

# ▶ Running AI Model Locally & Microsoft Sematic Kernel for C#

Carlos Bolivar  
*Principal Software Engineer*

# WHO AM I?

I work as a Principal Software Engineer at DASI (Aircraft Company) and have more than 18 years of experience in the software industry, primarily focused on Microsoft dotnet products.



# ► AGENDA

**01**

**BASIC CONCEPTS**

**02**

**RUN AI MODEL LOCALLY**

**03**

**SEMAC KERNEL**

**04**

**DEMO**

**05**

**Q&A**

# ► BASIC CONCEPTS



## What is a LLM?

A large language model (LLM) is a type of machine learning model designed for natural language processing tasks such as language generation.



## Key Features

- Natural Language Understanding (NLU)
- Text Generation
- Code Writing
- Summarization
- Conversational Abilities
- Question Answering
- Data extractions
- Text Search (RAG)



## Major LLM

- ChatGPT (OpenAI)
- DeepSeek R1
- Gemini (Google)
- Llama (Meta AI)
- Mistral (Mistral AI SAS)
- Claude (Anthropic)
- Phi-3 (Microsoft)
- Qwen (Alibaba)

# ► CLOUD-BASED PLATFORMS

## Multiples Models Deploy

<b>Azure AI Foundry</b>	<a href="https://ai.azure.com">https://ai.azure.com</a>
<b>Github Marketplace</b>	<a href="https://github.com/marketplace">https://github.com/marketplace</a>
<b>Together AI</b>	<a href="https://www.together.ai">https://www.together.ai</a>
<b>Hugging Face</b>	<a href="https://huggingface.co">https://huggingface.co</a>
<b>Amazon Bedrock</b>	<a href="https://aws.amazon.com/bedrock">https://aws.amazon.com/bedrock</a>

## Single Model Deploy

<b>OpenAI</b>	<a href="https://platform.openai.com/docs/overview">https://platform.openai.com/docs/overview</a>
<b>Google AI Studio</b>	<a href="https://aistudio.google.com">https://aistudio.google.com</a>
<b>Deepseek</b>	<a href="https://platform.deepseek.com">https://platform.deepseek.com</a>
<b>Llama</b>	<a href="https://www.llama.com">https://www.llama.com</a>
<b>Mistral</b>	<a href="https://mistral.ai">https://mistral.ai</a>

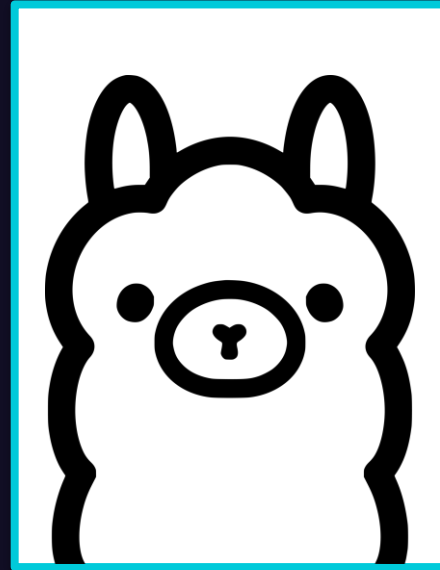
# ► ON-PERMISES & SELF-HOSTED PLATFORMS



## Ollama

Ollama is an open-source command-line tool that lets developers run large language models (LLMs) locally.

It provides a lightweight framework for downloading, running, and managing LLMs directly on your machine, with minimal setup.



<https://ollama.com/>

# ▶ RUN MODEL LOCALLY WITH OLLAMA

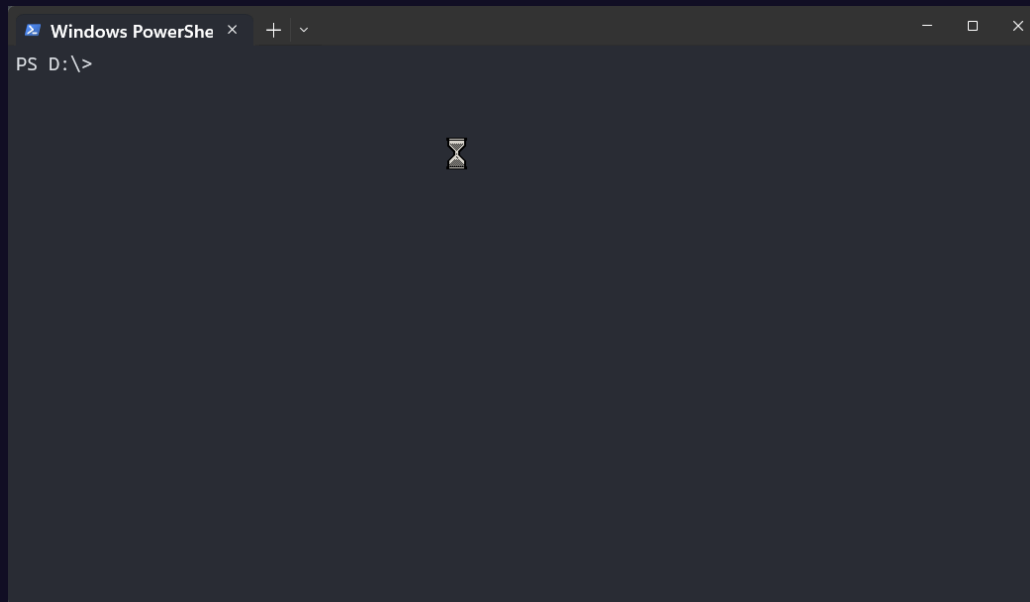


## MODEL LIBRARY

Model
Gemma 3
Claude
DeepSeek-R1
Llama 3
Phi 4
Mistral
.. And more



## RUNNING A MODEL LOCALLY

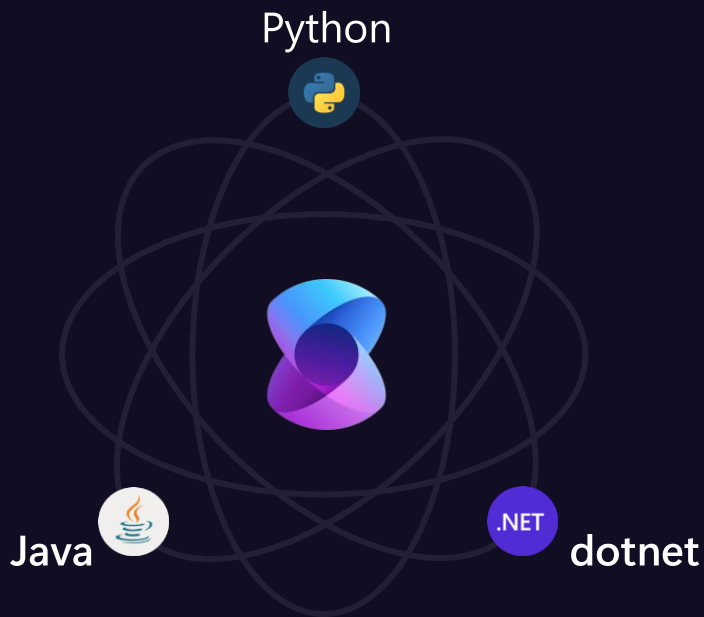


# ► MICROSOFT SEMATIC KERNEL



## WHAT IS IT?

An open-source SDK that lets you easily build agents that call your existing code. As a highly extensible framework, it can be used with models from OpenAI, Azure OpenAI, Hugging Face, and more!



<https://learn.microsoft.com/en-us/semantic-kernel/overview/>



# ► AI SERVICES OF SEMATIC KERNEL

Services	C#	Python	Java
Chat Completion	✓	✓	✓
Text Generation	✓	✓	✓
Embedding Generation (experimental)	✓	✓	✓
Text-to-image (experimental)	✓	✓	✗
Image-to-text (experimental)	✓	✗	✗
Text-to-audio	✓	✓	✗
Audio-to-text	✓	✓	✗



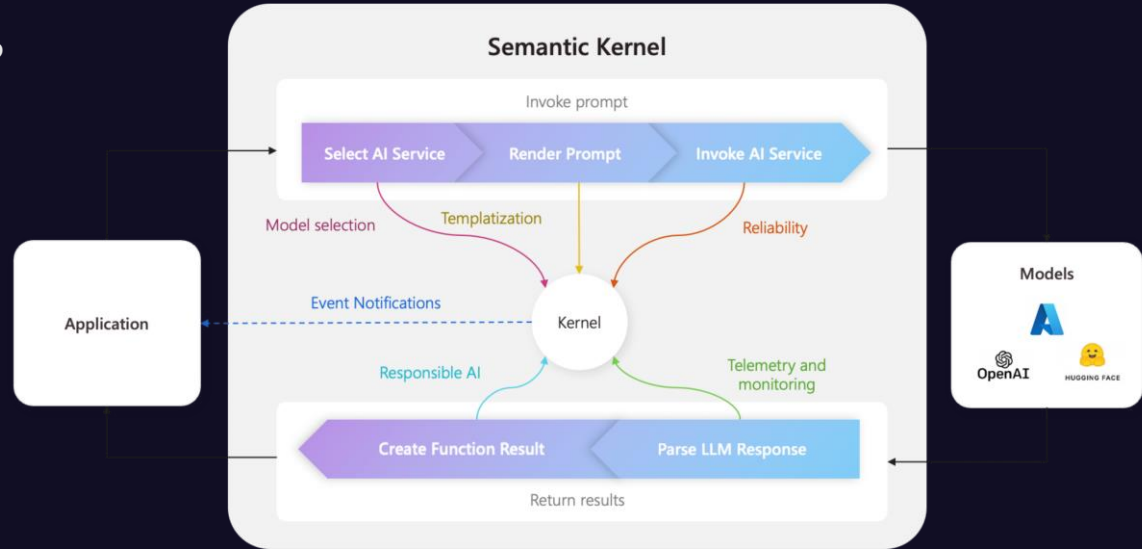
# ► SEMATIC KERNEL – THE KERNEL



## WHAT IS THE KERNEL?

The kernel is the central component of Semantic Kernel. It has all of the services and plugins necessary to run both native code and AI services, it is used by nearly every component within the Semantic Kernel SDK to power your agent.

*"All from a single place..."*



# ► SEMATIC KERNEL - COMPONENTS



## SERVICES

These include both AI services and additional services essential for operating your application.

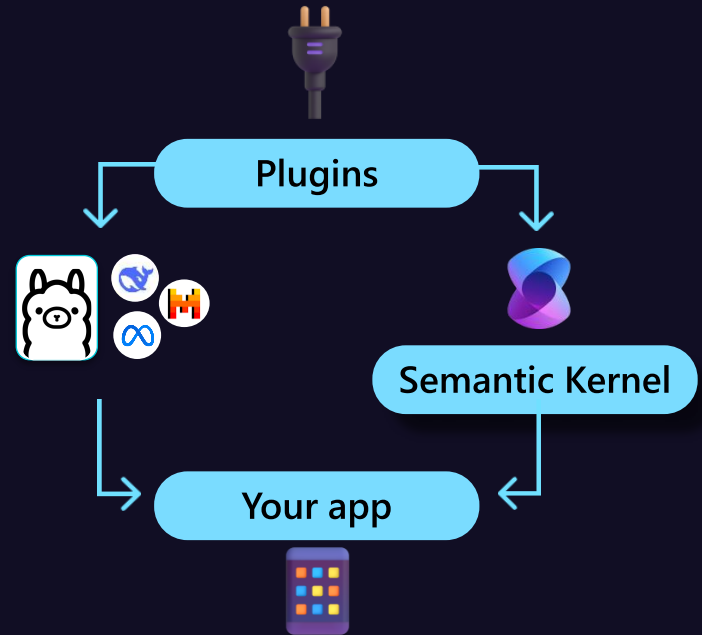
This was designed based on the Service Provider pattern used in .NET to facilitate dependency injection across all languages.



## PLUGINS

These are the components that are used by your AI services and prompt templates to perform work.

For example, can use plugins to retrieve data from a database or call an external API to perform actions.





► **DEMO**

# ► THANKS!

Do you have any questions?



<https://www.linkedin.com/in/cbolivar82/en>



<https://github.com/cbolivar82/ollama-sematic-kernel>



# SCAN ME!