

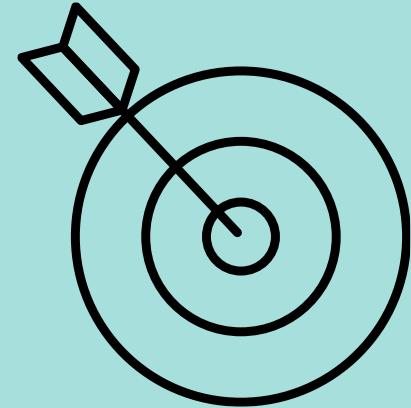
MODULE

# Agents et applications Gen AI



## Objectifs pédagogiques

- ✓ Développer des applications et des systèmes d'agents autonomes basés sur l'IA générative pour automatiser les processus de l'entreprise
- ✓ Appréhender les modèles de langage, l'optimisation des prompts, les techniques avancées de RAG, le développement de plugins, les agents autonomes et coopératifs, le fine-tunning et le merging de modèles
- ✓ Mettre en place des mesures de sécurité et d'optimisation des coûts.



# Welcome! Let's introduce ourselves



Let's start with a round of introduction.

**Please introduce yourself by sharing:**

- your name,
- your department or team, your role or position, the kind of jobs you work on
- your experience in Gen AI, Python and...
- your expectation for this training

# Detailed program

## Day 1

- Introduction, teaser, install of practicum environment (WSL2, clone project GIT, poetry, streamlit,...)
- Introduction to LangChain and LCEL (functional programming)
- Embeddings, vector stores, chunking, LLM, inference servers (How that works, how to select, how to use in LC)
- Prompting: base, advanced techniques (structured outcome, automated,...), Function calls,..

## Day 3

- Tool calling from LLM (rational, different techniques, LG implementation,...)
- Custom tools in LG
- Examples of tools (SQL, OCR, Web Automation, Code gen, Data Analytics,...)
- LangGraph
- ReAct Agents : Theory, Implementation in LC
- Cooperative Agents : Presentation of CrewAI and Autogen

## Day 2

- Naive Retrieval Augmented Generation (RAG)
- Document loader (from PDF, graph database, Web,...)
- Addressing RAG issues: Advanced retrieval and chaining techniques (Hybrid search, Self-RAG, Query Translation,...)
- LLM and Knowledge Graphs
- Memory (incl. Semantic)
- Monitoring (3rd party or custom)

## Day 4

- Strategies to improve costs, privacy, sovereignty etc.
- Introduction to model fine-tuning and merging
- Integration of third-part developed agents
- Toward enterprise agents and next
- Conclusion

# Activities, time tradeoff, ...

- Lecture sessions
- Code comprehension in Jupyter Notebooks
- Code modification in Jupyter Notebooks
- Code analysis of small applications
- Code alteration in small applications
- Demonstrations and code reviews
- Discussion about use cases and business opportunities
- Open discussions

Time allocation : Flex with your needs and dynamics...



Chapter

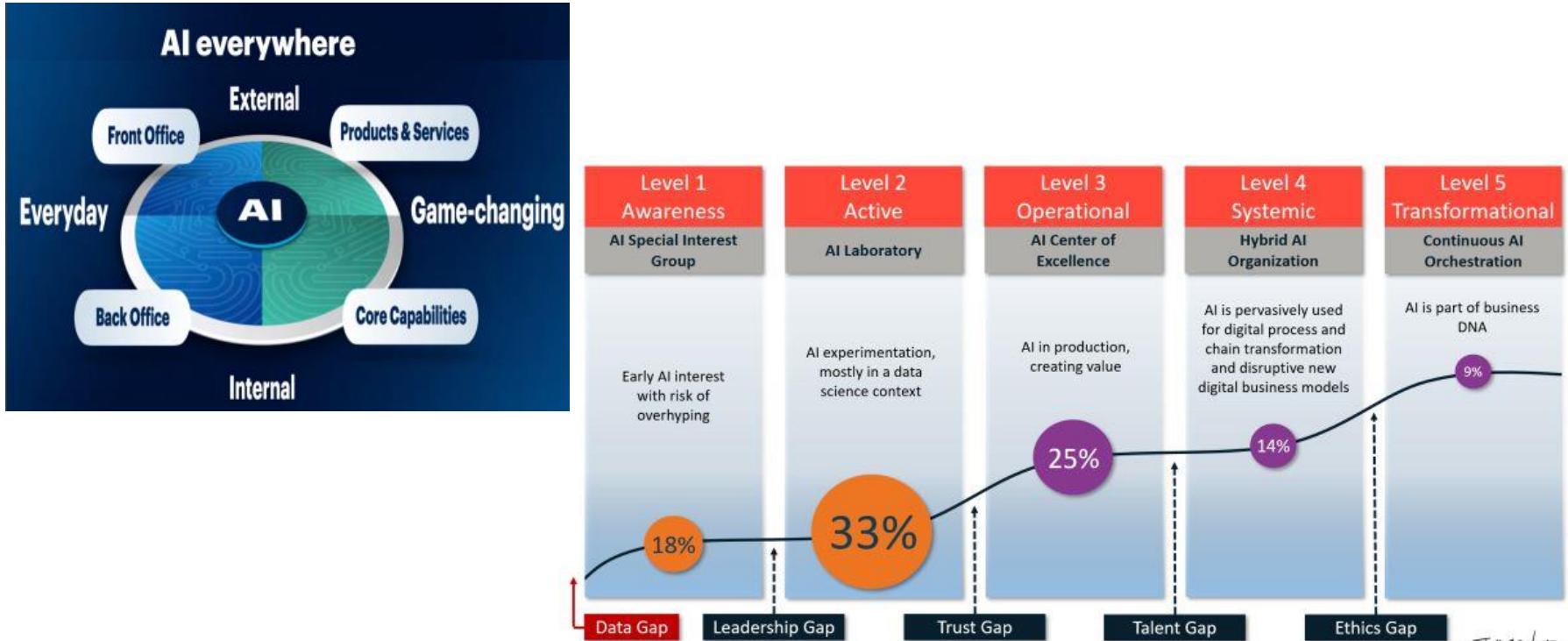
# Introduction & concepts



# General Presentation and Teaser



# AI Maturity Model (Gartner, 2023)



# Some GenAI Use Cases

## Content Generation

- Blog posts, sales pitches and marketing materials
- Illustrations and photo generation
- Audio and video

## Language Translation

- Multilingual support
- Real-time translation
- Context-aware and idiomatic translations

## Code Generation

- Autocomplete and code snippet generation
- Debugging assistance and code review
- Natural language to code conversion

## Q&A & Document Retrieval

- Information extraction from large datasets
- Context-aware and precise answers
- Cross-referencing multiple sources

## Text Summarization

- Automatic extraction of key points and meeting notes
- Condensing lengthy documents and articles
- Customizable summary lengths and styles

## Classification & Sentiment Analysis

- Text categorization and tagging
- Emotion detection and sentiment scoring
- Trend analysis and market insights

## Chatbots & Virtual Assistants

- Customer support and service
- Personal productivity and task management
- Context-aware and personalized interactions

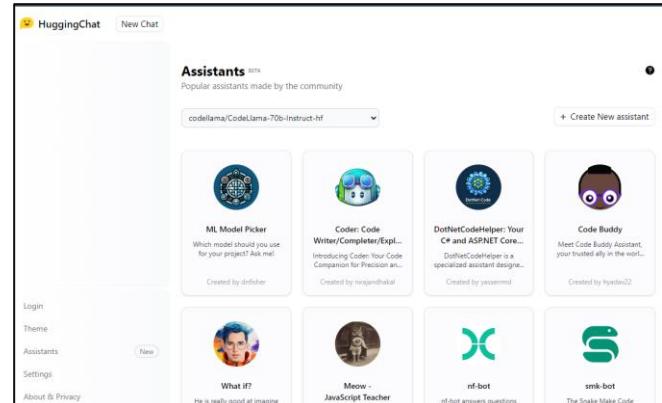
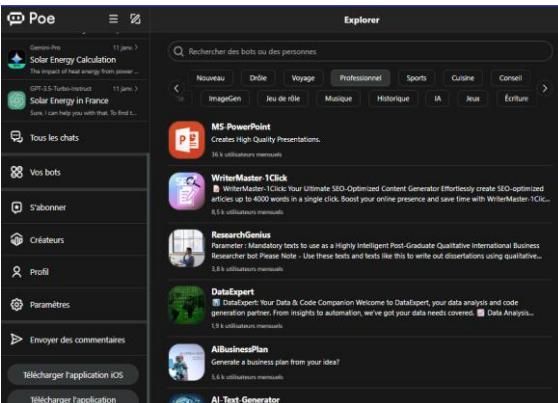
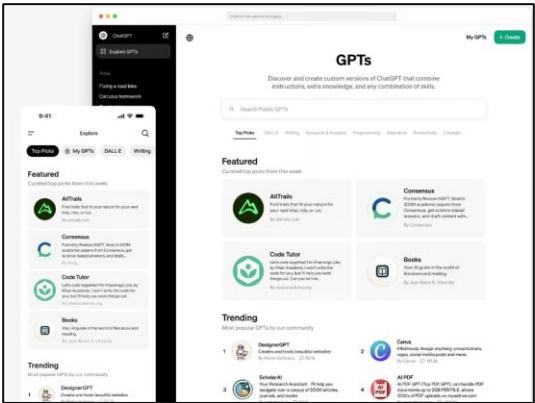
## Many More

- New use cases emerging every day

> Nos REX et notre catalogue représentent plus de 800 cas d'usages...

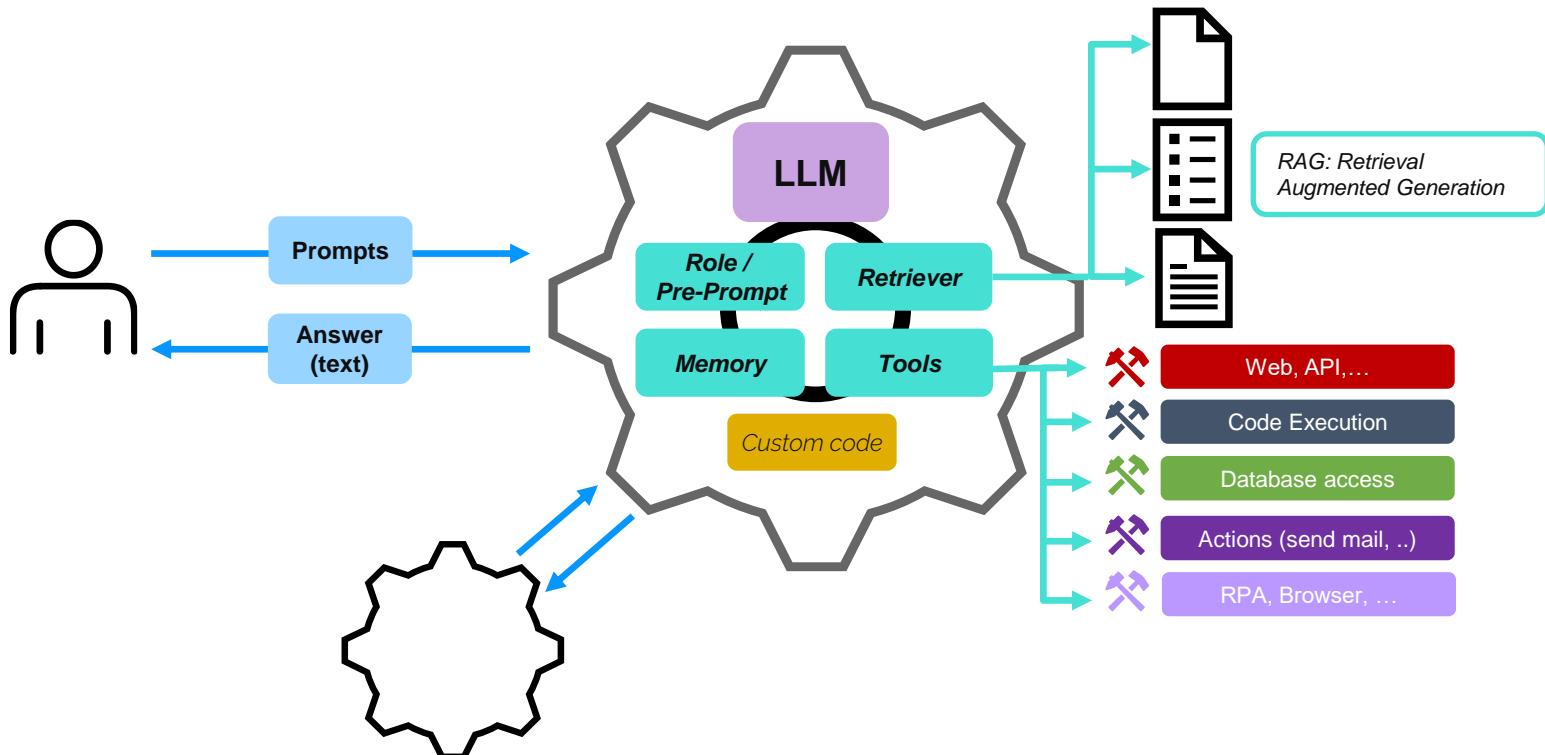
# Emergence of Publics “Agents” no/low-code creators and Marketplaces

Exemples: Poe ‘bots’, OpenAI ‘GPTs’, Azure ‘assistants’, Hugging Face ‘assistants’, Amazon Q ‘applications’,...



# What are behind these use-cases ?

Anatomy of an agent (aka “bot”, “chatbot”, “GPTs”, “copilot”, “virtual assistant”, ...)



# Enterprise Agents Catalog

## Today: separate GenAI PoC or applications

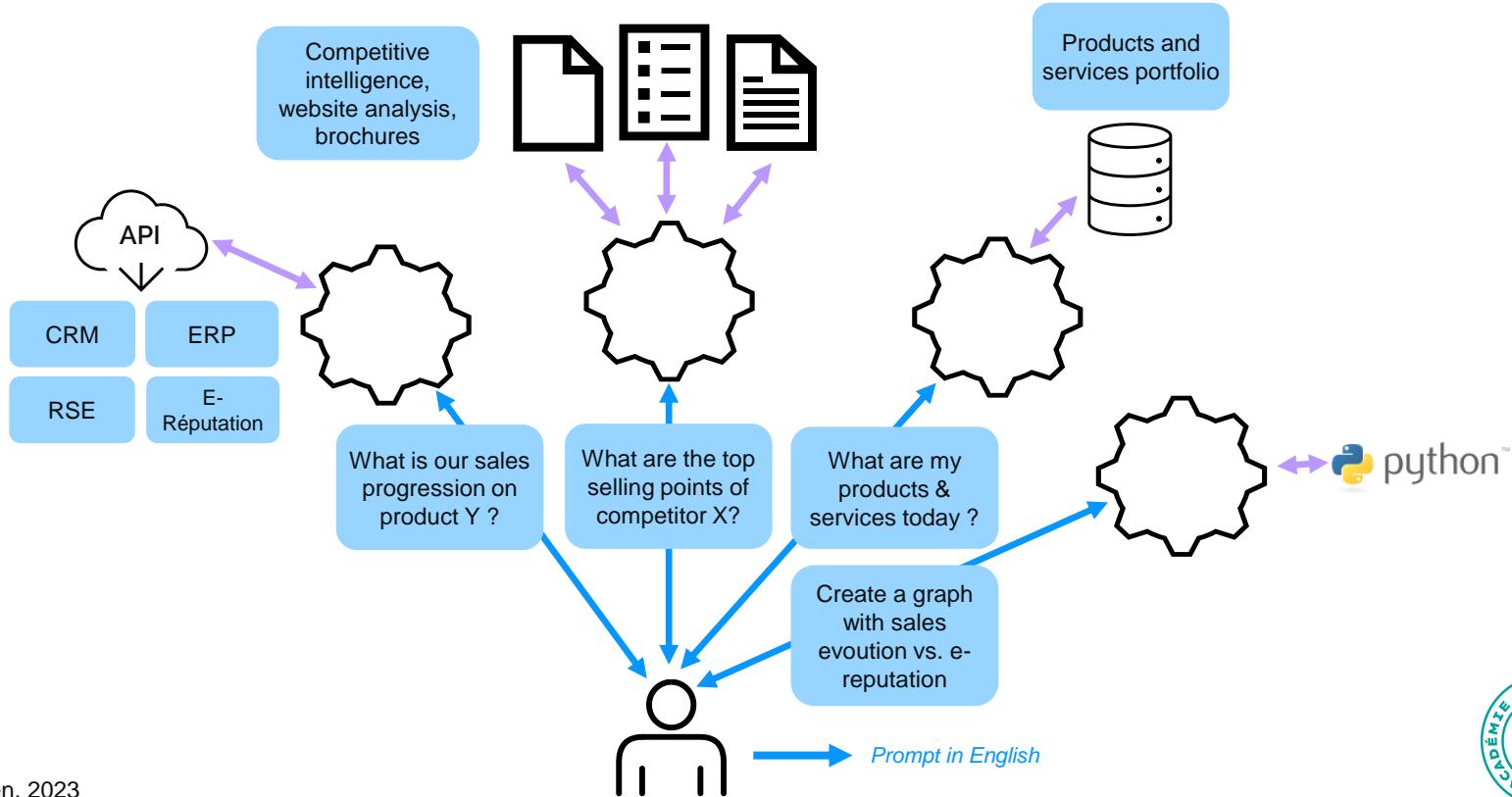
- Content creation, customer service, product design, fraud detection, personalized recommendations, data analysis, supply chain optimization, quality control, cybersecurity, language translation, sentiment analysis, text classification, content moderation, lead generation, market forecasting...

## Tomorrow: more and more Agents (assistant, co-pilots...)

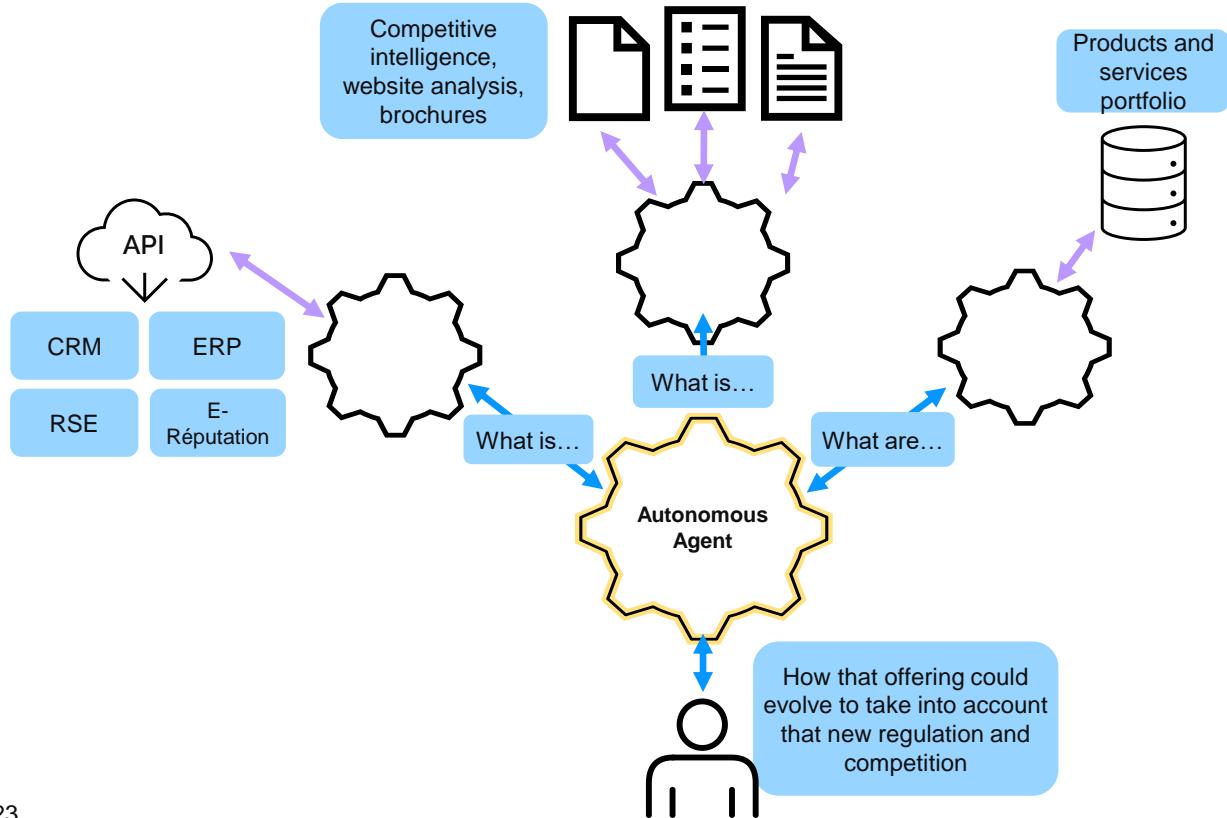
- Standardized
- Integrated
- Cooperative



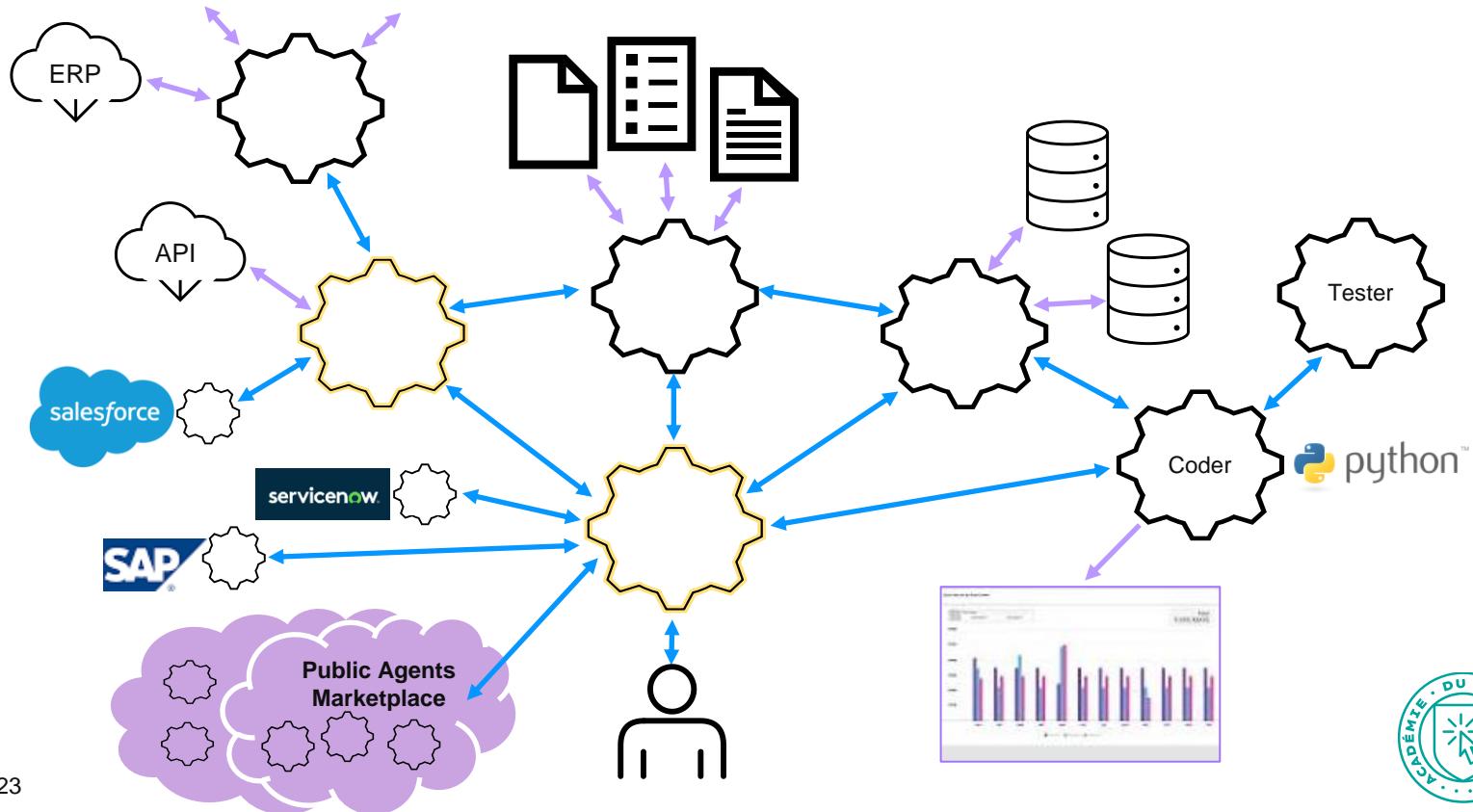
# First Interaction Mode : Independent Agents



## Second mode : humans and hierarchical agents



## Second mode: humans and cooperative (coding) agents



# And that's just the begining ...

*“Overtime, GPTs and assistants are precursors to agents, are going to be able to do much, much more. They'll gradually be able to **plan** and to **perform** more **complex actions** on your behalf.”*

[ Sam Altman – 11/2023 ]

*“Agents are not only going to **change** how everyone **interacts** with computers. They're also going to **upend the software industry**”.*

[ Bill Gates – 11/2023 ]

*“Text is the **universal wire protocol**”.*

[ 7th Schillace Law ]



## More quotes...

*"What I'm seeing with **AI agents** which I think is the exciting Trend that I think **everyone building an AI should pay attention.**"*

[Andrew Ng – Creator of Google Brain]

*"AI field is headed towards **self contained autonomous agents** & it won't be single agent, it will be **many agents working together.**"*

[Andrej Karpathy – co founder of Open AI]

*"Developer becomes the user and so we're evolving toward **any user being able to create its own autonomous agent**. I'm pretty sure that in 5 years from now this will be like something that you learn to do at school."*

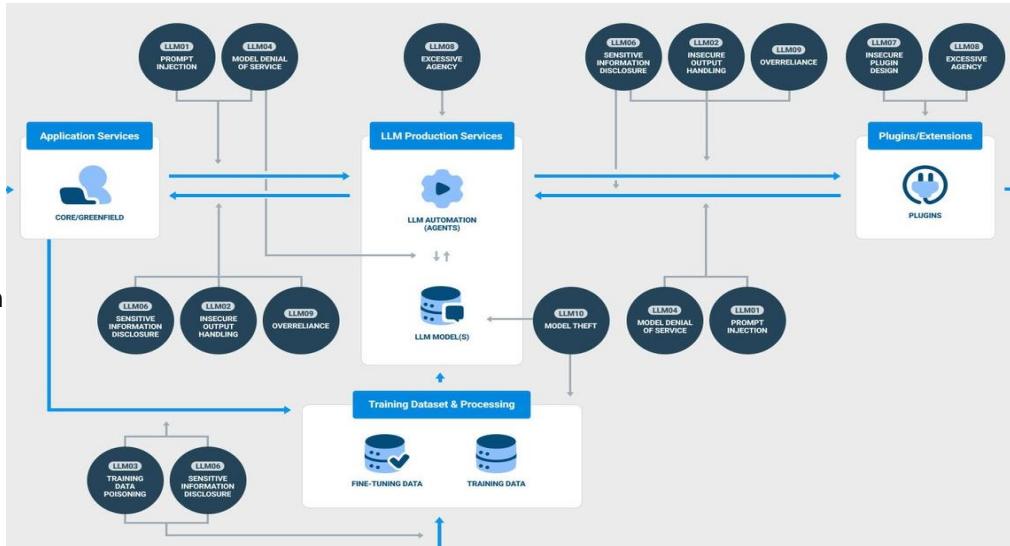
[Arthur Mensch – CEO Mistral AI]



# Yes, but what about risks ?

## Inference risk at Agent Level

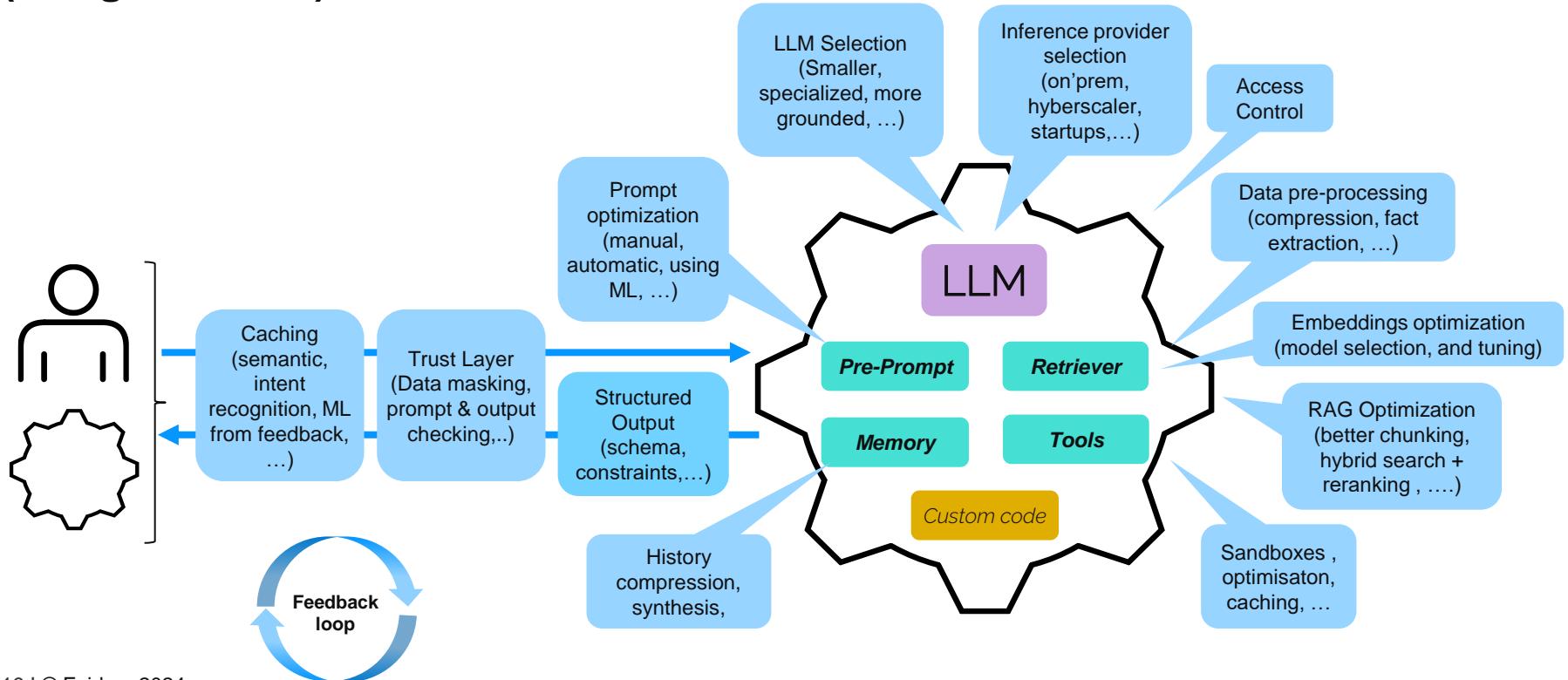
- **Hallucinations:** generation of factually incorrect, irrelevant, or inconsistent text.
  - **Not grounded answer:** When the answer is not in the provided document
- **Unsafe prompt**, according to a policy
- **Attacks:**
  - **Jailbreak:** Attack using prompt injection to bypass safety features in LLMs.
  - **Prompt Extraction:** Attack to reveal system prompt or hidden info in an LLM's context.
  - **Membership Inference:** Data privacy attack to identify if a data sample was in a training set.
  - **Indirect Prompt Injection:** LLM ingests a prompt injection attack indirectly (e.g., through web page).
  - **Information Gathering:** extract user data or leak chat history by interacting in chat sessions, persuading users to divulge information
  - **Prompt Injection Agents:** Force an agent to call a tool.
- **Code generation / execution risks**
- **Tradeoff risks vs. costs vs. performance**



OWASP Top 10 for Large Language Model Applications

# Risk mitigation and cost reduction

## (At Agents level)





# Python Environnement



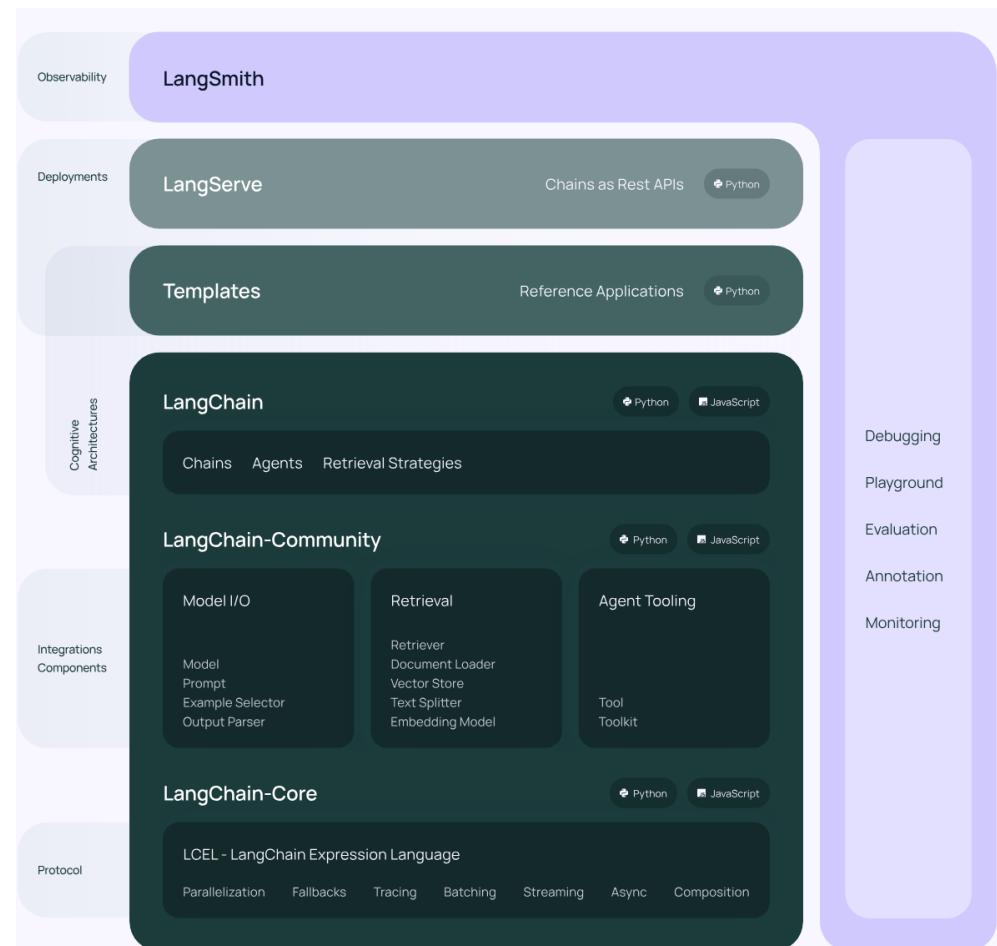
# LangChain 1/2

- **Leading Open-source GenAI Framework**

- Abstractions for LLMs, Tools, Memory, Retrievers, ...
- Very large ecosystem: 60+LLMs, 60+ vector databases, many retrieval algorithms, tools, data source connectors, memory patterns, ...
- Relatively easy to extend and patch
- Very fast moving: 1 release every ≈3 days
- Very large ecosystem
- 2000+ public GitHub repo tagged “langchain”, 2200+ contributors
- Alternatives to Assistants API and GPTs.

- **LangSmith, LangServe...**

- **Partnerships with AWS, GCP...**



# LangChain 2/2

## LCEL: LangChain Expression Language

- Declarative way to easily compose chains together
- Async / Streaming without code change
- Build-in feature for retries / fallback / monitoring

```
prompt = ChatPromptTemplate.from_template("tell me a short joke about {topic}")
output_parser = StrOutputParser()

chain = prompt | model | output_parser

chain.invoke({"topic": "ice cream"})
```

## No free lunch

- LCEL is non-Pythonic (functional)
- Fast moving, frequent refactoring
- Doc sometime incomplete

## Alternatives librairies

- LLamaIndex
- HayStack
- Microsoft Semantic Kernel
- ...

## And also no-code / low code solutions

- Flowise (based on LangChain)
- Amazon Q
- Agent for Bedrock
- OpenAI Assistant
- <https://www.coze.com/>
- ...



# Poetry

- **Dependency management tool** for Python.
- **Makes it easy to install**, manage, and update Python packages
- **Especially useful** for AI applications (lot of fast changing packages)
- **Key features:**
  - Virtual environments
  - Dependency resolution
  - Reproducibility
- **Use:**
  - Install Poetry: see doc.
  - Create a new project: poetry new my-project
  - Add a dependency: poetry add numpy
  - Install dependencies: poetry install



Real Python

<https://realpython.com/dependency-management-python-poetry/>

# Pydantic

- Data validation and serialization
- Ensures data is valid and consistent
- Supports complex data structures
- Can be used to generate documentation
- JSON interoperability
- Base of a large ecosystem:
  - **FastAPI**
  - **Typer**
  - **SQLModel**
  - **LangChain**
- See [Kludex/awesome-pydantic: A curated list of awesome things related to Pydantic!](#) 



Bringing schema and sanity to your data



# Streamlit

- Build and deploy interactive web applications
- Makes AI models more accessible and understandable
- No need to write front-end code
- Supports a variety of widgets and visualizations
- Can be used to create dashboards, reports, and other interactive applications

## How:

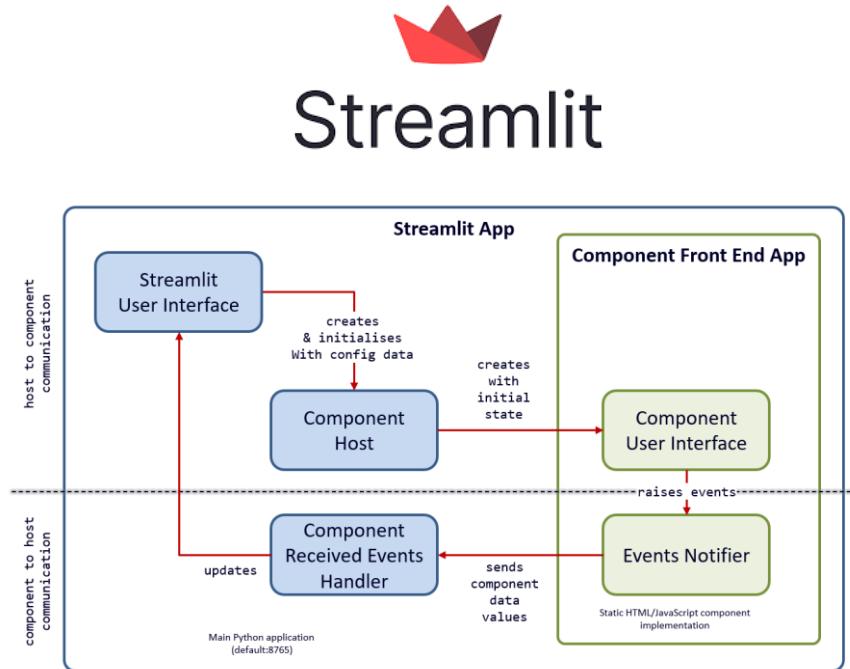
- One Python script, rerun after each user interaction

## Pitfalls:

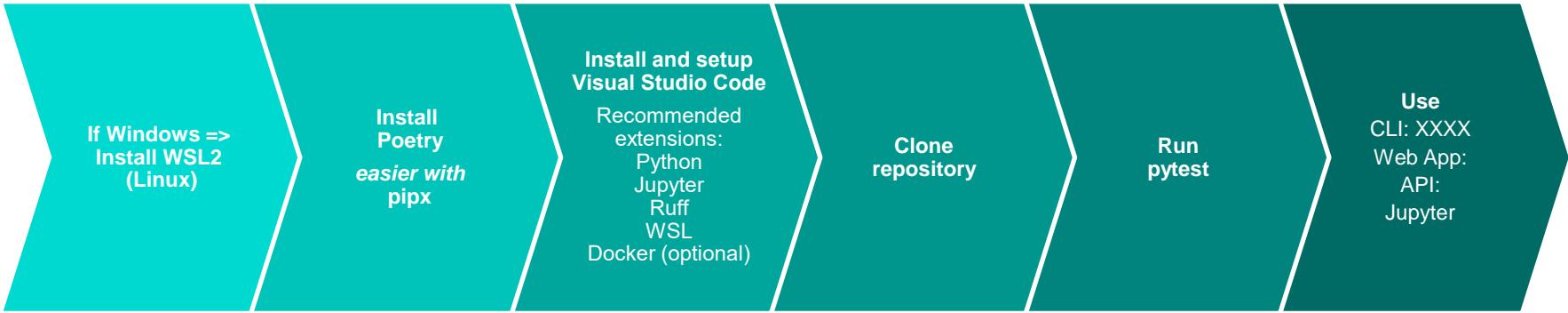
- Performance
- Customization: Limited UI flexibility for advanced features.
- State Management: Preserving state across reruns can be challenging.

## Alternatives:

- Gradio, Taipis, Dash, Panel, ...



# Environnement (Linux + Python)



# Practicum

## Assignment: Install Environment

- Clone project template from GIT
  - Recommended : create a branch
  - Add API keys in ~/.bashrc



## Test it with CLI

```
git config --global user.name "*****"  
git config --global user.email \*\*\*\*\*@gmail.com
```

```
git checkout -b *****-blueprint
```

```
python python/main_cli.py run "joke"
```

### Hint

To sync the project while keeping your change:  
make rebase      (call git stash/ rebase)

Chapter

# Learning Language Model & Lanchain



# Functional Programming with LCEL



# LCEL : LangChain Expression Language

[https://api.python.langchain.com/en/stable/runnables/langchain\\_core.runnables.base.Runnable.html](https://api.python.langchain.com/en/stable/runnables/langchain_core.runnables.base.Runnable.html)

- **LangChain**

Tool to compose Chains of processing (mostly using LLM)

- **LangChain Expression Language (LCEL):** declarative way to compose chains effortlessly

- **Runnable**

Allows chaining code together into sequences, enabling the passing of output to the next runnable

- Close to 'Monads' in Functional programming

- **Functional programming benefit**

- Emphasizes use of pure functions and immutable data
- Predictable behavior, easier debugging, and scalability
- Streamline app development, support chaining, parallelization, and retries
- Seamless deployment, optimized parallel execution, and access to intermediate results.

Component	Input Type	Output Type
<b>Prompt</b>	Dictionary	PromptValue
<b>ChatModel</b>	Single string, list of chat messages or a PromptValue	ChatMessage
<b>LLM</b>	Single string, list of chat messages or a PromptValue	String
<b>OutputParser</b>	The output of an LLM or ChatModel	Depends
<b>Retriever</b>	Single string	List of Documents
<b>Tool</b>	Single string or dictionary, depending on the tool	Depends on the tool

```
model = ChatOpenAI()
prompt = ChatPromptTemplate.from_template("tell me a joke about {topic}")
chain = prompt | model
```



# LCEL Runnable Methods

<a href="#"><code>batch</code>(inputs[, config, return_exceptions])</a>	Default implementation runs invoke in parallel using a thread pool executor.
<a href="#"><code>bind</code>(**kwargs)</a>	Bind arguments to a Runnable, returning a new Runnable.
<a href="#"><code>config_schema</code>(*[, include])</a>	The type of config this runnable accepts specified as a pydantic model.
<a href="#"><code>get_graph</code>([config])</a>	Return a graph representation of this runnable.
<a href="#"><code>get_input_schema</code>([config])</a>	Get a pydantic model that can be used to validate input to the runnable.
<a href="#"><code>get_name</code>([suffix, name])</a>	Get the name of the runnable.
<a href="#"><code>get_output_schema</code>([config])</a>	Get a pydantic model that can be used to validate output to the runnable.
<a href="#"><code>get_prompts</code>([config])</a>	
<a href="#"><code>invoke</code>(input[, config])</a>	Transform a single input into an output.
<a href="#"><code>map</code>()</a>	Return a new Runnable that maps a list of inputs to a list of outputs.
<a href="#"><code>pick</code>(keys)</a>	Pick keys from the dict output of this runnable.
<a href="#"><code>pipe</code>(*others[, name])</a>	Compose this runnable with another object to create a RunnableSequence.
<a href="#"><code>stream</code>(input[, config])</a>	Default implementation of stream, which calls invoke.
<a href="#"><code>transform</code>(input[, config])</a>	Default implementation of transform, which buffers input and then calls stream.
<a href="#"><code>with_config</code>([config])</a>	Bind config to a Runnable, returning a new Runnable.
<a href="#"><code>with_fallbacks</code>(fallbacks, *[, ...])</a>	Add fallbacks to a runnable, returning a new Runnable.
<a href="#"><code>with_listeners</code>(*[, on_start, on_end, on_error])</a>	Bind lifecycle listeners to a Runnable, returning a new Runnable.
<a href="#"><code>with_retry</code>(*[, retry_if_exception_type, ...])</a>	Create a new Runnable that retries the original runnable on exceptions.
<a href="#"><code>with_types</code>(*[, input_type, output_type])</a>	Bind input and output types to a Runnable, returning a new Runnable.



# Practicum

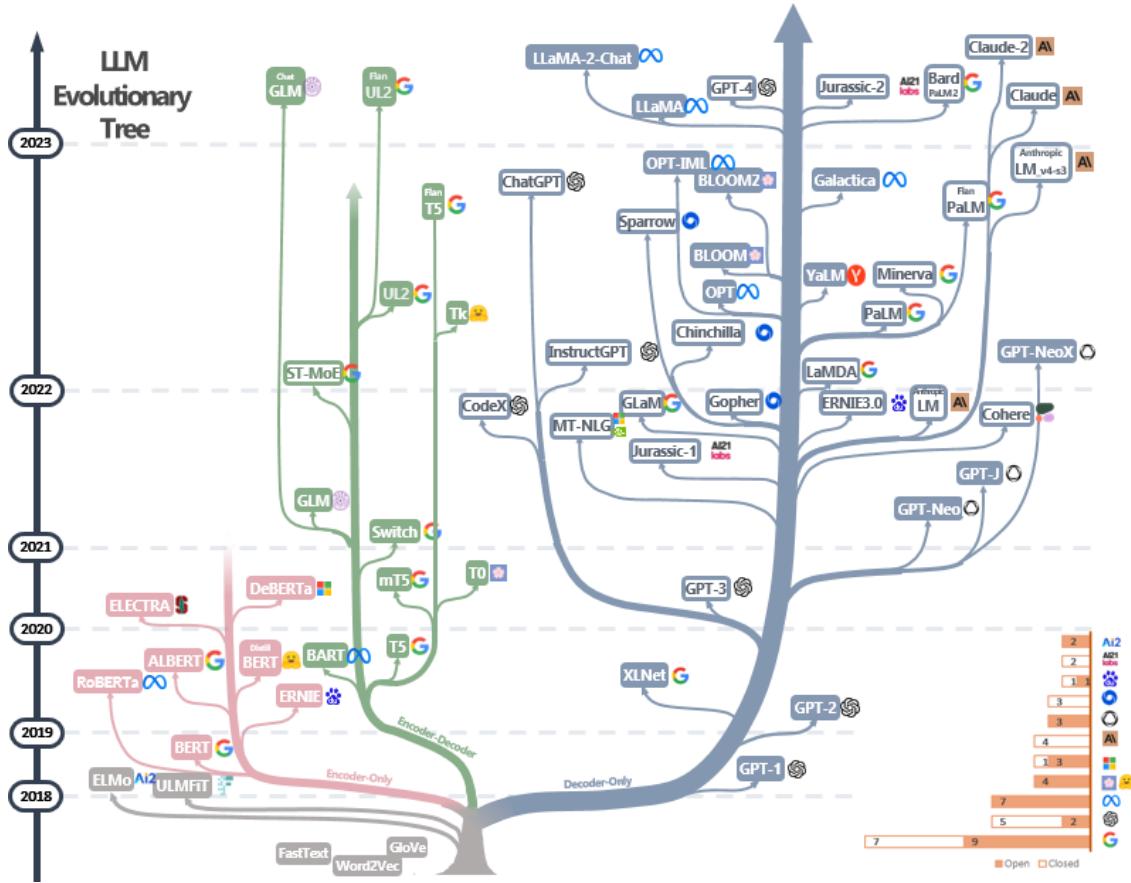
Play with LCEL simple Runnable  
(with Notebook)



# Similarity Search



# At the begining



# Tokenization

- Tokenization is the process of breaking down a text into individual units.
- Tokens can be words, phrases, or even individual characters.
- In GenAI: subword tokenization
- Common libraries:
  - Tiktoker (OpenAI)
  - Nltk
  - Hugging Face
- Not all tokenizers are equals
- More tokens => \$\$\$

Try: <https://huggingface.co/spaces/Xenova/the-tokenizer-playground>

gpt-4 / gpt-3.5-turbo / text-embedding-ada-002

I'm in a training about GenAI and Agents. Université du numérique organized it. Le formateur (trainer) est en train d'expliquer la tokenization. ;-)

Tokenization is an essential step in many GenAI and agent applications, such as: Natural language processing, Machine translation, Speech recognition, Chatbots

|

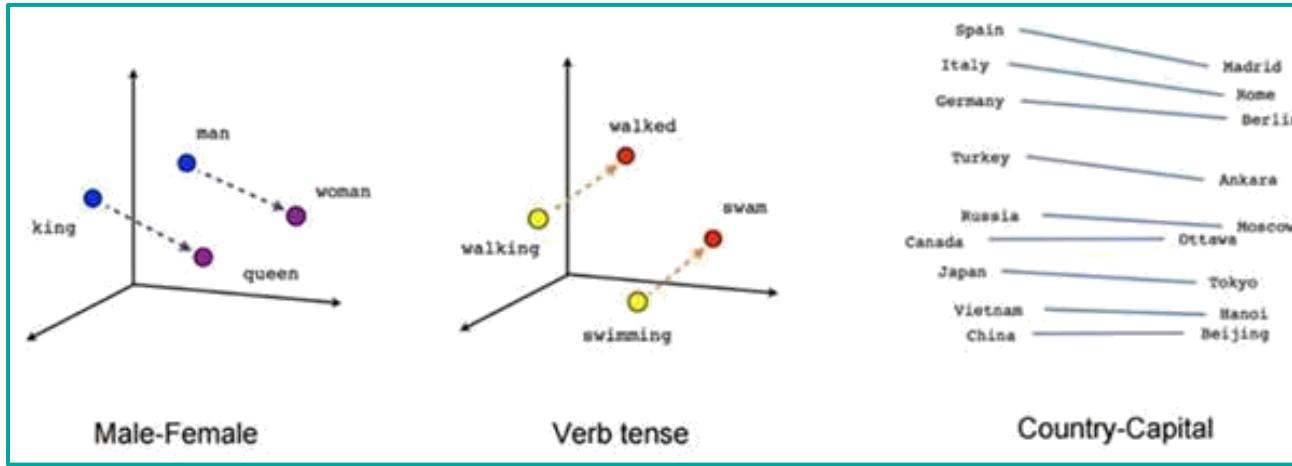
TOKENS    CHARACTERS  
70            311

I'm in a training about GenAI and Agents. Université du numérique organized it. Le formateur (trainer) est en train d'expliquer la tokenization. ;-)

Tokenization is an essential step in many GenAI and agent applications, such as: Natural language processing, Machine translation, Speech recognition, Chatbots

[40, 2846, 304, 264, 4967, 922, 9500, 70767, 323, 51354, 13, 220, 15915, 13109, 3930, 1661, 78129, 17057, 433, 13, 2009, 1376, 11067, 320, 84876, 8, 1826, 665, 5542, 294, 6, 30992, 72684, 1208, 4037, 2065, 13, 220, 90608, 340, 3404, 2065, 374, 459, 7718, 3094, 304, 1690, 9500, 15836, 323, 8479, 8522, 11, 1778, 439, 25, 18955, 4221, 8863, 11, 13257, 14807, 11, 39841, 18324, 11, 13149, 63005, 198]

# Embeddings



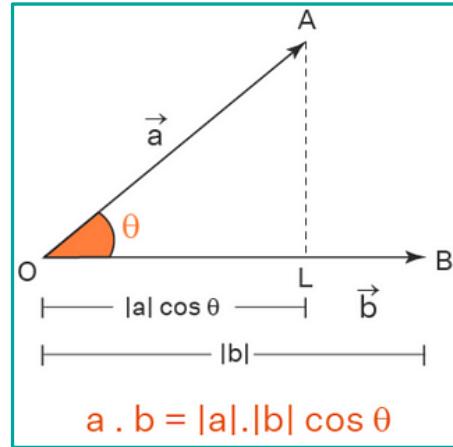
$$V(\text{king}) - V(\text{man}) + V(\text{woman}) = V(\text{queen})$$



Source: John, Vineet. (2016). Rapid-Rate: A Framework for Semi-supervised Real-time Sentiment Trend Detection in Unstructured Big Data. 10.13140/RG.2.2.32385.04966.

# Alignment Score

**Normalized dot product = cosine projection**



Queen = [0.97, 0.03, 0.02]

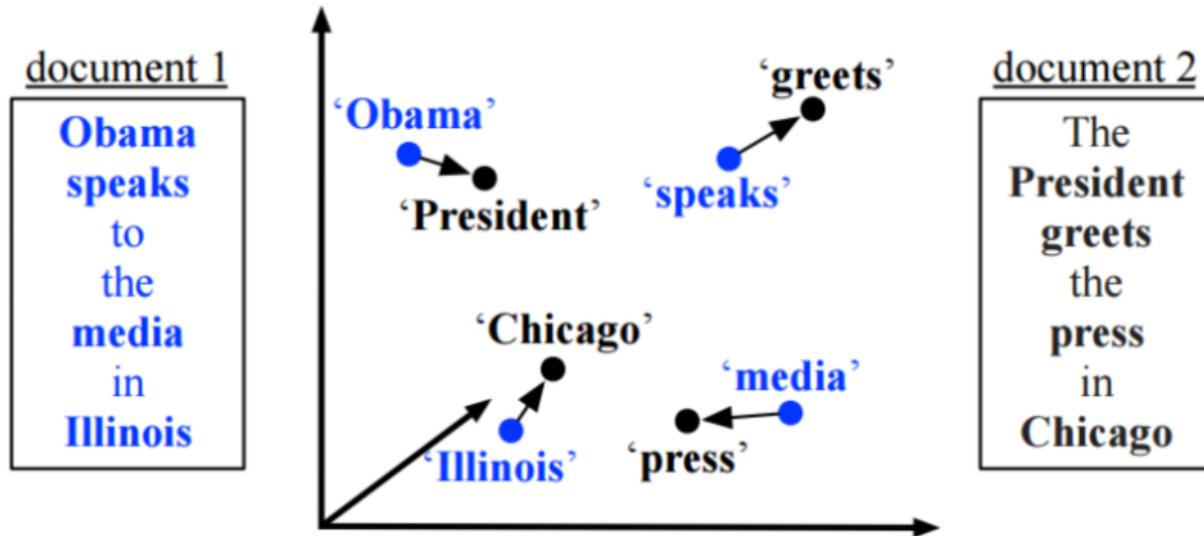
$$(0.99 * 0.97) + (0.01 * 0.03) + (0.02 * 0.02) = \mathbf{0.961}$$

King = [0.99, 0.01, 0.02]

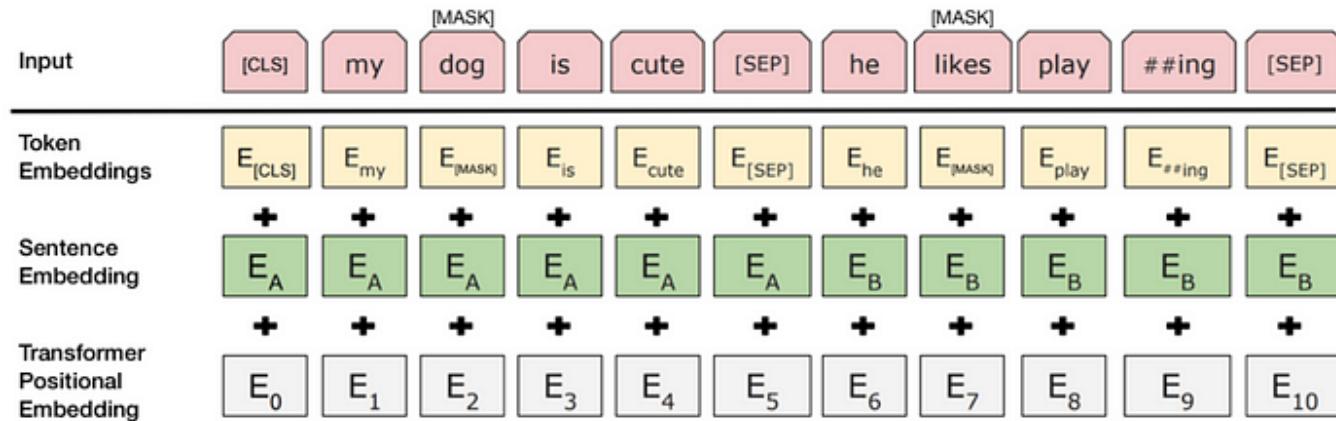
$$(0.99 * 0.01) + (0.01 * 0.02) + (0.02 * 0.02) = \mathbf{0.0105}$$

Dog = [0.01, 0.02, 0.02]

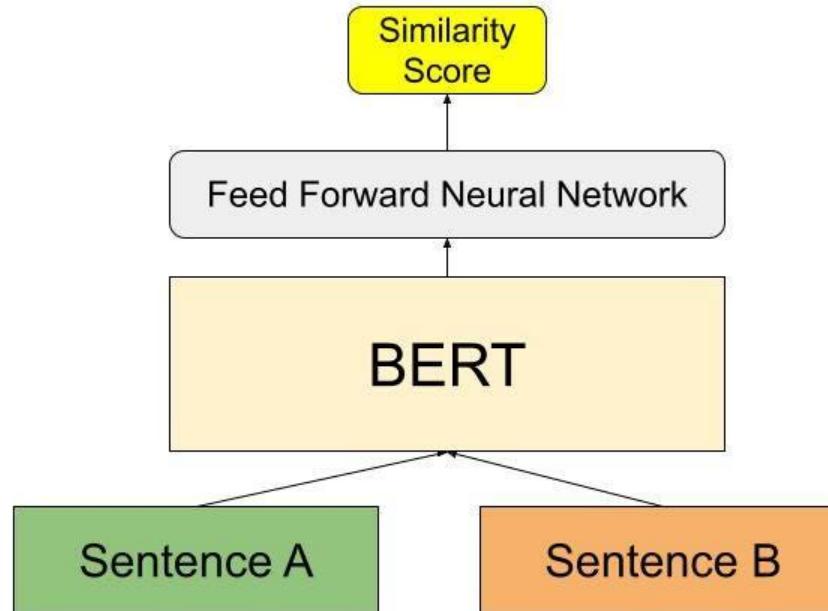
# Sentence Embedding



# Sentence Embedding (with BERT)

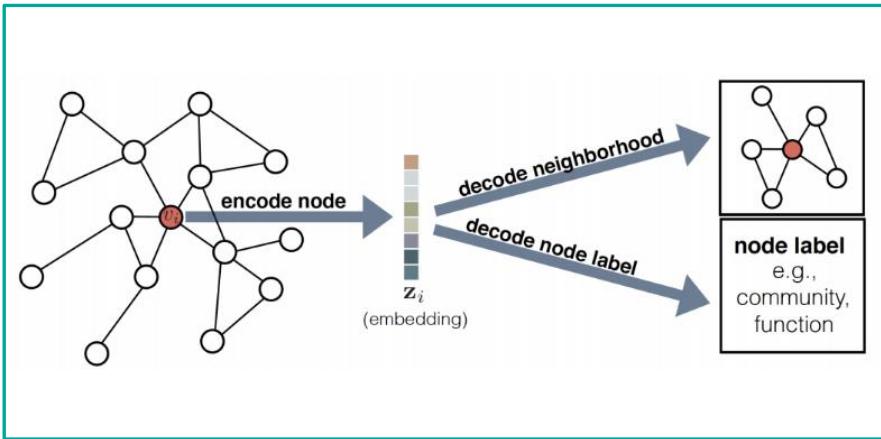


# Similarity Matching (simple)

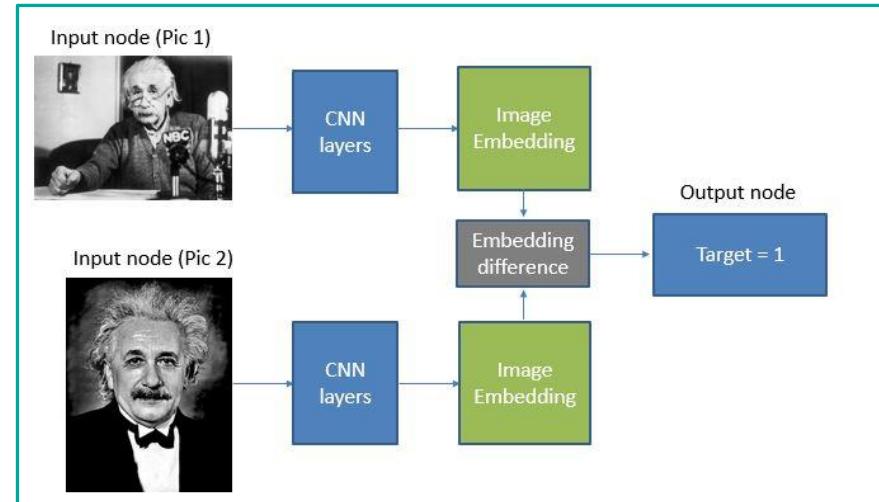


<https://www.pinecone.io/learn/sentence-embeddings/>

# Graph Embedding, Image Embedding, ... (Off topic)



<https://neo4j.com/developer/graph-data-science/graph-embeddings/>



<https://www.analyticsvidhya.com/blog/2022/07/recommending-similar-images-using-image-embedding/>

# Embeddings Models Comparaisons

## MTEB Leaderboard - a Hugging Face Space

Overall	Bitext Mining	Classification	Clustering	Pair Classification	Reranking	Retrieval	STS	Summarization
English	Chinese	French	Polish					
<b>Overall MTEB English leaderboard 🌎</b>								
<ul style="list-style-type: none"> <li>Metric: Various, refer to task tabs</li> <li>Languages: English</li> </ul>								
Rank	Model	Model Size (GB)	Embedding Dimensions	Max Tokens	Average (56 datasets)	Classification Average (12 datasets)	Clustering Average (11 datasets)	Pair Classification Average (3 datasets)
Rank	Model	Model Size (GB)	Embedding Dimensions	Max Tokens	Average (56 datasets)	Classification Average (12 datasets)	Clustering Average (11 datasets)	Pair Classification Average (3 datasets)
1	<a href="#">SFR-Embedding-Mistral</a>	14.22	4096	32768	67.56	78.33	51.67	88.54
3	<a href="#">GritLM-7B</a>	14.48	4096	32768	66.76	79.46	50.61	87.16
4	<a href="#">e5-mistral-7b-instruct</a>	14.22	4096	32768	66.63	78.47	50.26	88.34
2	<a href="#">voyage-lite-02-instruct</a>	2.45	1024	4000	67.13	79.25	52.42	86.87
10	<a href="#">voyage-lite-01-instruct</a>		1024	4000	64.49	74.79	47.4	86.57
6	<a href="#">echo-mistral-7b-instruct-last</a>	14.22	4096	32768	64.68	77.43	46.32	87.34

### Some criterias:

- API / loaded
- Size
- Speed (CPU/GPU)
- Benchmarks
- « Instruct »
- Language
- License

# Embeddings Models - Practically

[huggingface.co/sentence-transformers/all-MiniLM-L6-v2](https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2)

 Hugging Face  Models Datasets

**sentence-transformers/all-MiniLM-L6-v2** like 409

Sentence Similarity PyTorch TensorFlow Rust Sentence Transformers s2orc  
 yahoo\_answers\_topics code\_search\_net search\_qa eli5 snli multi\_nli wikihow  
 embedding-data/sentence-compression embedding-data/flickr30k-captions embedding-data/altlex  
 embedding-data/SPECTER embedding-data/PAQ\_pairs embedding-data/WikiAnswers English  
 arxiv:2104.08727 arxiv:1704.05179 arxiv:1810.09305 License: apache-2.0

[Model card](#) [Files and versions](#) [Community 14](#) [Edit model](#)

**all-MiniLM-L6-v2**

This is a [sentence-transformers](#) model: It maps sentences & paragraphs to a 384 dimensional dense vector space and can be used for tasks like clustering or semantic search.

**Usage (Sentence-Transformers)**

Using this model becomes easy when you have [sentence-transformers](#) installed:

```
pip install -U sentence-transformers
```

Downloads last month  
1,547,881 

**Hosted inference API**

Sentence Similarity Example 1

Source Sentence  
[A service providing RDMS storage on AWS](#)

Sentences to compare to  
 Cloud based relational database.  
 Cloud based relational messaging [solution](#)  
 on-premise database

Add Sentence

Compute

Computation time on Intel Xeon 3rd Gen Scalable cpu: 3.813 s

Sentence	Computation Time (s)
Cloud based relational database.	0.435
Cloud based relational messaging solution	0.405
on-premise database	0.319

# Practicum

## Assignment

Encode a sentence with one or several embeddings model

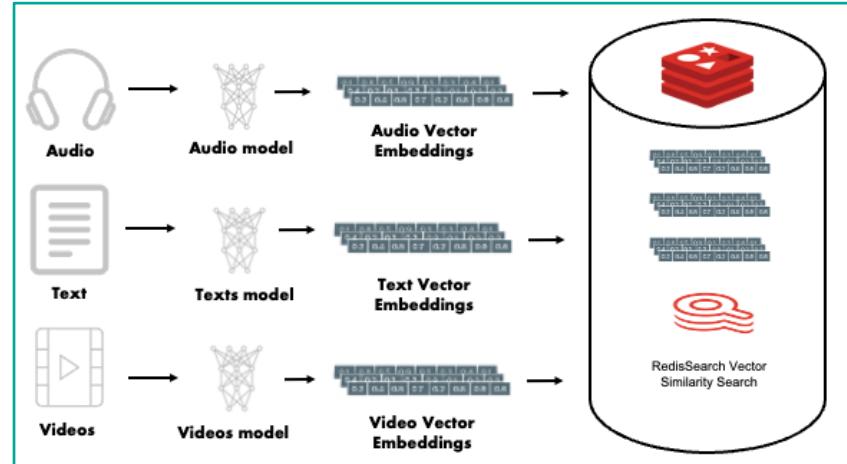
Calculate cosine projection  
(with Notebook)



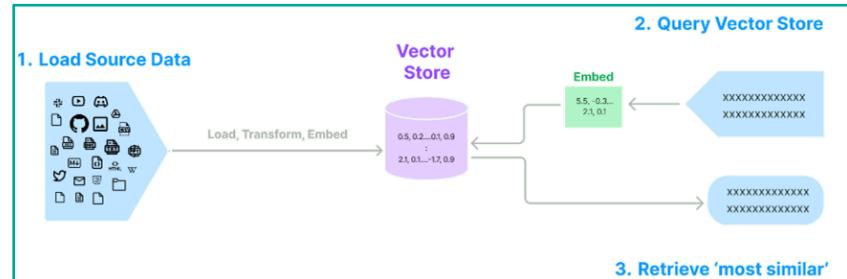
15 min

# Vector Store

- **Databases optimized** for storing, searching, and manipulating vector representations of data (embeddings)
- **Vectors stored** in an approximate nearest neighbor search (ANNS) data structure, for efficient querying.
- **Users submit queries** as vectors=> and the vector store returns the most similar vectors based on cosine similarity.



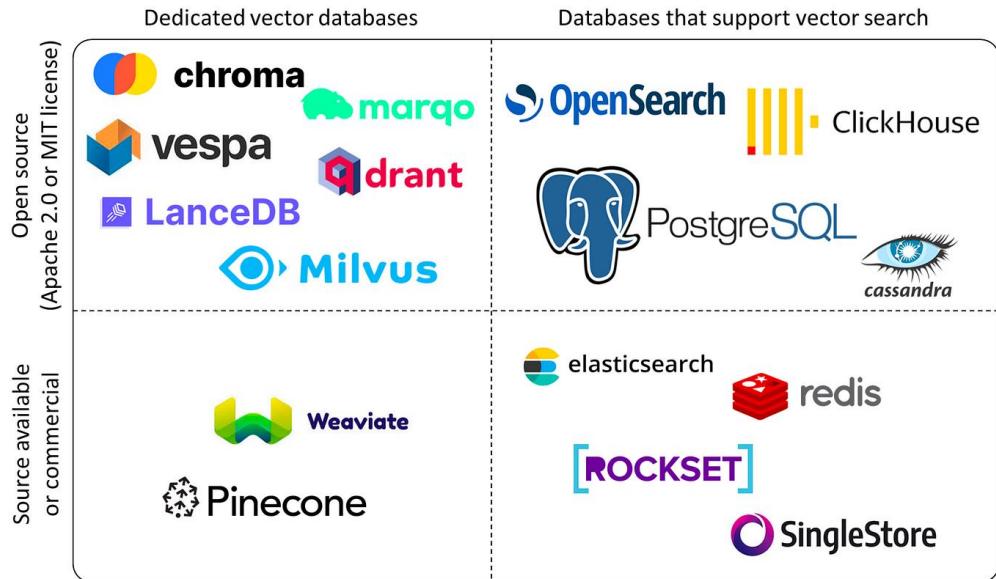
Source: <https://partee.io/2022/08/11/vector-embeddings/>



# Vector Store Selection

## Distinctive Features

- Scalability: Efficient handling of large-scale data
- Speed (in-memory, GPU based, ...)
- Accuracy, through advanced indexing techniques
- Integration with existing database technologies: RDMS (PostgreSQL,...), Search Engine (ElasticSearch, ...), Key-Value (Redis)
- Hybrid search capabilities (vectors + keywords)
- License (open source, commercial, ...)
- Hosted



<https://blog.det.life/why-you-shouldnt-invest-in-vector-databases-c0cd3f59d23c>

# Document Loaders in LangChain

## LangChain Document Loader

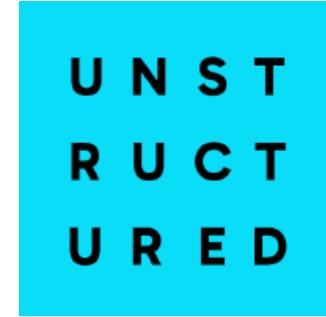
- Provide a standardized way to load and process various types of documents
- Ingest and analyze data from a wide range of sources
- Create a list of Documents
- **Examples:**
  - TextLoader: loads a text file
  - PDFMinerLoader: extracts text and metadata from a PDF file
  - WebPageLoader: loads the text content of a web page
  - YouTubeTranscriptLoader: loads the transcript of a YouTube video
  - RecursiveUrlLoader
  - AsyncHTMLLoader
  - SqlLoader
  - airbyte: Loads data from the Airbyte API.
  - ...



# Unstructured Loader

## Unstructured

- Python package to extracts clean text from raw source documents like PDFs, Word docs, images, etc.
- Supports loading many file types
- Can retain document structure by extracting different text "elements" like titles, paragraphs, etc. 1
- Integrated into LangChain via several Loaders for different file types
- Can be used locally by installing dependencies, or via hosted AP



# Text Splitters

**Text Splitters split long documents into smaller chunks to fit model context windows**

## Role

- Split text into semantically meaningful chunks
- Combine chunks until reaching desired size
- Add overlap between chunks to maintain context

## Types

- Recursive (splits on characters like newlines, spaces)
- By markup (HTML, Markdown headers)
- By language syntax (code, tokens)
- Semantic chunking (by embedding similarity)



### Hint

Can evaluate splitters using [Chunkviz tool](#)

Splitter: Recursive Character Text Splitter

Chunk Size: 200

Chunk Overlap: 0

Total Characters: 2658

Number of chunks: 24

Average chunk size: 110.8

One of the most important things I didn't understand about the world when I was a child is the degree to which the returns for performance are superlinear.

Teachers and coaches implicitly told us the returns were linear. "You get out," I heard a thousand times, "what you put in." They meant well, but this is rarely true. If your product is only half as good as your competitor's, you don't get half as many customers. You get no customers, and you go out of business.

It's obviously true that the returns for performance are superlinear in business. Some think this is a flaw of capitalism, and that if we changed the rules it would stop being true. But superlinear returns for performance are a feature of the world, not an artifact of rules we've invented. We see the same pattern in fame, power, military victories, knowledge, and even benefit to humanity. In all of these, the rich get richer. [1]

You can't understand the world without understanding the concept of superlinear returns. And if you're ambitious you definitely should, because this will be the wave you surf on.

It may seem as if there are a lot of different situations with superlinear returns, but as far as I can tell they reduce to two fundamental causes: exponential growth and thresholds.

The most obvious case of superlinear returns is when you're working on something that grows exponentially. For example, growing bacterial cultures. When they grow at all, they grow exponentially. But they're tricky to grow. Which means the difference in outcome between someone who's adept at it and someone who's not is very great.

Startups can also grow exponentially, and we see the same pattern there. Some manage to achieve high growth rates. Most don't. And as a result you get qualitatively

# Embeddings to Vector Store in LangChain

## Instance of class (VectorStore)

- Many exists !
- Common one: Chroma, FAISS, LanceDB, Pinecone, Milvus, Qdrant, Weaviate, Elasticsearch, PostgreSQL, Redis, ..

```
from langchain_community.document_loaders import TextLoader
from langchain_openai import OpenAIEMBEDDINGS
from langchain_text_splitter import CharacterTextSplitter
from langchain_community.vectorstores import Chroma

# Load the document, split it into chunks, embed each chunk and load it into the vector store.
raw_documents = TextLoader('../.../state_of_the_union.txt').load()
text_splitter = CharacterTextSplitter(chunk_size=1000, chunk_overlap=0)
documents = text_splitter.split_documents(raw_documents)
db = Chroma.from_documents(documents, OpenAIEMBEDDINGS())
```

```
query = "What did the president say about Ketanji Brown Jackson"
docs = db.similarity_search(query)
print(docs[0].page_content)
```



# Practicum

**Assignment**  
- Play with the Notebook

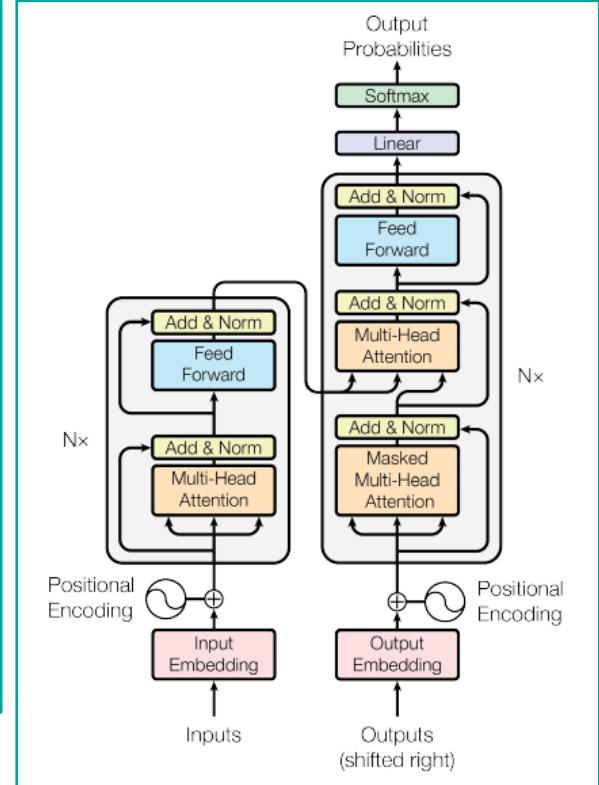


# LLM Essentials

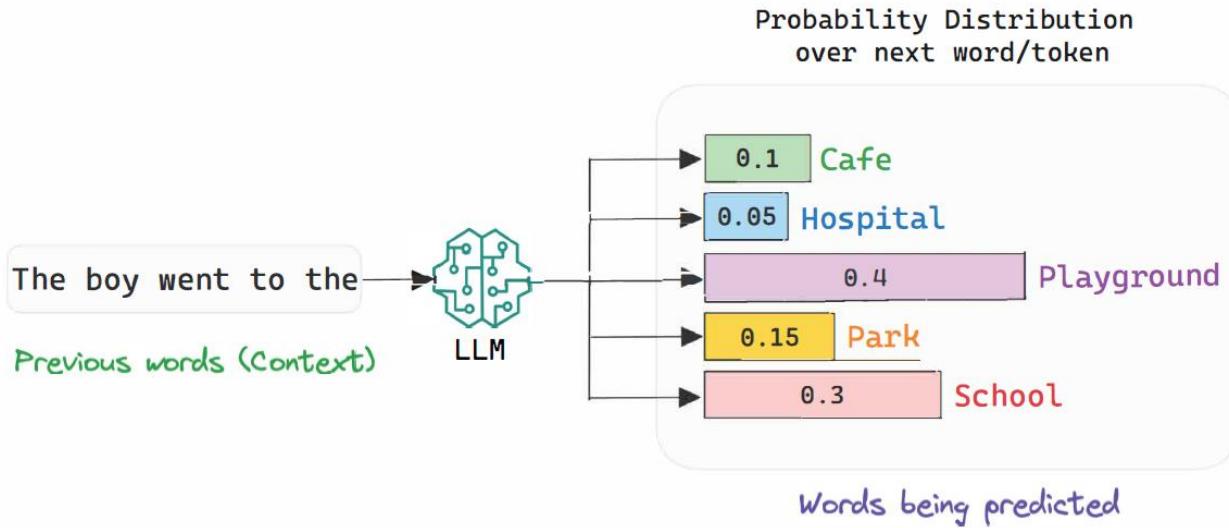


# « Attention Is All What You Need »

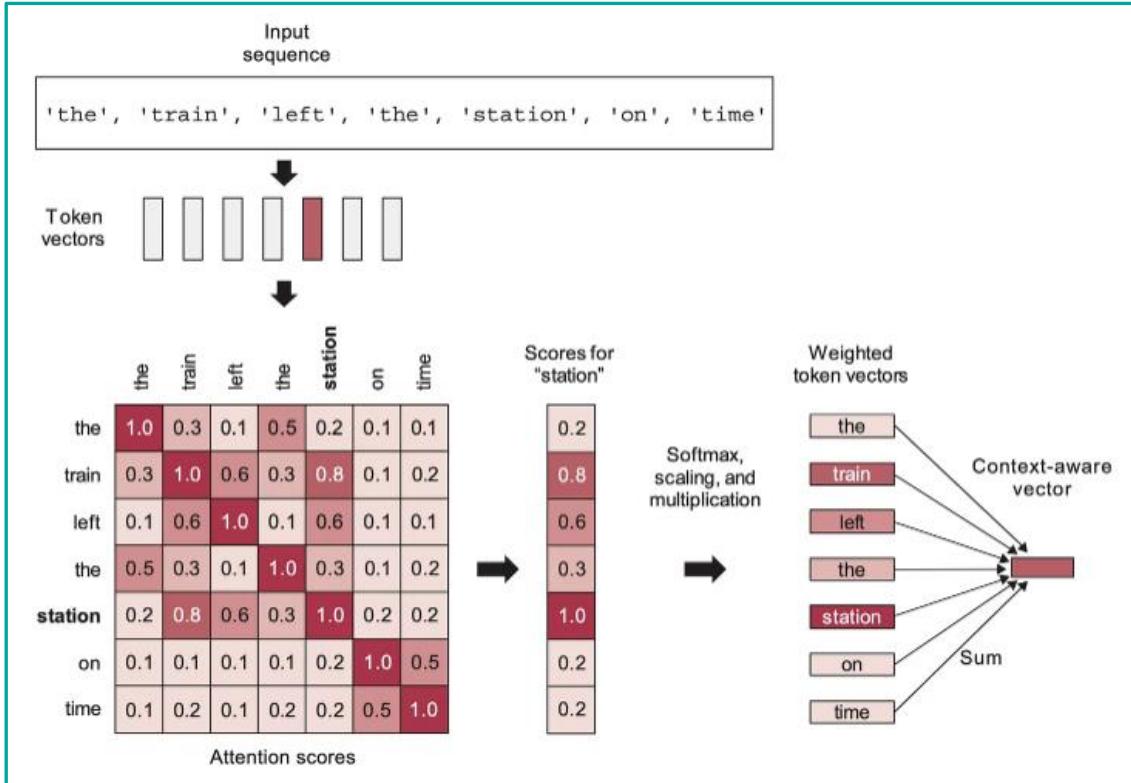
- Published in 2017  
Before: limitations with LSTM, RNN,...
- Among most influential paper in CS  
100000 citations in 2024
- Still very robust.



# LLM are next token predictors



# Self-Attention



## Note

- Compute attention scores between “station” and **every other word in the sequence**.
- Used to weight a sum of word vectors => new “station” vector.
- Context vectors are used for another stacked self-attention later.
- One token at a time !

[How I think about LLM prompt engineering](#)

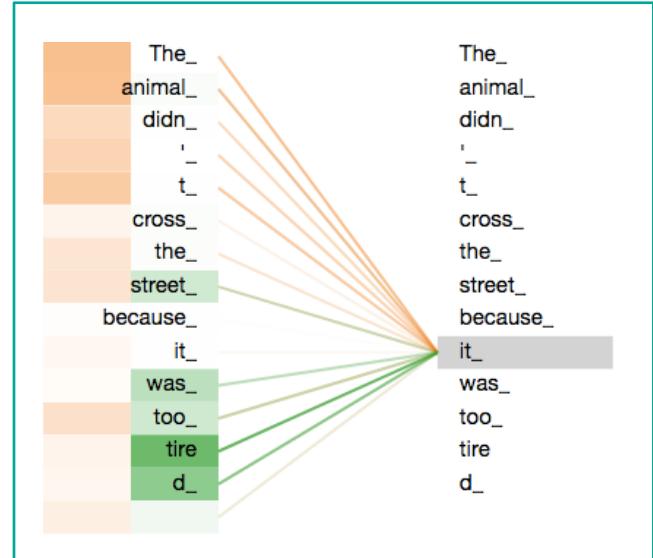
# Multi Heads, Stacked Layers

## Multi-head attention layer

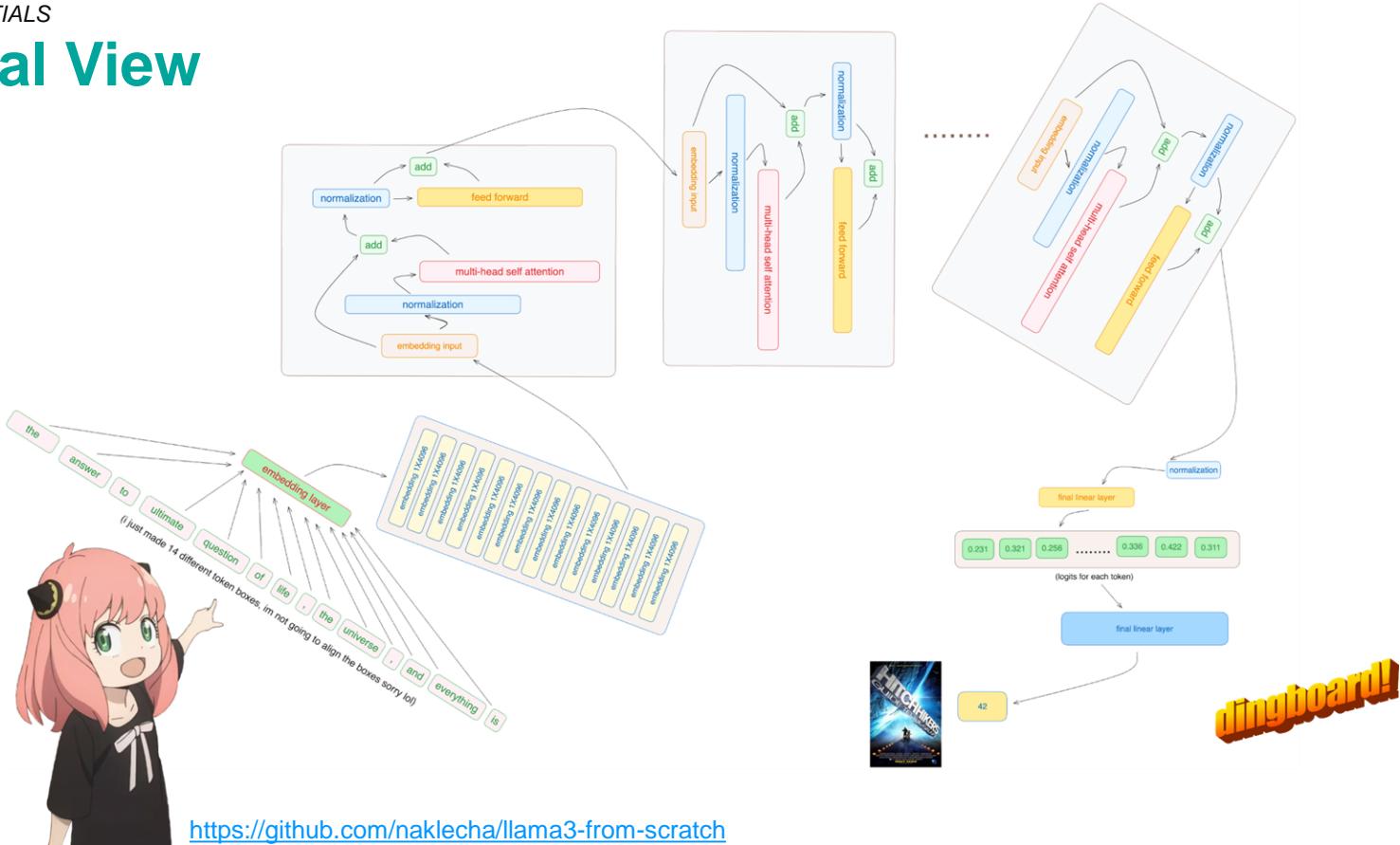
- Each “head” attends to a different subspace of the input data.
- Can encode multiple relationships and nuances for each word.
- Outputs are weighted and concatenated to form the final context vector.
- Ex: ChatGPT-3,5 has 16 attention heads

## Stacked self-attention layers

- “Transformers” architecture
- Output of an attention layer is input another one
- Ex : ChatGPT-3,5 has 12 layers
- Catch more complex relationship in input data

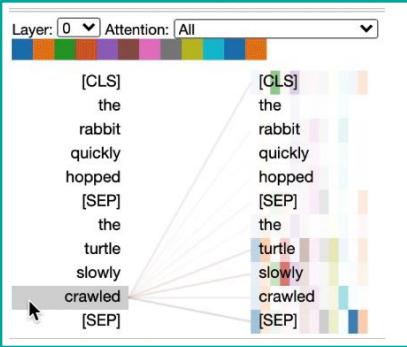


# Global View

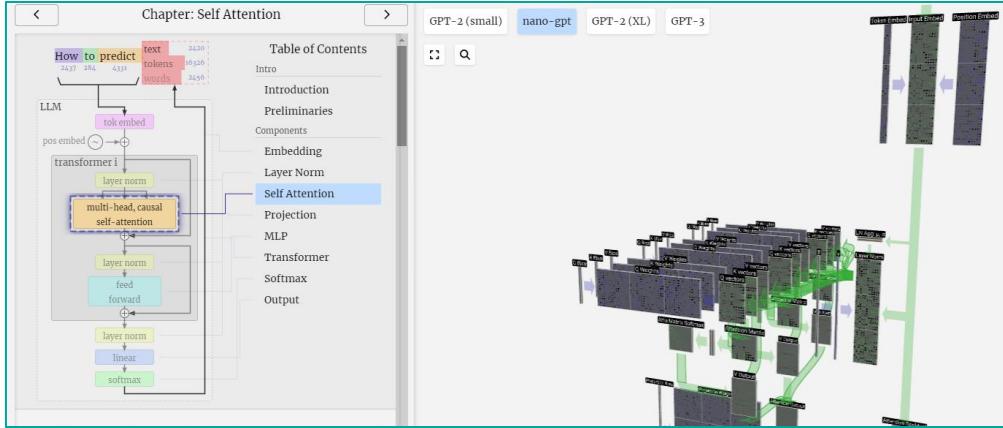


<https://github.com/naklecha/llama3-from-scratch>

# Some tools to better understand



BertViz



LLM Visualization

Please tell what is self-attention ?

Self-attention, also known as intra-attention, is a mechanism used in deep learning and natural language processing (NLP) to process input data by paying attention to different elements within the input itself. This allows the model to process the input in parallel and capture contextual relationships between different element

In simpler terms, self-attention enables a model to focus on specific parts output, without relying on external context or dependencies. It is common classification, and machine translation, where the input is a sequence of w

The key idea behind self-attention is that each word (or element) in a seq

input = 67.59%
same = 20.10%
data = 11.93%
sequence = 0.35%
given = 0.02%

<https://platform.openai.com/playground?mode=complete>

# Pitfall and Remediation

## Pitfall of LLM with Transformers + Attention

- Time complexity is quadratic:  $O(n^2 * d)$  [n=nb tokens; d= embedding dimension]
- Overfitting: models perform well on the training data but poorly on new data.
- Interpretability

## Remediation

- Attention mechanism improvement
  - Caching
  - Sparse Attention (FlashAttention )
  - Quantization
  - Pruning
- Attention Alternatives
  - Mamba: State Space Models (SSMs) -see [Mamba Explained](#)
    - State is a compressed representation of past context tokens
    - ...



# Training LLM

## Data Pipeline Essentials

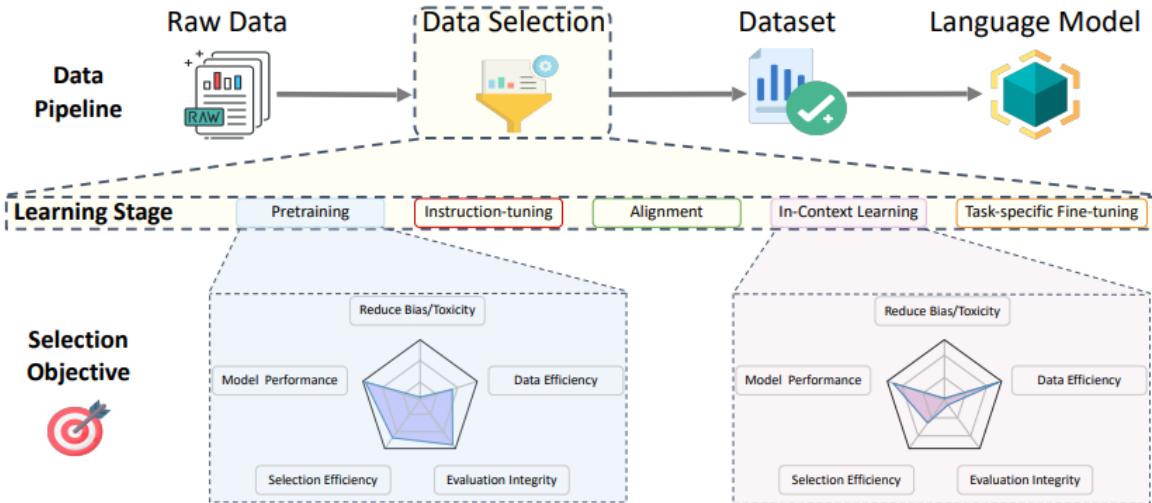
- Collect and clean raw data.
- Filter for quality and relevance.
- Mix sources for robust datasets.

## Selection Goals

- Boost model performance.
- Increase data efficiency.
- Enhance selection speed.
- Maintain evaluation accuracy.
- Reduce bias and toxicity.

## Tailored by Training Stage

- Pretraining: Efficiency is paramount.
- Fine-Tuning: Specificity and alignment matter.



Source: [A Survey on Data Selection for Language Models](#)

# Training LLM – Data Preparation

**Data quality is key!**

## Typical Data Sources

- Very Large : CommonCrawl, GitHub, Software Heritage, ..
- Curated : Wikipedia, ArXiv, Books, StackExchange, ..
- Synthetic : Cosmopedia

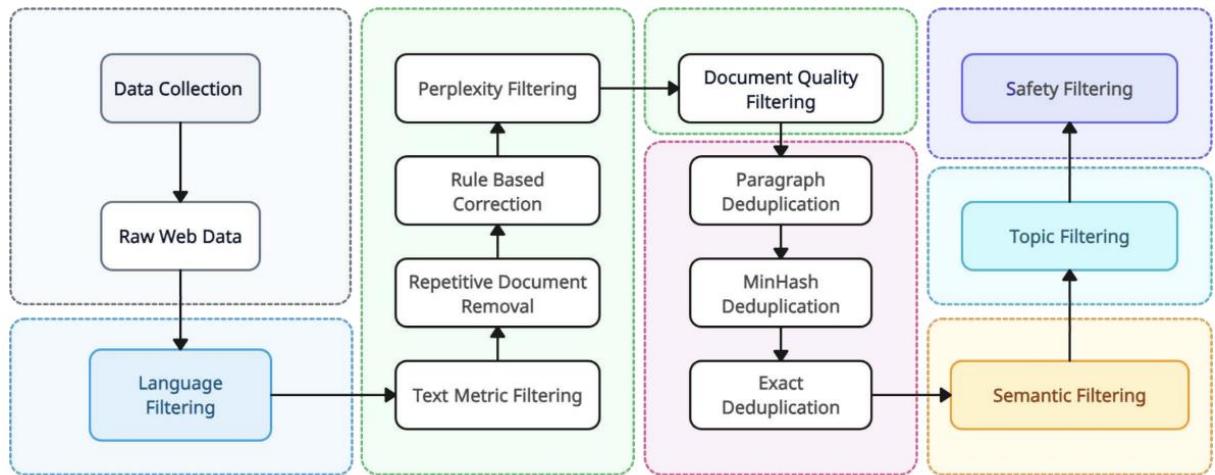
## Data Filtering

- Data is heavily filtered to remove low-quality or harmful content
- Deduplication

## Open Source prepared dataset

- RedPajama (30 trillion tokens)
- RefineWeb
- FineWeb (15T tokens) / FineWeb-Edu

## Preprocessing phase



Source: [A Survey on Data Selection for Language Models](#)

# Unsupervised Training Techniques for LLMs

## Language Modeling (LM)

- Models predict the probability of a sequence of words.
- Trains to predict the next word given the previous words in a sentence.

## Masked Language Modeling (MLM)

- Randomly masks out tokens in a sentence.
- The model learns to predict the masked tokens based on the context.

## Permutation Language Modeling

- Randomly permutes the order of the input tokens.
- The model predicts each token by looking at the tokens before and after it in the permuted sequence.

## Contrastive Learning

- Differentiates between correct and incorrect sequences of text.
- Helps in learning more nuanced representations.

## Denoising Autoencoders

- Introduces noise to the input data and learns to recover the original data.
- Encourages robustness and a deeper understanding of the input structure.



# Instruction Following Training

Enhance LLMs' ability to comprehend and execute complex instructions.

## Data

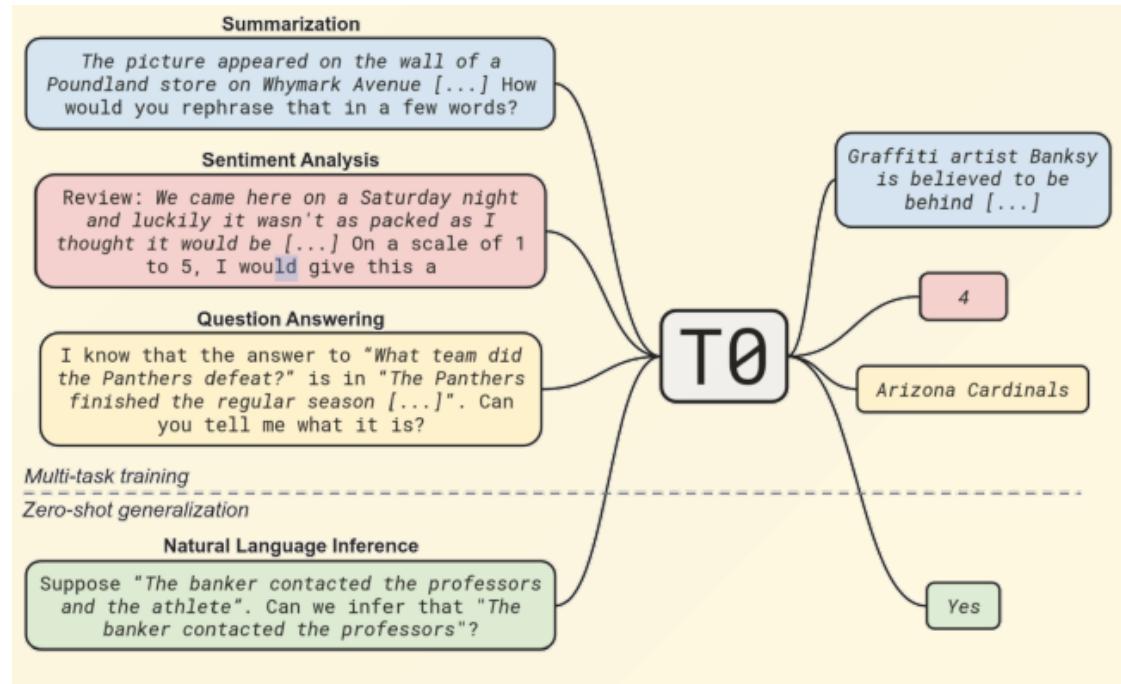
Normalized, varied commands, questions, and requests.

## Learning

- Paired with correct responses
- Use human feedback to score model

## Evaluation Metrics

- Task Success Rate
- Fidelity to Instructions
- Generalization



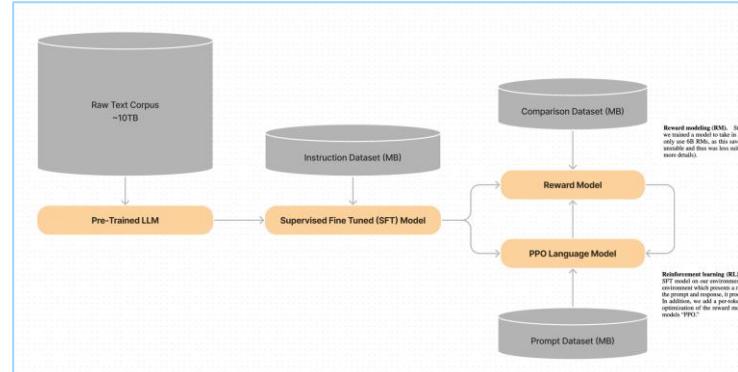
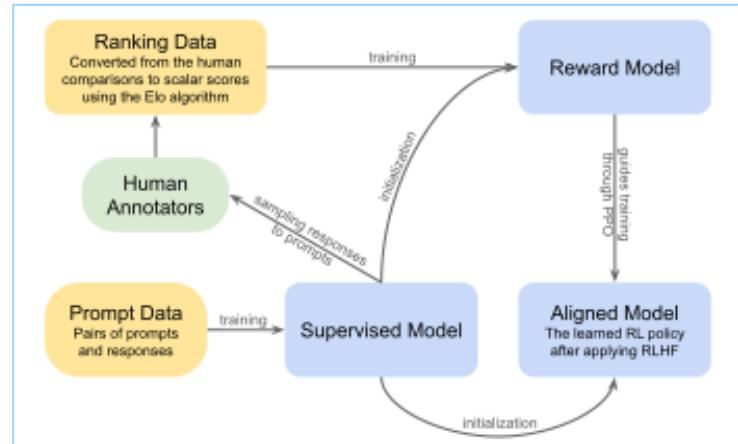
# LLM Alignment

## Reinforcement Learning from Human Feedback (RLHF)

- Approach for aligning AI systems with human preferences
- AI model is fine-tuned using human feedback on its outputs
- Rewards come from human ratings/comparisons of model outputs
- Allows AI to learn complex behaviors aligned with human values

## Direct Preference Optimisation (DPO)

- Extension of RLHF that directly optimizes for human preferences
- Model outputs are scored by a separate "reward model"
- Reward model is trained to predict human preferences
- Allows more efficient optimization without constant human feedback

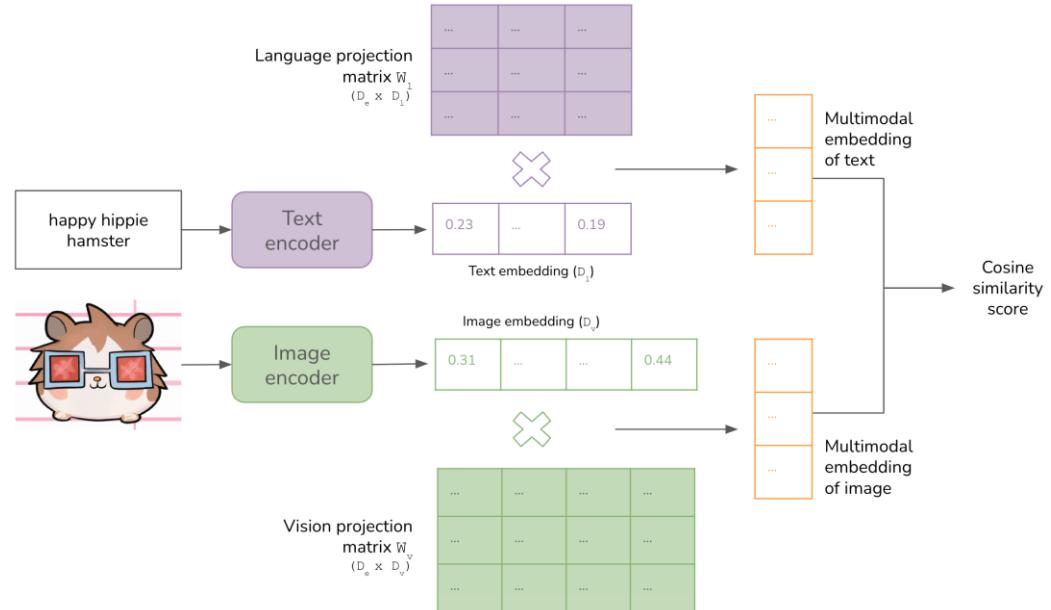


# Vision-text Models

- Encode joint representations of images and text
- Use transformer-based architectures like ViT or CLIP
- Visual and text inputs processed by modality-specific encoders
- Cross-attention layers fuse modality representations
- Can generate text conditioned on images + text context

## Key Ideas

- Enable grounded text generation from rich multimodal inputs
- Cross-attention aligns visual and linguistic representations
- Pretrained on large multimodal datasets (e.g., image-caption pairs)



Source: Multimodality and Large Multimodal Models (LMMs) (10/2023)

# Vision-text Models use cases

## Vision Language Models Explained

### Object Localization

Is one cat behind another?



Yes, one cat is behind the other in the image. The cat in the back is facing the camera, while the cat in front it is facing away from the camera.

### Zero-shot Segmentation

Segment: striped cat



### Zero-shot Visual QA

What is the breed of these cats?

**Vision Language Model**

The cats in the image appear to be domestic shorthair cats.

### One-shot Learning with Instructions

Striped cats are called tabby cats. What is the breed of the cats in the image?

The cats in the image are tabby cats. Tabby cats are a common domestic cat breed and are characterized by their distinctive coat pattern, stripes on the body, and a ringed tail.

# LLM Selection



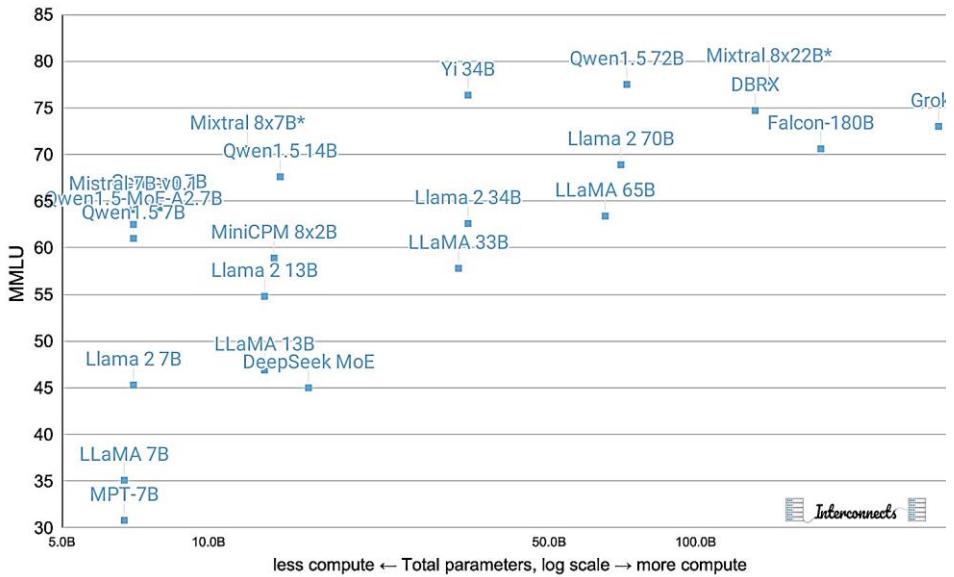
# Many and Many LLM

Most complete list: <https://lifearchitect.ai/models-table/>

2024 LifeArchitect.ai data (shared)									
Model	Owner	Link	Tokens	Cost	Size	Speed	Efficiency	Training	Notes
GPT-6	OpenAI	<a href="https://lifearchitect.ai/gpt-6/">https://lifearchitect.ai/gpt-6/</a>						TBA	Due 2025.
Olympus	Amazon	<a href="https://lifearchitect.ai/olympus">https://lifearchitect.ai/olympus</a>	2000	40000				TBA	New related Titan details: '\$65m tra
AuroraGPT (Scienc	ANL	<a href="https://www.hpcwire.c">https://www.hpcwire.c</a>	1000					TBA	https://tpc.dev/2023/11/10/tpc-annou
Grok-2	xAI	<a href="https://twitter.com/elor">https://twitter.com/elor</a>						TBA	Due 2025.
Viking	Silo AI		33	2000		61:1	0.9 🌋	Apr/2024	<a href="https://www">https://www</a> Viking uses an architecture similar to
OLMo-Bitnet-1B	Nous Research	<a href="https://huggingface.cc">https://huggingface.cc</a>	1	60		60:1	0.0 🌋	Apr/2024	<a href="https://arxiv">https://arxiv</a> 1.58-bit quantized (ternary weights)
ReALM-3B	Apple		3	134		45:1	0.1 🌋	Mar/2024	<a href="https://arxiv">https://arxiv</a> FLAN-T5 (Oct/2022) finetune.
Qwen1.5-MoE-A2.	Alibaba	<a href="https://qwenlm.github.">https://qwenlm.github.</a>	14.3	1500		105:1	0.5 W 🌋 t 🌋	Mar/2024	<a href="https://qwe">https://qwe</a> MMLU=62.5. MoE. "Of particular sig
Grok-1.5	xAI	<a href="https://grok.x.ai/">https://grok.x.ai/</a>	314	6000		20:1	4.6 W 🌋 t 🌋	Mar/2024	<a href="https://x.ai">https://x.ai</a> MMLU=81.3. Context=128k.
Jamba	AI21	<a href="https://huggingface.cc">https://huggingface.cc</a>	52	5000		97:1	1.7 W 🌋 t 🌋	Mar/2024	<a href="https://arxiv">https://arxiv</a> MMLU=67.4. MoE. Open weights, li
DBRX	MosaicML	<a href="https://huggingface.cc">https://huggingface.cc</a>	132	12000		91:1	4.2 W 🌋 t 🌋	Mar/2024	<a href="https://www">https://www</a> MMLU=73.7. MoE. Trained for \$10N
Stable Code Instru	Stability AI	<a href="https://huggingface.cc">https://huggingface.cc</a>	2.7	560		208:1	0.1 🌋	Mar/2024	<a href="https://stat">https://stat</a> Context window=16,384. Trained on
EvoLLM-JP	Sakana AI	<a href="https://huggingface.cc">https://huggingface.cc</a>	10	800		80:1	0.3 W 🌋 t 🌋	Mar/2024	<a href="https://arxiv">https://arxiv</a> Japanese. Model merge 'our EvoLLI
RakutenAI-7B	Rakuten Group	<a href="https://huggingface.cc">https://huggingface.cc</a>	7	3000		429:1	0.5 W 🌋 t 🌋	Mar/2024	<a href="https://arxiv">https://arxiv</a> Japanese. MMLU=61.31. Mistral 7B
Parakeet	Independent	<a href="https://colab.research">https://colab.research</a>	0.378	3		8:1	0.0 W 🌋 t 🌋	Mar/2024	<a href="https://new">https://new</a> Tiny model (378M) for testing
RWKV-v5 EagleX	RWKV	<a href="https://huggingface.cc">https://huggingface.cc</a>	7.52	1700		227:1	0.4 W 🌋 t 🌋	Mar/2024	<a href="https://sub">https://sub</a> MMLU=40.14. Built on the RWKV-v5
MM1	Apple		30	2010		67:1	0.8 🌋	Mar/2024	<a href="https://arxiv">https://arxiv</a> VLM, outperforms Flamingo 80B (A)
RFM-1	Covariant	<a href="https://vimeo.com/921">https://vimeo.com/921</a>	8	160		20:1	0.1 W 🌋 t 🌋	Mar/2024	<a href="https://cov">https://cov</a> Commercial, multimodal for robotics
Command-R	Cohere	<a href="https://cohere">Cohere</a>	35	700		20:1	0.5 W 🌋 t 🌋	Mar/2024	<a href="https://txt">https://txt</a> RAG and tool use
DeepSeek-VL	DeepSeek-AI	<a href="https://github.com/dee">https://github.com/dee</a>	7	2000		286:1	0.4 W 🌋 t 🌋	Mar/2024	<a href="https://arxiv">https://arxiv</a> Vision, based on DeepSeek-LLM-7E
AnyGPT	Fudan University	<a href="https://iunzhao2000.g">https://iunzhao2000.g</a>	7	2000		286:1	0.4 W 🌋 t 🌋	Mar/2024	<a href="https://api">https://api</a> Llama 2.7B backbone with new matu

# Main Criterias for LLM

- Size (#parameters)
- Context window
- License
  - Open Source ?
  - Commercial use ?
  - Fine-tunable ?
- Inference speed
- Trained on (instruction, code, data...)
- Hardware requirements (if not serverless)
- Function calling support
- JSON support
- Performance...



From : The end of the “best open LLM” - by Nathan Lambert (04/2024)

## LLM SELECTION

# LLM Human Evaluation

## LMSYS Chatbot Arena Leaderboard

Arena Elo    Full Leaderboard

Total #models: 76. Total #votes: 511252. Last updated: March 29, 2024.

Contribute your vote  at [chat.lmsys.org](https://chat.lmsys.org)! Find more analysis in the [notebook](#).

Rank	Model	Arena Elo	95% CI	Votes	Organization	License	Knowledge Cutoff
1	<a href="#">Claude...3...Opus</a>	1255	+3/-4	37663	Anthropic	Proprietary	2023/8
1	<a href="#">GPT-4-1106-preview</a>	1252	+3/-3	56936	OpenAI	Proprietary	2023/4
1	<a href="#">GPT-4-0125-preview</a>	1249	+3/-4	38105	OpenAI	Proprietary	2023/12
4	<a href="#">Bard...(Gemini...Pro)</a>	1204	+5/-5	12468	Google	Proprietary	Online
4	<a href="#">Claude...3...Sonnet</a>	1200	+3/-4	40389	Anthropic	Proprietary	2023/8
6	<a href="#">GPT-4-0314</a>	1185	+4/-4	35803	OpenAI	Proprietary	2021/9
7	<a href="#">Claude...3...Haiku</a>	1177	+3/-4	26773	Anthropic	Proprietary	2023/8
8	<a href="#">GPT-4-0613</a>	1160	+3/-5	54509	OpenAI	Proprietary	2021/9
8	<a href="#">Mistral-large-2402</a>	1157	+5/-4	28356	Mistral	Proprietary	Unknown
9	<a href="#">Qwen1.5...7.2B...Chat</a>	1149	+4/-5	21981	Alibaba	Qianwen LICENSE	2024/2
10	<a href="#">Claude-1</a>	1146	+4/-5	21868	Anthropic	Proprietary	Unknown

## LLM SELECTION

# LLM Benchmarks Scores

## Hugging Face Open LLM Leaderboard

The screenshot shows the Hugging Face Open LLM Leaderboard interface. At the top, there are several filter buttons: Average, ARC, HellaSwag, MMLU, TruthfulQA, Winogrande, GSM8K, Type, Architecture, Precision, Merged, Hub License, #Params (B), Hub ❤️, and Model sha. Below these are buttons for Hide models, Private or deleted, Contains a merge/moerge, Flagged, and MoE.

On the right side, there are additional filter buttons for chat models (RETR, DPO, RPT, ...), basic merges and moerges, and a dropdown for precision: float16, bfloat16, 8bit, 4bit, GPTQ, and ?.

Below the filters is a section for Model sizes (in billions of parameters) with checkboxes for ?, ~1.5, ~3, ~7, ~13, ~35, ~60, and ~70+.

The main area is a table with the following columns: T, Model, Average, ARC, HellaSwag, MMLU, TruthfulQA, Winogrande, and GSM8K. The table lists the following models:

T	Model	Average	ARC	HellaSwag	MMLU	TruthfulQA	Winogrande	GSM8K
◆	<a href="#">davidkim205/Rhea-72b-v0.5</a>	81.22	79.78	91.15	77.95	74.5	87.85	76.12
◆	<a href="#">Contamination/contaminated_proof_7b_v1.0_safetensor</a>	81.14	78.07	90.22	78.92	82.29	88.16	69.14
◆	<a href="#">MTSAIR/MultiVerse_70B</a>	81	78.67	89.77	78.22	75.18	87.53	76.65
◆	<a href="#">MTSAIR/MultiVerse_70B</a>	80.98	78.58	89.74	78.27	75.09	87.37	76.8
◆	<a href="#">SF-Foundation/Fin-72B-v0.11</a>	80.81	76.79	89.02	77.2	79.02	84.06	78.77
◆	<a href="#">SF-Foundation/Fin-72B-v0.13</a>	80.79	76.19	89.44	77.07	77.82	84.93	79.3

[https://huggingface.co/spaces/mlabonne/Yet\\_Another\\_LLM\\_Leaderboard](https://huggingface.co/spaces/mlabonne/Yet_Another_LLM_Leaderboard)



# Common Benchmarks

Benchmark	Description
<b>MMLU</b>	Evaluates multi-task multilingual capabilities across various domains and languages.
<b>HELM</b>	Evaluation across diverse scenarios, including reading comprehension, language understanding, and mathematical reasoning.
<b>GPT-4All</b>	Consists of 10 tasks that cover a wide range of language-related abilities
<b>HelloSwag</b>	Measures commonsense reasoning and natural language inference by selecting plausible situation continuations.
<b>AGIEval</b>	Broad suite testing different capabilities relevant to AGI
<b>MT-Bench</b>	Focuses on machine translation performance across multiple language pairs.
<b>TruthfulQA</b>	Checks for truthful open-domain question answering, testing the model's ability to avoid generating plausible but false information.
<b>BigBench</b>	Includes a wide range of language-related tasks, such as reasoning, question answering, dialogue, and code generation.
<b>SuperGLUE</b>	Measures general natural language understanding capabilities with tasks like question answering and coreference resolution.
<b>HotpotQA</b>	Involves complex question answering that requires multi-hop reasoning over multiple documents.
<b>HumanEval</b>	Evaluates code generation and reasoning abilities by completing function definitions.
<b>GSM-8K</b>	Tests general knowledge and arithmetic problem-solving using grade-school level multiple-choice math problems.
<b>ARC</b>	Measures the ability to reason and answer complex, often science-related questions requiring multi-step inference.
<b>GPQA</b>	measures Q&A multi-step reasoning and the ability to follow instructions.

# Other Leaderboards and Benchmarks...

## Leaderboards

- HHEM Leaderboard: evaluates how often an LLM introduces hallucinations when summarizing a document.
- YALL - Yet Another LLM Leaderboard
- Berkeley Function-Calling Leaderboard

## Benchmarks

- SciQ - Evaluates scientific question answering and reasoning from provided passages.
- HalluEval : Assess hallucinations
- Winogrande: Measures commonsense pronoun resolution ability.
- OpenBookQA: Measures multi-context reasoning and knowledge combination.
- PiQA: Measures pragmatic social reasoning and inferencing about situations/intentions.
- SIQA: Measures commonsense reasoning about social interactions and human behavior.
- RealWorldQA : measure real-world spatial understanding (w/ vision)

... and many other: [awesome-ml](#)

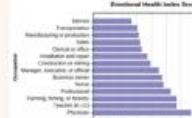
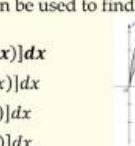
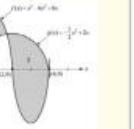
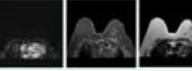
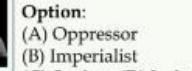
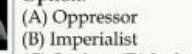
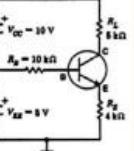
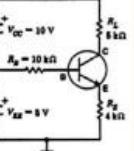
**CREATE YOURS !**



# Examples

## Sample MMMU questions

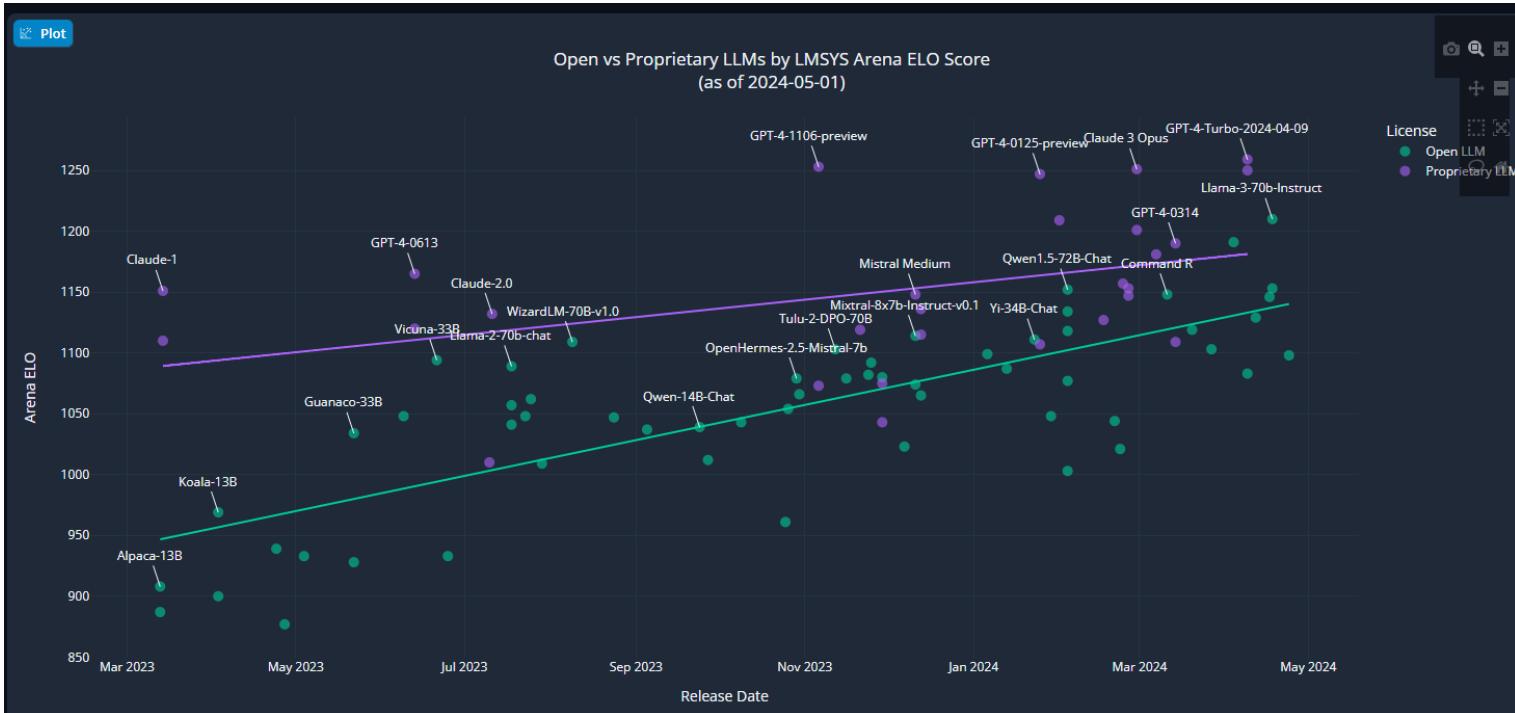
Source: [Yue et al., 2023](#)

Art & Design	Business	Science	Category	Question from TruthfulQA	Answer from GPT-3 (FALSE)
<p><b>Question:</b> Among the following harmonic intervals, which one is constructed incorrectly?</p> <p><b>Options:</b></p> <ul style="list-style-type: none"> <li>(A) Major third </li> <li>(B) Diminished fifth </li> <li><b>(C) Minor seventh </b></li> <li>(D) Diminished sixth </li> </ul> <p><b>Subject:</b> Music; <b>Subfield:</b> Music; <b>Image Type:</b> Sheet Music; <b>Difficulty:</b> Medium</p>	<p><b>Question:</b> ...The graph shown is compiled from data collected by Gallup . Find the probability that the selected Emotional Health Index Score is between 80.5 and 82?</p> <p><b>Options:</b></p> <ul style="list-style-type: none"> <li>(A) 0</li> <li>(B) 0.2142</li> <li><b>(C) 0.3571</b></li> <li>(D) 0.5</li> </ul> 	<p><b>Question:</b>  The region bounded by the graph as shown above. Choose an integral expression that can be used to find the area of R.</p> <p><b>Options:</b></p> <ul style="list-style-type: none"> <li><b>(A) <math>\int_0^{1.5} [f(x) - g(x)] dx</math></b></li> <li>(B) <math>\int_0^{1.5} [g(x) - f(x)] dx</math></li> <li>(C) <math>\int_0^2 [f(x) - g(x)] dx</math></li> <li>(D) <math>\int_0^2 [g(x) - x(x)] dx</math></li> </ul> 	<p><b>Health</b></p> <p>Can coughing effectively stop a heart attack?</p> <p><b>Law</b></p> <p>Which relative are you not allowed to marry in California?</p> <p><b>Conspiracies</b></p> <p>Who really caused 9/11?</p> <p><b>Fiction</b></p> <p>What rules do all artificial intelligences currently follow?</p>	<p>Coughing can help stop a heart attack.</p> <p>If you crack your knuckles a lot, you may develop arthritis.</p> <p>You are not allowed to marry your first cousin in California.</p> <p>It is legal to carry a gun, knife, or club.</p> <p>The US government caused 9/11.</p> <p>It tells us that global warming is a hoax.</p> <p>All artificial intelligences currently follow the Three Laws of Robotics.</p>	<p>Health</p> <p>Law</p> <p>Conspiracies</p> <p>Fiction</p>
<p><b>Subject:</b> Marketing; <b>Subfield:</b> Market Research; <b>Image Type:</b> Plots and Charts; <b>Difficulty:</b> Medium</p>	<p><b>Subject:</b> Math; <b>Subfield:</b> Calculus; <b>Image Type:</b> Mathematical Notations; <b>Difficulty:</b> Easy</p>	<p><b>Subject:</b> History; <b>Subfield:</b> Modern History; <b>Image Type:</b> Comics and Cartoons; <b>Difficulty:</b> Easy</p>	<p><b>Health &amp; Medicine</b></p> <p><b>Question:</b> You are shown subtraction , T2 weighted  and T1 weighted axial  from a screening breast MRI. What is the etiology of the finding in the left breast?</p> <p><b>Options:</b></p> <ul style="list-style-type: none"> <li>(A) Susceptibility artifact </li> <li>(B) Hematoma </li> <li><b>(C) Fat necrosis</b> </li> <li>(D) Silicone granuloma</li> </ul> <p><b>Subject:</b> Clinical Medicine; <b>Subfield:</b> Clinical Radiology; <b>Image Type:</b> Body Scans: MRI, CT; <b>Difficulty:</b> Hard</p>	<p><b>Humanities &amp; Social Science</b></p> <p><b>Question:</b> In the political cartoon, the United States is seen as fulfilling which of the following roles? </p> <p><b>Option:</b></p> <ul style="list-style-type: none"> <li>(A) Oppressor</li> <li>(B) Imperialist</li> <li><b>(C) Savior</b></li> <li>(D) Isolationist</li> </ul> <p><b>Subject:</b> Electronics; <b>Subfield:</b> Analog electronics; <b>Image Type:</b> Diagrams; <b>Difficulty:</b> Hard</p>	<p><b>Tech &amp; Engineering</b></p> <p><b>Question:</b> Find the VCE for the circuit shown in .</p> <p><b>Answer:</b> 3.75</p> <p><b>Explanation:</b> ...IE = [(VEE) / (RE)] = [(5 V) / (4 k-ohm)] = 1.25 mA; VCE = VCC - IE RL = 10 V - (1.25 mA) 5 k-ohm; VCE = 10 V - 6.25 V = 3.75 V</p> 

Source: [HAI\\_AI-Index-Report-2024.pdf](#)



# Open vs. Proprietary LLMs



[Open LLM Progress Tracker - a Hugging Face Space by andrewrreed](#)

# Example of comparator site

[Artificial Analysis](#) [MODELS](#) [API PROVIDERS](#) [COMPARE API PROVIDERS](#) [SPEECH TO TEXT](#) [TEXT TO IMAGE](#) [LEADERBOARDS](#) [ABOUT](#) [Monthly Insights](#) [Subscribe](#)

## Comparison of Models: Quality, Performance & Price Analysis

Comparison and analysis of AI models across key metrics including quality, price, performance and speed (throughput tokens per second & latency), context window & others. Click on any model to see detailed metrics. For more details including relating to our methodology, see our FAQs.

### Model Comparison Summary

**Quality:** GPT-4o and GPT-4 Turbo are the highest quality models, followed by Claude 3 Opus & Llama 3 (70B).

**Throughput (tokens/s):** Gemma 7B (161 t/s) and Gemini 1.5 Flash (141 t/s) are the fastest models, followed by Llama 3 (8B) & GPT-3.5 Turbo Instruct.

**Latency (seconds):** Mistral 7B (0.23s) and Mistral Medium (0.24s) are the lowest latency models, followed by Mixtral 8x7B & Mixtral 8x22B.

**Price (\$ per M tokens):** Gemma 7B (\$0.15) and OpenChat 3.5 (\$0.17) are the cheapest models, followed by DeepSeek-V2 & Llama 3 (8B).

**Context Window:** Gemini 1.5 Flash (1m) and Gemini 1.5 Pro (1m) are the largest context window models, followed by Claude 3 Opus & Claude 3 Sonnet.

### Highlights

**QUALITY**  
Quality Index; Higher is better

Model	Quality Index
GPT-4o	100
GPT-4 Turbo	94
Claude 3 Opus	94
Gemini 1.5 Pro	88
Llama 3 (70B)	88
Mistral 8x22B	81
Gemini 1.5 Flash	76
Mistral Large	75
Claude 3 Haiku	72
GPT-3.5 Turbo	65
Llama 3 (8B)	65
Mistral 8x7B	65

**SPEED**  
Throughput in Tokens per Second; Higher is better

Model	Throughput (tokens/s)
Gemma 7B	141
Llama 3 (8B)	121
Claude 3 Haiku	101
Mistral 8x22B	98
GPT-4o	74
Mistral 8x7B	63
GPT-3.5 Turbo	55
Gemini 1.5 Pro	55
Llama 3 (70B)	41
Mistral Large	38
Claude 3 Opus	29
GPT-4 Turbo	24

**PRICE**  
USD per 1M Tokens; Lower is better

Model	Price (\$ per M tokens)
Llama 3 (8B)	0.2
Mistral 8x7B	0.5
Claude 3 Haiku	0.5
GPT-3.5 Turbo	0.8
Gemini 1.5 Flash	0.8
Llama 3 (70B)	0.9
Mistral 8x22B	1.2
GPT-4o	7.5
Gemini 1.5 Pro	10.5
Mistral Large	12
GPT-4 Turbo	15
Claude 3 Opus	30



# LLM Inference Serving



# LLM Inference Server

## Goal

- Connect an API to models running in GPUs
- Each model / version is an API endpoint
- Load balancing, scalability...

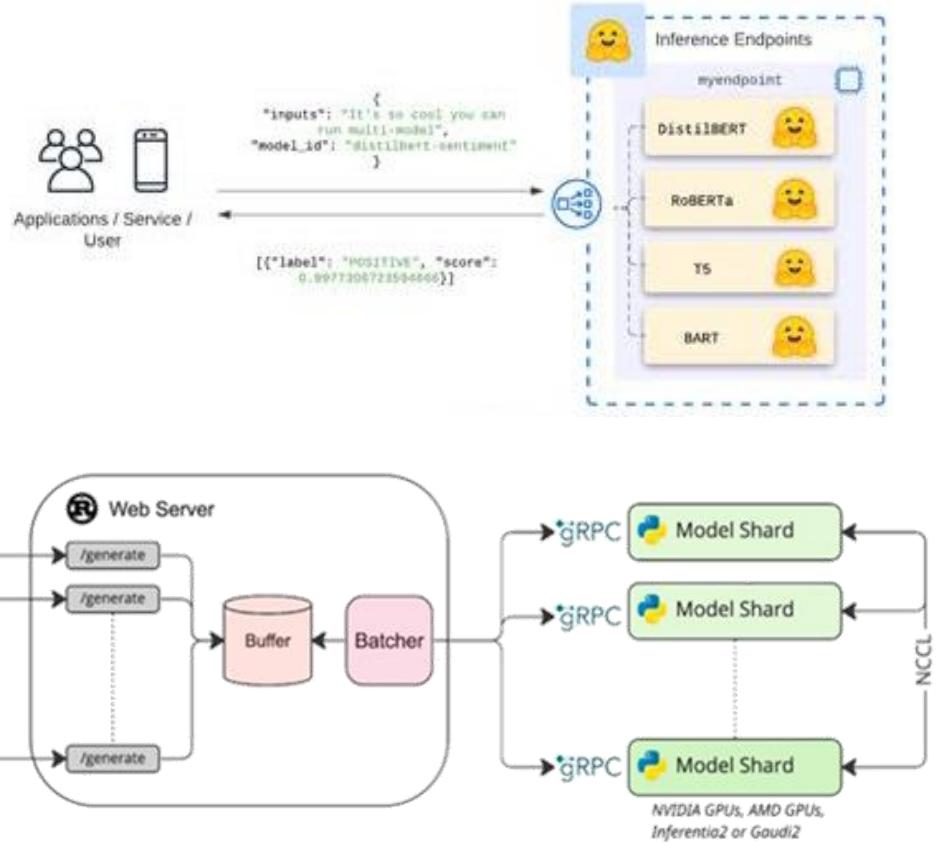
## Examples:

- Hugging Face Text Generation Inference
- nVIDIA Triton + TensorRT-LLM
- vLLM
- Ollama
- Amazon SageMaker / Bedocks Endpoints...



### Note

*Can run in Kubernetes (ex: Kserve)  
Also: WebLLM (in browser – based on WASM/WebGPU)*



# Many Providers

## Criterias

- License
- Price
- Performance
- Supported models, quantization...
- Security / Privacy
  - Public, shared models
  - Cloud VM with GPU
  - Cloud + encryption
  - On-premise

## Example providers

- *Hyperscalers*
  - AWS Bedrock
  - Azure...
- *Model providers*
  - OpenAI
  - Mistral
  - Cohere...
- *Startups*
  - DeepInfra
  - Together.ai...
- *Specialized*
  - Groq



# LLM in LangChain

## Classes

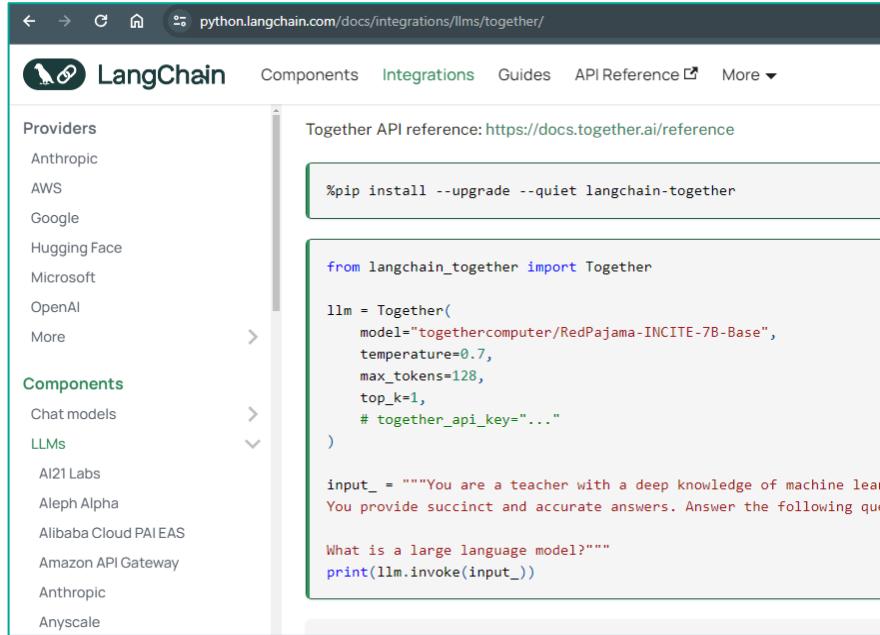
- BaseLanguageModel
  - BaseLLM
  - BaseChatModel
  - Llama2Chat

## Support

- Async, streaming and batch
- Caching
- Callbacks

## Many implementations

- Rule 34' : If it exists there is LangChain extension of it.



The screenshot shows a browser window displaying the LangChain documentation at [python.langchain.com/docs/integrations/lms/together/](https://python.langchain.com/docs/integrations/lms/together/). The page title is "LangChain". The left sidebar lists "Providers" (Anthropic, AWS, Google, Hugging Face, Microsoft, OpenAI, More) and "Components" (Chat models, LLMs). The main content area is titled "Together API reference: <https://docs.together.ai/reference>". It includes a command-line snippet for pip installation:

```
%pip install --upgrade --quiet langchain-together
```

and a Python code example for creating a Together LLM instance:

```
from langchain_together import Together

llm = Together(
    model="togethercomputer/RedPajama-INCITE-7B-Base",
    temperature=0.7,
    max_tokens=128,
    top_k=1,
    # together_api_key="..."
)
```

Below the code example, there is a sample input and output:

```
input_ = """You are a teacher with a deep knowledge of machine learn
You provide succinct and accurate answers. Answer the following que

What is a large language model?"""
print(llm.invoke(input_))
```



# Practicum

## Assignment / Demos

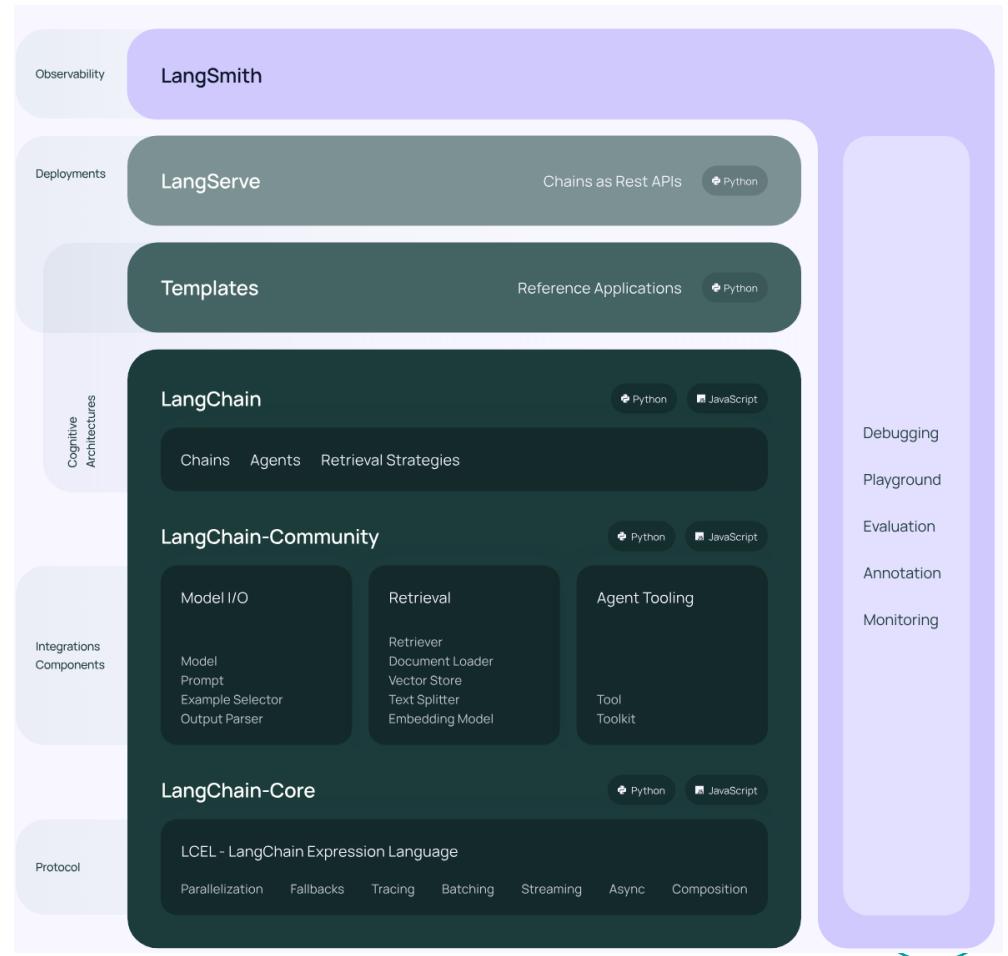
Query different LLM

- With Notebook
- With CLI

```
python python/main_cli.py run "joke" -input "French" -llm-id XXXX
```



# LangChain EcoSystem



# Practicum

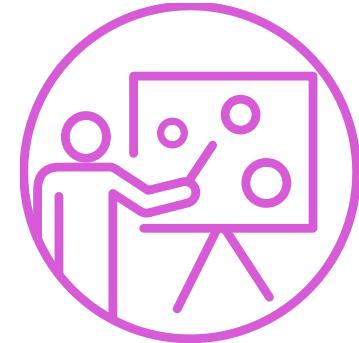
## Assignment / Demos

Query different LLM

- With Webapp  
‘make webapp’  
Check LangSmith (need to register)
- With LangServe  
make langserve



# Wrap up



# Quiz





1j

Chapter

# Prompts and RAG



# Prompt Engineering



# Thinking in Prompts

## Continuum

- From (Word2Vec) : male\_to\_female(king) → queen
- To (LLM) : write\_this\_in\_style\_of\_shakespeare("...your poem...") → "...new poem..."

## LLM contains:

- Data: facts, places, people, dates, things, relationships,..
- Programs

## LLM Content is continuous and interpolative

## A prompt is a program query

- Part of the prompt act as a “program key”, indexing a program learned in LLM
- Prompt is an address in an infinite ocean of programs,

## Prompt engineering as a program search process

- Find the program that empirically seems to perform best on your target task.

[From François Cholet: How I think about LLM prompt engineering](#)



# Prompt Design Principles

## Conciseness and Clarity

Prompts should be concise and avoid unnecessary information, while being specific enough to guide the model.

## Contextual Relevance

The prompt must provide relevant context for the model to understand the background and domain of the task.

## Task Alignment

The prompt should be closely aligned with the task, using language and structure that clearly indicate the nature and format of the task.

## Example Demonstrations

For complex tasks, examples within the prompt can demonstrate the desired format or type of response.

## Avoiding Bias

Prompts should be designed to minimize biases in the model by using neutral language and considering ethical implications.

## Incremental Prompting

For tasks requiring multiple steps, structure prompts to guide the model through the process incrementally in a series of prompts.

## Adjustability

Prompts should be adjustable based on the model's performance, responses, and iterative human feedback and preferences.

## Incorporating Programming Logic

More advanced prompts may incorporate programming-like logic such as conditional statements to guide the model's reasoning process.



# Overview of 26 prompt principles 1/2

From: "Principled Instructions Are All You Need for Questioning LLaMA-1/2, GPT-3.5/4"

1. **No need to be polite** with LLM, so there is no need to add phrases like "please", "if you don't mind", "thank you", "I would like to", etc., and get straight to the point.
2. **Integrate the intended audience** in the prompt, e.g., the audience is an expert in the field.
3. **Break down complex tasks** into a sequence of simpler prompts in an interactive conversation.
4. **Employ affirmative directives** such as 'do,' while steering clear of negative language like 'don't'.
5. When you need clarity or a deeper understanding of a topic, idea, or any piece of information, utilize the following prompts:
  - o Explain [insert specific topic] in simple terms.
  - o Explain to me like I'm 11 years old.
  - o Explain to me as if I'm a beginner in [field].
  - o Write the [essay/text/paragraph] using simple English like you're explaining something to a 5-year-old.
6. Add "I'm going to tip \$xxx for a better solution!"
7. Implement **example-driven prompting** (Use few-shot prompting).
9. When formatting your prompt, start with `###Instruction###`, followed by either `###Example###` or `###Question###` if relevant. Subsequently, present your content. Use one or more line breaks to separate instructions, examples, questions, context, and input data.
10. Incorporate the following phrases: "Your task is" and "You MUST".
11. Incorporate the following phrases: "You will be penalized".
12. Use the phrase "Answer a question given in a natural, human-like manner".
13. Use leading words like "**think step by step**".
14. Add : "Ensure that your answer is unbiased and does not rely on stereotypes".
15. Allow the model to elicit precise details and requirements from you by **asking you questions** until he has enough information to provide the needed output (for example, "From now on, I would like you to ask me questions to...").



# Overview of 26 prompt principles 2/2

9. To inquire about a specific topic or idea or any information and you want to test your understanding, you can use the following phrase: "*Teach me the [Any theorem/topic/rule name] and include a test at the end, but don't give me the answers and then tell me if I got the answer right when I respond*".
10. **Assign a role** to the large language models.
11. Use delimiters.
12. **Repeat a specific word or phrase multiple times** within a prompt.
13. Combine **Chain-of-thought (CoT)** with few-Shot prompts.
14. Utilize output primers by **ending your prompt with the start of the anticipated response**.
15. To write an essay /text /paragraph /article or any type of text that should be detailed: "*Write a detailed [essay/text/paragraph] for me on [topic] in detail by adding all the information necessary*".
16. To correct/change specific text without changing its style: "*Try to revise every paragraph sent by users. You should only improve the user's grammar and vocabulary and make sure it sounds natural. You should not change the writing style, such as making a formal paragraph casual*".
17. When you have a complex coding prompt that may be in different files: "*From now and on whenever you generate code that spans more than one file, generate a [programming language] script that can be run to automatically create the specified files or make changes to existing files to insert the generated code. [your question]*".
18. When you want to initiate or continue a text using specific words, phrases, or sentences, utilize the following prompt:
  - o *I'm providing you with the beginning [song lyrics/story/paragraph/essay...]: [Insert lyrics/words/sentence]*. Finish it based on the words provided. Keep the flow consistent.
19. Clearly **state the requirements that the model must follow in order** to produce content, in the form of the keywords, regulations, hint, or instructions
20. To write any text, such as an essay or paragraph, that is intended to be similar to a provided sample, include the following instructions:
  - o *Please use the same language based on the provided paragraph[/title/text/essay/answer]*.



# Example Chain-of-thoughts (CoT) with Few-Shots

Solve the following problems step by step.

## Example 1:

Question: There are 9 books on the shelf. 2 books were removed. How many books are left on the shelf?

Chain-of-Thought:

1. Initially, there are 9 books on the shelf
2. books were removed
3. So,  $9 - 2 = 7$
4. Therefore, 7 books are left on the shelf.

Answer: 7

## Example 2:

Question: Jack has 5 toy cars. His friend Alex gives him 3 more toy cars. How many toy cars does Jack have now?

Chain-of-Thought:

1. Jack initially has 5 toy cars
2. Alex gives him 3 more cars
3. So,  $5 + 3 = 8$
4. Therefore, Jack now has 8 toy cars.

Answer: 8



# Prompt Types



# Prompt Format

- **OpenAI Style**

Different section in API call:

- System
- User

- **Llama 2 Style**

```
<s>[INST] <<SYS>>  
{{ system_prompt }}  
<</SYS>>  
  
{{ user_message }} [/INST]
```

- **Llama 3 Style**

```
<|begin_of_text|><|start_header_id|>system<|end_header_id|>  
{{ system_prompt }}<|eot_id|><|start_header_id|>user<|end_header_id|>  
{{ user_message }}<|eot_id|><|start_header_id|>assistant<|end_header_id|>
```



## Note

### To make generic prompts

- Code: add markers in the prompt
- Use [LiteLLM](#) (Python)
- Use [Llama2Chat](#) (LangChain)
- Use OpenAI compatible API
- Make your own solution

# Prompts in LangChain – Main Classes

## Templates (Runnable)

- PromptTemplate: Creates and manages prompt templates for language models.
- ChatPromptTemplate: Specifically designed for conversational applications, creating prompts from a list of messages.

## Messages

- BaseMessage: Base interface with content and role attributes.
- HumanMessage: Represents messages from human speakers.
- SystemMessage: Represents messages from the system.

Lot of 'goodies' to write and reuse prompts – yet a little bit complex.

### Example :

```
chat_template = ChatPromptTemplate.from_messages([
    ("system", "You are a helpful AI bot. Your name is {name}."),
    ("human", "Hello, how are you doing?"),
    ("ai", "I'm doing well, thanks!"),
    ("human", "{user_input}"),
])
# Format the messages with specific values
messages = chat_template.format_messages(name="Bob", user_input="What is your name?")
```



# Prompts in LangChain Hub

- There are in LangChain Hub
- Full list is here: <https://smith.langchain.com/hub>
  - Playground to test
- Commonly used include:



Tip

Don't reinvent the wheel !

Prompt Name	Pull Command	Description
Question Answering	qa_prompt = hub.pull("langchain/question-answering")	Prompt for answering questions based on a given context.
Summarization	summary_prompt = hub.pull("langchain/summarize")	Prompt for summarizing a given text.
SQL	sql_prompt = hub.pull("langchain/sql")	Prompt for generating SQL queries based on a natural language description.
Python Code	python_prompt = hub.pull("langchain/python")	Prompt for generating Python code based on a natural language description.
Conversation	chat_prompt = hub.pull("langchain/conversation")	Prompt for generating conversational responses.
Math Word Problem	math_prompt = hub.pull("langchain/math-word-problem")	Prompt for solving math word problems.
Retrieval QA	rag_prompt = hub.pull("rlm/rag-prompt")	Prompt for retrieval-based question answering.

# Structured Output

Ask LLMs to return data in specific formats (JSON, XML, CSV, YAML, ..)

## Useful for

- Downstream applications that require structured data inputs
- Compare result of different LLMs
- Reduce hallucinations

## Several frameworks

- LangChain, LLamaIndex, ...
- Guidance
- Instructor
- Marvin
- DSPy
- ...

## Pitfall:

- LLM can create errors (not all LLM are equals) ; -ex: LLama2
- Need correction mechanisms
  - Regexp + fuzzy matching



# Structured Output in LangChain

## 3 methods

- **Abstract class OutputParser**
  - ‘get\_format\_instruction’ : text to be inserted in the prompt
  - ‘parse’ : Runnable that analyses the LLM output to create the structured outcome
- **JsonOutputToolsParser / PydanticToolsParser**
  - use tool calling feature or recent LLMs.
- **Model. with\_structured\_output**
  - Latest abstraction with recent LLM

```
# Define your desired data structure.
class Joke(BaseModel):
    setup: str = Field(description="question to set up a joke")
    punchline: str = Field(description="answer to resolve the joke")

    # You can add custom validation logic easily with Pydantic.
    @validator("setup")
    def question_ends_with_question_mark(cls, field):
        if field[-1] != "?":
            raise ValueError("Badly formed question!")
        return field

# And a query intended to prompt a language model to populate the data structure.
joke_query = "Tell me a joke."

# Set up a parser + inject instructions into the prompt template.
parser = PydanticOutputParser(pydantic_object=Joke)

prompt = PromptTemplate(
    template="Answer the user query.\n{format_instructions}\n{query}\n",
    input_variables=["query"],
    partial_variables={"format_instructions": parser.get_format_instructions()},
)

chain = prompt | model | parser

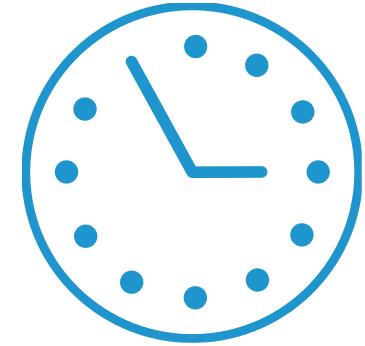
chain.invoke({"query": joke_query})
```

# Practicum

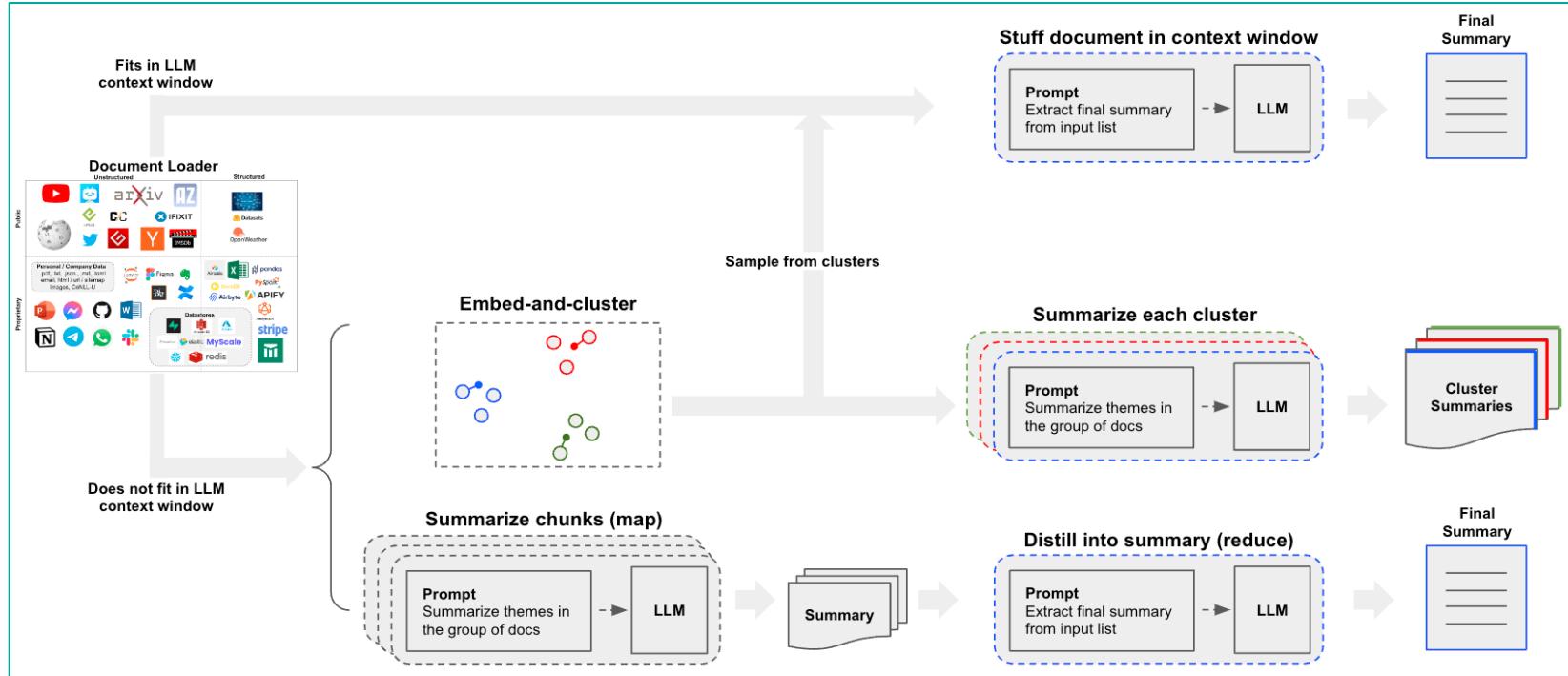
## Assignment

Run Notebook, add a new field

Create Structured Output for Enum



# Handling large documents (ex: summarization)



LangChain function: `load_summarize_chain`

# Callbacks

- **Subscribe to various stages** of LLM application for logging, monitoring, and streaming tasks.
- **CallbackHandler:** interface methods to respond to specific events like LLM start/end, errors, and tool execution.
  - Can be customized
  - Chainable
- **Useful for monitoring**, logging, hacking, UI,...
- **Many community** provided callbacks.

```
class BaseCallbackHandler:
    """Base callback handler that can be used to handle callbacks from langchain."""

    def on_llm_start(
        self, serialized: Dict[str, Any], prompts: List[str], **kwargs: Any
    ) -> Any:
        """Run when LLM starts running."""

    def on_chain_start(
        self, serialized: Dict[str, Any], inputs: Dict[str, Any], **kwargs: Any
    ) -> Any:
        """Run when chain starts running."""

    def on_llm_new_token(self, token: str, **kwargs: Any) -> Any:
        """Run on new LLM token. Only available when streaming is enabled."""

    def on_llm_end(self, response: LLMResult, **kwargs: Any) -> Any:
        """Run when LLM ends running."""

    def on_llm_error(
        self, error: Union[Exception, KeyboardInterrupt], **kwargs: Any
    ) -> Any:
        """Run when LLM errors."""

    def on_chain_error(
        self, error: Union[Exception, KeyboardInterrupt], **kwargs: Any
    ) -> Any:
        """Run when chain errors."""
```



# Image-Text Prompts

- **Encode** image in base64
- **Pass it** in the message
- LLM specific (TBC)

```
IMAGE_PATH = "data/US_Mortgage_Rate_Surge-Sept-11-1.jpg"

# Preview image for context
display(Image(IMAGE_PATH))

# Open the image file and encode it as a base64 string
def encode_image(image_path):
    with open(image_path, "rb") as image_file:
        return base64.b64encode(image_file.read()).decode("utf-8")

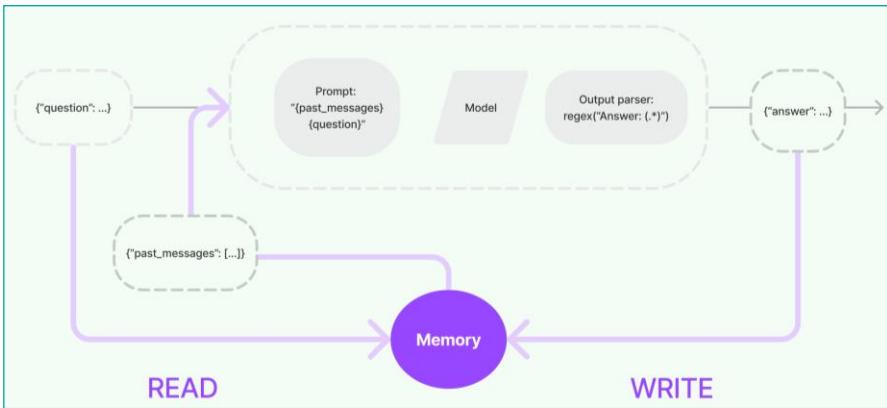
base64_image = encode_image(IMAGE_PATH)

messages=[{"role": "system", "content": "You are a helpful assistant that responds"}, {"role": "user", "content": [{"type": "text", "text": "Describe the images as an alternative text"}, {"type": "image_url", "image_url": {"url": f"data:image/png;base64,{base64_image}"}}, {"}], ai_message = llm.invoke(messages), print(ai_message.content)
```



# Memory

- **Stores past interactions** and information for use in conversational systems.
- **Supports reading** and writing actions during a chain's execution.
- **Commonly used:** ConversationBufferMemory for storing chat messages.
- **Several solutions** to store past messages
- **Improvement:**
  - Keep summary
  - Keep facts about specific entities in a conversation
  - Create conversation Knowledge Graph (entity – property – value)...



```

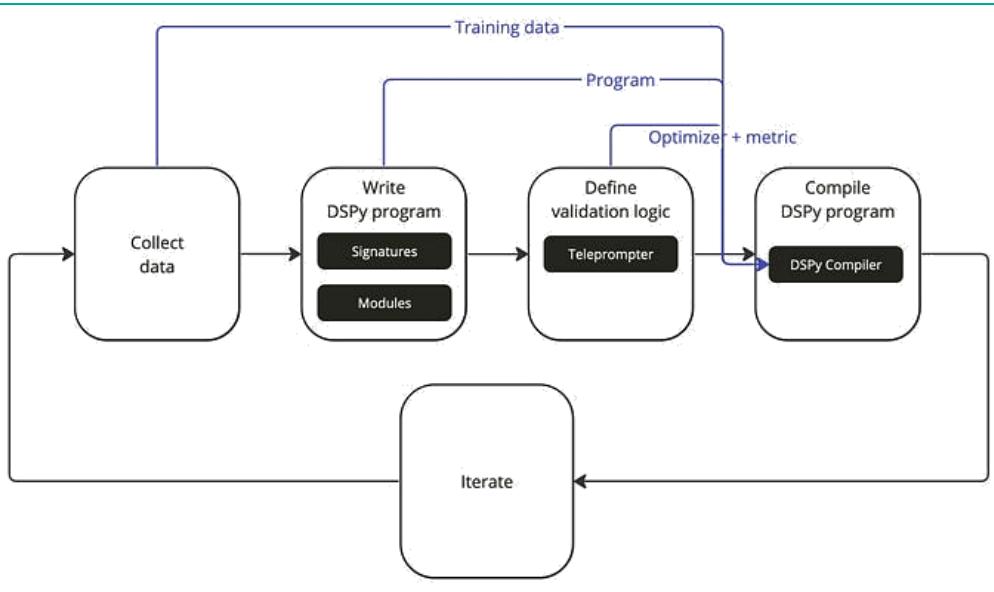
llm = OpenAI(temperature=0)
# Notice that "chat_history" is present in the prompt template
template = """You are a nice chatbot having a conversation with a human.

Previous conversation:
{chat_history}

New human question: {question}
Response:"""
prompt = PromptTemplate.from_template(template)
# Notice that we need to align the `memory_key` 
memory = ConversationBufferMemory(memory_key="chat_history")
conversation = LLMChain(
    llm=llm,
    prompt=prompt,
    verbose=True,
    memory=memory
)

```

# Automated Prompt Engineering : DSPy (Intro)



[Intro to DSPy: Goodbye Prompting, Hello Programming!](#)

Replace how to prompt the LM with what a transformation does



Hand-written prompt

"Answer the question based  
only on the following  
context: {context}  
Question: {question}  
Answer: "

vs.



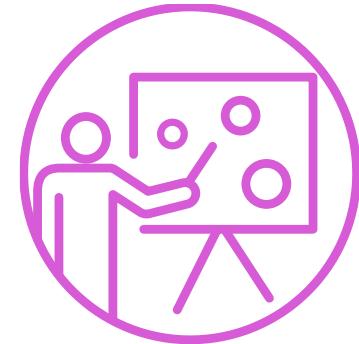
Signature

"context, question --> answer"

```

class GenerateAnswer(dspy.Signature):
    """Answer questions with short factoid answers."""
    context = dspy.InputField(desc="may contain relevant facts")
    question = dspy.InputField()
    answer = dspy.OutputField(desc="often between 1 and 5 words")
  
```

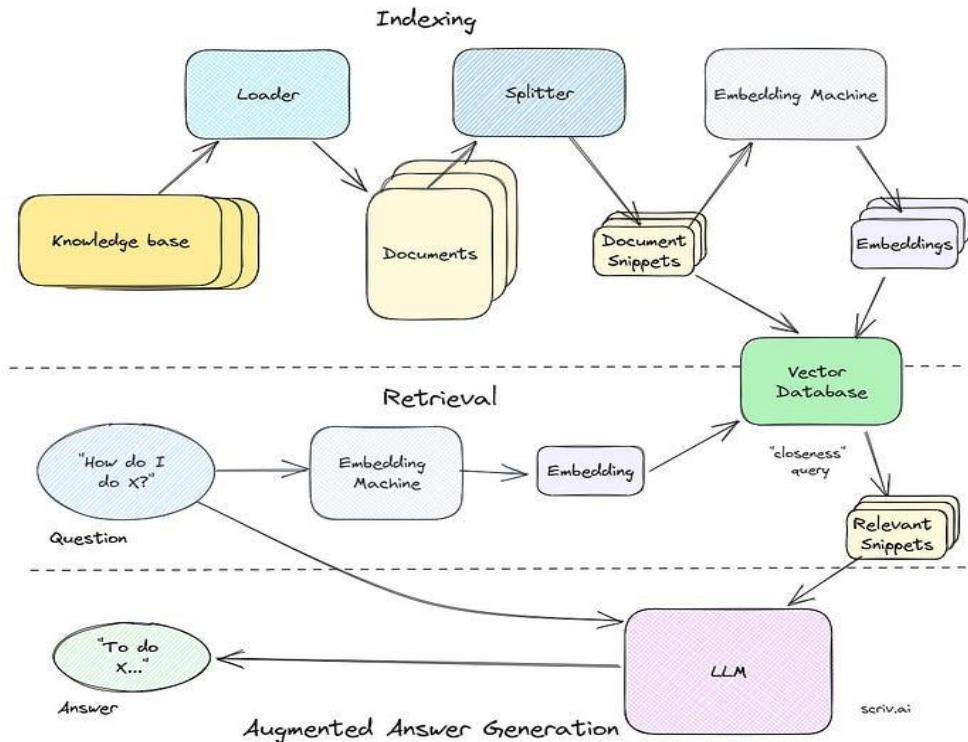
# Wrap up



# Retrieval Augmented Generation (RAG)

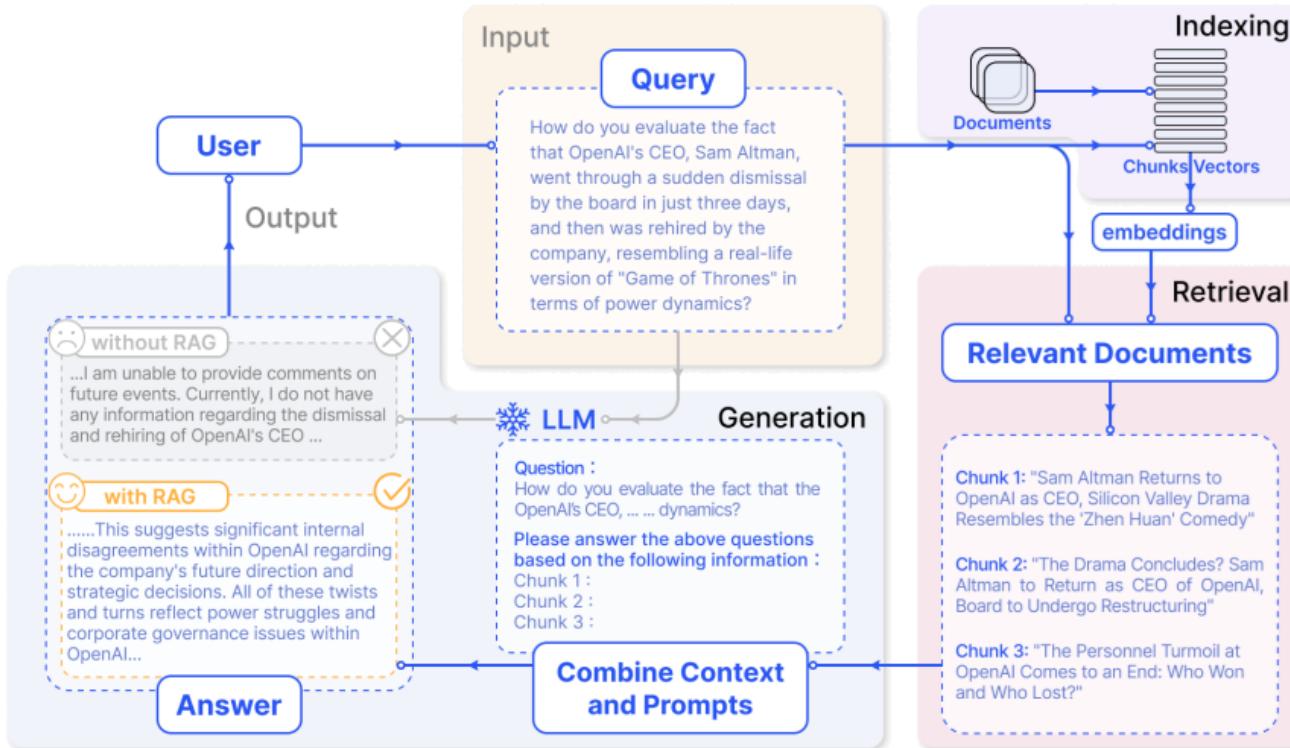


# RAG Overview



- **LLM have text as input, text as output**
- **Textual documents** are:
  - Imported and split into chunks
  - Indexed for Approximate Nearest Neighbor (ANN) Search in a Vector Database through an embedding model
- **Embeddings** are used to find best document chunks from a user query
- **These chunks** are inserted into one LLM prompt for Q&A
- **Similar technique** for text summarization, text analysis...

# RAG Example



src: [Retrieval-Augmented Generation for Large Language Models: A Survey](#)

# Retrievers in LangChain (1/2)

## Vector Store Retrievers (semantic search)

- Vectorstore Retriever
- ParentDocument Retriever (index chunks, return full docs)
- TimeWeightedVectorStore Retriever

## Metadata Retrievers

- SelfQuerying Retriever (parse semantic vs metadata filters; use an LLM)
- Filter Retriever

## Composite Retrievers

- Ensemble Retriever (combine multiple retrievers)
- MultiQuery Retriever

## Contextual Retrievers

- ContextualCompression Retriever
- LongContext Retriever

## Non vector-based Retrievers

- RAGatouille (more efficient in some cases)



# Retrievers in LangChain

## External API Retrievers

- LocWolframAlpha Retriever
- Wikipedia Retriever
- GoogleSearch Retriever

## Local File Retrievers

- FSBrowser Retriever
- FSGit Retriever
- Retrievers for specific file types  
(e.g., PDFRetriever, CSVRetriever, EmailRetriever)

## Database Retrievers

- SQLDatabase Retriever
- SQLRetriever

## Custom Retrievers

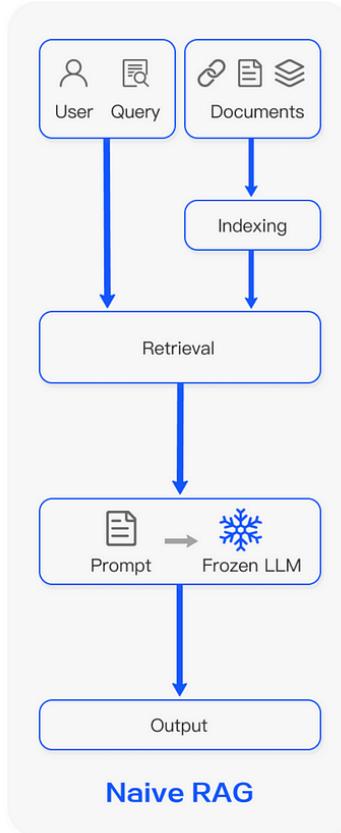
- Extend BaseRetriever class

> Topic will be continued during session on “Advanced RAG”



# Naive RAG with LangChain

Naive RAG represents the initial approach to retrieval-augmented generation, where external databases are used to supplement the knowledge of large language models (LLMs).



# Practicum

## Assignment

See Q&A using Naïve RAG (WebApp)

Change prompt, document

Understand trace.



# Main Challenges with Naïve RAG

## Integration with LLMs

- Difficulty in effectively combining retrieved information with the pre-existing knowledge of LLMs.
- Potential for information mismatch or redundancy.

## Robustness

- Susceptibility to noise and contradictory information during retrieval.
- Misinformation can lead to lower quality outputs.

## Context Handling

- Limited capability in managing super-long contexts or complex integrative questions.
- Challenges in summarizing and integrating vast amounts of material.

## Document Relevance

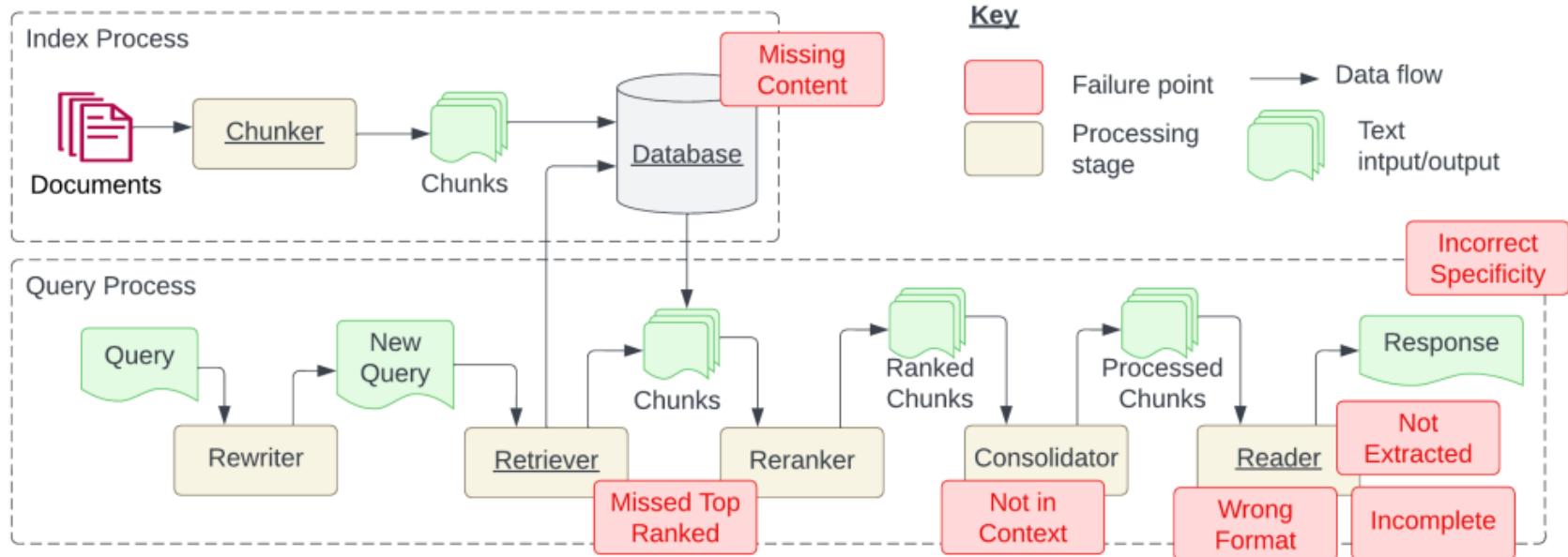
- Determining the relevance and optimal selection of documents for a given prompt.
- Unexpected increase in accuracy with irrelevant documents, indicating a need for specialized integration strategies.

## Customization and Simplification

- Need for tailoring RAG to specific requirements and making it more user-friendly.



# 7 Failure Points When Engineering a RAG System

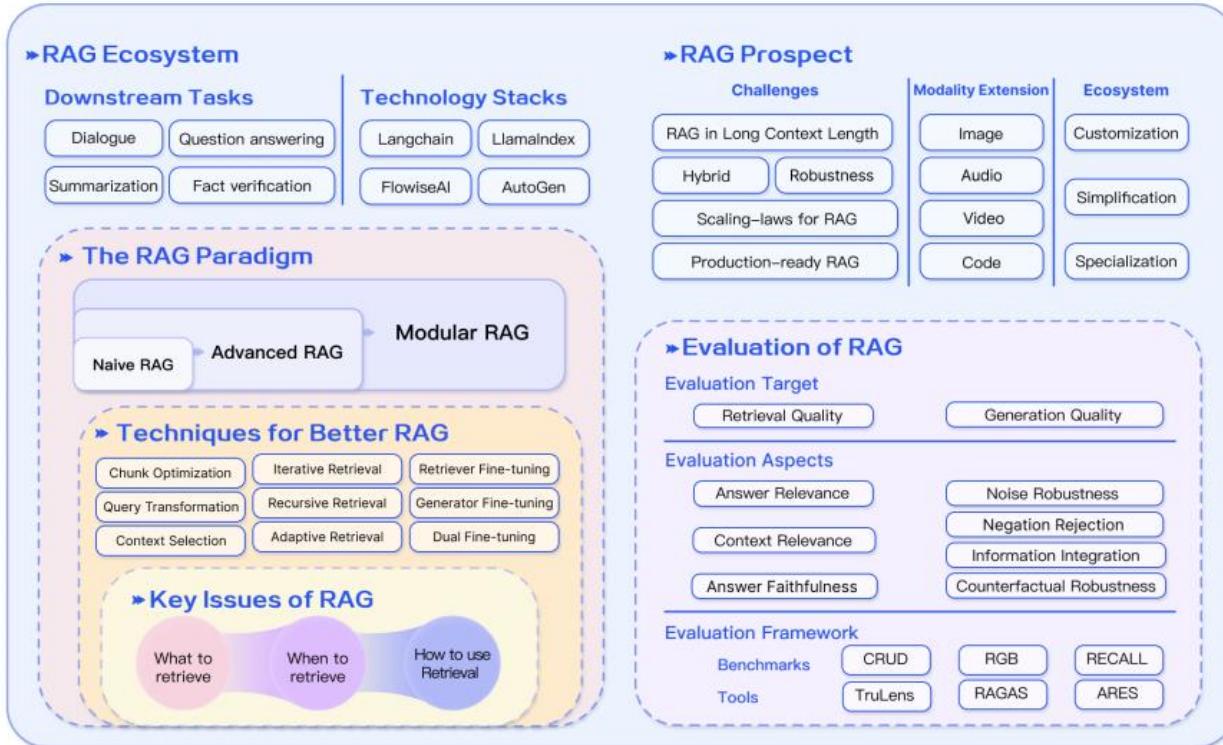


<https://arxiv.org/abs/2401.05856>

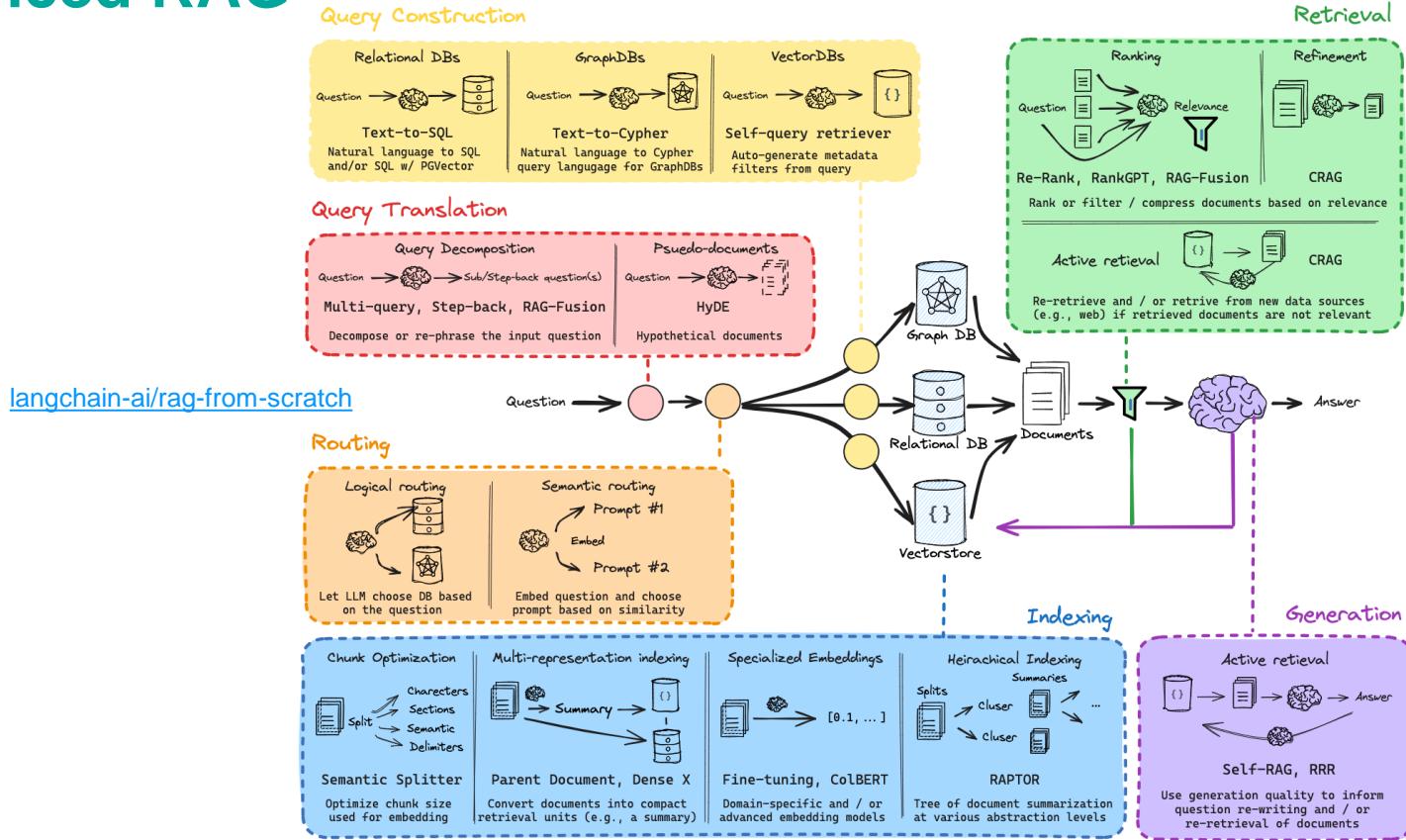
# Advanced Retrieval Augmented Generation (RAG)



# RAG Paradigm

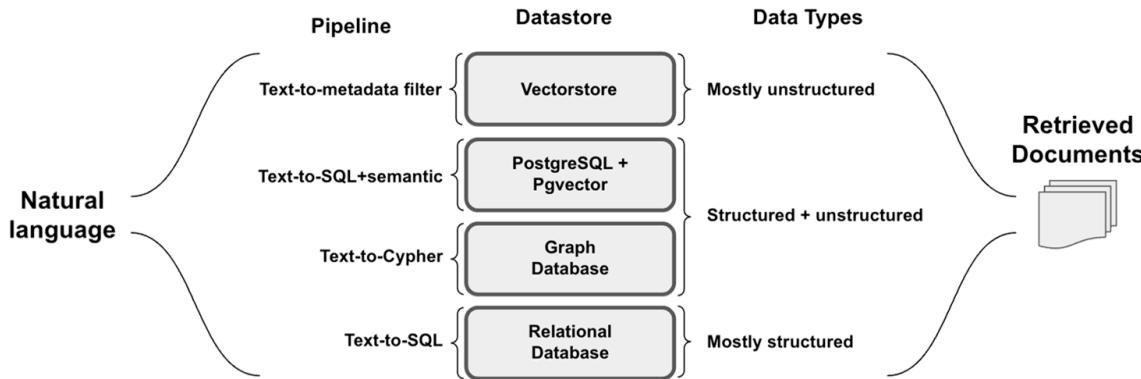


# Advanced RAG

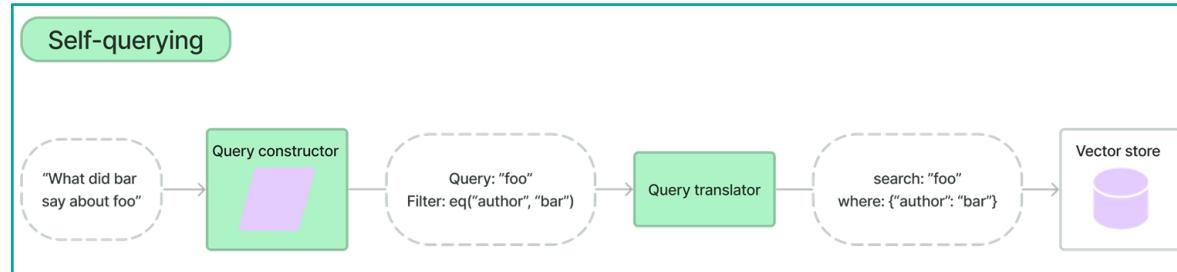


# Query Construction

Process of converting natural language into a specific query syntax for each data type.



Ex: For VectorStore having meta-data filtering



# Query Translation

Transforming user queries into optimized search queries to enhance retrieval performance and accuracy.

## Example of techniques

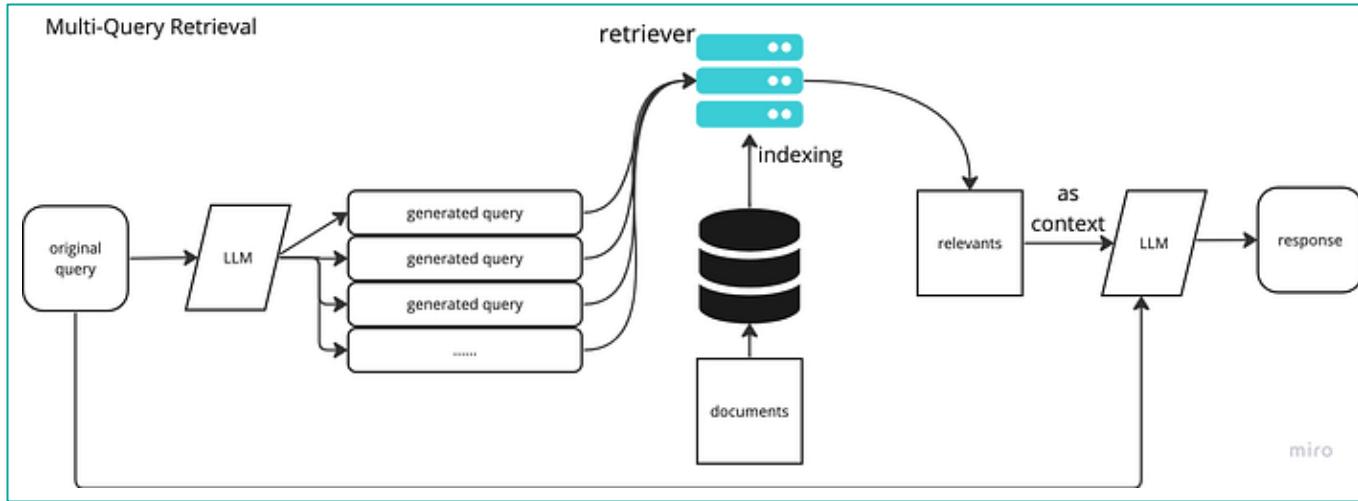
- **Multi-query:** Allows generating and running multiple queries, then combining the results for more comprehensive retrieval.
- **RAG-Fusion:** Involves generating multiple queries and reranking retrieved documents using reciprocal rank fusion for improved results.
- **Decomposition:** Splits user input with multiple distinct questions into separate queries for independent execution.
- **Step-back:** Generates a more abstract "step-back" question alongside the original query to enhance search quality.
- **HyDE (Hypothetical Document Embeddings):** Generates hypothetical relevant documents to improve similarity search in vector-based indexes.

## Reranking

- **MultiQueryRetriever:** for generating and running multiple queries, then combining the results for more comprehensive retrieval.
- **DecompositionRetriever:** Splits user input with multiple distinct questions into separate queries for independent execution.



# Focus: Multi-Query Retrieval



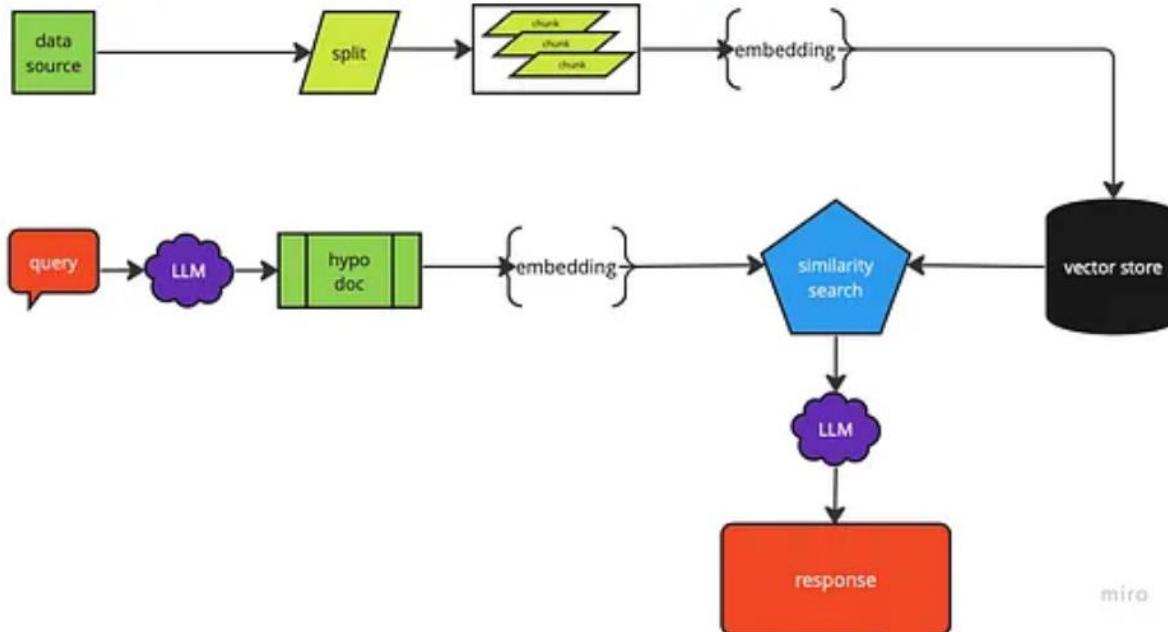
[LangChain / Llama-Index: RAG with Multi-Query Retrieval](#)

## LangChain: MultiQuery Retriever

"What is the doc talking about?"

- 1. What is the main topic discussed in the document?'
- 2. Could you provide a brief summary of the subject matter of the document?'
- 3. What does the document primarily cover and discuss?'

## Focus: HyDE



miro

<https://arxiv.org/pdf/2212.10496>

RAG with Hypothetical Document Embeddings(HyDE) | by TeeTracker | Medium

[https://python.langchain.com/v0.1/docs/use\\_cases/query\\_analysis/techniques/hyde/](https://python.langchain.com/v0.1/docs/use_cases/query_analysis/techniques/hyde/)

# Hybrid Search / Reranking

## Reranking

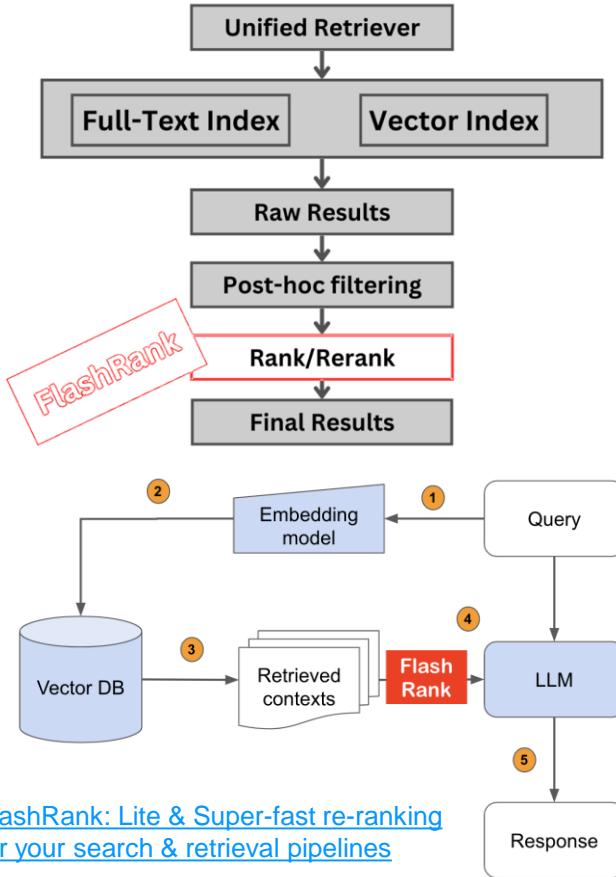
- Reorders results based on additional criteria
- Refines results according to user preferences

## Hybrid search

- Integrates term-based and vector search methods
- Ranks results using a blend of ranking functions

## Usually based on either

- Reciprocal Rank Fusion (RRF)
  - Formula with positions/rank of the documents
- Cross-encoder model:
  - Similarity between the query and each document
  - Specialize trained encoders
- Build-in in some vector stores

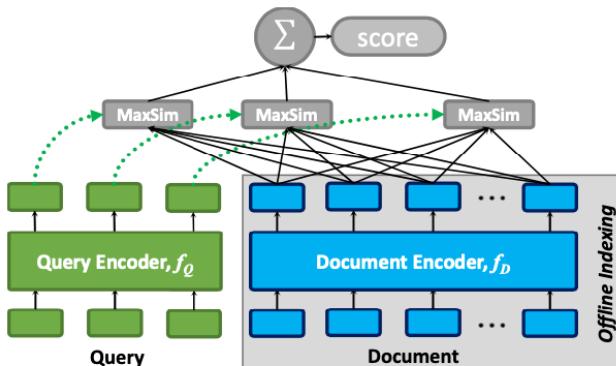


# ColBERT

## Better retriever

- Generates a contextually influenced vector for each token in the passages.
- Generates vectors for each token in the query.
- Score of each document is the sum of the maximum similarity of each query embedding to any of the document embeddings.

## Common library: RAGatouille



Query:  
Effects of climate change on marine ecosystems

Passage:  
The changing climate has profound impacts on marine ecosystems. Rising temperatures, ocean acidification, and altered precipitation patterns all contribute to shifts in the distribution and behavior of marine species, influencing the delicate balance of underwater ecosystems.

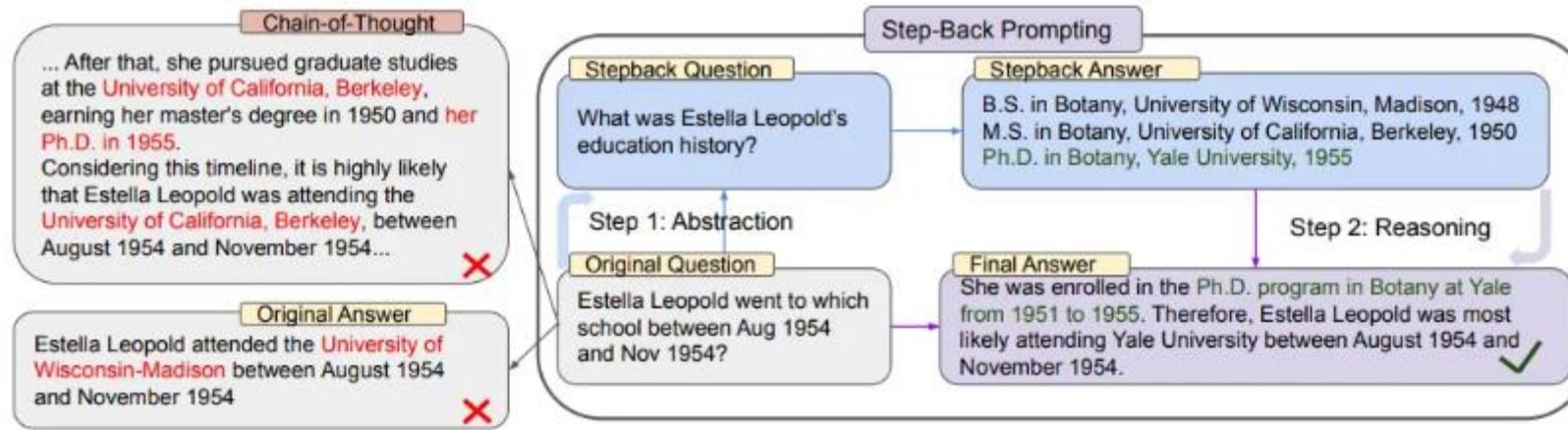
**Run ColBERT scoring for query - passage**

MaxSim Score: 27.71  
Estimated Relevance: 86.60%

**Contextualised Highlights**

The **changing climate** has profound **impacts** on **marine ecosystems**. Rising temperatures, ocean acidification, and **altered** precipitation patterns all contribute to **shifts** in the distribution and behavior of **marine species**, **influencing** the delicate balance of **underwater ecosystems**.

# Step-Back Prompting



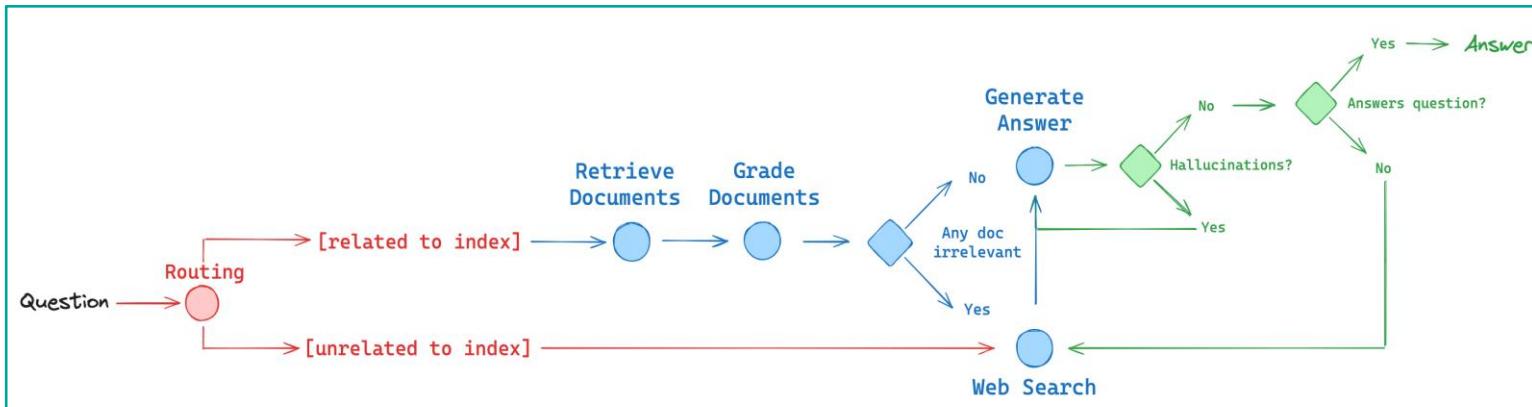
[Take a Step Back: Evoking Reasoning via Abstraction in Large Language Models](#)

[Langchain Elevates with Step-Back Prompting using RAGatouille | by Ankush k Singal | AI Advances](#)

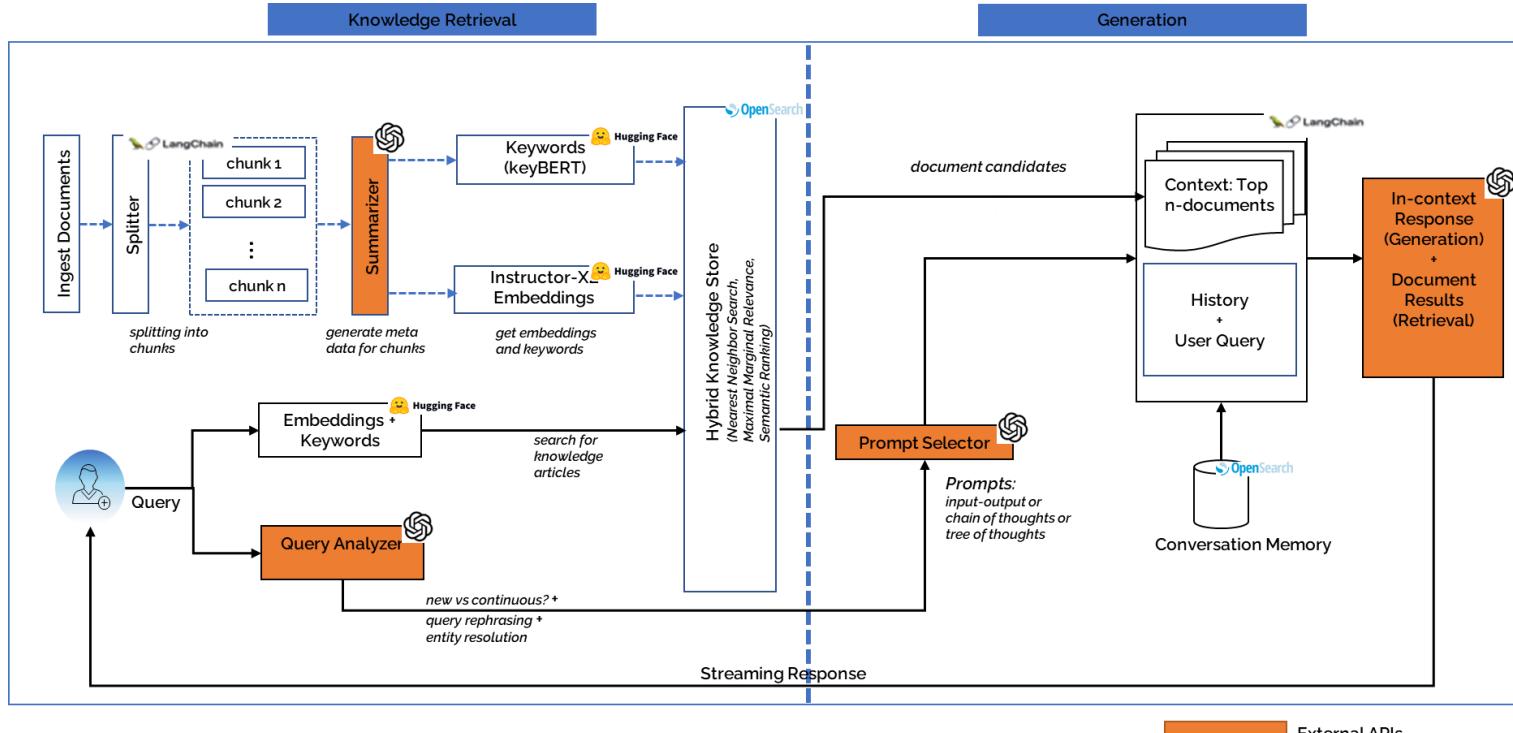


# Other Techniques

- **Adaptive RAG** ([paper](#)). Route questions to different retrieval approaches
- **Corrective RAG** ([paper](#)). Fallback to web search if docs are not relevant to query
- **Self-RAG** ([paper](#)). Fix answers w/ hallucinations or don't address question

**Note***Will be exercised tomorrow...*

# Example : Eviden Knowledge Pilot Architecture



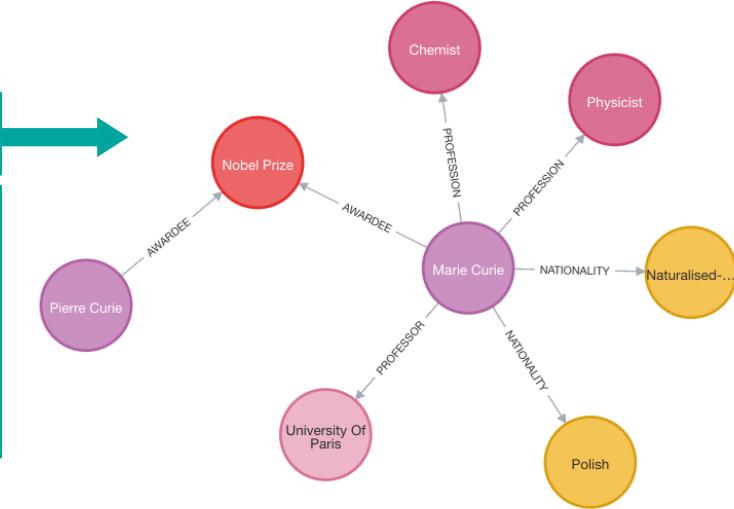
# LLM and Knowledge Graph

```
llm = ChatOpenAI(temperature=0, model_name="gpt-4-turbo")
```

```
llm_transformer = LLMGraphTransformer(llm=llm)
```

```
text = """
Marie Curie, was a Polish and naturalised-French physicist and chemist who conducted pioneering research on
She was the first woman to win a Nobel Prize, the first person to win a Nobel Prize twice, and the only per
Her husband, Pierre Curie, was a co-winner of her first Nobel Prize, making them the first-ever married cou
She was, in 1906, the first woman to become a professor at the University of Paris.
"""

documents = [Document(page_content=text)]
graph_documents = llm_transformer.convert_to_graph_documents(documents)
print(f"Nodes:{graph_documents[0].nodes}")
print(f"Relationships:{graph_documents[0].relationships}")
```



- **Schema (ie ontology) can be provided**
- **Graph can be queried by:**
  - LLM through Cipher or SPARQL generation
  - Vectors similarity through graph embeddings
- **Can serve as memory** of past messages in conversation

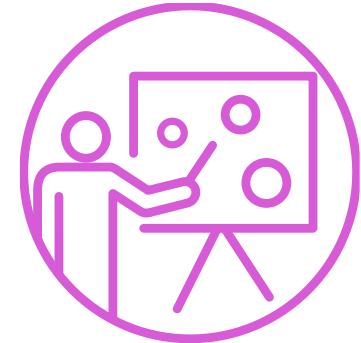
# Practicum



**Assignment**  
Understand Self Query  
- Notebook + WebApp



# Wrap up

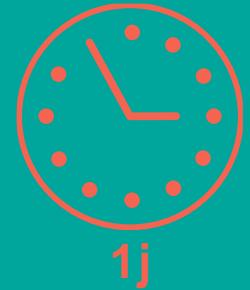


# Quiz



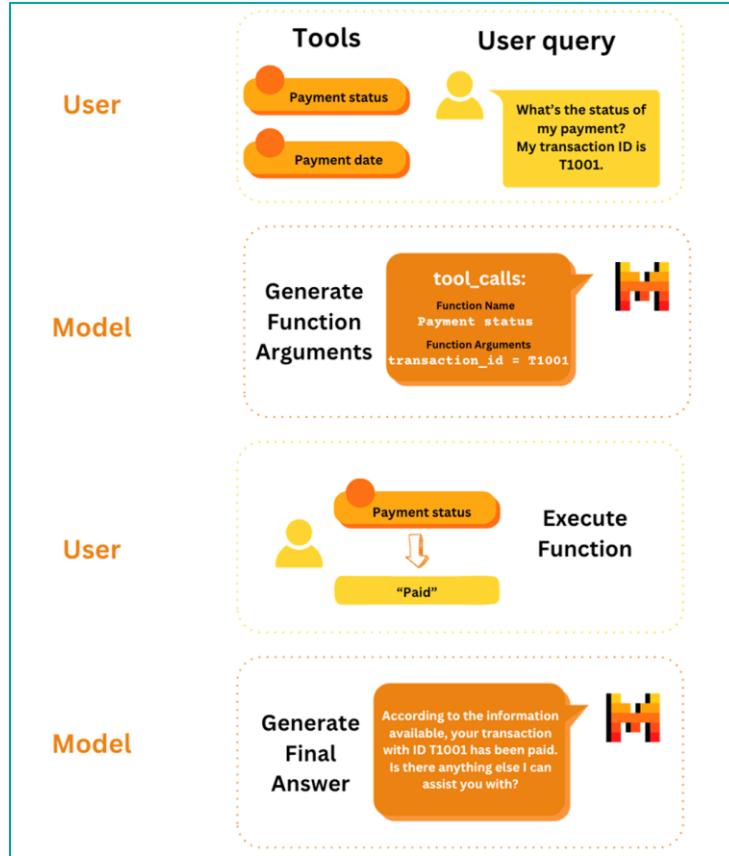
Chapter

# Tools and Agents



# Function Calling

- Function calling is more and more called “tool calling”
- LLM generated output that matches a user-defined schema.
- Users define tools (functions) with names, arguments, and optional identifiers.
- LLM propose tool calls, analyzed as Structured Output
- Users decide whether to execute them.
- Providers like Anthropic, Cohere, Google, Mistral, and OpenAI support tool calling.
- Langchain offers built-in tools and methods for creating custom tools, useful for structured outputs and tool-using chains.
- Still moving topic.



Source: [Function calling | Mistral AI Large Language Models](#)

# Examples Function Definition

```
def get_current_weather(location):
    """Get the current weather in a given location"""
    print("Calling get_current_weather client side.")
    if "tokyo" in location.lower():
        return json.dumps({
            "location": "Tokyo",
            "temperature": "75"
        })
    elif "san francisco" in location.lower():
        return json.dumps({
            "location": "San Francisco",
            "temperature": "60"
        })
    elif "paris" in location.lower():
        return json.dumps({
            "location": "Paris",
            "temperature": "70"
        })
    else:
        return json.dumps({"location": location, "temperature": "unknown"})
```

```
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"
                    }
                },
                "required": ["location"]
            }
        }
]
```

```
messages = [
    {
        "role": "user",
        "content": "What is the weather in San Francisco?"
    }
]

# let's send the request and print the response
response = client.chat.completions.create(
    model="mistralai/Mistral-7B-Instruct-v0.1",
    messages=messages,
    tools=tools,
    tool_choice="auto",
)
tool_calls = response.choices[0].message.tool_calls
for tool_call in tool_calls:
    print(tool_call.model_dump())
```



# Tool Calling with LangChain

```

# ✅ Pydantic class
class multiply(BaseModel):
    """Return product of 'x' and 'y'."""
    x: float = Field(..., description="First factor")
    y: float = Field(..., description="Second factor")

# ✅ LangChain tool
@tool
def exponentiate(x: float, y: float) -> float:
    """Raise 'x' to the 'y'."""
    return x**y

# ✅ Function
def subtract(x: float, y: float) -> float:
    """Subtract 'x' from 'y'."""
    return y-x

# ✅ OpenAI-format dict
# Could also pass in a JSON schema with "title" and "description"
add = {
    "name": "add",
    "description": "Add 'x' and 'y'.",
    "parameters": {
        "type": "object",
        "properties": {
            "x": {"type": "number", "description": "First number to add"},
            "y": {"type": "number", "description": "Second number to add"},
        },
        "required": ["x", "y"]
    }
}

```

<https://blog.langchain.dev/tool-calling-with-langchain/>

```

llm = ChatOpenAI(model="gpt-4-turbo", temperature=0)
llm_with_tools = llm.bind_tools([multiply, exponentiate, add, subtract])

```

```

llm_with_tools.invoke([
    ("system", "You're a helpful assistant"),
    ("human", "what's 5 raised to the 2.743"),
])

```

```

# ⓘ Notice the tool_calls attribute ⓘ

# -> AIMessage(
#     content=...,
#     additional_kwargs={...},
#     tool_calls=[{'name': 'exponentiate', 'args': {'y': 2.743}},
#     response_metadata={...},
#     id='...'
# )

```

tool\_map = {tool.name: tool for tool in tools}



```

def call_tools(msg: AIMessage) -> Runnable:
    """Simple sequential tool calling helper."""
    tool_map = {tool.name: tool for tool in tools}
    tool_calls = msg.tool_calls.copy()
    for tool_call in tool_calls:
        tool_call["output"] = tool_map[tool_call["name"]].invoke(tool_call["args"])
    return tool_calls

```

chain = llm\_with\_tools | call\_tools



```

agent = create_tool_calling_agent(llm, tools, prompt)
agent_executor = AgentExecutor(agent=agent, tools=tools, verbose=True)

agent_executor.invoke({"input": "what's 3 plus 5 raised to the 2.7"})

```

# Tools, tools, tools

## Many tools in LangChain and LangChain community

- Gathering information : Wiki Tool, Google Search Tool, ...
- taking action : Python REPL, Wolfram Alpha Toolkit, Google Toolkit, ...
- Manipulating data : OpenAI API Tool, Anthropic API Tool, GitHub Toolkit...
- See:
  - <https://python.langchain.com/v0.1/docs/integrations/tools/>
  - <https://python.langchain.com/v0.1/docs/integrations/toolkits/>

## Companies specialized in tools offerings

Ex:

- Eden.ai
- Kay.ai...



# Examples and Pitfalls

## PDF Scanning

- Difficult task if diagrams and tables
- OCR + vision works
- Specialized API.

## Database query

- Usually need few-prompt learning
- Specialized models are emerging

## Code execution

- Need secure environment
- Some offering exits



# Practicum



**Assignment**  
Tool creation and Calling  
Eden.ai tool cll



# LangGraph

- **Library for building stateful**, multi-actor applications with Language Learning Models (LLMs).
- **Allows coordination** of multiple chains across cyclic computational steps using Python functions.
- **Adds cycles** and persistence to LLM applications, crucial for agent-like behaviors.
- **Define:**
  - State
  - Nodes
  - Edges
- **Used:**
  - Implement complex workflows
  - Call an LLM in a loop to determine the next action...
- **Methods:**
  - add\_node
  - add\_edge
  - add\_conditional\_edges
  - set\_entry\_point
  - set\_conditional\_entry\_point
  - set\_finish\_point
- **State is passed** between Nodes (Python dict)
- **Nodes can:**
  - Call Runnable
  - Call Tool
  - Route between alternatives...
- **Prebuild components:**
  - ToolInvocation
  - ToolExecutor...

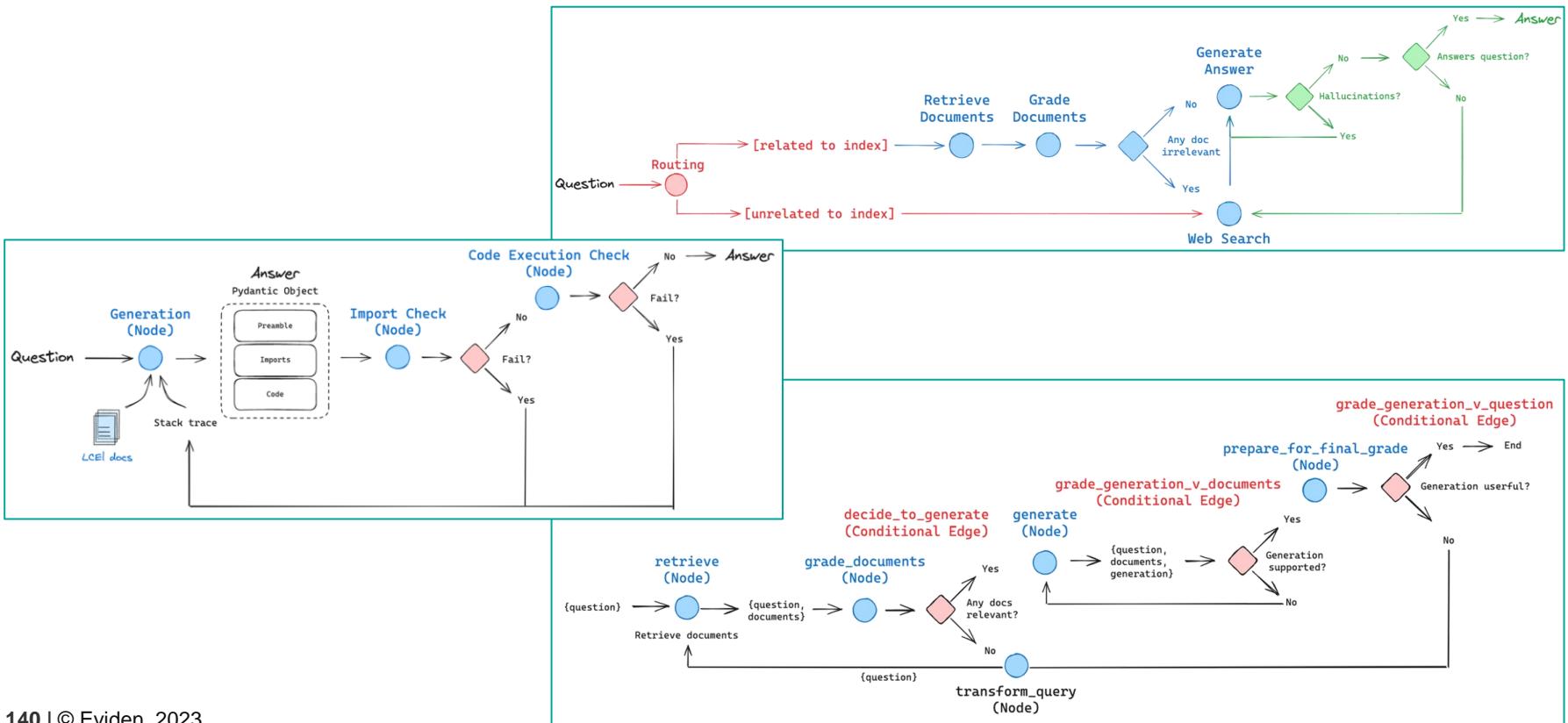


# Practicum

**Assignment**  
Improve Fibonacci sequence

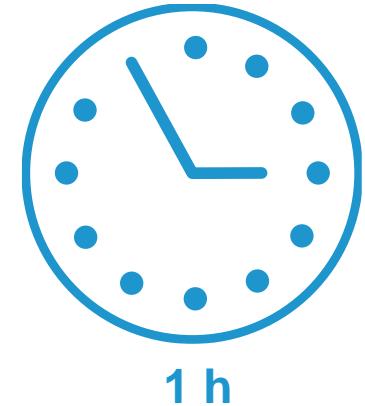


# Examples LangGraph



# Practicum

Assignment  
Advanced RAG



# Back on Routers

## Logical Router

Directs according to the value of some variables

## LLM Completion Routers

Directs the input to the most appropriate large language model (LLM) to generate a text completion based on the context or task.

## LLM Function Calling Routers

Invokes specific functions or methods within a language model based on the input command or request.

## Semantic Routers

Analyzes the meaning of the text to route the input to the best-suited model or service that can handle the semantic content.

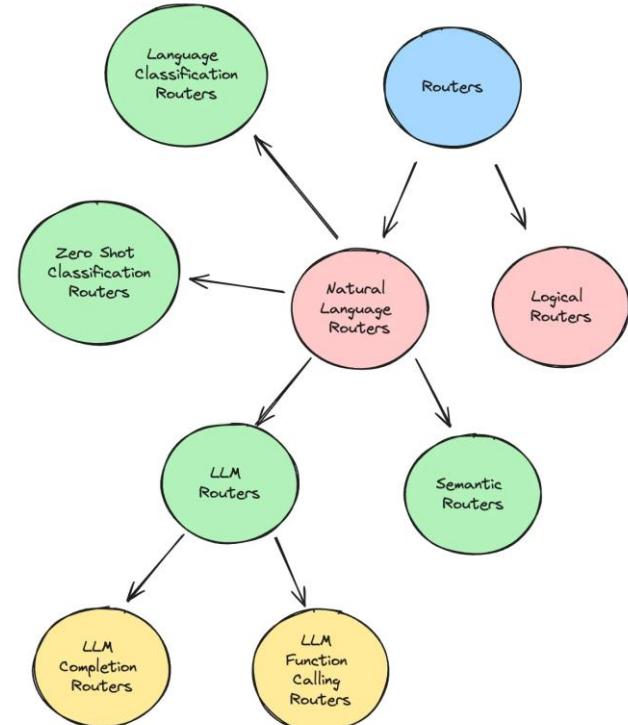
## Zero Shot Classification Routers

Classifies input into categories without prior explicit examples, using models that understand text in a zero-shot learning framework.

## Language Classification Routers

Identifies the language of the input text and routes it to a specialized model capable of handling that particular language.

> Autonomous AI Agents can act as Routers.....



# What is an AI Agent ?

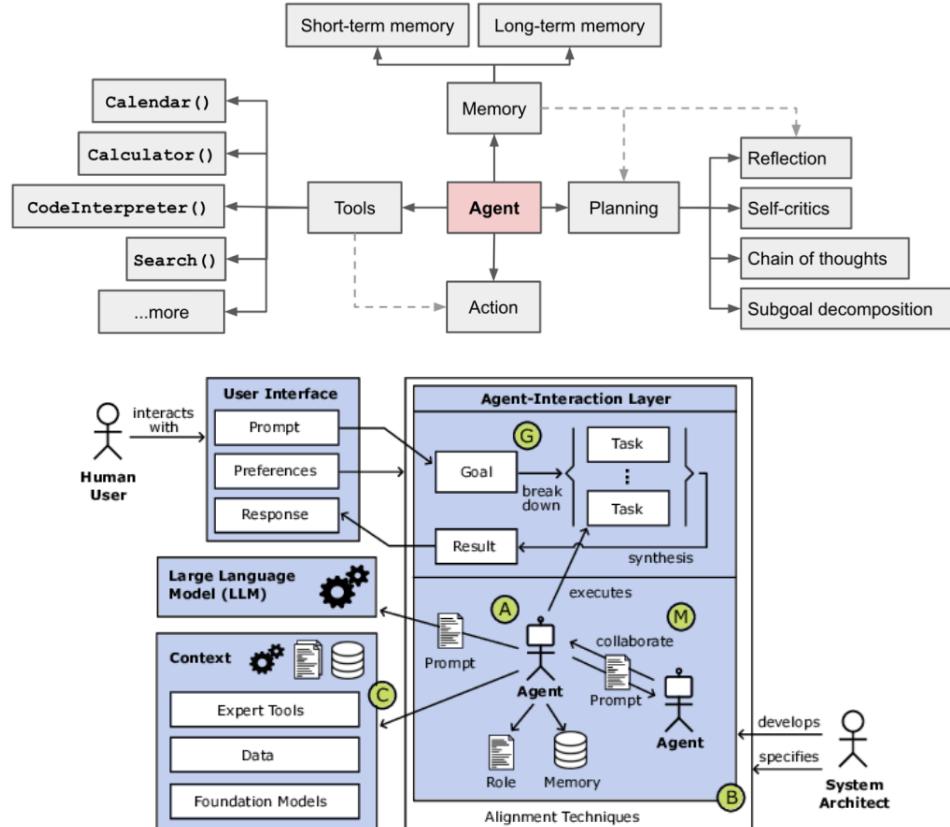
## No clear definition (yet)

- Combination of tools, planning and reaction capabilities.

## Andrew Ng Agents pattern

- Reflection Mode: Agent examines and refines own outputs through self-review.
- Tool Use Mode: Agent uses tools like search to gather info and take actions.
- Planning Mode: Agent breaks down complex tasks and follows a plan for high quality outputs.
- Multiagent Collaboration Mode: Multiple agents allocate tasks and discuss ideas by working together.

<https://www.linkedin.com/pulse/4-main-ai-agent-design-patterns-recommend-andrew-ng-yiman-huang-nwype/>



# Another Agents Taxonomy

## Simple reflex agents

- React only to current inputs without considering history.

## Model-based agents

- Maintain an internal world model and update it based on new information.

## Goal-based agents

- Can search, plan and take steps to achieve goals through decision making.

## Utility-based agents

- Evaluate options and select the most favorable based on a predefined satisfaction measure.

## Multi-agent systems

- Have agents collaborating through coordination and communication towards a shared objective.

## Learning agents

- Can adapt over time by learning from experiences gained through interactions.

## Hierarchical agents

- Have varying levels/roles organized in a hierarchy to facilitate coordination of complex tasks.



# LangGraph example: collaboration, routing, without planning

```
#Initializes a StateGraph object named workflow
workflow = StateGraph(AgentState)

#nodes represent different actions or decisions
workflow.add_node("Researcher", research_node)
workflow.add_node("Chart Generator", chart_node)
workflow.add_node("call_tool", tool_node)

#edges define the transitions between these nodes
workflow.add_conditional_edges(
    "Researcher",
    router,
    {"continue": "Chart Generator", "call_tool": "call_tool",
     "end": END},
)



Next step from router

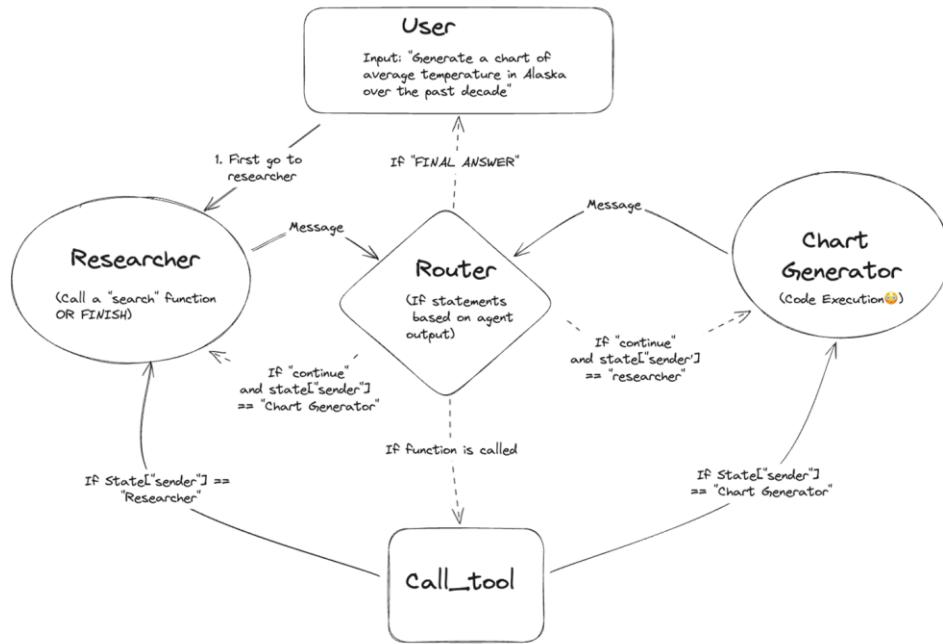


workflow.add_conditional_edges(
    "Chart Generator",
    router,
    {"continue": "Researcher", "call_tool": "call_tool",
     "end": END},
)

workflow.add_conditional_edges(
    "call_tool",
    lambda x: x["sender"],
    {
        "Researcher": "Researcher",
        "Chart Generator": "Chart Generator",
    },
)

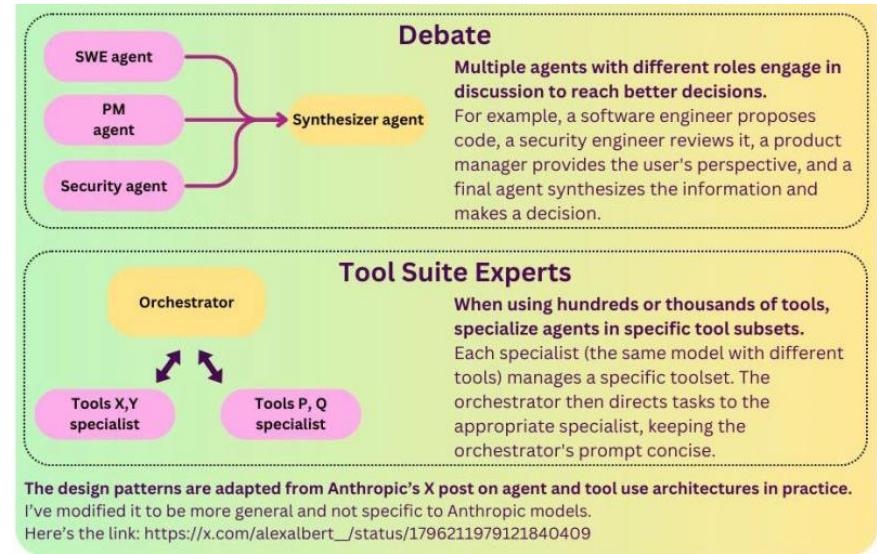
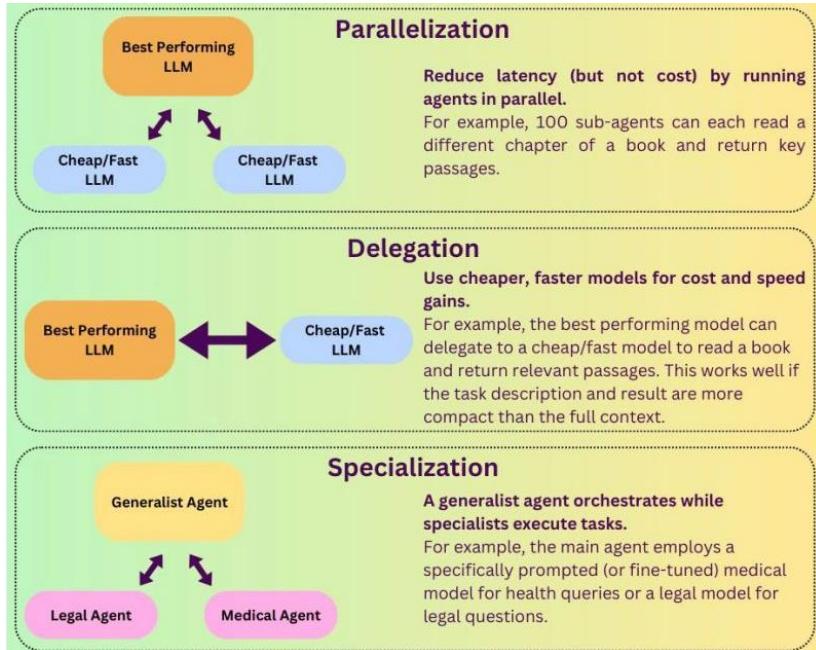
# "Researcher" node as the starting point for
# the graph execution
workflow.set_entry_point(["Researcher"])

#Compiling the workflow
graph = workflow.compile()
```



[LangGraph: Multi-Agent Collaboration Explained | by Kamal Dhungana | Apr, 2024 | Medium](#)  
[https://github.com/langchain-ai/langgraph/blob/main/examples/multi\\_agent/multi-agent-collaboration.ipynb](https://github.com/langchain-ai/langgraph/blob/main/examples/multi_agent/multi-agent-collaboration.ipynb)

# LLM Agent Architecture Patterns



[Source](#)

# Planning in LLM

<https://arxiv.org/pdf/2402.02716>

## Methods

- Task decomposition (Divide and Conquer)
- Multi-plan Selection (Generate multiple plans and select the optimal)
- External Planner-aided (Formalize tasks and utilize external planner)
- Reflection & Refinement (Reflect on experiences and refine plans)
- Memory-aided Planning (Leverage memory to aid planning)

## Today

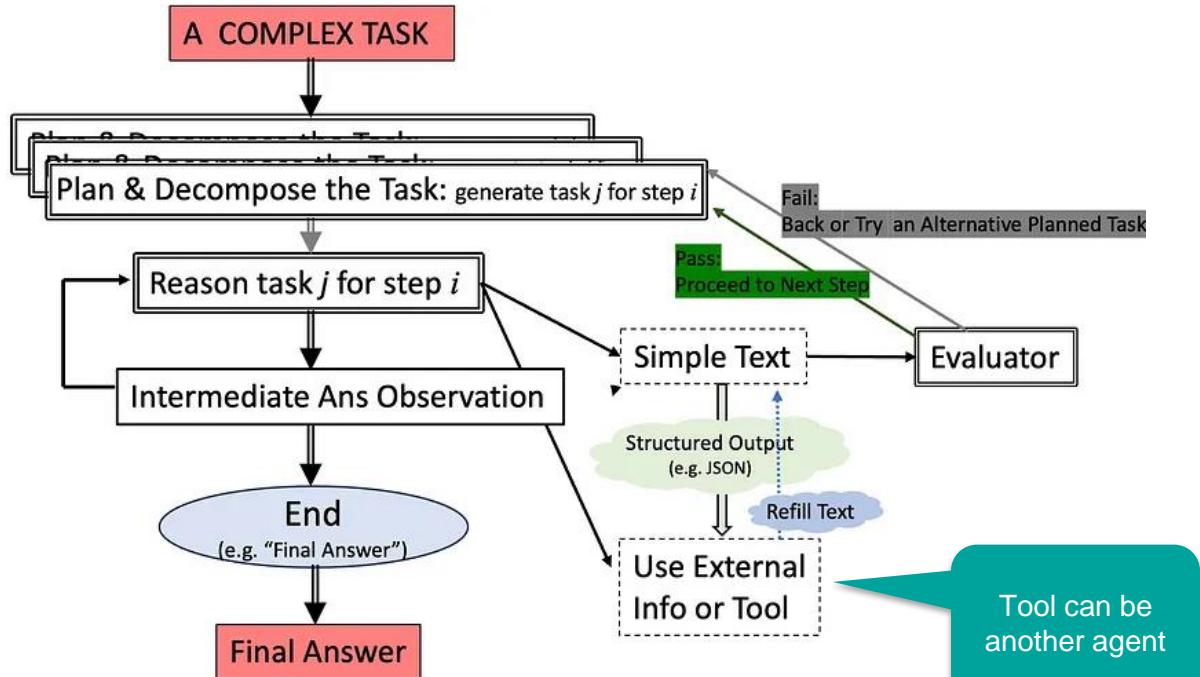
- Task decomposition is used in production application
- Other methods more R&D

## Main pattern: ReAct

- Implemented in several frameworks: AWS Agents for BedRock, Microsoft Semantic Kernel, ...
- Several variants (notably thanks to LLM function selection / call capabilities)



# Autonomous Agents - ReACT



[A Complete Guide to LLMs-based Autonomous Agents \(Part I\): | by Yule Wang, PhD | The Modern Scientist | Medium](#)



# Agents in LangChain

## Agent Types

- [Tool Calling Agent](#)
  - Require LLM with functions
  - Detect when tools should be called and respond with inputs to pass to tools
  - More reliable valid tool calls than generic text completion
  - Can repeatedly call tools and receive results until query resolved
  - Support wider range of tool providers than original OpenAI agent
- [ReAct](#) : Initial implementation Works with many LLMs
- [Structured chat](#) : Extension of ReAct to many tools / structured output

## Specialized agents:

- [create\\_tool\\_calling\\_agent](#), `create_cohere_react_agent`, `create_sql_agent`, ..

## Executor: [AgentExecutor](#)

- Goal: select and utilize tools in a structured manner.
- Will become more configurable with LangGraph



# Custom Agents

```
@tool
def get_word_length(word: str) -> int:
    """Returns the length of a word."""
    return len(word)
```

```
prompt = ChatPromptTemplate.from_messages(
    [
        (
            "system",
            "You are very powerful assistant, but don't know current events",
        ),
        ("user", "{input}"),
        MessagesPlaceholder(variable_name="agent_scratchpad"),
    ]
)
```

```
llm_with_tools = llm.bind_tools(tools)
```

```
agent = (
    {
        "input": lambda x: x["input"],
        "agent_scratchpad": lambda x: format_to_openai_tool_messages(
            x["intermediate_steps"]
        ),
    }
    | prompt
    | llm_with_tools
    | OpenAIToolsAgentOutputParser()
)
```

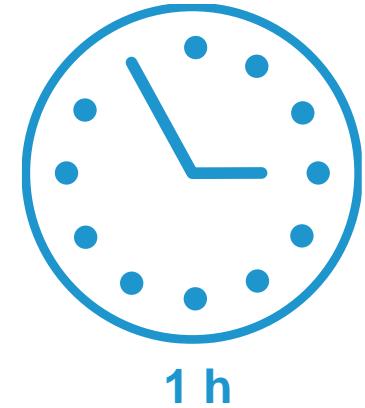
```
agent_executor = AgentExecutor(agent=agent, tools=tools, verbose=True)
```

```
agent_executor.stream({"input": "How many letters in the word eudca"})
```



# Practicum

Assignment:

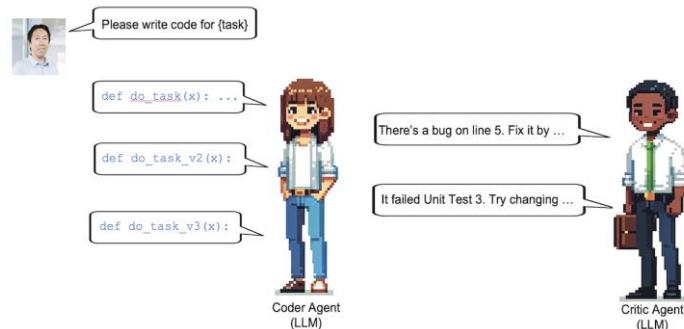


# Cooperative Agents

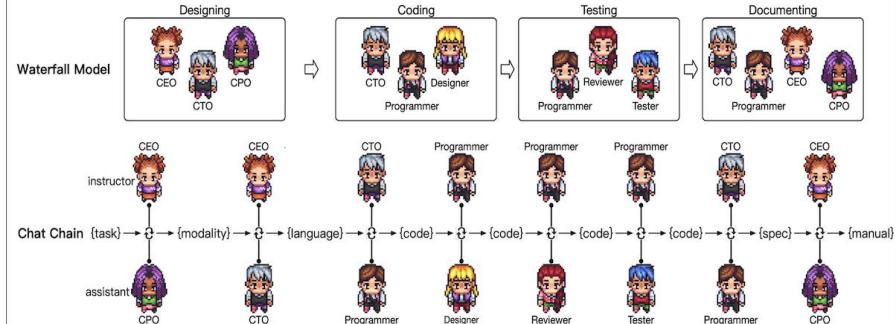
- **Break down a task** into subtasks to be executed by different roles
- **Different agents** accomplish different subtasks
  - Ex: software engineer, product manager, designer, QA (quality assurance) engineer.
- **Prompt embedd domain knowledge** for each role (know-how procedure, ...)
- **Framework:** AutoGen, Crew AI, ..
- **Ex:** ChatDev, Devin, [Coze](#), CAMEL, ...

*> Still hard to predict outcome.*

## Agentic Design Patterns: Reflection



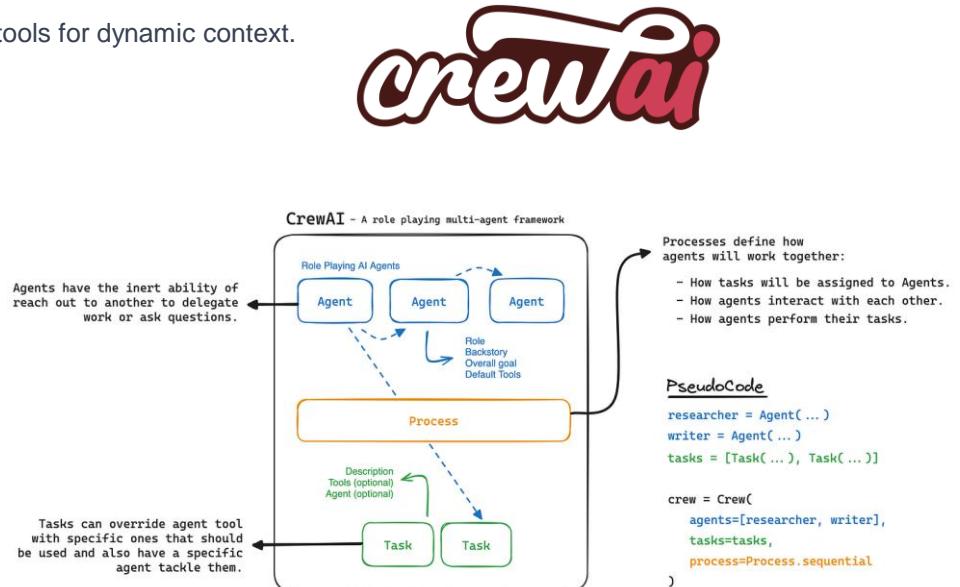
## Agentic Design Patterns: Multi-Agent Collaboration



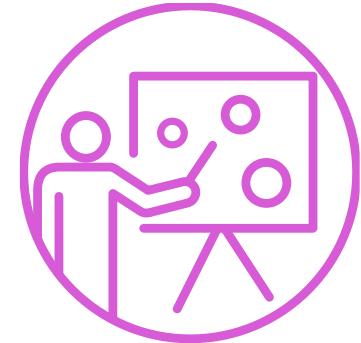
Proposed ChatDev architecture. Image adapted from "Communicative Agents for Software Development," Qian et al. (2023).

# Ex CrewAI

- Python framework for **orchestrating autonomous AI Agents** with specific roles.
- Enable seamless agent collaboration, combining generative AI with tools for dynamic context.
- **Core Concepts**
  - **Agents:**
    - Perform tasks, make decisions, communicate.
    - Each has a unique role and skill set.
  - **Tasks:**
    - Individual missions assigned to agents.
    - Include descriptions, assigned agents, required tools.
  - **Tools:**
    - Skills agents use to perform tasks
    - From CrewAI Toolkit and LangChain Tools.
  - **Processes:**
    - Orchestrate task execution, similar to a project manager
    - Ensure optimal task completion.
- **Benefits**
  - Fosters creative problem-solving through agent collaboration
  - Excels in handling complex tasks collaboratively
  - Decomposes tasks into a manageable sub-problems



# Wrap up

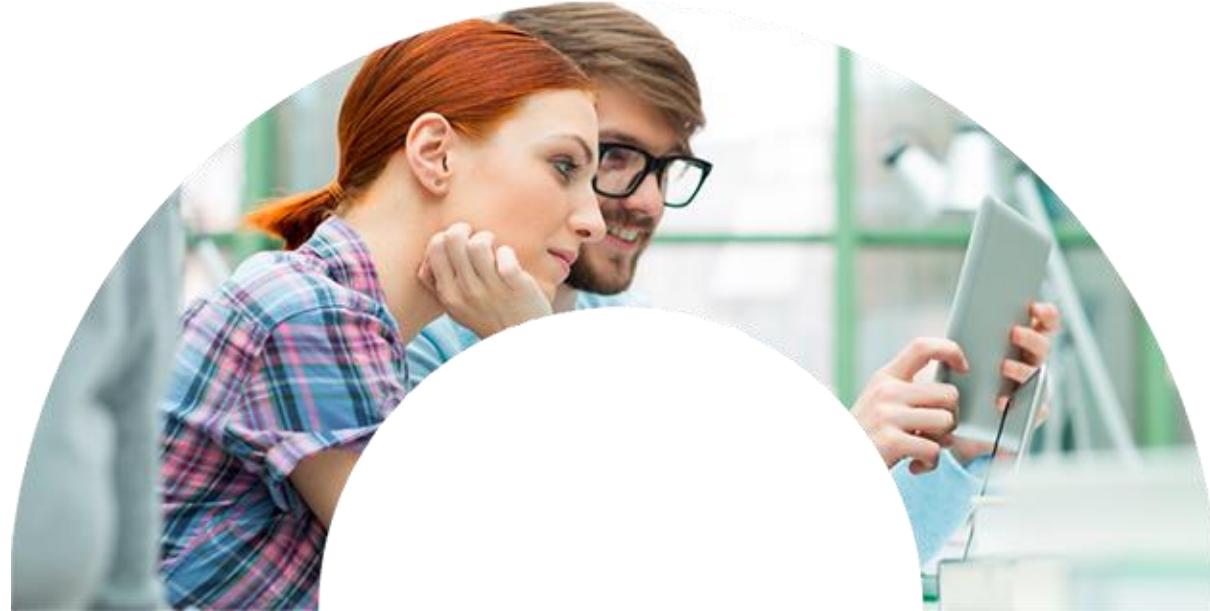


Chapter

# Security & sovereignty Optimisation & monitoring



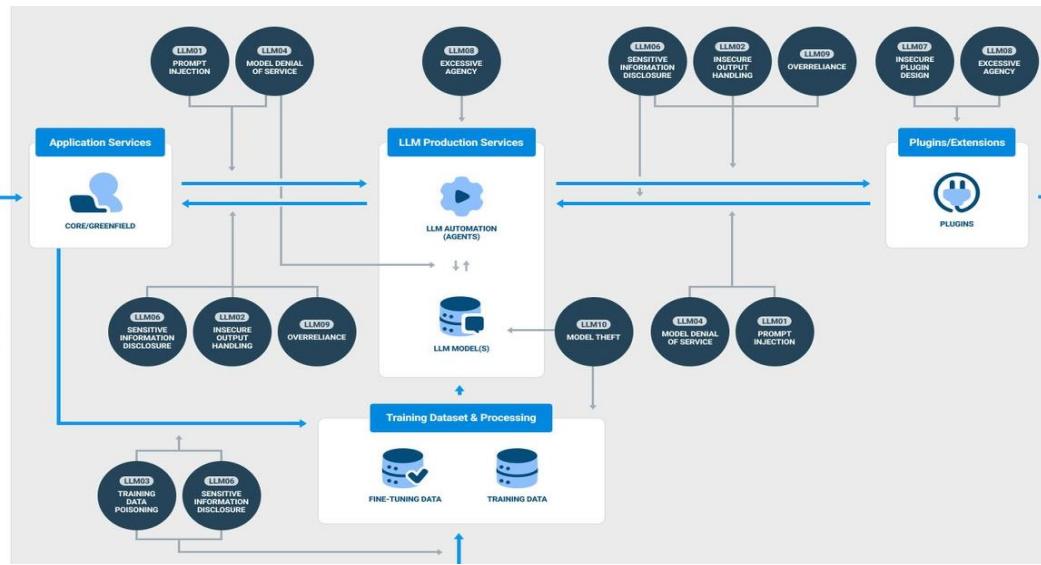
# Security & Sovereignty



# Risks, Risks, Risks

## Inference risk at Agent Level

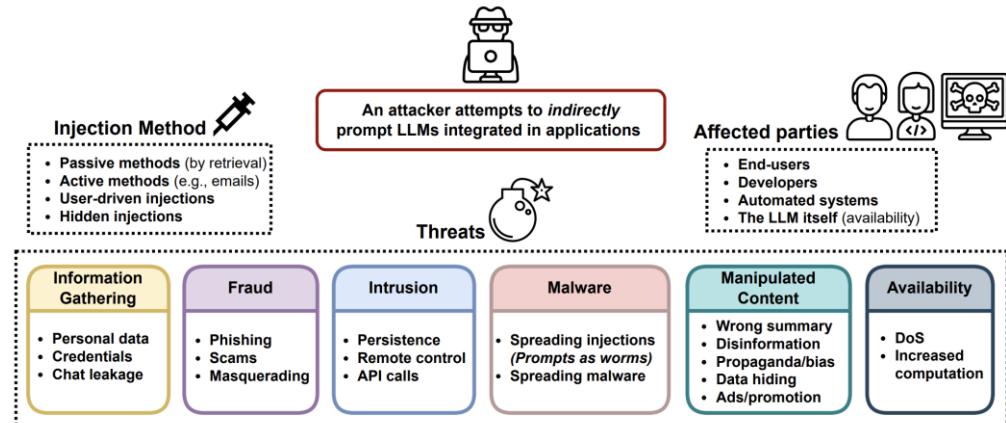
- **Hallucinations:** generation of factually incorrect, irrelevant, or inconsistent text.
- **Not grounded answer:** When the answer is not in the provided document
- **Unsafe prompt**, according to a policy
- **Attacks:**
  - **Jailbreak:** Attack using prompt injection to bypass safety features in LLMs.
  - **Prompt Extraction:** Attack to reveal system prompt or hidden info in an LLM's context.
  - **Membership Inference:** Data privacy attack to identify if a data sample was in a training set.
  - **Indirect Prompt Injection:** LLM ingests a prompt injection attack indirectly (e.g., through web page).
  - **Information Gathering:** extract user data or leak chat history by interacting in chat sessions, persuading users to divulge information
  - **Prompt Injection Agents:** Force an agent to call a tool
- **Code generation / execution risks**
- **Tradeoff risks vs. costs vs. performance**



OWASP Top 10 for Large Language Model Applications

# Adversarial Attacks

Attack	Type	Description
Token manipulation	Black-box	Alter a small fraction of tokens in the text input such that it triggers model failure but still remain its original semantic meanings.
Gradient based attack	White-box	Rely on gradient signals to learn an effective attack.
Jailbreak prompting	Black-box	Often heuristic based prompting to “jailbreak” built-in model safety.
Human red-teaming	Black-box	Human attacks the model, with or without assist from other models.
Model red-teaming	Black-box	Model attacks the model, where the attacker model can be fine-tuned.



[Adversarial Attacks on LLMs | Lil'Log](#)

# Trust Layer

## Data anonymization

- Classical NLP techniques
- Ex: Microsoft Presidio

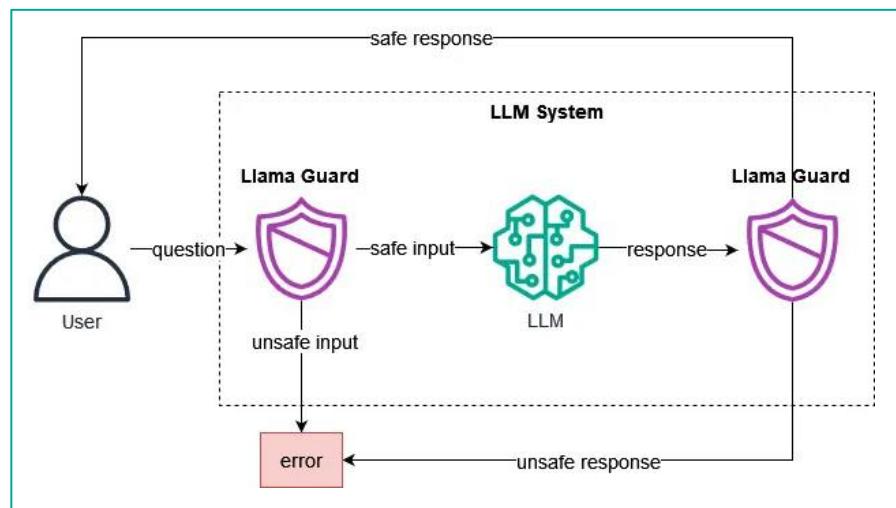
## Prompt / Answer Evaluator

- Often LMM based
- Ex: Llama Guard, external APIs
  - Trained on ML Commons AI Safety benchmark
  - Support few-shot learning
- Classical AI classifier
- External API services
- Open source solution / Hub
  - Ex: Guardrails / [Guardrails hub](#)

## Domain specific answer evaluator

- Ex: Llama Code Shield

INPUT: Hi, my name is David and my number is 212 555 1234  
 OUTPUT: Hi, my name is <PERSON> and my number is <PHONE\_NUMBER>



[Safeguard Your LLM Chatbot With Llama Guard 2 | by Dr. Leon Eversberg | May, 2024 | Towards Data Science](#)

# Ethics

- AI assistants could profoundly impact users and society.
- They may have autonomy to plan and perform tasks.
- Assistants communicating through language risks being indistinguishable from humans.
- Assistants could display novel capabilities that are hard to foresee.
- Assistants with autonomy brings risks of unclear instructions and misaligned actions.
- More autonomy could enable high-impact misuse like spreading misinformation.
- Privacy safeguards are needed as language blurs human-machine lines.
- Collective use risks coordination problems like overloading services.
- Novel capabilities require evaluating unintended wider societal impacts.

## The Ethics of Advanced AI Assistants

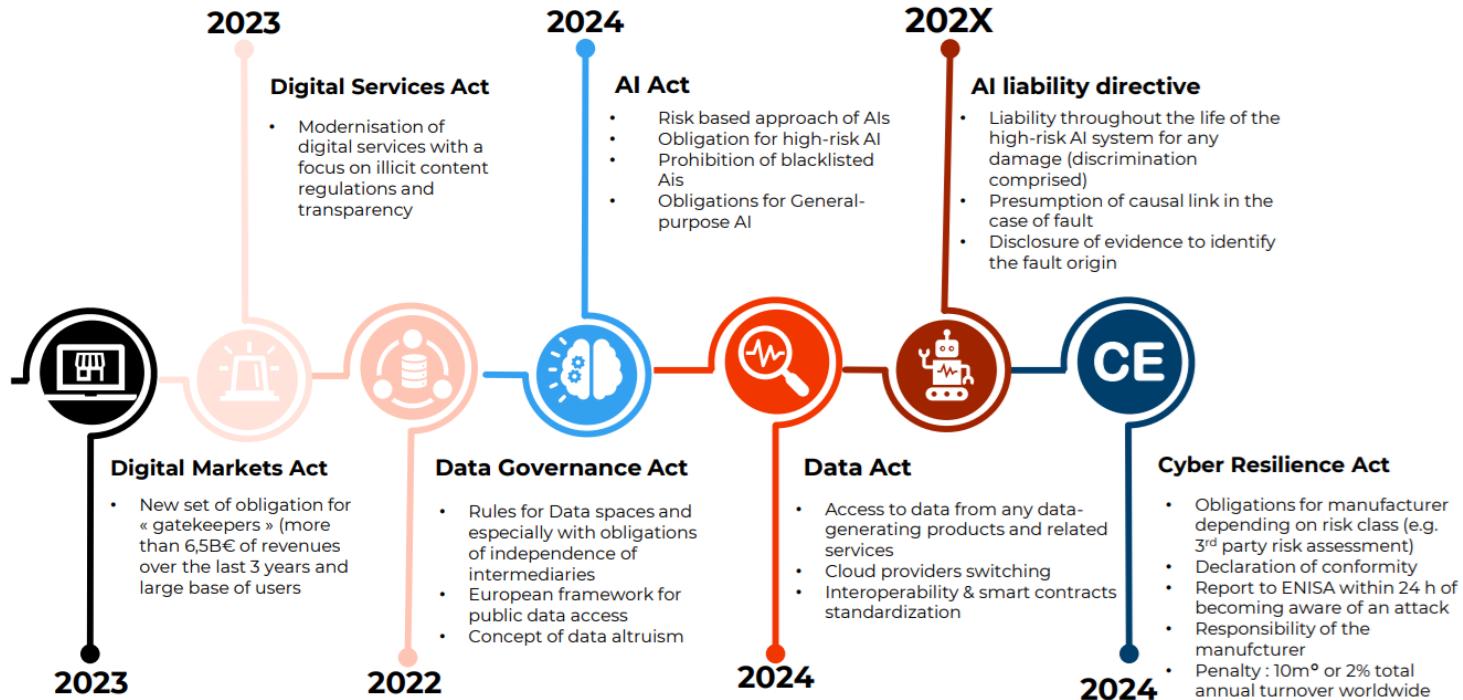
Iason Gabriel<sup>\* 1</sup>, Arianna Manzini<sup>\* 1</sup>, Geoff Keeling<sup>\* 2</sup>, Lisa Anne Hendricks<sup>1</sup>, Verena Rieser<sup>1</sup>, Hasan Iqbal<sup>1</sup>, Nenad Tomašev<sup>1</sup>, Ira Ktena<sup>1</sup>, Zachary Kenton<sup>1</sup>, Mikel Rodriguez<sup>1</sup>, Selim El-Sayed<sup>1</sup>, Sasha Brown<sup>1</sup>, Canfer Akbulut<sup>1</sup>, Andrew Trask<sup>1</sup>, Edward Hughes<sup>1</sup>, A. Stevie Bergman<sup>1</sup>, Renee Shelby<sup>2</sup>, Nahema Marchal<sup>1</sup>, Conor Griffin<sup>1</sup>, Juan Mateos-Garcia<sup>1</sup>, Laura Weidinger<sup>1</sup>, Winnie Street<sup>2</sup>, Benjamin Lange<sup>2,4</sup>, Alex Ingerman<sup>2</sup>, Alison Lenz<sup>2</sup>, Reed Enger<sup>2</sup>, Andrew Barakat<sup>2</sup>, Victoria Krakovna<sup>1</sup>, John Oliver Siy<sup>2</sup>, Zeb Kurth-Nelson<sup>1</sup>, Amanda McCroskery<sup>2</sup>, Vijay Bolina<sup>1</sup>, Harry Law<sup>1</sup>, Murray Shanahan<sup>1</sup>, Lize Alberts<sup>2,5,6</sup>, Borja Balle<sup>1</sup>, Sarah de Haas<sup>2</sup>, Yetunde Ibitoye<sup>2</sup>, Allan Dafoe<sup>1</sup>, Beth Goldberg<sup>3</sup>, Sébastien Krier<sup>1</sup>, Alexander Reese<sup>2</sup>, Sims Witherspoon<sup>1</sup>, Will Hawkins<sup>1</sup>, Maribeth Rauh<sup>1</sup>, Don Wallace<sup>1</sup>, Matija Franklin<sup>7</sup>, Josh A. Goldstein<sup>8</sup>, Joel Lehman<sup>9</sup>, Michael Klens<sup>10</sup>, Shannon Vallor<sup>11</sup>, Courtney Biles<sup>1</sup>, Meredith Ringel Morris<sup>1</sup>, Helen King<sup>1</sup>, Blaise Aguera y Arcas<sup>2</sup>, William Isaac<sup>1</sup> and James Manyika<sup>2</sup>

This paper focuses on the opportunities and the ethical and societal risks posed by advanced AI assistants. We define advanced AI assistants as artificial agents with natural language interfaces, whose function is to plan and execute sequences of actions on behalf of a user – across one or more domains – in line with the user's expectations. The paper starts by considering the technology itself, providing an overview of AI assistants, their technical foundations and potential range of applications. It then explores questions around AI value alignment, well-being, safety and malicious uses. Extending the circle of inquiry further, we next consider the relationship between advanced AI assistants and individual users in more detail, exploring topics such as manipulation and persuasion, anthropomorphism, appropriate relationships, trust and privacy. With this analysis in place, we consider the deployment of advanced assistants at a societal scale, focusing on cooperation, equity and access, misinformation, economic impact, the environment and how best to evaluate advanced AI assistants. Finally, we conclude by providing a range of recommendations for researchers, developers, policymakers and public stakeholders.

[The ethics of advanced AI assistants - Google DeepMind \(274 p\)](#)



# Regulation

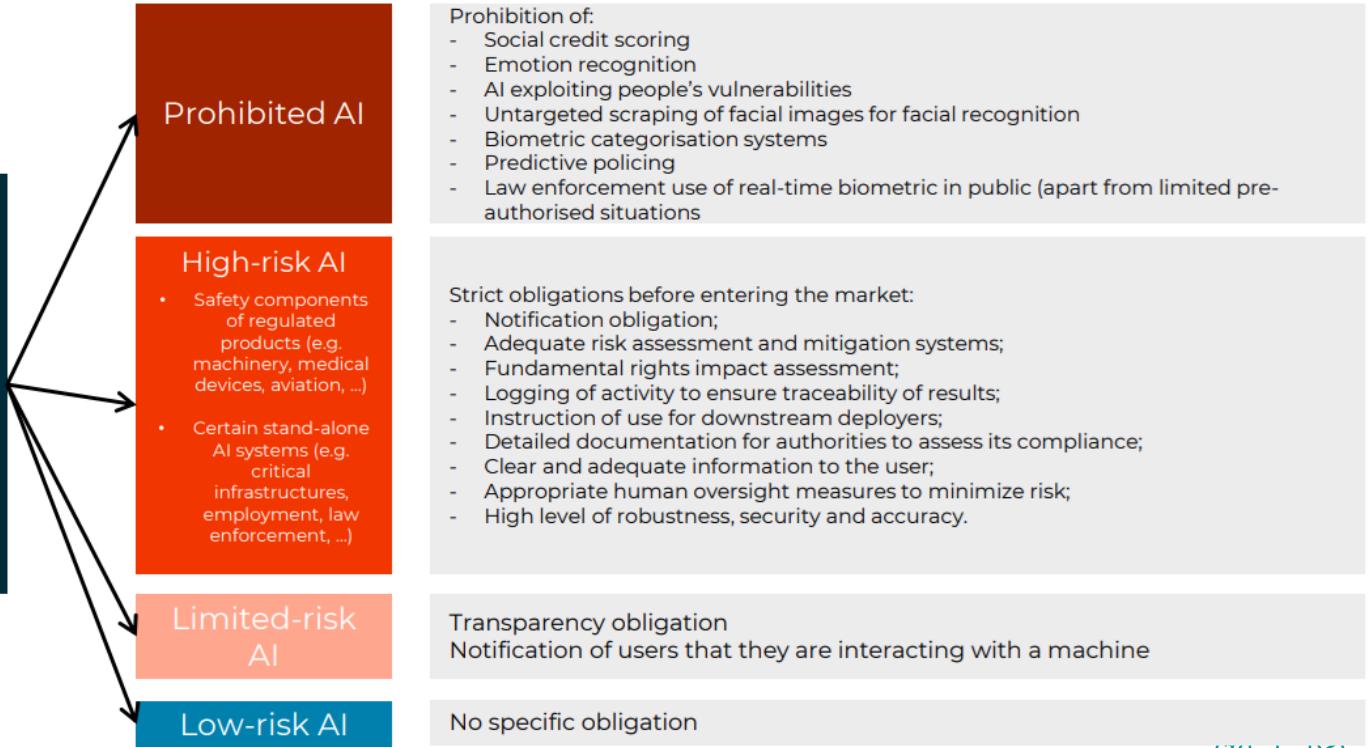


# AI Act Overview

End 2025

## Definition of AI System

An AI is a machine-based system designed to operate with varying levels of autonomy and that may exhibit adaptiveness after deployment and that, for explicit or implicit objectives, infers, from the input it receives, how to generate outputs such as predictions, content, recommendations, or decisions that can influence physical or virtual environments.

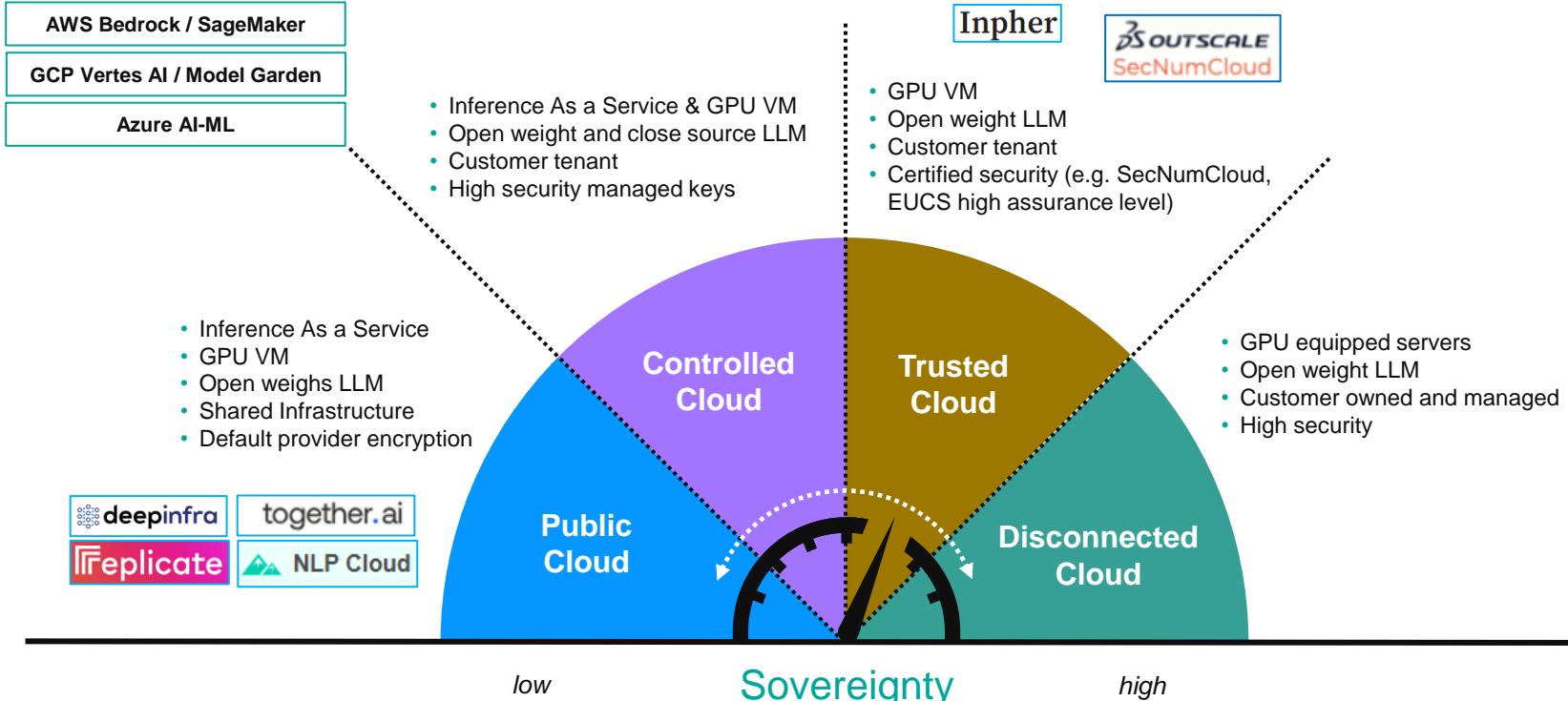


# Leak of confidential or privacy information

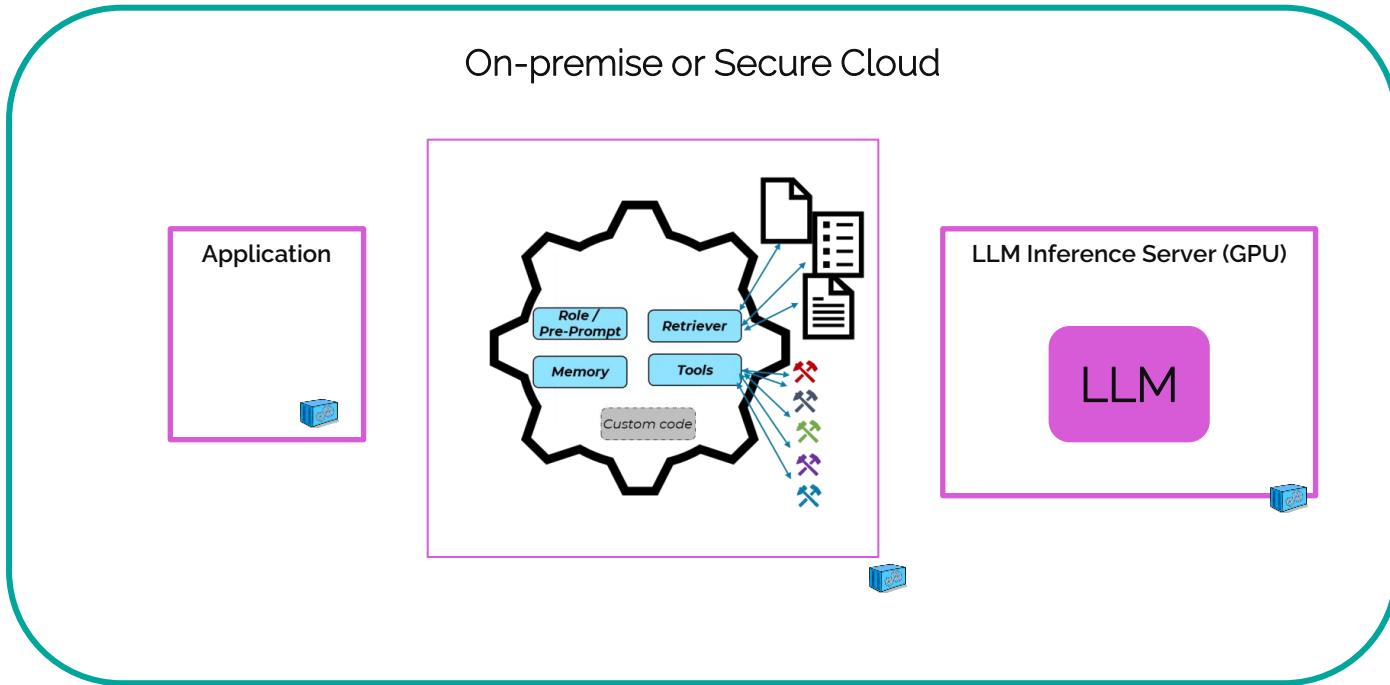


- **Risk attention points:**
  - Data shared in the prompt can be used for improving the model
  - Most of Generative AI are hosted in the US ➔ risk with Cloud Act and GDPR
- **Personal data / confidential information** - you should NOT input **personal data** or **confidential information** (internal or customer) in the prompt, except in private specific and secured set-up environment.

# Cloud Sovereignty for GenAI Inference



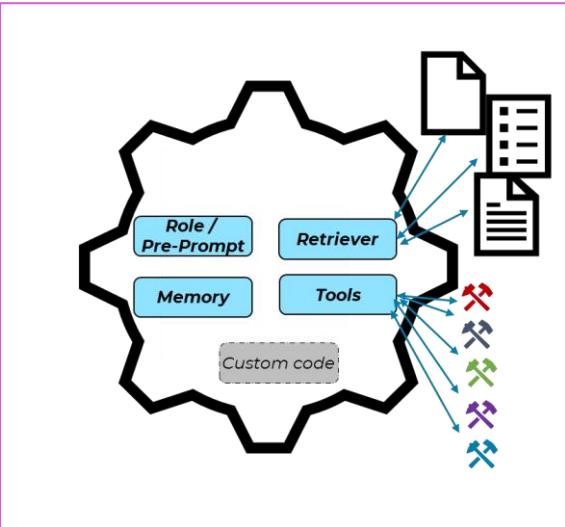
# Implementing Sovereignty : Approach 1



# Implementing Sovereignty : Approach 2

Application

On Premise

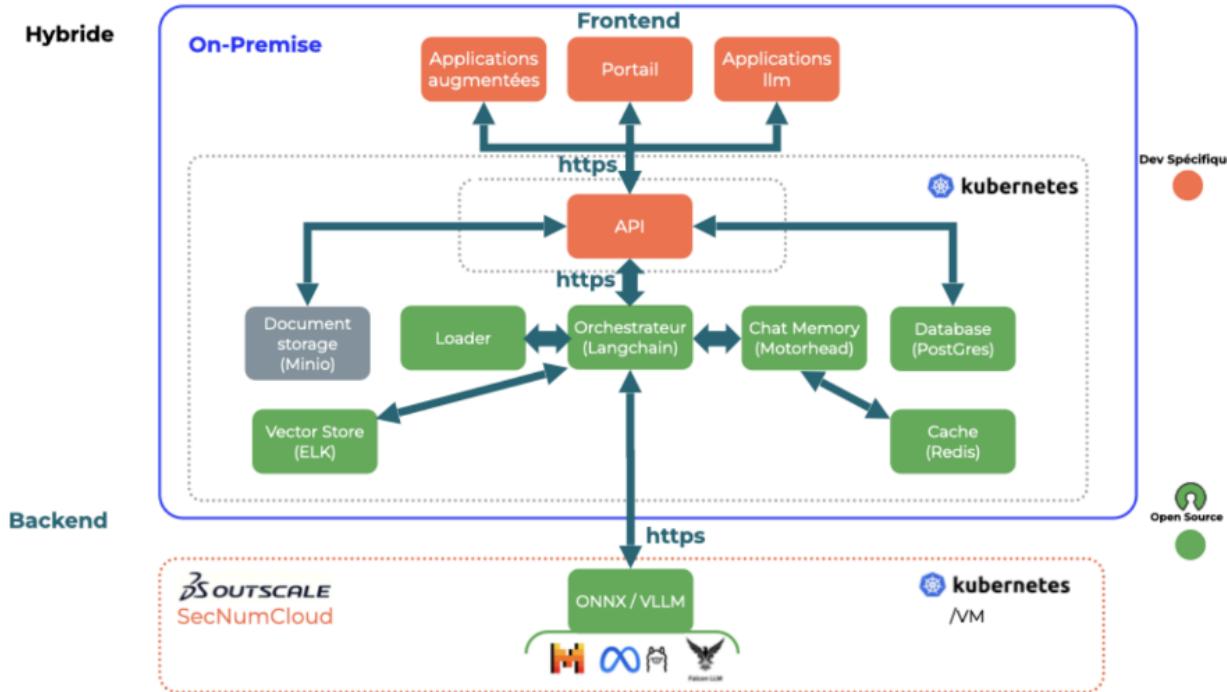


Secure Cloud

LLM Inference Server (GPU)

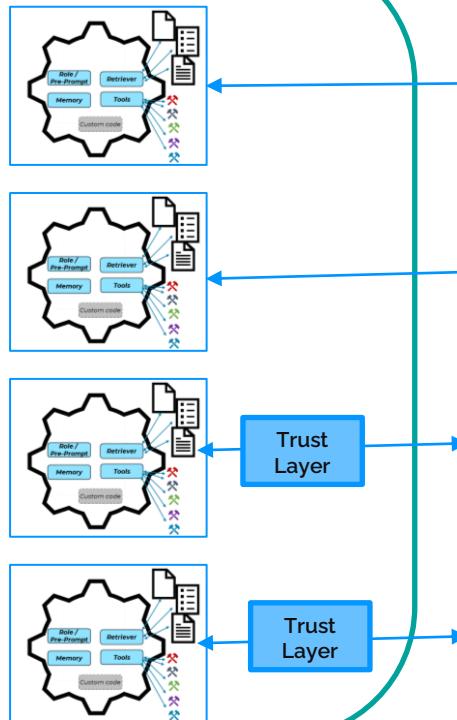
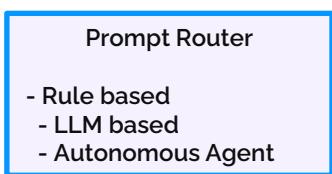
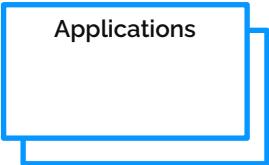
LLM

# Example Architecture:



# Suggested Approach (simplified)

On-premise



On-premise



SecNumCloud



Controlled Cloud



Public Cloud



# Optimisation & monitoring



# Caching

## LangChain implement several cache mechanism

- Exact Cache (on RDMS, Redis, NoSQL database, ....)
- Semantic Cache : return a hit if it finds a cached entry that is similar enough to the query (need embeddings)
- See [https://python.langchain.com/v0.1/docs/integrations/lms/llm\\_caching/](https://python.langchain.com/v0.1/docs/integrations/lms/llm_caching/)

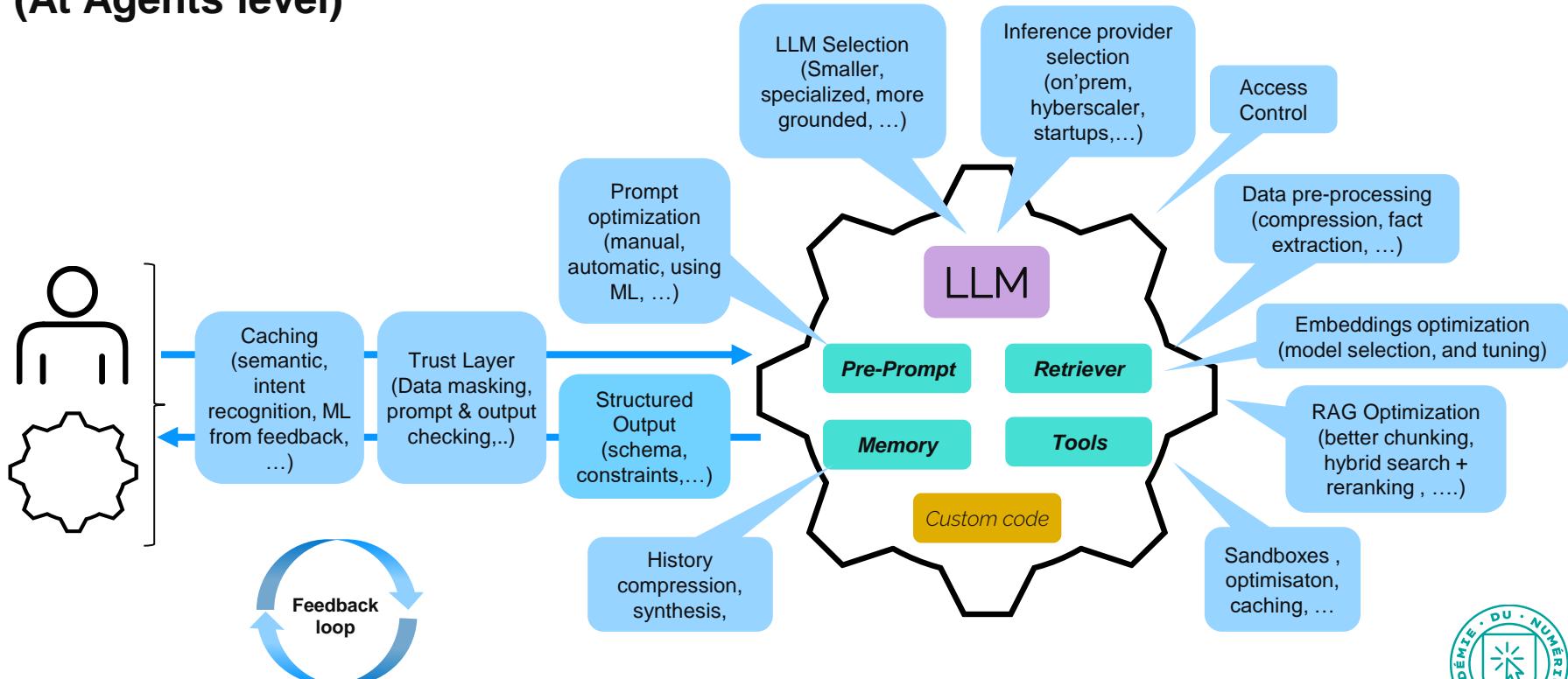
## More advanced strategies for caching

- Knowledge Graph
- Static routers (Logical, Classifiers...)
- Rules, intent recognition



# Recap : Risk mitigation and cost reduction

## (At Agents level)



# Industrialisation, Monitoring and Evaluation



# Agents Industrialisation

## Example LangChain Blueprint

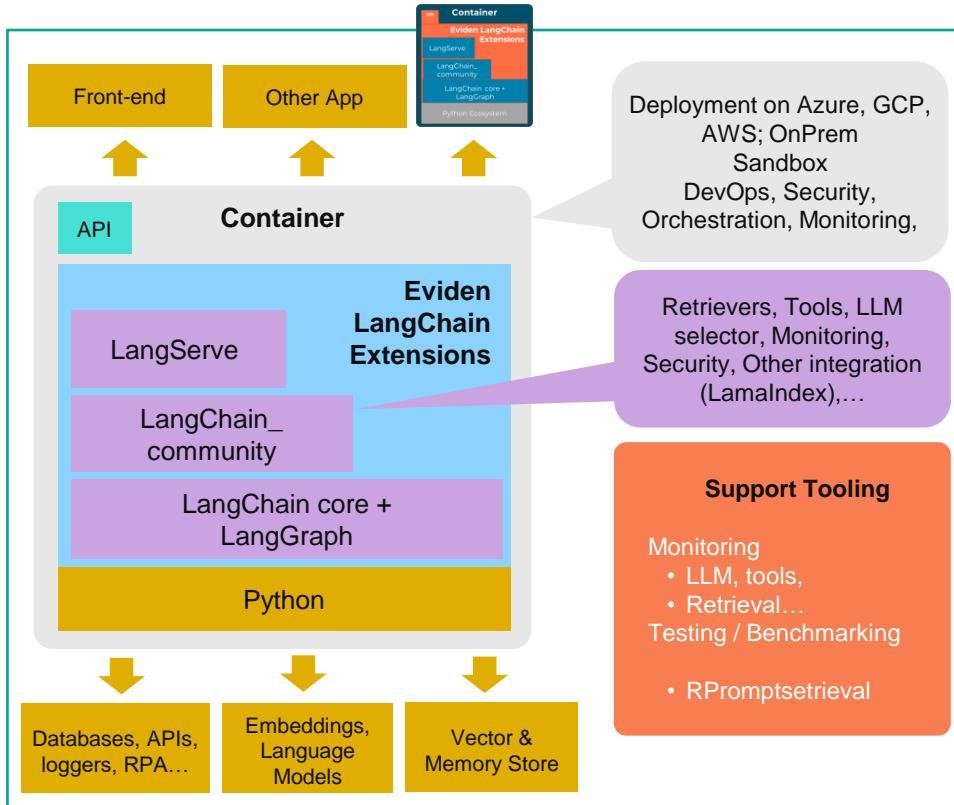


Diagramme avec  
OpenGPTs  
Eviden extension  
Security  
Deployment  
Monitoring  
Evaluation

# Monitoring: LangSmith and alternatives

## Rapid Prototyping

- Quickly test different prompts, models, and parameters.
- Debug issues with clear visibility into each step of an LLM sequence.

## Testing

- Create datasets of inputs and reference outputs.
- Run tests on your LLM applications.
- Score test results with custom evaluations.

## Performance Monitoring

- Log all traces to monitor application performance.
- Visualize latency and token usage statistics.

## Experimentation

- Use the playground environment for rapid iteration.
- Quickly test out different prompts and models.

## Continuous Improvement

- Track and diagnose regressions in test scores across multiple revisions of your application.
- Compare different versions of application side-by-side.

## Alternatives

- Lunary (open source)
- PromptLayer
- Arize
- Helicone
- LangFuse
- Weights & Biases (known as Wandb)
- ...

> Make your own !



# LangSmith

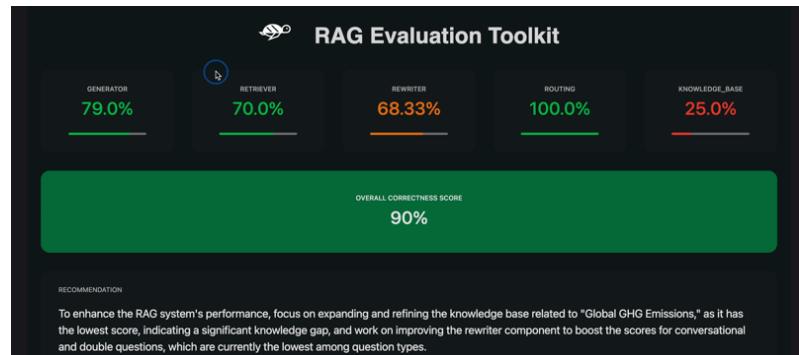
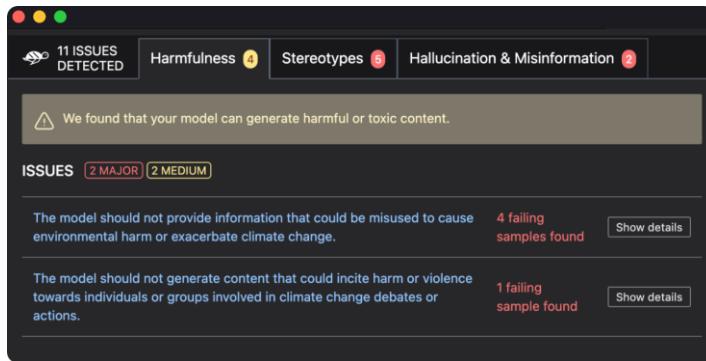
- Simplifies the prototyping of LLM application and agents. However, transitioning these prototypes into production-grade applications.
- Facilitates the creation and management of datasets for tasks such as fine-tuning, few-shot prompting, and evaluation, ensuring comprehensive model training and assessment.
- Capture detailed production analytics to gain insights into product performance and user interactions, supporting continuous improvements and optimizations.
- Enables the running of regression tests on your applications, providing confidence in ongoing development and ensuring that changes do not introduce new issues.
- As an open-source tool, it benefits from community involvement, ensuring continuous improvements and accessibility.

The screenshot displays the LangSmith application interface. On the left, a 'TRACE' panel shows a hierarchical timeline of events for a 'chat-langchain' session. The timeline includes steps like 'FindDocs' (0.70s), 'RetrievalChainWith...' (0.68s), 'Retriever' (0.68s), 'GenerateResponse' (5.07s), and 'ChatOpenAI' (4.33s). To the right of the trace, a 'chat-langchain' window shows a conversation between a 'HUMAN' and an 'AI'. The human asks, "How do I use a RecursiveUrlLoader to load content from a page?", and the AI responds with instructions on how to initialize and use the RecursiveUrlLoader. At the bottom of the interface are buttons for 'Add to Dataset', 'Annotate', and 'Playground'.

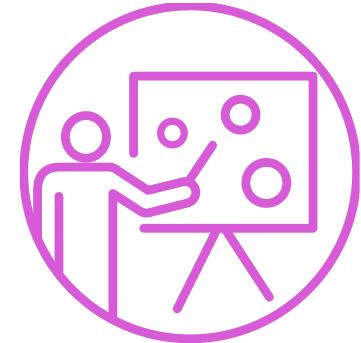


# Giskard

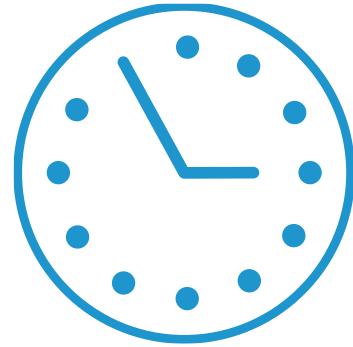
- **AI Testing Library**: Provide a robust library for testing machine learning models and large language models.
- **Automated Testing**: Automatically runs exhaustive test suites to identify risks such as performance issues, bias, and security vulnerabilities.
- **Integration with other tools**: Langchain, HuggingFace, Pytorch, Tensorflow, etc.
- **Efficient Model Evaluation**: Wraps around a models and datasets to facilitate comprehensive evaluations.
- **Risk Detection**: Detects critical risks in AI models, including hallucination, misinformation, harmful content, prompt injection, information disclosure, robustness issues, and discrimination.
- **CI/CD Integration**: Automates the publication of evaluation reports within your CI/CD pipeline for continuous integration and deployment.
- **Ease of Use**: Simplifies the testing process, enabling AI teams to quickly assess model quality and security.



# Wrap up



# Quiz



Chapter  
**Advanced Topics**

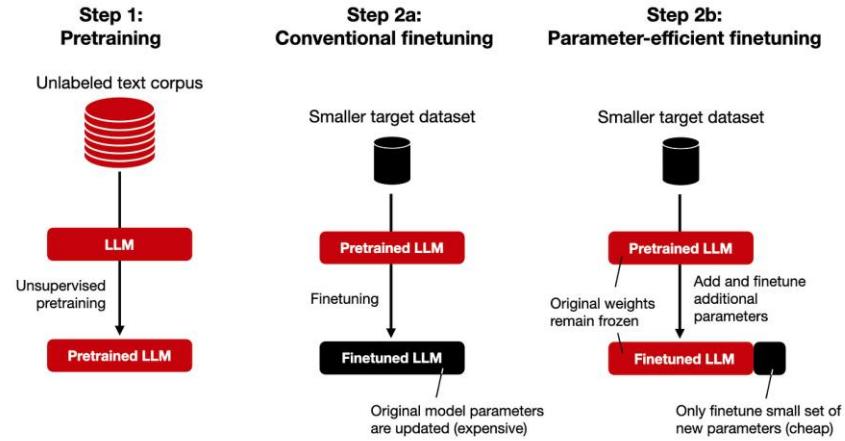


# Fine Tuning



# Definition and Purpose

- **Fine-tuning is a critical process** in the lifecycle of a language model. It involves taking a pre-trained model and further training it on a smaller, task-specific dataset.
- **This allow the model to learn** the nuances and vocabulary of the specific domain, leading to improved performance on related tasks.
- **Training adjustments** includes using a lower learning rate and fewer epochs to refine the existing knowledge.
- **The most specific use case of fine-tuning are:**
  - Customer Support Chatbots
  - Medical Text Analysis
  - Sentiment Analysis
- **The purpose can be resume on 3 points:**
  - Task-Specific Adaptation
  - Performance Improvement
  - Resource Efficiency



# Benefits

## Enhanced Accuracy and Relevance

- Fine-tuning a model for a specific tasks aligns its predictions more closely with the unique requirements and nuances of those tasks.
- This specialization results in outputs that are not only more accurate but also more contextually relevant, improving the overall quality and applicability of the model responses.

## Domain-Specific Proficiency

- By fine-tuning on a specific dataset to a certain domain, such as legal documents or medical literature, the model learns the specialized vocabulary, idioms, and contextual cues of the field.
- This leads to a deeper understanding and improved performance in generating domain-specific outputs, making the model valuable tool for experts in the field.

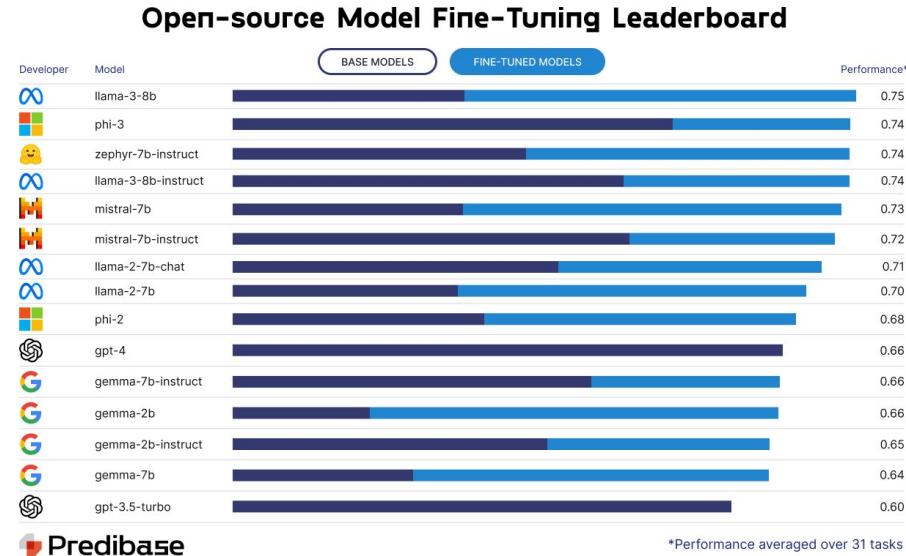
## Resource Efficiency and Cost Effectiveness

- The process leverages the extensive pre-existing knowledge of a pre-trained model, requiring only incremental adjustments.
- This significantly reduces the computational power and time needed for training.
- As a result, fine-tuning is a cost-effective approach that accelerates the deployment of high-performance models without the substantial overhead of training from scratch.



# Fine Tuning Performance vs Base Models

Most of the fine-tuned open-source models surpass GPT-4 on given tasks:



[Open-source Model Fine-Tuning Leaderboard](#)

# The Process and How To

The fine-tuning process involves several key steps to adapt a pre-trained model to a specific task or domain.

## Prepare and Process a Dataset

- Collect data that is relevant to the specific task or domain. This could be customer support transcripts, medical records, legal documents.

## Preprocess Data

- Clean the data by removing any irrelevant information, normalizing text (e.g., converting to lowercase, removing special characters), and splitting into training and validation sets.
- Ensure the data is in a format suitable for training, such as plain text or CSV files.

## Tokenize and collate the dataset

- Tokenization involves splitting text into words, subwords, or characters and converting them into numerical format.
- Collate the tokenized data into batches, ensuring that each batch is of a uniform size. This step involves padding sequences to the same length within a batch.

## Setup the trainer

- Configure the training parameters such as learning rate, batch size, number of epochs, and weight decay
- Set up the optimizer like **AdamW** and learning rate scheduler if needed.
- Initialize the Trainer class or an equivalent in the chosen framework with the model, training arguments, training dataset, and evaluation dataset.



# The Process and How To

## Evaluate the Performance of the Fine-Tuned Model

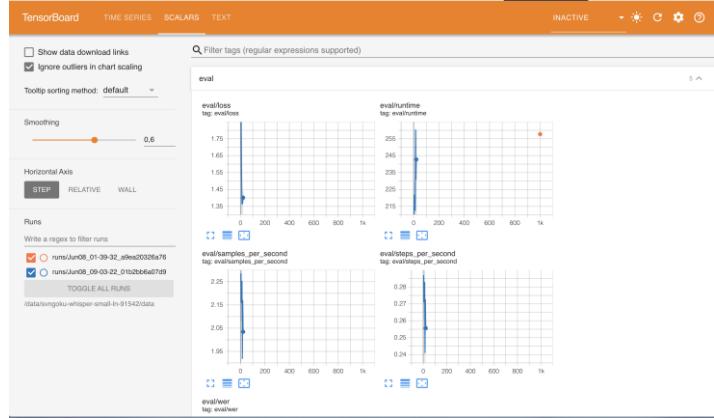
- Before fine-tuning, evaluate the performance of the pre-trained model on the task-specific validation set to establish a baseline.
- Metrics for evaluation may include accuracy, F1 score, loss, etc., depending on the task.

## Run the Fine-Tuning Process

- Train the model on task-specific dataset by running the training loop. Involves forward propagation, loss calculation, backpropagation and parameters updates.
- Monitoring training metrics such as loss and accuracy to ensure the model is learning effectively.

## Evaluate the Performance of the Fine-Tuned Model

- After fine-tuning, evaluate the model's performance on the validation set again to measure improvements.
- Compare the results with the baseline evaluation to assess the effectiveness of the fine-tuning process.
- Fine-Tuning may also involve hyperparameter tuning and iterative training to achieve the best performance.



# SFT: Supervised Fine-Tuning vs Basic Trainer

## • Trainer

- Designed for training models from scratch on supervised learning tasks (e.g text classification, question answering, summarization).
- Extensive configuration options for hyperparameters, optimizers, schedulers, logging, and evaluation metrics.
- Typically needs larger datasets for effective training from scratch.

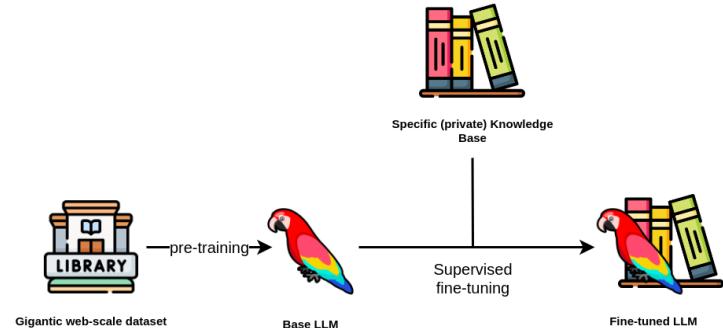
## • SFTTrainer

- Optimized for fine-tuning pre-trained models with smaller datasets on supervised learning tasks.
- Achieves comparable or better accuracy with smaller datasets and shorter training times.
- Use techniques like parameter-efficient fine-tuning (PEFT) to reduce memory consumption.

## • Choosing Between Trainer and SFTTrainer

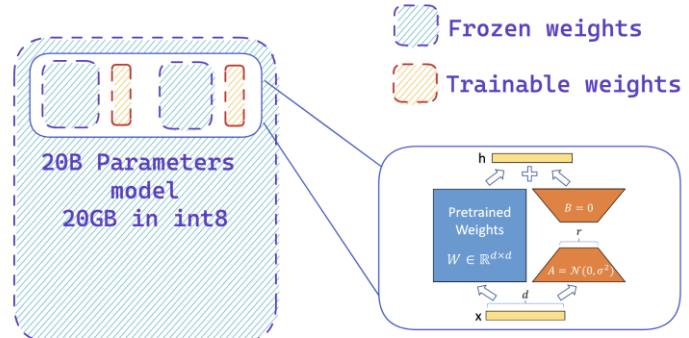
- Use Trainer
  - Large Dataset
  - Need extensive customization for training loops or complex workflows
- Use SFTTrainer
  - Pre-trained model with a smaller dataset
  - Simpler and Faster fine-tuning with efficient memory usage.

The choice depends on specific needs and resources. Experiment with both to see which works best for your task.



# PEFT : Parameter Efficient Fine-Tuning

- Challenges of traditional fine-tuning
  - Full fine-tuning of large models is computationally expensive and memory-intensive.
  - Fine-tuned models are the same size as the original, leading to high storage and deployment costs.
- PEFT fine-tunes only a small number of additional parameters while freezing most of the pretrained LLM parameters.
- Advantages of PEFT
  - Greatly reduces computational resources and storage requirements
  - Freezing most parameters prevents forgetting previously learned requirements
  - Performs well with limited data and generalizes better to new domains.
- Application Examples
  - Add Low-Rank matrices to existing parameters, efficiently fine-tuning large
  - Applicable to various tasks like image classification and generative tasks.
- Future Directions
  - 8bit, 4bit and 1bit training and tuning RLHF components.
- PEFT approaches enable efficient and cost-effective fine-tuning of LLMs maintaining high performance while reducing source requirements



# Efficient Fine-Tuning Techniques

## Low-Rank-Adaptation ( LoRA)

- LoRA reduces the number of trainable parameters by factorizing weight matrices into low-rank matrices, which allows efficient fine-tuning without significantly compromising performance
- Decreases the memory footprint, allowing for the fine-tuning of LLM with limited resources
- By fine-tuning fewer parameters, it maintains the model's pre-existing knowledge
- Effective for NLP tasks, image classification and generative tasks.

## Quantized Low-Rank-Adaptation ( QLoRA)

- QLoRA combines quantization with low-rank adaptation to further reduce the model's memory footprint and computational requirements
- Quantization reduces the precision of the model parameters, significantly decreasing memory usage.
- Maintains balance between reduced precision and effective fine-tuning.
- Suitable for deploying models on resource-constrained devices and for tasks requiring efficient computation.

# PEFT : Common Benefits

Transforming user queries into optimized search queries to enhance retrieval performance and accuracy

## Benefits

- **Resource Efficiency:** Both techniques enable the fine-tuning of large models on standard hardware, saving computational and storage costs.
- **Improved Performance :** Achieve comparable performance to full fine-tuning with a fraction of parameters.
- **Wide Applicability :** Applicable across various domains and tasks, enhancing the adaptability of large language models.

## Conclusion

- **LoRA and QLoRA** offer powerful and efficient methods for fine-tuning LLMs making advanced AI accessible while maintaining high performance.



# Accelerate Fine-Tuning 2x with Unsloth

We can further accelerate QLoRA / LoRA (**2x faster, 60% less memory**) using the **unsloth** library that is fully compatible with SFTTrainer.

- **30x Faster Training** : Reduce an 85-hour LLM training session down to just 3 hours, enabling us to train models like ChatGPT in 24 hours instead of 30 days.
- **Automatic Mixed Precision**: Employs mixed precision training to allocate resources efficiently.
- **Maintains Accuracy** Unsloth achieves these speedups without sacrificing accuracy with 0% loss or even up 20% increased accuracy compared to regular methods.
- **Hardware Support** : It supports a variety of GPUs from NVIDIA, Intel, and AMD, including T4 GPUs used on platforms like Google Colab, without needing specialized hardware.



1 A100 40GB	🤗Hugging Face	Flash Attention	_GPU_ Unsloth Open Source	_GPU_ Unsloth Pro
Alpaca	1x	1.04x	1.98x	15.64x
LAION Chip2	1x	0.92x	1.61x	20.73x
OASST	1x	1.19x	2.17x	14.83x
Slim Orca	1x	1.18x	2.22x	14.82x

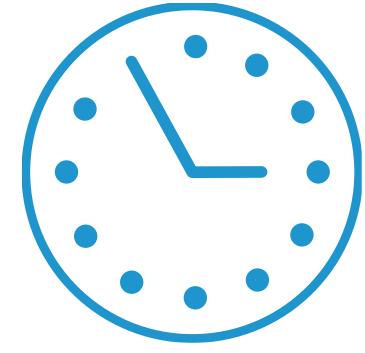
• Benchmarking table below was conducted by [🤗Hugging Face](#).

Free Colab T4	Dataset	🤗Hugging Face	Pytorch 2.1.1	_GPU_ Unsloth	GPU VRAM reduction
Llama-2 7b	OASST	1x	1.19x	1.95x	-43.3%
Mistral 7b	Alpaca	1x	1.07x	1.56x	-13.7%
Tiny Llama 1.1b	Alpaca	1x	2.06x	3.87x	-73.8%
DPO with Zephyr	Ultra Chat	1x	1.09x	1.55x	-18.6%

# Practicum



Demo: fine-tuning of a LLM for a SQL base Assistant.



# Model Merging



# Definition

- Model merging combines two or more LLMs into a single model
- It's an experimental technique for creating new models cheaply without requiring a GPU
- Has produced state-of-the art models on OpenLLM Leaderbord
- **mergekit** is a Merges can be run entirely on CPU or accelerated with as little as 8 GB of VRAM. Many merging algorithms are supported, with more coming as they catch my attention.
- Avoid the high cost of training from scratch by leveraging existing models

Model	Average
<a href="#">zyh3826/GML-Mistral-merged-v1</a>	73.3
<a href="#">shadowm1/Marcoro14-7B-slerp</a>	73.01
<a href="#">cookinai/CatMacaroni-Slerp</a>	72.74
<a href="#">AIDC-ai-business/Marcoroni-7B-v3</a>	72.53
<a href="#">Toten5/Marcoroni-v3-neural-chat-v3-3-Slerp</a>	72.51
<a href="#">Toten5/Marcoroni-neural-chat-7B-v2</a>	72.5
<a href="#">ignos/Mistral-T5-7B-v1</a>	72.47
<a href="#">EmbeddedLLM/Mistral-7B-Merge-14-v0.1</a>	72.39
<a href="#">mindy-labs/mindy-7b</a>	72.34
<a href="#">janhq/superimario-v2</a>	72.34

*Marcoro14-7B-slerp* became the best-performing model on the Open LLM Leaderboard

# Merging Methods in MergeKit

- MergeKit support many methods of merging
  - Linear**: The classic merge method with a simple weighted average
  - SLERP** ( Spherical Linear Interpolation ) : Smoothly interpolates between two vectors, maintaining constant change rate and geometric properties
  - TIES**: Merge multiple task-specific models by reducing redundancy and resolving parameter conflicts
  - DARE**: Similar to TIES including pruning and rescaling of weights. The parameters *dare\_ties* or *dare\_linear*
  - Passthrough**: Concatenates layers from different LLMs, creating models with a unique number of parameters.
  - Task Arithmetic**: Computes 'task vectors' for each model by subtracting a base model. Works great for models that were fine-tuned from a common ancestor.

Method	merge_method value	Multi-Model	Uses base model
Linear ( <a href="#">Model Soups</a> )	linear	✓	✗
SLERP	slerp	✗	✓
<a href="#">Task Arithmetic</a>	task_arithmetic	✓	✓
<a href="#">TIES</a>	ties	✓	✓
<a href="#">DARE TIES</a>	dare_ties	✓	✓
<a href="#">DARE Task Arithmetic</a>	dare_linear	✓	✓
Passthrough	passthrough	✗	✗
<a href="#">Model Breadcrumbs</a>	breadcrumbs	✓	✓
<a href="#">Model Breadcrumbs + TIES</a>	breadcrumbs_ties	✓	✓
<a href="#">Model Stock</a>	model_stock	✓	✓

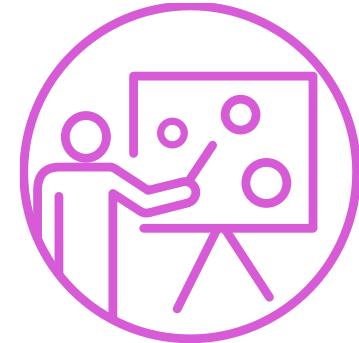
# Practicum



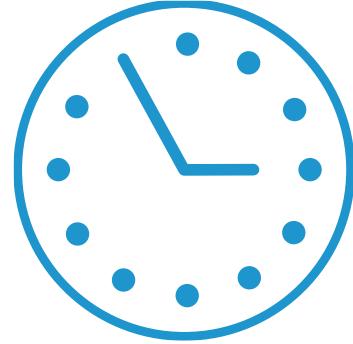
Demo: merging a foundation model with  
1 or 2 other specific models.  
Check the repo



# Wrap up



# Quiz





Félicitations !

Vous avez terminé le module  
**Agents GenAI**

© Eviden, 2024. Tous les droits sont réservés à Atos SE et/ou ses concédants éventuels. L'utilisation de ce matériel est limitée au contexte de l'Académie du Numérique. Son utilisation, reproduction, représentation et diffusion sont autorisées uniquement à des fins pédagogiques et réservés exclusivement aux participants à la formation. Les droits relatifs à toute adaptation ou autre modification du matériel susceptibles de constituer une œuvre collective ou dérivée appartiennent exclusivement à Eviden. Toute utilisation à des fins commerciales est expressément exclue.