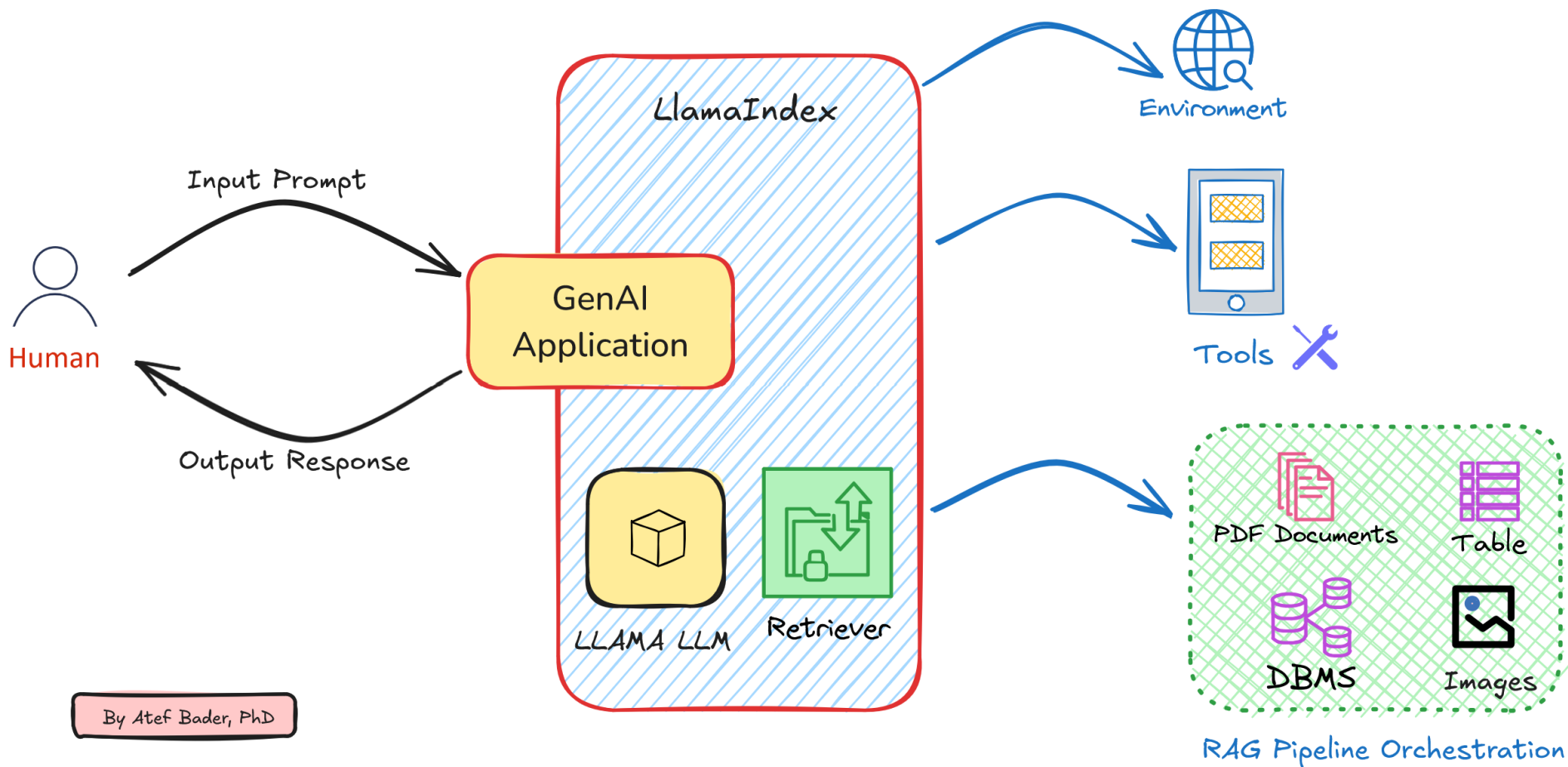# GenAI Applications:

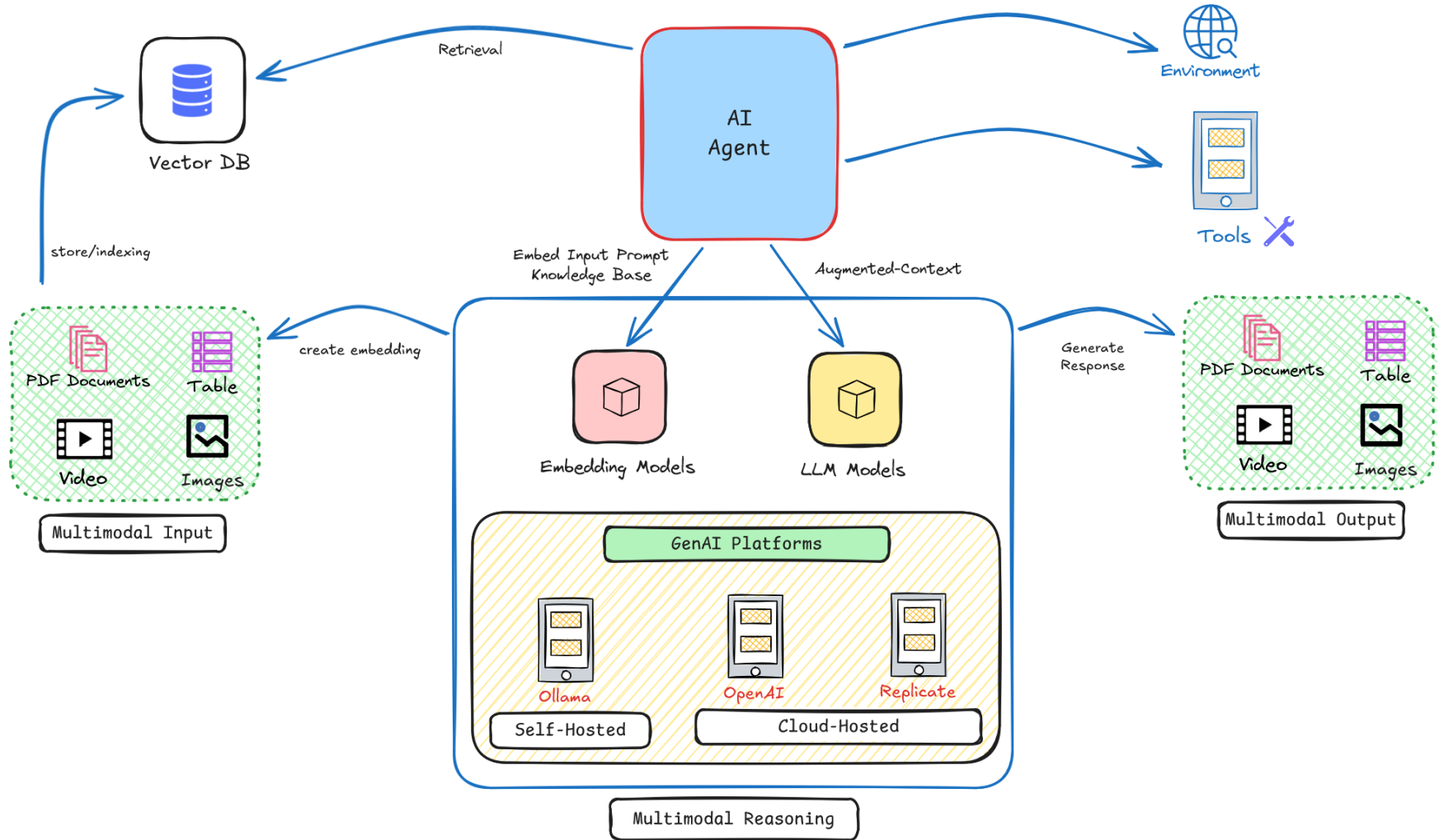**RAG Pipeline Orchestrations, Agents, Workflows, Chats, etc.**

**Dr. Atef Bader**

Human

Input Prompt

Output Response

GenAI Application

LlamaIndex

LLAMA LLM

Retriever

Environment

Tools 🛠️

PDF Documents

Table

DBMS

Images

RAG Pipeline Orchestration

By Atef Bader, PhD

2

By Atef Bader, PhD

GenAI Application

Environment

Tools

create embedding

Embedding Model

Retrieval

Augement-Context

Multimodel Input

PDF Documents

Table

Video

Images

store/indexing

Vector DB

LLM Model

Generate Response

Multimodel Output

PDF Documents

Table

Video

Images

# RAG Pipeline Orchestration

Retrieval

AI
Agent

Environment

Tools

Vector DB

store/indexing

Embed Input Prompt
Knowledge Base

Augmented-Context

create embedding

PDF Documents

Table

Video

Images

Multimodal Input

Embedding Models

LLM Models

GenAI Platforms

Ollama

OpenAI

Replicate

Self-Hosted

Cloud-Hosted

Multimodal Reasoning

Generate
Response

PDF Documents

Table

Video

Images

Multimodal Output

# Multimodal AI Agent Engineering

# Human

**Persona, Social, Ergonomics, Dialog**

# Interaction

- HCI Laws
- UI/UX
- User Research
- Design Process
- Design Principles

- Libraries & Frameworks
- Evaluation & Usability
- UI Patterns
- Prompt Engineering & NLP
- Generative AI

**Human-Computer Interaction**

# Computer

**Architecture, I/O Devices, OS Platform, Graphics**

# Human-LLM Application Engineering

## Human-Computer Interaction

**Human**

Persona, Social, Ergonomics, Dialog

**Interaction**

- Program a Computer
- Algorithm in Pseudo code
- HCI Laws          - Libraries & Frameworks
- UI/UX             - Evaluation & Usability
- User Research     - UI Patterns

**Computer**

Architecture, I/O Devices, OS Platform, Graphics

## Human-LLM Interaction

**Human**

Persona, Social, Ergonomics, Dialog

**Interaction**

- Program an LLM
- Prompt Engineering & Chain-of-Thought (CoT) in NL
- Generative AI for Text, Code, Digital Content
- Retrieval-Augmented Generation (RAG)
- Embeddings & Vector Database
- Conversational Agents
- Causal Inference, Baysian/Belief Networks

**LLM**

GPU,  AI Platforms, Models

**+**

# Human Brain

-
  - The human brain is made up of about 86 billion nerve cells, along with many other types of cells.



The human brain is made up of about 86 billion nerve cells, along with many other types of cells. They interact and link together in unique ways, creating distinct brain regions with specific functions. Uncovering the complex makeup and interactions of these many cells could lead to a new understanding of how the brain functions in health and disease, and new tools to study the complex activities and functions of these cells.

To better understand the identities and roles of brain cells, NIH's Brain Research Through Advancing Innovative Neurotechnologies® (BRAIN) Initiative launched an international network of collaborating researchers called the BRAIN Initiative Cell Census Network. Its aim is to create a comprehensive inventory of all the cells in the human, nonhuman primate, and mouse brains, including cell locations, interconnections,

An international network of researchers created the most complete cell atlases yet of the human brain.
*vitstudio / Shutterstock*

and activities. The study of brain cells across species can pinpoint features that are uniquely human and give insights into which animals to study for different scientific questions. The latest findings were reported in a series of more than 20 papers published in *Science, Science Advances,* and *Science Translational Medicine* on October 13, 2023.

One paper examined three human brains to find over 3,000 types of brain cells—more than previously known. The team identified specific types of cells in distinct clusters in different brain regions. These findings could help shed light on conditions that are known to affect specific brain areas, such as cancer or neurodegenerative diseases.

The researchers have created the most detailed cell atlas yet of the adult human brain. The atlas reveals information about each cell's gene activity and epigenome—the changes to a cell's DNA and chromosomes that alter genetic activity. The findings also show that, besides variation among brain regions, there is variation between individuals. More people will need to be studied to fully understand the patterns of healthy and diseased brains.
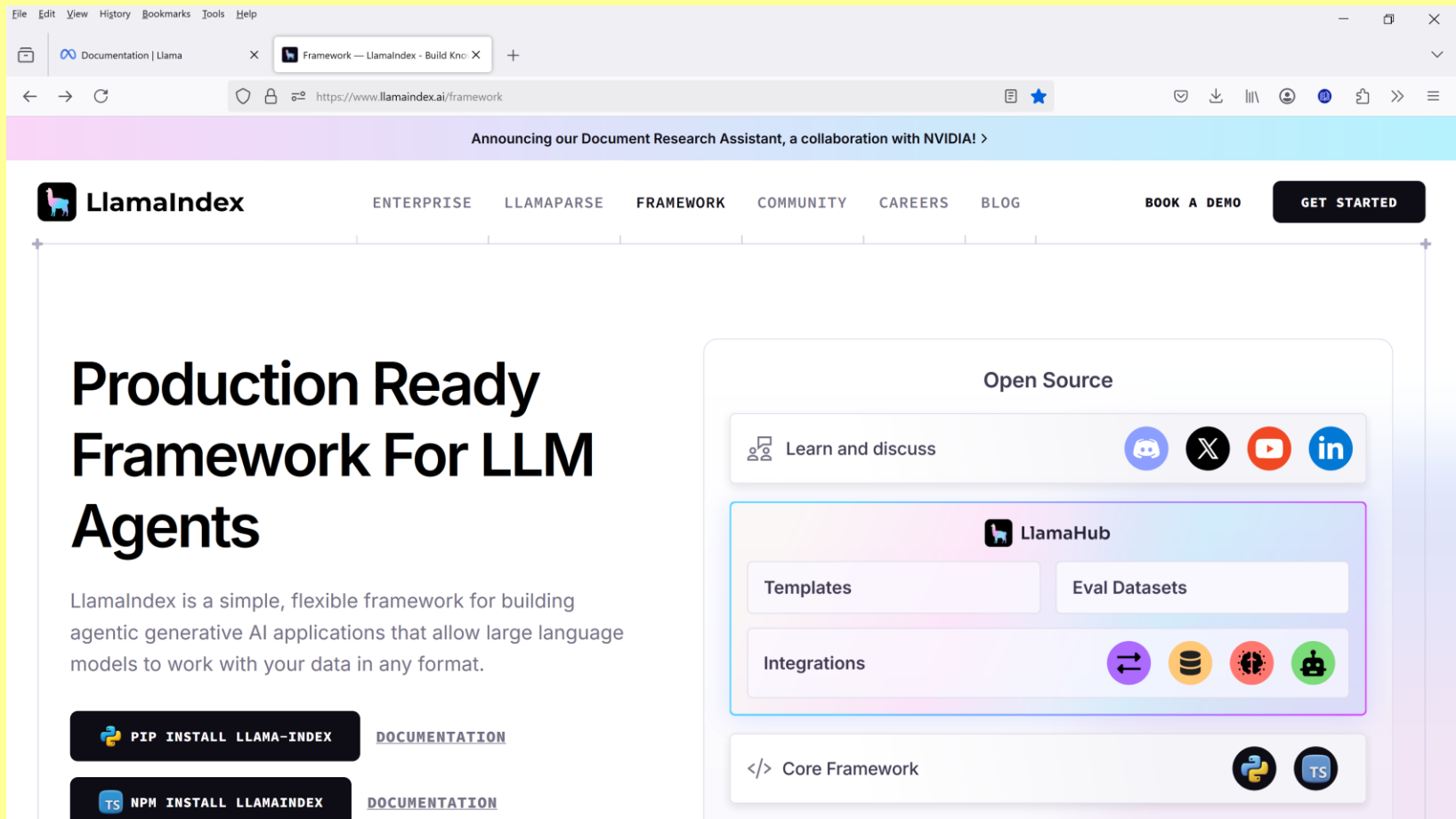
Another paper compared the cellular and molecular properties of the brains of humans and several nonhuman primates: the

# LlamaIndex : Framework for GenAI Applications

- https://www.llamaindex.ai/framework

# LLAMA – LLM Model

- [https://www.llama.com/docs/overview/](https://www.llama.com/docs/overview/)

# HuggingFace : Repo for Open-Source Models

- https://huggingface.co/meta-llama

# Replicate : Platform for Cloud-Hosted Open-Source/Public Models

- [https://replicate.com/explore](https://replicate.com/explore)

# OpenAI : Platform for Cloud-Hosted Proprietary Models

- [https://openai.com/api/](https://openai.com/api/)

# Ollama : Container for Self-Hosted Open-Source Models

- https://ollama.com/

# Why LlamaIndex?

- LlamaIndex is the framework for Context-Augmented LLM Applications

- LlamaIndex provides tools like:
  - **Data connectors** ingest your existing data from their native source and format. These could be APIs, PDFs, SQL, and (much) more.
  - **Data indexes** structure your data in intermediate representations that are easy and performant for LLMs to consume.
  - **Engines** provide natural language access to your data. For example:
    - **Query engines** are powerful interfaces for question-answering (e.g. a RAG flow).
    - **Chat engines** are conversational interfaces for multi-message, "back and forth" interactions with your data.
  - **Agents** are LLM-powered knowledge workers augmented by tools, from simple helper functions to API integrations and more.
  - **Observability/Evaluation integrations** that enable you to rigorously experiment, evaluate, and monitor your app in a virtuous cycle.
  - **Workflows** allow you to combine all of the above into an event-driven system far more flexible than other, graph-based approaches.

# Why LlamaIndex?

- You can use LlamaIndex for Context-Augmented LLM Applications, when using:

  - ## OpenAI

  - ## Replicate

  - ## Ollama

- Here is a tutorial:

  - https://docs.llamaindex.ai/en/stable/getting_started/starter_example_local/

# Benchmark & Testbed

| Platform | Model | | Framework |
|---|---|---|---|
| Ollama | llama3.2:1b   [1.3GB]<br>llama3.2:3b   [2.0 GB]<br>llama3.2:3b-instruct-fp16   [6.4GB] | HuggingFaceEmbedding<br>BAAI/bge-small-en-v1.5 | LlamaIndex<br>LangChain/LangGraph |
| Replicate | meta/meta-llama-3-70b-instruct<br>meta-llama-3.1-405b-instruct<br>meta/meta-llama-3-8b-instruct<br>meta-llama-3.1-405b-instruct<br>meta/meta-llama-3-8b-instruct<br><br>deepseek-ai/deepseek-r1 | HuggingFaceEmbedding<br>BAAI/bge-small-en-v1.6 | LlamaIndex<br>LangChain/LangGraph |
| OpenAI | gpt-4o-mini<br>gpt-3.5-turbo | text-embedding-3-small<br>text-embedding-ada-002 | LlamaIndex<br>LangChain/LangGraph |

# References

- https://ollama.com/
- https://replicate.com/google/imagen-3/examples?input=python
- https://replicate.com/google/imagen-3-fast
- https://deepmind.google/technologies/imagen-3/
- https://huggingface.co/meta-llama
- https://www.llamaindex.ai/
- https://www.llama.com/docs/overview/