

# Generative Conversational AI

## TECHNOLOGIES

Diego Gosmar





# Diego Gosmar

**Chief AI Officer**



**Start-up Mentor & Advisor**



**Ambassador**





## Diego Gosmar



Chief AI Officer, XCALLY.ai

Verified email at ieee.org - [Homepage](#)

Artificial Intelligence Conversational AI Conversation Analysis

FOLLOWING



Google Scholar

TITLE :

CITED BY YEAR

- [Conversational hyperconvergence: an onlife evolution model for conversational AI agency](#)

1 2024

D Gosmar  
AI and Ethics, 1-15

- [Conversational AI Multi-Agent Interoperability, Universal Open APIs for Agentic Natural Language Multimodal Communications](#)

2024

D Gosmar, DA Dahl, E Coin  
arXiv preprint arXiv:2407.19438

- [Insight AI Risk Detection Model-Vulnerable People Emotional Situation Support](#)

2024

D Gosmar, E Peretto, O Coleman  
Proceedings of the 28th International Conference on Evaluation and ...

ChatGPT, GenAI  
Real applications  
Image Recognition  
RNN to Transformers

Neural Networks  
Speech Analytics  
Conversational AI  
Sentiment Analysis

BLUE  
Edition

AI and Ethics

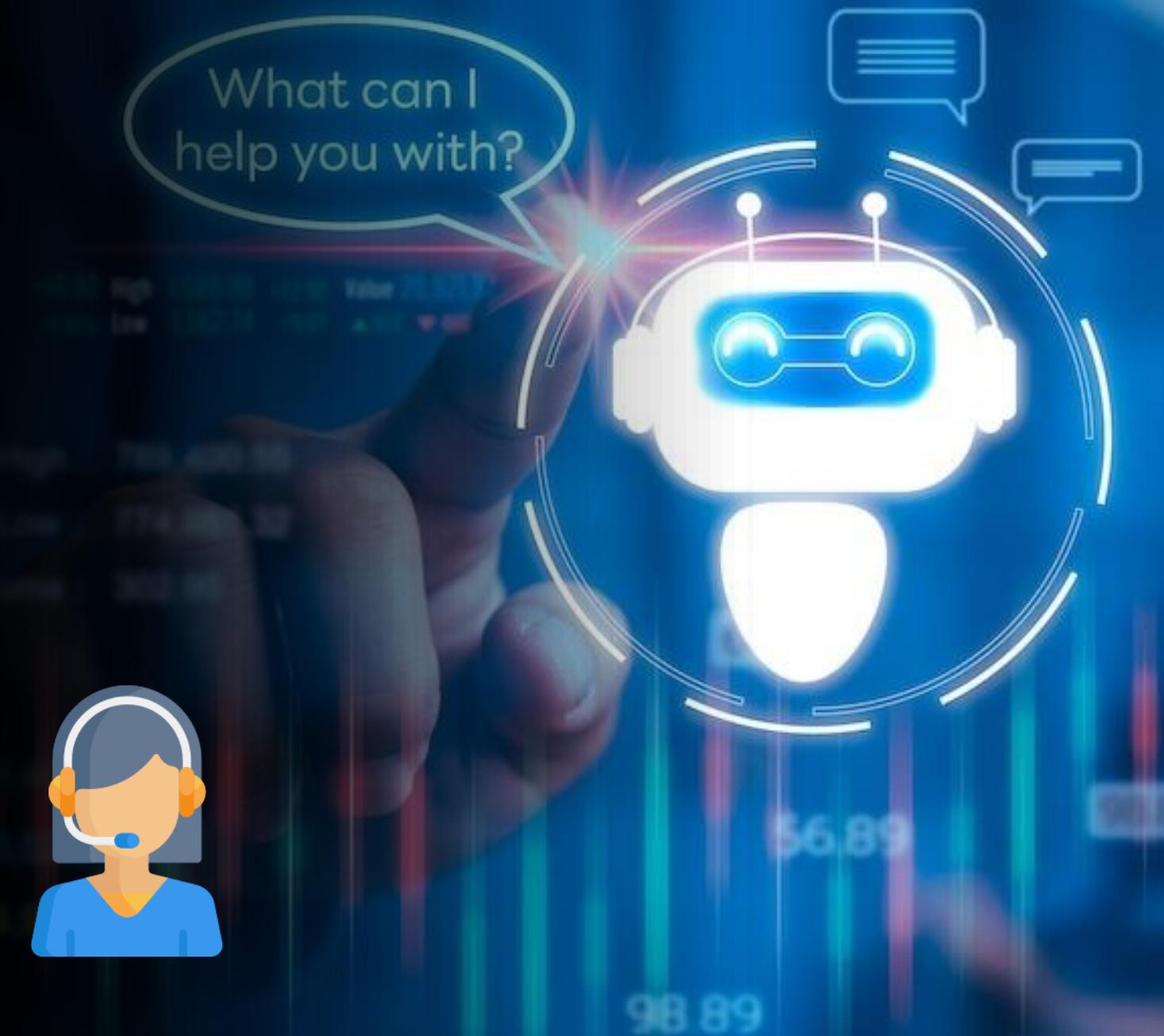




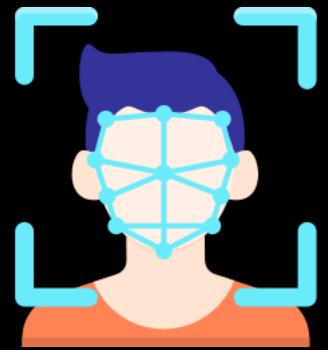
- AI intro
- Generative AI and Conversational AI intro
- LLM and RAG intro
- AI new models, self-supervised pretrained and fine-tuning
- LLM Key concepts and Embedding
- Demo exercises
- RAG
- Local examples
- Demo exercises
- RAG Framework example
- Demo exercises
- RAG with Open AI
- Demo exercises
- LLM, RAG management and monitoring
- OVON use case (smart library)

# Conversational AI

## What is it?



# Discriminative AI vs Generative AI



- **Discriminative AI** is capable of classifying data and making predictions based on predefined models, data patterns and historical data.
- **Generative AI** is capable of generating new information and content from provided datasets.

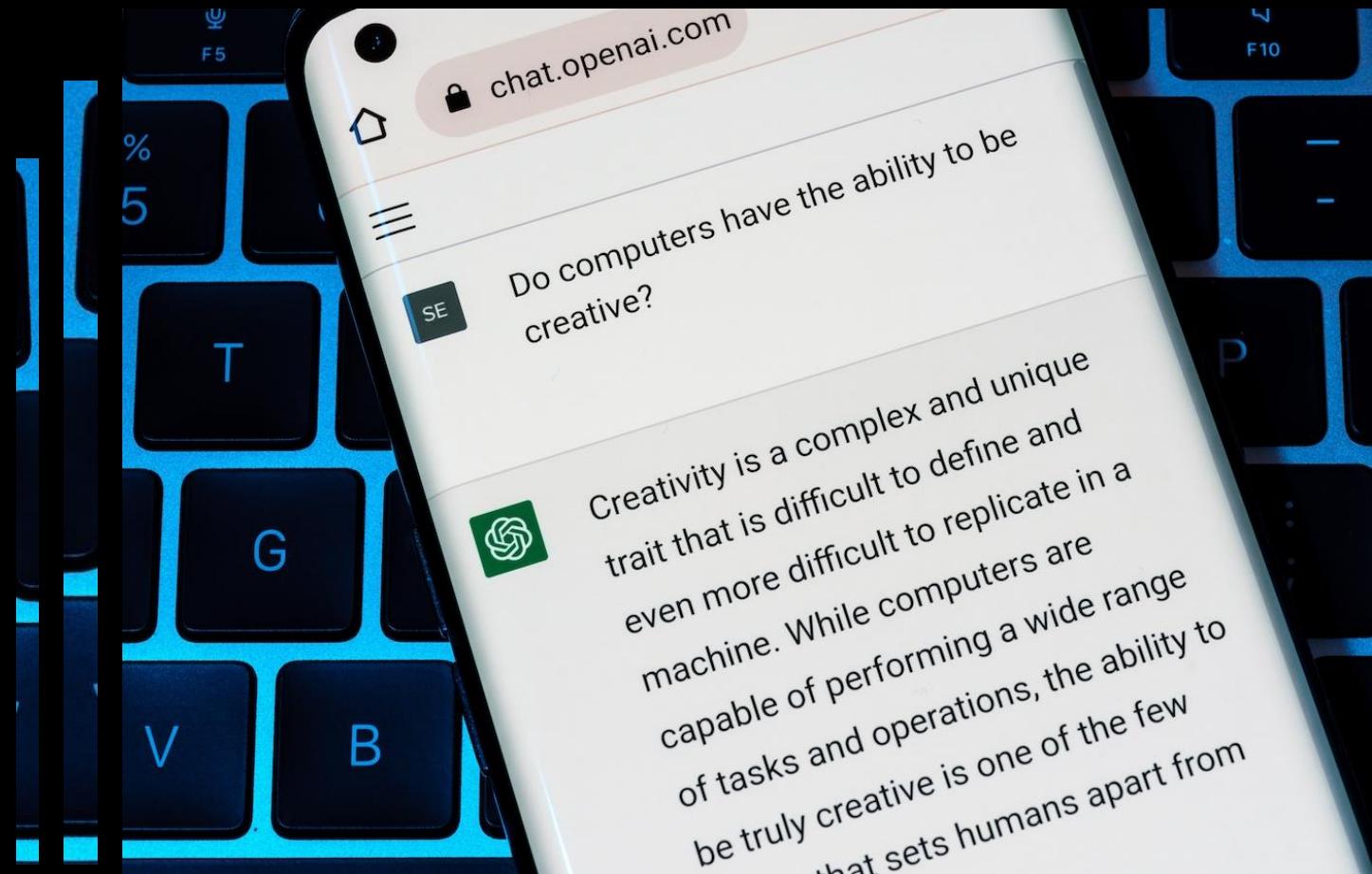
# Generative AI

Machine Learning models/AI solutions trained on a large amount of DATA in order to generate NEW DATA themselves

Often Multimodal, with different kind of input/output:

i.e. text, voice, images, video

LMM

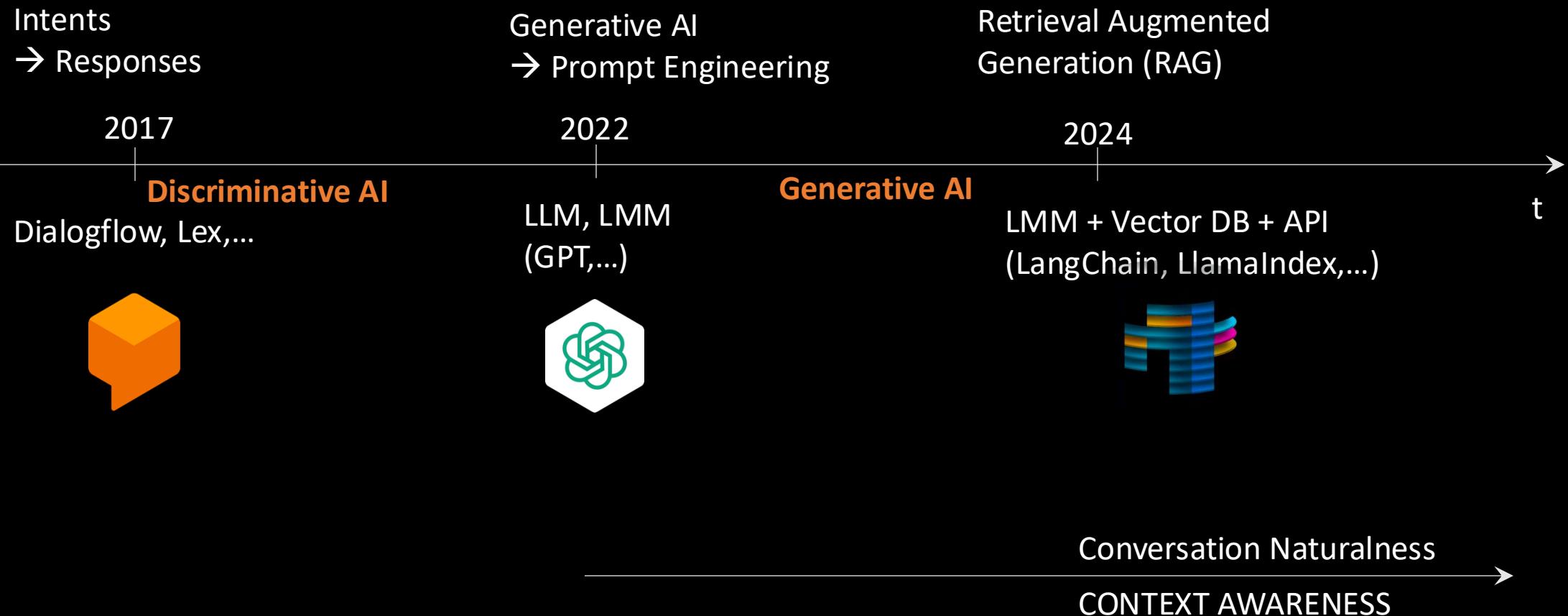


# Generative AI



# Conversational AI

## Evolution Timeline



**RAG** (Retrieval Augmented Generation)

**Fine Tuning** (Open-Weight LLM)

**LAM** (Large Action Models)

# RAG, Fine Tuning and LAM

CUSTOM AGENTS for everything!

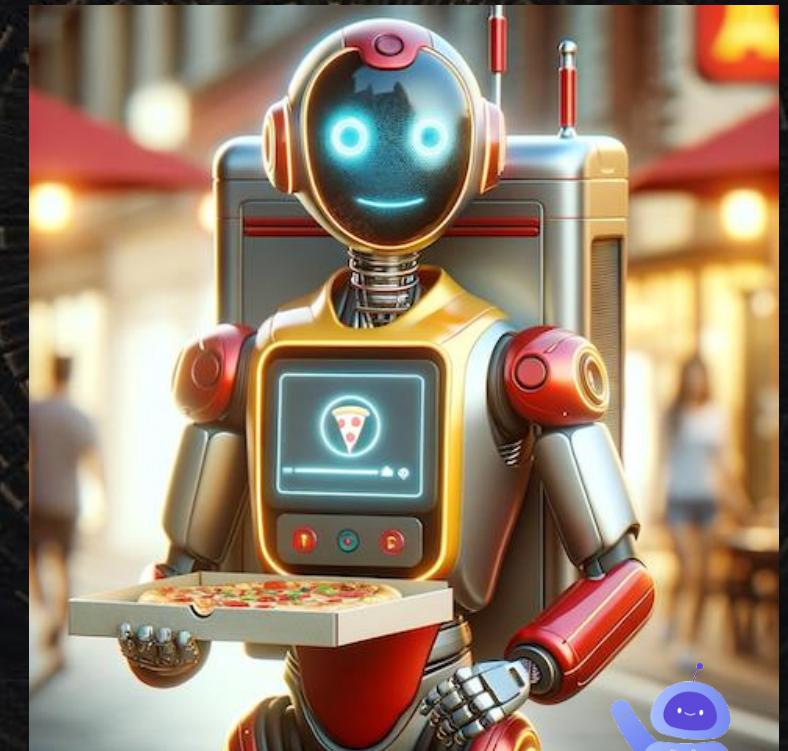
HEALTHCARE



FINANCE



FOOD ORDERING



# RAG

## Retrieval-Augmented Generation

### Prompt:

*what's the amount related to our latest higher invoices in 2024?*

### Context:

Invoice #1, amount: 1200 \$  
Invoice #2, amount: 750 \$  
...  
Invoice #N, amount: 2350 \$



### Role:

*You are an assistant agent specialized in billing information*

### Output:

*Sure, the following ones are the highest 2023 invoices:  
Invoice #5: 1850 \$  
Invoice #43: 2470 \$  
Etc... etc...*



**RAG** (Retrieval Augmented Generation)

**Fine Tuning** (Open-Weight LLM)

**LAM** (Large Action Models)

Function-Calling capabilities

# Fine-Tuning EXERCISES



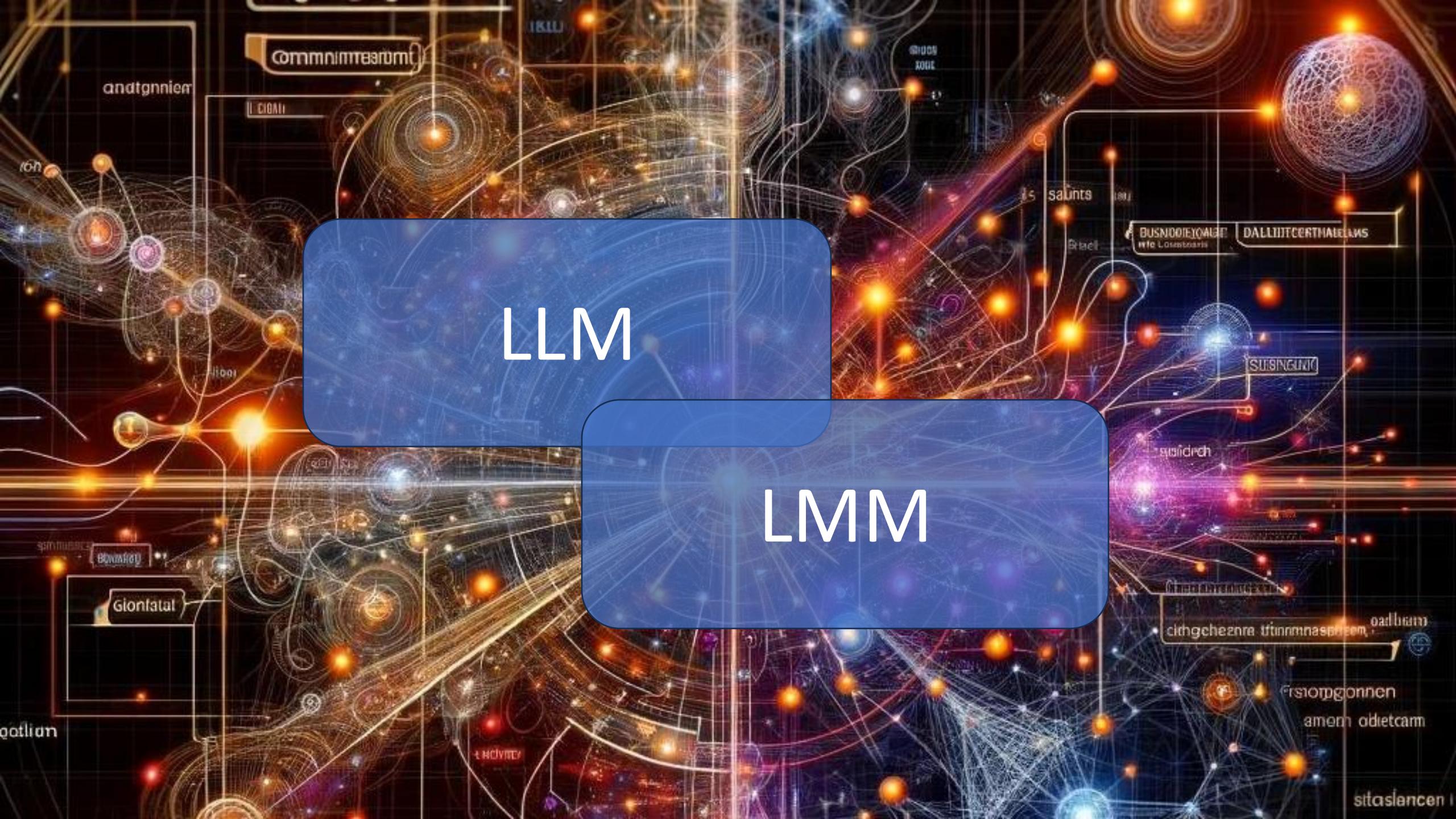
# Homeworks



<https://huggingface.co/models>

i.e. EleutherAI/gpt-neo-125m





LLM

LMM

# Language Model

A language model is a probabilistic model of a natural language



**Text generation:** (predict next words, utterances, sentences...)

Text Summarization

Translation

...

LMM: beyond text!

# LLM: training and tuning

## KEY concepts



- **Model Architecture**: This refers to the specific implementation and mathematical structure of the model. Most top-performing LLMs today use a variation of the "decoder-only" **Transformer architecture**.
- **Training Dataset**: The dataset comprises examples and documents used to train the model. These documents can be in natural languages (like English or Chinese), programming languages (like Python), or any structured text (like tables or equations). The training process involves the model learning patterns from these documents.

Exploratory



# LLM: training and tuning

## KEY concepts

- **Tokenizer:** This is a tool that converts text from the training dataset into numbers, which the model can process. It breaks down text into smaller units called **tokens**, which could be words, sub-words, or characters. The tokenizer's vocabulary size varies, indicating the number of different tokens it recognizes: typically between **32k and 200k**



*The size of an LLM dataset is often measured as the number of tokens it contains once split in a sequence of these individual, "atomistic" units: these days it ranges from **several hundred billion tokens** to **several trillion tokens!***

# LLM: training and tuning

## KEY concepts

- **Training Process:** Requires significant computing power and skilled personnel. It involves setting up the model architecture and running the training algorithm with the chosen **hyperparameters** on the training dataset. The outcome is a set of learned model **weights**, which are used for making predictions on new inputs.
- **Model Weights:** the **weights** can then be used for **inference**, i.e., for prediction on new inputs, for instance, to generate text.  
Pretrained LLMs can also be specialized or adapted for a specific task after pretraining, particularly **when the weights are openly released**. They are then used as a starting point for use cases and applications through a process called **fine-tuning**.



# LLM: training and tuning

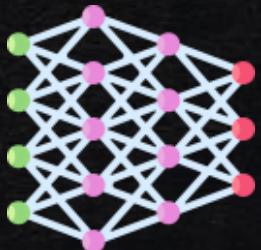
## KEY concepts

- **Pretraining and Fine-Tuning:** Pretrained LLMs can be further specialized for specific tasks through fine-tuning, which involves additional training on a more specialized dataset.  
This process is **less resource-intensive** than training a model from scratch and allows for the **adaptation of open-source pretrained models for various applications**, even with limited computing resources.



i.e. HuggingFace H4 Zephyr 7b beta model

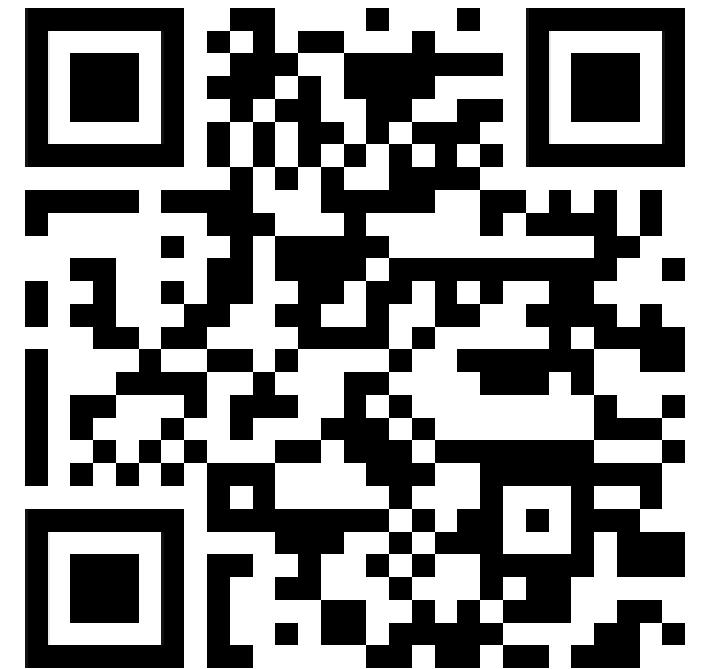
7.24 billion parameters



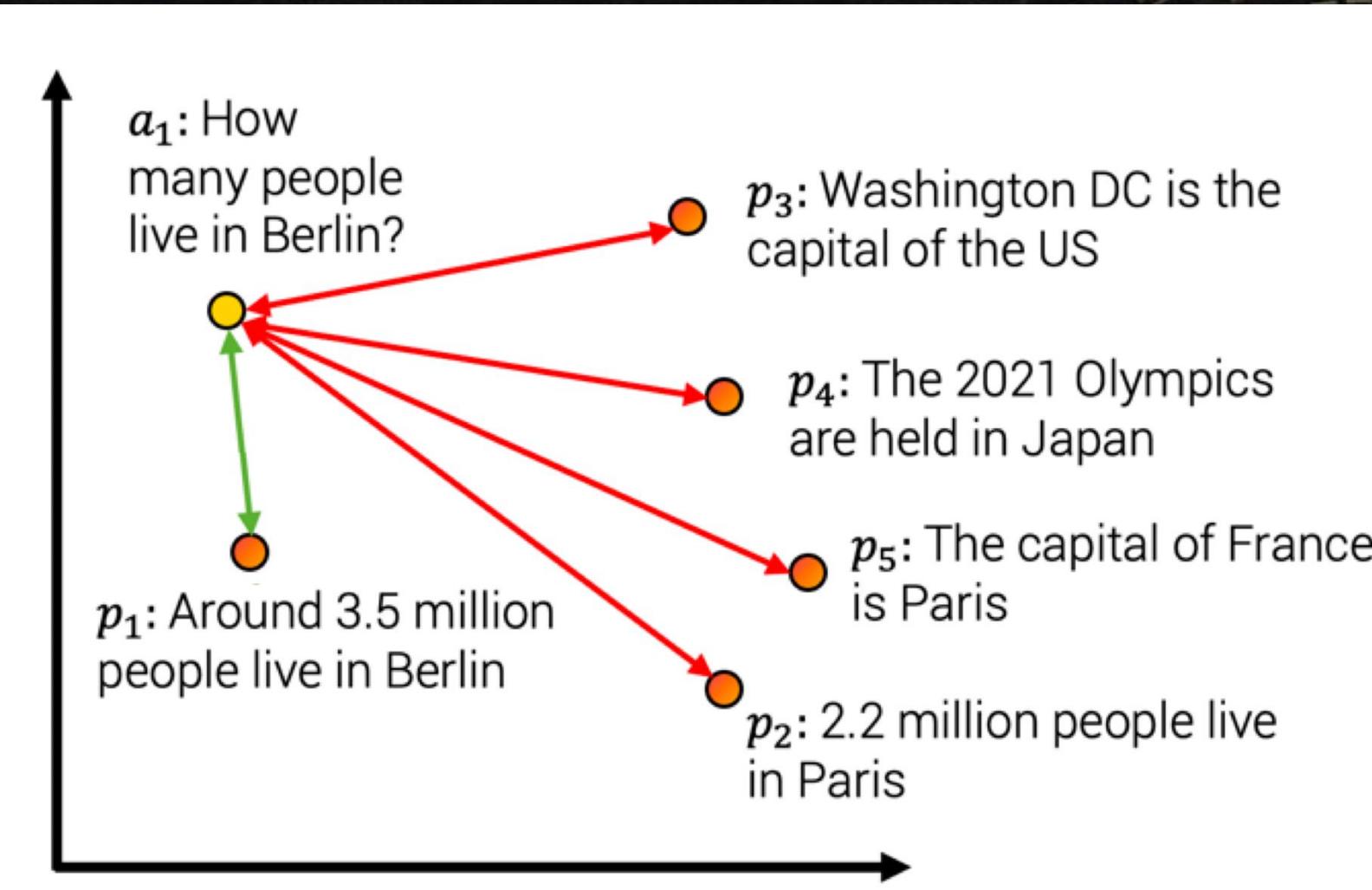
Fine-tuned version of [mistralai/Mistral-7B-v0.1](#)

Datasets: ● [stingning/ultrachat](#)

774K synthetic dialogues generated by ChatGPT



# Embedding



# EXERCISES



# LLM: training and tuning

## Visualize LLM Tokens

GPT2

GPT3.5 / GPT4

LLaMA

chakra of artificial intelligence." This book explores the field of machine learning and its significance in the realm of artificial intelligence. It reflects Gosmar's expertise and interest in innovative technologies, particularly in the context of customer experience and contact center solutions. The book delves into the complex aspects of machine learning, likely offering insights from Gosmar's extensive experience in the technology sector.

Tokens: 89 Characters: 516

Clear

Show example



# Tokenizer VS Embedding



**Tokenizer → Dictionary → Syntactic domain**

Syllabus, Word, part of text mapping rules to numbers!

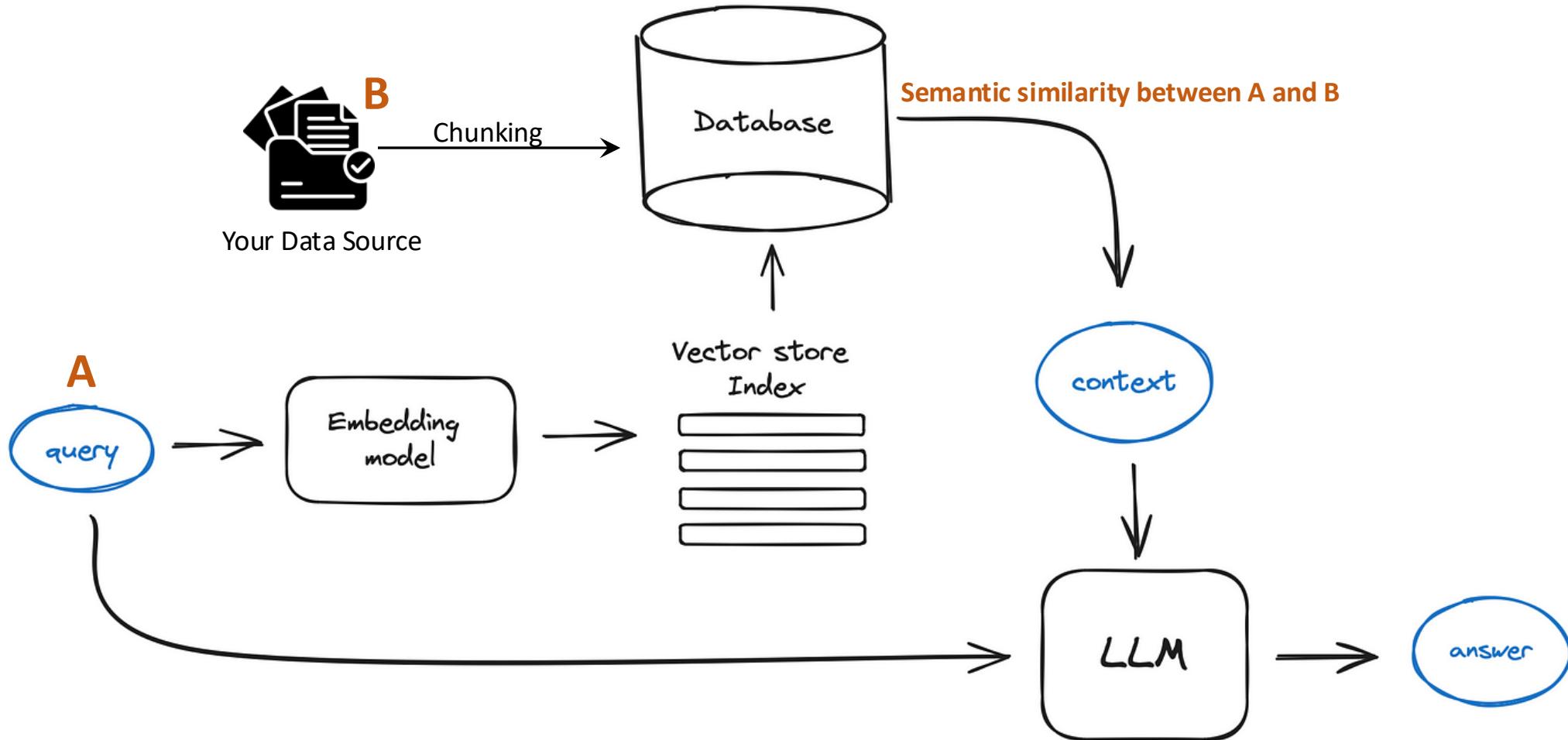
**Embedding → Vectorization → Semantic domain**

Text meaning, correlations, context,...

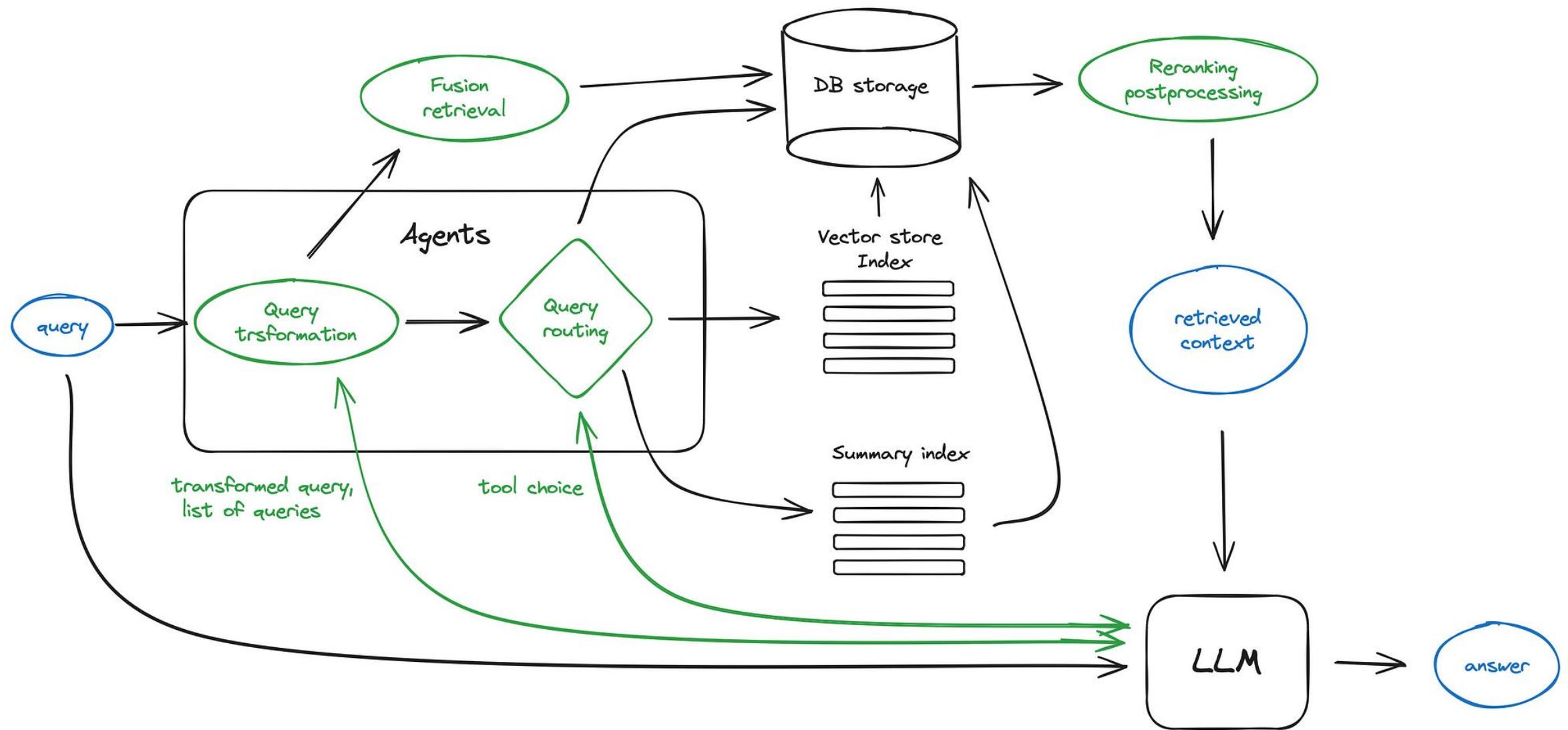
# RAG



## Naive RAG



# Advanced RAG

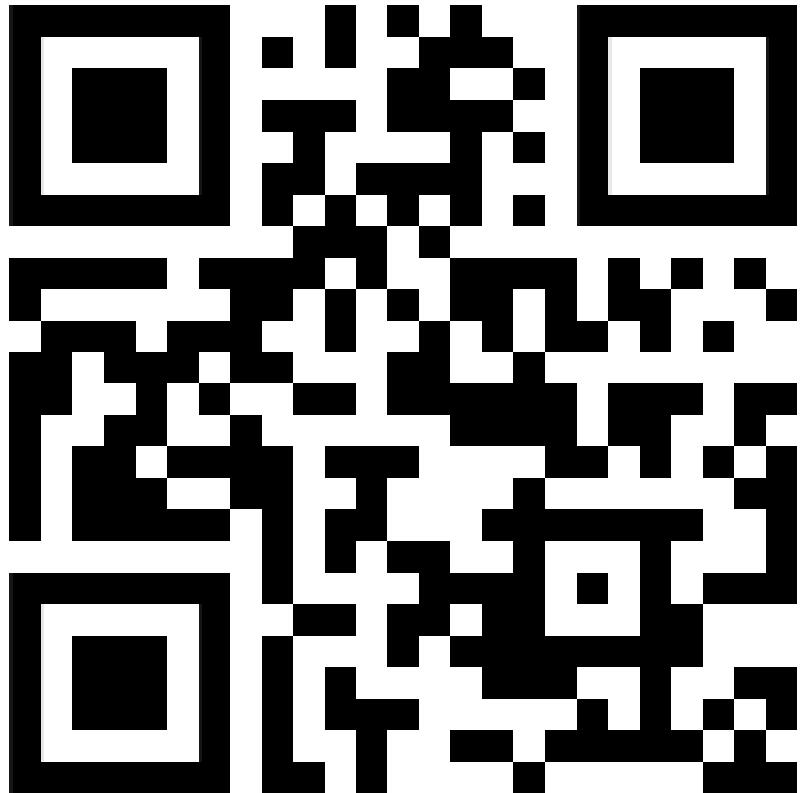


# EXERCISES



# ollama setup

<https://ollama.com>



# ollama cfg

- Run a model
- Pull a model
- Customize a Role for a Model
- ollama API

SEE NOTES



# RAG: ollama web app demo-kit

ACME LTD

## ORDER

Acme Ltd 123 Business  
Road, Tech

City Phone: 555-1234

Email:  
[contact@acmelt.com](mailto:contact@acmelt.com)

Order Date: Jan 5, 2024

Order Number:  
ACME-2024-005

To: AutoSpare Inc.  
456 Parts Avenue,  
Industrial Park

Phone: 555-5678

Email:  
[sales@autospareinc.com](mailto:sales@autospareinc.com)

Shipment Address:

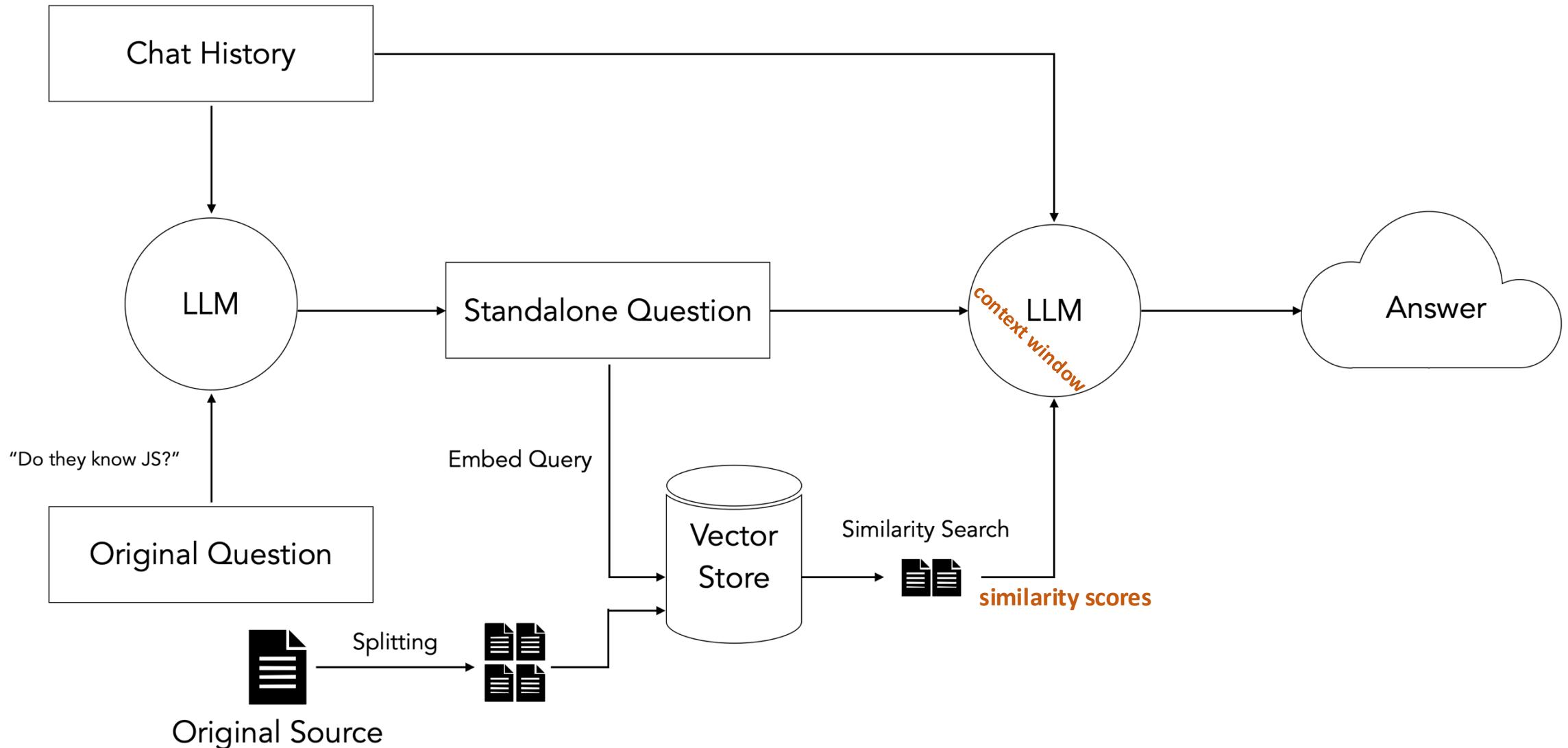
Acme Ltd Warehouse  
789 Storage Lane, Tech City

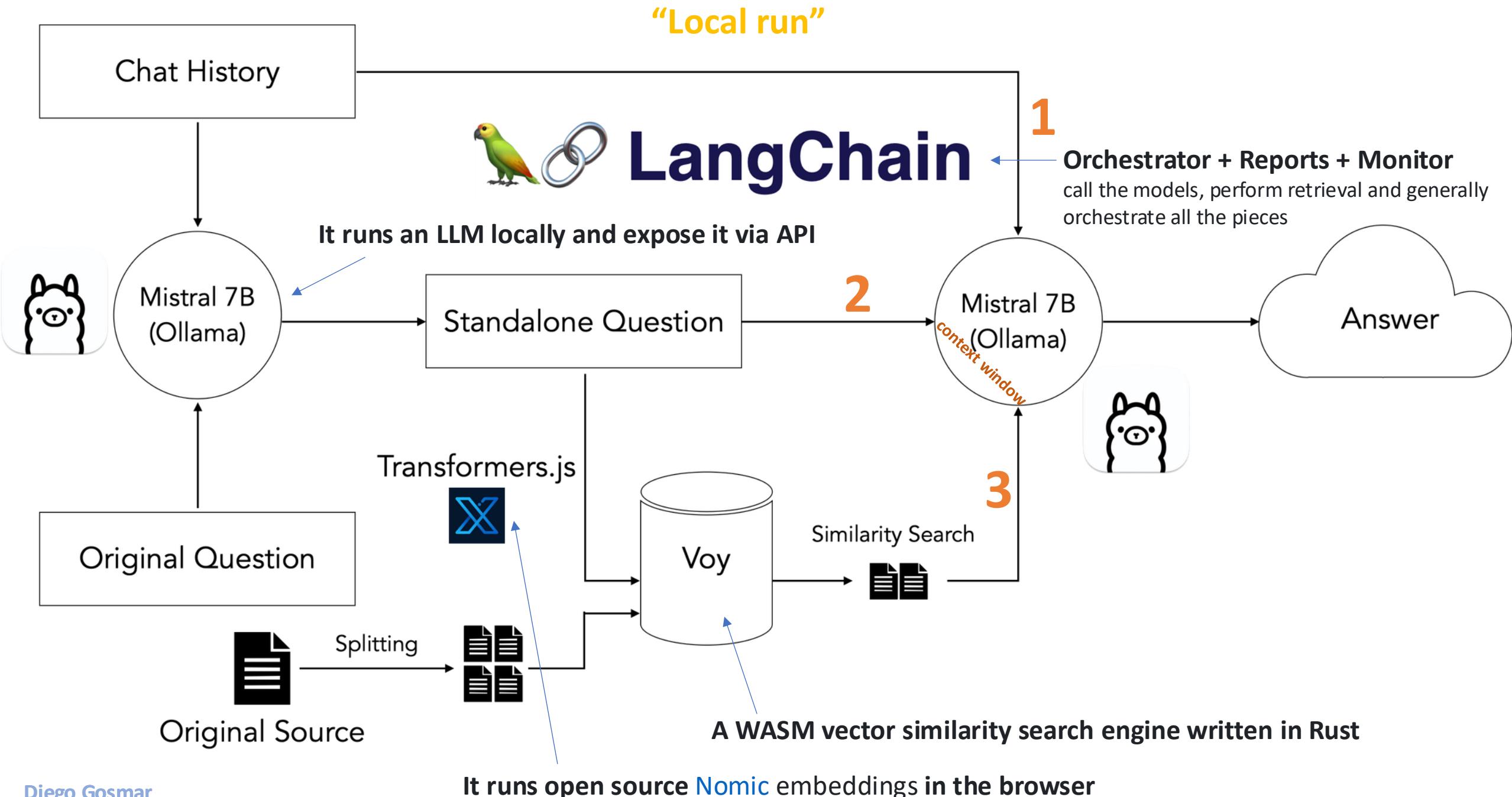
Billing Address:

Acme Ltd  
123 Business Road, Tech City

Description	Quantity	Unit Price	Cost
Car Brake Pads (Model A1)	50	30,00 US\$	1.500,00 US\$
Engine Oil Filter (Type B2)	100	5,00 US\$	500,00 US\$
Headlights (Model C3)	20	45,00 US\$	900,00 US\$
Windshield Wipers	200	10,00 US\$	2.000,00 US\$
Car Batteries (Model E5)	10	80,00 US\$	800,00 US\$
		Subtotal	5.700,00 US\$
	Tax	10,00%	570,00 US\$
		Total	6.270,00 US\$

# Architecture





# RAG: ollama web app demo-kit

## Local fork and setup:

```
clone https://github.com/jacobleee93/fully-local-pdf-chatbot.git
yarn
npm run-script dev
```

```
OLLAMA_ORIGINS=http://localhost:3000 OLLAMA_HOST=127.0.0.1:11435 ollama serve
OLLAMA_HOST=127.0.0.1:11435 ollama pull mistral
```

Access via: <http://localhost:3000/>



# RAG: ollama web app demo-kit

**Here are some links for the various pieces used in the app:**

- Demo app: <https://webml-demo.vercel.app/>
- Demo app GitHub repo: <https://github.com/jacobleee93/fully-local-pdf-chatbot>
- Voy: <https://github.com/tantaraio/voy>
- Ollama: <https://github.com/jmorganca/ollama/>
- LangChain.js: <https://js.langchain.com/>
- Transformers.js: <https://huggingface.co/docs/transformers.js/index>

# XCALLY demo

# Conversational AI & RAG

1/2

**Smart Order**

**Assistant**

- Info** Assistant settings
- Files** Associate files to the assistant

**Name\*** Smart Order

**Model\*** gpt-4-turbo-preview

**Instructions**

You are a professional AI agent specialized in Order information. You reply in details to any information about orders. In case the users ask us anything not related to orders, reply them kindly explaining that you are AI agent specialized in Order information, therefore you cannot fulfil their request.

```

graph LR
    Start((Start)) --> TTSWelcome[TTS Welcome message  
Hello]
    TTSWelcome --> ASR[ASR]
    ASR --> ChatGPT[OpenAI ChatGPT]
    ChatGPT -- Got off: false --> DefaultMessage[Default message if no speech]
    DefaultMessage --> Hangup((Hangup))
    ChatGPT -- Got off: true --> TTSEnding[true  
TTS ending message]
    TTSEnding --> Hangup
  
```

**VoiceBotDemo**  
demo for a smart agent assistant. - Created At 12/06/24 23:35:53

**Edit OpenAI ChatGPT**

**Label** OpenAI ChatGPT

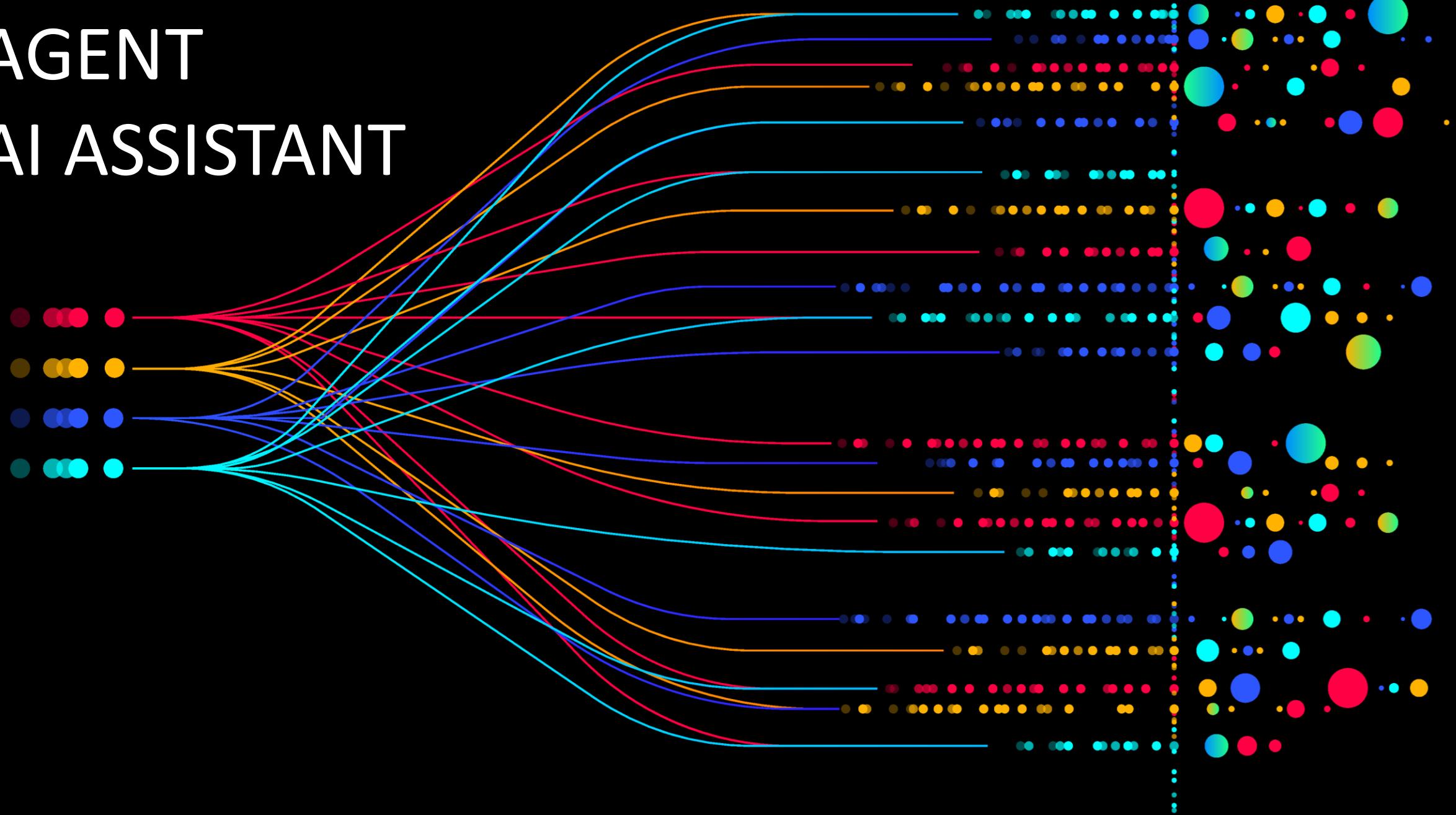
Chat\_Bot\_05-06-2024\_13:36:50  
Chat\_Bot\_05-06-2024\_12:49:21  
Chat\_Bot\_05-06-2024\_11:07:20  
Chat\_Bot\_18-04-2024\_09:46:07  
Chat\_Bot\_11-04-2024\_09:53:09

**SAVE**



# AGENT AI ASSISTANT

IVE AI





TCHAC - john.doe@email.com

Currently in use by



#30 - TCHAC

Open

00:00:52

john.doe@email.com - 24/06/24 09:58:09



New chat request  
from: <http://localhost:7002>  
name: John Doe  
email: john.doe@email.com

john.doe@email.com - 24/06/24 09:58:34



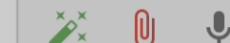
Can you tell me something about the error code 3



Get help from XCALLY Assistant

ode 3ZE0VZ1K?

Insert your text message here



SEND ANSWER



The error code 3ZE0VZ1K indicates that the appliance has detected that the heating element is calcified. To resolve this issue, you should descale the appliance and operate it with the water softening system. [\[4:0†source\]](#) .

CONTACT

INTERACTION

CUSTOMER >

Name

john.doe@email.com

Email

john.doe@email.com

Phone

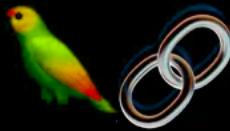
Mobile

Fax

MORE

# What more?

chain = prompt | llm | output\_parser



# LangChain



# LlamaIndex

GPT4ALL

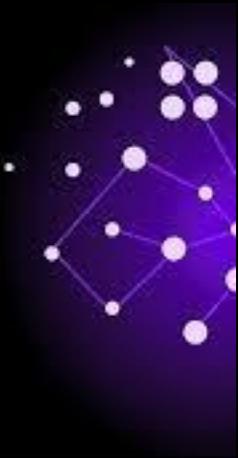
A free-to-use, locally running, privacy-aware chatbot. **No GPU or internet required.**

graph databases VS vector databases

KNOWLEDGE GRAPH

An innovative approach to knowledge retrieval

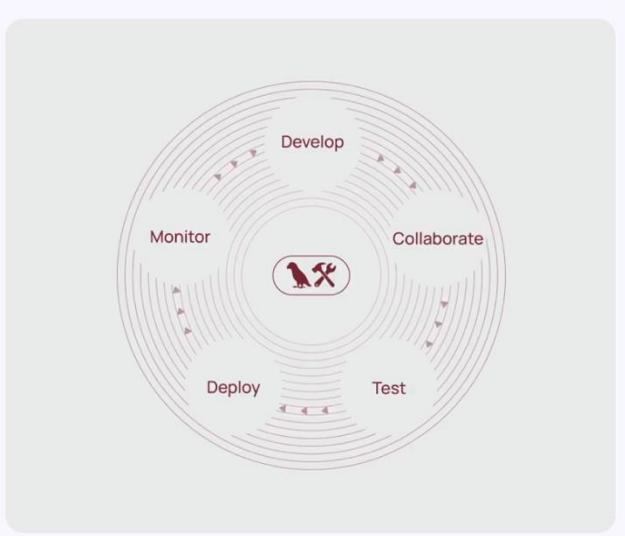
WRITER



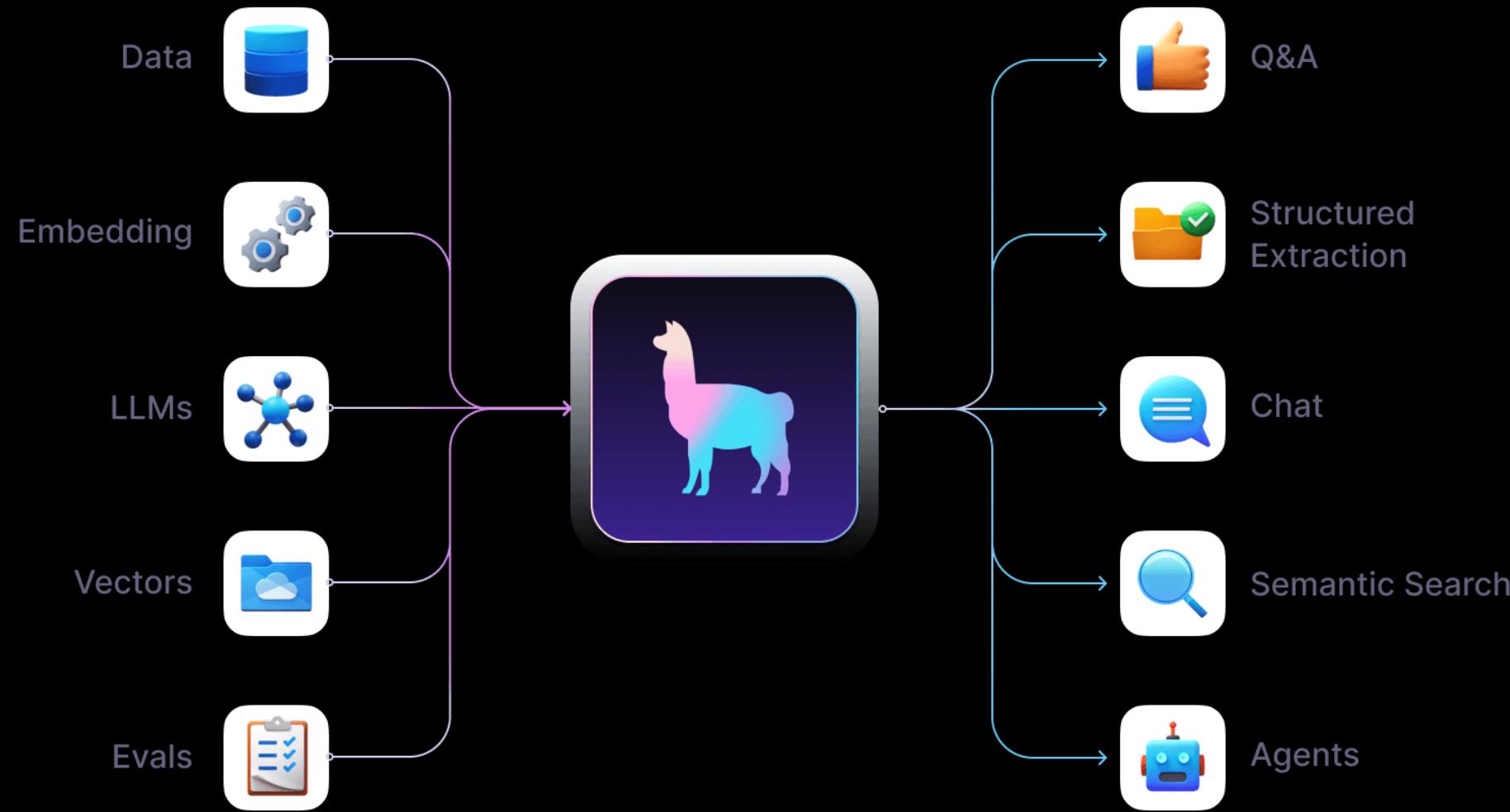
Observe performance with LangSmith

Ship faster with LangSmith's debug, test, deploy, and monitoring workflows. Don't rely on "vibes" – add engineering rigor to your LLM-development workflow, whether you're building with LangChain or not.

[Learn more about LangSmith ↗](#)



# LLmalndex: data framework for building LLM applications





Personal > Projects > myproject

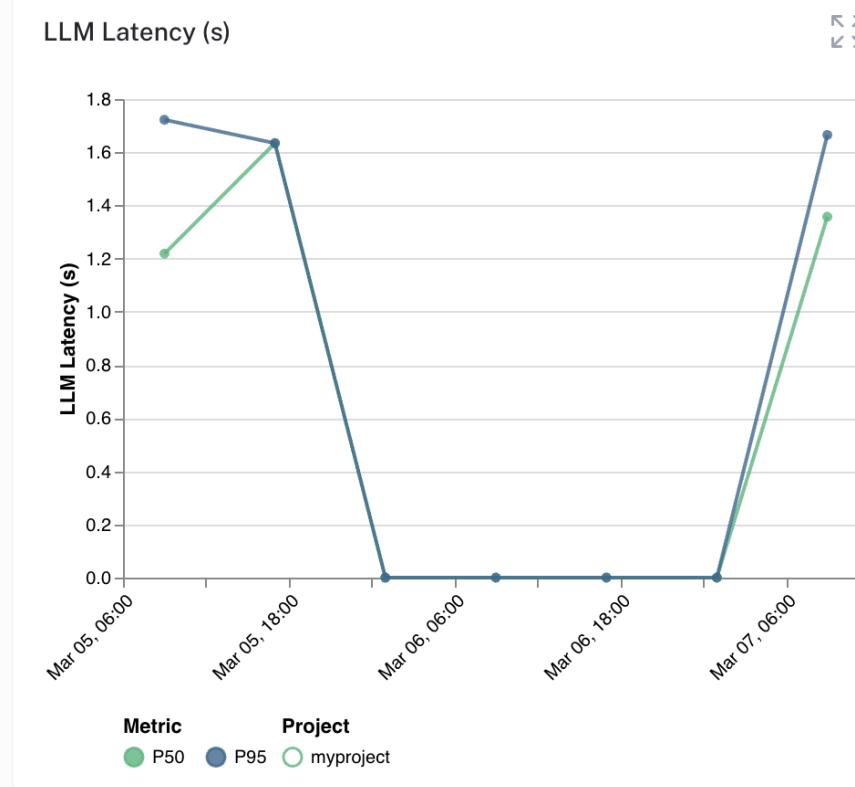
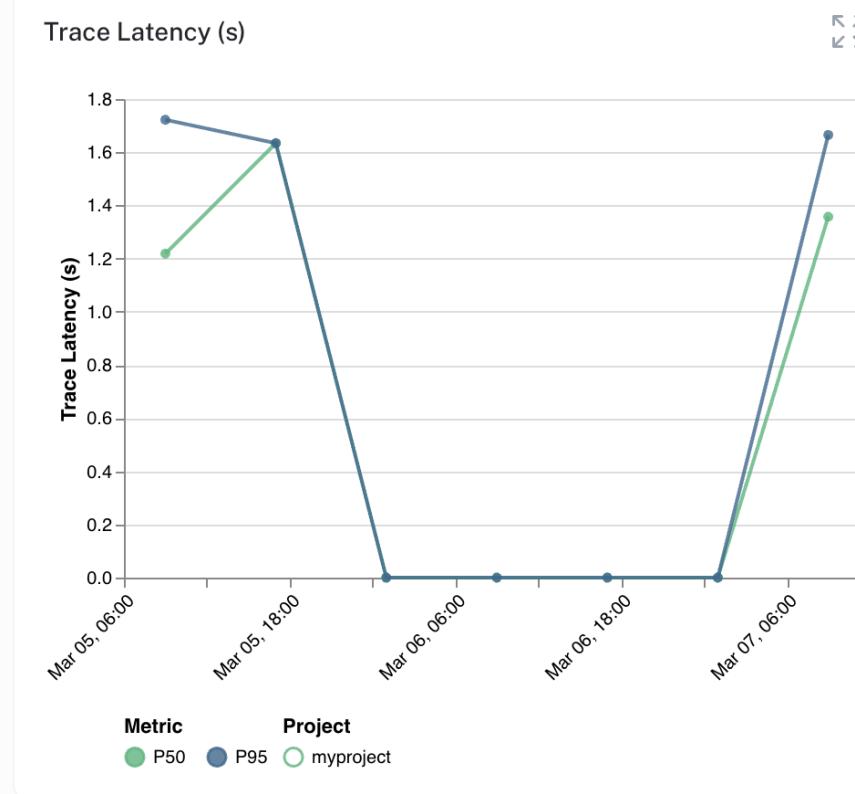
--> One API Key and project x Manager Customer



## myproject

Edit

Traces LLM Calls All Runs Monitor Setup



### Stats

Last 7 days

RUN COUNT

17

TOTAL TOKENS

1,534 / \$0.001203 ⓘ

MEDIAN TOKENS

102

ERROR RATE

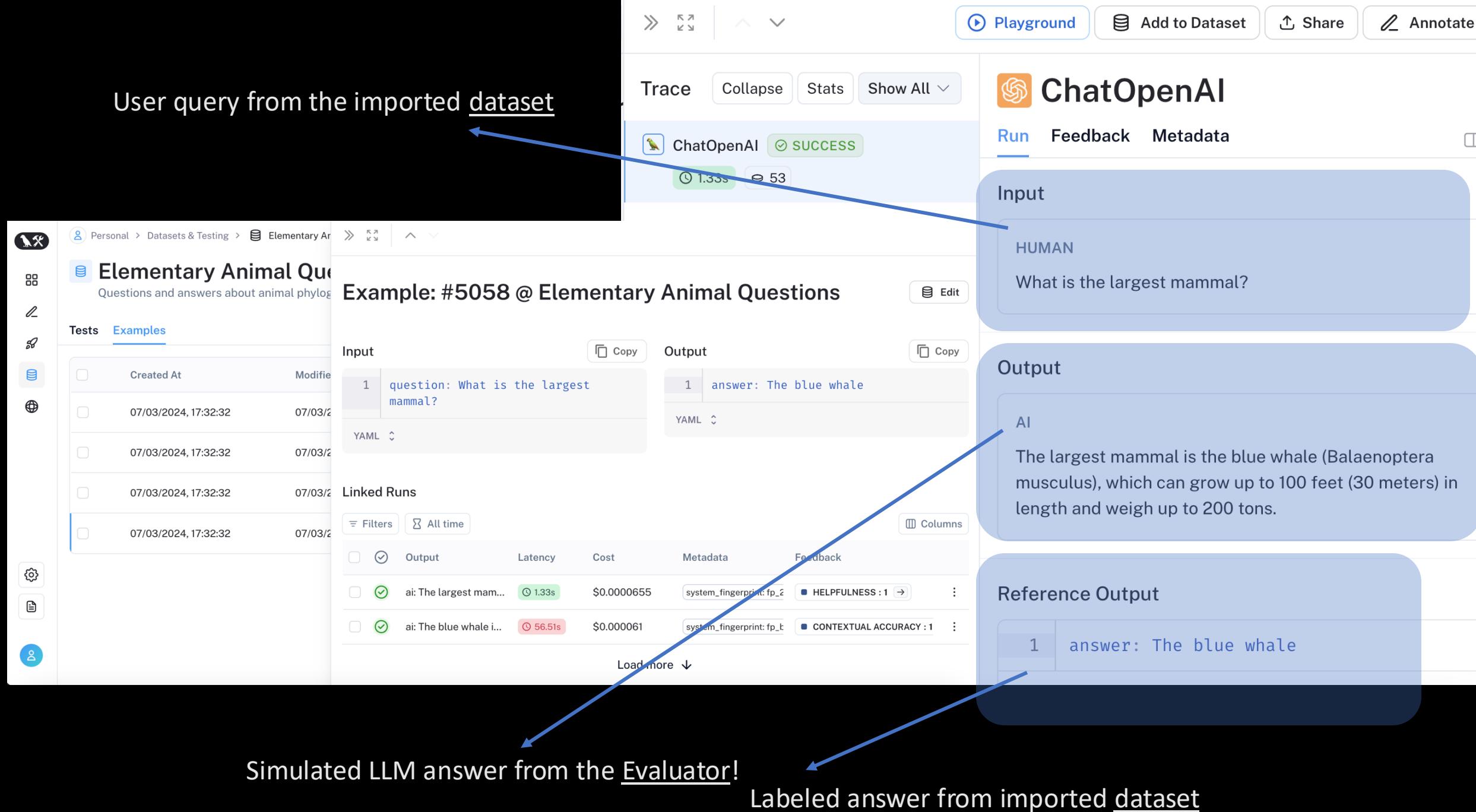
0%

% STREAMING

0%

LATENCY

P50: 1.38s P99: 1.78s



The screenshot shows a user interface for a machine learning model. At the top, there's a navigation bar with icons for file operations like back, forward, and search. Below it, the title "Elementary Animal Questions" is displayed, followed by a subtitle "Questions and answers about animal phylogeny". A sub-header "Example: #5058 @ Elementary Animal Questions" is centered above the main content area.

The main area is divided into two sections: "Input" and "Output".

- Input:** Shows a single row of data:
  - Index: 1
  - Text: "question: What is the largest mammal?"
- Output:** Shows a single row of data:
  - Index: 1
  - Text: "answer: The blue whale"

Below these sections, there's a "YAML" button. Underneath the "Output" section, there's another "YAML" button.

On the left side of the main content area, there's a sidebar titled "Tests Examples". It lists several entries with columns for "Created At", "Modified", and "Metadata". One entry is highlighted with a blue border, showing "07/03/2024, 17:32:32" in both columns. The "Metadata" column for this entry shows "ai: The largest mammal..." with a latency of 1.33s and a cost of \$0.0000655. The "Feedback" column for this entry shows "HELPFULNESS : 1" and "CONTEXTUAL ACCURACY : 1".

At the bottom of the main content area, there are buttons for "Filters" and "All time", and a "Load more" button.

YAML feedback:

key: COT Contextual Accuracy

score: 1

value: CORRECT

comment: |-

The answer correctly identifies the blue whale as the largest mammal, which is the correct answer according to the context. The additional information about the size and weight of the blue whale is also correct and does not conflict with the context or the question.

GRADE: CORRECT

evaluator\_info:

\_run:

run\_id: 91f6b719-2ea2-436c-a012-8a29f7a091ca

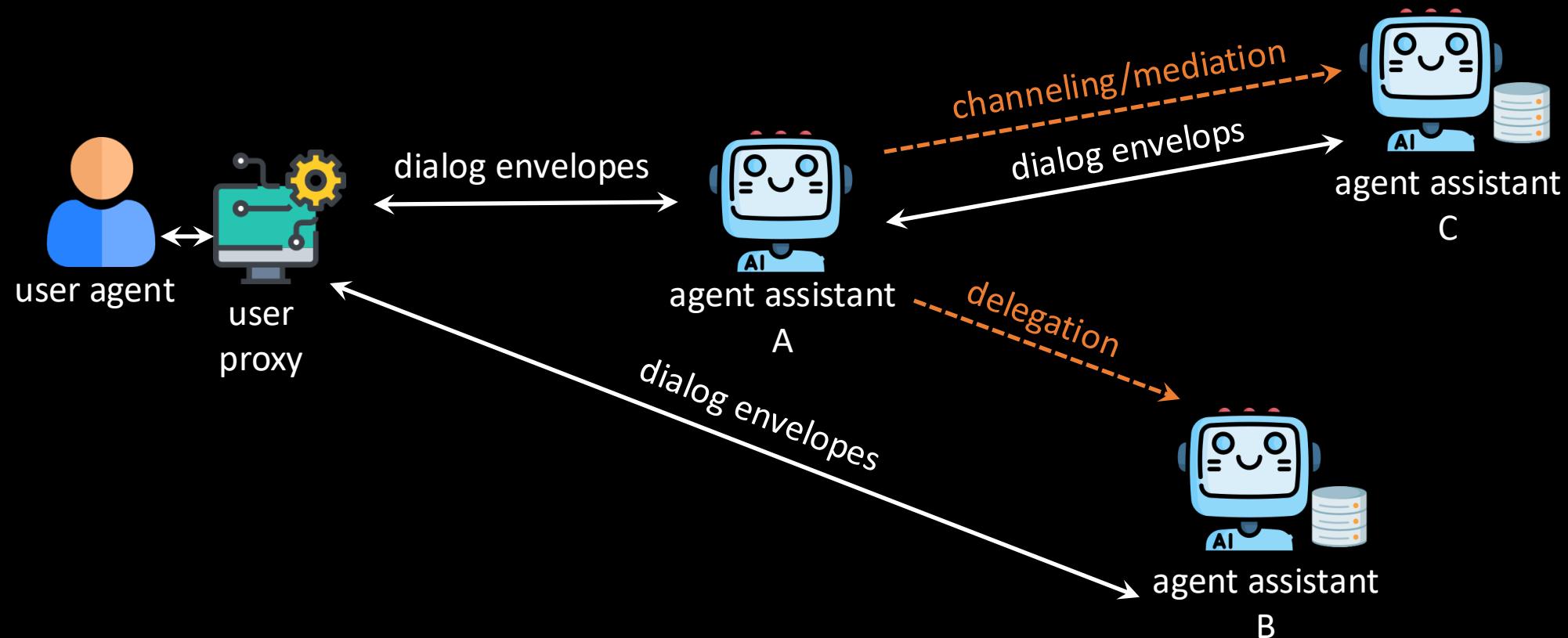


OPEN VOICE  
**INTEROPERABILITY**



# Dialog Agent Message Envelopes

Universal JSON structure



Agent A: 1<sup>st</sup> level dialog agent assistant

Agent B/C: usually with specific knowledge, and resource connections (i.e. CRM, specific KB etc...)



The Open Voice Network

Open Voice Interoperability Initiative - LF AI & Data Foundation  
Architecture Work Group

3 July 2024

Draft Version 0.9.2 Status: Published

*Editor-in-Chief: David Attwater*

*Contributors: Emmett Coin, Deborah Dahl, Jim Larson, Allan  
Wylie, Rainer Turner and Diego Gosmar*

<https://github.com/open-voice-interoperability/docs>

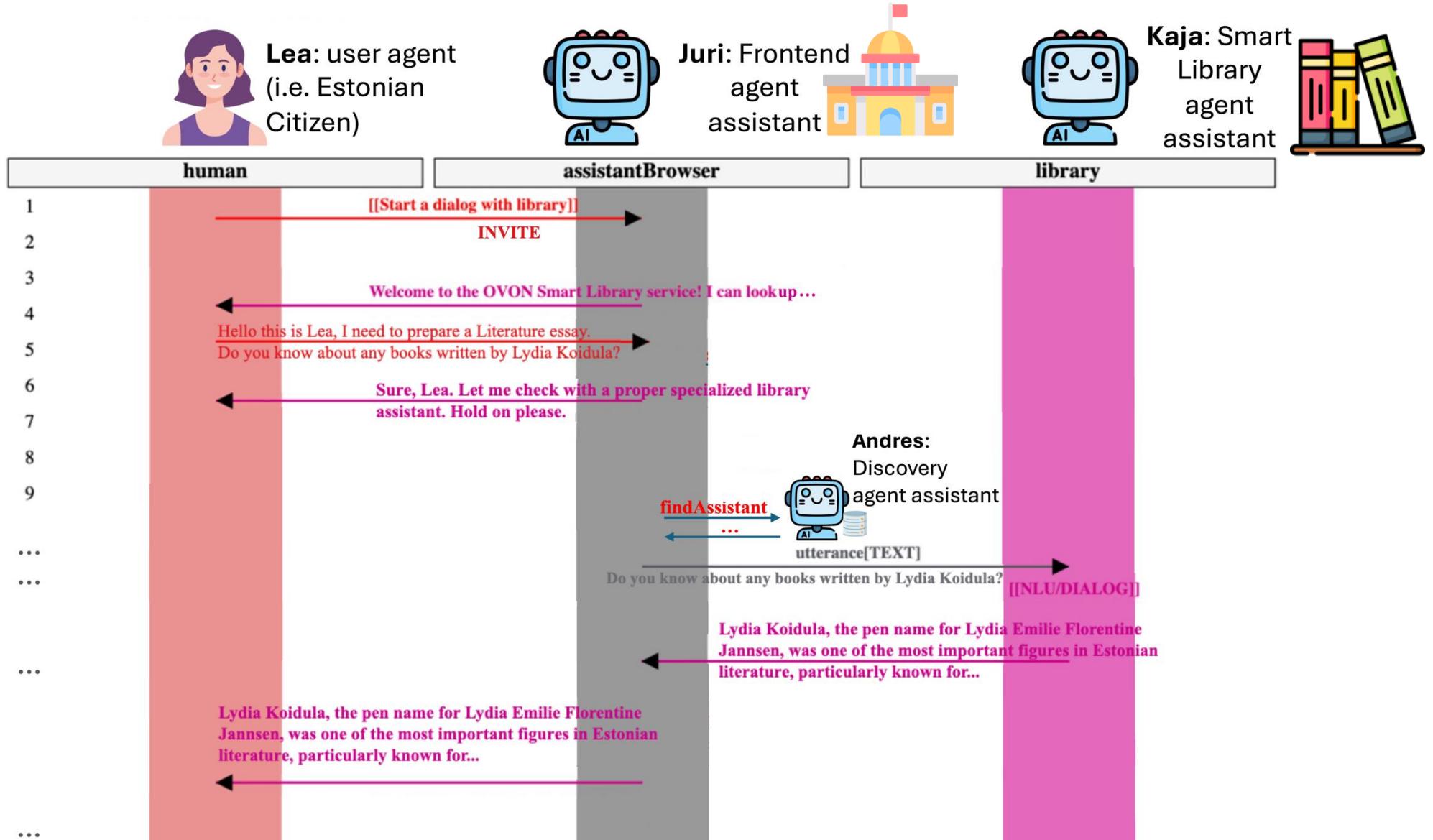


**Conversational AI Multi-Agent Interoperability,  
Universal Open APIs for Agentic Natural Language  
Multimodal Communications**

Diego Gosmar, Deborah A. Dahl, Emmett Coin

<https://arxiv.org/abs/2407.19438>

# OPEN VOICE INTEROPERABILITY



# Example OVON REQUEST UTTERANCE & WHISPER

To be sent as part of the body of an  
HTTPS POST request to the endpoint:  
<https://ovon.xally.com/smartlibrary>

Replace values with your desired  
questions about books:

**utterance**: main request

**whisper**: additional infos you'd like  
to send/request



Formatted according to the:  
*Conversational Envelope and Dialog Event specs*

```
"ovon": {
  "schema": {
    "version": "0.9.0",
    ...
  }
  "eventType": "invite",
  "parameters": {
    "to": {
      "url": "https://ovon.xally.com"
    }
  }
  ...
  {
    "eventType": "utterance",
    "parameters": {
      ...
      "mimeType": "text/plain",
      "tokens": [ { "value": "Do you know about any books written by Lydia Koidula?" } ]
    }
  }
  ...
  {
    "eventType": "whisper",
    "parameters": {
      ...
      "features": {
        "text": {
          "mimeType": "text/plain",
          "tokens": [ { "value": "Lydia Koidula, the pen name for Lydia Emilie Florentine Jannsen, was one of the most important figures in Estonian literature, particularly known for..." } ]
        }
      }
    }
  }
}
```

# “It’s a Wrap-up”

---

Thank you!

<https://www.linkedin.com/in/diegogosmar>