Tech Community Sharing
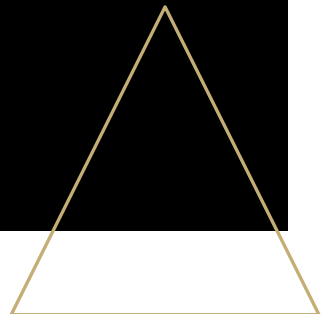
# Azure OpenAI 集成
# 云原生应用

# Overview

Architecture

Semantic Kernel
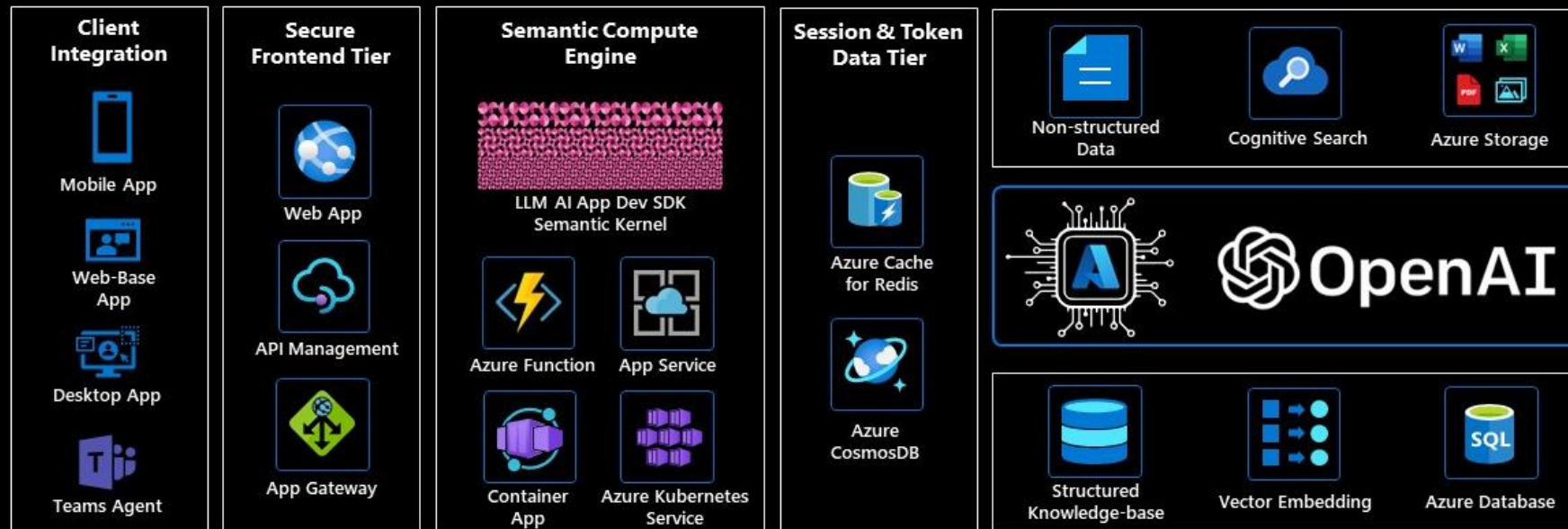
Q & A

Details

Demo

# Enterprise GPT-ize Intelligent Application Architecture



A prompt engineering approach with LLM AI Dev Framework

Azure-OpenAI-App-Innovation-Workshop

Azure Cognitive Search can quickly index unstructured data such as PDF and WORD files, allowing existing data to be used immediately.

Utilizing Azure Database with vector storage and processing capabilities and combining them with AOAI's embedding vector generation model, the enterprise's existing structured knowledge base can be integrated easily
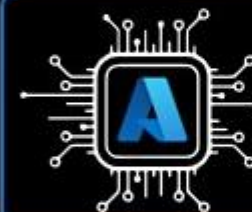
## Dual-engine

Non-structured Data

Cognitive Search

Azure Storage
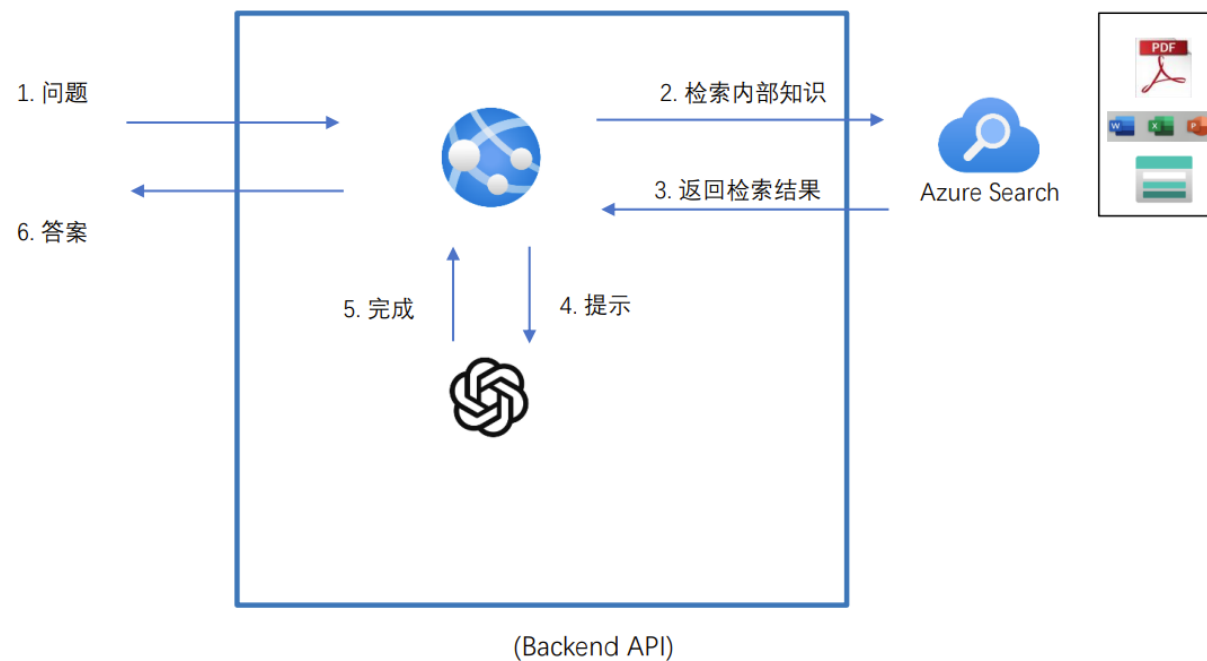
OpenAI

Structured Knowledge-base

Vector Embedding

Azure Database

# GPT + enterprise knowledge base/data
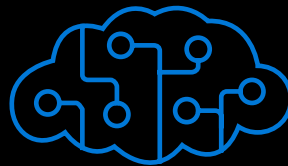
# Azure Cognitive Search

| ingest | abundant | Browse |
|--------|----------|--------|

Cognitive skills

101010
010101
101010

Data in any
format, any
Azure store

Annotations

Search

# Cognitive Search Architecture



Cognitive Skill

Document cracking

Skillset:
Extensible and scalable pipeline

Search index

**client data**

.pdf
.doc
.jpeg
...

**Azure has built-in cognitive skills OCR,
keyword extraction,
Sentiment analyzer, computer vision,
Form Recognizer,
ink recognizer,
entity linking,**

Text Analytics
translation
Content Moderator
personalize
QnA manufacturer
speech

**Annotated Documents**

**Search Index**

**Third Party Enrichers**
custom classification model,
custom entity extraction,

Azure machine Learning

Web

Mobile

Bots

PowerBI

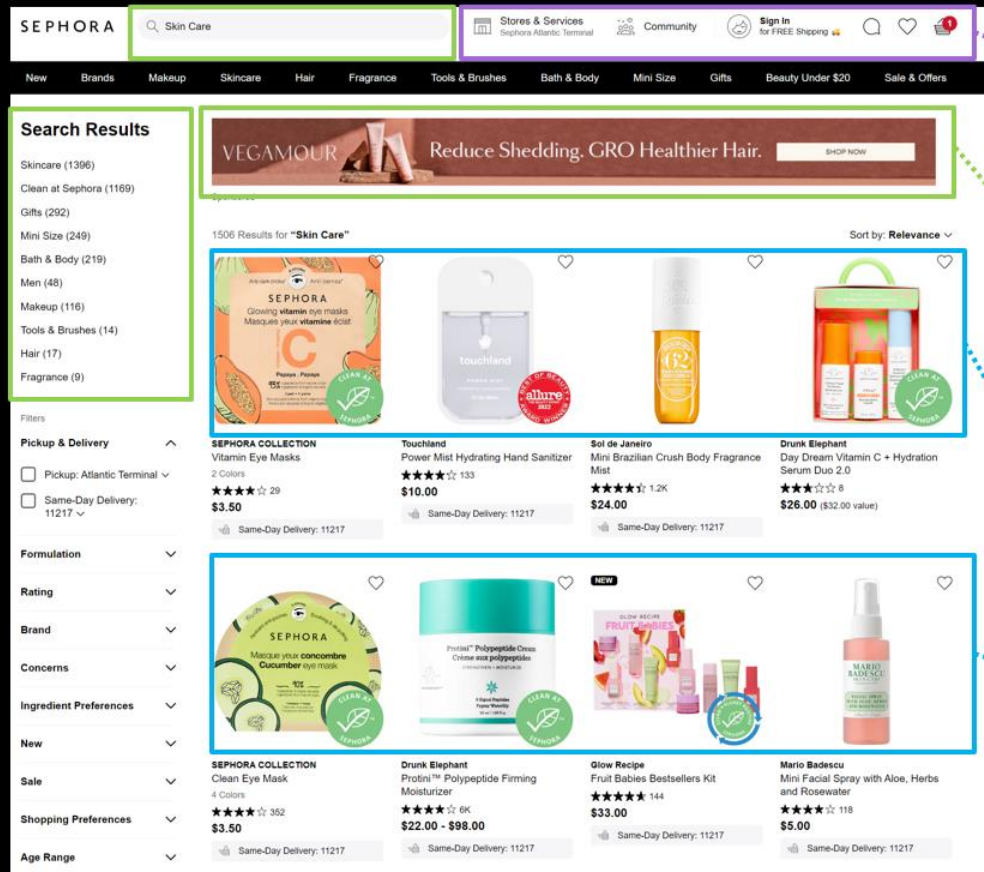**Session & Token Data Tier**

Azure Cache for Redis

Azure CosmosDB

- Based on Redis and CosmosDB.
- Add contextual caching, session persistence, prompt persistence, and other capabilities to the application.
- Leave room for future model or engine optimization based on prompts.

# How is Azure Cache for Redis Leveraged?



**Session Store**
Easily resume a user's journey on any mobile or web application by quickly sorting the data of their last session regardless of where they are across the globe

**Message Broker Store**
Azure Cache for Redis can communicate across all Microservices. In the example, it communicated with its ad server to give a recommended ad for the shopper
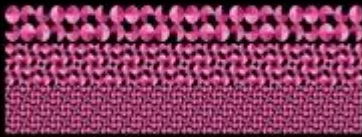
**Distributed Cache**
All images need to be consistently accessed and are static therefore, minimizes load times by keeping all assets into a cache no matter where the user is across the globe

- The API encapsulation, load balancing, and gateway at the front end further enhance application security and reliability.
- Allow the intelligent entity to integrate with various front-end apps in a more secure and stable manner.

Semantic Compute Engine

LLM AI App Dev SDK
Semantic Kernel

Azure Function    App Service

Container App    Azure Kubernetes Service

- Utilize a flexible semantic compute engine as prompt engine
- Support various deployment forms such as PaaS, Serverless, and containers
- Enable enterprises to optimize the prompt engine with modern LLM AI application development framework
- Elastically scale computing resources according to real-time business needs.

# Function as a Service

Use serverless code to handle events.

- Writing cloud apps made easy
- Extend functionality based on customer needs
- Develop functions in languages such as C#, Node.js, F#, Python, Java, and more
- Easily schedule event-driven tasks across services
- Exposes the function as an HTTP API endpoint

# Azure Functions Common application scenarios

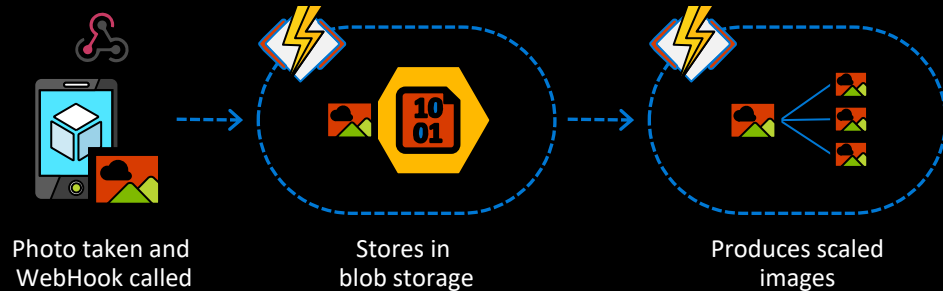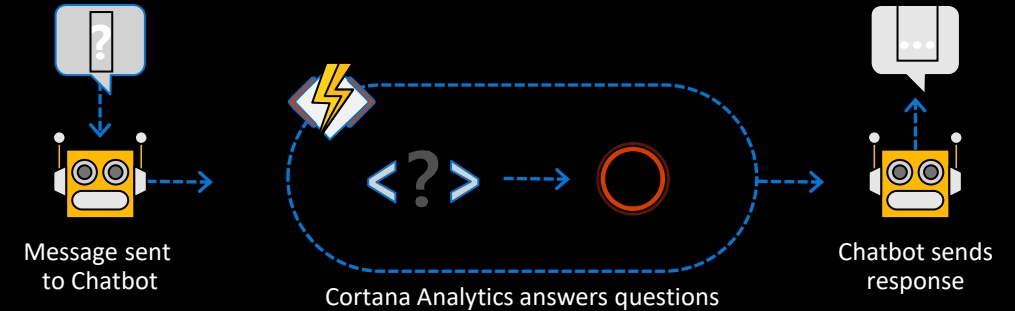Any scenario that requires a response to an event



**Real-time stream processing**

Millions of devices feed into Stream Analytics → Transform to structured data → Store data in SQL DB

**Timer-based processing**

Every 15 minutes → Find and clean invalid data → Clean table

**Backend (Mobile/IoT/Web)**

Photo taken and WebHook called → Stores in blob storage → Produces scaled images

**Real-time bot messaging**

Message sent to Chatbot → Cortana Analytics answers questions → Chatbot sends response

# Semantic Kernel - LLM AI Dev Framework

## 💪 THE CORE

Everything you need to manage complex prompts, chains, long-running tasks, and **planning.**

SKILLS

PLANNER
1 2 3 ...

RUNTIME
... > 3 > 2 > 1 >

MEMORIES

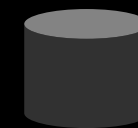CONNECTORS

## 👶 THE BFD
But wait, there's more! All the longform **memories** readily available in the MS Graph. Plus, the 900 Power **connectors**.
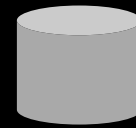
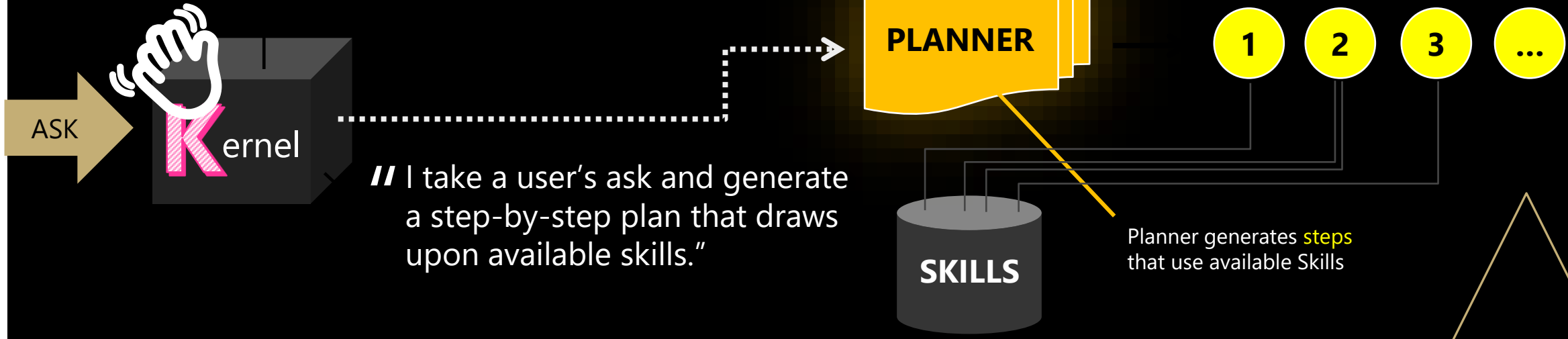SUBSTRATE    LINKEDIN    GITHUB

POWER PLATFORM

# Meet the lightweight Kernel of Semantic Kernel

```
using Microsoft.SemanticKernel;

var myKernel = Kernel.Builder.Build();
```

" I've been designed to reduce hallucinations, orchestrate complicated LLM AI prompts combined with native code, use multiple AI models, and ... I have a special skill to **PLAN**."

**STEPS**

ASK

**K**ernel

**PLANNER**

**1** **2** **3** **...**

" I take a user's ask and generate a step-by-step plan that draws upon available skills."

**SKILLS**

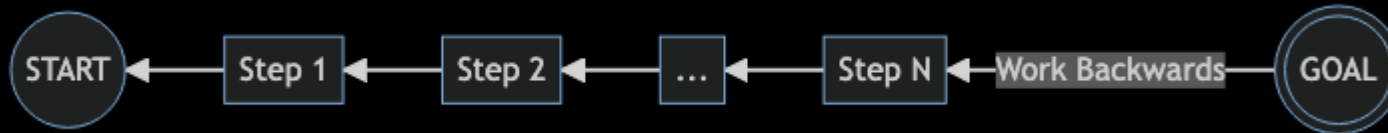Planner generates steps that use available Skills

# A more "goal-oriented" approach to problem solving

# Planner

- The planner will operate within the skills it has available. In the event that a desired skill does not exist, the planner can suggest you create the skill. Or, depending upon the level of complexity the kernel can help you write the missing skill.

START ← Step 1 ← Step 2 ← ... ← Step N ← Work Backwards — GOAL

# Skill

- an LLM AI prompt — also called a "semantic" function
- native computer code -- also called a "native" function

```
SkillName (directory name)
|
|_____ Function1Name (directory name)    skprompt.txt
|                                         config.json
|
|_____ Function2Name (directory name)
```

```
MyAppSource
|
|_____ MySkillsDirectory
|       |
|       |_____ MySemanticSkill (a directory)
|       |       |
|       |       |_____ MyFirstSemanticFunction (a directory)
|       |       |_____ MyOtherSemanticFunction (a directory)
|       |
|       |_____ MyNativeSkill.cs (a file)
|       |_____ MyOtherNativeSkill.cs (a file)
```

COLLECTION OF SKILLS

Skill A

Skill B

Skill C

Function A-1      Function A-2      Function B-1      Function C-1      Function C-2      Function C-3      Function C-4

# Memory

**Conventional key-value pairs**: Just like you would set an environment variable in your shell, the same can be done when using SK. The lookup is "conventional" because it's a one-to-one match between a key and your query.

**Conventional local-storage**: When you save information to a file, it can be retrieved with its filename. When you have a lot of information to store in a key-value pair, you're best off keeping it on disk.

**Semantic memory search**: You can also represent text information as a long vector of numbers, known as "embeddings." This lets you execute a "semantic" search that compares meaning-to-meaning with your query.

# Connector

## MS Graph Connector Kit

- Add an event to your calendar
- Send an email for you
- Add a file to your OneDrive
- Create a share link to a file in your OneDrive
- Query your organization hierarchy
- Manage your MS To Do list

## Out-of-box

- Issue a Bing search query
- Read OpenXML streams (e.g. Word docs)
- Use SQLite as a lightweight database

# SK makes app developers' work lives easier

Fast Integration: SK is designed to be embedded in any kind of application, making it easy for developers to add LLM AI functionality to test inside their apps.

Power Prompting: Plain prompts that are fed as API calls can only get you so far. SK provides the abstractions and machinery to unlock your OpenAI or Azure OpenAI API key.

Novel-But-Familiar: For 100% determininism, native code is always available as a first-class partner on your prompt engineering quests. You get the best of both worlds.

# How to use Semantic Kernel R1 in just 1 minute

**⏱ 1 MINUTE**

Install the nuget package and go

#r "nuget: Microsoft.SemanticKernel, *-*"

**⏱ HAVE MORE TIME?**

Go deeper with the GitHub repo aka.ms/skrepo

**⏱ MORE MINUTES**

Learn more about its history aka.ms/sk

# Demo

# Semantic Kernel Service API

semantic-kernel/samples/dotnet/KernelHttpServer

>>func start –csharp

[2023-04-02T13:18:19.291Z] Found C:\Users\huolu\Desktop\semantic-kernel\samples\dotnet\KernelHttpServer\KernelHttpServer.csproj. Using for user secrets file configuration.

Functions:

        ExecutePlan: [POST] http://localhost:7071/api/planner/execute/{maxSteps?}

        InvokeFunction: [POST] http://localhost:7071/api/skills/{skillName}/invoke/{functionName}

        Ping: [GET] http://localhost:7071/api/ping

- Written in C# against Azure Function Runtime v4
- Expose some Semantic Kernel APIs that you can call via HTTP POST requests

Azure Functions Core Tools

# GitHub Repo Q&A Bot



semantic-kernel/samples/apps/github-qna-webapp-react
>>yarn install
>>yarn start

# Simple Chat Summary App

# Workaround(Optional)

```
samples > apps > github-qna-webapp-react > src > hooks > TS useSemanticKernel.ts > useSemanticKernel
1    // Copyright (c) Microsoft. All rights reserved.
2
3    import React from 'react';
4    import { SemanticKernel } from './SemanticKernel';
5
6    export const useSemanticKernel = (uri: string) => {
7        const [semanticKernel] = React.useState(new SemanticKernel('http://localhost:7071'));
8        return semanticKernel;
9    };
10
```

url

# Q & A

- [Semantic Kernel - Home (sharepoint.com)](Semantic Kernel - Home (sharepoint.com))

- [Azure-OpenAI-App-Innovation-Workshop/Workshop Content EN at main · xuhaoruins/Azure-OpenAI-App-Innovation-Workshop (github.com)](Azure-OpenAI-App-Innovation-Workshop/Workshop Content EN at main · xuhaoruins/Azure-OpenAI-App-Innovation-Workshop (github.com))