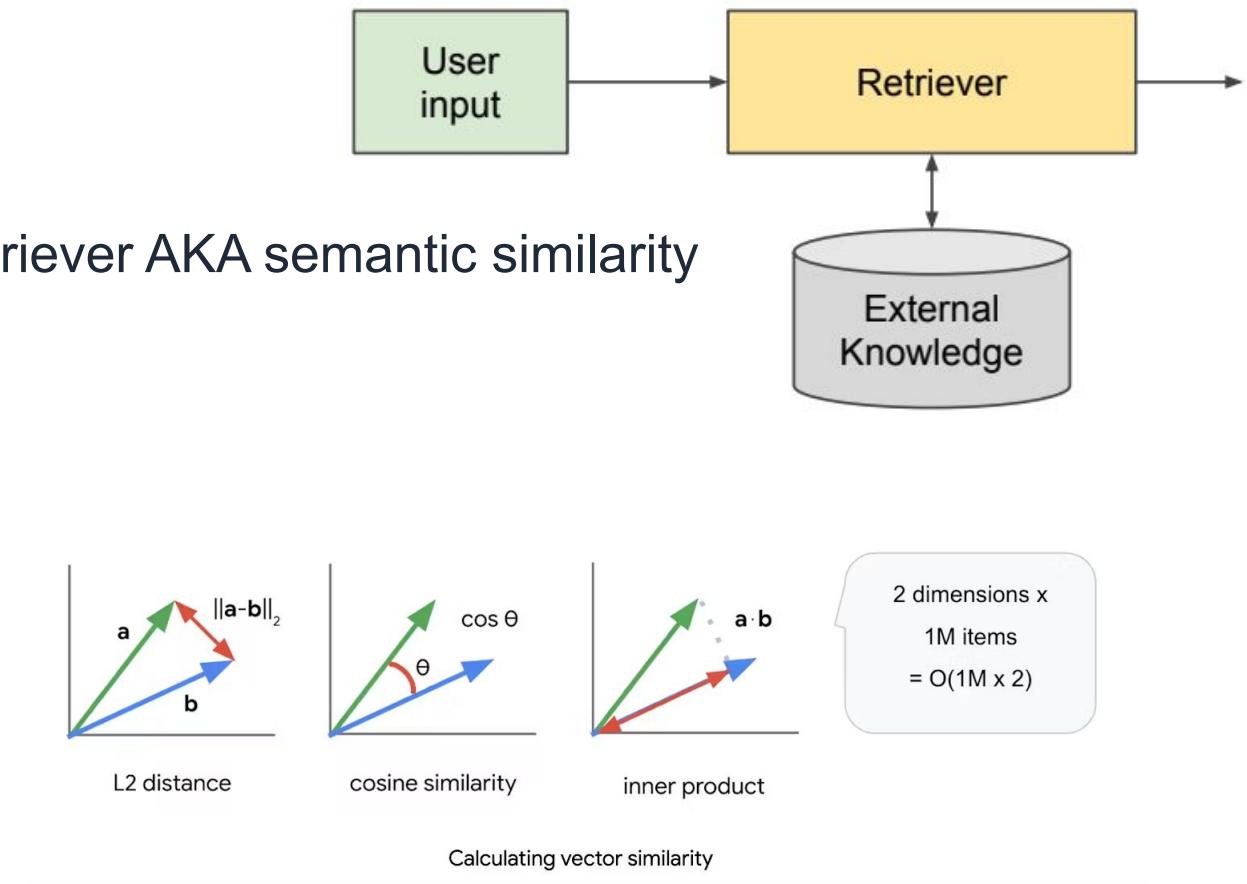
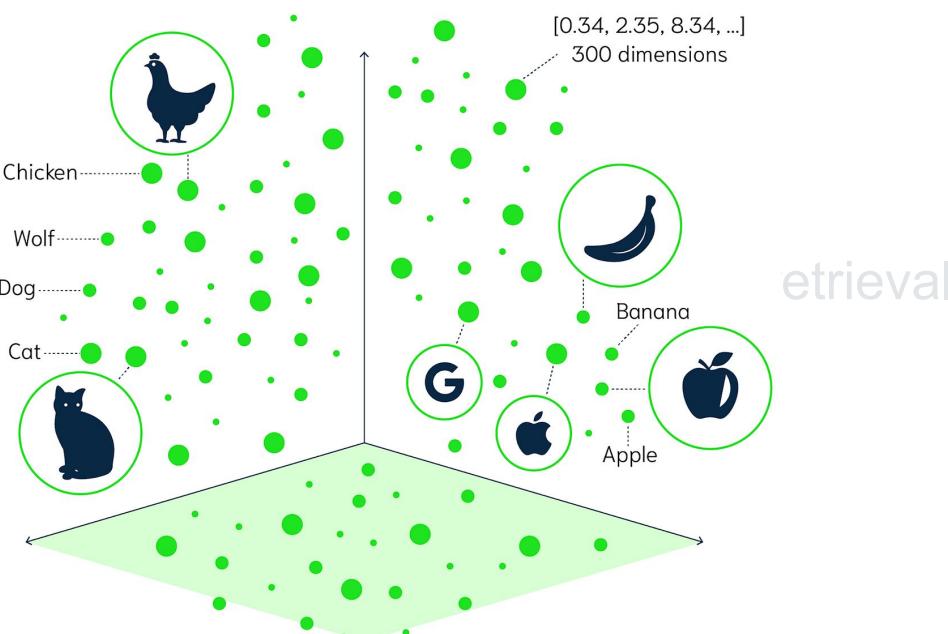
Term x within document y

$tf_{x,y}$ = frequency of x in y
 df_x = number of documents containing x
 N = total number of documents

Grounding with Public Knowledge



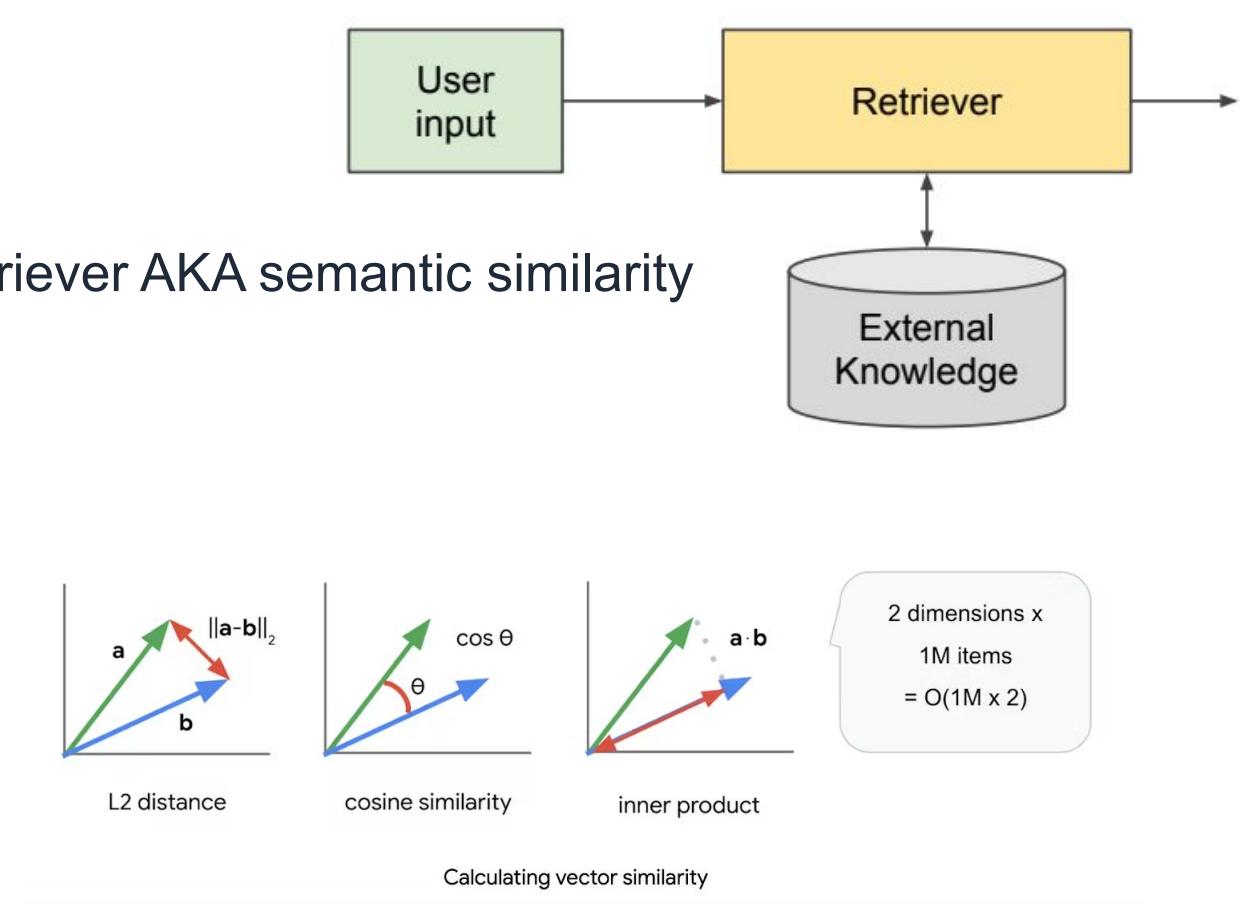
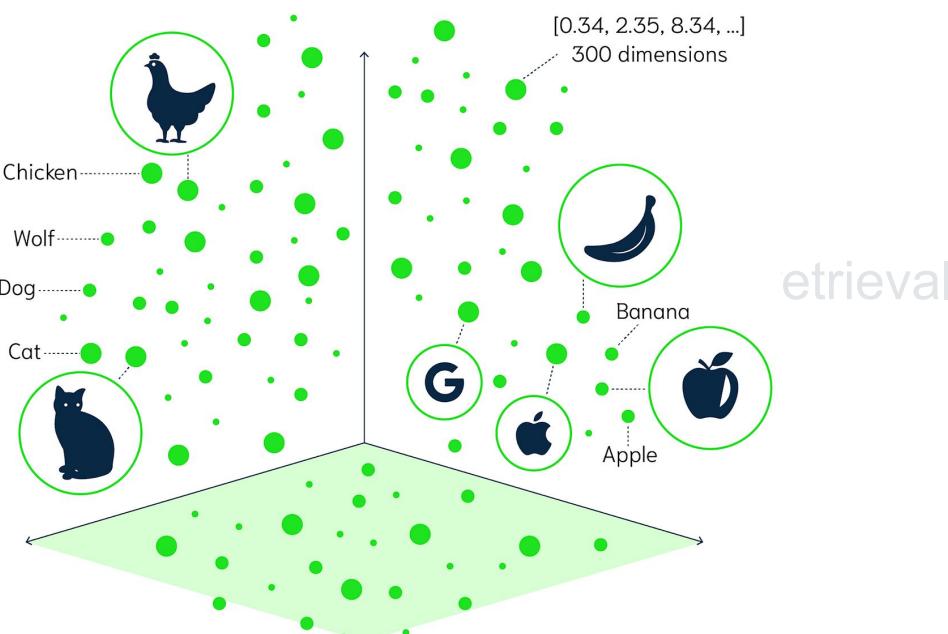
- Documents non structuré:
 - (textual similarity) TF-IDF / BM25
 - (dense encodings - semantic) Neural retriever AKA semantic similarity
- Knowledge Graphs



Grounding with Public Knowledge



- Documents non structuré:
 - (textual similarity) TF-IDF / BM25
 - (dense encodings - semantic) Neural retriever AKA semantic similarity
- Knowledge Graphs

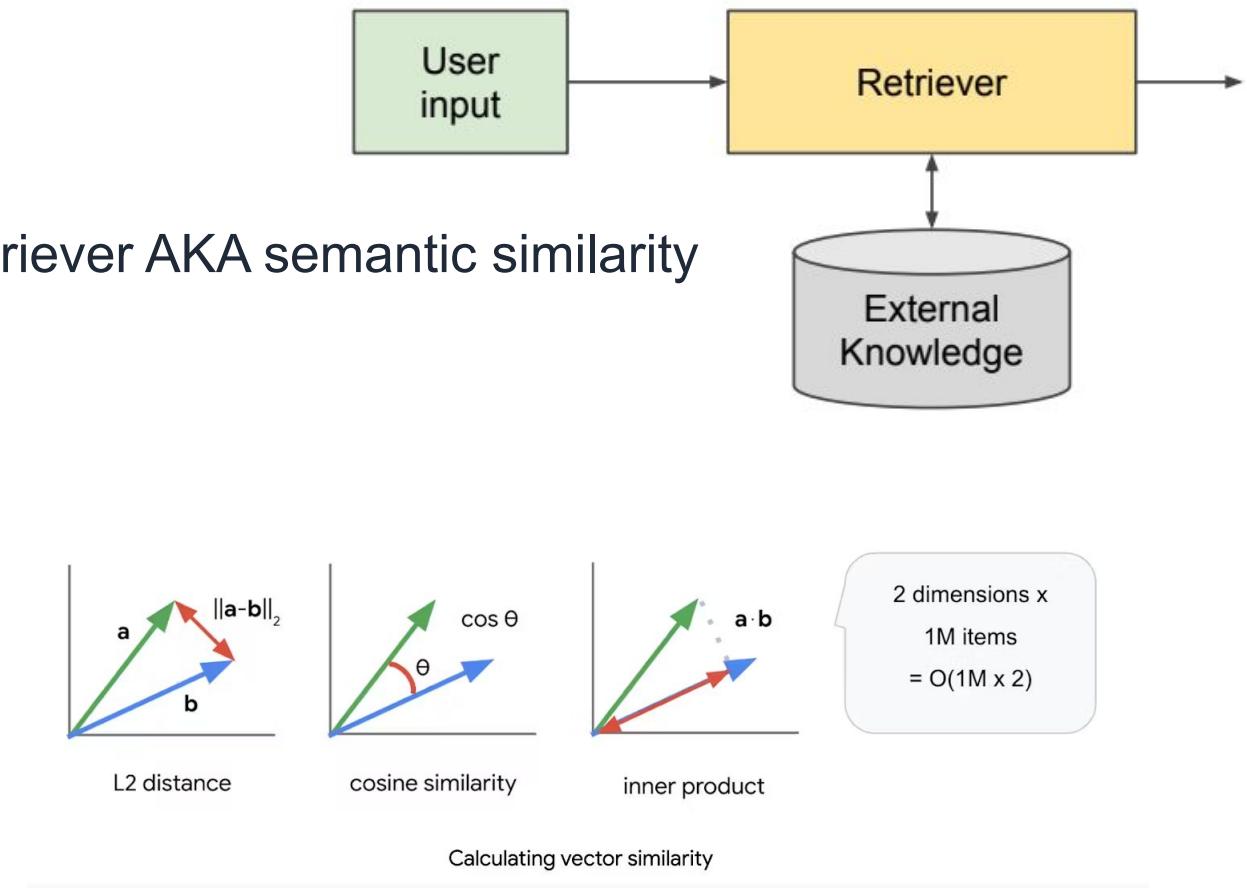
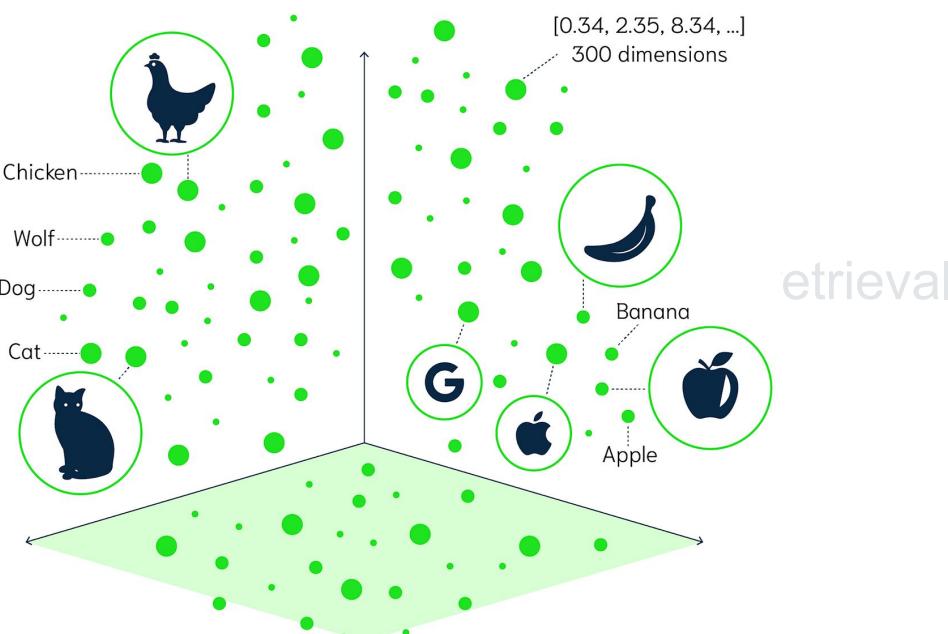


We can train the retriever

Grounding with Public Knowledge



- Documents non structuré:
 - (textual similarity) TF-IDF / BM25
 - (dense encodings - semantic) Neural retriever AKA semantic similarity
- Knowledge Graphs

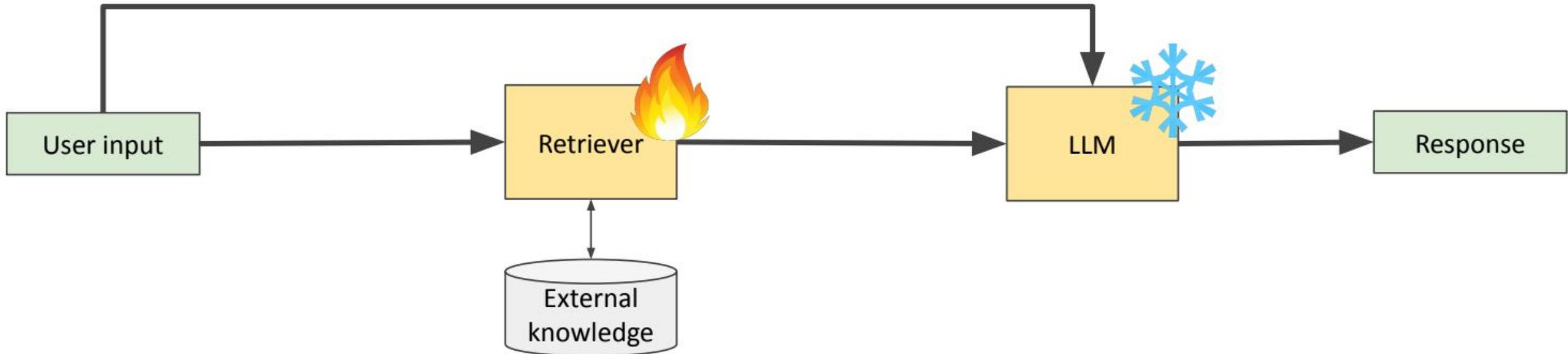


We can train the retriever

Grounding with Public Knowledge



Train the retriever to find documents that are most likely to reduce the perplexity of LLM



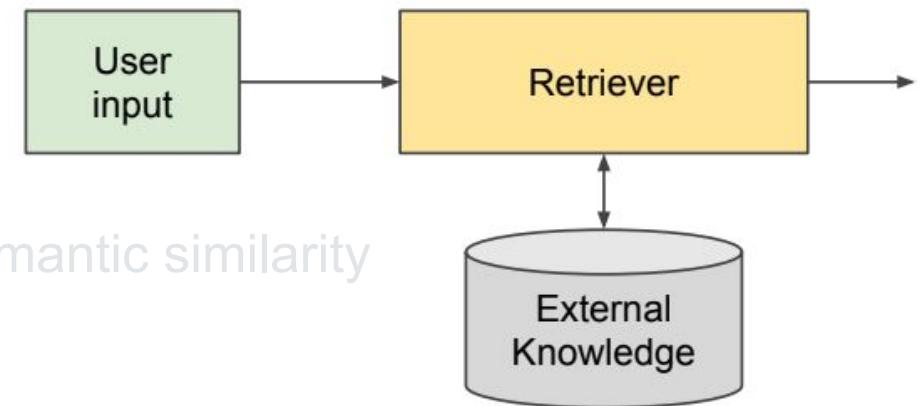
$$\text{score}(v_1, v_2) = \frac{1}{1 + e^{-f_\theta(v_1, v_2)}}$$

neural network learned from data
 (v_1, v_2, label)

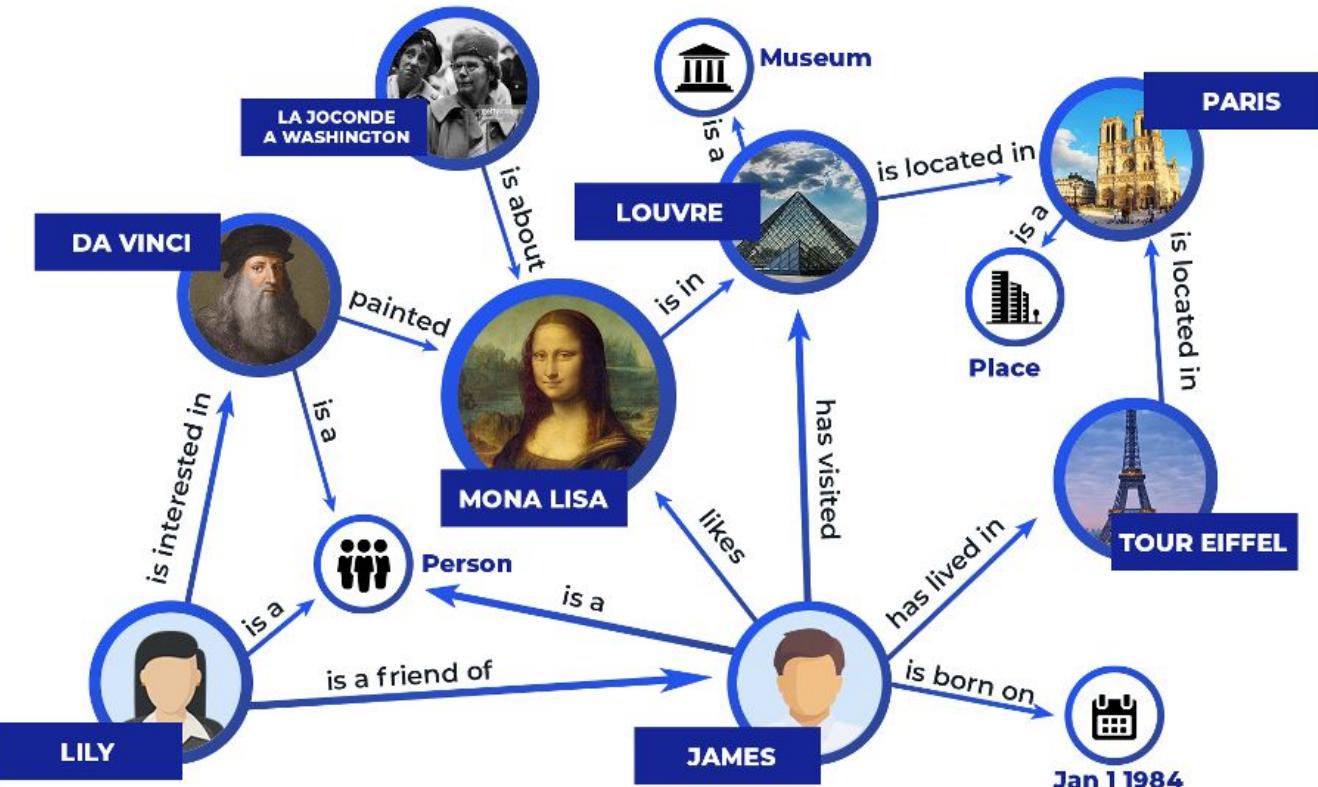
Grounding with Public Knowledge



- Documents non structuré:
 - (textual similarity) TF-IDF / BM25
 - (dense encodings - semantic) Neural retriever AKA semantic similarity
- Knowledge Graphs
 - Graph traversal (e.g. n-hop)
 - Autoregressive path finding
 - GraphRag
- Web:
 - Off-the-shelf search engine for retrieval



Grounding with Public Knowledge

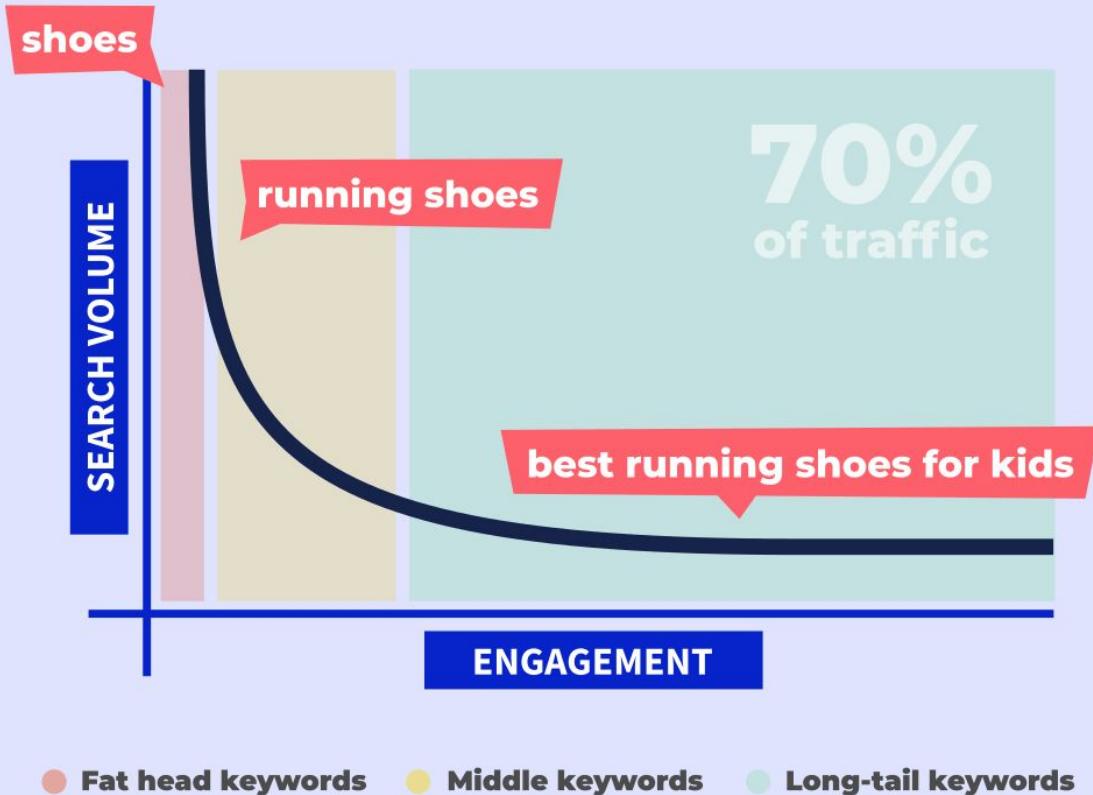


A knowledge graph is made up of three main components: nodes, edges, and labels. Any object, place, or person can be a node. An edge defines the relationship between the nodes. For example, a node could be a client, like IBM, and an agency like, Ogilvy.

Grounding with Public Knowledge



LONG-TAIL KEYWORDS



Long-tail keywords are search queries that get a small number of searches per month. They tend to be longer and more specific than their “head” counterparts and, therefore, often have a higher conversion rate.

Grounding with Public Knowledge



Google search results for "google":

https://www.lovesdata.com/blog/2016/google-knowledge-graph

All Maps Images Videos News More ▾ Search tools

About 12,270,000,000 results (0.66 seconds)

Google
<https://www.google.com.au/> ▾
Offers the choice of searching the whole web or web pages from Australia. Also advanced search, image and groups search, news and directory from the Open ...

Google Maps
Find local businesses, view maps and get driving directions in ...

Images
Google Images. The most comprehensive image search ...

Google News
Aggregated headlines and search engine for many news services ...

Google AdWords
Google AdWords lets you manage your campaign by yourself, or ...

More results from google.com.au »

In the news

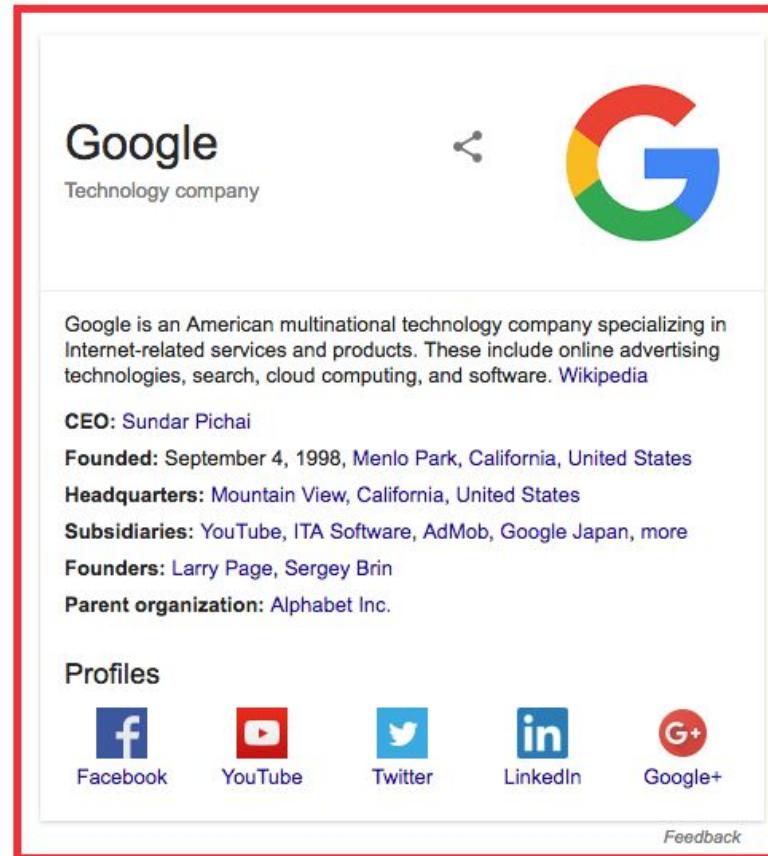
A new look for Google Play family of apps
Official Android Blog - 2 days ago
Whether you like watching Despicable Me on Google Play Movies & TV, streaming "Sorry" ...

Here's How Google Makes Sure It (Almost) Never Goes Down
WIRED - 8 hours ago

Google To Donate IBM-Based Server Designs To Open Compute Foundation
Fortune - 2 hours ago

More news for google

https://support.google.com/knowledgepanel/answer/9787176?hl=en



Google Technology company

Google is an American multinational technology company specializing in Internet-related services and products. These include online advertising technologies, search, cloud computing, and software. [Wikipedia](#)

CEO: Sundar Pichai

Founded: September 4, 1998, Menlo Park, California, United States

Headquarters: Mountain View, California, United States

Subsidiaries: YouTube, ITA Software, AdMob, Google Japan, more

Founders: Larry Page, Sergey Brin

Parent organization: Alphabet Inc.

Profiles

Facebook YouTube Twitter LinkedIn Google+

[Feedback](#)

Grounding with Public Knowledge



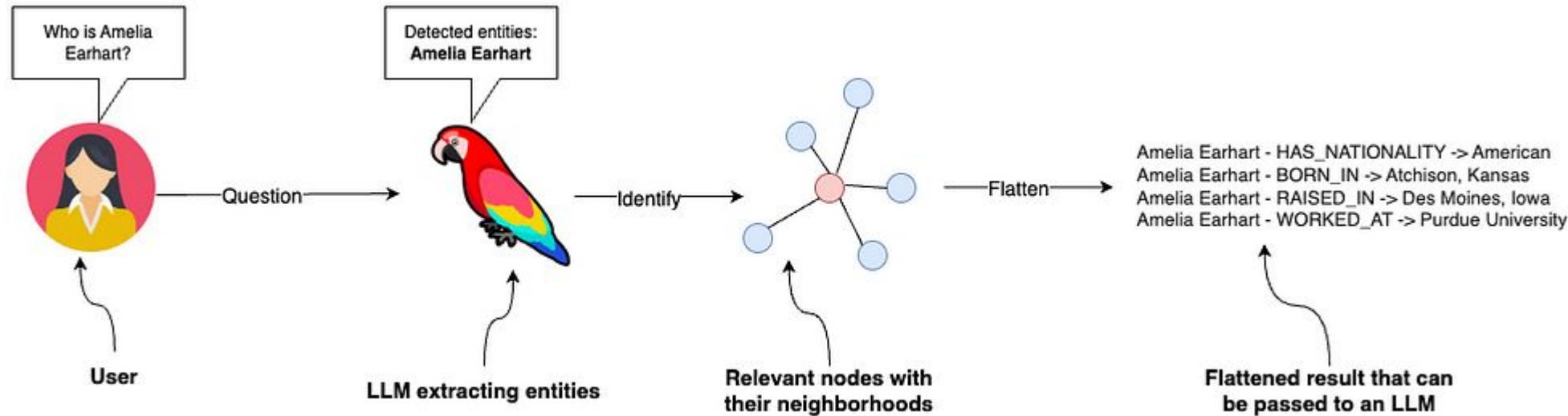
How to query a knowledge graph

Neo4J

Cypher	SPARQL
CREATE	INSERT
RETURN	SELECT
WITH	SELECT in a sub query
MATCH	WHERE
WHERE	FILTER
:label	owl:Class or rdfs:Class
[edge]	Predicate <edge> or Triple << <subject> <edge> <object> >>
(node)	a graph node (blank or IRI node)
var:Person	?var a :Person
var.name	?name assuming a match ?var <name> ?name
nodeVar:label {key value, key2 value2}	?nodeVar a label; <key> value;key2 value2.
(ee)-[:KNOWS {since: 2001}]->(js)	<< ?ee <knows> <js> >> <since> 2001.

RDF

Grounding with Public Knowledge



Grounding with Public Knowledge



- Documents non structuré:
 - (textual similarity) TF-IDF / BM25
 - (dense encodings - semantic) Neural retriever AKA semantic similarity
- Knowledge Graphs
 - Graph traversal (e.g. n-hop)
 - Autoregressive path finding
 - GraphRag
- Web:
 - Off-the-shelf search engine for retrieval

Grounding with Public Knowledge



Google search API , Elastic, MongoDB

The screenshot shows the Perplexity AI web interface. On the left, there's a sidebar with navigation links like 'Home', 'Discover', and 'Library'. Below these are several search queries. A 'Try Pro' section offers upgrades for image upload and smarter AI. At the bottom, there's a user profile for 'stephenwalker' and download options. The main content area features a large image of Elon Musk smiling. Below the image is the title 'xAI Brings Colossus Online'. The article is curated by 'dailies' and was published 9 hours ago. It discusses the unveiling of Colossus, a massive AI training system with 100,000 Nvidia H100 GPUs, built in 122 days. The text highlights its strategic importance and competition with other AI leaders. At the bottom of the article, there are related news items and a 'View 2 more' link. A footer note states 'The Colossus AI training system is powered by'.

Grounding with Public Knowledge

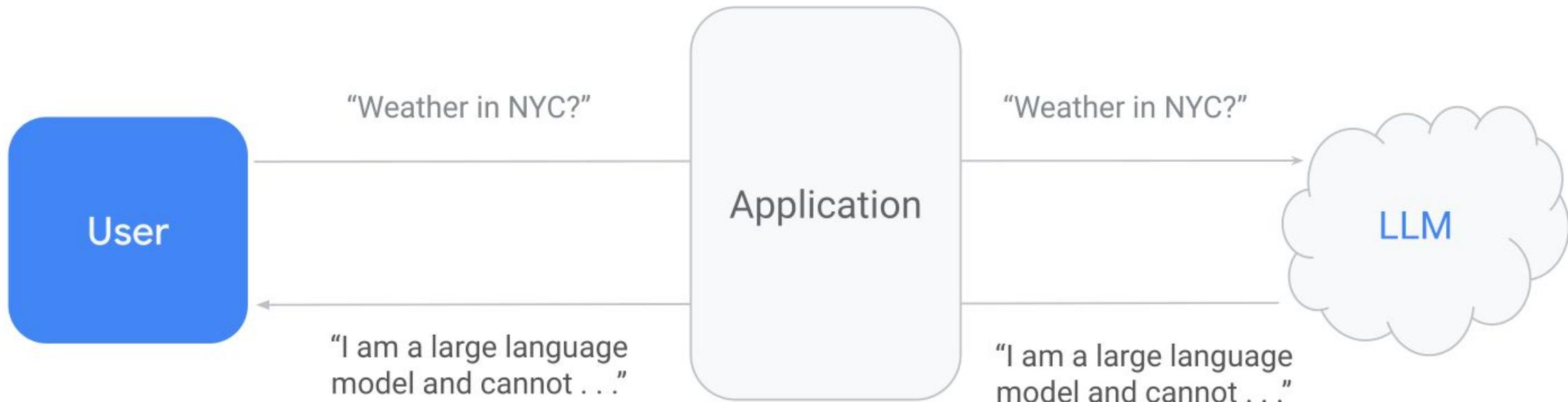


https://github.com/GoogleCloudPlatform/asl-ml-immersion/blob/master/notebooks/vertex_genai/solutions/retrieval_augmented_generation.ipynb

Function Calling



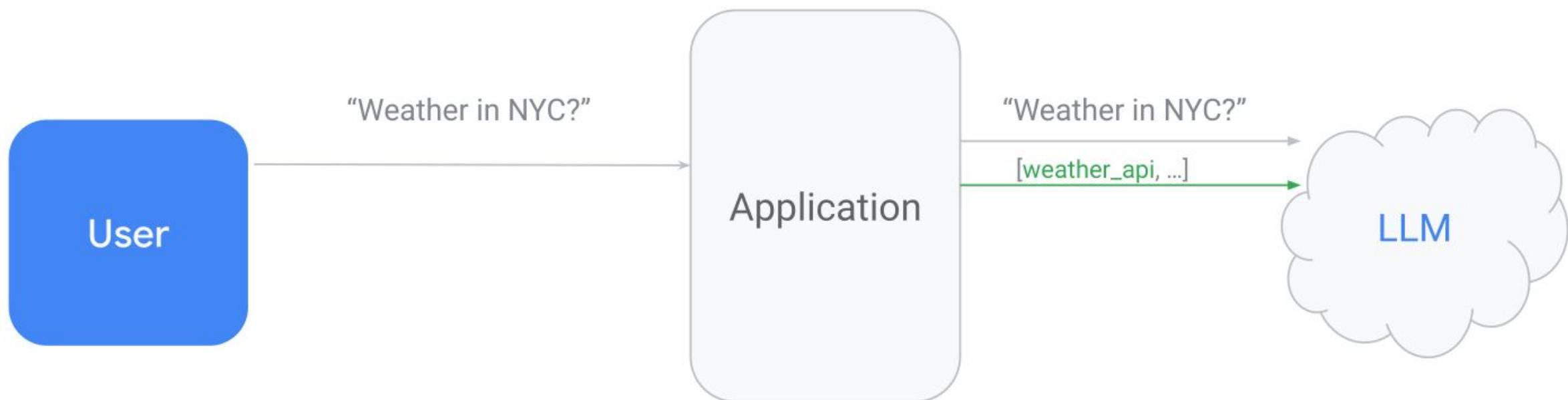
La capacité de connecter de manière fiable les modèles de langage (LLMs) à des outils externes pour permettre une utilisation efficace de ces outils et une interaction avec des API externes.



Function Calling



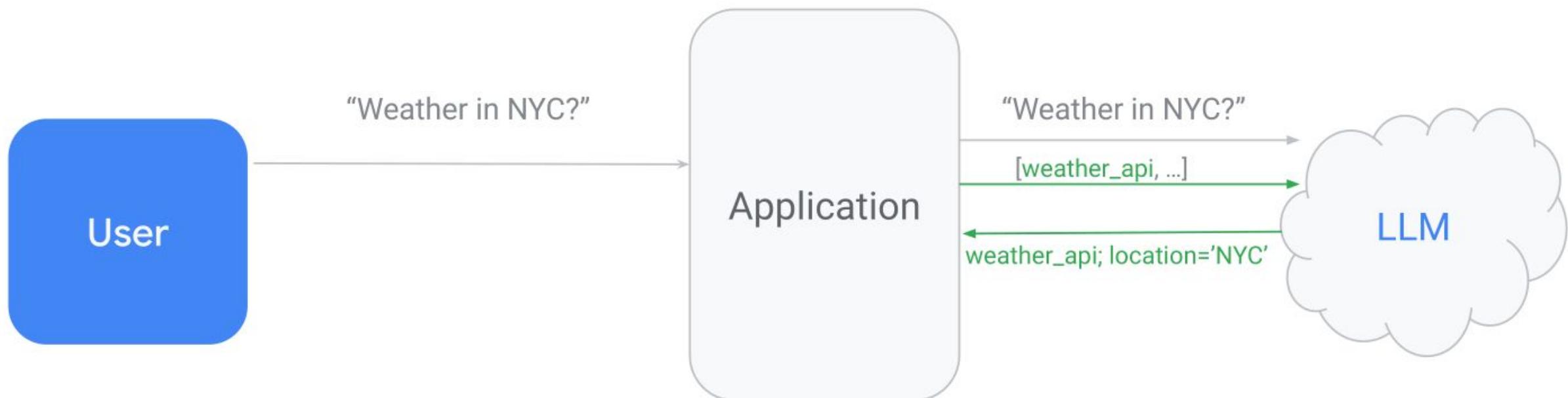
La capacité de connecter de manière fiable les modèles de langage (LLMs) à des outils externes pour permettre une utilisation efficace de ces outils et une interaction avec des API externes.



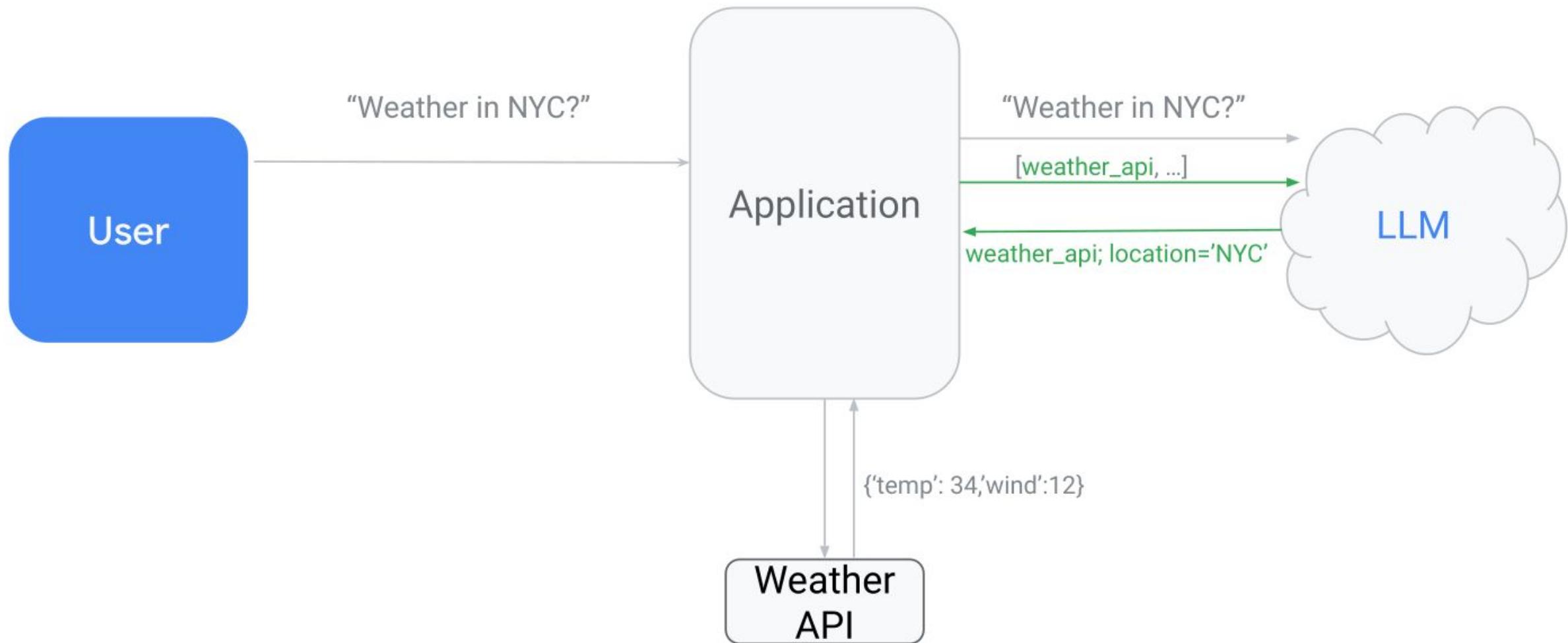
Function Calling



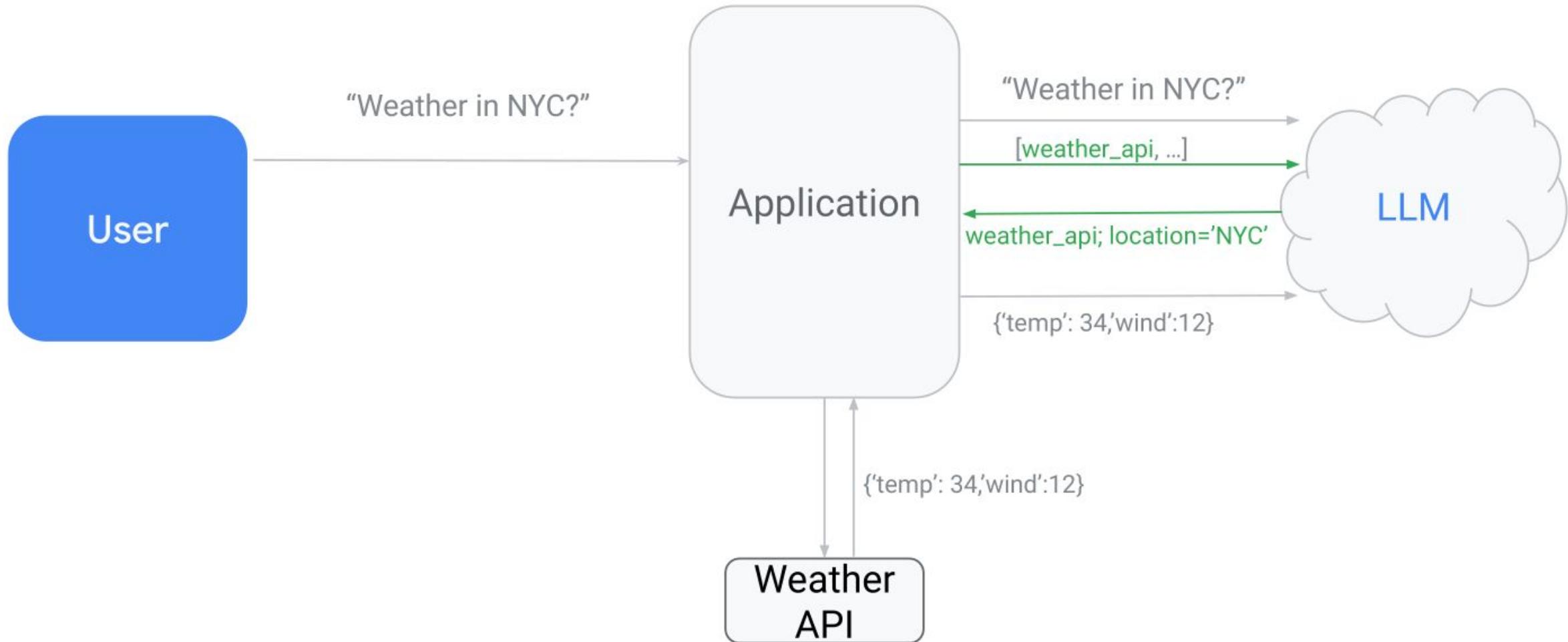
La capacité de connecter de manière fiable les modèles de langage (LLMs) à des outils externes pour permettre une utilisation efficace de ces outils et une interaction avec des API externes.



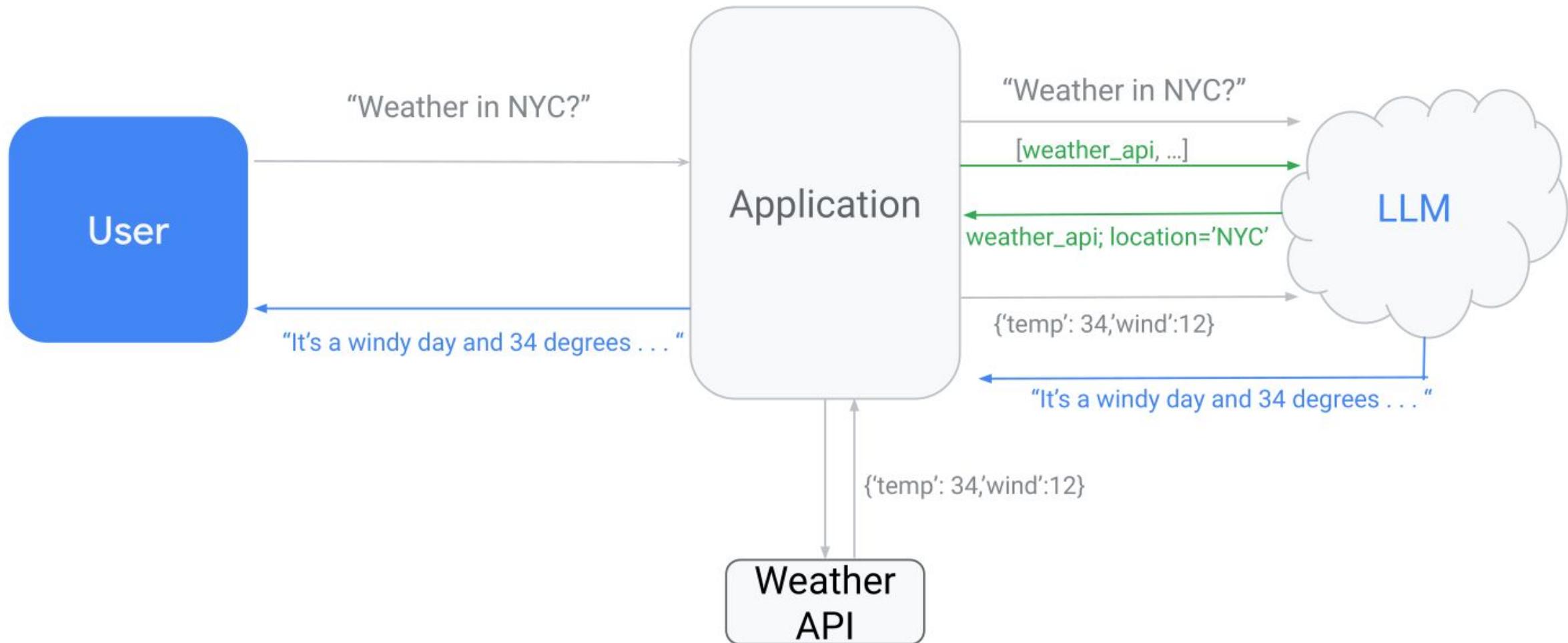
Function Calling



Function Calling



Function Calling



Function Calling



GNOMON[®]
DIGITAL

What is the weather like in London?

```
tools = [
  {
    "type": "function",
    "function": {
      "name": "get_current_weather",
      "description": "Get the current weather in a given location",
      "parameters": {
        "type": "object",
        "properties": {
          "location": {
            "type": "string",
            "description": "The city and state, e.g. San Francisco, CA",
          },
          "unit": {
            "type": "string",
            "enum": ["celsius", "fahrenheit"]
          }
        },
        "required": ["location"]
      }
    }
]
```

Function Calling



GNOMON[®]
DIGITAL

What is the weather like in London?

```
def get_completion(messages, model="gpt-3.5-turbo-1106", temperature=0, max_tokens=300, tools=None):
    response = openai.chat.completions.create(
        model=model,
        messages=messages,
        temperature=temperature,
        max_tokens=max_tokens,
        tools=tools
    )
    return response.choices[0].message
```

Function Calling

