# Hands-On Retrieval Augmented Generation (RAG)

Introduction to LanceDB and LangChain

# About Me

- MSc Student, Information Technology Engineering, University of Tehran

- BSc Computer Engineering, University of Isfahan

- Graduate Researcher at
  - NLP Lab
    - Supervisor: Dr. Hesham Faili
  - Computational Cognitive Psychology Lab
    - Supervisor: Dr. Hadi Moradi

- LLM Researcher

- Research Field: LLM base Language Agents and Psychology
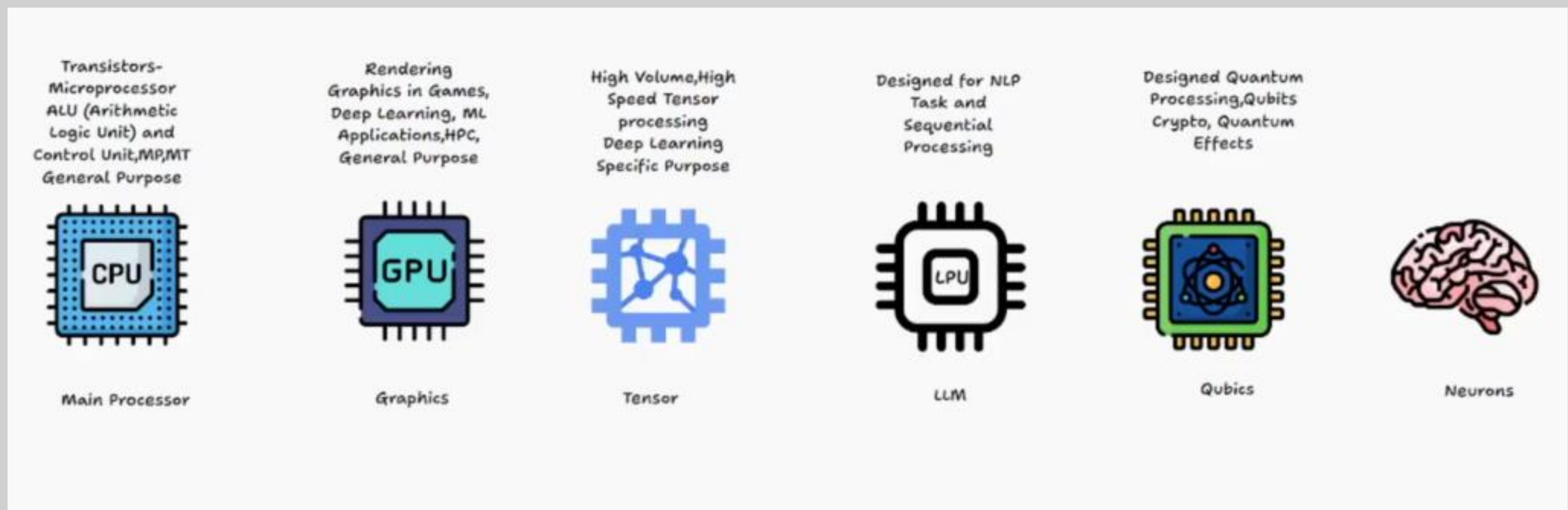
- AI/Software Engineer

# Workshop Plan



- Introduction to RAG

- How to work with LanceDB
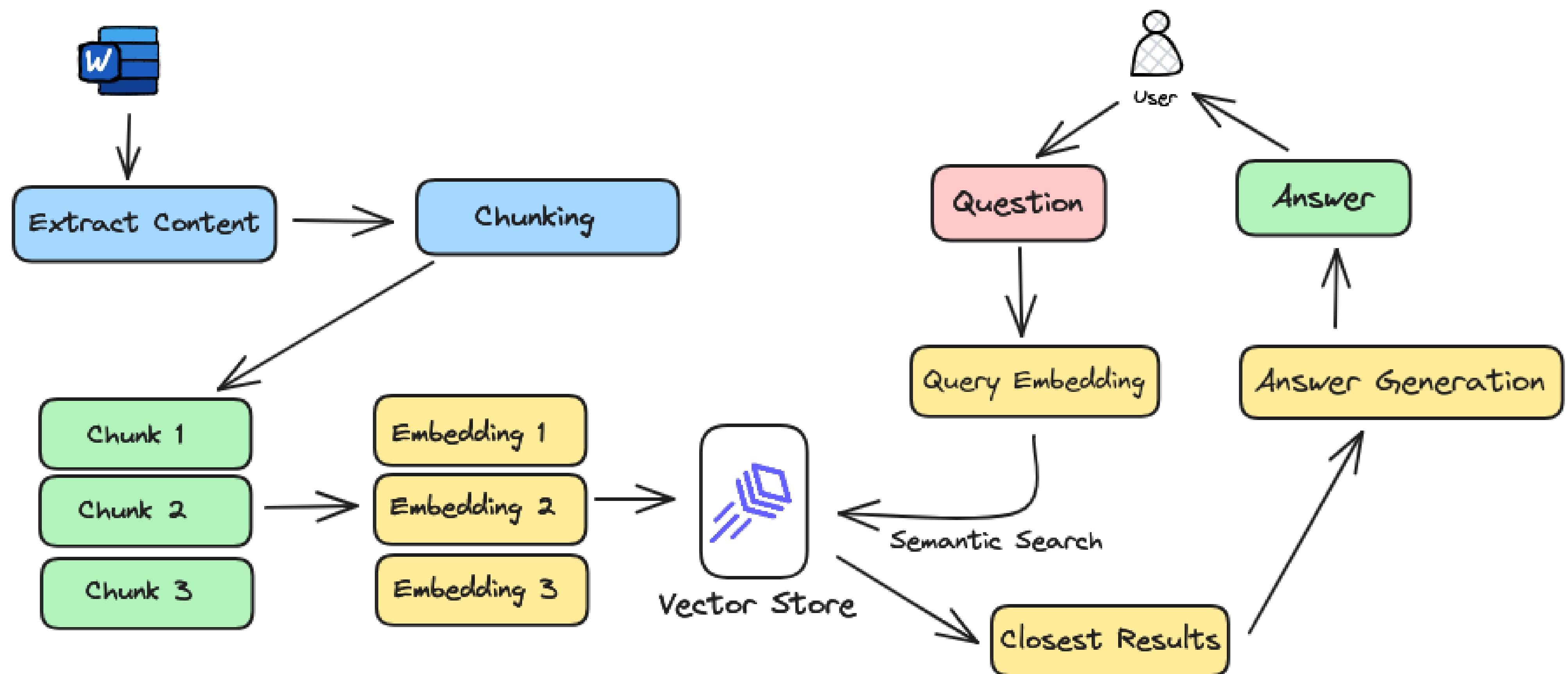
- How to work with LLMs using APIs and Langchain framwork
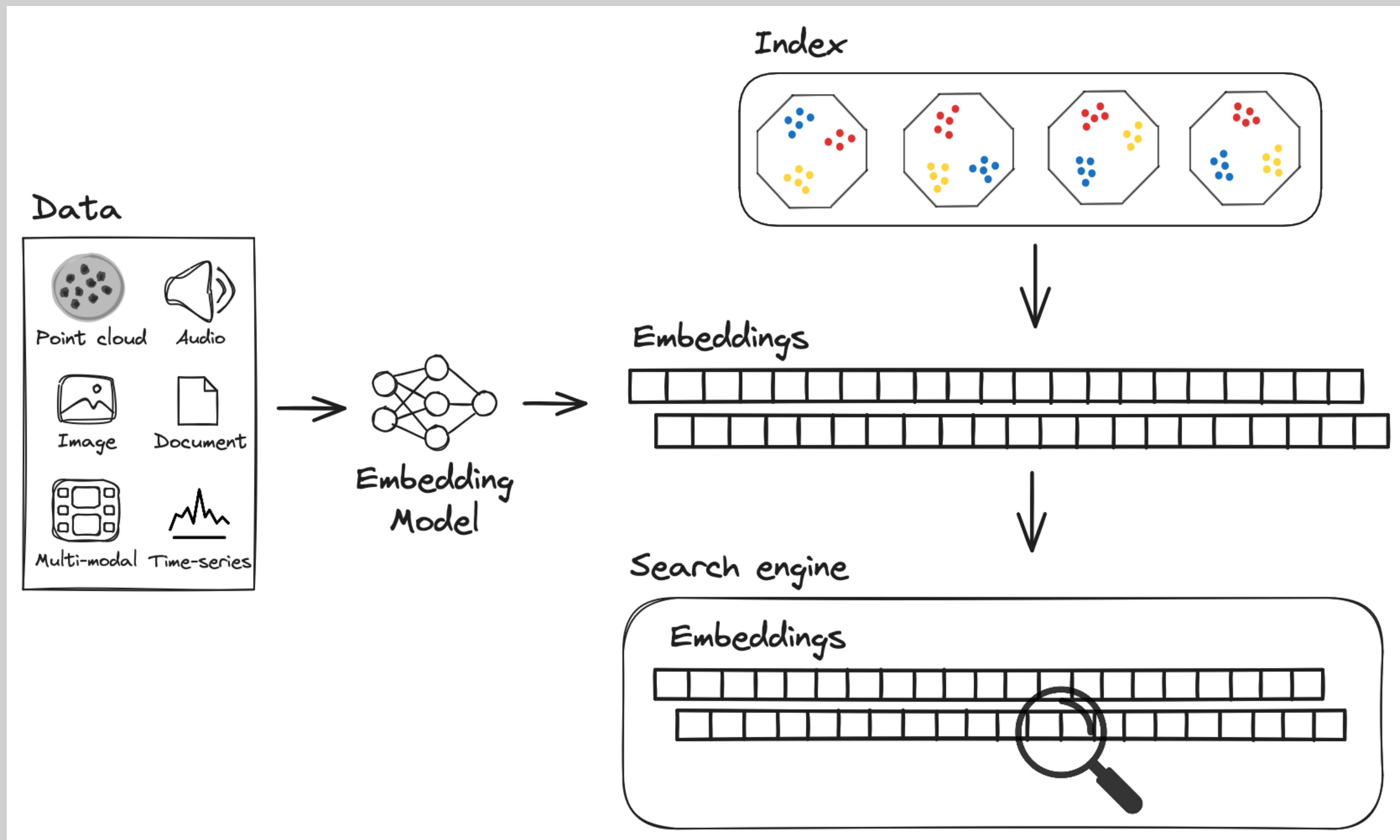
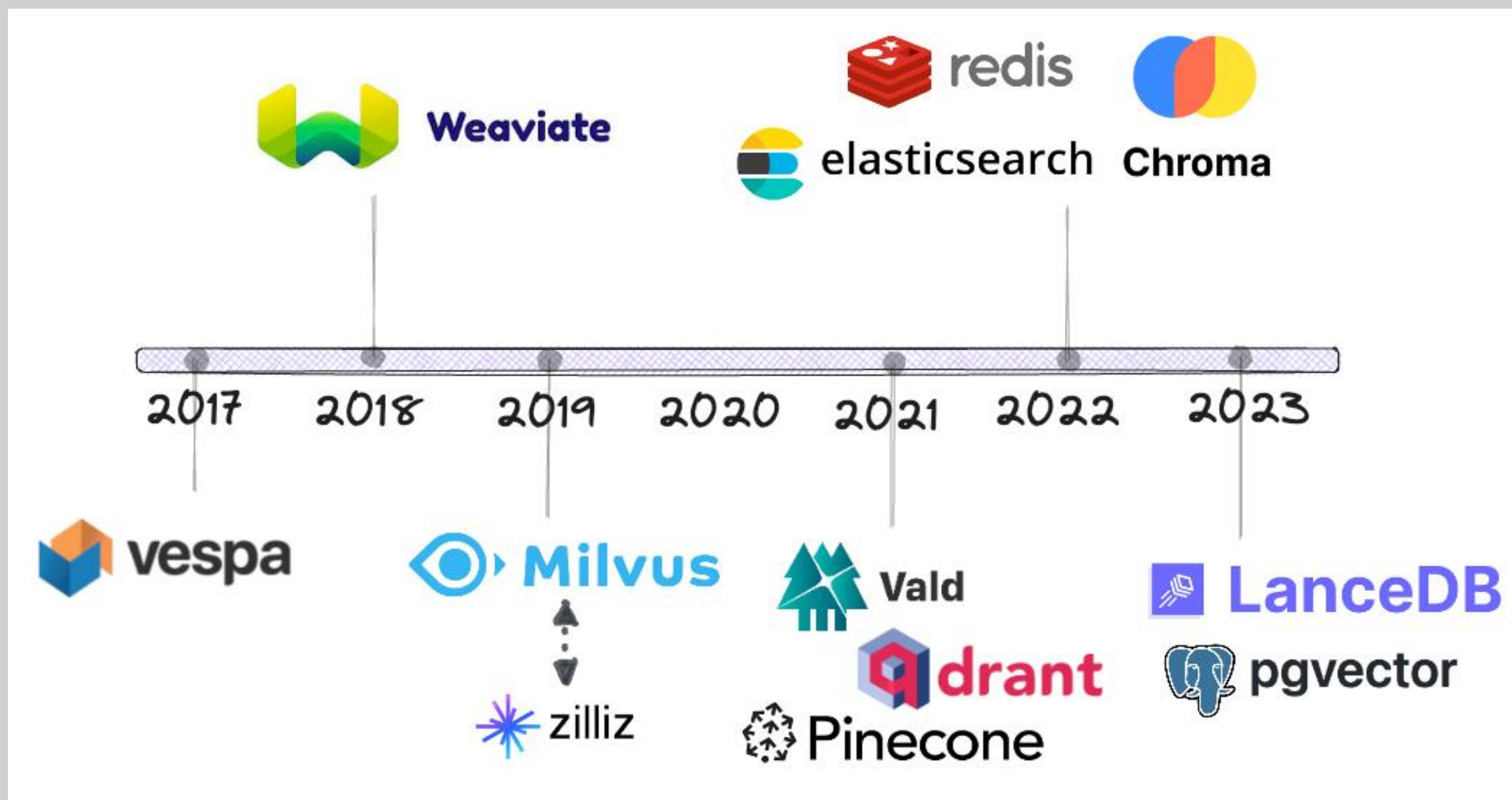# LLMs, but not just input/output



Source: JCharisTech

Milad Mohammadi | Hands-On RAG | Intelligent Information Retrieval: Workshop 2 | Fall 2024

# RAG



Source: LanceDB github

# Vector Search



Source: LanceDB github

Milad Mohammadi | Hands-On RAG | Intelligent Information Retrieval: Workshop 2 | Fall 2024

# Vector Databased



Source: The Data Quarry blog

Milad Mohammadi | Hands-On RAG | Intelligent Information Retrieval: Workshop 2 | Fall 2024

# LanceDB



Source: LanceDB github

Milad Mohammadi | Hands-On RAG | Intelligent Information Retrieval: Workshop 2 | Fall 2024

# Let's code with LanceDB

Milad Mohammadi | Hands-On RAG | Intelligent Information Retrieval: Workshop 2 | Fall 2024

# LanceDB Documentation



Source: LanceDB github

# LLM based systems and products



**AI Products** re-bundle user workflow around AI systems (e.g. feedback to AI)

**AI Systems** combine multiple AI components (including LLMs) to do things like access private data, automate workflows, ...

**LLMs** are the core ("intel inside")

Vertically integrated

End-to-end feedback loop

AI Product

AI System

LLM

GPT-4, Claude, Customs LLMs

Source: BoxCars AI: Tabrez Syed

# Model Providers

Milad Mohammadi | Hands-On RAG | Intelligent Information Retrieval: Workshop 2 | Fall 2024

# TogetherAI



Source: TogetherAI

# What does 'client' mean here?



Source: TowardsDataScience: Shaw Talebi

# Stack

# LangChain

# Let's code with OpenAI client and Langchain

Milad Mohammadi | Hands-On RAG | Intelligent Information Retrieval: Workshop 2 | Fall 2024

# Advanced RAG techniques



Source:
Latest and
Greatest
blog

div.beehiiv
.com

# Agentic RAG



Source: X.com/
@helloiamleonie

Milad Mohammadi | Hands-On RAG | Intelligent Information Retrieval: Workshop 2 | Fall 2024

# Good Luck with LLMs!

**you'll need it ;)**

For further discussions or questions, feel free to reach out via email:

miladmohammadi@ut.ac.ir