



FOSDEM'24



Using Generative AI and Content Service Platforms together

Angel Borroy
Developer Evangelist



LangChain

neo4j

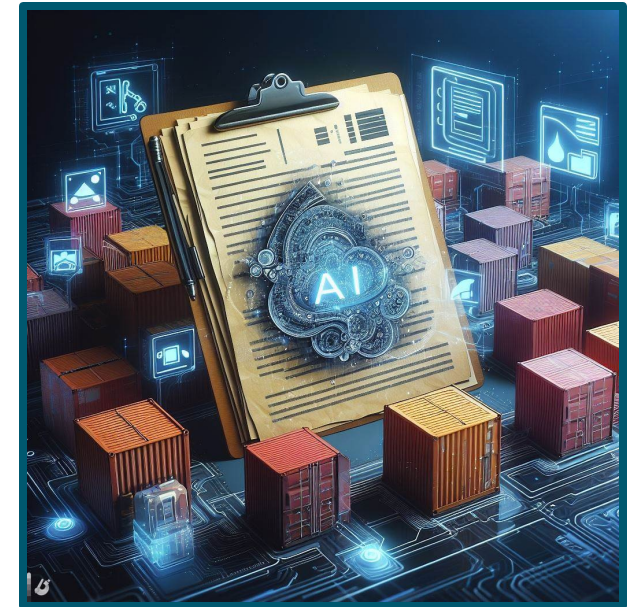


docker

February 1, 2024

Agenda

- GenAI Stack
- Features
- Integration
 - Existing Content
 - New Content
- What Else?



Content credentials
Generated with AI · 25 January 2024 at 10:09 am



GenAI Stack





Components

ollama

- Local management of *open source* LLMs
- Catalog of preconfigured LLMs, such as Llama2 or Mistral



neo4j

- Graph and native vector search capabilities
- Ground LLMs for more precise GenAI predictions and outcomes



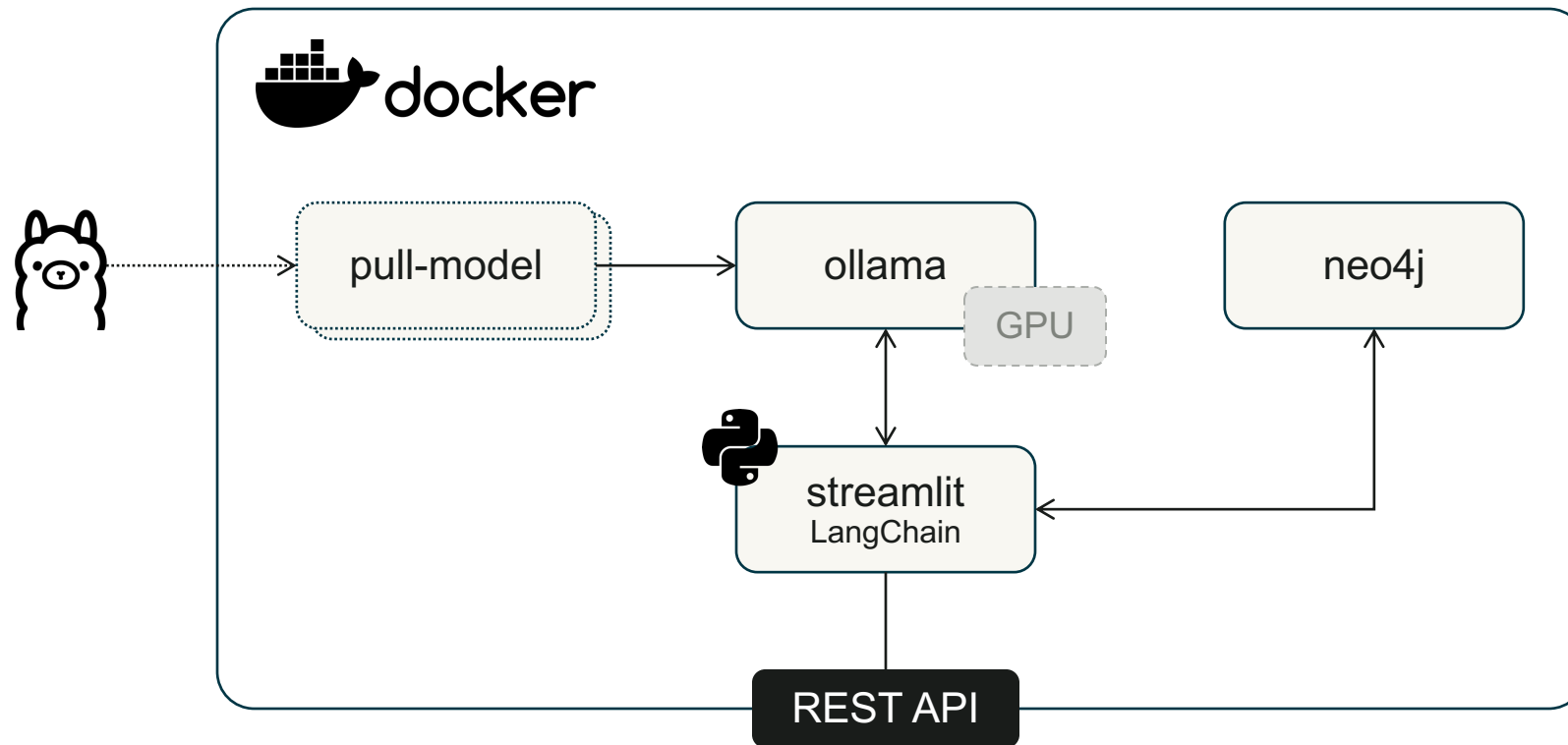
LangChain

- Communication between the LLM, your application, and the database
- Python framework for developing applications powered by LLMs



Deployment

<https://github.com/docker/genai-stack>



Docker GenAI

<https://ollama.ai/library>



40+ GB Size
64+ GB RAM

4 GB RAM
2 GB Size

[Llama2 Community License Agreement](#)
[Llama2 Community License Agreement](#)
[Apache License 2.0](#)
[Deepseek License Agreement](#)
[Llama2 Community License Agreement](#)
[Falcon TII License 1.0](#)

llama2
 codellama
 mixtral
 deepseek-coder
 vicuna
 falcon

13+B

3B

phi
 dolphin-phi
 orca-mini
 deepseek-coder

[MIT License](#)
[MIT License](#)
[cc-by-nc-sa-4.0](#)
[Deepseek License Agreement](#)



[Llama2 Community License Agreement](#)
[Llama2 Community License Agreement](#)
[Llama2 Community License Agreement](#)
[cc-by-nc-sa-4.0](#)
[Apache License 2.0](#)
[Microsoft Research License](#)

llama2
 codellama
 vicuna
 orca-mini
 llava
 orca2

13B

7B

llama2
 codellama
 vicuna
 mistral
 mistral-openorca
 llava
 orca-mini
 deepseek-coder
 orca2
 falcon

[Llama2 Community License Agreement](#)
[Llama2 Community License Agreement](#)
[Llama2 Community License Agreement](#)
[Apache License 2.0](#)
[Apache License 2.0](#)
[Apache License 2.0](#)
[cc-by-nc-sa-4.0](#)
[Deepseek License Agreement](#)
[Microsoft Research License](#)
[Falcon TII License 1.0](#)

8 GB Size
16 GB RAM

8 GB RAM
4 GB Size



deepseek coder



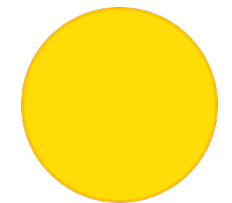


Ethical AI Rating

Mistral 7B



Nextcloud



Is the software (both for inferencing and training) open source?

<https://github.com/mistralai/mistral-src>



Is the trained model freely available for self-hosting?

<https://ollama.ai/library/mistral>



Is the training data available and free to use?

No



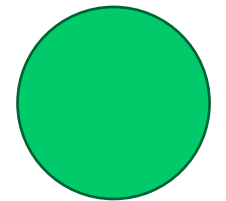


Ethical AI Rating

LlaVA with Visual Encoder



Nextcloud



Is the software (both for inferencing and training) open source?

<https://github.com/haotian-liu/LLaVA>



Is the trained model freely available for self-hosting?

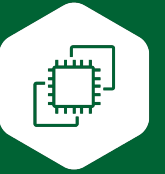
<https://ollama.ai/library/llava>



Is the training data available and free to use?

<https://github.com/haotian-liu/LLaVA/blob/main/docs/Data.md>





MacBook Pro 2021

Apple M1 Pro

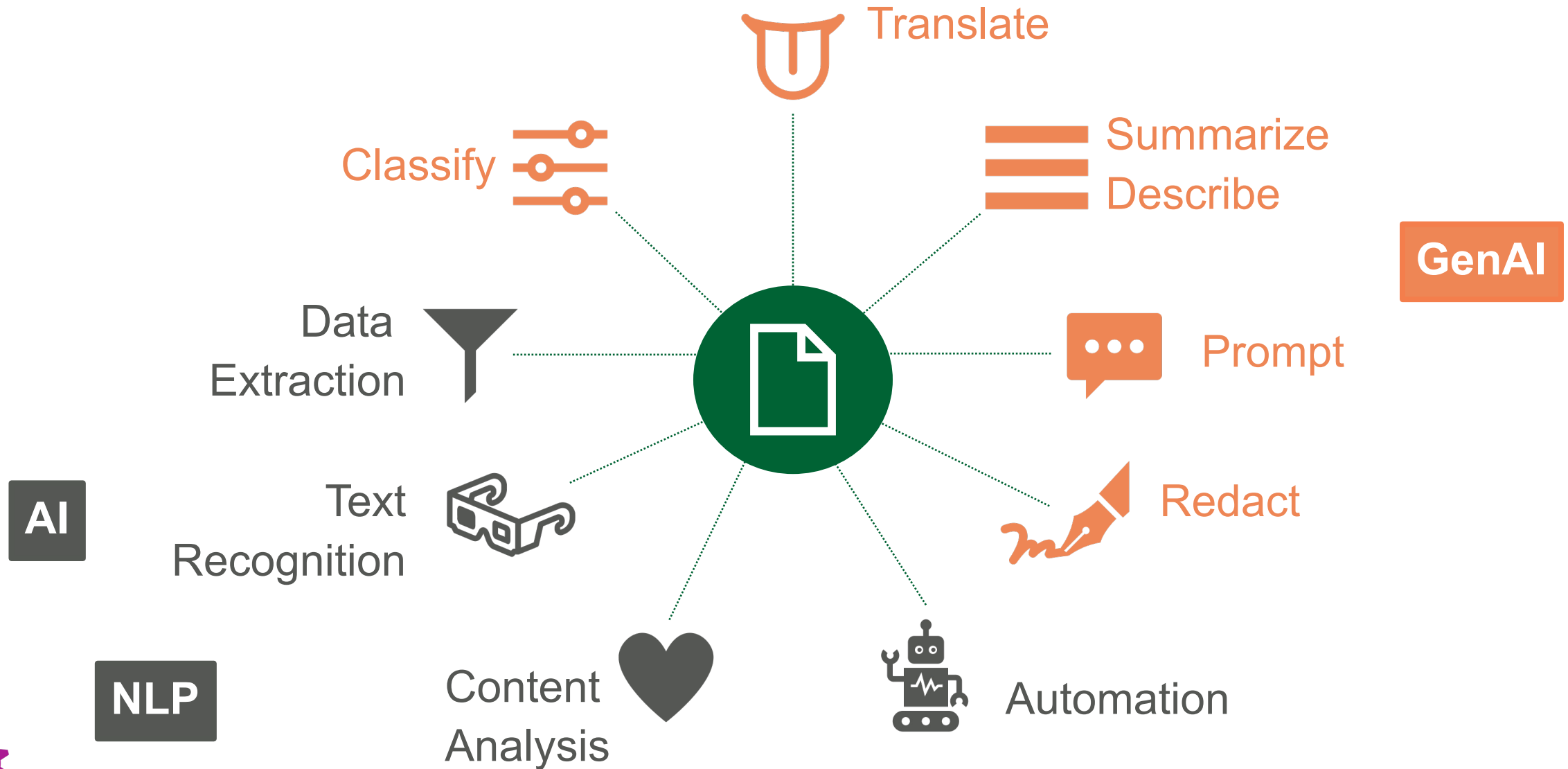
32 GiB RAM

GPU

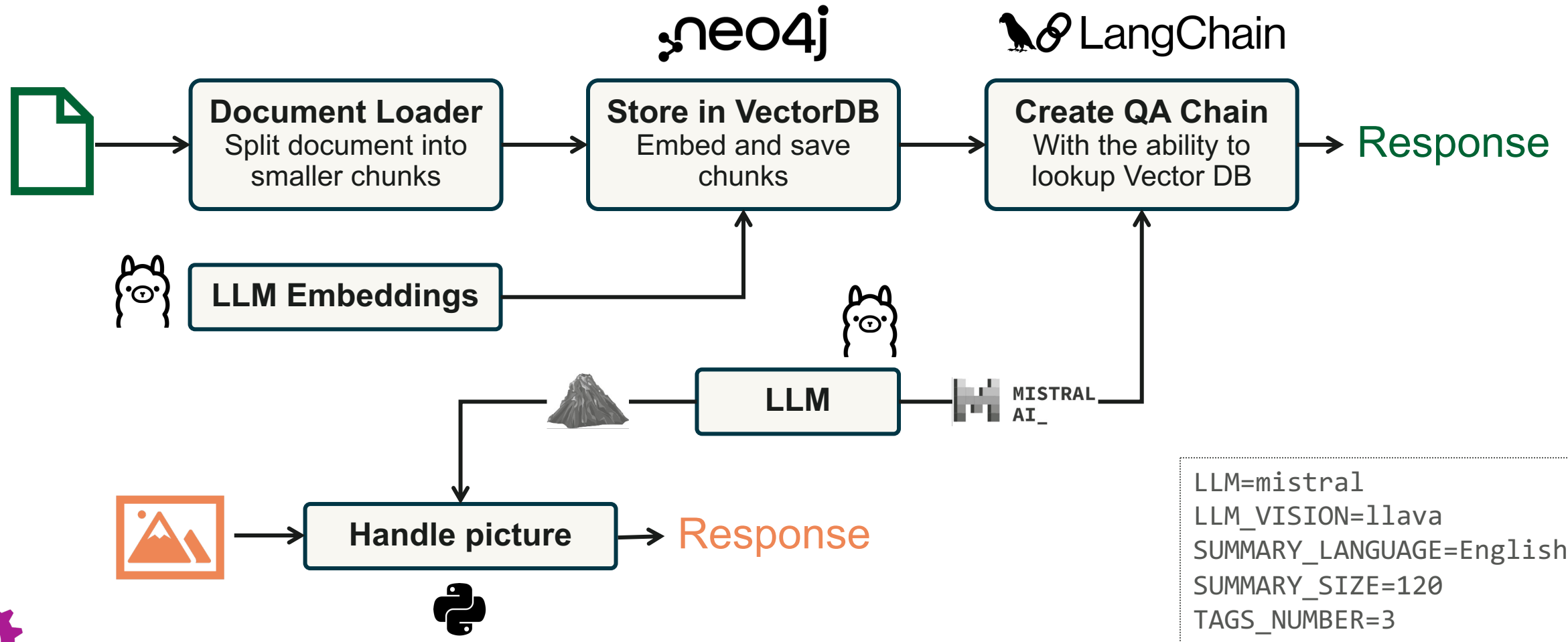
Features




Document Centric Features




Python application that provides REST API endpoints



```
curl --location 'http://localhost:8506/summary' --form 'file=@"/file.pdf"'
{
  "summary": " The text discusses...",
  "tags": " Golang, Merkle, Difficulty",
  "model": "mistral"
}
```

A small white icon of a document with a folded top-right corner, located at the bottom right of the code block.

```
curl --location \
'http://localhost:8506/classify?termList="Japanese,Spanish,Vietnamese"' \
--form 'file=./file.pdf'
{
  "term": " English",
  "model": "mistral"
}
```

A small white icon of a document with a folded top-right corner, located at the bottom right of the code block.

```
curl --location \  
'http://localhost:8506/prompt?prompt="What is the name of the son?"' \  
--form 'file=./file.pdf' \  
{  
  "answer": "The name of the son is Musuko.",  
  "model": "mistral"  
}
```



```
curl --location 'http://localhost:8506/describe' \  
--form 'image=@"image.jpg"' \  
{  
  "description": "The image features a man standing confidently. He is wearing  
                 glasses, a beanie hat, and a jacket.",  
  "model": "llava"  
}
```



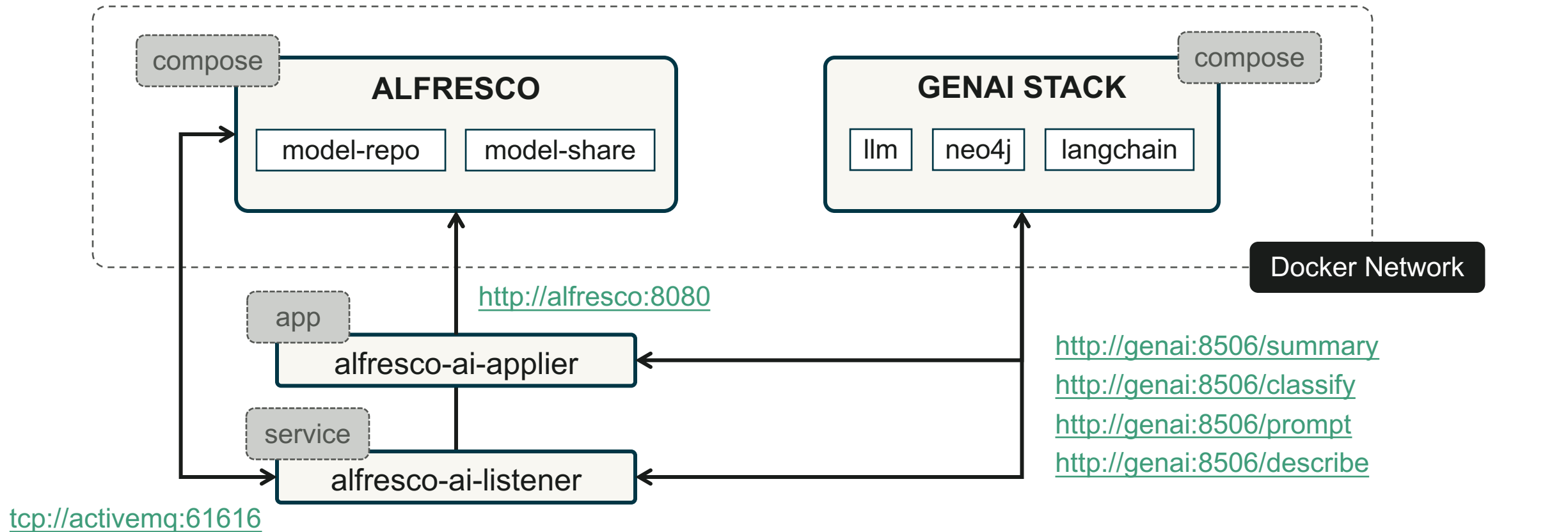
Integration



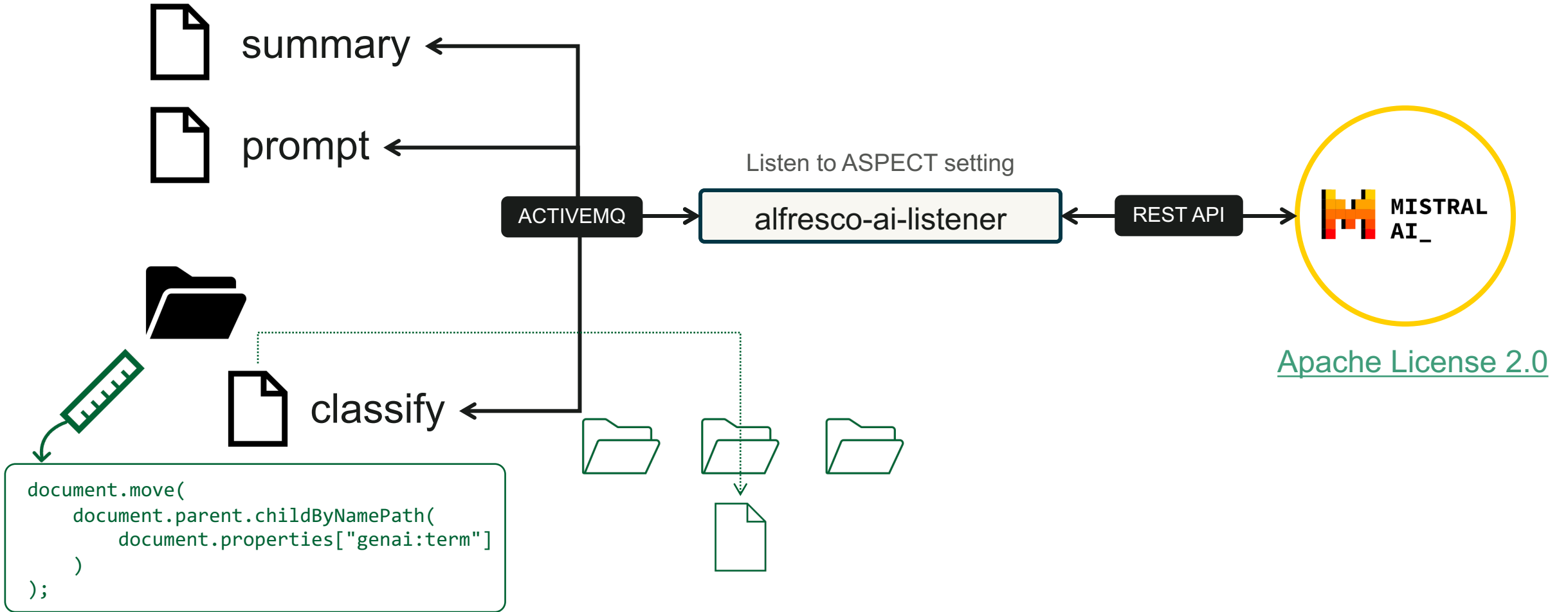
Integration with Content Service Platform



Deployment



Content Service Platform Integration



What Else?



Find your way

Docker AI/ML Hack-a-thon

- [Readme AI](#)
- [Techdocs](#)
- [Docker Image Analyzer](#)
- [Docker Log Sentiment Analyzer](#)
- [GitChats AI](#)

ollama alternatives

- <https://gpt4all.io/index.html>
- <https://localai.io>
- <https://www.secondstate.io/run-llm>
- <https://huggingface.co/docs/hub/spaces-sdks-docker-first-demo>



Hyland™

Thanks!

