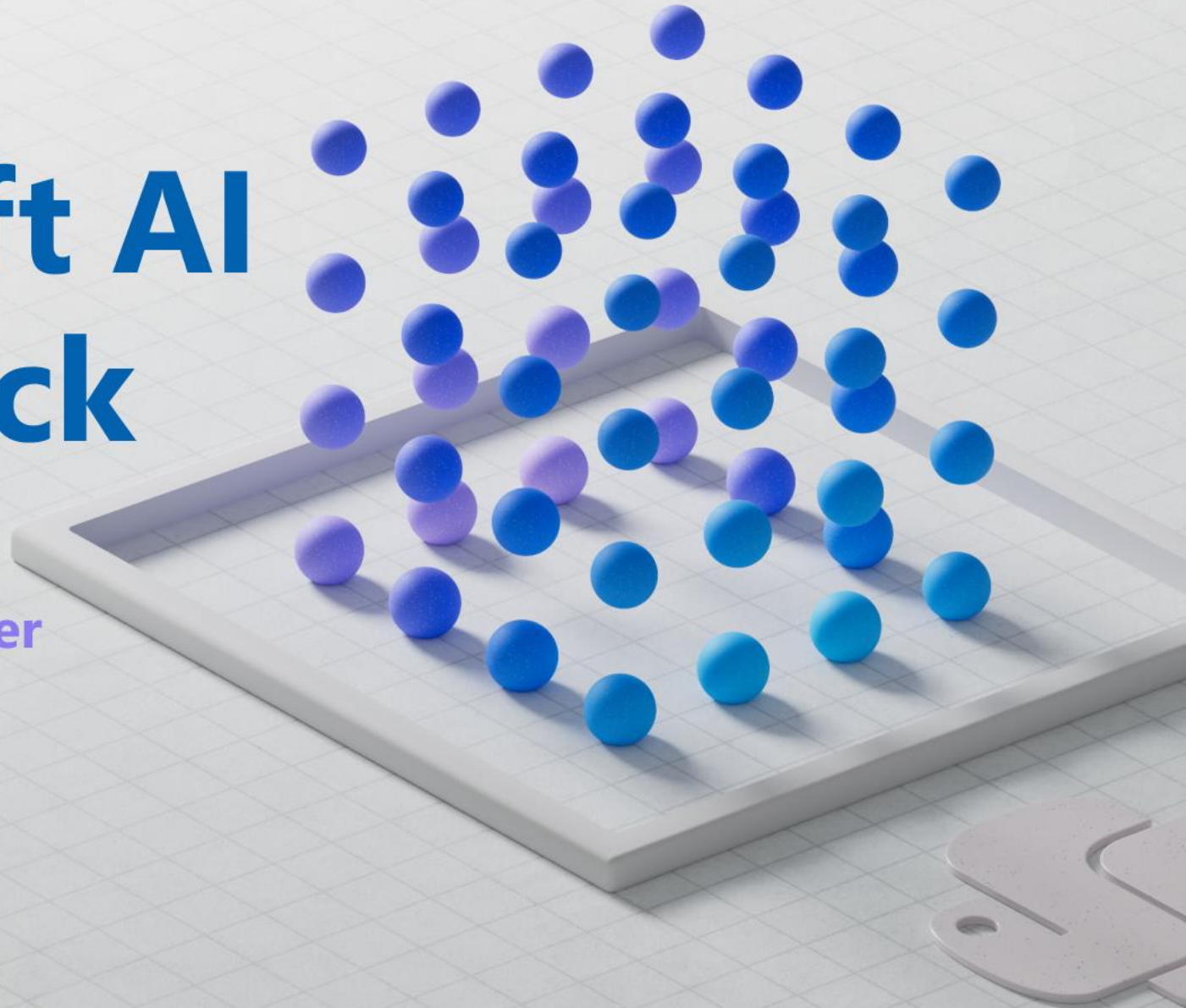January 29th - February 12th

# The Microsoft AI Chat App Hack

**Build, innovate, and #HackTogether**
aka.ms/hacktogether/chatapp
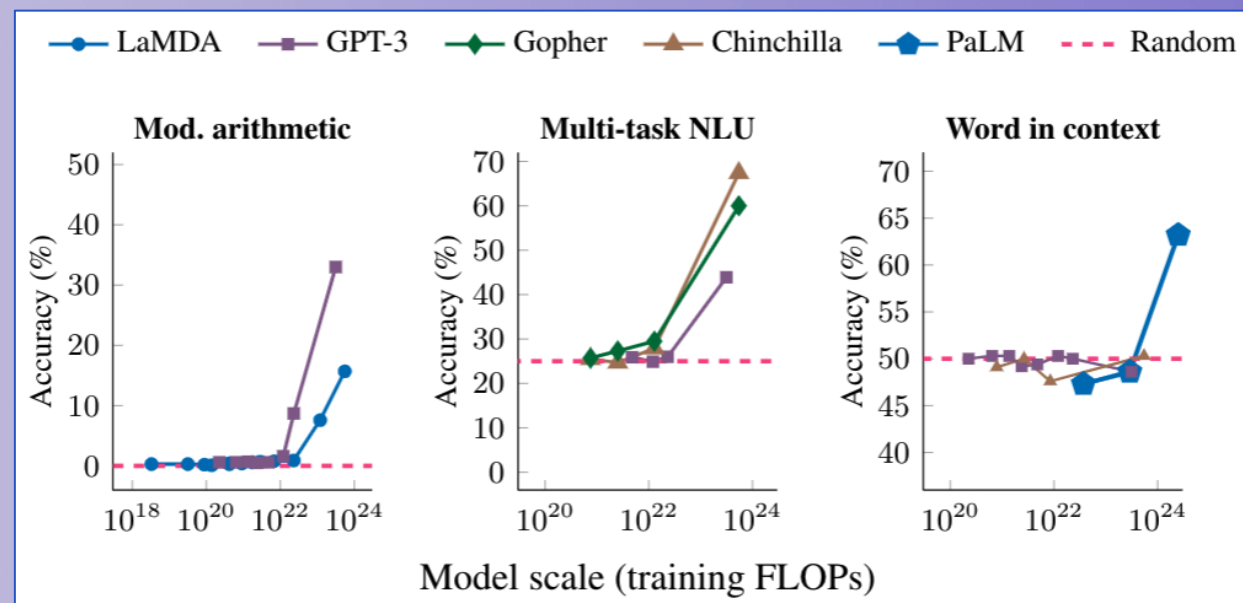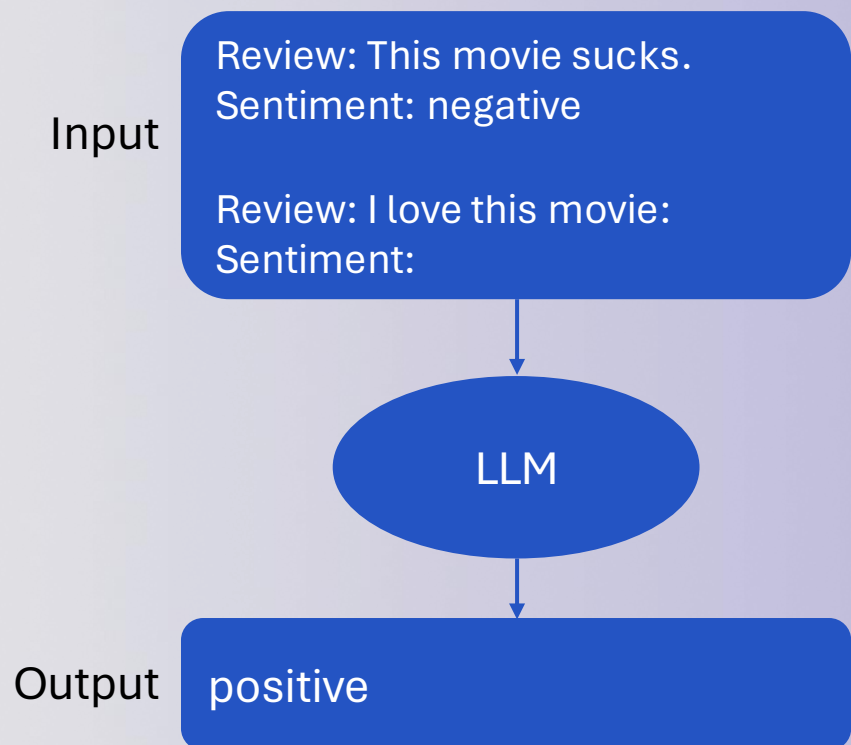
Microsoft

# Goal

Build an AI Chat App using the **RAG** approach to combine an **LLM** with your own data.

# LLMs

# LLM: Large Language Model

An LLM is a model that is so large that it achieves general-purpose language understanding and generation.

Input

Review: This movie sucks.
Sentiment: negative

Review: I love this movie:
Sentiment:

LLM

Output   positive



LaMDA —■— GPT-3 —◆— Gopher —▲— Chinchilla —⬟— PaLM --- Random

**Mod. arithmetic**

Accuracy (%)

**Multi-task NLU**

Accuracy (%)

**Word in context**

Accuracy (%)

Model scale (training FLOPs)

# LLMs in use today

| Model | # of Parameters | Creator | Uses |
| --- | --- | --- | --- |
| GPT 3.5 | 175 B | OpenAI | ChatGPT, Copilots, APIs |
| GPT 4 | Undisclosed | OpenAI | |
| PaLM | 540 B | Google | Bard |
| Claude 2 | 130 B | Anthropic | APIs |
| LlaMA | 70 B | Meta | OSS |
| Mistral-7B | 7 B | Mistral AI | OSS |

Microsoft

# GPT: Generative Pre-trained Transformer

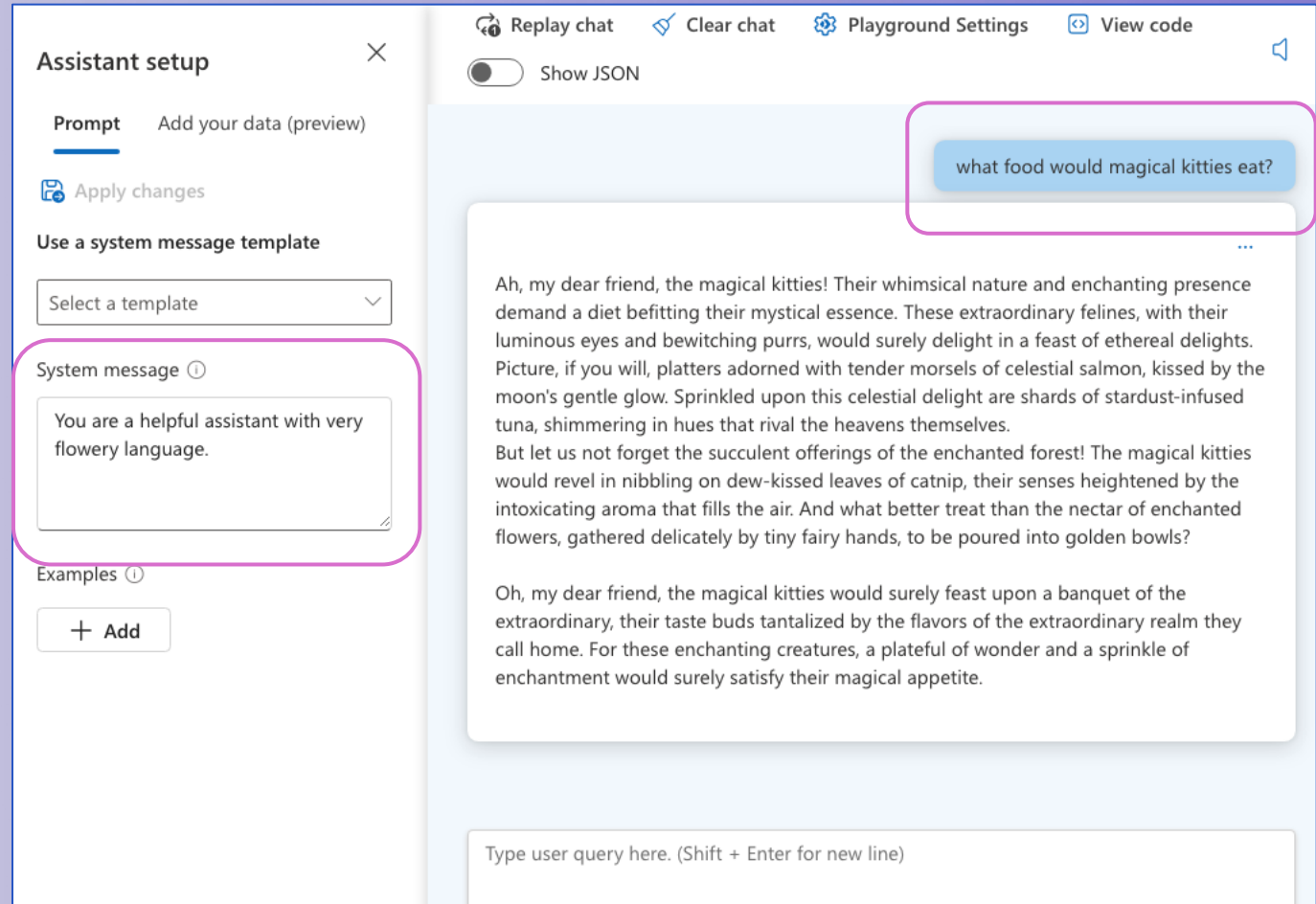GPT models are LLMs based on Transformer architecture from "Attention is all you need" paper

Learn more:
- Andrej Karpathy: 🎥 State of GPT
- Andrej Karpathy: 🎥 Let's build

GPT: from scratch, in code

# Using OpenAI GPT models: Python 🐍

```python
response = openai.ChatCompletion.create(
  stream=True,
  messages = [
   {
    "role": "system",
    "content": "You are a helpful assistant with very flowery language"
   },
   {
    "role": "user",
    "content": "What food would magical kitties eat?"
   }
])

for event in response:
  print(event.choices[0].delta.content)
```

# The limitations of LLMs



Outdated public knowledge

No internal knowledge

# RAG: Retrieval Augmented Generation

Microsoft

Do my company perks cover underwater activities?

Yes, your company perks cover underwater activities such as scuba diving lessons [1]

User Question

Document Search

PerksPlus.pdf#page=2: Some of the lessons covered under PerksPlus include: · Skiing and snowboarding lessons · Scuba diving lessons · Surfing lessons · Horseback riding lessons These lessons provide employees with the opportunity to try new things, challenge themselves, and improve their physical skills.....

Large Language Model

# The benefit of RAG



Up-to-date public knowledge

Access to internal knowledge

# RAG components

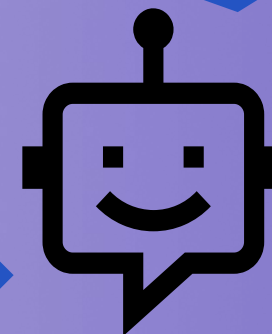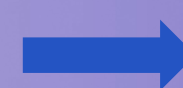| Component | Examples |
|---|---|
| **Retriever**: A knowledge base that can efficiently retrieve sources that match a user query | Azure AI Search, Azure CosmosDB, PostgreSQL, Qdrant, Pinecone |
| **LLM**: A model that can answer questions based on the query based on the provided sources, and can include citations | GPT 3.5, GPT 4 |
| **Glue**: A way to chain the retriever to the LLM (optional) | Langchain, Llamaindex, Semantic Kernel |
| **Features** | Chat history, Feedback buttons, Text-to-speech, User login, File upload, Access control, etc. |

# Copilot Studio – On Your Data



**Retriever:**
  Uploaded files

**LLM**: GPT 3.5

https://copilotstudio.preview.microsoft.com/

# Azure Studio – On Your Data



**Retriever:**
   Azure AI Search
   Azure Blob Storage
   Azure CosmosDB for MongoDB vCore
   URL/Web address
   Uploaded files

**LLM**: GPT 3.5/4

**Features**:
   User authentication
   Chat history persistence

https://learn.microsoft.com/azure/ai-services/openai/concepts/use-your-data

# Open source RAG chat app solution



**Retriever:**
   Azure AI Search

**LLM**: GPT 3.5/4

**Features**:
   Multi-turn chats
   User authentication with ACLs
   Chat with image documents

**https://github.com/Azure-Samples/azure-search-openai-demo/**     **aka.ms/ragchat**

# Deep dive:
# RAG chat app
# solution

# Prerequisites

- Azure account and subscription
  - A free account can be used, but will have limitations.
- Access to Azure OpenAI or an openai.com account
  - Request access to Azure OpenAI today!

  https://aka.ms/oaiapply

https://github.com/Azure-Samples/azure-search-openai-demo/#azure-account-requirements
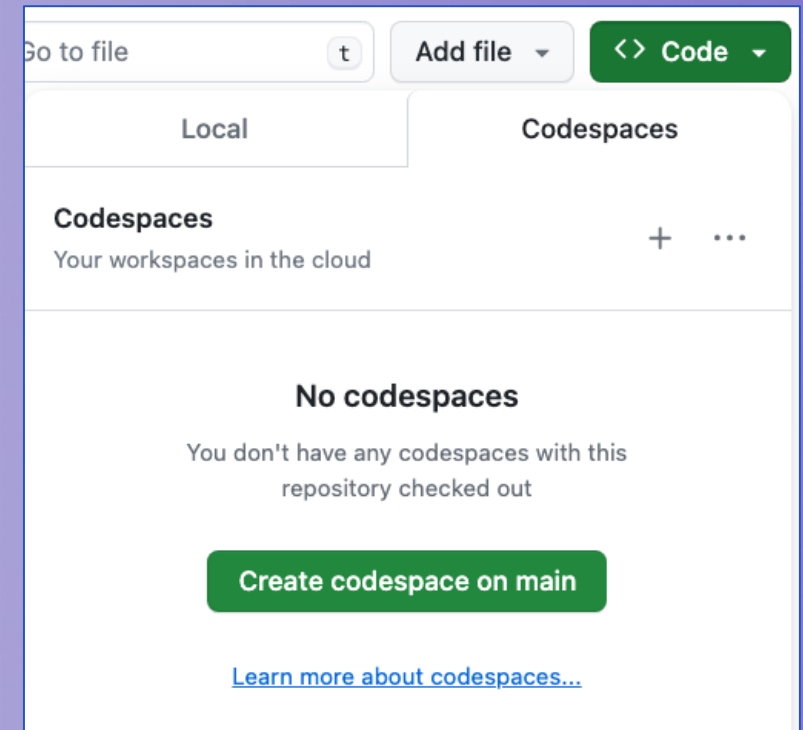
# Opening the project: 3 options

- GitHub Codespaces →

- VS Code with Dev Containers extension

- Your Local Environment
  - Python 3.9+
  - Node 14+
  - Azure Developer CLI

https://github.com/Azure-Samples/azure-search-openai-demo/?tab=readme-ov-file#project-setup

# Deploying with the Azure Developer CLI

Login to your Azure account:

```
azd auth login --use-device-code
```

Create a new azd environment: (to track deployment parameters)
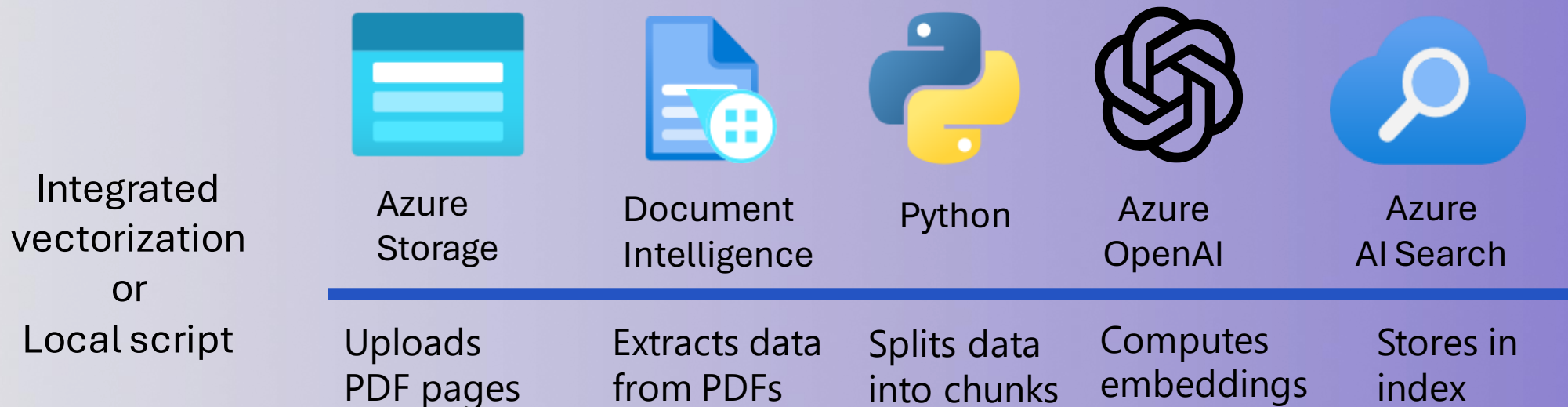
```
azd env new
```

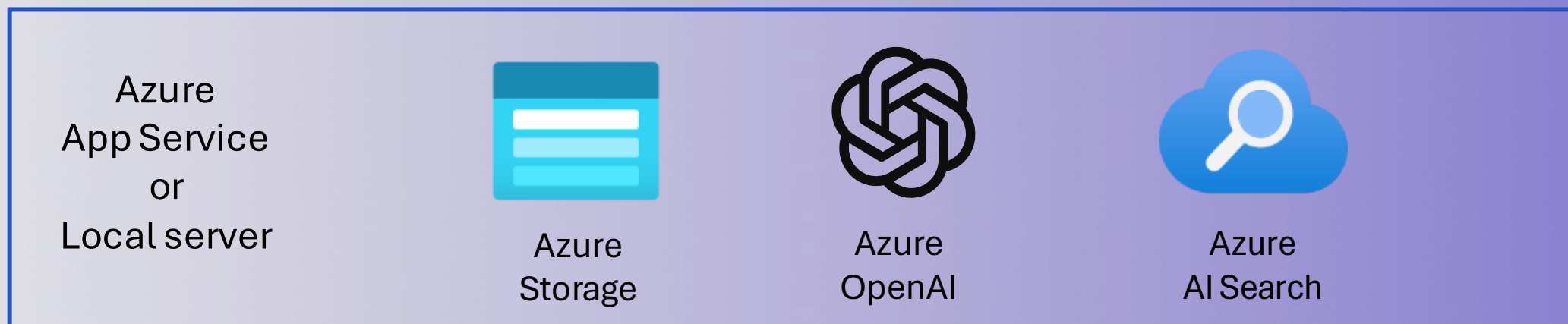Provision resources and deploy app:

```
azd up
```

azd up is a combination of azd provision and  azd deploy

# Code walkthrough

**Typescript frontend**

(React, FluentUI)

**Python backend**

(Quart, Uvicorn)

*chat.tsx*
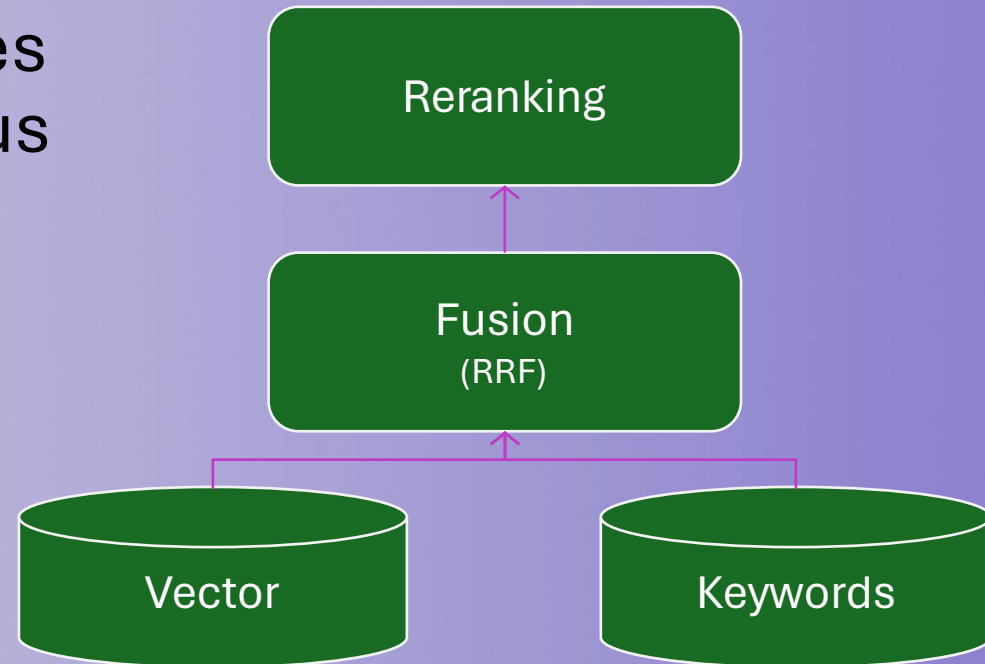makeApiRequest() → *api.ts*
chatApi()

*app.py*
chat() → *chatreadretrieveread.py*
run()

get_search_query()
compute_text_embedding()
search()
get_messages_from_history()
chat.completions.create()

# Search approach

For optimal retrieval, search() uses hybrid retrieval (text + vectors) plus the semantic ranker option.

https://aka.ms/ragrelevance



🔊 Learn more at this week's session: Azure AI Search Best Practices

# Next steps

- Register for the hackathon → **aka.ms/hacktogether/chatapp**

- Introduce yourself in our discussion forum

- Deploy the repo with the sample data
  - See steps on low cost deployment → **aka.ms/ragchat/free**

- Post in forum if you have any questions or issues deploying.

- Join tomorrow's session: Customizing your RAG Chat App!

Microsoft