



HOW TO BUILD A CHATBOT

Session 5 -
Building a Chatbot

SESSION 5

AGENDA



1

Demo of Target Solution

2

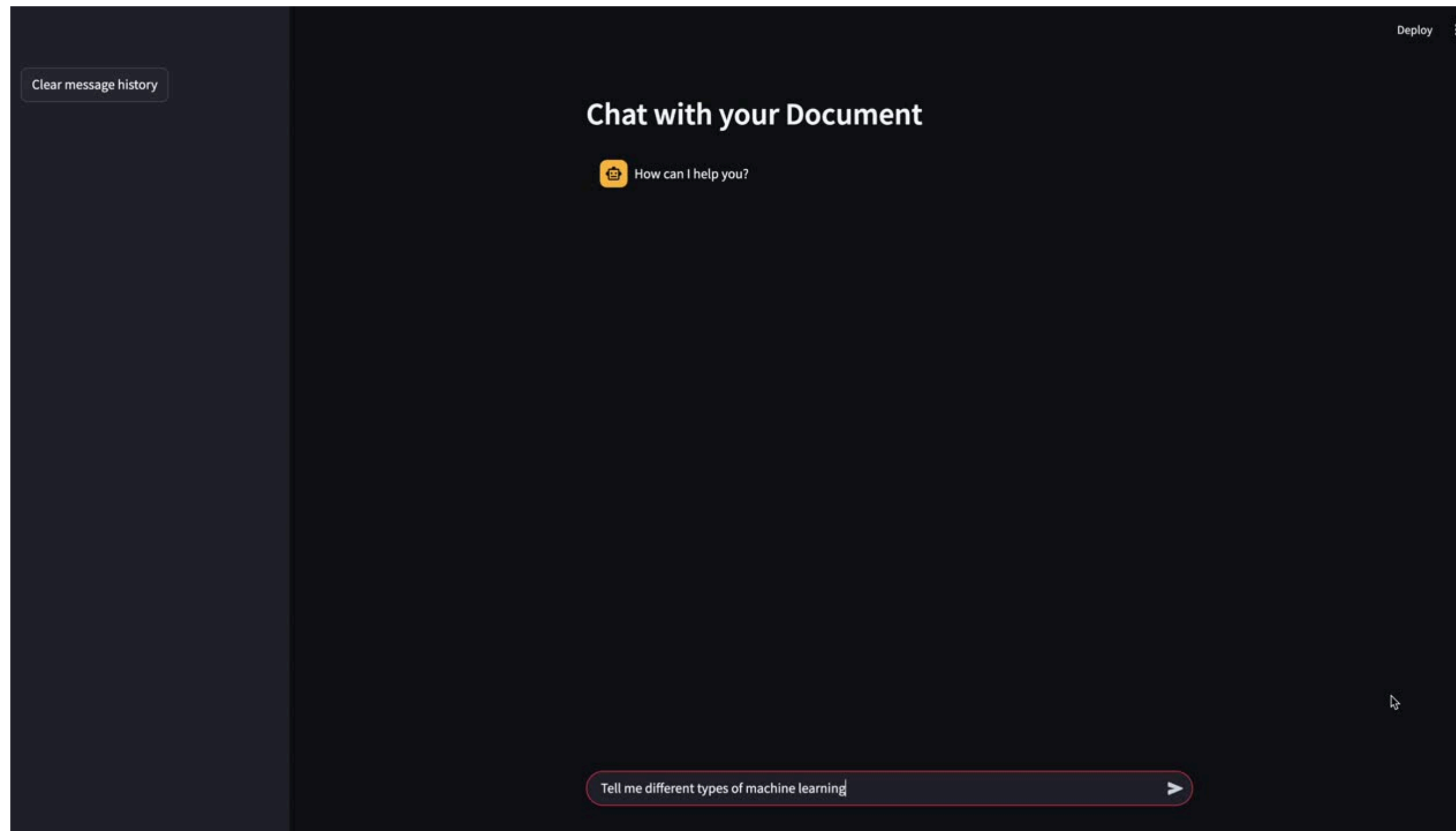
Target Architecture

3

Building Blocks

DEMO OF TARGET SOLUTION

Chatbot App in Action



TARGET ARCHITECTURE

Chatbot App:

- Web app built with Streamlit, accessible via browser.
- Python-based with FastAPI and LangChain.

LLM Serving:

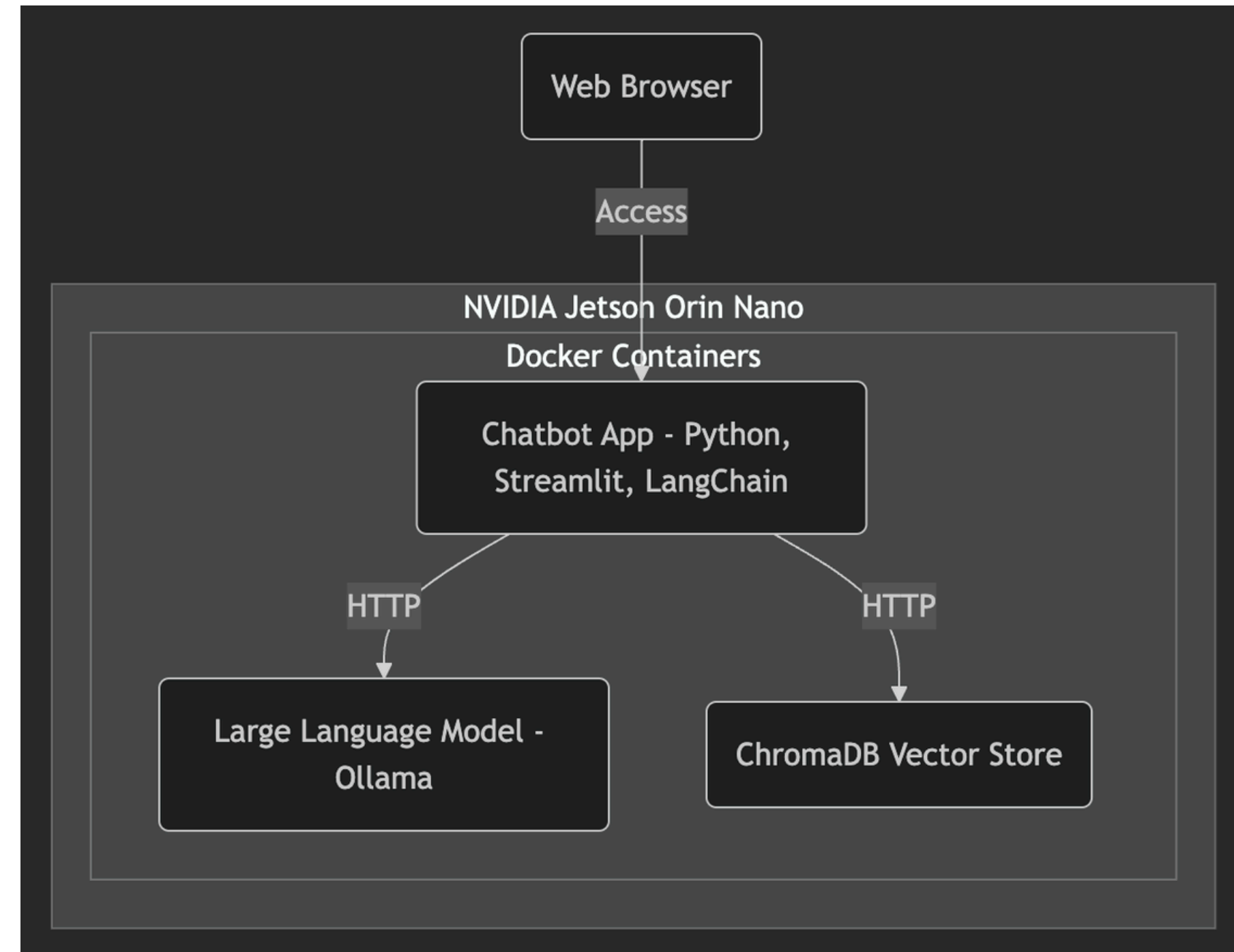
- Ollama for managing large language models.

Knowledge Storage:

- Vector database for knowledge management.

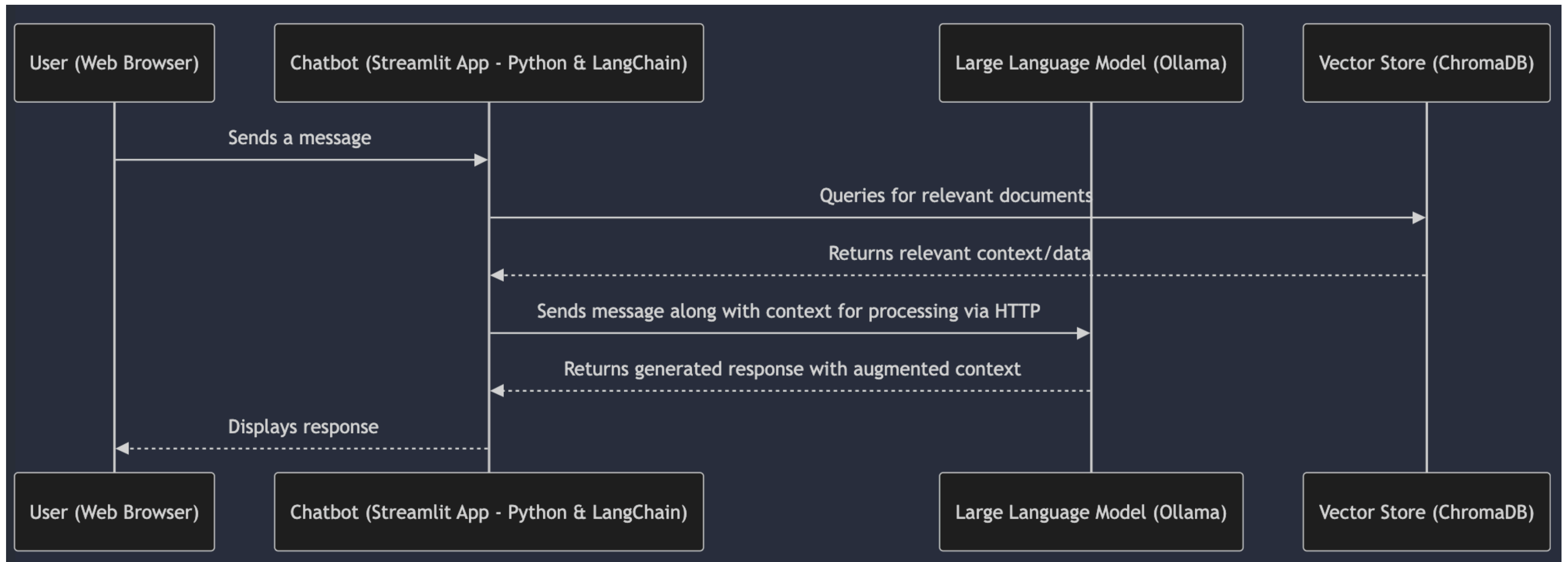
Deployment:

- Docker containers for application deployment.



TARGET ARCHITECTURE

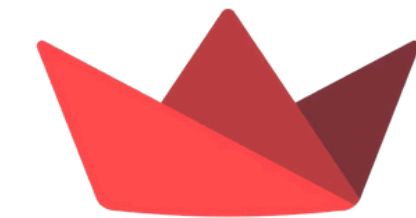
User Interaction Workflow



BUILDING BLOCKS

Frontend - Streamlit Webapp

- Streamlit is an open-source Python framework for building interactive web apps.
- It offers built-in widgets for easy UI development.
- Apps update in real-time without manual refreshes.
- Streamlit apps can be deployed on cloud platforms easily.

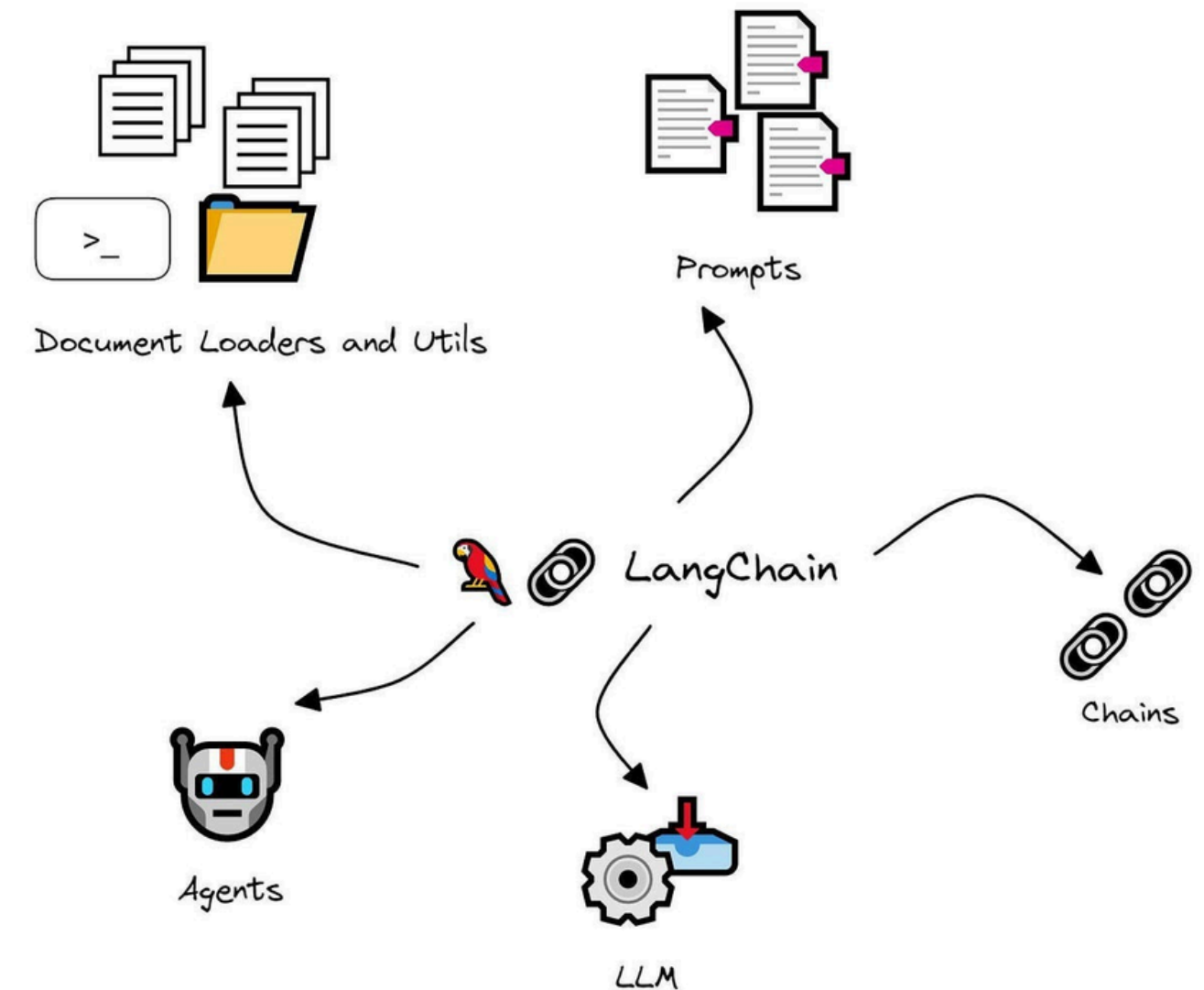


Streamlit

BUILDING BLOCKS

RAG Chatbot with LangChain

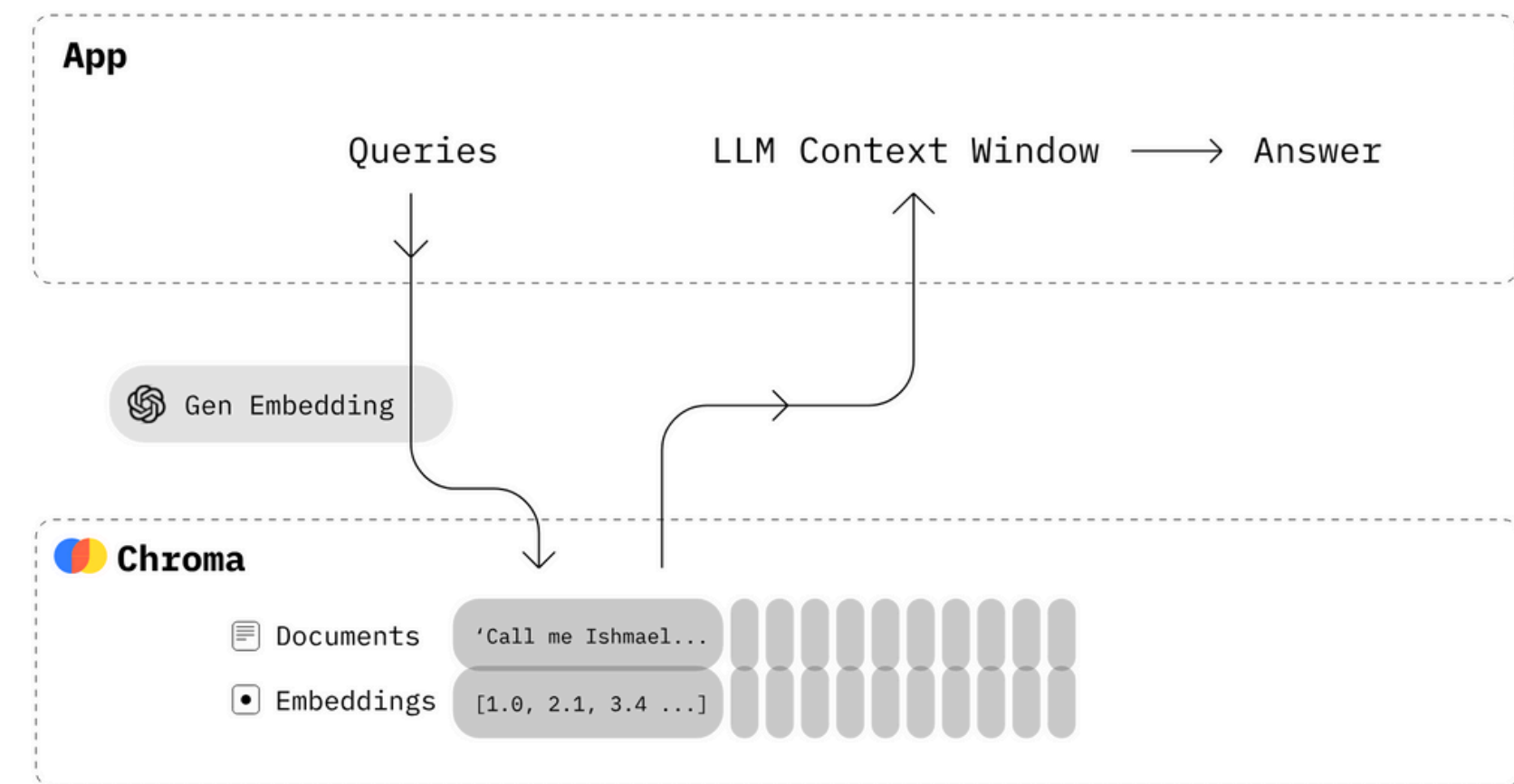
- Build LLM based apps
- Supports APIs, databases, and custom logic for flexible workflows.
- Enables context persistence across multiple interactions.



BUILDING BLOCKS

Chroma as Vector Database

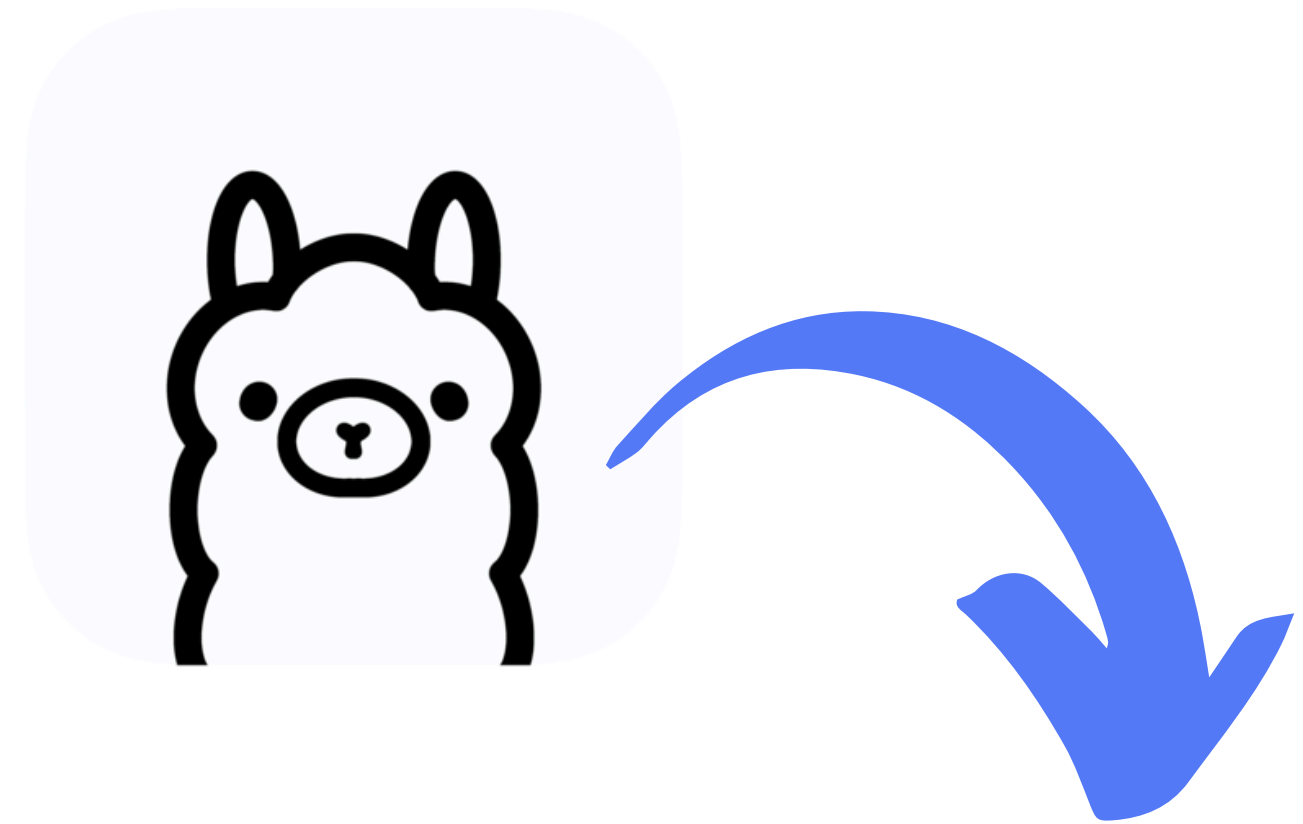
- Specialized for storing and querying high-dimensional embeddings.
- Works with popular ML frameworks like LangChain.
- Enables fast similarity searches for embeddings-based applications.



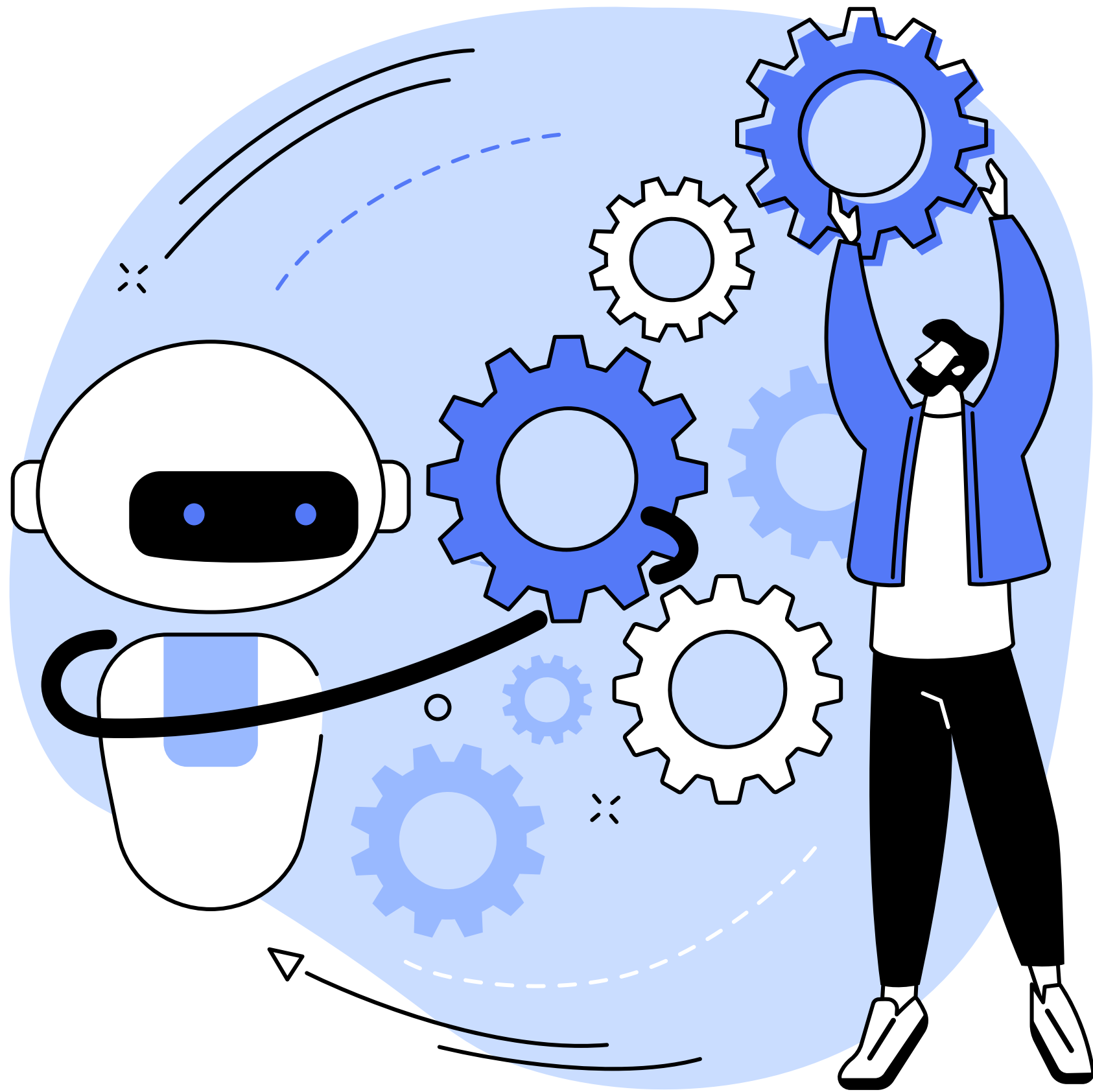
BUILDING BLOCKS

Ollama as LLM Runtime

- Run large language models on local machines efficiently.
- Designed for high-speed inference with minimal resource usage.
- Simple setup for running and experimenting with LLMs on your device.



```
nvidia@jao-60:/$ jetson-containers run $(autotag ollama) ollama run mistral
Namespace(packages=['ollama'], prefer=['local', 'registry', 'build'], disable=[''], user='dustynv',
output='/tmp/autotag', quiet=False, verbose=False)
-- L4T_VERSION=36.2.0 JETPACK_VERSION=6.0 CUDA_VERSION=12.2
-- Finding compatible container image for ['ollama']
cu122/ollama:r36.2.0
+ docker run --runtime nvidia -it --rm --network host --volume /tmp/argus_socket:/tmp/argus_socket
--volume /etc/enctune.conf:/etc/enctune.conf --volume /etc/nv_tegra_release:/etc/nv_tegra_relea
se --volume /tmp/nv_jetson_model:/tmp/nv_jetson_model --volume /var/run/dbus:/var/run/dbus --volu
me /var/run/avahi-daemon/socket:/var/run/avahi-daemon/socket --volume /var/run/docker.sock:/var/r
un/docker.sock --volume /mnt/NVME/jetson-containers/dev/data:/data --device /dev/snd --device /de
v/bus/usb --device /dev/video0 --device /dev/video1 cu122/ollama:r36.2.0 ollama run mistral
pulling manifest
pulling e8a35b5937a5... 82% | 3.4 GB/4.1 GB 38 MB/s
```



IT'S YOUR TURN

