



An Introduction to OpenAI

2023.08

Jeff Fattic

App Innovation Lead | Microsoft

Randy Pagels

DevOps Architect | Xpirit USA

<https://github.com/microsoft/globalopenaihack>
SemanticKernel

[OPTIONAL]

- 1. Clone this repo**
- 2. Create Settings.json in Tutorial00**



Artificial Intelligence

The AI Journey – How we got here!
Some AI vocabulary you should know
OpenAI + Azure
MSFT Responsible AI
Real World Case Studies
Learning More!

Brief history of AI

Artificial Intelligence

Machine Learning

Deep Learning

Generative
AI



Artificial Intelligence

The field of computer science that seeks to create intelligent machines that can replicate or exceed human intelligence



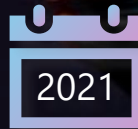
Machine Learning

Subset of AI that enables machines to learn from existing data and improve upon that data to make decisions or predictions



Deep Learning

A machine learning technique in which layers of neural networks are used to process data and make decisions



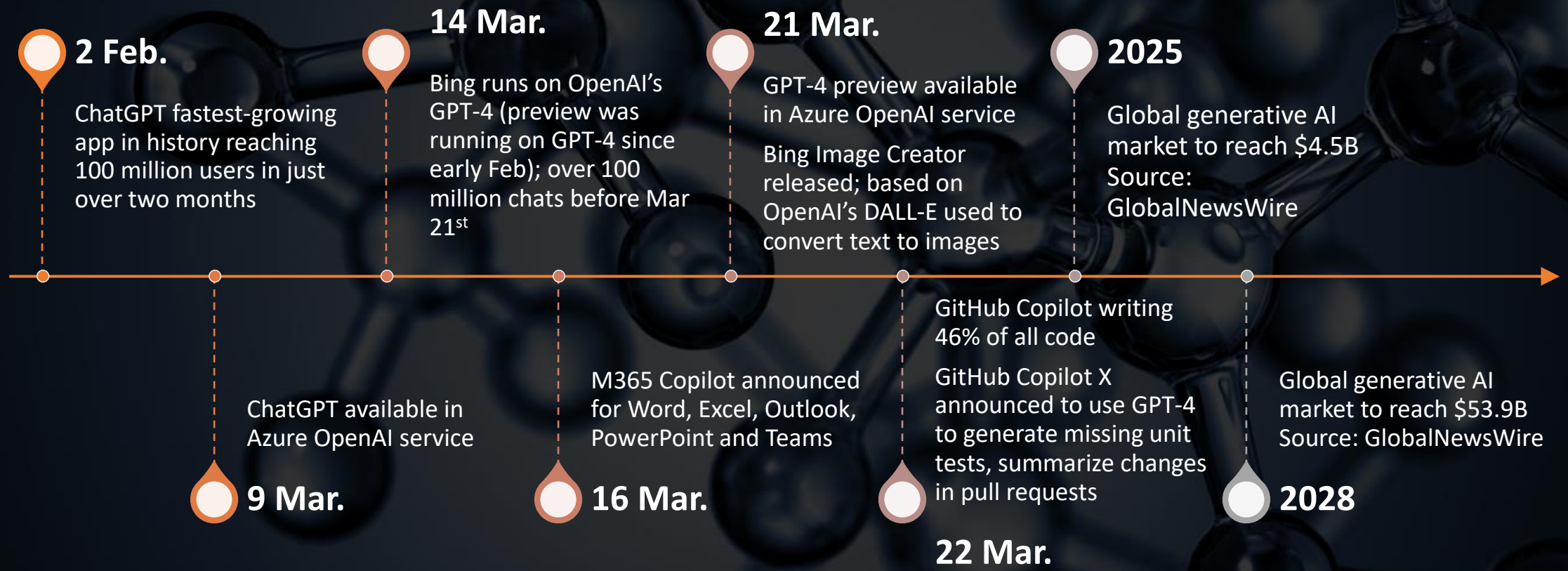
Generative AI

Create new written, visual, and auditory content given prompts or existing data

And then ChatGPT...



And more recently...



Some AI Vocabulary

OpenAI

ChatGPT

Prompts

Prompt
Engineering

Completions

Tokens

Fine-Tuning

Embeddings

Copilot

Semantic
Kernel

What is Azure OpenAI Service?

- One of several Azure AI services. Azure OpenAI Service allows developers to use large language generative AI models for enterprise-grade applications.
- With built-in responsible AI and enterprise-grade security, the service is designed to detect and mitigate harmful use.



GPT-3 / GPT-4

Generate & understand text



Codex (deprecated)

Generate & understand code



DALL-E 2

Generate images from text prompts

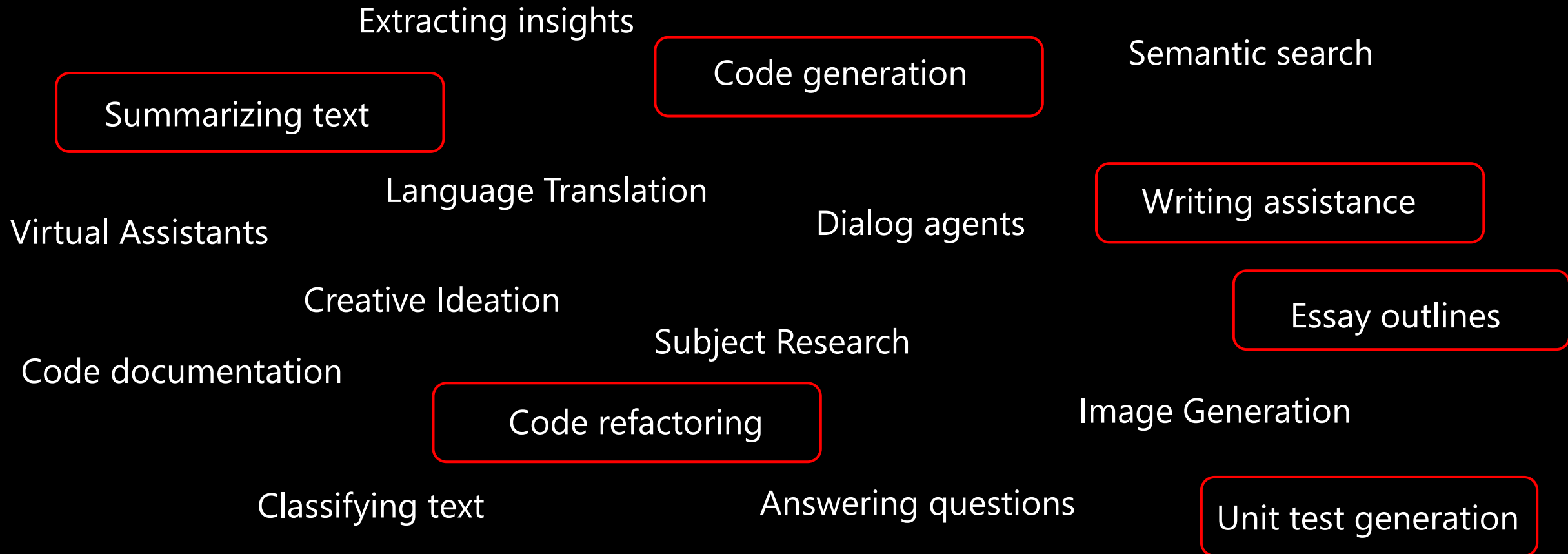


ChatGPT

Personalize conversations

Azure OpenAI Service Capabilities

- Use cases hitting every industry: education, healthcare, entertainment, financial services, retail, manufacturing, and more!



Anatomy of Azure OpenAI Service (Completions)

The screenshot shows the Azure OpenAI Completions playground interface. The left sidebar contains navigation links: Azure OpenAI, Playground, Chat, Completions (selected), DALL-E (Preview), Management, Deployments, Models, Data files, Quotas, and Content filters (Preview). The main area is titled 'Completions playground' and shows the 'davinci-003' deployment with a 'Load an example' button and a 'View code' link. The example text is: 'The only thing we have to fear is fear itself. This phrase was coined by Franklin D. Roosevelt in his first inaugural address in 1933. He was speaking about the fear and uncertainty that was gripping the nation during the Great Depression. He was trying to reassure the American people that they could overcome their fears and that the nation would be able to weather the storm.' The 'Parameters' panel on the right includes: Temperature (1), Max length (tokens) (100), Stop sequences (empty), Top probabilities (0.5), Frequency penalty (0), Presence penalty (0), Best of (1), Pre-response text (empty), and Post-response text (empty). Red arrows point from the parameter names to their respective controls. The bottom of the interface shows a 'Generate' button, a 'Tokens: 74' indicator, and a 'Privacy & cookies' link.

Temperature
Controls randomness
Generate more creative completions.

Tokens
Limits number of tokens per model.
A token approximately 4 characters.

Top probabilities
Controls randomness

Frequency penalty
Decreases likelihood of repeating exact same text in a response.

Presence penalty
Increases likelihood of introducing new topics in a response.

Best of
Generate multiple responses

Anatomy of Azure OpenAI Service (Chat)

Chat playground

Deploy to

Import setup Export setup Show panels

Assistant setup

System message Add your data (preview)

Save changes

Specify how the chat should act

Use a template to get started, or just start writing your own system message below. Want some tips? [Learn more](#)

Use a system message template

Select a template

System message

You are an AI assistant that helps people find information.

Examples

Add examples to show the chat what responses you want. It will try to mimic any responses you add here so make sure they match the rules you laid out in the system message.

+ Add an example

Chat session

Clear chat View code Show raw JSON

Tell me a Joke.
What's the capitol of France?
Can you suggest a good movie?

Sure, here's a joke for you: Why don't scientists trust atoms? Because they make up everything.

The capital of France is Paris.

As for a movie suggestion, it depends on what kind of movie you're in the mood for. If you're looking for a comedy, I recommend "The Grand Budapest Hotel." If you're in the mood for a drama, "The Shawshank Redemption" is a classic choice. If you want something more recent, "Parasite" won Best Picture at the 2020 Oscars and is definitely worth a watch.

User message

Type user query here. (Shift + Enter for new line)

Configuration

Deployment Parameters

Max response 800

Temperature 0.7

Top P 0.95

Stop sequence

Stop sequences

Frequency penalty 0

Presence penalty 0

[Learn more](#)

Current token count

Input tokens progress indicator

185/4000

Token Count

132/4000

The API supports a maximum of 4000 tokens shared:

- Prompt
- system message
- Examples
- message history
- user query
- model's response

Microsoft Cloud Runs on Trust



- Your data is **your** data
 - Data is stored encrypted in your Azure subscription
- Your data is **not** used to train the foundational AI models
 - Azure OpenAI Service provisioned in your Azure subscription
 - Model fine tuning stays in your Azure subscription and never moves into the foundation AI models
- Your data is **protected** by the most comprehensive enterprise compliance and security controls
 - Encrypted with Customer Managed Keys
 - Private virtual networks and role-based access control
 - Soc2, ISO, HIPPA, CSA STAR compliant

Microsoft Responsible AI

[Responsible AI principles from Microsoft](#)



Fairness

AI systems should **treat people fairly**



Reliability & Safety

AI systems should **perform reliably** and safely



Privacy & Security

AI systems should be secure and **respect privacy**



Inclusiveness

AI systems should **empower everyone** and engage people



Transparency

AI systems should **be understandable**



Accountability

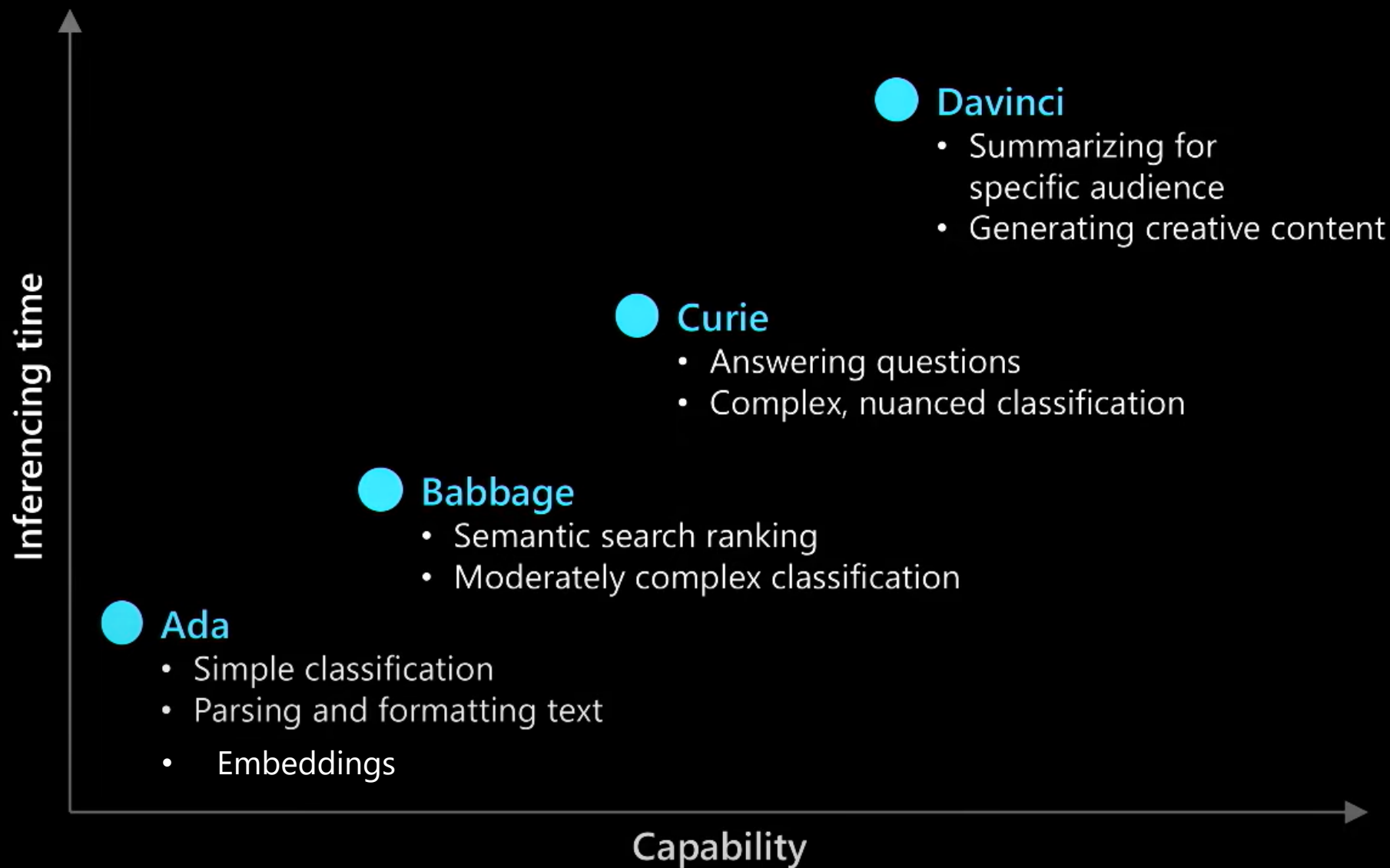
People should **be accountable** for AI systems



Comparison of GPT versions

- **GPT-3.5** (4,000 Tokens)
 - Use-case specific models to optimize inference time and performance
 - Suitable for a large range of use cases
- **ChatGPT**
 - Should be first choice for most use cases
 - Most economical GPT model in Azure OpenAI Service
 - For all workloads, not just chat
- **GPT-4** (32,000 Tokens)
 - Improved problem solving and reasoning capabilities by 10x
 - Iterative refinement:
 - Paste in code errors & GPT-4 will fix for you
 - Iterate on stories
 - Increased token limit – works well for long content

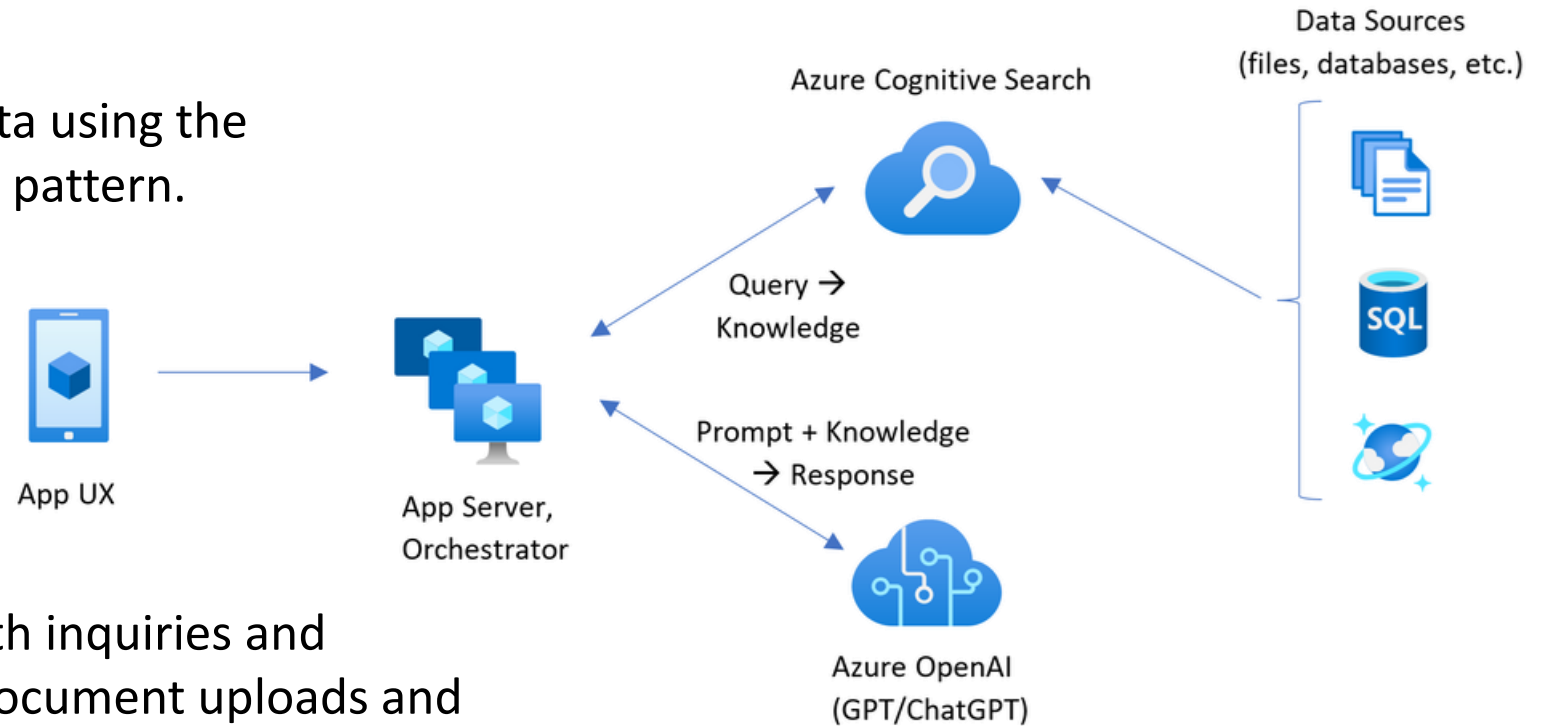
GPT-3 models



Example: Retrieval Augmented Generation pattern

Uses Azure OpenAI Service to access the ChatGPT model (gpt-35-turbo), and Azure Cognitive Search for data indexing and retrieval.

ChatGPT-like experiences over data using the **Retrieval Augmented Generation** pattern.

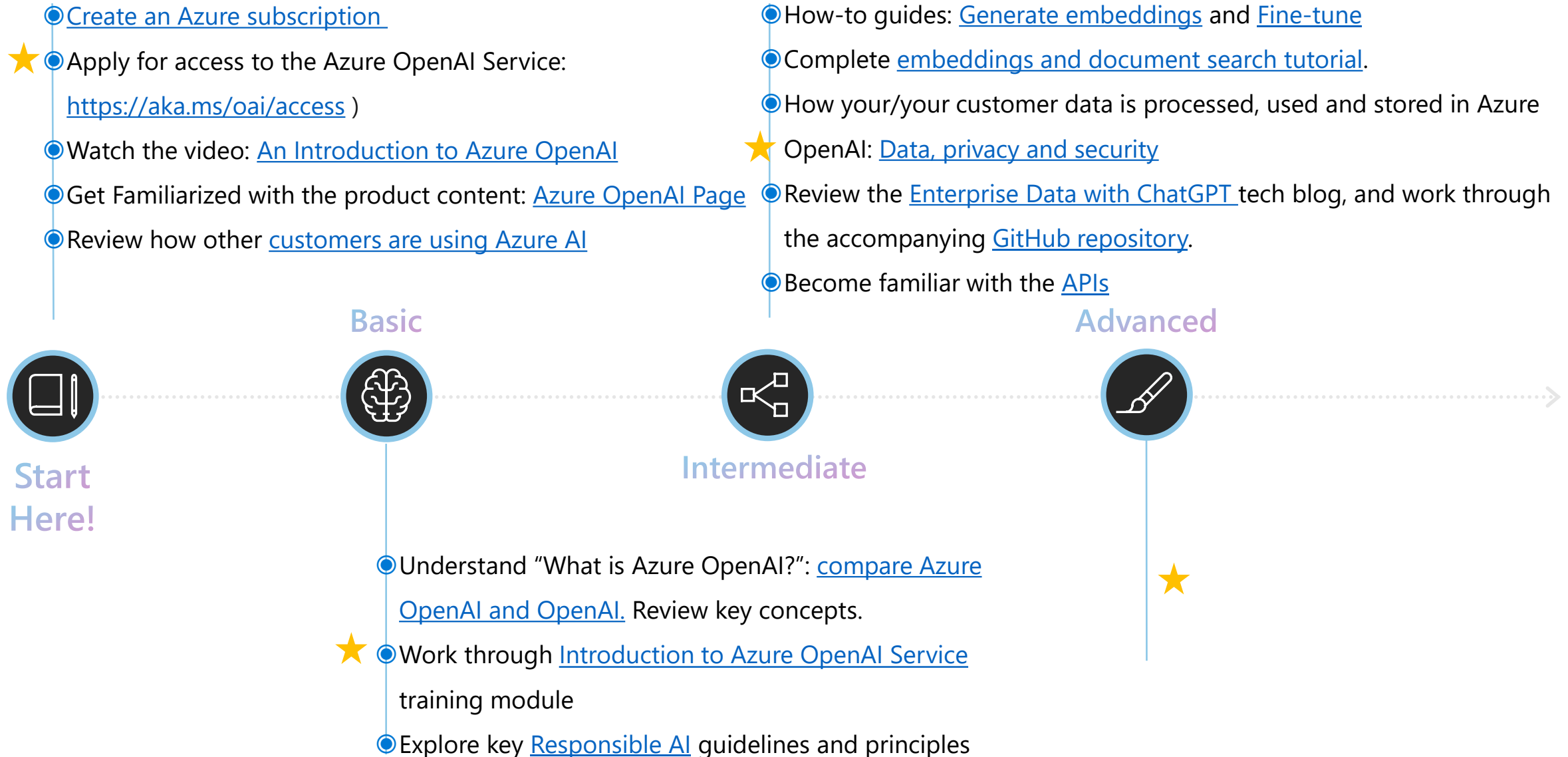


Example

Applications to empower staff with inquiries and patients with interactive health document uploads and simplification.

<https://github.com/Azure-Samples/azure-search-openai-demo>

Learning: Azure OpenAI Service





Demos

<https://github.com/microsoft/globalopenaihack>

SemanticKernel