

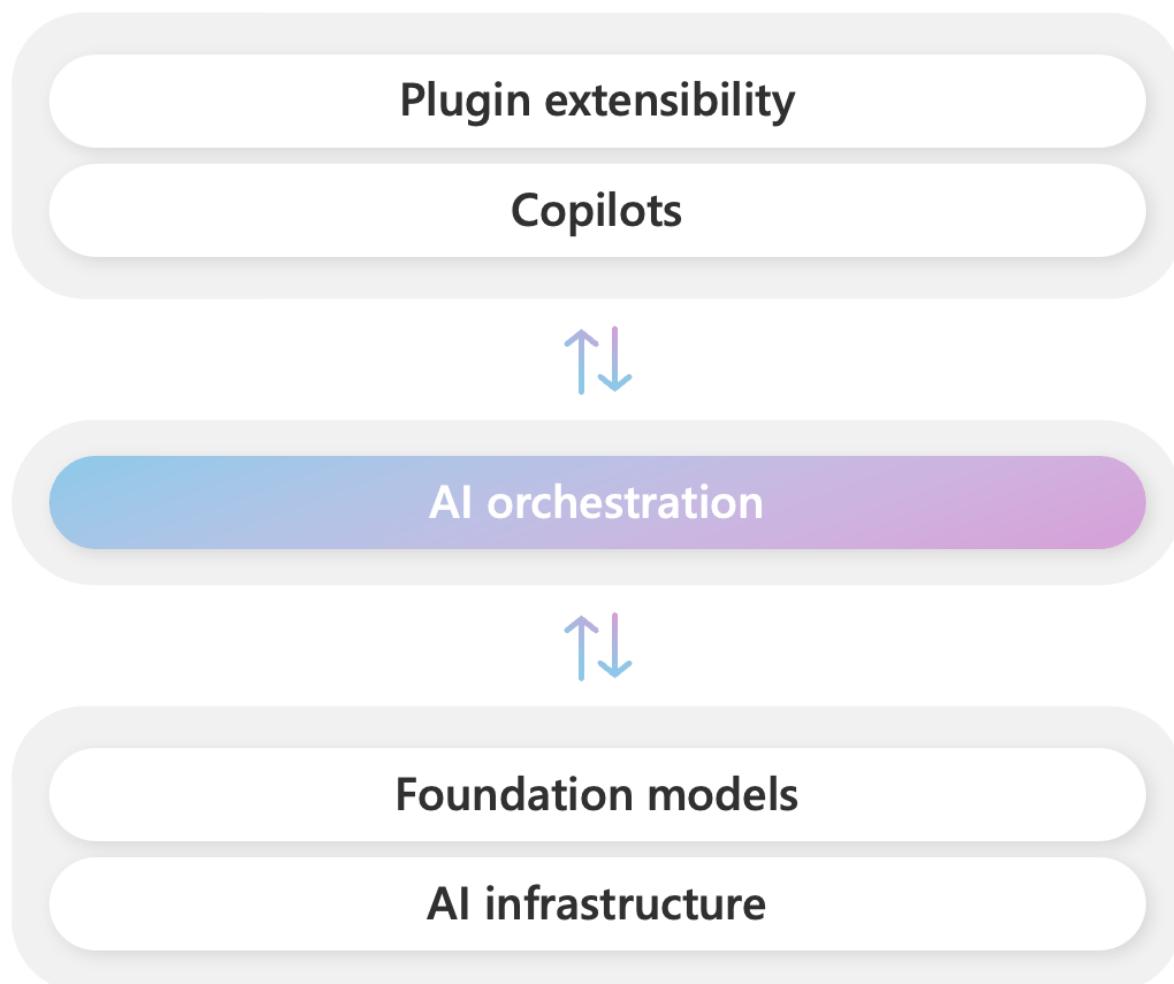
What is Semantic Kernel?

Article • 07/11/2023



Semantic Kernel is an open-source SDK that lets you easily combine AI services like [OpenAI](#), [Azure OpenAI](#), and [Hugging Face](#) with conventional programming languages like C# and Python. By doing so, you can create AI apps that combine the best of both worlds.

Semantic Kernel is at the center of the copilot stack



During Kevin Scott's talk [The era of the AI Copilot](#), he showed how Microsoft powers its [Copilot system](#) with a stack of AI models and plugins. At the center of this stack is an AI orchestration layer that allows us to combine AI models and plugins together to create brand new experiences for users.

To help developers build their *own* Copilot experiences on top of AI plugins, we have released Semantic Kernel, a lightweight open-source SDK that allows you to orchestrate AI plugins. With Semantic Kernel, you can leverage the *same* AI orchestration patterns that power Microsoft 365 Copilot and Bing in your *own* apps, while still leveraging your existing development skills and investments.

💡 Tip

If you are interested in seeing a sample of the copilot stack in action (with Semantic Kernel at the center of it), check out [Project Miyagi](#). Project Miyagi reimagines the design, development, and deployment of intelligent applications on top of Azure with all of the latest AI services and tools.

Semantic Kernel makes AI development extensible

Semantic Kernel has been engineered to allow developers to flexibly integrate AI services into their existing apps. To do so, Semantic Kernel provides a set of connectors that make it easy to add [memories](#) and [models](#). In this way, Semantic Kernel is able to add a simulated "brain" to your app.

Additionally, Semantic Kernel makes it easy to add skills to your applications with [AI plugins](#) that allow you to interact with the real world. These plugins are composed of [prompts](#) and [native functions](#) that can respond to triggers and perform actions. In this way, plugins are like the "body" of your AI app.

Because of the extensibility Semantic Kernel provides with connectors and [plugins](#), you can use it to orchestrate AI plugins from both OpenAI and Microsoft on top of nearly any model. For example, you can use Semantic Kernel to orchestrate plugins built for ChatGPT, Bing, and Microsoft 365 Copilot on top of models from OpenAI, Azure, or even Hugging Face.



Models and Memory

Connectors

Semantic Kernel

Plugins



Triggers and actions

As a developer, you can use these pieces individually or together. For example, if you just need an abstraction over OpenAI and Azure OpenAI services, you could use the SDK to *just* run pre-configured prompts within your plugins, but the *real* power of Semantic Kernel comes from combining these components together.

Why do you need an AI orchestration SDK?

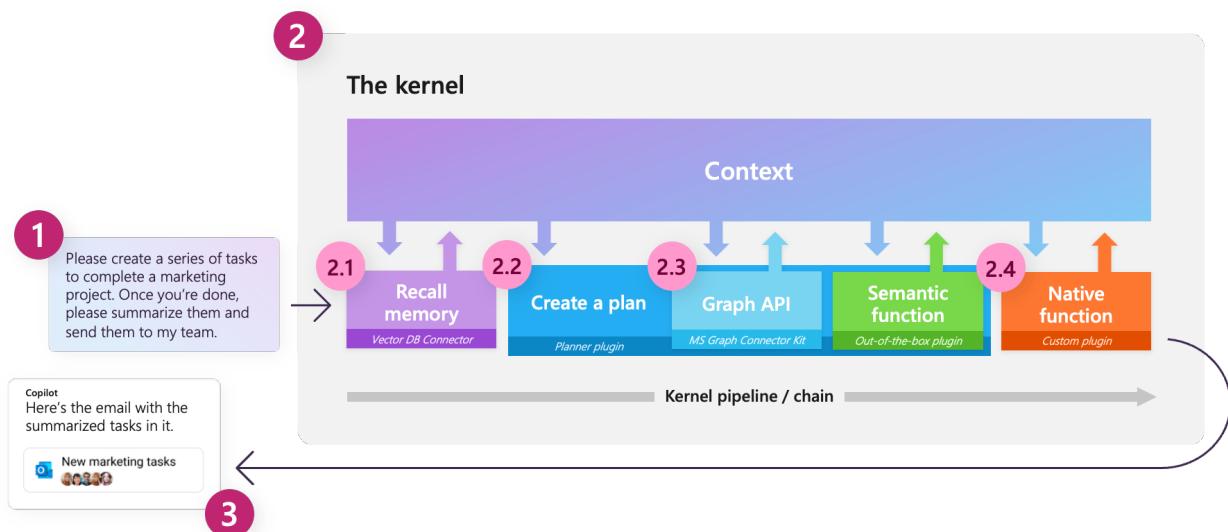
If you wanted, you could use the APIs for popular AI services directly and directly feed the results into your existing apps and services. This, however, requires you to learn the APIs for each service and then integrate them into your app. Using the APIs directly also does not allow you to easily draw from the recent advances in AI research that require solutions *on top* of these services. For example, the existing APIs do not provide planning or AI memories out-of-the-box.

To simplify the creation of AI apps, open source projects like [LangChain](#) have emerged. Semantic Kernel is Microsoft's contribution to this space and is designed to support enterprise app developers who want to integrate AI into their existing apps.

Seeing AI orchestration with Semantic Kernel

By using multiple AI models, plugins, and memory all together within Semantic Kernel, you can create sophisticated pipelines that allow AI to automate complex tasks for users.

For example, with Semantic Kernel, you could create a pipeline that helps a user send an email to their marketing team. With [memory](#), you could retrieve information about the project and then use [planner](#) to autogenerate the remaining steps using available plugins (e.g., ground the user's ask with Microsoft Graph data, generate a response with GPT-4, and send the email). Finally, you can display a success message back to your user in your app using a custom plugin.



Step	Component	Description
1	Ask	It starts with a goal being sent to Semantic Kernel by either a user or developer.
2	Kernel	The kernel orchestrates a user's ask. To do so, the kernel runs a pipeline / chain that is defined by a developer. While the chain is run, a common context is provided by the kernel so data can be shared between functions.
2.1	Memories	With a specialized plugin, a developer can recall and store context in vector databases. This allows developers to simulate memory within their AI apps.
2.2	Planner	Developers can ask Semantic Kernel to auto create chains to address novel needs for a user. Planner achieves this by mixing-and-matching plugins

Step	Component	Description
		that have already been loaded into the kernel to create additional steps. This is similar to how ChatGPT, Bing, and Microsoft 365 Copilot combines plugins together in their experiences.
2.3	Connectors	To get additional data or to perform autonomous actions, you can use out-of-the-box plugins like the Microsoft Graph Connector kit or create a custom connector to provide data to your own services.
2.4	Custom plugins	As a developer, you can create custom plugins that run inside of Semantic Kernel. These plugins can consist of either LLM prompts (semantic functions) or native C# or Python code (native function). This allows you to add new AI capabilities and integrate your existing apps and services into Semantic Kernel.
3	Response	Once the kernel is done, you can send the response back to the user to let them know the process is complete.

Semantic Kernel is open-source

To make sure all developers can take advantage of our learnings building Copilots, we have released Semantic Kernel as an [open-source project](#) on GitHub. Today, we provide the SDK in .NET and Python flavors (TypeScript and Java are coming soon). For a full list of what is supported in each language, see [supported languages](#).

The screenshot shows the GitHub repository page for `microsoft/semantic-kernel`. The repository is public, has 134 watchers, 1.1k forks, and 8.9k stars. It contains 18 branches and 9 tags. The main branch has 633 commits. The repository description is: "Integrate cutting-edge LLM technology quickly and easily into your apps". It uses the `aka.ms/semantic-kernel` URL. Key features listed include `sdk`, `ai`, `artificial-intelligence`, `openai`, and `llm`. The repository has a MIT license, a Code of conduct, a Security policy, and 134 people watching it. There are 1.1k forks. A "Report repository" link is also present. The releases section shows one release, `python-0.2.7.dev`, which is the latest version, released 2 weeks ago. The packages section indicates "No packages published".

Given that new breakthroughs in LLM AIs are landing on a daily basis, you should expect this SDK evolve. We're excited to see what you build with Semantic Kernel and we look forward to your feedback and contributions so we can build the best practices together in the SDK.

[Open the Semantic Kernel repo](#)

Watch the latest Semantic Kernel roadmap video

<https://learn-video.azurefd.net/vod/player?id=b60199d8-1fa4-41b3-a9e0-68a068681b2a&locale=en-us&embedUrl=%2Fsemantic-kernel%2Foverview%2F>

Contribute to Semantic Kernel

We welcome contributions and suggestions from the Semantic Kernel community! One of the easiest ways to participate is to engage in discussions in the [GitHub repository](#). Bug reports and fixes are welcome!

For new features, components, or extensions, please [open an issue](#) and discuss with us before sending a PR. This will help avoid rejections since it will allow us to discuss the impact to the larger ecosystem.

[Learn more about contributing](#)

Get started using the Semantic Kernel SDK

Now that you know what Semantic Kernel is, follow the [get started](#) link to try it out. Within minutes you can create prompts and chain them with out-of-the-box plugins and native code. Soon afterwards, you can give your apps memories with embeddings and summon even more power from external APIs.

[Get started with Semantic Kernel](#)

Start learning how to use Semantic Kernel

Article • 07/11/2023



In just a few steps, you can start running the getting started guides for Semantic Kernel in either C# or Python. After completing the guides, you'll know how to...

- Configure your local machine to run Semantic Kernel
- Run AI prompts from the kernel
- Make AI prompts dynamic with variables
- Create a simple AI chatbot
- Automatically combine functions together with planner
- Store and retrieve memory with embeddings

If you are an experienced developer, you can skip the guides and directly access the packages from the Nuget feed or PyPI.

C#

Instructions for accessing the `SemanticKernel` Nuget feed is available [here ↗](#). It's as easy as:

Nuget

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

Requirements to run the guides

Before running the guides in C#, make sure you have the following installed on your local machine.

- ✓ `git` or the [GitHub app ↗](#)
- ✓ [VSCode ↗](#) or [Visual Studio ↗](#)
- ✓ An OpenAI key via either [Azure OpenAI Service](#) or [OpenAI ↗](#)
- ✓ [.Net 7 SDK ↗](#) - for C# notebook guides
- ✓ In VS Code the [Polyglot Notebook ↗](#) - for notebook guides

If you are using the Python guides, you just need `git` and `python`. These guides have been tested on python versions 3.8-3.11.

Download and run the guides

To setup the guides, follow the steps below.

💡 Tip

Have your OpenAI or Azure OpenAI keys ready to enter when prompted by the Jupyter notebook.

1. Use your web browser to visit aka.ms/sk/repo on GitHub.
2. Clone or fork the repo to your local machine.

❗ Note

If you are new to using GitHub and have never cloned a repo to your local machine, please review [this guide](#).

❗ Note

If you are a new contributor to open source, please [fork the repo](#) to start your journey.

If you have trouble cloning or forking the repo, you can watch the video below.
<https://aka.ms/SK-Local-Setup>

3. While the repository is open in VS Code, navigate to the `/samples/notebooks` folder.
4. Choose either the `dotnet` or `python` folder based on your preferred programming language.
5. Open the `00-getting-started.ipynb` notebook.
6. Activate each code snippet with the "play" button on the left hand side.

If you need help running the `00-getting-started.ipynb` notebook, you can watch the video below.

<https://aka.ms/SK-Getting-Started-Notebook>

7. Repeat for the remaining notebooks.

Navigating the guides

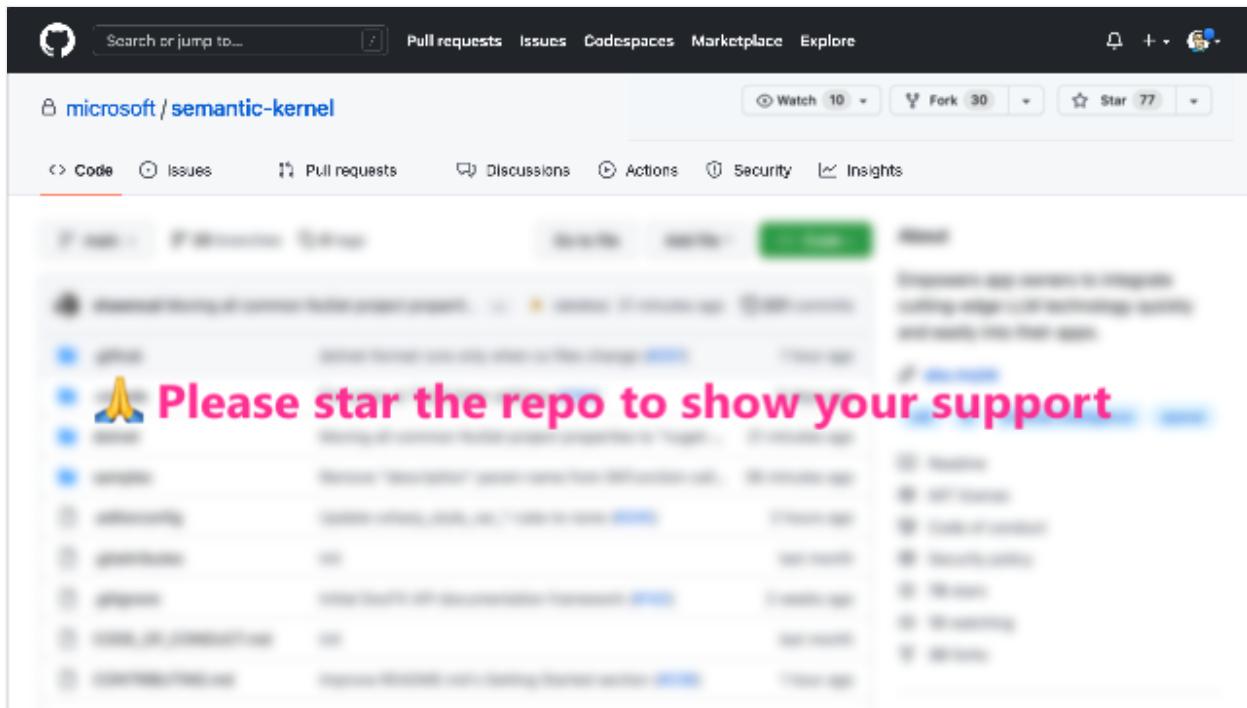
The guides are designed to be run in order to build on the concepts learned in the previous notebook. If you are interested in learning a particular concept, however, you can jump to the notebook that covers that concept. Below are the available guides; each one can also be opened within the docs website by clicking on the **Open guide** link.

File	Link	Description
<i>00-getting-started.ipynb</i>	Open guide	Run your first prompt
<i>01-basic-loading-the-kernel.ipynb</i>	Open guide	Changing the configuration of the kernel
<i>02-running-prompts-from-file.ipynb</i>	Open guide	Learn how to run prompts from a file
<i>03-semantic-function-inline.ipynb</i>	Open guide	Configure and run prompts directly in code
<i>04-context-variables-chat.ipynb</i>	Open guide	Use variables to make prompts dynamic
<i>05-using-the-planner.ipynb</i>	Open guide	Dynamically create prompt chains with planner
<i>06-memory-and-embeddings.ipynb</i>	Open guide	Store and retrieve memory with embeddings

[Start the first guide](#)

Like what you see?

If you are a fan of Semantic Kernel, please give the repo a  star to show your support.



Keep learning

The guides are an easy way run sample code and learn how to use Semantic Kernel. If you want to learn more about the concepts behind Semantic Kernel, keep reading the docs. Based on your experience level, you can jump to the section that best fits your needs.

Experience level	Next step
For beginners who are just starting to learn about AI	Learn prompt engineering
For people who are well versed in prompt engineering	Orchestrate AI plugins
For people familiar with using AI plugins	Store and retrieve memory
For those who want to see how it all works together	Run the sample apps

[Learn how to Orchestrate AI](#)

Hackathon materials for Semantic Kernel

Article • 07/11/2023



With these materials you can run your own Semantic Kernel Hackathon, a hands-on event where you can learn and create AI solutions using Semantic Kernel tools and resources.

By participating and running a Semantic Kernel hackathon, you will have the opportunity to:

- Explore the features and capabilities of Semantic Kernel and how it can help you solve problems with AI
- Work in teams to brainstorm and develop your own AI plugins or apps using Semantic Kernel SDK and services
- Present your results and get feedback from other participants
- Have fun!

Download the materials

To run your own hackathon, you will first need to download the materials. You can download the zip file here:

[Download hackathon materials](#)

Once you have unzipped the file, you will find the following resources:

- Hackathon sample agenda
- Hackathon prerequisites
- Hackathon facilitator presentation
- Hackathon team template
- Helpful links

Name	Date Modified	Size	Kind
Hackathon facilitator presentation	Today at 9:41 PM	--	Folder
Semantic Kernel - Hackathon Presentation.pptx	Apr 6, 2023 at 5:42 PM	12.6 MB	PowerP...(pptx)
Hackathon prerequisites	Today at 9:41 PM	--	Folder
Semantic Kernel - Hackathon Prerequisites.docx	Apr 6, 2023 at 5:42 PM	70 KB	Micros...(docx)
Semantic Kernel - Hackathon Prerequisites.pdf	Apr 6, 2023 at 5:42 PM	146 KB	PDF Document
Hackathon sample agenda	Today at 9:41 PM	--	Folder
Semantic Kernel - Hackathon Agenda.docx	Apr 6, 2023 at 5:42 PM	71 KB	Micros...(docx)
Semantic Kernel - Hackathon Agenda.pdf	Apr 6, 2023 at 5:42 PM	137 KB	PDF Document
Hackathon team template	Today at 9:41 PM	--	Folder
Semantic Kernel - Hackathon Team Template.pptx	Apr 6, 2023 at 5:42 PM	2.8 MB	PowerP...(pptx)
Participant Exit Survey Link.url	Apr 6, 2023 at 5:42 PM	121 bytes	Web sit...ocation
Semantic Kernel Discord Community Link.url	Apr 6, 2023 at 5:42 PM	115 bytes	Web sit...ocation
Semantic Kernel Documentation Site.url	Apr 6, 2023 at 5:42 PM	113 bytes	Web sit...ocation
Semantic Kernel Recipes.url	Apr 6, 2023 at 5:42 PM	145 bytes	Web sit...ocation
Semantic Kernel Repo.url	Apr 6, 2023 at 5:42 PM	112 bytes	Web sit...ocation

Preparing for the hackathon

Before the hackathon, you and your peers will need to download and install software needed for Semantic Kernel to run. Additionally, you should already have API keys for either OpenAI or Azure OpenAI and access to the Semantic Kernel repo. Please refer to the prerequisites document in the facilitator materials for the complete list of tasks participants should complete before the hackathon.

You should also familiarize yourself with the available documentation and tutorials. This will ensure that you are knowledgeable of core Semantic Kernel concepts and features so that you can help others during the hackathon. The following resources are highly recommended:

- [What is Semantic Kernel?](#)
- [Semantic Kernel recipes](#)
- [Semantic kernel recipes videos ↗](#)
- [Semantic Kernel LinkedIn training video ↗](#)

Additionally, you can check out the available sample AI plugins and apps that demonstrate how Semantic Kernel can be used for various scenarios. You can use them as inspiration and even modify them for your own projects during the hackathon. You can find them here:

- [Semantic Kernel samples](#)

Running the hackathon

The hackathon will consist of six main phases: welcome, overview, brainstorming, development, presentation, and feedback.

Here is an approximate agenda and structure for each phase but feel free to modify this based on your team:

Length (Minutes)	Phase	Description
15	Welcome	The hackathon facilitator will welcome the participants, introduce the goals and rules of the hackathon, and answer any questions.
60	Overview	The facilitator will guide you through either a pre-recorded video or a live presentation that will give you an overview of AI and why it is important for solving problems in today's world. Along with an overview of Semantic Kernel and its features, such as the kernel , planner , plugins , memories and more. You will also see demos of how Semantic Kernel can be used for different scenarios.
120	Brainstorming	The facilitator will help you form teams based on your interests or skill levels. You will then brainstorm ideas for your own AI plugins or apps using design thinking techniques.
360+	Development	You will use Semantic Kernel SDKs tools, and resources to develop, test, and deploy your projects. This could be for the rest of the day or over multiple days based on the time available and problem to be solved.
60	Presentation	Each team will present their results using a PowerPoint template provided. You will have about 15 minutes per team to showcase your project, demonstrate how it works, and explain how it solves a problem with AI. You will also receive feedback from other participants.
30	Feedback	Each team can share their feedback on the hackathon and Semantic Kernel with the group and fill out the Hackathon Exit Survey ↗ .

Following up after the hackathon

We hope you enjoyed running a Semantic Kernel Hackathon and the overall experience! We would love to hear from you about what worked well, what didn't, and what we can improve for future content. Please take a few minutes to fill out the [hackathon facilitator survey ↗](#) and share your feedback and suggestions with us.

If you want to continue developing your AI plugins or projects after the hackathon, you can find more resources and support for Semantic Kernel.

- [Semantic Kernel Discord community ↗](#)
- [Semantic Kernel blog ↗](#)
- [Semantic Kernel GitHub repo ↗](#)

Thank you for your engagement and creativity during the hackathon. We look forward to seeing what you create next with Semantic Kernel!

Additional learning for Semantic Kernel

Article • 07/11/2023



Want to learn more about Semantic Kernel? Check out these in-depth tutorials and videos. We will add more content over time from our team and community, so check back often!

Cook with Semantic Kernel

Learn how to supercharge your problem-solving creativity with Semantic Kernel running on your own machine just like your own “Easy Bake Oven.” We’ll use plenty of cooking analogies to land the core ideas of LLM AI running on Semantic Kernel so be prepared to get hungry!

[Start the tutorial](#)

The screenshot shows a Microsoft Teams video call interface. On the right, a man wearing glasses and a black shirt is speaking. On the left, a Microsoft Teams video camera icon is visible. The main screen displays a Microsoft Visual Studio Code (VS Code) interface. The title bar says "notebook.ipynb M". The left sidebar shows a file tree for "e2-first-dish" containing files like "SK-RECIPES", "config", "Settings.cs", and various "e1" through "e8" recipe files. The main code editor area contains the following text:

```
#!import ../config/Settings.cs
bool useAzureOpenAI = false;

await Settings.AskAzureEndpoint(useAzureOpenAI);
await Settings.AskModel(useAzureOpenAI);
await Settings.AskApiKey(useAzureOpenAI);
```

Kernel syntax examples

This project contains a collection of .NET examples for various scenarios using Semantic Kernel components. There are already 40 examples that show how to achieve basic tasks like creating a chain, adding a plugin, and running a chain. There are also more advanced examples that show how to use the Semantic Kernel API to create a custom plugin, stream data, and more.

[Start the tutorial](#)

-  Example01_NativeFunctions.cs
-  Example02_Pipeline.cs
-  Example03_Variables.cs
-  Example04_CombineLLMPromptsAndNativeCode.cs
-  Example05_InlineFunctionDefinition.cs
-  Example06_TemplateLanguage.cs
-  Example07_BingAndGoogleSkills.cs
-  Example08_RetryHandler.cs

More tutorials coming soon

If you have a tutorial you would like to share, please let us know on our [Discord server](#). We would love to highlight your content with the community here!

Supported Semantic Kernel languages

Article • 07/18/2023



ⓘ Note

Skills are currently being renamed to plugins. This article has been updated to reflect the latest terminology, but some images and code samples may still refer to skills.

Semantic Kernel plans on providing support to the following languages:

- ✓ C#
- ✓ Python
- ✓ Java ([coming soon ↗](#))

While the overall architecture of the kernel is consistent across all languages, we made sure the SDK for each language follows common paradigms and styles in each language to make it feel native and easy to use.

Available features

Today, not all features are available in all languages. The following tables show which features are available in each language. The  symbol indicates that the feature is partially implemented, please see the associated note column for more details. The  symbol indicates that the feature is not yet available in that language; if you would like to see a feature implemented in a language, please consider [contributing to the project](#) or [opening an issue](#).

AI Services

Services	C#	Python	Java	Notes
TextGeneration				Example: Text-Davinci-003
TextEmbeddings				Example: Text-Embeddings-Ada-002
ChatCompletion				Example: GPT4, Chat-GPT

Services	C#	Python	Java	Notes
Image Generation	✓	✗	✗	Example: Dall-E

AI service endpoints

Endpoints	C#	Python	Java	Notes
OpenAI	✓	✓	✓	
AzureOpenAI	✓	✓	✓	
Hugging Face Inference API	⌚	✗	✗	Coming soon to Python, not all scenarios are covered for .NET
Hugging Face Local	✗	✓	✗	
Custom	✓	⌚	✗	Requires the user to define the service schema in their application

Tokenizers

Tokenizers	C#	Python	Java	Notes
GPT2	✓	✓	✓	
GPT3	✓	✗	✗	
tiktoken	⌚	✗	✗	Coming soon to Python and C#. Can be manually added to Python via <code>pip install tiktoken</code>

Core plugins

Plugins	C#	Python	Java	Notes
TextMemorySkill	✓	✓	⌚	
ConversationSummarySkill	✓	✓	✗	
FileIOSkill	✓	✓	✓	
HttpSkill	✓	✓	✓	
MathSkill	✓	✓	✓	
TextSkill	✓	✓	⌚	

Plugins	C#	Python	Java	Notes
TimeSkill	✓	✓	✓	
WaitSkill	✓	✓	✓	

Planners

Planners	C#	Python	Java	Notes
Plan Object Model	✓	✓	✗	
BasicPlanner	✗	✓	✗	
ActionPlanner	✓	⌚	⌚	In development
SequentialPlanner	✓	⌚	⌚	In development
StepwisePlanner	✓	✗	✗	

Connectors

Memory Connectors	C#	Python	Java	Notes
Azure Cognitive Search	✓	✓	✗	
Chroma	✓	✓	✗	
Pinecone	✓	✓	✗	
Weaviate	✓	✓	✗	Currently supported on Python 3.9+, 3.8 coming soon
Qdrant	✓	⌚	✗	In feature branch for review
Postgres	✓	✓	✗	Vector optimization requires pgvector ↗
Sqlite	✓	✗	✗	Vector optimization requires sqlite-vss ↗
Milvus	⌚	⌚	✗	

Plugins

Plugins	C#	Python	Java	Notes
MsGraph	✓	✗	✗	Contains plugins for OneDrive, Outlook, ToDos, and Organization Hierarchies
Document and data loading plugins (i.e. pdf, csv, docx, pptx)	✓	✗	✗	Currently only supports Word documents
OpenAPI	✓	✗	✗	
Web search plugins (i.e. Bing, Google)	✓	✓	✗	
Text chunkers	⟳	⟳	✗	

Notes about the Python SDK

During the initial development phase, many Python best practices have been ignored in the interest of velocity and feature parity. The project is now going through a refactoring exercise to increase code quality.

To make the Kernel as lightweight as possible, the core pip package should have a minimal set of external dependencies. On the other hand, the SDK should not reinvent mature solutions already available, unless of major concerns.

Contributing to Semantic Kernel

Article • 06/21/2023

You can contribute to Semantic Kernel by submitting issues, starting discussions, and submitting pull requests (PRs). Contributing code is greatly appreciated, but simply filing issues for problems you encounter is also a great way to contribute since it helps us focus our efforts.

Reporting issues and feedback

We always welcome bug reports, API proposals, and overall feedback. Since we use GitHub, you can use the [Issues ↗](#) and [Discussions ↗](#) tabs to start a conversation with the team. Below are a few tips when submitting issues and feedback so we can respond to your feedback as quickly as possible.

Reporting issues

New issues for the SDK can be reported in our [list of issues ↗](#), but before you file a new issue, please search the list of issues to make sure it does not already exist. If you have issues with the Semantic Kernel documentation (this site), please file an issue in the [Semantic Kernel documentation repository ↗](#).

If you *do* find an existing issue for what you wanted to report, please include your own feedback in the discussion. We also highly recommend up-voting ( reaction) the original post, as this helps us prioritize popular issues in our backlog.

Writing a Good Bug Report

Good bug reports make it easier for maintainers to verify and root cause the underlying problem. The better a bug report, the faster the problem can be resolved. Ideally, a bug report should contain the following information:

- A high-level description of the problem.
- A *minimal reproduction*, i.e. the smallest size of code/configuration required to reproduce the wrong behavior.
- A description of the *expected behavior*, contrasted with the *actual behavior* observed.
- Information on the environment: OS/distribution, CPU architecture, SDK version, etc.

- Additional information, e.g. Is it a regression from previous versions? Are there any known workarounds?

[Create issue](#)

Submitting feedback

If you have general feedback on Semantic Kernel or ideas on how to make it better, please share it on our [discussions board](#). Before starting a new discussion, please search the list of discussions to make sure it does not already exist.

We recommend using the [ideas category](#) if you have a specific idea you would like to share and the [Q&A category](#) if you have a question about Semantic Kernel.

You can also start discussions (and share any feedback you've created) in the Discord community by joining the [Semantic Kernel Discord server](#).

[Start a discussion](#)

Help us prioritize feedback

We currently use up-votes to help us prioritize issues and features in our backlog, so please up-vote any issues or discussions that you would like to see addressed.

If you think others would benefit from a feature, we also encourage you to ask others to up-vote the issue. This helps us prioritize issues that are impacting the most users. You can ask colleagues, friends, or the [community on Discord](#) to up-vote an issue by sharing the link to the issue or discussion.

Submitting pull requests

We welcome contributions to Semantic Kernel. If you have a bug fix or new feature that you would like to contribute, please follow the steps below to submit a pull request (PR). Afterwards, project maintainers will review code changes and merge them once they've been accepted.

Recommended contribution workflow

We recommend using the following workflow to contribute to Semantic Kernel (this is the same workflow used by the Semantic Kernel team):

1. Create an issue for your work.

- You can skip this step for trivial changes.
- Reuse an existing issue on the topic, if there is one.
- Get agreement from the team and the community that your proposed change is a good one by using the discussion in the issue.
- Clearly state in the issue that you will take on implementation. This allows us to assign the issue to you and ensures that someone else does not accidentally work on it.

2. Create a personal fork of the repository on GitHub (if you don't already have one).

3. In your fork, create a branch off of main (`git checkout -b mybranch`).

- Name the branch so that it clearly communicates your intentions, such as "issue-123" or "githubhandle-issue".

4. Make and commit your changes to your branch.

5. Add new tests corresponding to your change, if applicable.

6. Build the repository with your changes.

- Make sure that the builds are clean.
- Make sure that the tests are all passing, including your new tests.

7. Create a PR against the repository's **main** branch.

- State in the description what issue or improvement your change is addressing.
- Verify that all the Continuous Integration checks are passing.

8. Wait for feedback or approval of your changes from the code maintainers.

9. When area owners have signed off, and all checks are green, your PR will be merged.

Dos and Don'ts while contributing

The following is a list of Dos and Don'ts that we recommend when contributing to Semantic Kernel to help us review and merge your changes as quickly as possible.

Do's:

- **Do** follow the standard [.NET coding style](#) and [Python code style](#) ↗
- **Do** give priority to the current style of the project or file you're changing if it diverges from the general guidelines.

- **Do** include tests when adding new features. When fixing bugs, start with adding a test that highlights how the current behavior is broken.
- **Do** keep the discussions focused. When a new or related topic comes up it's often better to create new issue than to side track the discussion.
- **Do** clearly state on an issue that you are going to take on implementing it.
- **Do** blog and/or tweet about your contributions!

Don'ts:

- **Don't** surprise the team with big pull requests. We want to support contributors, so we recommend filing an issue and starting a discussion so we can agree on a direction before you invest a large amount of time.
- **Don't** commit code that you didn't write. If you find code that you think is a good fit to add to Semantic Kernel, file an issue and start a discussion before proceeding.
- **Don't** submit PRs that alter licensing related files or headers. If you believe there's a problem with them, file an issue and we'll be happy to discuss it.
- **Don't** make new APIs without filing an issue and discussing with the team first. Adding new public surface area to a library is a big deal and we want to make sure we get it right.

Breaking Changes

Contributions must maintain API signature and behavioral compatibility. If you want to make a change that will break existing code, please file an issue to discuss your idea or change if you believe that a breaking change is warranted. Otherwise, contributions that include breaking changes will be rejected.

The continuous integration (CI) process

The continuous integration (CI) system will automatically perform the required builds and run tests (including the ones you should also run locally) for PRs. Builds and test runs must be clean before a PR can be merged.

If the CI build fails for any reason, the PR issue will be updated with a link that can be used to determine the cause of the failure so that it can be addressed.

Contributing to documentation

We also accept contributions to the [Semantic Kernel documentation repository](#). To learn how to make contributions, please start with the Microsoft [docs contributor guide](#).

Using the kernel to orchestrate AI

Article • 07/11/2023



The term "kernel" can have different meanings in different contexts, but in the case of the Semantic Kernel, the kernel refers to an instance of the processing engine that fulfills a user's request with a collection of [plugins](#).

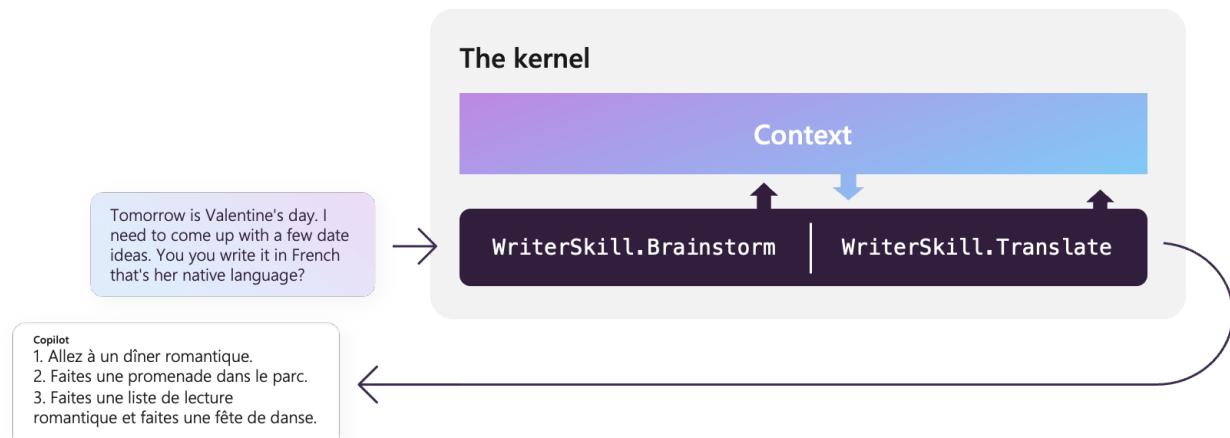
Kernel: "The core, center, or essence of an object or system." —[Wiktionary](#)

ⓘ Note

Skills are currently being renamed to plugins. This article has been updated to reflect the latest terminology, but some images and code samples may still refer to skills.

How does the kernel work?

The kernel has been designed to encourage [function composition](#) which allows developers to combine and interconnect the input and outputs of plugins into a single pipeline. It achieves this by providing a context object that is passed to each plugin in the pipeline. The context object contains the outputs of all previous plugins so the current plugin can use them as inputs.



In this way, the kernel is very similar to the UNIX kernel and its pipes and filters architecture; only now, instead of chaining together programs, we are chaining together AI prompts and native functions.

Using the kernel

To start using the kernel, you must first create an instance of it.

```
C#
```

```
C#  
  
using Microsoft.SemanticKernel;  
  
// Set Simple kernel instance  
IKernel kernel = KernelBuilder.Create();
```

Customizing the kernel

The kernel can also be customized to change how it behaves. During instantiation you can adjust the following properties:

- The initial set of plugins
- The prompt template engine
- The logging engine
- The memory store

Below is an example of how to instantiate the kernel with a custom logger in both C# and Python.

Tip

To see these code samples in action, we recommend checking out the [Loading the kernel](#) guide to explore different ways to load the kernel in a C# or Python notebook.

```
C#
```

You can use the `KernelBuilder` class to create a kernel with custom configuration. The following code snippet shows how to create a kernel with a logger.

```
C#
```

```
using Microsoft.Extensions.Logging;  
using Microsoft.Extensions.Logging.Abstractions;
```

```
ILogger myLogger = NullLogger.Instance;
IKernel kernel_2 = Kernel.Builder
    .WithLogger(myLogger)
    .Build();
```

Additionally, you can create a kernel using a configuration object. This is useful if you want to define a configuration elsewhere in your code and then inject it into the kernel.

C#

```
var config = new KernelConfig();
IKernel kernel_3 = Kernel.Builder
    .WithConfiguration(config)
    .Build();
```

Adding an AI model to the kernel

To add an AI model to the kernel, we can use the out-of-the-box connectors that Semantic Kernel provides. Remember, connectors are what allow you to give you AI a "brain." In this case, we're giving the kernel the ability to think by adding a model. Later, when you learn about [memory](#), you'll see how to give the kernel the ability to remember with connectors.

The following code snippets show how to add a model to the kernel in C# and Python.

C#

```
C#
IKernel kernel = KernelBuilder.Create();

kernel.Config.AddAzureTextCompletionService(
    "my-finetuned-Curie", // Azure OpenAI *Deployment ID*
    "https://contoso.openai.azure.com/", // Azure OpenAI *Endpoint*
    "...your Azure OpenAI Key...", // Azure OpenAI *Key*
    "Azure_curie" // alias used in the prompt
    templates' config.json
);

kernel.Config.AddOpenAITextCompletionService(
    "text-davinci-003", // OpenAI Model Name
    "...your OpenAI API Key...", // OpenAI API key
    "...your OpenAI Org ID...", // *optional* OpenAI Organization ID
);
```

```
        "OpenAI_davinci"          // alias used in the prompt
    templates' config.json
);
```

Importing and registering plugins

Now that you have an instance of the kernel with an AI model, you can start adding plugins to it which will give it the ability to fulfill user requests. Below you can see how to import and register plugins, both from a file and inline.

Tip

Many of the code samples below come from the quick start notebooks. To follow along (and to learn more about how the code works), we recommend checking out the [Running prompts from files](#) and [Running semantic functions inline](#) guides.

The following example leverages the sample [FunSkill plugin](#) that comes with the Semantic Kernel repo.

C#

```
C#  
  
var skillsDirectory =
Path.Combine(System.IO.Directory.GetCurrentDirectory(), "...", "...",
"skills");  
  
var funSkillFunctions =
kernel.ImportSemanticSkillFromDirectory(skillsDirectory, "FunSkill");
var jokeFunction = funSkillFunctions["Joke"];
```

Alternatively, you can register a semantic function inline. To do so, you'll start by defining the semantic prompt and its configuration. The following samples show how you could have registered the same `Joke` function from the `FunSkill` plugin inline.

C#

Create the semantic prompt as a string.

```
C#
```

```

string skPrompt = @"WRITE EXACTLY ONE JOKE or HUMOROUS STORY ABOUT THE
TOPIC BELOW

JOKE MUST BE:
- G RATED
- WORKPLACE/FAMILY SAFE
NO SEXISM, RACISM OR OTHER BIAS/BIGOTRY

BE CREATIVE AND FUNNY. I WANT TO LAUGH.
+++++
{{\$input}}
+++++
";

```

Create the prompt configuration.

```
C#
var promptConfig = new PromptTemplateConfig
{
    Completion =
    {
        MaxTokens = 1000,
        Temperature = 0.9,
        TopP = 0.0,
        PresencePenalty = 0.0,
        FrequencyPenalty = 0.0,
    }
};
```

Register the semantic function.

```
C#
var promptTemplate = new PromptTemplate(
    skPrompt,
    promptConfig,
    kernel
);

var functionConfig = new SemanticFunctionConfig(promptConfig,
    promptTemplate);
var jokeFunction = kernel.RegisterSemanticFunction("FunSkill", "Joke",
    functionConfig);
```

Running a function from a plugin

Whether you imported a plugin from a file or registered it inline, you can now run the plugin to fulfill a user request. Below you can see how to run the joke function you registered above by passing in an input.

C#

C#

```
var result = await jokeFunction.InvokeAsync("time travel to dinosaur  
age");  
  
Console.WriteLine(result);
```

After running the above code examples, you should receive an output like the following.

Output

A time traveler went back to the dinosaur age and was amazed by the size of the creatures. He asked one of the dinosaurs, "How do you manage to get around with such short legs?"

The dinosaur replied, "It's easy, I just take my time!"

Chaining functions within plugins

If you have multiple functions, you can chain them together to create a pipeline. The kernel will automatically pass the outputs of each plugin to the next plugin in the pipeline using the `$input` variable. You can learn more about how to chain functions in the [chaining plugins](#) article.

C#

Create and register the semantic functions.

C#

```
string myJokePrompt = """  
Tell a short joke about {{$INPUT}}.  
""";  
string myPoemPrompt = """  
Take this="{{$INPUT}}" and convert it to a nursery rhyme.  
""";  
string myMenuPrompt = """  
Make this poem="{{$INPUT}}" influence the three items in a coffee shop  
menu.
```

The menu reads in enumerated form:

```
""";  
  
var myJokeFunction = kernel.CreateSemanticFunction(myJokePrompt,  
maxTokens: 500);  
var myPoemFunction = kernel.CreateSemanticFunction(myPoemPrompt,  
maxTokens: 500);  
var myMenuFunction = kernel.CreateSemanticFunction(myMenuPrompt,  
maxTokens: 500);
```

Create and run the pipeline.

C#

```
var myOutput = await kernel.RunAsync(  
    "Charlie Brown",  
    myJokeFunction,  
    myPoemFunction,  
    myMenuFunction);  
  
Console.WriteLine(myOutput);
```

The above code will output something like the following.

Output

1. Colossus of Memnon Latte - A creamy latte with a hint of sweetness, just like the awe-inspiring statue.
2. Gasp and Groan Mocha - A rich and indulgent mocha that will make you gasp and groan with delight.
3. Heart Skipping a Beat Frappuccino - A refreshing frappuccino with a hint of sweetness that will make your heart skip a beat.

Take the next step

This article only scratches the surface of what you can do with the kernel. To learn more about additional features, check out the following articles.

Your goal	Next step
Learn what plugins are and what they can do	Understand AI plugins
Create more advanced pipelines with Semantic Kernel	Chaining functions together

Your goal	Next step
Automatically creating pipelines with Planner	Auto create plans with planner
Simulating memory within Semantic Kernel	Give you AI memories

[Understanding plugins](#)

Using AI plugins in Semantic Kernel

Article • 07/12/2023



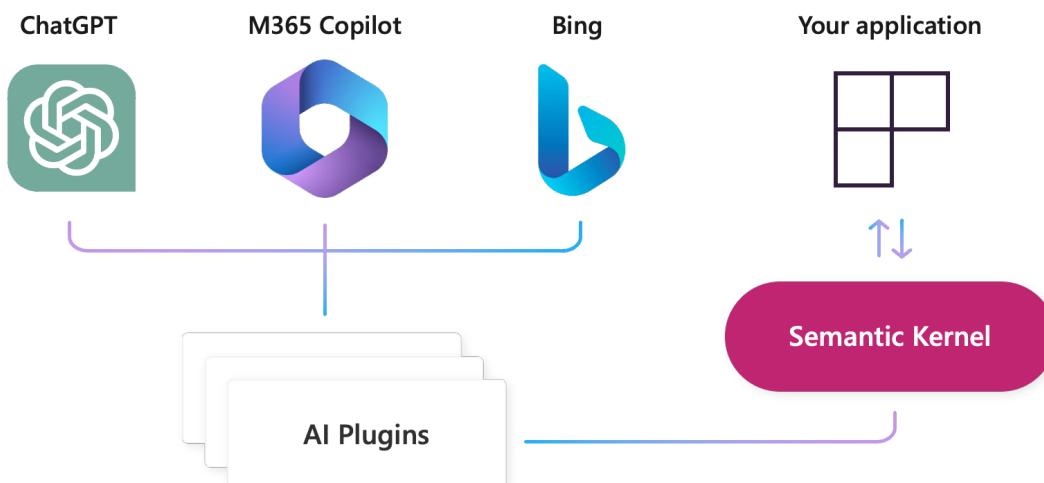
With plugins, you can encapsulate AI capabilities into a single unit of functionality. Plugins are the building blocks of the Semantic Kernel and can interoperate with plugins in ChatGPT, Bing, and Microsoft 365.

ⓘ Note

Skills are currently being renamed to plugins. This article has been updated to reflect the latest terminology, but some images and code samples may still refer to skills.

What is a plugin?

To drive alignment across the industry, we're in progress of adopting the [OpenAI plugin specification](#) as the standard for plugins. This will help create an ecosystem of interoperable plugins that can be used across all of the major AI apps and services like ChatGPT, Bing, and Microsoft 365.



For developers using Semantic Kernel, this means any plugins you build will soon be usable in ChatGPT, Bing, and Microsoft 365, allowing you to increase the reach of your AI capabilities without rewriting your code. It also means that plugins built for ChatGPT, Bing, and Microsoft 365 will be able to integrate with Semantic Kernel.

To see our status on adopting the OpenAI plugin specification, see [our roadmap](#) at the bottom of the page. Each of the integration points have accompanying issues on GitHub that you can follow to stay up-to-date on our progress. We are also open to feedback on our unification plan, so feel free to leave a comment on the GitHub issues so that we can best support your needs.

What does a plugin look like?

At a high-level, a plugin is a group of functions that can be exposed to AI apps and services. The functions within plugins can then be orchestrated by an AI application to accomplish user requests. Within Semantic Kernel, you can invoke these functions either manually (see [chaining functions](#)) or automatically with a [planner](#).

Just providing functions, however, is not enough to make a plugin. To power automatic orchestration with a [planner](#), plugins *also* need to provide details that semantically describe how they behave. Everything from the function's input, output, and side effects need to be described in a way that the AI can understand, otherwise, planner will provide unexpected results.

For example, in the [WriterSkill plugin](#), each function has a semantic description that describes what the function does. Planner can then use this description to choose the best function to call based on a user's ask.

In the picture on the right, planner would likely use the `ShortPoem` and `StoryGen` functions to satisfy the users ask thanks to the provided semantic descriptions.

Writer plugin

Function	Description for model
Brainstorm	Given a goal or topic description generate a list of ideas.
EmailGen	Write an email from the given bullet points.
ShortPoem	Turn a scenario into a short and entertaining poem.
StoryGen	Generate a list of synopsis for a novel or novella with sub-chapters.
Translate	Translate the input into a language of your choice.

Can you write me a short poem about living in Dublin, Ireland and then create a story based on the poem?



Planner

Copilot
Sure! Here's a story based on living along the Grand Canal in Dublin, Ireland...



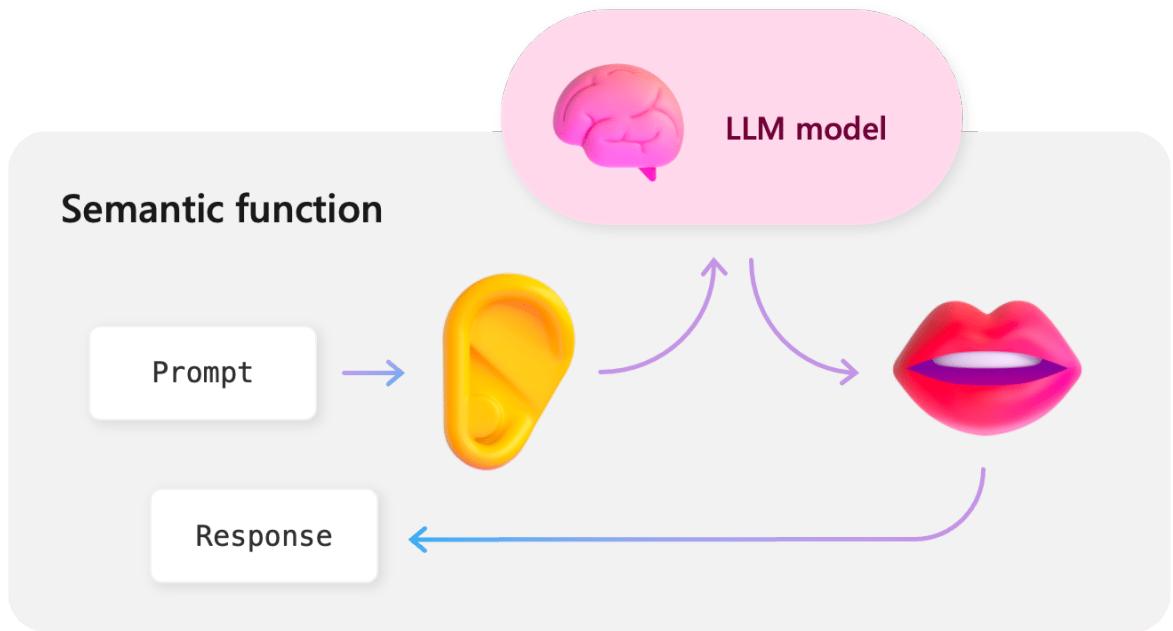
Adding functions to plugins

Now that you know what a plugin is, let's take a look at how to create one. Within a plugin, you can create two types of functions: semantic functions and native functions. The following sections describe how to create each type. For further details, please refer to the [Creating semantic functions](#) and [Creating native functions](#) articles.

Semantic functions

If plugins represent the "body" of your AI app, then semantic functions would represent the ears and mouth of your AI. They allow your AI app to listen to users asks and respond back with a natural language response.

To connect the ears and the mouth to the "brain," Semantic Kernel uses connectors. This allows you to easily swap out the AI services without rewriting code.



Below is an sample called `Summarize` that can be found in the [samples folder](#) in the GitHub repository.

```

Prompt

[SUMMARIZATION RULES]
DON'T WASTE WORDS
USE SHORT, CLEAR, COMPLETE SENTENCES.
DO NOT USE BULLET POINTS OR DASHES.
USE ACTIVE VOICE.
MAXIMIZE DETAIL, MEANING
FOCUS ON THE CONTENT

[BANNED PHRASES]
This article
This document
This page
This material
[END LIST]

Summarize:
Hello how are you?
+++++
Hello

Summarize this
{{$input}}
+++++

```

To semantically describe this function (as well as define the configuration for the AI service), you must also create a `config.json` file in the same folder as the prompt. This file

describes the function's input parameters and description. Below is the `config.json` file for the `Summarize` function.

```
JSON

{
  "schema": 1,
  "type": "completion",
  "description": "Summarize given text or any text document",
  "completion": {
    "max_tokens": 512,
    "temperature": 0.0,
    "top_p": 0.0,
    "presence_penalty": 0.0,
    "frequency_penalty": 0.0
  },
  "input": {
    "parameters": [
      {
        "name": "input",
        "description": "Text to summarize",
        "defaultValue": ""
      }
    ]
  }
}
```

Both `description` fields are used by `planner`, so it's important to provide a detailed, yet concise, description so planner can make the best decision when orchestrating functions together. We recommend testing multiple descriptions to see which one works best for the widest range of scenarios.

You can learn more about creating semantic functions in the [Creating semantic functions](#) article. In this article you'll learn the best practices for the following:

- ✓ How to create semantic functions
- ✓ Adding input parameters
- ✓ Calling functions within semantic functions

[Learn more about creating semantic functions](#)

Native functions

With native functions, you can have the kernel call C# or Python code directly so that you can manipulate data or perform other operations. In this way, native functions are like the hands of your AI app. They can be used to save data, retrieve data, and perform

any other operation that you can do in code that is ill-suited for LLMs (e.g., performing calculations).

Native function



Save information



Look things up



Perform actions

Instead of providing a separate configuration file with semantic descriptions, planner is able to use annotations in the code to understand how the function behaves. Below are examples of the annotations used by planner in both C# and Python for out-of-the-box native functions.

C#

The following code is an excerpt from the `DocumentSkill` plugin, which can be found in the [document plugin](#) folder in the GitHub repository. It demonstrates how you can use the `SKFunction` and `SKFunctionInput` attributes to describe the function's input and output to planner.

C#

```
[SKFunction("Read all text from a document")]
[SKFunctionInput(Description = "Path to the file to read")]
public async Task<string> ReadTextAsync(string filePath, SKContext
context)
{
    this._logger.LogInformation("Reading text from {0}", filePath);
    using var stream = await
this._fileSystemConnector.GetFileContentStreamAsync(filePath,
context.CancellationToken).ConfigureAwait(false);
    return this._documentConnector.ReadText(stream);
}
```

You can learn more about creating native functions in the [Creating native functions](#) article. In this article you'll learn the best practices for the following:

- ✓ How to create simple native functions

- ✓ Calling Semantic Kernel functions from within native functions
- ✓ Different ways to invoke native functions

[Learn more about creating native functions](#)

Plugins will soon be interoperable

Both OpenAI and Microsoft are converging on the same [plugin definition ↗](#), so as a developer, you can build one plugin that works seamlessly across ChatGPT, Bing, Microsoft 365 Business Chat, and soon, Semantic Kernel.

To support this effort, the Semantic Kernel team is in the progress of updating what was previously known as "skills" into plugins. After this unification process, we will allow developers to bring outside AI plugins into Semantic Kernel *and* deploy Semantic Kernel plugins so they can be consumed by other apps and services.

The following tables describe the teams progress towards plugin unification. The  symbol means the feature is currently being planned for.

Bringing ChatGPT plugins to Semantic Kernel

Feature	Description	Link	Status
Use OpenAPI files in Semantic Kernel	With the OpenAPI plugin, you can take an OpenAPI file from a ChatGPT plugin and bring it into Semantic Kernel	Blog post ↗	
Import a ai-plugin.json manifest in Semantic Kernel	This will make it even easier to bring in a ChatGPT plugin. By providing an ai-plugin.json file, we can pull in the relevant OpenAPI file and also provide the plugin description planner for more effective plan generation.	Coming soon	

Bringing Semantic Kernel plugins to ChatGPT

Feature	Description	Link	Status
Deploy Semantic Kernel plugins for ChatGPT	Today, you can deploy your existing semantic and native functions to Azure. Once they've been deployed, you can manually craft the necessary OpenAPI specification and ai-plugin.json file so that your functions can be used by ChatGPT.	Documentation	

Feature	Description	Link	Status
Automatically generate OpenAPI spec and ai-plugin.json file	To make it easier to deploy your plugins to ChatGPT, we will help you automatically create the files that ChatGPT requires to use a plugin. You can then host these files yourself so you can provide them to ChatGPT.	<i>Coming soon</i>	

Take the next step

Now that you understand the basics of plugins, you can now go deeper into the details of creating semantic and native functions for your plugin.

[Create a semantic function](#)

Creating semantic functions

Article • 07/12/2023



In previous articles, we demonstrated [how to load a semantic function](#). We also showed how to run the function either [by itself](#) or [in a chain](#). In both cases, we used out-of-the-box sample functions that are included with Semantic Kernel to demonstrate the process.

In this article, we'll demonstrate how to actually *create* a semantic function so you can easily import them into Semantic Kernel. As an example in this article, we will demonstrate how to create a semantic function that gathers the intent of the user. This semantic function will be called `GetIntent` and will be part of a plugin called `OrchestratorPlugin`.

By following this example, you'll learn how to create a semantic function that can use multiple context variables and functions to elicit an AI response. If you want to see the final solution, you can check out the following samples in the public documentation repository.

Language	Link to final solution
C#	Open solution in GitHub ↗
Python	Open solution in GitHub ↗

ⓘ Note

Skills are currently being renamed to plugins. This article has been updated to reflect the latest terminology, but some images and code samples may still refer to skills.

💡 Tip

We recommend using the **Semantic Kernel Tools** extension for Visual Studio Code to help you create semantic functions. This extension provides an easy way to

create and test functions directly from within VS Code.



Semantic Kernel Tools

Microsoft microsoft.com | 3,604 installs | ★★★★★ (2) | Free

Prompt engineering tools to create AI plugins with Semantic Kernel

[Install](#)

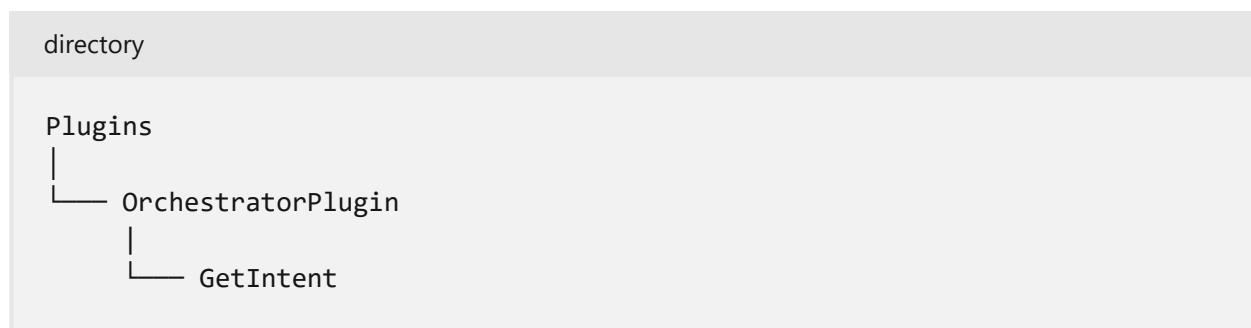
[Trouble Installing?](#)



Creating a home for your semantic functions

Before creating the `OrchestratorPlugin` or the `GetIntent` function, you must first define a folder that will hold all of your plugins. This will make it easier to import them into Semantic Kernel later. We recommend putting this folder at the root of your project and calling it *plugins*.

Within your *plugins* folder, you can then create a folder called *OrchestratorPlugin* for your plugin and a nested folder called *GetIntent* for your function.

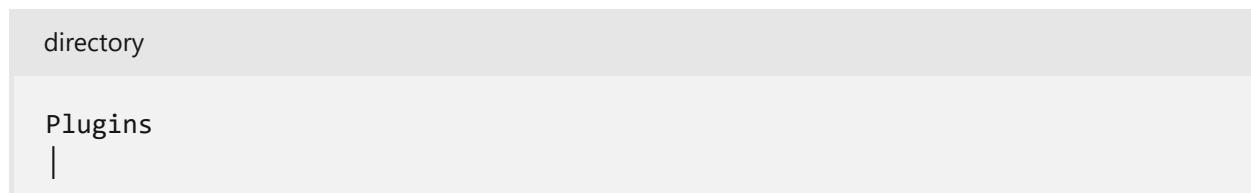


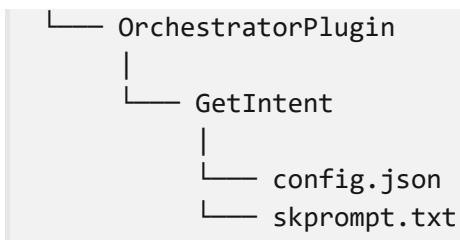
To see a more complete example of a plugins directory, check out the [Semantic Kernel sample plugins](#) folder in the GitHub repository.

Creating the files for your semantic function

Once inside of a semantic functions folder, you'll need to create two files: `skprompt.txt` and `config.json`. The `skprompt.txt` file contains the prompt that will be sent to the AI service and the `config.json` file contains the configuration along with semantic descriptions used by planner.

Go ahead and create these two files in the *GetIntent* folder. In the following sections, we'll walk through how to configure them.





Writing a prompt in the `skprompt.txt` file

The `skprompt.txt` file contains the request that will be sent to the AI service. The [prompt engineering](#) section of the documentation provides a detailed overview of how to write prompts, but at a high level, prompts are requests written in natural language that are sent to an AI service.

In most cases, you'll send your prompt to a text or chat completion service which will return back a response that attempts to complete the prompt. For example, if you send the prompt `I want to go to the`, the AI service might return back `beach`. This is a very simple example, but it demonstrates the basic idea of how prompts work.

In the case of the `GetIntent` function, we want to create a prompt that asks the AI service what the intent of a user is. The following prompt will do just that:

```
txt

Bot: How can I help you?
User: {{$input}}


-----
The intent of the user in 5 words or less:
```

Notice that we're using a variable called `$input` in the prompt. This variable is later defined in the `config.json` file and is used to pass the user's input to the AI service.

Go ahead and copy the prompt above and save it in the `skprompt.txt` file.

Configuring the function in the `config.json` file

Next, we need to define the configuration for the `GetIntent` function. The [configuring prompts](#) article provides a detailed overview of how to configure prompts, but at a high level, prompts are configured using a JSON file that contains the following properties:

- `type` – The type of prompt. In this case, we're using the `completion` type.

- `description` – A description of what the prompt does. This is used by planner to automatically orchestrate plans with the function.
- `completion` – The settings for completion models. For OpenAI models, this includes the `max_tokens` and `temperature` properties.
- `input` – Defines the variables that are used inside of the prompt (e.g., `$input`).

For the `GetIntent` function, we'll use the following configuration:

JSON

```
{
  "schema": 1,
  "type": "completion",
  "description": "Gets the intent of the user.",
  "completion": {
    "max_tokens": 500,
    "temperature": 0.0,
    "top_p": 0.0,
    "presence_penalty": 0.0,
    "frequency_penalty": 0.0
  },
  "input": {
    "parameters": [
      {
        "name": "input",
        "description": "The user's request.",
        "defaultValue": ""
      }
    ]
  }
}
```

Copy the configuration above and save it in the `config.json` file.

Testing your semantic function

At this point, you can import and test your function with the kernel by using the following code.

C#

C#

```
var pluginsDirectory =
Path.Combine(System.IO.Directory.GetCurrentDirectory(), "path", "to",
"your", "plugins", "folder");
```

```
// Import the OrchestratorPlugin from the plugins directory.  
var orchestratorPlugin = kernel  
    .ImportSemanticSkillFromDirectory(pluginsDirectory,  
    "OrchestratorPlugin");  
  
// Get the GetIntent function from the OrchestratorPlugin and run it  
var result = await orchestratorPlugin["GetIntent"]  
    .InvokeAsync("I want to send an email to the marketing team  
celebrating their recent milestone.");  
  
Console.WriteLine(result);
```

You should get an output that looks like the following:

Output

Send congratulatory email.

Making your semantic function more robust

While our function works, it's not very useful when combined with native code. For example, if we had several native functions available to run based on an intent, it would be difficult to use the output of the `GetIntent` function to choose which native function to actually run.

We need to find a way to constrain the output of our function so that we can use the output in a switch statement.

Templatizing a semantic function

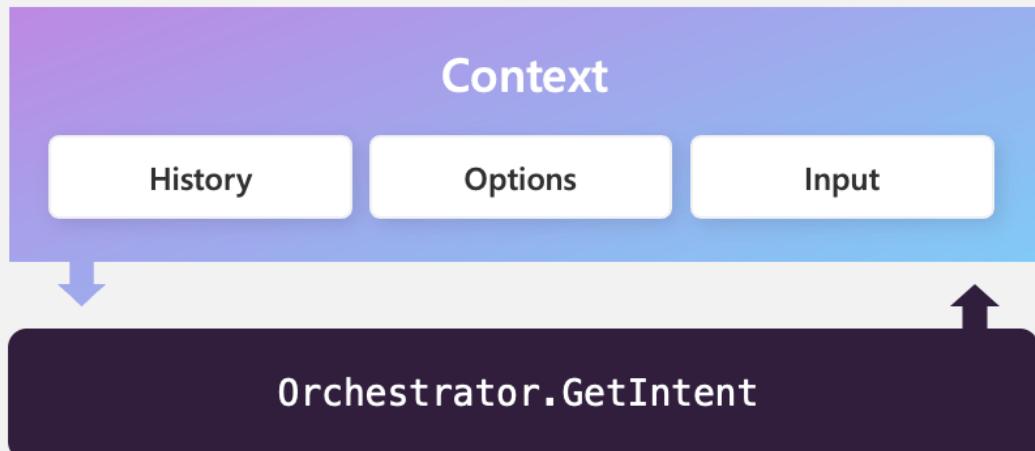
One way to constrain the output of a semantic function is to provide a list of options for it to choose from. A naive approach would be to hard code these options into the prompt, but this would be difficult to maintain and would not scale well. Instead, we can use Semantic Kernel's templating language to dynamically generate the prompt.

The [prompt template syntax](#) article in the [prompt engineering](#) section of the documentation provides a detailed overview of how to use the templating language. In this article, we'll show you just enough to get started.

The following prompt uses the `>{{$options}}` variable to provide a list of options for the LLM to choose from. We've also added a `>{{$history}}` variable to the prompt so that the previous conversation is included.

By including these variables, we are able to help the LLM choose the correct intent by allowing it to leverage variables within the Semantic Kernel context object.

The kernel



txt

```
{{$history}}
User: {{$input}}
```

Provide the intent of the user. The intent should be one of the following:
{{\$options}}

INTENT:

When you add a new variable to the prompt, you also should update the `config.json` file to include the new variable. While these properties aren't used now, it's good to get into the practice of adding them so that they can be used by the [planner](#) later. The following configuration adds the `$options` and `$history` variable to the `input` section of the configuration.

JSON

```
{
  "schema": 1,
  "type": "completion",
  "description": "Gets the intent of the user.",
  "completion": {
    "max_tokens": 500,
    "temperature": 0.0,
    "top_p": 0.0,
    "presence_penalty": 0.0,
```

```
        "frequency_penalty": 0.0
    },
    "input": {
        "parameters": [
            {
                "name": "input",
                "description": "The user's request.",
                "defaultValue": ""
            },
            {
                "name": "history",
                "description": "The history of the conversation.",
                "defaultValue": ""
            },
            {
                "name": "options",
                "description": "The options to choose from.",
                "defaultValue": ""
            }
        ]
    }
}
```

You can now update your code to provide a list of options to the `GetIntent` function by using context.

C#

```
C#
```

```
// Import the OrchestratorPlugin from the plugins directory.
var pluginsDirectory =
Path.Combine(System.IO.Directory.GetCurrentDirectory(), "path", "to",
"your", "plugins", "folder");
var orchestrationPlugin = kernel
    .ImportSemanticSkillFromDirectory(pluginsDirectory,
"OrchestratorPlugin");

// Create a new context and set the input, history, and options
variables.
var context = kernel.CreateNewContext();
context["input"] = "Yes";
context["history"] = @"Bot: How can I help you?
User: My team just hit a major milestone and I would like to send them a
message to congratulate them.
Bot:Would you like to send an email?";
context["options"] = "SendEmail, ReadEmail, SendMeeting, RsvpToMeeting,
SendChat";

// Run the GetIntent function with the context.
var result = await
orchestrationPlugin["GetIntent"].InvokeAsync(context);
```

```
Console.WriteLine(result);
```

Now, instead of getting an output like `Send congratulatory email.`, we'll get an output like `SendEmail`. This output could then be used within a switch statement in native code to run the correct function.

Calling functions *within* a semantic function

We now have a more useful semantic function, but you might run into token limits if you had a long list of options and a long conversation history. To get around this, we can call other functions within our semantic function to help break up the prompt into smaller pieces.

To learn more about calling functions within a semantic function, see the [calling functions within a semantic function](#) section in the [prompt engineering](#) section of the documentation.

The following prompt uses the `Summarize` function in the [SummarizeSkill plugin](#) to summarize the conversation history before asking for the intent.

txt

```
{{SummarizeSkill.Summarize $history}}
User: {{$input}}
```

```
-----
Provide the intent of the user. The intent should be one of the following:
{{$options}}
```

INTENT:

You can now update your code to load the `SummarizeSkill` plugin so the kernel can find the `Summarize` function.

C#

```
C#
var pluginsDirectory =
Path.Combine(System.IO.Directory.GetCurrentDirectory(), "path", "to",
"your", "plugins", "folder");
```

```
// Import the OrchestratorPlugin and SummarizeSkill from the plugins
// directory.
var orchestrationPlugin =
kernel.ImportSemanticSkillFromDirectory(pluginsDirectory,
"OrchestratorPlugin");
var summarizationPlugin =
kernel.ImportSemanticSkillFromDirectory(pluginsDirectory,
"SummarizeSkill");

// Create a new context and set the input, history, and options
variables.
var context = kernel.CreateNewContext();
context["input"] = "Yes";
context["history"] = @"Bot: How can I help you?
User: My team just hit a major milestone and I would like to send them a
message to congratulate them.
Bot:Would you like to send an email?";
context["options"] = "SendEmail, ReadEmail, SendMeeting, RsvpToMeeting,
SendChat";

// Run the Summarize function with the context.
var result = await
orchestrationPlugin["GetIntent"].InvokeAsync(context);

Console.WriteLine(result);
```

Take the next step

Now that you can create a semantic function, you can now learn how to [create a native function](#).

[Create a native function](#)

Adding native functions to the kernel

Article • 07/14/2023



In the [how to create semantic functions](#) article, we showed how you could create a semantic function that retrieves a user's intent, but what do you do once you have the intent? In this article we'll show how to create native functions that can route the intent and perform a task.

As an example, we'll add an additional function to the `OrchestratorPlugin` we created in the [semantic functions](#) article to route the user's intent. We'll also add a new plugin called `MathPlugin` that will perform simple arithmetic for the user.

By the end of this article, you'll have a kernel that can correctly answer user questions like `What is the square root of 634?` and `What is 42 plus 1513?`. If you want to see the final solution, you can check out the following samples in the public documentation repository.

Language	Link to final solution
C#	Open solution in GitHub ↗
Python	Open solution in GitHub ↗

ⓘ Note

Skills are currently being renamed to plugins. This article has been updated to reflect the latest terminology, but some images and code samples may still refer to skills.

Finding a home for your native function

You can place plugins in the same directory as the other plugins. For example, to create functions for a plugin called `MyNewPlugin`, you can create a new file called `MyCSharpPlugin.cs` in the same directory as your semantic functions.

directory

```
MyPluginsDirectory
└── MyNewPlugin
    ├── MyFirstSemanticFunction
    │   └── skprompt.txt
    │   └── config.json
    ├── MyOtherSemanticFunctions
    |   ...
    └── MyNewPlugin.cs
```

In our example, we'll create two files one for the native `OrchestratorPlugin` functions and another for the `MathPlugin`. Depending on the language you're using, you'll create either C# or Python files for each.

C#

```
directory
Plugins
└── OrchestratorPlugin
    └── GetIntent
        └── skprompt.txt
        └── config.json
        └── OrchestratorPlugin.cs
└── MathPlugin
    └── MathPlugin.cs
```

It's ok if you have a plugin folder with native functions and semantic functions. The kernel will load both functions into the same plugin namespace. What's important is that you don't have two functions with the same name within the same plugin namespace. If you do, the last function loaded will overwrite the previous function.

We'll begin by creating the `MathPlugin` functions. Afterwards, we'll call the `MathPlugin` functions from within the `OrchestratorPlugin`. At the end of this example you will have the following supported functions.

Plugin	Function	Type	Description
OrchestratorPlugin	GetIntent	Semantic	Gets the intent of the user

Plugin	Function	Type	Description
OrchestratorPlugin	GetNumbers	Semantic	Gets the numbers from a user's request
OrchestratorPlugin	RouteRequest	Native	Routes the request to the appropriate function
MathPlugin	Sqrt	Native	Takes the square root of a number
MathPlugin	Add	Native	Adds two numbers together

Creating simple native functions

Open up the *MathPlugin.cs* or *MathPlugin.py* file you created earlier and follow the instructions below to create the two necessary functions: `Sqrt` and `Add`.

Defining a function that takes a single input

Add the following code to your file to create a function that takes the square root of a number.

C#

```
C#  
  
using Microsoft.SemanticKernel.SkillDefinition;  
using Microsoft.SemanticKernel.Orchestration;  
  
namespace Plugins;  
  
public class MathPlugin  
{  
    [SKFunction("Takes the square root of a number")]  
    public string Sqrt(string number)  
    {  
        return Math.Sqrt(Convert.ToDouble(number)).ToString();  
    }  
}
```

Notice that the input and return types are strings. This is because the kernel will pass the input as a string and expect a string to be returned. You can convert the input to any type you want within the function.

Also notice how we've added a description to each function with attributes. This description will be used by the [planner](#) to automatically create a plan using these

functions. In our case, we're telling planner that this function [Takes the square root of a number](#).

Using context parameters to take multiple inputs

Adding numbers together requires multiple numbers as input. Since we cannot pass multiple numbers into a native function, we'll need to use context parameters instead.

Add the following code to your `MathPlugin` class to create a function that adds two numbers together.

```
C#  
  
C#  
  
[SKFunction("Adds two numbers together")]
[SKFunctionContextParameter(Name = "input", Description = "The first
number to add")]
[SKFunctionContextParameter(Name = "number2", Description = "The second
number to add")]
public string Add(SKContext context)
{
    return (
        Convert.ToDouble(context["input"]) +
        Convert.ToDouble(context["number2"])
    ).ToString();
}
```

Notice that instead of taking a string as input, this function takes an `SKContext` object as input. This object contains all of the variables in the Semantic Kernel's context. We can use this object to retrieve the two numbers we want to add. Also notice how we provide descriptions for each of the context parameters. These descriptions will be used by the [planner](#) to automatically provide inputs to this function.

The `SKContext` object only supports strings, so we'll need to convert the strings to doubles before we add them together.

Running your native functions

You can now run your functions using the code below. Notice how we pass in the multiple numbers required for the `Add` function using a `SKContext` object.

```
C#
```

```
C#  
  
using Microsoft.SemanticKernel;  
using Plugins;  
  
// ... instantiate your kernel  
  
var mathPlugin = kernel.ImportSkill(new MathPlugin(), "MathPlugin");  
  
// Run the Sqrt function  
var result1 = await mathPlugin["Sqrt"].InvokeAsync("64");  
Console.WriteLine(result1);  
  
// Run the Add function with multiple inputs  
var context = kernel.CreateNewContext();  
context["input"] = "3";  
context["number2"] = "7";  
var result2 = await mathPlugin["Add"].InvokeAsync(context);  
Console.WriteLine(result2);
```

The code should output 8 since it's the square root of 64 and 10 since it's the sum of 3 and 7.

Creating a more complex native function

We can now create the native routing function for the `OrchestratorPlugin`. This function will be responsible for calling the right `MathPlugin` function based on the user's request.

Calling Semantic Kernel functions within a native function

In order to call the right `MathPlugin` function, we'll use the `GetIntent` semantic function we defined in the [previous tutorial](#). Add the following code to your `OrchestratorPlugin` class to get started creating the routing function.

```
C#
```

```
C#  
  
using Microsoft.SemanticKernel;  
using Microsoft.SemanticKernel.Orchestration;  
using Microsoft.SemanticKernel.SkillDefinition;  
using Newtonsoft.Json.Linq;  
  
namespace Plugins;
```

```

public class OrchestratorPlugin
{
    IKernel _kernel;

    public OrchestratorPlugin(IKernel kernel)
    {
        _kernel = kernel;
    }

    [SKFunction("Routes the request to the appropriate function.")]
    public async Task<string> RouteRequest(SKContext context)
    {
        // Save the original user request
        string request = context["input"];

        // Add the list of available functions to the context
        context["options"] = "Sqrt, Add";

        // Retrieve the intent from the user request
        var GetIntent = _kernel.Skills.GetFunction("OrchestratorPlugin",
            "GetIntent");
        await GetIntent.InvokeAsync(context);
        string intent = context["input"].Trim();

        // Call the appropriate function
        switch (intent)
        {
            case "Sqrt":
                // Call the Sqrt function
            case "Add":
                // Call the Add function
            default:
                return "I'm sorry, I don't understand.";
        }
    }
}

```

From the code above, we can see that the `OrchestratorPlugin` function does the following:

1. Saves the kernel to a private variable during initialization so it can be used later to call the `GetIntent` function.
2. Sets the list of available functions to the context so it can be passed to the `GetIntent` function.
3. Uses a `switch` statement to call the appropriate function based on the user's intent.

Unfortunately, we have a challenge. Despite knowing the user's intent, we don't know which numbers to pass to the `MathPlugin` functions. We'll need to add another semantic

function to the `OrchestratorPlugin` to extract the necessary numbers from the user's input.

Using semantic functions to extract data for native functions

To pull the numbers from the user's input, we'll create a semantic function called `GetNumbers`. Create a new folder under the `OrchestratorPlugin` folder named `GetNumbers`. Then create a `skprompt.txt` and `config.json` file within the folder. Add the following code to the `skprompt.txt` file.

txt

Extract the numbers from the input and output them in JSON format.

INPUT: Take the square root of 4
OUTPUT: {"number1":4}

INPUT: Subtract 3 dollars from 2 dollars
OUTPUT: {"number1":2,"number2":3}

INPUT: I have a 2x4 that is 3 feet long. Can you cut it in half?
OUTPUT: {"number1":3, "number2":2}

INPUT: {{\$input}}
OUTPUT:

Add the following code to the `config.json` file.

JSON

```
{  
    "schema": 1,  
    "type": "completion",  
    "description": "Gets the numbers from a user's request.",  
    "completion": {  
        "max_tokens": 500,  
        "temperature": 0.0,  
        "top_p": 0.0,  
        "presence_penalty": 0.0,  
        "frequency_penalty": 0.0  
    },  
    "input": {  
        "parameters": [  
            {  
                "name": "input",  
                "description": "The user's request."  
            }  
        ]  
    }  
}
```

```
        "defaultValue": ""
    }
]
}
}
```

Putting it all together

We can now call the `GetNumbers` function from the `OrchestratorPlugin` function. Replace the `switch` statement in the `OrchestratorPlugin` function with the following code.

C#

```
C#  
  
var GetNumbers = _kernel.Skills.GetFunction("OrchestratorPlugin",
"GetNumbers");
SKContext getNumberContext = await GetNumbers.InvokeAsync(request);
JObject numbers = JObject.Parse(getNumberContext["input"]);  
  
// Call the appropriate function
switch (intent)
{
    case "Sqrt":
        // Call the Sqrt function with the first number
        var Sqrt = _kernel.Skills.GetFunction("MathPlugin", "Sqrt");
        SKContext sqrtResults = await
Sqrt.InvokeAsync(numbers["number1"]!.ToString());  
  
        return sqrtResults["input"];
    case "Add":
        // Call the Add function with both numbers
        var Add = _kernel.Skills.GetFunction("MathPlugin", "Add");
        context["input"] = numbers["number1"]!.ToString();
        context["number2"] = numbers["number2"]!.ToString();
        SKContext addResults = await Add.InvokeAsync(context);  
  
        return addResults["input"];
    default:
        return "I'm sorry, I don't understand.";
}
```

Finally, you can invoke the `OrchestratorPlugin` function from your main file using the code below.

C#

C#

```
using Microsoft.SemanticKernel;
using Plugins;

// ... instantiate your kernel

var pluginsDirectory =
Path.Combine(System.IO.Directory.GetCurrentDirectory(), "plugins");

// Import the semantic functions
kernel.ImportSemanticSkillFromDirectory(pluginsDirectory,
"OrchestratorPlugin");
kernel.ImportSemanticSkillFromDirectory(pluginsDirectory,
"SummarizeSkill");

// Import the native functions
var mathPlugin = kernel.ImportSkill(new MathPlugin(), "MathPlugin");
var orchestratorPlugin = kernel.ImportSkill(new
OrchestratorPlugin(kernel), "OrchestratorPlugin");

// Make a request that runs the Sqrt function
var result1 = await orchestratorPlugin["RouteRequest"].InvokeAsync("What
is the square root of 634?");
Console.WriteLine(result1);

// Make a request that runs the Add function
var result2 = await orchestratorPlugin["RouteRequest"].InvokeAsync("What
is 42 plus 1513?");
Console.WriteLine(result2);
```

Take the next step

You now have the skills necessary to create both semantic and native functions to create custom plugins, but up until now, we've only called one function at a time. In the next article, you'll learn how to chain multiple functions together.

[Chaining functions](#)

Chaining functions together

Article • 07/12/2023



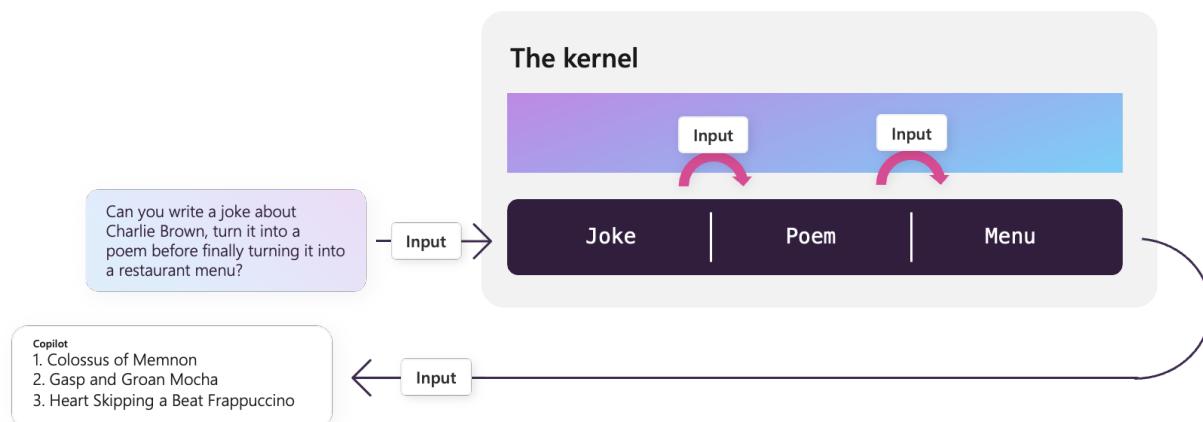
In previous articles, we showed how you could invoke a Semantic Kernel function (whether semantic or native) individually. Oftentimes, however, you may want to string multiple functions together into a single pipeline to simplify your code.

[Later in this article](#), we'll put this knowledge to use by demonstrating how you could refactor the code from the [native functions](#) to make it more readable and maintainable. If you want to see the final solution, you can check out the following samples in the public documentation repository.

Language	Link to final solution
C#	Open solution in GitHub ↗
Python	Open solution in GitHub ↗

Passing data to semantic functions with `$input`

Semantic Kernel was designed in the spirit of UNIX's piping and filtering capabilities. To replicate this behavior, we've added a special variable called `$input` into the kernel's context object that allows you to stream output from one semantic function to the next.



For example we can make three inline semantic functions and string their outputs into the next by adding the `$input` variable into each prompt.



Create and register the semantic functions.

C#

```
string myJokePrompt = """  
Tell a short joke about {{$input}}.  
""";  
string myPoemPrompt = """  
Take this "{$input}" and convert it to a nursery rhyme.  
""";  
string myMenuPrompt = """  
Make this poem "{$input}" influence the three items in a coffee shop  
menu.  
The menu reads in enumerated form:  
""";  
  
var myJokeFunction = kernel.CreateSemanticFunction(myJokePrompt,  
maxTokens: 500);  
var myPoemFunction = kernel.CreateSemanticFunction(myPoemPrompt,  
maxTokens: 500);  
var myMenuFunction = kernel.CreateSemanticFunction(myMenuPrompt,  
maxTokens: 500);
```

Run the functions sequentially. Notice how all of the functions share the same context.

C#

```
var context = kernel.CreateNewContext("Charlie Brown");  
await myJokeFunction.InvokeAsync(context);  
await myPoemFunction.InvokeAsync(context);  
await myMenuFunction.InvokeAsync(context);  
  
Console.WriteLine(context);
```

Which would result in something like:

Output

1. Colossus of Memnon Latte - A creamy latte with a hint of sweetness, just like the awe-inspiring statue.
2. Gasp and Groan Mocha - A rich and indulgent mocha that will make you gasp and groan with delight.
3. Heart Skipping a Beat Frappuccino - A refreshing frappuccino with a hint of sweetness that will make your heart skip a beat.

Using the `RunAsync` method to simplify your code

Running each function individually can be very verbose, so Semantic Kernel also provides the `RunAsync` method in C# or `run_async` method in Python that automatically calls a series of functions sequentially, all with the same context object.

```
C#
```

```
C#
```

```
var myOutput = await kernel.RunAsync(
    new ContextVariables("Charlie Brown"),
    myJokeFunction,
    myPoemFunction,
    myMenuFunction);

Console.WriteLine(myOutput);
```

Passing more than just `$input` with native functions

In the previous articles, we've already seen how you can update and retrieve additional properties from the context object within native functions. We can use this same technique to pass additional data between functions within a pipeline.

We'll demonstrate this by updating the code written in the [native functions](#) article to use the `RunAsync` method instead.

Adding a function that changes variables in the context

In the previous example, we used the `RouteRequest` function to individually call each of the Semantic Kernel functions, and in between calls, we updated the context object with the new data. We can simplify this code by creating a new native function that performs the same context update operations. We'll call this function `ExtractNumbersFromJson` and it will take the JSON string from the `input` variable and extract the numbers from it.

Add the following code to your `OrchestratorPlugin` class.

```
C#
```

```
C#
```

```
[SKFunction("Extracts numbers from JSON")]
public SKContext ExtractNumbersFromJson(SKContext context)
{
    JObject numbers = JObject.Parse(context["input"]);

    // loop through numbers and add them to the context
    foreach (var number in numbers)
    {
        if (number.Key == "number1")
        {
            // add the first number to the input variable
            context["input"] = number.Value.ToString();
            continue;
        }
        else
        {
            // add the rest of the numbers to the context
            context[number.Key] = number.Value.ToString();
        }
    }
    return context;
}
```

Using the `RunAsync` method to chain our functions

Now that we have a function that can extracts numbers, we can update our `RouteRequest` function to use the `RunAsync` method to call the functions in a pipeline. Update the `RouteRequest` function to the following:

C#

```
C#
```

```
[SKFunction("Routes the request to the appropriate function.")]
public async Task<string> RouteRequest(SKContext context)
{
    // Save the original user request
    string request = context["input"];

    // Add the list of available functions to the context
    context["options"] = "Sqrt, Add";

    // Retrieve the intent from the user request
    var GetIntent = _kernel.Skills.GetFunction("OrchestratorPlugin",
    "GetIntent");
    await GetIntent.InvokeAsync(context);
    string intent = context["input"].Trim();
```

```

    // Prepare the functions to be called in the pipeline
    var GetNumbers = _kernel.Skills.GetFunction("OrchestratorPlugin",
"GetNumbers");
    var ExtractNumbersFromJson =
_kernel.Skills.GetFunction("OrchestratorPlugin",
"ExtractNumbersFromJson");
    ISKFunction MathFunction;

    // Retrieve the correct function based on the intent
    switch (intent)
    {
        case "Sqrt":
            MathFunction = _kernel.Skills.GetFunction("MathPlugin",
"Sqrt");
            break;
        case "Add":
            MathFunction = _kernel.Skills.GetFunction("MathPlugin",
"Add");
            break;
        default:
            return "I'm sorry, I don't understand.";
    }

    // Run the functions in a pipeline
    var output = await _kernel.RunAsync(
        request,
        GetNumbers,
        ExtractNumbersFromJson,
        MathFunction);

    return output["input"];
}

```

After making these changes, you should be able to run the code again and see the same results as before. Only now, the `RouteRequest` is easier to read and you've created a new native function that can be reused in other pipelines.

Starting a pipeline with additional context variables

So far, we've only passed in a string to the `RunAsync` method. However, you can also pass in a context object to start the pipeline with additional information. This can be useful to pass additional information to any of the functions in the pipeline.

It's *also* useful in persisting the initial `$input` variable across all functions in the pipeline without it being overwritten. For example, in our current pipeline, the user's original request is overwritten by the output of the `GetNumbers` function. This makes it difficult to

retrieve the original request later in the pipeline to create a natural sounding response. By storing the original request as another variable, we can retrieve it later in the pipeline.

Passing a context object to `RunAsync`

To pass a context object to `RunAsync`, you can create a new context object and pass it as the first parameter. This will start the pipeline with the variables in the context object. We'll be creating a new variable called `original_input` to store the original request. Later, we'll show where to add this code in the `RouteRequest` function.

C#

```
C#  
  
// Create a new context object  
var pipelineContext = new ContextVariables(request);  
pipelineContext["original_request"] = request;
```

Creating a semantic function that uses the new context variables

Now that we have a variable with the original request, we can use it to create a more natural sounding response. We'll create a new semantic function called `CreateResponse` that will use the `original_request` variable to create a response in the `OrchestratorPlugin`.

Start by creating a new folder called `CreateResponse` in your `OrchestratorPlugin` folder. Then create the `config.json` and `skprompt.txt` files and paste the following code into the `config.json` file. Notice how we now have two input variables, `input` and `original_request`.

JSON

```
{  
  "schema": 1,  
  "type": "completion",  
  "description": "Creates a response based on the original request and  
the output of the pipeline",  
  "completion": {  
    "max_tokens": 256,  
    "temperature": 0.0,
```

```
        "top_p": 0.0,
        "presence_penalty": 0.0,
        "frequency_penalty": 0.0
    },
    "input": {
        "parameters": [
            {
                "name": "input",
                "description": "The user's request.",
                "defaultValue": ""
            },
            {
                "name": "original_request",
                "description": "The original request from the user.",
                "defaultValue": ""
            }
        ]
    }
}
```

Next, copy and paste the following prompt into *skprompt.txt*.

```
txt
```

```
The answer to the users request is: {{$input}}
The bot should provide the answer back to the user.

User: {{$original_request}}
Bot:
```

You can now update the `RouteRequest` function to include the `CreateResponse` function in the pipeline. Update the `RouteRequest` function to the following:

```
C#
```

```
[SKFunction("Routes the request to the appropriate function.")]
public async Task<string> RouteRequest(SKContext context)
{
    // Save the original user request
    string request = context["input"];

    // Add the list of available functions to the context
    context["options"] = "Sqrt, Add";

    // Retrieve the intent from the user request
    var GetIntent = _kernel.Skills.GetFunction("OrchestratorPlugin",
    "GetIntent");
    var CreateResponse =
```

```

_kernel.Skills.GetFunction("OrchestratorPlugin", "CreateResponse");
    await GetIntent.InvokeAsync(context);
    string intent = context["input"].Trim();

    // Prepare the functions to be called in the pipeline
    var GetNumbers = _kernel.Skills.GetFunction("OrchestratorPlugin",
"GetNumbers");
    var ExtractNumbersFromJson =
_kernel.Skills.GetFunction("OrchestratorPlugin",
"ExtractNumbersFromJson");
    ISKFunction MathFunction;

    // Prepare the math function based on the intent
    switch (intent)
    {
        case "Sqrt":
            MathFunction = _kernel.Skills.GetFunction("MathPlugin",
"Sqrt");
            break;
        case "Add":
            MathFunction = _kernel.Skills.GetFunction("MathPlugin",
"Add");
            break;
        default:
            return "I'm sorry, I don't understand.";
    }

    // Create a new context object with the original request
    var pipelineContext = new ContextVariables(request);
    pipelineContext["original_request"] = request;

    // Run the functions in a pipeline
    var output = await _kernel.RunAsync(
        pipelineContext,
        GetNumbers,
        ExtractNumbersFromJson,
        MathFunction,
        CreateResponse);

    return output["input"];
}

```

Testing the new pipeline

Now that we've updated the pipeline, we can test it out. Run the following code in your main file.

C#

C#

```

using Microsoft.SemanticKernel;
using Plugins;

// ... instantiate your kernel

var pluginsDirectory =
Path.Combine(System.IO.Directory.GetCurrentDirectory(), "plugins");

// Import the semantic functions
kernel.ImportSemanticSkillFromDirectory(pluginsDirectory,
"OrchestratorPlugin");
kernel.ImportSemanticSkillFromDirectory(pluginsDirectory,
"SummarizeSkill");

// Import the native functions
var mathPlugin = kernel.ImportSkill(new MathPlugin(), "MathPlugin");
var orchestratorPlugin = kernel.ImportSkill(new
OrchestratorPlugin(kernel), "OrchestratorPlugin");

// Make a request that runs the Sqrt function
var result1 = await orchestratorPlugin["RouteRequest"]
    .InvokeAsync("What is the square root of 524?");
Console.WriteLine(result1);

// Make a request that runs the Add function
var result2 = await orchestratorPlugin["RouteRequest"]
    .InvokeAsync("How many sheep would I have if I started with 3 and
then got 7 more?");
Console.WriteLine(result2);

```

You should get a response like the following. Notice how the response is now more natural sounding.

Output

```

The square root of 524 is 22.891046284519195.
You would have 10 sheep.

```

Take the next step

You are now becoming familiar with orchestrating both semantic and non-semantic functions. Up until now, however, you've had to manually orchestrate the functions. In the next section, you'll learn how to use planner to orchestrate functions automatically.

[Automatically create chains with planner](#)

Automatically orchestrate AI with planner

Article • 07/12/2023



So far, we have manually orchestrated all of the functions on behalf of the user. This, however, is not a scalable solution because it would require the app developer to predict all possible requests that could be made by the user. So instead, we will learn how to automatically orchestrate functions on the fly using planner. If you want to see the final solution, you can check out the following samples in the public documentation repository.

Language	Link to final solution
C#	Open solution in GitHub ↗
Python	Open solution in GitHub ↗

What is planner?

Planner is a function that takes a user's ask and returns back a plan on how to accomplish the request. It does so by using AI to mix-and-match the plugins registered in the kernel so that it can recombine them into a series of steps that complete a goal.

This is a powerful concept because it allows you to create atomic functions that can be used in ways that you as a developer may not have thought of.

For example, if you had task and calendar event plugins, planner could combine them to create workflows like "remind me to buy milk when I go to the store" or "remind me to call my mom tomorrow" without you explicitly having to write code for those scenarios.



With great power comes great responsibility, however. Because planner can combine functions in ways that you may not have thought of, it is important to make sure that you only expose functions that you want to be used in this way. It's also important to make sure that you apply [responsible AI ↗](#) principles to your functions so that they are used in a way that is fair, reliable, safe, private, and secure.

Planner is an extensible part of Semantic Kernel. This means we have several planners to choose from and that you could create a custom planner if you had specific needs. Below is a table of the out-of-the-box planners provided by Semantic Kernel and their language support. The ✘ symbol indicates that the feature is not yet available in that language; if you would like to see a feature implemented in a language, please consider [contributing to the project](#) or [opening an issue](#).

Planner	Description	C#	Python	Java
BasicPlanner	A simplified version of SequentialPlanner that strings together a set of functions.	✘	✓	✘
ActionPlanner	Creates a plan with a single step.	✓	✗	✗
SequentialPlanner	Creates a plan with a series of steps that are interconnected with custom generated input and output variables.	✓	✗	✗

Testing out planner

For the purposes of this article, we'll build upon the same code we wrote in the [previous section](#). Only this time, instead of relying on our own `OrchestratorPlugin` to chain the `MathPlugin` functions, we'll use planner to do it for us!

At the end of this section, we'll have built a natural language calculator that can answer simple word problems for users.

Adding more functions to `MathPlugin`

Before we use planner, let's add a few more functions to our `MathPlugin` class so we can have more options for our planner to choose from. The following code adds a `Subtract`, `Multiply`, and `Divide` function to our plugin.

```
C#  
  
C#  
  
using Microsoft.SemanticKernel.Orchestration;  
using Microsoft.SemanticKernel.SkillDefinition;  
  
namespace Plugins;  
  
public class MathPlugin  
{  
    [SKFunction("Take the square root of a number")]  
    public string Sqrt(string number)  
    {  
        return Math.Sqrt(Convert.ToDouble(number)).ToString();  
    }  
  
    [SKFunction("Add two numbers")]  
    [SKFunctionContextParameter(Name = "input", Description = "The first  
number to add")]  
    [SKFunctionContextParameter(Name = "number2", Description = "The  
second number to add")]  
    public string Add(SKContext context)  
    {  
        return (Convert.ToDouble(context["input"]) +  
Convert.ToDouble(context["number2"])).ToString();  
    }  
  
    [SKFunction("Subtract two numbers")]  
    [SKFunctionContextParameter(Name = "input", Description = "The first  
number to subtract from")]  
    [SKFunctionContextParameter(Name = "number2", Description = "The  
second number to subtract away")]  
    public string Subtract(SKContext context)  
    {  
        return (Convert.ToDouble(context["input"]) -
```

```

Convert.ToDouble(context["number2"])).ToString();
}

[SKFunction("Multiply two numbers. When increasing by a percentage,
don't forget to add 1 to the percentage.")]
[SKFunctionContextParameter(Name = "input", Description = "The first
number to multiply")]
[SKFunctionContextParameter(Name = "number2", Description = "The
second number to multiply")]
public string Multiply(SKContext context)
{
    return (Convert.ToDouble(context["input"]) *
Convert.ToDouble(context["number2"])).ToString();
}

[SKFunction("Divide two numbers")]
[SKFunctionContextParameter(Name = "input", Description = "The first
number to divide from")]
[SKFunctionContextParameter(Name = "number2", Description = "The
second number to divide by")]
public string Divide(SKContext context)
{
    return (Convert.ToDouble(context["input"]) /
Convert.ToDouble(context["number2"])).ToString();
}

```

Instantiating planner

To instantiate planner, all you need to do is pass it a kernel object. Planner will then automatically discover all of the plugins registered in the kernel and use them to create plans. The following code initializes both a kernel and a `SequentialPlanner`. At the end of this article we'll review the other types of Planners that are available in Semantic Kernel.

C#

```

C#

using Microsoft.SemanticKernel;
using Plugins;

// ... instantiate your kernel

// Add the math plugin
var mathPlugin = kernel.ImportSkill(new MathPlugin(), "MathPlugin");

```

```
// Create planner
var planner = new SequentialPlanner(kernel);
```

Creating and running a plan

Now that we have planner, we can use it to create a plan for a user's ask and then invoke the plan to get a result. The following code asks our planner to solve a math problem that is difficult for an LLM to solve on its own because it requires multiple steps and it has numbers with decimal points.

C#

```
C#  
  
// Create a plan for the ask
var ask = "If my investment of 2130.23 dollars increased by 23%, how
much would I have after I spent $5 on a latte?";
var plan = await planner.CreatePlanAsync(ask);

// Execute the plan
var result = await plan.InvokeAsync();

Console.WriteLine("Plan results:");
Console.WriteLine(result.Result);
```

After running this code, you should get the correct answer of 2615.1829 back, but how?

How does planner work?

Behind the scenes, planner uses an LLM prompt to generate a plan. You can see the prompt that is used by `SequentialPlanner` by navigating to the [skprompt.txt file](#) in the Semantic Kernel repository. You can also view the prompt used by the [basic planner](#) in Python.

Understanding the prompt powering planner

The first few lines of the prompt are the most important to understanding how planner works. They look like this:

txt

Create an XML plan step by step, to satisfy the goal given.

To create a plan, follow these steps:

0. The plan should be as short as possible.
1. From a <goal> create a <plan> as a series of <functions>.
2. Before using any function in a plan, check that it is present in the most recent [AVAILABLE FUNCTIONS] list. If it is not, do not use it. Do not assume that any function that was previously defined or used in another plan or in [EXAMPLES] is automatically available or compatible with the current plan.
3. Only use functions that are required for the given goal.
4. A function has a single 'input' and a single 'output' which are both strings and not objects.
5. The 'output' from each function is automatically passed as 'input' to the subsequent <function>.
6. 'input' does not need to be specified if it consumes the 'output' of the previous function.
7. To save an 'output' from a <function>, to pass into a future <function>, use <function.{FunctionName} ... setContextVariable: "
<UNIQUE_VARIABLE_KEY>"/>
8. To save an 'output' from a <function>, to return as part of a plan result, use <function.{FunctionName} ... appendToResult:
"RESULT__<UNIQUE_RESULT_KEY>"/>
9. Append an "END" XML comment at the end of the plan.

With these steps, planner is given a set of rules that it can use to generate a plan in XML. Afterwards, the prompt provides a few examples of valid plans before finally providing the `$available_functions` and user's goal.

```
txt

[AVAILABLE FUNCTIONS]

{{\$available_functions}}

[END AVAILABLE FUNCTIONS]

<goal>{{\$input}}</goal>
```

Giving planner the *best* data

When you render the prompt, one of the main things you might notice is that all of the descriptions we provided for our functions are included in the prompt. For example, the description for `MathPlugin.Add` is included in the prompt as `Add two numbers`.

```
txt
```

[AVAILABLE FUNCTIONS]

```
MathPlugin.Add:  
  description: Add two numbers  
  inputs:  
    - input: The first number to add  
    - number2: The second number to add  
  
MathPlugin.Divide:  
  description: Divide two numbers  
  inputs:  
    - input: The first number to divide from  
    - number2: The second number to divide by
```

Because of this, it's incredibly important to provide the best descriptions you can for your functions. If you don't, planner will not be able to generate a plan that uses your functions correctly.

You can also use the descriptions to provide explicit instructions to the model on how to use your functions. Below are some techniques you can use to improve the use of your functions by planner.

- **Provide help text** – It's not always clear when or how to use a function, so giving advice helps. For example, the description for `MathPlugin.Multiply` reminds the bot to add 1 whenever it increases a number by a percentage.
- **Describe the output**. – While there is not an explicit way to tell planner what the output of a function is, you can describe the output in the description.
- **State if inputs are required**. – If a function requires an input, you can state that in the input's description so the model knows to provide an input. Conversely, you can tell the model that an input is optional so it knows it can skip it if necessary.

Viewing the plan produced by planner

Because the plan is returned as plain text (either as XML or JSON), we can print the results to inspect what plan planner actually created. The following code shows how to print the plan and the output for C# and Python.

C#

```
C#
```

```
Console.WriteLine(plan);
```

Output

```
{
  "state": [
    {
      "Key": "INPUT",
      "Value": ""
    }
  ],
  "steps": [
    {
      "state": [
        {
          "Key": "INPUT",
          "Value": ""
        }
      ],
      "steps": [],
      "parameters": [
        {
          "Key": "number2",
          "Value": "1.23"
        },
        {
          "Key": "INPUT",
          "Value": "2130.23"
        }
      ],
      "outputs": [
        "INVESTMENT_INCREASE"
      ],
      "next_step_index": 0,
      "name": "Multiply",
      "skill_name": "MathPlugin",
      "description": "Multiply two numbers"
    },
    {
      "state": [
        {
          "Key": "INPUT",
          "Value": ""
        }
      ],
      "steps": [],
      "parameters": [
        {
          "Key": "number2",
          "Value": "5"
        },
        {
          "Key": "INPUT",
          "Value": "$INVESTMENT_INCREASE"
        }
      ],
    }
  ]
}
```

```

    "outputs": [
        "RESULT_FINAL_AMOUNT"
    ],
    "next_step_index": 0,
    "name": "Subtract",
    "skill_name": "MathPlugin",
    "description": "Subtract two numbers"
}
],
"parameters": [
{
    "Key": "INPUT",
    "Value": ""
}
],
"outputs": [
    "RESULT_FINAL_AMOUNT"
],
"next_step_index": 0,
"name": "",
"skill_name": "Microsoft.SemanticKernel.Planning.Plan",
"description": "If my investment of 2130.23 dollars increased by 23%, how much would I have after I spent $5 on a latte?"
}

```

Notice how in the example, planner can string together functions and pass parameters to them. This effectively allows us to deprecate the `OrchestratorPlugin` we created previously because we no longer need the `RouteRequest` native function or the `GetNumbers` semantic function. Planner does both.

When to use planner?

As demonstrated by this example, planner is extremely powerful because it can automatically recombine functions you have already defined, and as AI models improve and as the community develops better planners, you will be able to rely on them to achieve increasingly more sophisticated user scenarios.

There are, however, considerations you should make before using a planner. The following table describes the top considerations you should make along with mitigations you can take to reduce their impact.

Considerations	Description	Mitigation
Performance	It takes time for a planner to consume the full list of tokens and to generate a plan for a user, if you rely on the planner after a user provides input,	While building UI, it's important to provide feedback to the user to let them know something is happening with loading experiences. You can also

Considerations	Description	Mitigation
	you may unintentionally hang the UI while waiting for a plan.	use LLMs to stall for time by generating an initial response for the user while the planner completes a plan. Lastly, you can use predefined plans for common scenarios to avoid waiting for a new plan.
Cost	both the prompt and generated plan consume many tokens. To generate a very complex plan, you may need to consume <i>all</i> of the tokens provided by a model. This can result in high costs for your service if you're not careful, especially since planning typically requires more advanced models like GPT 3.5 or GPT 4.	The more atomic your functions are, the more tokens you'll require. By authoring higher order functions, you can provide planner with fewer functions that use fewer tokens. Lastly, you can use predefined plans for common scenarios to avoid spending money on new plans.
Correctness	Planner can generate faulty plans. For example, it may pass variables incorrectly, return malformed schema, or perform steps that don't make sense.	To make planner robust, you should provide error handling. Some errors, like malformed schema or improperly returned schema, can be recovered by asking planner to "fix" the plan.

Using predefined plans

There are likely common scenarios that your users will frequently ask for. To avoid the performance hit and the costs associated with planner, you can pre-create plans and serve them up to a user.

This is similar to the front-end development adage coined by Aaron Swartz: "[Bake, don't fry](#)." By pre-creating, or "baking," your plans, you can avoid generating them on the fly (i.e., "frying"). You won't be able to get rid of "frying" entirely when creating AI apps, but you can reduce your reliance on it so you can use healthier alternatives instead.

To achieve this, you can generate plans for common scenarios offline, and store them as XML in your project. Based on the intent of the user, you can then serve the plan back up so it can be executed. By "baking" your plans, you also have the opportunity to create additional optimizations to improve speed or lower costs.

Next steps

You now have the skills necessary to automatically generate plans for your users. You can use these skills to create more advanced AI apps that can handle increasingly

complex scenarios. In the next section, you'll learn how to author plugins that can be used by planner *and* ChatGPT.

[Create and run ChatGPT plugins](#)

Create and run ChatGPT plugins using Semantic Kernel

Article • 07/14/2023



In this article, we'll show you how to take a Semantic Kernel plugin and expose it to ChatGPT with Azure Functions. As an example, we'll demonstrate how to transform the `MathPlugin` we created in previous articles into a ChatGPT plugin.

At the [end of this article](#), you'll also learn how to load a ChatGPT plugin into Semantic Kernel and use it with a planner.

Once we're done, you'll have an Azure Function that exposes each of your plugin's native functions as HTTP endpoints so they can be used by Semantic Kernel or ChatGPT. If you want to see the final solution, you can check out the sample in the public documentation repository.

Language	Link to final solution
C#	Open solution in GitHub ↗
Python	<i>Coming soon</i>

Prerequisites

To complete this tutorial, you'll need the following:

- [Azure Functions Core Tools ↗](#) version 4.x.
- [.NET 6.0 SDK. ↗](#)

To publish your plugin once you're complete, you'll also need an Azure account with an active subscription. [Create an account for free ↗](#) and one of the following tools for creating Azure resources:

- [Azure CLI version 2.4](#) or later.
- The [Azure Az PowerShell module](#) version 5.9.0 or later.

You do **not** need to have access to OpenAI's plugin preview to complete this tutorial. If you do have access, however, you can upload your final plugin to OpenAI and use it in ChatGPT at the very end.

What are ChatGPT plugins?

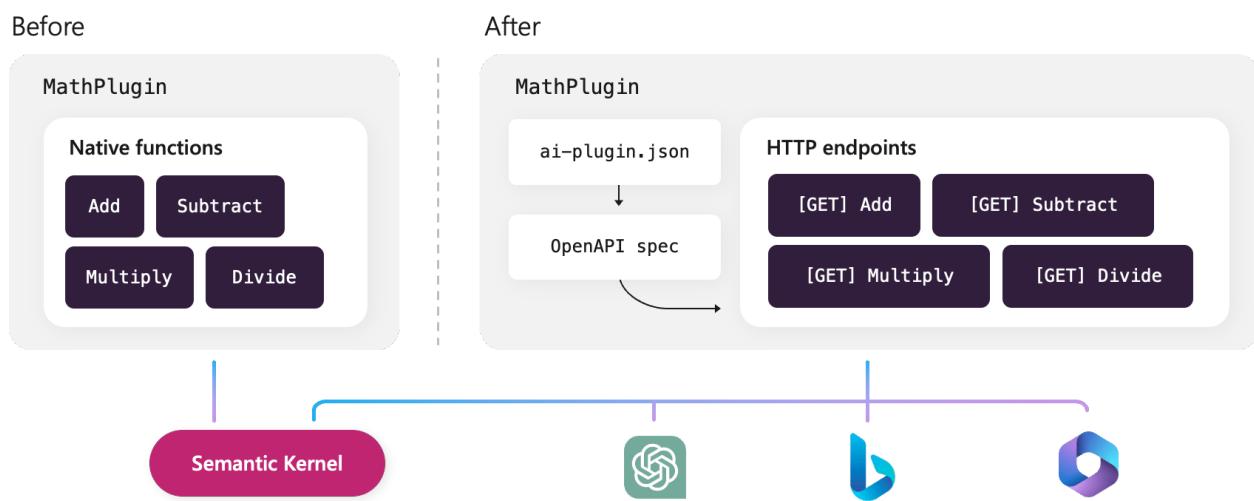
In the [plugin article](#) we described how all plugins are moving towards the common standard defined by OpenAI. This standard, which is called a ChatGPT plugin in this article, uses a plugin manifest file that points to an accompanying [OpenAPI specification](#). Plugins defined in this way can then be used by any application that supports the OpenAI specification, including Semantic Kernel and ChatGPT.

ⓘ Important

OpenAPI is different than OpenAI. OpenAPI is a specification for describing REST APIs, while OpenAI is a company that develops AI models and APIs. While the two are not related, OpenAI has adopted the OpenAPI specification for describing plugin APIs.

Transforming our `MathPlugin` into a ChatGPT plugin

So far, however, we've only shown how to create plugins that are *natively* loaded into Semantic Kernel instead of being exposed through an OpenAPI specification. This has helped us demonstrate the core concepts of plugins without adding the additional complexity of standing up an HTTP endpoint. With minimal changes, however, we can take the plugins we've already created and expose them to ChatGPT.



There are three steps we must take to turn our existing `MathPlugin` into a ChatGPT plugin:

1. Create HTTP endpoints for each native function.
2. Create an OpenAPI specification and plugin manifest file that describes our plugin.
3. Test the plugin in either Semantic Kernel or ChatGPT.

Provide HTTP endpoints for each function

Before we can expose our plugin to other applications, we need to create an HTTP endpoint for each of our native functions. This will allow us to call our native functions from any other service. You can achieve this multiple ways, but in this article we'll use Azure Functions.

Create a new Azure Function project

There are several ways to create an Azure Function, but in this article we'll use the Azure Functions Core Tools. If you are using Visual Studio or Visual Studio Code, you can also use the Azure Functions extension to create a new Azure Function project. For more information, see [Create your first function using Visual Studio](#) or [Create your first function using Visual Studio Code](#).

1. To create a new Azure Function project, run the following command in your terminal:

```
Bash
```

```
func init MathPlugin --worker-runtime dotnet-isolated --target-framework net6.0
```

2. Navigate into the `MathPlugin` directory:

```
Bash
```

```
cd MathPlugin
```

3. Open the `MathPlugin.csproj` file.

4. Update the `MathPlugin.csproj` file to include the following package references:

```
XML
```

```
<PackageReference Include="Microsoft.Azure.Functions.Worker" Version="1.17.0" />
<PackageReference Include="Microsoft.Azure.Functions.Worker.Extensions.Http" Version="3.0.13" />
<PackageReference Include="Microsoft.Azure.Functions.Worker.Sdk" Version="1.11.0" />
```

5. Open the `host.json` file.

6. Update the `host.json` file to include the following configuration:

JSON

```
{  
    "version": "2.0",  
    "logging": {  
        "applicationInsights": {  
            "samplingSettings": {  
                "isEnabled": true,  
                "excludedTypes": "Request"  
            }  
        }  
    },  
    "extensions": {  
        "http": {  
            "routePrefix": ""  
        }  
    }  
}
```

ⓘ Note

Later, this will allow us to serve up the plugin manifest file from the `.well-known` directory as required by the OpenAI specification.

7. Run the following command to restore the packages:

Bash

```
dotnet restore
```

Add the native functions to the Azure Function project

We can now add our native functions to the Azure Function project.

1. Run the following command in your terminal to create placeholder for the `Add` function:

Bash

```
func new --name Add --template "HTTP trigger" --authlevel "anonymous"
```

2. Open the `Add.cs` file.

3. Replace the `Run` function with the following code:

```
C#  
  
[Function("Add")]
public HttpResponseMessage Run([HttpTrigger(AuthorizationLevel.Anonymous,
"get", "post")] HttpRequestData req)
{
    bool result1 = double.TryParse(req.Query["number1"], out double
number1);
    bool result2 = double.TryParse(req.Query["number2"], out double
number2);

    if (result1 && result2)
    {
        HttpResponseMessage response =
req.CreateResponse(HttpStatusCode.OK);
        response.Headers.Add("Content-Type", "text/plain");
        double sum = number1 + number2;
        response.WriteString(sum.ToString());

        return response;
    }
    else
    {
        HttpResponseMessage response =
req.CreateResponse(HttpStatusCode.BadRequest);
        response.Headers.Add("Content-Type", "application/json");
        response.WriteString("Please pass two numbers on the query
string or in the request body");

        return response;
    }
}
```

4. Repeat the previous steps to create HTTP endpoints for the `Subtract`, `Multiply`, `Divide`, and `Sqrt` functions. When replacing the `Run` function, be sure to update the function name and logic for each function accordingly.

Validate the HTTP endpoints

At this point, you should have five HTTP endpoints in your Azure Function project. You can test them by following these steps:

1. Run the following command in your terminal:

```
Bash
```

```
func start
```

2. Open a new terminal window and run the following commands:

Bash

```
curl "http://localhost:7071/Add?number1=1&number2=2"
curl "http://localhost:7071/Subtract?number1=1&number2=2"
curl "http://localhost:7071/Multiply?number1=1&number2=2"
curl "http://localhost:7071/Divide?number1=1&number2=2"
curl "http://localhost:7071/Sqrt?number=9"
```

3. You should see the following responses:

Output

```
3
-1
2
0.5
3
```

Create the manifest files

Now that we have HTTP endpoints for each of our native functions, we need to create the files that will tell ChatGPT and other applications how to call them. We'll do this by creating an OpenAPI specification and plugin manifest file.

Add an OpenAPI spec to your Azure Function project

An OpenAPI specification describes the HTTP endpoints that are available in your plugin. Instead of manually creating an OpenAPI specification, you can use NuGet packages provided by Azure Functions to automatically create and serve up these files. To add an OpenAPI specification to your Azure Function project, follow these steps:

1. Run the following commands in your Azure Function project directory:

Bash

```
dotnet add package Microsoft.Azure.WebJobs.Extensions.OpenApi --version 1.5.1
dotnet add package Microsoft.Azure.Functions.Worker.Extensions.OpenApi
--version 1.5.1
```

2. Open the `Add.cs` file.

3. Add the following `using` statements:

C#

```
using Microsoft.Azure.WebJobs.Extensions.OpenApi.Core.Attributes;
using Microsoft.OpenApi.Models;
```

4. Add the following attributes to the `Run` function:

C#

```
[OpenApiOperation(operationId: "Add", tags: new[] { "ExecuteFunction" }, Description = "Adds two numbers.")]
[OpenApiParameter(name: "number1", Description = "The first number to add", Required = true, In = ParameterLocation.Query)]
[OpenApiParameter(name: "number2", Description = "The second number to add", Required = true, In = ParameterLocation.Query)]
[OpenApiResponseWithBody(statusCode: HttpStatusCode.OK, contentType: "text/plain", bodyType: typeof(string), Description = "Returns the sum of the two numbers.")]
[OpenApiResponseWithBody(statusCode: HttpStatusCode.BadRequest, contentType: "application/json", bodyType: typeof(string), Description = "Returns the error of the input.")]
```

5. Repeat the previous steps for the `Subtract`, `Multiply`, `Divide`, and `Sqrt` functions.

When adding the attributes, update the operation and parameter descriptions accordingly.

ⓘ Important

The `Description` fields for both the operation and the parameters are the most important attributes because they will be used by the planner to determine which function to call. We recommend reusing the same description values from the previous walkthroughs.

Function	Description	Number 1	Number 2
Add	Add two numbers.	The first number to add	The second number to add
Subtract	Subtract two numbers.	The first number to subtract from	The second number to subtract away

Function	Description	Number 1	Number 2
Multiply	Multiply two numbers. When increasing by a percentage, don't forget to add 1 to the percentage.		
Divide	Divide two numbers.	The first number to divide from	The second number to divide by
Sqrt	Take the square root of a number.	The number to calculate the square root of	N/A

Validate the OpenAPI spec

You can then test the OpenAPI document by following these steps:

1. Run the following command in your terminal:

```
Bash
```

```
func start
```

2. Navigate to the following URL in your browser:

```
Bash
```

```
http://localhost:7071/swagger/ui
```

3. You should see the following page:

The screenshot shows the Swagger UI interface for an Azure Functions API. At the top, it displays the URL `http://localhost:7071/swagger.json` and a green "Explore" button. Below this, the title "OpenAPI Document on Azure Functions 1.0.0" is shown, along with the base URL `localhost:7071` and the JSON file path `http://localhost:7071/swagger.json`. A sub-header states "This is the OpenAPI Document on Azure Functions". A dropdown menu labeled "Schemes" is set to "HTTP". The main content area is titled "ExecuteFunction" and lists five API endpoints under the "GET" method:

- /Multiply
- /Add
- /Divide
- /Sqrt
- /Subtract

4. Navigating to `http://localhost:7071/swagger.json` will allow you to download the OpenAPI specification.

Add the plugin manifest file

The last step is to serve up the plugin manifest file. Based on the OpenAI specification, the manifest file is always served up from the `./well-known/ai-plugin.json` file and contains the following information:

Field	Type	Description	Required
schema_version	String	Manifest schema version	<input checked="" type="checkbox"/>
name_for_model	String	Name the model will use to target the plugin (no spaces allowed, only letters and numbers). 50 character max.	<input checked="" type="checkbox"/>
name_for_human	String	Human-readable name, such as the full company name. 20 character max.	<input checked="" type="checkbox"/>
description_for_model	String	Description better tailored to the model, such as token context length considerations or keyword	<input checked="" type="checkbox"/>

Field	Type	Description	Required
		usage for improved plugin prompting. 8,000 character max.	
description_for_human	String	Human-readable description of the plugin. 100 character max.	<input checked="" type="checkbox"/>
auth	ManifestAuth	Authentication schema	<input checked="" type="checkbox"/>
api	Object	API specification	<input checked="" type="checkbox"/>
logo_url	String	URL used to fetch the logo. Suggested size: 512 x 512. Transparent backgrounds are supported. Must be an image, no GIFs are allowed.	<input checked="" type="checkbox"/>
contact_email	String	Email contact for safety/moderation	<input checked="" type="checkbox"/>
legal_info_url	String	Redirect URL for users to view plugin information	<input checked="" type="checkbox"/>
HttpAuthorizationType	HttpAuthorizationType	"bearer" or "basic"	<input checked="" type="checkbox"/>
ManifestAuthType	ManifestAuthType	"none", "user_http", "service_http", or "oauth"	
interface BaseManifestAuth	BaseManifestAuth	type: ManifestAuthType; instructions: string;	
ManifestNoAuth	ManifestNoAuth	No authentication required: BaseManifestAuth & { type: 'none', }	
ManifestAuth	ManifestAuth	ManifestNoAuth, ManifestServiceHttpAuth, ManifestUserHttpAuth, ManifestOAuthAuth	

To create an Azure Function that serves up this manifest, follow these steps:

1. Run the following command in your terminal:

Bash

```
func new --name AIPluginJson --template "HTTP trigger" --authlevel "anonymous"
```

2. Open the *AIPluginJson.cs* file.

3. Replace the `Run` method with the following code:

```
C#  
  
[Function("GetAiPluginJson")]
public HttpResponseMessage Run([HttpTrigger(AuthorizationLevel.Anonymous,
    "get", Route = ".well-known/ai-plugin.json")] HttpRequestData req)
{
    var currentDomain = $"{req.Url.Scheme}://{req.Url.Host}:{req.Url.Port}";  
  
    HttpResponseMessage response = req.CreateResponse(HttpStatusCode.OK);
    response.Headers.Add("Content-Type", "application/json");  
  
    var json = $@"
        ""schema_version"": ""v1"",
        ""name_for_human"": ""Simple calculator"",
        ""name_for_model"": ""calculator"",
        ""description_for_human"": ""This plugin performs basic math
operations"",
        ""description_for_model"": ""Help the user perform math. You can
add, subtract, multiple, divide, and perform square roots"",
        ""auth"": {},
        ""type"": ""none""
    },
    ""api"": {
        ""type"": ""openapi"",
        ""url"": "{currentDomain}/swagger.json"
    },
    ""logo_url"": "{currentDomain}/logo.png",
    ""contact_email"": "support@example.com",
    ""legal_info_url"": "http://www.example.com/legal"
};  
  
    response.WriteString(json);  
  
    return response;
}
```

Validate the plugin manifest file

You can then test that the plugin manifest file is being served up by following these steps:

1. Run the following command in your terminal:

```
Bash
```

```
func start
```

2. Navigate to the following URL in your browser:

Bash

```
http://localhost:7071/.well-known/ai-plugin.json
```

3. You should now see the plugin manifest file.

```
{  
    "schema_version": "v1",  
    "name_for_human": "Simple calculator",  
    "name_for_model": "calculator",  
    "description_for_human": "This plugin performs basic math operations.",  
    "description_for_model": "Help the user perform math. You can add, subtract,  
multiple, divide, and perform square roots.",  
    "auth": {  
        "type": "none"  
    },  
    "api": {  
        "type": "openapi",  
        "url": "http://localhost:7071/api/swagger.json"  
    },  
    "logo_url": "http://localhost:7071/logo.png",  
    "contact_email": "support@example.com",  
    "legal_info_url": "http://www.example.com/legal"  
}
```

Testing the plugin end-to-end

You now have a complete plugin that can be used in Semantic Kernel and ChatGPT. Since there is currently a waitlist for creating plugins for ChatGPT, we'll first demonstrate how you can test your plugin with Semantic Kernel.

Running the plugin with Semantic Kernel

By testing your plugin in Semantic Kernel, you can ensure that it is working as expected before you get access to the plugin developer portal for ChatGPT. While testing in Semantic Kernel, we recommend using the Stepwise Planner to invoke your plugin since it is the only planner that supports JSON responses.

To test the plugin in Semantic Kernel, follow these steps:

1. Create a new C# project.
2. Add the necessary Semantic Kernel NuGet packages:

Bash

```
dotnet add package Microsoft.SemanticKernel  
dotnet add package Microsoft.SemanticKernel.Planning.StepwisePlanner
```

```
dotnet add package Microsoft.SemanticKernel.Skills.OpenAPI
```

3. Paste the following code into your *program.cs* file:

C#

```
using Microsoft.Extensions.Logging;
using Microsoft.SemanticKernel;
using Microsoft.SemanticKernel.Planning;

// ... create a new Semantic Kernel instance here

// Add the math plugin using the plugin manifest URL
const string pluginManifestUrl = "http://localhost:7071/.well-known/ai-
plugin.json";
var mathPlugin = await
kernel.ImportChatGptPluginSkillFromUrlAsync("MathPlugin", new
Uri(pluginManifestUrl));

// Create a stepwise planner and invoke it
var planner = new StepwisePlanner(kernel);
var question = "I have $2130.23. How much would I have after it grew by
24% and after I spent $5 on a latte?";
var plan = planner.CreatePlan(question);
var result = await plan.InvokeAsync(kernel.CreateNewContext());

// Print the results
Console.WriteLine("Result: " + result);

// Print details about the plan
if (result.Variables.TryGetValue("stepCount", out string? stepCount))
{
    Console.WriteLine("Steps Taken: " + stepCount);
}
if (result.Variables.TryGetValue("skillCount", out string? skillCount))
{
    Console.WriteLine("Skills Used: " + skillCount);
}
```

4. After running the code, you should see the following output:

Output

```
Result: After the amount grew by 24% and $5 was spent on a latte, you
would have $2636.4852 remaining.
Steps Taken: 3
Skills Used: 2 (MathPlugin.Multiply(1), MathPlugin.Subtract(1))
```

Running the plugin in ChatGPT

If you would like to test your plugin in ChatGPT, you can do so by following these steps:

1. Request access to plugin development by filling out the [waitlist form](#).
2. Once you have access, follow the steps [provided by OpenAI](#) to register your plugin.

Next steps

Congratulations! You have successfully created a plugin that can be used in Semantic Kernel and ChatGPT. Once you have fully tested your plugin, you can deploy it to Azure Functions and register it with OpenAI. For more information, see the following resources:

- [Deploying Azure Functions](#)
- [Submit a plugin to the OpenAI plugin store](#)

Use the out-of-the-box plugins in the kernel

Article • 07/12/2023



ⓘ Note

Skills are currently being renamed to plugins. This article has been updated to reflect the latest terminology, but some images and code samples may still refer to skills.

To provide a degree of standardization across Semantic Kernel implementations, the GitHub repo has several plugins available out-of-the-box depending on the language you are using. These plugins are often referred to as **Core plugins**. Additionally, each library also includes a handful of other plugins that you can use. The following section covers each set of plugins in more detail.

Core plugins

The core plugins are planned to be available in all languages since they are core to using Semantic Kernel. Below are the core plugins currently available in Semantic Kernel along with their current support for each language. The ✘ symbol indicates that the feature is not yet available in that language; if you would like to see a feature implemented in a language, please consider [contributing to the project](#) or [opening an issue](#).

Plugin	Description	C#	Python	Java
ConversationSummarySkill	To summarize a conversation	✓	✓	*
FileIOSkill	To read and write to the filesystem	✓	✓	✗
HttpSkill	To call APIs	✓	✓	✗
MathSkill	To perform mathematical operations	✓	✓	✗
TextMemorySkill	To store and retrieve text in memory	✓	✓	✗
TextSkill	To deterministically manipulating text strings	✓	✓	*

Plugin	Description	C#	Python	Java
TimeSkill	To acquire the time of day and any other temporal information	✓	✓	*
WaitSkill	To pause execution for a specified amount of time	✓	✗	✗

You can find the full list of core plugins for each language by following the links below:

- [C# core plugins ↗](#)
- [Python core plugins ↗](#)

Using core plugins in Semantic Kernel

If you want to use one of the core plugins, you can easily import them into your project. For example, if you want to use the `TimeSkill` in either C# or Python, you can import it as follows.

C#

When using a core plugin, be sure to include a `using Microsoft.SemanticKernel.CoreSkills`:

```
C#
using Microsoft.SemanticKernel.CoreSkills;
// ... instantiate a kernel and configure it first
kernel.ImportSkill(new TimeSkill(), "time");

const string ThePromptTemplate = @"
Today is: {{time.Date}}
Current time is: {{time.Time}}
```

Answer to the following questions using JSON syntax, including the data used.

Is it morning, afternoon, evening, or night
(morning/afternoon/evening/night)?
Is it weekend time (weekend/not weekend)?;

```
var myKindOfDay = kernel.CreateSemanticFunction(ThePromptTemplate,
maxTokens: 150);
```

```
var myOutput = await myKindOfDay.InvokeAsync();
Console.WriteLine(myOutput);
```

The output should be similar to the following:

resulting-output

```
{
    "date": "Wednesday, 21 June, 2023",
    "time": "12:17:02 AM",
    "period": "night",
    "weekend": "not weekend"
}
```

Chaining core plugins together in Semantic Kernel

Most of the core plugins were built so that they can be easily chained together. For example, the `TextSkill` can be used to trim whitespace from a string, convert it to uppercase, and then convert it to lowercase.

C#

```
C#  
  
using Microsoft.SemanticKernel;
using Microsoft.SemanticKernel.Orchestration;
using Microsoft.SemanticKernel.CoreSkills;  
  
var kernel = Kernel.Builder.Build();  
  
var myText = kernel.ImportSkill(new TextSkill());  
  
SKContext myOutput = await kernel.RunAsync(
    " i n f i n i t e   s p a c e   ",
    myText["TrimStart"],
    myText["TrimEnd"],
    myText["Uppercase"]);  
  
Console.WriteLine(myOutput);
```

Note how the input streams through a pipeline of three functions executed serially. Expressed sequentially as in a chain of functions:

" i n f i n i t e s p a c e " **TextSkill.TrimStart**

→

→

TextSkill.TrimEnd

→

TextSkill.Uppercase

→

The output reads as:

I N F I N I T E S P A C E

Take the next step

[Deploy your plugins to Azure](#)

What are Prompts?

Article • 05/23/2023



Prompts play a crucial role in communicating and directing the behavior of Large Language Models (LLMs) AI. They serve as inputs or queries that users can provide to elicit specific responses from a model.

The subtleties of prompting

Effective prompt design is essential to achieving desired outcomes with LLM AI models. Prompt engineering, also known as prompt design, is an emerging field that requires creativity and attention to detail. It involves selecting the right words, phrases, symbols, and formats that guide the model in generating high-quality and relevant texts.

If you've already experimented with ChatGPT, you can see how the model's behavior changes dramatically based on the inputs you provide. For example, the following prompts produce very different outputs:

Prompt

Please give me the history of humans.

Prompt

Please give me the history of humans in 3 sentences.

The first prompt produces a long report, while the second prompt produces a concise response. If you were building a UI with limited space, the second prompt would be more suitable for your needs. Further refined behavior can be achieved by adding even more details to the prompt, but it's possible to go too far and produce irrelevant outputs. As a prompt engineer, you must find the right balance between specificity and relevance.

When you work directly with LLM models, you can also use other controls to influence the model's behavior. For example, you can use the `temperature` parameter to control the randomness of the model's output. Other parameters like top-k, top-p, frequency penalty, and presence penalty also influence the model's behavior.

Prompt engineering: a new career

Because of the amount of control that exists, prompt engineering is a critical skill for anyone working with LLM AI models. It's also a skill that's in high demand as more organizations adopt LLM AI models to automate tasks and improve productivity. A good prompt engineer can help organizations get the most out of their LLM AI models by designing prompts that produce the desired outputs.

Becoming a great prompt engineer with Semantic Kernel

Semantic Kernel is a valuable tool for prompt engineering because it allows you to experiment with different prompts and parameters across multiple different models using a common interface. This allows you to quickly compare the outputs of different models and parameters, and iterate on prompts to achieve the desired results.

Once you've become familiar with prompt engineering, you can also use Semantic Kernel to apply your skills to real-world scenarios. By combining your prompts with native functions and connectors, you can build powerful AI-powered applications.

Lastly, by deeply integrating with Visual Studio Code, Semantic Kernel also makes it easy for you to integrate prompt engineering into your existing development processes.

- ✓ Create prompts directly in your preferred code editor.
- ✓ Write tests for them using your existing testing frameworks.
- ✓ And deploy them to production using your existing CI/CD pipelines.

Additional tips for prompt engineering

Becoming a skilled prompt engineer requires a combination of technical knowledge, creativity, and experimentation. Here are some tips to excel in prompt engineering:

- **Understand LLM AI models:** Gain a deep understanding of how LLM AI models work, including their architecture, training processes, and behavior.
- **Domain knowledge:** Acquire domain-specific knowledge to design prompts that align with the desired outputs and tasks.
- **Experimentation:** Explore different parameters and settings to fine-tune prompts and optimize the model's behavior for specific tasks or domains.
- **Feedback and iteration:** Continuously analyze the outputs generated by the model and iterate on prompts based on user feedback to improve their quality and relevance.
- **Stay updated:** Keep up with the latest advancements in prompt engineering techniques, research, and best practices to enhance your skills and stay ahead in

the field.

Prompt engineering is a dynamic and evolving field, and skilled prompt engineers play a crucial role in harnessing the capabilities of LLM AI models effectively.

Take the next step

[Create your first prompt](#)

Writing prompts in Semantic Kernel

Article • 05/23/2023



To write an LLM AI prompt that Semantic Kernel is uniquely fit for, all you need is a concrete goal in mind — something you would like an AI to get done for you. For example:

I want to make a cake. Give me the best chocolate cake recipe you can think of.

Congratulations! You have imagined a delicious ask for Semantic Kernel to run to completion. This ask can be given to the Planner to get decomposed into steps. Although to make the Planner work reliably, you'll need to use the most advanced model available to you. So let's start from writing basic prompts to begin with.

ⓘ Note

Skills are currently being renamed to plugins. This article has been updated to reflect the latest terminology, but some images and code samples may still refer to skills.

💡 Tip

Want to easily follow along as you write your first prompts? Download the [Semantic Kernel VS Code Extension](#) which allows you to easily create and run prompts from within VS Code.



Semantic Kernel Tools

Microsoft microsoft.com | 3,604 installs | ★★★★★ (2) | Free

Prompt engineering tools to create AI plugins with Semantic Kernel

[Install](#)

[Trouble Installing?](#)

Writing a simple prompt

Writing prompts is like making a wish. Let's imagine we are entrepreneurs trying to make it in downtown Manhattan and we need to drive more leads to our store. We write the prompt:

Plain-Prompt

Write me a marketing slogan for my apparel shop in New York City with a focus on how affordable we are without sacrificing quality.

The result of this prompt from an actual LLM AI model is:

Response-From-LLM-AI-Model

New York Style, Low-Cost Smile:
Shop at NYC's Best Apparel Store!

Let's try another example where we are eager to play with the summarizing capability of LLM AIs and want to show off its superpower when applied to text that we explicitly define:

Plain-Prompt

Summarize the following text in two sentences or less.

---Begin Text---

Jan had always wanted to be a writer, ever since they were a kid. They spent hours reading books, writing stories, and imagining worlds. They grew up and pursued their passion, studying literature and journalism, and submitting their work to magazines and publishers. They faced rejection after rejection, but they never gave up hope. Jan finally got their breakthrough, when a famous editor discovered their manuscript and offered them a book deal.

---End Text---

The result of this prompt from an actual LLM AI model is:

Response-From-LLM-AI-Model

A possible summary is:

Jan's lifelong dream of becoming a writer came true when a famous editor offered them a book deal, after years of rejection and perseverance.

And there we have it. Two simple prompts that aren't asking the model for too much: 1/ we're asking the model to give us a marketing slogan, and separately 2/ we're asking the model to summarize a body of text down to two sentences.

Both of these simple prompts qualify as "functions" that can be packaged as part of an [Semantic Kernel plugin](#). The only problem is that they can do only one thing — as defined by the prompt — and with no flexibility. We set up the first plain prompt in Semantic Kernel within a directory named `SloganMaker` into a file named `skprompt.txt`:

```
SloganMaker/skprompt.txt
```

```
Write me a marketing slogan for my apparel shop in New  
York City with a focus on how affordable we are without  
sacrificing quality.
```

Similarly, we place the second plain prompt into a directory named `SummarizeBlurb` as a file named into a file named `skprompt.txt`.

```
SummarizeBlurb/skprompt.txt
```

```
Summarize the following text in two sentences or less.  
---Begin Text---  
Jan had always wanted to be a writer, ever since they  
were a kid. They spent hours reading books, writing  
stories, and imagining worlds. They grew up and pursued  
their passion, studying literature and journalism, and  
submitting their work to magazines and publishers. They  
faced rejection after rejection, but they never gave up  
hope. Jan finally got their breakthrough, when a famous  
editor discovered their manuscript and offered them a  
book deal.  
---End Text---
```

Each of these directories comprise a Semantic Kernel function. When both of the directories are placed inside an enclosing directory called `TestPlugin` the result is a brand new plugin.

```
Semantic-Plugins-And-Their-Functions
```

```
TestPlugin  
|  
|   SloganMaker  
|   |  
|   |   skprompt.txt  
|   |   [config.json](../howto/configuringfunctions)  
|  
|   SummarizeBlurb
```

```
|  
└── skprompt.txt  
└── [config.json](../howto/configuringfunctions)
```

This plugin can do one of two things by calling one of its two functions:

- `TestPlugin.SloganMaker()` generates a slogan for a specific kind of shop in NYC
- `TestPlugin.SummarizeBlurb()` creates a short summary of a specific blurb

Next, we'll show you how to make a more powerful plugin by introducing Semantic Kernel prompt templates. But before we do so, you may have noticed the `config.json` file. That's a special file for customizing how you want the function to run so that its performance can be tuned. If you're eager to know what's inside that file you can go [here](#) but no worries — you'll be running in no time. So let's keep going!

Writing a more powerful "templated" prompt

Let's say we want to go into the advertising business with AI powering the slogan-side of our offerings. We'd like to encapsulate how we create slogans to be repeatable and across any industry. To do so, we take our first prompt and write it as such as a "templated prompt":

```
SloganMakerFlex/skprompt.txt
```

```
Write me a marketing slogan for my {{$INPUT}} in New  
York City with a focus on how affordable we are without  
sacrificing quality.
```

Such "templated" prompts include variables and function calls that can dynamically change the content and the behavior of an otherwise plain prompt. Prompt templates can help you to generate more diverse, relevant, and effective prompts, and to reuse and combine them for different tasks and domains.

In a templated prompt, the double `{} curly braces {}` signify to Semantic Kernel that there's something special for it to notice within the LLM AI prompt. To pass an input to a prompt, we refer to the default input variable `$INPUT` — and by the same token if we have other variables to work with, they will start with a dollar sign `$` as well.

Our other plain prompt for summarizing text into two sentences can take an `input` by simply replacing the existing body of text and replacing it with `$input` as follows:

```
SummarizeBlurbFlex/skprompt.txt
```

```
Summarize the following text in two sentences or less.
```

```
---Begin Text---
```

```
{{\$INPUT}}
```

```
---End Text---
```

We can name these two functions `SloganMakerFlex` and `SummarizeBlurbFlex` — as two new Semantic Kernel functions that can belong to a new `TestPluginFlex` plugin that now takes an input. To package these two function to be used by Semantic Kernel in the context of a plugin, we arrange our file hierarchy the same as we did before:

File-Structure-For-Plugin-Definition-With-Functions

```
TestPluginFlex
|
└── SloganMakerFlex
    |
    └── skprompt.txt
    └── config.json
|
└── SummarizeBlurbFlex
    |
    └── skprompt.txt
    └── config.json
```

Recall that the difference between our new "flex" plugins and our original "plain" plugins is that we've gained the added flexibility of being able to pass a single parameter like:

- `TestPluginFlex.SloganMakerFlex('detective agency')` generates a slogan for a 'detective agency' in NYC
- `TestPluginFlex.SummarizeBlurbFlex('<insert long text here>')` creates a short summary of a given blurb

Templated prompts can be further customized beyond a single `$INPUT` variable to take on more inputs to gain even greater flexibility. For instance, if we wanted our `SloganMaker` plugin to not only take into account the kind of business but also the business' location and specialty, we would write the function as:

```
SloganMakerFlex/skprompt.txt
```

```
Write me a marketing slogan for my {{\$INPUT}} in {{\$CITY}} with
a focus on {{\$SPECIALTY}} we are without sacrificing quality.
```

Note that although the use of `$INPUT` made sense as a generic input for a templated prompt, you're likely to want to give it a name that makes immediate sense like

`$BUSINESS` — so let's change the function accordingly:

```
SloganMakerFlex/skprompt.txt
```

```
Write me a marketing slogan for my {{$BUSINESS}} in {{$CITY}} with  
a focus on {{$SPECIALTY}} we are without sacrificing quality.
```

We can replace our `TestPluginFlex` plugin with this new definition for `SloganMakerFlex` to serve the minimum capabilities of a copywriting agency.

In Semantic Kernel, we refer to prompts and templated prompts as *functions* to clarify their role as a fundamental unit of computation within the kernel. We specifically refer to *semantic* functions when LLM AI prompts are used; and when conventional programming code is used we say *native* functions. To learn how to make a native function you can skip ahead to [building a native functions](#) if you're anxious.

Get your kernel ready

First off, you'll want to create an instance of the kernel and configure it to run with Azure OpenAI or regular OpenAI. If you're using Azure OpenAI:

```
C#  
  
using Microsoft.SemanticKernel;  
  
var kernel = Kernel.Builder.Build();  
  
kernel.Config.AddAzureOpenAITextCompletion(  
    "Azure_davinci", // LLM AI model alias  
    "text-davinci-003", // Azure OpenAI *Deployment ID*  
    "https://contoso.openai.azure.com/", // Azure OpenAI *Endpoint*  
    "...your Azure OpenAI Key..." // Azure OpenAI *Key*  
);
```

If you're using regular OpenAI:

```
C#  
  
using Microsoft.SemanticKernel;  
  
var kernel = Kernel.Builder.Build();  
  
kernel.Config.AddOpenAITextCompletion(  
    "OpenAI_davinci", // LLM AI model alias  
    "text-davinci-003", // OpenAI Model Name  
    "...your OpenAI API Key...", // OpenAI API key  
    "...your OpenAI Org ID..." // *optional* OpenAI
```

```
Organization ID  
);
```

Invoking a semantic function from C#

When running a semantic function from your app's root source directory `MyAppSource` your file structure will looks like:



When running the kernel in C# you will:

1. Import your desired semantic function by specifying the root plugins directory and the plugin's name
2. Get ready to pass your semantic function parameters with a `ContextVariables` object
3. Set the corresponding context variables with `<your context variables>.Set`
4. Select the semantic function to run within the plugin by selecting a function

In code, and assuming you've already instantiated and configured your kernel as `kernel` as described [above](#):

```
C#  
  
using Microsoft.SemanticKernel;  
using Microsoft.SemanticKernel.KernelExtensions;  
using Microsoft.SemanticKernel.Orchestration;  
  
// ... instantiate a kernel as kernel  
  
var myPlugin = kernel.ImportSemanticSkillFromDirectory("MyPluginsDirectory",
```

```
"TestPluginFlex");

var myContext = new ContextVariables();
myContext.Set("BUSINESS", "Basketweaving Service");
myContext.Set("CITY", "Seattle");
myContext.Set("SPECIALTY", "ribbons");

var myResult = await kernel.RunAsync(myContext, myPlugin["SloganMakerFlex"]);

Console.WriteLine(myResult);
```

The output will read similar to:

```
"Ribbons with Seattle Style: Quality You Can Count On!"
```

Invoking a semantic function inline from C#

It's possible to bypass the need to package your semantic functions explicitly in `skprompt.txt` files by choosing to create them on-the-fly as inline code at runtime. Let's take `summarizeBlurbFlex`:

```
summarizeBlurbFlex

Summarize the following text in two sentences or less.
---Begin Text---
{{$INPUT}}
---End Text---
```

and define the function inline in C# — assuming you've already instantiated and configured your kernel as `kernel` as described [above](#):

```
C#

using Microsoft.SemanticKernel;
using Microsoft.SemanticKernel.SemanticFunctions;

// ... instantiate a kernel as kernel

string summarizeBlurbFlex = """
Summarize the following text in two sentences or less.
---Begin Text---
{{$INPUT}}
---End Text---
""";
```



```
var myPromptConfig = new PromptTemplateConfig
{
    Description = "Take an input and summarize it super-succinctly.",
```

```

Completion =
{
    MaxTokens = 1000,
    Temperature = 0.2,
    TopP = 0.5,
}
};

var myPromptTemplate = new PromptTemplate(
    summarizeBlurbFlex,
    myPromptConfig,
    kernel
);

var myFunctionConfig = new SemanticFunctionConfig(myPromptConfig,
myPromptTemplate);

var myFunction = kernel.RegisterSemanticFunction(
    "TestPluginFlex",
    "summarizeBlurbFlex",
    myFunctionConfig);

var myOutput = await kernel.RunAsync("This is my input that will get
summarized for me. And when I go off on a tangent it will make it harder.
But it will figure out that the only thing to summarize is that this is a
text to be summarized. You think?",
    myFunction);

Console.WriteLine(myOutput);

```

Note that the configuration was given inline to the kernel with a `PromptTemplateConfig` object instead of a `config.json` file with the maximum number of tokens to use `MaxTokens`, the variability of words it will use as `TopP`, and the amount of randomness to consider in its response with `Temperature`. Keep in mind that when using C# these parameters will be *PascalCased* (each word is explicitly capitalized in a string) to be consistent with C# conventions, but in the `config.json` the parameters are *lowercase*. To learn more about these function parameters read how to [configure functions](#).

A more succinct way to make this happen is with default settings across the board:

C#

```

using Microsoft.SemanticKernel;
using Microsoft.SemanticKernel.KernelExtensions;
using Microsoft.SemanticKernel.Orchestration;

// ... instantiate a kernel as kernel

string summarizeBlurbFlex = """
Summarize the following text in two sentences or less.
---Begin Text---

```

```
{{\$INPUT}}
---End Text---
""";  
  
var mySummarizeFunction = kernel.CreateSemanticFunction(summarizeBlurbFlex,  
maxTokens: 1000);  
  
var myOutput = await kernel.RunAsync(  
    new ContextVariables("This is my input that will get summarized for me.  
And when I go off on a tangent it will make it harder But it will figure out  
that the only thing to summarize is that this is a text to be summarized.  
You think?"),  
    mySummarizeFunction);  
  
Console.WriteLine(myOutput);
```

Links to learn more about prompts

- [Azure OAI Text Generation Tutorial](#)
- [Transparency Note On Azure OAI](#)
- [Mini-Course on Azure OAI](#)
- [OpenAI's Best Practices Guide ↗](#)

Take the next step

You're now ready to take advantage of the *Kernel's* pipelining capability.

[Compose functions to connect them end-to-end](#)

What are Models?

Article • 05/23/2023



💡 Tip

The article provides a brief overview of GPT models, including their variants, how they work, and how they can be fine-tuned. It also mentions similar LLM AI models and compares models based on their number of parameters.

👉Summary generated by plugin **SummarizeSkill.Summarize** ↗

A *model* refers to a specific instance or version of an LLM AI, such as GPT-3 or Codex, that has been trained and fine-tuned on a large corpus of text or code (in the case of the Codex model), and that can be accessed and used through an API or a platform. OpenAI and Azure OpenAI offer a variety of models that can be customized and controlled through parameters or options, and that can be applied and integrated to various domains and tasks.

About available OpenAI and Azure OpenAI GPT models

There are four Generative Pre-trained Transformer (GPT) models currently available from OpenAI and Azure OpenAI. They are composed of four variants: Ada, Babbage, Curie, and Davinci. They differ in the number of parameters, the amount of data they were trained on, and the types of tasks they can perform.

Ada is the smallest and simplest model, with 350 million parameters and 40GB of text data. It can handle basic natural language understanding and generation tasks, such as classification, sentiment analysis, summarization, and simple conversation.

Babbage is a larger model, with 3 billion parameters and 300GB of text data. It can handle more complex natural language tasks, such as reasoning, logic, arithmetic, and word analogy.

Curie is a very large model, with 13 billion parameters and 800GB of text data. It can handle advanced natural language tasks, such as text-to-speech, speech-to-text, translation, paraphrasing, and question answering.

Davinci is the largest and most powerful model, with 175 billion parameters and 45TB of text data. It can handle almost any natural language task, as well as some multimodal tasks, such as image captioning, style transfer, and visual reasoning. It can also generate coherent and creative texts on any topic, with a high level of fluency, consistency, and diversity.

Model	Parameters	Tasks
text-ada-001	350 million	Basic NLU** and NLG**
text-babbage-001	3 billion	Complex NLU and NLG
text-curie-001	13 billion	Advanced NLU and NLG
text-davinci-003	175 billion	Almost any NLU, NLG, and multimodal task

**Natural Language Understanding (NLU) / Natural Language Generating (NLG)

How does a GPT model work?

A GPT model is a type of neural network that uses the transformer architecture to learn from large amounts of text data. The model has two main components: an encoder and a decoder. The encoder processes the input text and converts it into a sequence of vectors, called embeddings, that represent the meaning and context of each word. The decoder generates the output text by predicting the next word in the sequence, based on the embeddings and the previous words. The model uses a technique called attention to focus on the most relevant parts of the input and output texts, and to capture long-range dependencies and relationships between words. The model is trained by using a large corpus of texts as both the input and the output, and by minimizing the difference between the predicted and the actual words. The model can then be fine-tuned or adapted to specific tasks or domains, by using smaller and more specialized datasets.

What is a baseline comparison rubric for LLM AIs?

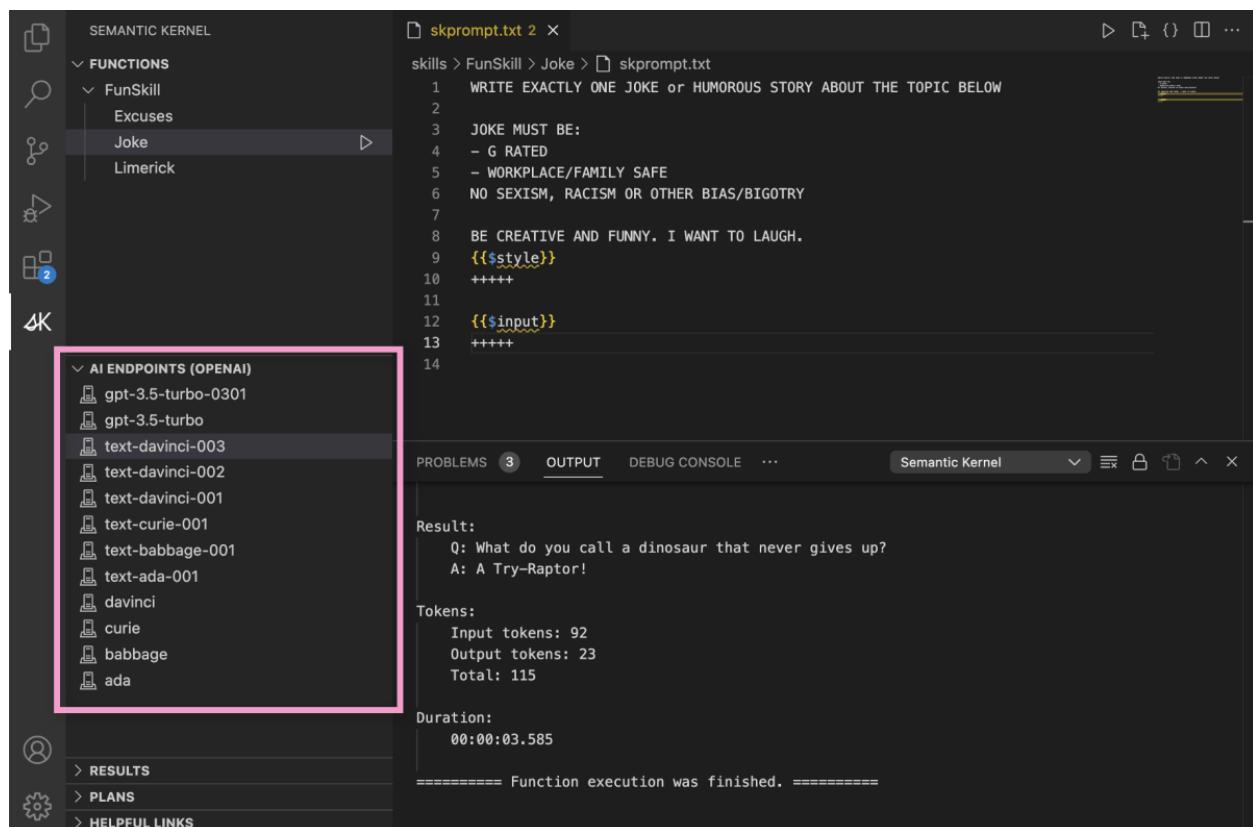
LLM AI Model	Parameters	Year
BERT	340 million	2018
GPT-2	1.5 billion	2019
Meena	2.6 billion	2020

LLM AI Model	Parameters	Year
GPT-3	175 billion	2020
LaMDA	137 billion	2022
BLOOM	176 billion	2022

LLM AI models are generally compared by the number of parameters — where bigger is usually better. The number of parameters is a measure of the size and the complexity of the model. The more parameters a model has, the more data it can process, learn from, and generate. However, having more parameters also means having more computational and memory resources, and more potential for overfitting or underfitting the data. Parameters are learned or updated during the training process, by using an optimization algorithm that tries to minimize the error or the loss between the predicted and the actual outputs. By adjusting the parameters, the model can improve its performance and accuracy on the given task or domain.

Easily test different models using Semantic Kernel tools

If you want to easily test how different models perform, you can use the [Semantic Kernel VS Code Extension](#) to quickly run a prompt on AI models from OpenAI, Azure OpenAI, and even Hugging Face.



The screenshot shows the Semantic Kernel extension in the VS Code interface. On the left, the sidebar displays the Semantic Kernel tree, including functions like FunSkill, Joke, and Limerick, and AI endpoints such as gpt-3.5-turbo-0301, gpt-3.5-turbo, text-davinci-003, and text-davinci-001. A pink box highlights the 'AI ENDPOINTS (OPENAI)' section. The main editor area contains a file named 'skprompt.txt' with the following content:

```

skills > FunSkill > Joke > skprompt.txt
1 WRITE EXACTLY ONE JOKE OR HUMOROUS STORY ABOUT THE TOPIC BELOW
2
3 JOKE MUST BE:
4 - G RATED
5 - WORKPLACE/FAMILY SAFE
6 NO SEXISM, RACISM OR OTHER BIAS/BIGOTRY
7
8 BE CREATIVE AND FUNNY. I WANT TO LAUGH.
9 {{\$style}}
10 ++++++
11
12 {{\$input}}
13 ++++++

```

The bottom right panel shows the execution results:

Result:
Q: What do you call a dinosaur that never gives up?
A: A Try-Raptor!

Tokens:
Input tokens: 92
Output tokens: 23
Total: 115

Duration:
00:00:03.585

===== Function execution was finished. =====

Take the next step

[Learn about configuring prompts](#)

Configuring prompts

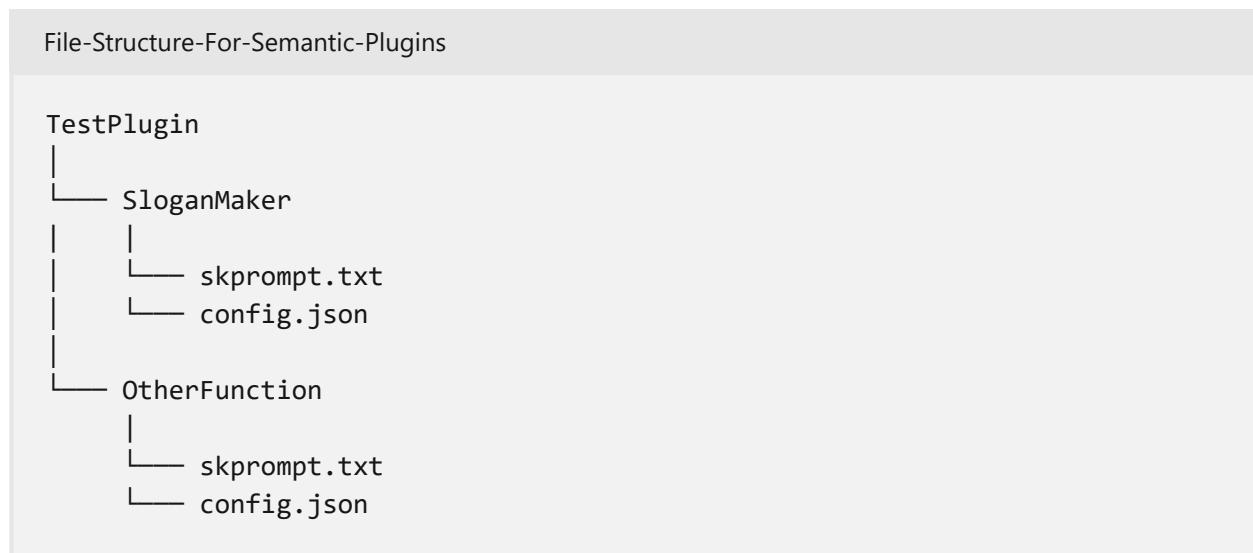
Article • 05/23/2023



When creating a prompt, there are many parameters that can be set to control how the prompt behaves. In Semantic Kernel, these parameters both control how a function is used by [planner](#) and how it is run by an [LLM AI model](#).

Semantic Kernel allows a developer to have complete control over these parameters by using a `config.json` file placed in the same directory as the `skprompt.txt` file.

For example, if you were to create a plugin called `TestPlugin` with two semantic functions called `SloganMaker` and `OtherFunction`, the file structure would look like this:



The `config.json` file for the `SloganMaker` function would look like this:

config.json-example

```
{
  "schema": 1,
  "type": "completion",
  "description": "a function that generates marketing slogans",
  "completion": {
    "max_tokens": 1000,
    "temperature": 0.0,
    "top_p": 0.0,
    "presence_penalty": 0.0,
    "frequency_penalty": 0.0
  }
  "input": {
    "parameters": [
```

```
{  
    "name": "input",  
    "description": "The product to generate a slogan for",  
    "defaultValue": ""  
}  
]  
}  
}
```

ⓘ Note

The `config.json` file is currently optional, but if you wish to exercise precise control of a function's behavior be sure to include it inside each function directory.

Parameters used by planner

The `description` field in the root object and `input` object are used by [planner](#) to determine how to use a function. The root `description` tells planner what the function does, and the input `description` tells planner how to populate the input parameters.

Because these parameters impact the behavior of planner, we recommend running tests on the values you provide to ensure they are used by planner correctly.

When writing `description` and `input`, we recommend using the following guidelines:

- The `description` fields should be short and concise so that it does not consume too many tokens when used in planner prompt.
- Consider the `descriptions` of other functions in the same plugin to ensure that they are sufficiently unique. If they are not, planner may not be able to distinguish between them.
- If you have trouble getting planner to use a function, try adding recommendations or examples for when to use the function.

Completion parameters in config.json

In addition to providing parameters for planner, the `config.json` file also allows you to control how a function is run by an [LLM AI model](#). The `completion` object in the root object of the `config.json` file allows you to set the parameters used by the model.

The following table describes the parameters available for use in the `completion` object for the OpenAI and Azure OpenAI APIs:

Completion Parameter	Type	Required?	Default	Description
<code>max_tokens</code>	integer	Optional	16	The maximum number of tokens to generate in the completion. The token count of your prompt plus <code>max_tokens</code> can't exceed the model's context length. Most models have a context length of 2048 tokens (except <code>davinci-codex</code> , which supports 4096).
<code>temperature</code>	number	Optional	1	What sampling temperature to use. Higher values mean the model will take more risks. Try 0.9 for more creative applications, and 0 (argmax sampling) for ones with a well-defined answer. We generally recommend altering this or <code>top_p</code> but not both.
<code>top_p</code>	number	Optional	1	An alternative to sampling with temperature, called nucleus sampling, where the model considers the results of the tokens with <code>top_p</code> probability mass. So 0.1 means only the tokens comprising the top 10% probability mass are considered. We generally recommend altering this or temperature but not both.
<code>presence_penalty</code>	number	Optional	0	Number between -2.0 and 2.0. Positive values penalize new tokens based on whether they appear in the text so far, increasing the model's likelihood to talk about new topics.
<code>frequency_penalty</code>	number	Optional	0	Number between -2.0 and 2.0. Positive values penalize new tokens based on their existing frequency in the text so far, decreasing the model's likelihood to repeat the same line verbatim.

To learn more about the various parameters available for tuning how a function works, visit the [Azure OpenAI reference](#).

Default setting for OpenAI and Azure OpenAI

If you do not provide completion parameters in the `config.json` file, Semantic Kernel will use the default parameters for the OpenAI API. Learn more about the current defaults by reading the [Azure OpenAI API reference](#).

Take the next step

[Understanding tokens](#)

What are Tokens?

Article • 05/23/2023



💡 Tip

Key topics:

- Tokens: basic units of text/code for LLM AI models to process/generate language.
- Tokenization: splitting input/output texts into smaller units for LLM AI models.
- Vocabulary size: the number of tokens each model uses, which varies among different GPT models.
- Tokenization cost: affects the memory and computational resources that a model needs, which influences the cost and performance of running an OpenAI or Azure OpenAI model.

👉 Topics list generated by plugin **SummarizeSkill.Topics** ↗

Tokens are the basic units of text or code that an LLM AI uses to process and generate language. Tokens can be characters, words, subwords, or other segments of text or code, depending on the chosen tokenization method or scheme. Tokens are assigned numerical values or identifiers, and are arranged in sequences or vectors, and are fed into or outputted from the model. Tokens are the building blocks of language for the model.

How does tokenization work?

Tokenization is the process of splitting the input and output texts into smaller units that can be processed by the LLM AI models. Tokens can be words, characters, subwords, or symbols, depending on the type and the size of the model. Tokenization can help the model to handle different languages, vocabularies, and formats, and to reduce the computational and memory costs. Tokenization can also affect the quality and the diversity of the generated texts, by influencing the meaning and the context of the tokens. Tokenization can be done using different methods, such as rule-based, statistical, or neural, depending on the complexity and the variability of the texts.

OpenAI and Azure OpenAI uses a subword tokenization method called "Byte-Pair Encoding (BPE)" for its GPT-based models. BPE is a method that merges the most frequently occurring pairs of characters or bytes into a single token, until a certain number of tokens or a vocabulary size is reached. BPE can help the model to handle rare or unseen words, and to create more compact and consistent representations of the texts. BPE can also allow the model to generate new words or tokens, by combining existing ones. The way that tokenization is different dependent upon the different model Ada, Babbage, Curie, and Davinci is mainly based on the number of tokens or the vocabulary size that each model uses. Ada has the smallest vocabulary size, with 50,000 tokens, and Davinci has the largest vocabulary size, with 60,000 tokens. Babbage and Curie have the same vocabulary size, with 57,000 tokens. The larger the vocabulary size, the more diverse and expressive the texts that the model can generate. However, the larger the vocabulary size, the more memory and computational resources that the model requires. Therefore, the choice of the vocabulary size depends on the trade-off between the quality and the efficiency of the model.

What does tokenization have to do with the cost of running a model?

Tokenization affects the amount of data and the number of calculations that the model needs to process. The more tokens that the model has to deal with, the more memory and computational resources that the model consumes. Therefore, the cost of running an OpenAI or Azure OpenAI model depends on the tokenization method and the vocabulary size that the model uses, as well as the length and the complexity of the input and output texts. Based on the number of tokens used for interacting with a model and the different rates for different models, your costs can widely differ. For example, as of February 2023, the rate for using Davinci is \$0.06 per 1,000 tokens, while the rate for using Ada is \$0.0008 per 1,000 tokens. The rate also varies depending on the type of usage, such as playground, search, or engine. Therefore, tokenization is an important factor that influences the cost and the performance of running an OpenAI or [Azure OpenAI model ↗](#).

Using Semantic Kernel tools to measure token use

If you want to measure how much consumption each of your prompt uses, you can use the [Semantic Kernel VS Code Extension](#) to see how many input and output tokens are necessary to run a prompt.

The screenshot shows the Semantic Kernel IDE interface. On the left is a sidebar with icons for file operations, search, and navigation. The main area has two tabs: 'SEMANTIC KERNEL' and 'OUTPUT'. The 'SEMANTIC KERNEL' tab displays a tree view of functions and AI endpoints. Under 'FUNCTIONS', 'FunSkill' is expanded, showing 'Excuses', 'Joke' (which is selected), and 'Limerick'. Under 'AI ENDPOINTS (OPENAI)', 'text-davinci-003' is selected. The 'OUTPUT' tab shows a code editor with a file named 'skprompt.txt' containing a template for generating a joke. The terminal window below shows the execution results:

```
skills > FunSkill > Joke > skprompt.txt
1 WRITE EXACTLY ONE JOKE or HUMOROUS STORY ABOUT THE TOPIC BELOW
2
3 JOKE MUST BE:
4 - G RATED
5 - WORKPLACE/FAMILY SAFE
6 NO SEXISM, RACISM OR OTHER BIAS/BIGOTRY
7
8 BE CREATIVE AND FUNNY. I WANT TO LAUGH.
9 {{${style}}}
10 +++++
11
12 {{${input}}}
13 +++++
14
```

PROBLEMS 3 OUTPUT DEBUG CONSOLE ... Semantic Kernel

Result:
Q: What do you call a dinosaur that never gives up?
A: A Try-Raptor!

Tokens:
Input tokens: 92
Output tokens: 23
Total: 115

Duration:
00:00:03.585

===== Function execution was finished. =====

Prompt template syntax

Article • 05/23/2023



The Semantic Kernel prompt template language is a simple and powerful way to define and compose AI functions **using plain text**. You can use it to create natural language prompts, generate responses, extract information, **invoke other prompts** or perform any other task that can be expressed with text.

The language supports three basic features that allow you to (#1) include variables, (#2) call external functions, and (#3) pass parameters to functions.

You don't need to write any code or import any external libraries, just use the curly braces `{{...}}` to embed expressions in your prompts. Semantic Kernel will parse your template and execute the logic behind it. This way, you can easily integrate AI into your apps with minimal effort and maximum flexibility.

Variables

To include a variable value in your text, use the `{{$variableName}}` syntax. For example, if you have a variable called `name` that holds the user's name, you can write:

```
Hello {{$name}}, welcome to Semantic Kernel!
```

This will produce a greeting with the user's name.

Function calls

To call an external function and embed the result in your text, use the `{{$namespace.functionName}}` syntax. For example, if you have a function called `weather.getForecast` that returns the weather forecast for a given location, you can write:

```
The weather today is {{$weather.getForecast}}.
```

This will produce a sentence with the weather forecast for the default location stored in the `input` variable. The `input` variable is set automatically by the kernel when invoking a function. For instance, the code above is equivalent to:

```
The weather today is {{weather.getForecast $input}}.
```

Function parameters

To call an external function and pass a parameter to it, use the `{{namespace.functionName $varName}}` syntax. For example, if you want to pass a different input to the weather forecast function, you can write:

```
The weather today in {{$city}} is {{weather.getForecast $city}}. The weather today  
in {{$region}} is {{weather.getForecast $region}}.
```

This will produce two sentences with the weather forecast for two different locations, using the city stored in the `city` variable and the region name stored in the `region` variable.

Design Principles

The template language uses of the `$` symbol on purpose, to clearly distinguish between variables, which are retrieved from local temporary memory, and functions that retrieve content executing some code.

Branching features such as "if", "for", and code blocks are not part of SK's template language. This reflects SK's design principle of using natural language as much as possible, with a clear separation from conventional programming code.

By using a simple language, the kernel can also avoid complex parsing and external dependencies, resulting in a fast and memory efficient processing.

Semantic function example

Here's a very simple example of a semantic function defined with a prompt template, using the syntax described.

```
skprompt.txt
```

```
My name: {{msgraph.GetMyName}}  
My email: {{msgraph.GetMyEmailAddress}}  
Recipient: {{$recipient}}  
Email to reply to:  
=====  
{{$sourceEmail}}  
=====  
Generate a response to the email, to say: {{$input}}
```

Include the original email quoted after the response.

If we were to write that function in C#, it would look something like:

C#

```
async Task<string> GenResponseToEmailAsync(
    string whatToSay,
    string recipient,
    string sourceEmail)
{
    try {
        string name = await this._msgraph.GetMyName();
    } catch {
        ...
    }

    try {
        string email = await this._msgraph.GetMyEmailAddress();
    } catch {
        ...
    }

    try {
        // Use AI to generate an email using the 5 given variables
        // Take care of retry logic, tracking AI costs, etc.
        string response = await ...

        return response;
    } catch {
        ...
    }
}
```

What are Memories?

Article • 05/23/2023



Memories are a powerful way to provide broader context for your ask. Historically, we've always called upon memory as a core component for how computers work: think the RAM in your laptop. For with just a CPU that can crunch numbers, the computer isn't that useful unless it knows what numbers you care about. Memories are what make computation relevant to the task at hand.

We access memories to be fed into Semantic Kernel in one of three ways — with the third way being the most interesting:

1. Conventional key-value pairs: Just like you would set an environment variable in your shell, the same can be done when using Semantic Kernel. The lookup is "conventional" because it's a one-to-one match between a key and your query.
2. 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.
3. 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.

How does semantic memory work?

Embeddings are a way of representing words or other data as vectors in a high-dimensional space. Vectors are like arrows that have a direction and a length. High-dimensional means that the space has many dimensions, more than we can see or imagine. The idea is that similar words or data will have similar vectors, and different words or data will have different vectors. This helps us measure how related or unrelated they are, and also perform operations on them, such as adding, subtracting, multiplying, etc. Embeddings are useful for AI models because they can capture the meaning and context of words or data in a way that computers can understand and process.

So basically you take a sentence, paragraph, or entire page of text, and then generate the corresponding embedding vector. And when a query is performed, the query is transformed to its embedding representation, and then a search is performed through

all the existing embedding vectors to find the most similar ones. This is similar to when you make a search query on Bing, and it gives you multiple results that are proximate to your query. Semantic memory is not likely to give you an exact match — but it will always give you a set of matches ranked in terms of how similar your query matches other pieces of text.

Why are embeddings important with LLM AI?

Since a prompt is a text that we give as input to an AI model to generate a desired output or response, we need to consider the length of the input text based on the token limit of the model we choose to use. For example, GPT-4 can handle up to 8,192 tokens per input, while GPT-3 can only handle up to 4,096 tokens. This means that texts that are longer than the token limit of the model will not fit and may be cut off or ignored.

It would be nice if we could use an entire 10,000-page operating manual as context for our prompt, but because of the token limit constraint, that is impossible. Therefore, embeddings are useful for breaking down that large text into smaller pieces. We can do this by summarizing each page into a shorter paragraph and then generating an embedding vector for each summary. An embedding vector is like a compressed representation of the text that preserves its meaning and context. Then we can compare the embedding vectors of our summaries with the embedding vector of our prompt and select the most similar ones. We can then add those summaries to our input text as context for our prompt. This way, we can use embeddings to help us choose and fit large texts as context within the token limit of the model.

Take the next step

[Learn about embeddings](#)

What are Embeddings?

Article • 05/23/2023



💡 Tip

Memory: Embeddings

- Embeddings are vectors or arrays of numbers that represent the meaning and the context of tokens processed by the model.
- They are used to encode and decode input and output texts, and can vary in size and dimension. / Embeddings can help the model understand the relationships between tokens, and generate relevant and coherent texts.
- They are used for text classification, summarization, translation, and generation, as well as image and code generation.

👉 Notes generated by plugin **SummarizeSkill.Notegen** ↗

Embeddings are the representations or encodings of [tokens](#), such as sentences, paragraphs, or documents, in a high-dimensional vector space, where each dimension corresponds to a learned feature or attribute of the language. Embeddings are the way that the model captures and stores the meaning and the relationships of the language, and the way that the model compares and contrasts different tokens or units of language. Embeddings are the bridge between the discrete and the continuous, and between the symbolic and the numeric, aspects of language for the model.

What are embeddings to a programmer?

Embeddings are vectors or arrays of numbers that represent the meaning and the context of the tokens that the model processes and generates. Embeddings are derived from the parameters or the weights of the model, and are used to encode and decode the input and output texts. Embeddings can help the model to understand the semantic and syntactic relationships between the tokens, and to generate more relevant and coherent texts. Embeddings can also enable the model to handle multimodal tasks, such as image and code generation, by converting different types of data into a common representation. Embeddings are an essential component of the transformer architecture

that GPT-based models use, and they can vary in size and dimension depending on the model and the task.

How are embeddings used?

Embeddings are used for:

- **Text classification:** Embeddings can help the model to assign labels or categories to texts, based on their meaning and context. For example, embeddings can help the model to classify texts as positive or negative, spam or not spam, news or opinion, etc.
- **Text summarization:** Embeddings can help the model to extract or generate the most important or relevant information from texts, and to create concise and coherent summaries. For example, embeddings can help the model to summarize news articles, product reviews, research papers, etc.
- **Text translation:** Embeddings can help the model to convert texts from one language to another, while preserving the meaning and the structure of the original texts. For example, embeddings can help the model to translate texts between English and Spanish, French and German, Chinese and Japanese, etc.
- **Text generation:** Embeddings can help the model to create new and original texts, based on the input or the prompt that the user provides. For example, embeddings can help the model to generate texts such as stories, poems, jokes, slogans, captions, etc.
- **Image generation:** Embeddings can help the model to create images from texts, or vice versa, by converting different types of data into a common representation. For example, embeddings can help the model to generate images such as logos, faces, animals, landscapes, etc.
- **Code generation:** Embeddings can help the model to create code from texts, or vice versa, by converting different types of data into a common representation. For example, embeddings can help the model to generate code such as HTML, CSS, JavaScript, Python, etc.

Take the next step

[Learn about vector databases](#)

What is a vector database?

Article • 06/23/2023



A vector database is a type of database that stores data as high-dimensional vectors, which are mathematical representations of features or attributes. Each vector has a certain number of dimensions, which can range from tens to thousands, depending on the complexity and granularity of the data. The vectors are usually generated by applying some kind of transformation or embedding function to the raw data, such as text, images, audio, video, and others. The embedding function can be based on various methods, such as machine learning models, word embeddings, feature extraction algorithms.

The main advantage of a vector database is that it allows for fast and accurate similarity search and retrieval of data based on their vector distance or similarity. This means that instead of using traditional methods of querying databases based on exact matches or predefined criteria, you can use a vector database to find the most similar or relevant data based on their semantic or contextual meaning.

For example, you can use a vector database to:

- find images that are similar to a given image based on their visual content and style
- find documents that are similar to a given document based on their topic and sentiment
- find products that are similar to a given product based on their features and ratings

To perform similarity search and retrieval in a vector database, you need to use a query vector that represents your desired information or criteria. The query vector can be either derived from the same type of data as the stored vectors (e.g., using an image as a query for an image database), or from different types of data (e.g., using text as a query for an image database). Then, you need to use a similarity measure that calculates how close or distant two vectors are in the vector space. The similarity measure can be based on various metrics, such as cosine similarity, euclidean distance, hamming distance, jaccard index.

The result of the similarity search and retrieval is usually a ranked list of vectors that have the highest similarity scores with the query vector. You can then access the

corresponding raw data associated with each vector from the original source or index.

Use Cases for Vector Databases

Vector databases have many use cases across different domains and applications that involve natural language processing (NLP), computer vision (CV), recommendation systems (RS), and other areas that require semantic understanding and matching of data.

One use case for storing information in a vector database is to enable large language models (LLMs) to generate more relevant and coherent text based on an [AI plugin](#).

However, large language models often face challenges such as generating inaccurate or irrelevant information; lacking factual consistency or common sense; repeating or contradicting themselves; being biased or offensive. To overcome these challenges, you can use a vector database to store information about different topics, keywords, facts, opinions, sources, related to your desired domain or genre. Then, you can use a large language model and pass information from the vector database with your AI plugin to generate more informative and engaging content that matches your intent and style.

For example, if you want to write a blog post about the latest trends in AI, you can use a vector database to store the latest information about that topic and pass the information along with the ask to a LLM in order to generate a blog post that leverages the latest information.

Available connectors to vector databases

Today, we have several connectors to vector databases that you can use to store and retrieve information. These include:

- [Azure Cognitive Search](#) ↗
- [Chroma](#) ↗
- [COSMOS DB](#) ↗
- [Redis](#) ↗
- [Pinecone](#) ↗
- [Postgres](#) ↗
- [Qdrant](#) ↗
- [Sqlite](#) ↗
- [Weaviate](#) ↗ .net and [Python](#) ↗

Take the next step

Create and deploy plugins

Learn how to deploy Semantic Kernel to Azure as a web app service

Article • 05/24/2023



In this how-to guide we will provide steps to deploy Semantic Kernel to Azure as a web app service. Deploying Semantic Kernel as web service to Azure provides a great pathway for developers to take advantage of Azure compute and other services such as Azure Cognitive Services for responsible AI and vectorized databases.

You can use one of the deployment options to deploy based on your use case and preference.

Considerations

1. Azure currently limits the number of Azure OpenAI resources per region per subscription to 3. Azure OpenAI is not available in every region. (Refer to this [availability map](#)) Bearing this in mind, you might want to use the same Azure OpenAI instance for multiple deployments of Semantic Kernel to Azure.
2. F1 and D1 App Service SKU's (the Free and Shared ones) are not supported for this deployment.
3. Ensure you have sufficient permissions to create resources in the target subscription.
4. Using web frontends to access your deployment: make sure to include your frontend's URL as an allowed origin in your deployment's CORS settings. Otherwise, web browsers will refuse to let JavaScript make calls to your deployment.

Use Case	Deployment Option
<u>Use existing: Azure OpenAI Resources</u> Use this option to use an existing Azure OpenAI instance and connect the Semantic Kernel web API to it.	 Deploy to Azure PowerShell File ↗ Bash File ↗

Use Case	Deployment Option
<p><u>Create new: Azure OpenAI Resources</u></p> <p>Use this option to deploy Semantic Kernel in a web app service and have it use a new instance of Azure OpenAI.</p> <p>Note: access to new Azure OpenAI resources is currently limited due to high demand.</p>	<p> Deploy to Azure</p> <p> PowerShell File ↗</p> <p> Bash File ↗</p>
<p><u>Use existing: OpenAI Resources</u></p> <p>Use this option to use your OpenAI account and connect the Semantic Kernel web API to it.</p>	<p> Deploy to Azure</p> <p> PowerShell File ↗</p> <p> Bash File ↗</p>

Script Parameters

Below are examples on how to run the PowerShell and bash scripts. Refer to each of the script files for the complete list of available parameters and usage.

PowerShell

- Creating new Azure OpenAI Resources

PowerShell

```
.\DeploySK.ps1 -DeploymentName YOUR_DEPLOYMENT_NAME -Subscription  
YOUR_SUBSCRIPTION_ID
```

- Using existing Azure OpenAI Resources

After entering the command below, you will be prompted to enter your Azure OpenAI API key. (You can also pass in the API key using the -ApiKey parameter)

PowerShell

```
.\DeploySK-Existing-AzureOpenAI.ps1 -DeploymentName YOUR_DEPLOYMENT_NAME -  
Subscription YOUR_SUBSCRIPTION_ID -Endpoint "YOUR_AZURE_OPENAI_ENDPOINT"
```

- Using existing OpenAI Resources

After entering the command below, you will be prompted to enter your OpenAI API key. (You can also pass in the API key using the -ApiKey parameter)

PowerShell

```
.\DeploySK-Existing-OpenAI.ps1 -DeploymentName YOUR_DEPLOYMENT_NAME -  
Subscription YOUR_SUBSCRIPTION_ID
```

Bash

- Creating new Azure OpenAI Resources

Bash

```
./DeploySK.sh -d DEPLOYMENT_NAME -s SUBSCRIPTION_ID
```

- Using existing Azure OpenAI Resources

Bash

```
./DeploySK-Existing-AzureOpenAI.sh -d YOUR_DEPLOYMENT_NAME -s  
YOUR_SUBSCRIPTION_ID -e YOUR_AZURE_OPENAI_ENDPOINT -o  
YOUR_AZURE_OPENAI_API_KEY
```

- Using existing OpenAI Resources

Bash

```
./DeploySK-Existing-AI.sh -d YOUR_DEPLOYMENT_NAME -s YOUR_SUBSCRIPTION_ID -  
o YOUR_OPENAI_API_KEY
```

Azure Portal Template

If you choose to use Azure Portal as your deployment method, you will need to review and update the template form to create the resources. Below is a list of items you will need to review and update.

1. Subscription: decide which Azure subscription you want to use. This will house the resource group for the Semantic Kernel web application.
2. Resource Group: the resource group in which your deployment will go. Creating a new resource group helps isolate resources, especially if you are still in active development.
3. Region: select the geo-region for deployment. Note: Azure OpenAI is not available in all regions and is currently limited to three instances per region per subscription.

4. Name: used to identify the app. App Service SKU: select the pricing tier based on your usage. Click [here](#) to learn more about Azure App Service plans.
5. Package URL: there is no need to change this unless you want to deploy a customized version of Semantic Kernel. (See [this page](#) for more information on publishing your own version of the Semantic Kernel web app service)
6. Completion, Embedding and Planner Models: these are by default using the appropriate models based on the current use case - that is Azure OpenAI or OpenAI. You can update these based on your needs.
7. Endpoint: this is only applicable if using Azure OpenAI and is the Azure OpenAI endpoint to use.
8. API Key: enter the API key for the instance of Azure OpenAI or OpenAI to use.
9. Semantic Kernel API Key: the default value of "[newGuid()]" in this field will create an API key to protect your Semantic Kernel endpoint. You can change this by providing your own API key. If you do not want to use API authorization, you can make this field blank.
10. CosmosDB: whether to deploy a CosmosDB resource to store chats. Otherwise, volatile memory will be used.
11. Qdrant: whether to deploy a Qdrant database to store embeddings. Otherwise, volatile memory will be used.
12. Speech Services: whether to deploy an instance of the Azure Speech service to provide speech-to-text for input.

What resources are deployed?

Below is a list of the key resources created within the resource group when you deploy Semantic Kernel to Azure as a web app service.

1. Azure web app service: hosts Semantic Kernel
2. Application Insights: application logs and debugging
3. Azure Cosmos DB: used for chat storage (optional)
4. Qdrant vector database (within a container): used for embeddings storage (optional)
5. Azure Speech service: used for speech-to-text (optional)

Verifying the deployment

To make sure your web app service is running, go to
https://YOUR_INSTANCE_NAME.azurewebsites.net/healthz

To get your instance's URL, go to your deployment's resource group (by clicking on the "Go to resource group" button seen at the conclusion of your deployment if you use the

"Deploy to Azure" button). Then click on the resource whose name ends with "-skweb".

This will bring you to the Overview page on your web service. Your instance's URL is the value that appears next to the "Default domain" field.

Changing your configuration, monitoring your deployment and troubleshooting

After your deployment is complete, you can change your configuration in the Azure Portal by clicking on the "Configuration" item in the "Settings" section of the left pane found in the Semantic Kernel web app service page.

Scrolling down in that same pane to the "Monitoring" section gives you access to a multitude of ways to monitor your deployment.

In addition to this, the "Diagnose and solve problems" item near the top of the pane can yield crucial insight into some problems your deployment may be experiencing.

If the service itself is functioning properly but you keep getting errors (perhaps reported as 400 HTTP errors) when making calls to the Semantic Kernel, check that you have correctly entered the values for the following settings:

- AIService:AzureOpenAI
- AIService:Endpoint
- AIService:Models:Completion
- AIService:Models:Embedding
- AIService:Models:Planner

AIService:Endpoint is ignored for OpenAI instances from openai.com but MUST be properly populated when using Azure OpenAI instances.

How to clean up resources

When you want to clean up the resources from this deployment, use the Azure portal or run the following [Azure CLI](#) command:

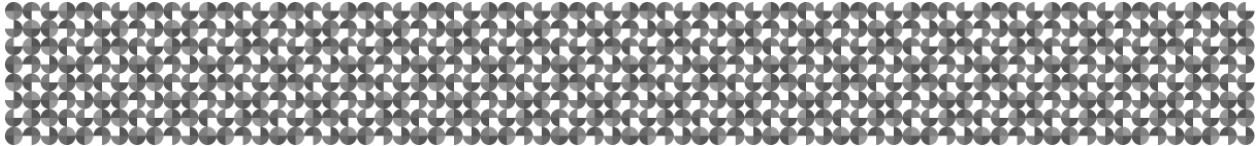
PowerShell

```
az group delete --name YOUR_RESOURCE_GROUP
```

Take the next step

Learn [how to make changes to your Semantic Kernel web app](#), such as adding new skills.

If you have not already done so, please star the GitHub repo and join the Semantic Kernel community! [Star the Semantic Kernel repo ↗](#)



Learn how to make changes to the Semantic Kernel web app service

Article • 05/26/2023



This guide provides steps to make changes to the skills of a deployed instance of the Semantic Kernel web app. Currently, changing semantic skills can be done without redeploying the web app service but changes to native skills do require re-deployments. This document will guide you through the process of doing both.

Prerequisites

1. An instance of the Semantic Kernel web app service deployed in your Azure subscription. You can follow the how-to guide [here](#) for details.
2. Have your web app's name handy. If you used the deployment templates provided with the Chat Copilot, you can find the web app's name by going to the [Azure Portal](#) and selecting the resource group created for your Semantic Kernel web app service. Your web app's name is the one of the resource listed that ends with "skweb".
3. Locally tested [skills](#) or [planner](#) ready to be added to your Semantic Kernel web app service.

How to publish changes to the Semantic Kernel web app service

There are two main ways to deploy changes to the Semantic Kernel web app service. If you have been working locally and are ready to deploy your changes to Azure as a new web app service, you can follow the steps in the first section. If you have already deployed your Semantic Kernel web app service and want to make changes to add Semantic skills, you can follow the steps in the second section.

1. Deploying your Chat Copilot App to Azure as a web application

After working locally, i.e. you cloned the code from the GitHub [repo](#) and have made changes to the code for your needs, you can deploy your changes to Azure as a web application.

You can use the standard methods available to [deploy an ASP.net web app](#) in order to do so.

Alternatively, you can follow the steps below to manually build and upload your customized version of the Semantic Kernel service to Azure.

First, at the command line, go to the '..semantic-kernel/samples/apps/copilot-chat-app/webapi' directory and enter the following command:

PowerShell

```
dotnet publish CopilotChatApi.csproj --configuration Release --arch x64 --os win
```

This will create a directory which contains all the files needed for a deployment:

Windows Command Prompt

```
../semantic-kernel/samples/apps/copilot-chat-app/webapi/bin/Release/net6.0/win-x64/publish'
```

Zip the contents of that directory and store the resulting zip file on cloud storage, e.g. Azure Blob Container. Put its URI in the "Package Uri" field in the web deployment page you access through the "Deploy to Azure" buttons or use its URI as the value for the PackageUri parameter of the deployment scripts found on this [page](#).

Your deployment will then use your customized deployment package. That package will be used to create a new Azure web app, which will be configured to run your customized version of the Semantic Kernel service.

2. Publish skills directly to the Semantic Kernel web app service

This method is useful for making changes when adding new semantic skills only.

How to add Semantic Skills

1. Go to https://YOUR_APP_NAME.scm.azurewebsites.net, replacing YOUR_APP_NAME in the URL with your app name found in Azure Portal. This will take you to the [Kudu](#) console for your app hosting.
2. Click on Debug Console and select CMD.
3. Navigate to the 'site\wwwroot\Skills'
4. Create a new folder using the (+) sign at the top and give a folder name to store your Semantic Skills e.g. SemanticSkills.
5. Now you can drag and drop your Semantic Skills into this folder
6. Next navigate to 'site\wwwroot'
7. Click on the pencil icon to edit the appsettings.json file.
8. In the appsettings.json file, update the SemanticSkillDirectory with the location of the skills you have created.

JSON

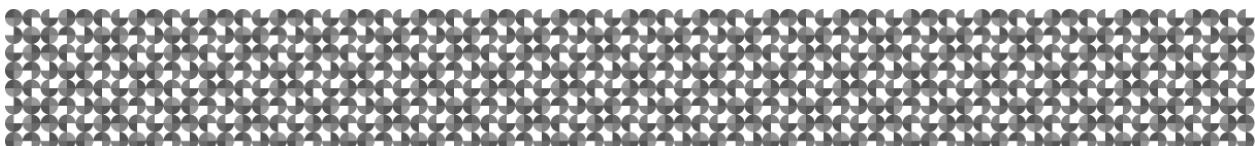
```
"Service": {  
    "SemanticSkillsDirectory": "/SemanticSkills",  
    "KeyVaultUri": ""  
},
```

9. Click on "Save" to save the changes to the appsettings.json file.
10. Now your web app is configured to use your Semantic Skills.

Take the next step

To explore how you build a front-end web app explore the [Chat Copilot app](#).

If you have not already done so, please star the GitHub repo and join the Semantic Kernel community! [Star the Semantic Kernel repo ↗](#)



Protect your Azure OpenAI API keys with Azure API Management

Article • 06/07/2023



With [Azure API Management](#), you can protect your AI API keys and manage access to your AI APIs. This is helpful if you want to give your users or developers access to Azure OpenAPI APIs without giving them direct access to your keys or to manage access to your AI APIs directly from Azure.

By the end of this article, you'll have a working API Management instance along with sample code that shows you how to use Azure API Management with Semantic Kernel.

Prerequisites

To complete this tutorial, you'll first need access to the following:

- Azure OpenAI API key
- API Management instance

To create a new API Management instance, see [Create an API Management instance](#).

You can use the default configuration when creating your instance.

Setup Azure API Management instance with Azure OpenAI API

The following steps describe how you can setup your Azure OpenAI API with Azure API Management.

1. Start by [downloading the API definition for Azure OpenAI](#) and saving it to your local machine.

[Download API reference](#)

2. Open the file in a text editor and change the `servers` property so that the `url` and `endpoint` properties point to your API endpoint.

For example, if your API endpoint is `https://contoso.openai.azure.com`, you would change the `servers` property to the following:

JSON

```
"servers": [
  {
    "url": "https://contoso.openai.azure.com/openai",
    "variables": {
      "endpoint": {
        "default": "contoso.openai.azure.com"
      }
    }
  ],
],
```

3. You can now import the file from step 2 into Azure API Management. The full process is described in the [import and publish](#) article.

4. Now create a named value for your Azure OpenAI API key.

This will allow you to better protect your Azure OpenAI API key. To create a named value, see [using named values in Azure API Management policies](#). Take note of the **Display name** you give your named value as you'll need it in the next step.

The screenshot shows the Azure portal's configuration interface for creating a named value. The 'Name' field is set to 'AzureOpenAIKey'. The 'Display name' field is highlighted with an orange border and contains the value 'azure-openai-key'. The 'Type' field is set to 'Secret'. A 'Show value' button is visible at the bottom left. The entire form is contained within a light gray box.

5. Finally, [edit the inbound policy](#) so your API adds the `api-key` header to your request.

- a. Navigate to your API Management instance and select **APIs** from the left-hand menu.
- b. Select your API from the list of APIs.
- c. Select the **Design** tab.
- d. Select **Add policy** within the Inbound processing section.
- e. Select **Set headers**.

- f. In the Set headers page, enter `openai-key` for the Name field
- g. For the Value field, take the display name of your named value in the previous step and wrap it in double curly braces. For example, if the display name of your named value was `azure-openai-key`, you would enter `{{azure-openai-key}}`.
- h. Leave the action as `append`.

Inbound processing

Modify the request before it is sent to the backend service.

Set headers

Set policy to add, remove or change headers that are passed to the backend service.

[Learn more about "set-header" policy.](#)

NAME	VALUE	ACTION	DELETE
api-key	<code>{{azure-openai-key}}</code>	append	

[+ Add header](#)

- i. Select **Save**.

You can also use key vault secrets to protect your Azure OpenAI API key even more. To learn more, see [using named values in Azure API Management policies](#).

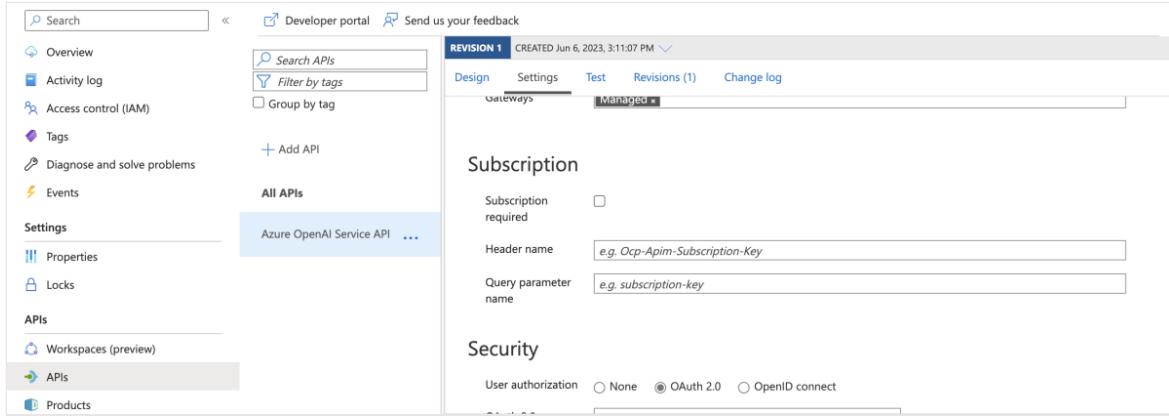
You've now created and initially setup your Azure OpenAI API with Azure API Management. You can now test your API by selecting the **Test** tab in Azure API Management. For more information, see [Test an API](#). Next, we'll configure authentication for your new API so only authorized users can access it.

Configure authentication for your new API

By default, your new Azure API Management service uses subscriptions to control access to your APIs, but this would be functionally equivalent to giving your users direct access to your API keys. So instead, we'll demonstrate how to setup OAuth 2.0 to control access to your Azure OpenAI APIs.

2. Configure OAuth 2.0 by following the steps in the [protect an API in Azure API Management using OAuth 2.0 authorization with Azure Active Directory](#) article. Once you're done with these steps you'll have an Azure AD application with the necessary scopes.
3. Next, [configure a JWT validation policy to pre-authorize requests](#). This will ensure that only users with the correct permissions can access your API and underlying Azure OpenAI API.
4. Lastly, turn off **Subscription required** for your API.

Navigate to your API Management instance and select **APIs** from the left-hand menu. Then select your API from the list of APIs. In the **Settings** tab, find the **Subscription required** checkbox and uncheck it. Only complete this step after configuring OAuth 2.0 and JWT validation to ensure that only authorized users can access your API.



Congrats, you've now configured authentication for your Azure OpenAI API with Azure API Management. You can now provide users with access to your API by creating users in Azure AD and assigning them the correct permissions. For more information, see [Assign a user or group to an enterprise app in Azure Active Directory](#).

Access your API from Semantic Kernel

Once users have the correct permissions, they can access your API from within Semantic Kernel. There are a few steps to complete to test this connection within a console application.

ⓘ Important

The following code is for illustrative purposes.

1. First, create a new class that implements the `TokenCredential` class. This class will be used to provide Semantic Kernel with the user authentication token for your service.

C#

```
using Azure.Core;

public class BearerTokenCredential : TokenCredential
{
    private readonly AccessToken _accessToken;

    // Constructor that takes a Bearer token string and its expiration
```

```
date
    public BearerTokenCredential(AccessToken accessToken)
    {
        _accessToken = accessToken;
    }

    public override AccessToken GetToken(TokenRequestContext
requestContext, CancellationToken cancellationToken)
    {
        return _accessToken;
    }

    public override ValueTask<AccessToken>
GetTokenAsync(TokenRequestContext requestContext, CancellationToken
cancellationToken)
    {
        return new ValueTask<AccessToken>(_accessToken);
    }
}
```

2. Next, within your console application, you'll create an interactive logon to get the user's token.

C#

```
string[] scopes = new string[] {
"https://cognitiveservices.azure.com/.default" };
var credential = new InteractiveBrowserCredential();
var requestContext = new TokenRequestContext(scopes);
var accessToken = await credential.GetTokenAsync(requestContext);
```

3. Finally, you can create a new instance of Kernel and pass in the `BearerTokenCredential` class you created in step 1 along with the access token you retrieved in step 2.

C#

```
IKernel kernel = new KernelBuilder()
.WithAzureTextCompletionService(
    "text-davinci-003",
    "https://apim...api.net/",
    new BearerTokenCredential(accessToken)
)
.Build();
```

Overview of sample apps

Article • 07/18/2023



Multiple learning samples are provided in the [Semantic Kernel GitHub repository](#) to help you learn core concepts of Semantic Kernel.

Requirements to run the apps

- ✓ [Azure Functions Core Tools](#) - used for running the kernel as a local API
- ✓ [Yarn](#) - used for installing the app's dependencies

Try the TypeScript/React sample apps

ⓘ Important

The local API service must be active for the sample apps to run.

Sample App	Illustrates
Simple chat summary	Use ready-to-use plugins and get those plugins into your app easily. <i>Be sure that the local API service is running for this sample to work.</i>
Book creator	Use planner to deconstruct a complex goal and envision using planner in your app. <i>Be sure that the local API service is running for this sample to work.</i>
Authentication and APIs	Use a basic plugin pattern to authenticate and connect to an API and imagine integrating external data into your app's LLM AI. <i>Be sure that the local API service is running for this sample to work.</i>
GitHub Repo Q&A Bot	Use embeddings to store local data and functions to question the embedded data. <i>Be sure that the local API service is running for this sample to work.</i>
Chat Copilot reference app	Build your own chatbot based on the Semantic Kernel.

Next step

[Run the simple chat summary app](#)

Local API service for app samples

Article • 05/23/2023



This service API is written in C# against Azure Function Runtime v4 and exposes some Semantic Kernel APIs that you can call via HTTP POST requests for the [learning samples](#).

ⓘ Important

Each function will call OpenAI which will use tokens that you will be billed for.

Walkthrough video

<https://aka.ms/SK-Local-API-Setup> ↗

Requirements to run the local service

- ✓ [Azure Functions Core Tools](#) - used for running the kernel as a local API

Running the service API locally

The [local API service](#) ↗ is located in the Semantic Kernel GitHub repository.

Run `func start --csharp` from the command line. This will run the service API locally at `http://localhost:7071`.

Two endpoints will be exposed by the service API:

- **InvokeFunction:** [POST]
`http://localhost:7071/api/skills/{skillName}/invoke/{functionName}`
- **Ping:** [GET] `http://localhost:7071/api/ping`

Take the next step

Now that your service API is running locally, try out some of the sample apps so you can learn core Semantic Kernel concepts!

[Jump into all the sample apps](#)

Simple chat summary sample app

Article • 05/23/2023



The Simple Chat Summary sample allows you to see the power of [functions](#) used in a chat sample app. The sample highlights the [Summarize ↗](#), [Topics ↗](#) and [Action Items ↗](#) functions in the [Summarize Plugin ↗](#). Each function calls OpenAI to review the information in the chat window and produces insights.

ⓘ Important

Each function will call OpenAI which will use tokens that you will be billed for.

Walkthrough video

[https://aka.ms/SK-Samples-SimChat-Video ↗](https://aka.ms/SK-Samples-SimChat-Video)

Requirements to run this app

- ✓ [Local API service](#) is running
- ✓ [Yarn ↗](#) - used for installing the app's dependencies

Running the app

The [Simple chat summary sample app ↗](#) is located in the Semantic Kernel GitHub repository.

1. Follow the [Setup](#) instructions if you do not already have a clone of Semantic Kernel locally.
2. Start the [local API service](#).
3. Open the ReadMe file in the Simple Chat Summary sample folder.
4. Open the Integrated Terminal window.
5. Run `yarn install` - if this is the first time you are running the sample. Then run `yarn start`.
6. A browser will open with the sample app running

Exploring the app

Setup Screen

Start by entering in your [OpenAI key ↗](#) or if you are using [Azure OpenAI Service](#) the key and endpoint. Then enter in the model you would like to use in this sample.

Interact Screen

A preloaded chat conversation is available. You can add additional items in the chat or modify the [chat thread ↗](#) before running the sample.

AI Summaries Screen

Three semantic functions are called on this screen

1. [Summarize ↗](#)
2. [Topics ↗](#)
3. [Action Items ↗](#)

Next step

[Run the book creator app](#)

Book creator sample app

Article • 05/23/2023



The Book creator sample allows you to enter in a topic then the [Planner](#) creates a plan for the functions to run based on the ask. You can see the plan along with the results. The [Writer Skill](#) functions are chained together based on the asks.

ⓘ Important

Each function will call OpenAI which will use tokens that you will be billed for.

Walkthrough video

<https://aka.ms/SK-Samples-CreateBook-Video>

Requirements to run this app

- ✓ [Local API service](#) is running
- ✓ [Yarn](#) - used for installing the app's dependencies

Running the app

The [Book creator sample app](#) is located in the Semantic Kernel GitHub repository.

1. Follow the [Setup](#) instructions if you do not already have a clone of Semantic Kernel locally.
2. Start the [local API service](#).
3. Open the ReadMe file in the Book creator sample folder.
4. Open the Integrated Terminal window.
5. Run `yarn install` - if this is the first time you are running the sample. Then run `yarn start`.
6. A browser will open with the sample app running

Exploring the app

Setup Screen

Start by entering in your [OpenAI key](#) or if you are using [Azure OpenAI Service](#) the key and endpoint. Then enter in the model you would like to use in this sample.

Topics Screen

On this screen you can enter in a topic for the children's book that will be created for you. This will use functions and AI to generate book ideas based on this topic.

Book Screen

By clicking on the asks, multiple steps will be found from the Planner and the process will run to return results.

Next step

[Run the authentication and API app](#)

Authentication and API calls sample app

Article • 05/23/2023

The Authenticated API's sample allows you to use authentication to connect to the Microsoft Graph using your personal account. If you don't have a Microsoft account or do not want to connect to it, you can review the code to see the patterns needed to call out to APIs. The sample highlights connecting to Microsoft Graph and calling APIs for Outlook, OneDrive, and ToDo. Each function will call Microsoft Graph and/or OpenAI to perform the tasks.

ⓘ Important

Each function will call OpenAI which will use tokens that you will be billed for.

Walkthrough video

[https://aka.ms/SK-Samples-AuthAPI-Video ↗](https://aka.ms/SK-Samples-AuthAPI-Video)

Requirements to run this app

- ✓ Local API service is running
- ✓ Yarn ↗ - used for installing the app's dependencies

Running the app

The [Authentication and API sample app ↗](#) is located in the Semantic Kernel GitHub repository.

1. Follow the [Setup](#) instructions if you do not already have a clone of Semantic Kernel locally.
2. Start the [local API service](#).
3. Open the ReadMe file in the Authentication and API sample folder.
4. You will need to register your application in the Azure Portal. Follow the steps to register your app [here](#).
 - Your Redirect URI will be <http://localhost:3000>
 - It is recommended you use the `Personal Microsoft accounts` account type for this sample

5. Once registered, copy the Application (client) ID from the Azure Portal and paste in the GUID into the `env` file next to `REACT_APP_GRAPH_CLIENT_ID=`
6. Open the Integrated Terminal window.
7. Run `yarn install` - if this is the first time you are running the sample. Then run `yarn start`.
8. A browser will open with the sample app running

Exploring the app

Your Info Screen

You can sign in with your Microsoft account by clicking 'Sign in with Microsoft'. This will give the sample app access to Microsoft Graph on your behalf and will be used for the functions to run on the Interact screen.

Setup Screen

Start by entering in your [OpenAI key](#) or if you are using [Azure OpenAI Service](#) the key and endpoint. Then enter in the model you would like to use in this sample.

Interact Screen

When you select each of the 3 actions, native functions will be called to perform actions through the Microsoft Graph API and connector.

The actions on this screen are:

1. Summarize and create a new Word document and save it to OneDrive
2. Get a shareable link and email the link to myself
3. Add a reminder to follow-up with the email sent above

Take the next step

[Run the GitHub Repo Q&A Bot app](#)

GitHub Repo Q&A Bot sample app

Article • 07/18/2023



The GitHub Repo Q&A Bot sample allows you to enter in a GitHub repo then those files are created as [embeddings](#). You can then question the stored files from the embedding local storage.

ⓘ Important

Each function will call OpenAI which will use tokens that you will be billed for.

Walkthrough video

[https://aka.ms/SK-GitHub-QA-Bot-Video ↗](https://aka.ms/SK-GitHub-QA-Bot-Video)

Requirements to run this app

- ✓ [Local API service](#) is running
- ✓ [Yarn ↗](#) - used for installing the app's dependencies

Running the app

The [GitHub Repo Q&A Bot sample app ↗](#) is located in the Semantic Kernel GitHub repository.

1. Follow the [Setup](#) instructions if you do not already have a clone of Semantic Kernel locally.
2. Start the [local API service](#).
3. Open the ReadMe file in the GitHub Repo Q&A Bot sample folder.
4. Open the Integrated Terminal window.
5. Run `yarn install` - if this is the first time you are running the sample. Then run `yarn start`.
6. A browser will open with the sample app running

Exploring the app

Setup Screen

Start by entering in your [OpenAI key](#) or if you are using [Azure OpenAI Service](#) the key and endpoint. Then enter in the model for completion and embeddings you would like to use in this sample.

GitHub Repository Screen

On this screen you can enter in public GitHub repo and the sample will download the repo using a function and add the files as embeddings.

Q&A Screen

By default the Markdown files are stored as embeddings. You can ask questions in the chat and get answers based on the embeddings.

Next step

Run the Chat Copilot reference app!

[Run the Chat Copilot reference app](#)

Chat Copilot

Article • 07/18/2023



The Chat Copilot reference app allows you to build your own integrated large language model chatbot. This is an enriched intelligence app, with multiple dynamic components including command messages, user intent, and memories.

The chat prompt and response will evolve as the conversation between the user and the application proceeds. This chat experience uses a chat plugin containing multiple functions that work together to construct the final prompt for each exchange.

ⓘ Important

Each function will call OpenAI which will use tokens that you will be billed for.

Requirements to run this app

- ✓ [Visual Studio Code](#)
- ✓ [Git](#)
- ✓ [.NET 6.0](#)
- ✓ [Node.js](#)
- ✓ [Yarn](#)

Running the app

The [Chat Copilot reference app](#) is located in the Semantic Kernel GitHub repository.

1. Clone [Semantic Kernel](#) locally.
2. To enable authentication, [register an Azure Application](#). We recommend using the following properties:
 - Select **Single-page application (SPA)** as platform type, and set the Web redirect URI to `http://localhost:3000`
 - Select **Accounts in any organizational directory and personal Microsoft Accounts** as supported account types for this sample.

(!) Note

Make a note of the Application (client) ID from the Azure Portal; we will use it in step 4.

3. Install requirements

Windows

Open a PowerShell terminal as an administrator and navigate to the `/samples/apps/copilot-chat-app/scripts` directory in the Semantic Kernel project.

PowerShell

```
cd /samples/apps/copilot-chat-app/scripts
```

Next, run the following command to install the required dependencies:

PowerShell

```
./Install-Requirements.ps1
```

4. Run the configuration script

PowerShell

If you are using Azure OpenAI, run the following command. Replace the `{AZURE_OPENAI_ENDPOINT}`, `{AZURE_OPENAI_API_KEY}`, and `{APPLICATION_CLIENT_ID}` values in the following command before running it:

PowerShell

```
./Configure.ps1 -AzureOpenAI -Endpoint {AZURE_OPENAI_ENDPOINT} -  
ApiKey {AZURE_OPENAI_API_KEY} -ClientId {APPLICATION_CLIENT_ID}
```

If you are using OpenAI, run the following command. Replace the `{OPENAI_API_KEY}` and `{APPLICATION_CLIENT_ID}` values in the following command before running it:

PowerShell

```
./Configure.ps1 -openai -ApiKey {OPENAI_API_KEY} -ClientId  
{APPLICATION_CLIENT_ID}
```

5. Run the start script



```
PowerShell

PowerShell

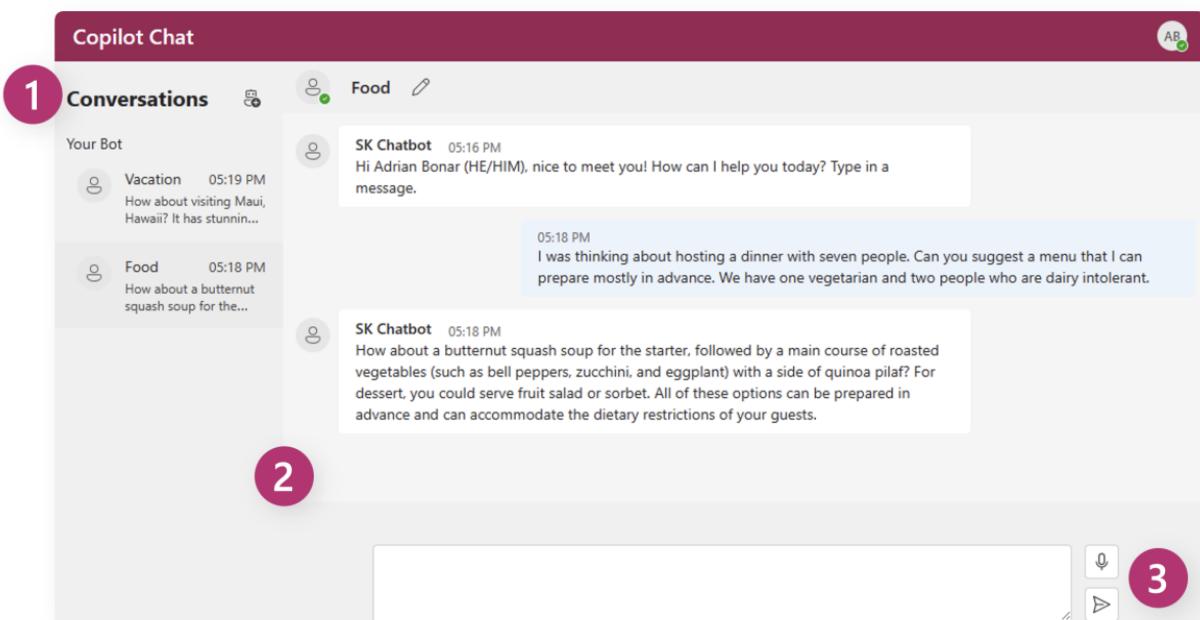
setx ASPNETCORE_ENVIRONMENT "Development"

./Start.ps1
```

6. Congrats! A browser should automatically launch and navigate to <https://localhost:3000> with the sample app running.

Exploring the app

With the Chat Copilot reference app running, you can start interacting with the chatbot. The app will start with a default conversation thread where you can ask it questions.



Feature	Name	Description
1	Conversation Pane	The left portion of the screen shows different conversation threads the user is holding with the chatbot. To start a new conversation, click the '+Bot' symbol.

Feature	Name	Description
2	Conversation Thread	Chatbot responses will appear in the main conversation thread, along with a history of your prompts. Users can scroll up and down to review a complete conversation history.
3	Prompt Entry Box	The bottom of the screen contains the prompt entry box, where users can type their prompts, and click the "Send" icon to the right of the box when ready to send it to the bot.

Next step

If you've tried all the apps and are excited to see more, please star the GitHub repo and join the Semantic Kernel community!

[Star the Semantic Kernel repo](#)

Visual Code Studio Semantic Kernel Extension

Article • 05/23/2023

The Semantic Kernel Tools help developers to write semantic functions for [Semantic Kernel](#).



ⓘ Note

Skills are currently being renamed to plugins. This article has been updated to reflect the latest terminology, but some images and code samples may still refer to skills.

These tools simplify Semantic Kernel development

With the Semantic Kernel Tools, you can easily create new semantic functions and test them without needing to write any code. Behind the scenes, the tools use the Semantic Kernel SDK so you can easily transition from using the tools to integrating your semantic functions into your own code.

In the following image you can see how a user can easily view all of their semantic functions, edit them, and run them from within Visual Studio Code using any of the supported AI endpoints.

```

skprompt.txt 2 ×
skills > FunSkill > Joke > skprompt.txt
1 WRITE EXACTLY ONE JOKE OR HUMOROUS STORY ABOUT THE TOPIC BELOW
2
3 JOKE MUST BE:
4 - G RATED
5 - WORKPLACE/FAMILY SAFE
6 NO SEXISM, RACISM OR OTHER BIAS/BIGOTRY
7
8 BE CREATIVE AND FUNNY. I WANT TO LAUGH.
9 {{style}}
10 +++++
11
12 {{$input}}
13 +++++
14

```

PROBLEMS 3 OUTPUT DEBUG CONSOLE ... Semantic Kernel

Result:
Q: What do you call a dinosaur that never gives up?
A: A Try-Raptor!

Tokens:
Input tokens: 92
Output tokens: 23
Total: 115

Duration:
00:00:03.585

===== Function execution was finished. =====

Installing the Semantic Kernel Extension

To get started with Semantic Kernel Tools, follow these simple steps:

1. Ensure that you have [Visual Studio Code](#) installed on your computer.
2. Open Visual Studio Code and press Shift+Control+X to bring up the [Extensions marketplace](#).
3. In the Extensions menu, search for "[Semantic Kernel Tools](#)".
4. Select Semantic Kernel Tools from the search results and click the Install button.



5. Wait for the installation to complete, then restart Visual Studio Code.

Connecting the extension to your AI endpoint

First you must configure an AI endpoint to be used by the Semantic Kernel

Open the Command Palette i.e., View -> Command Palette or Ctrl+Shift+P

1. If you have access to an Azure subscription that contains the Azure OpenAI resource you want to use:
 - Endpoint type, select "completion"
 - Allow the extension to sign in to Azure Portal
 - Select the subscription to use
 - Select the resource group which contains the Azure OpenAI resource
 - Select the Azure OpenAI resource
 - Select the Azure OpenAI model
2. Type "Add Azure OpenAI Endpoint" and you will be prompted for the following information:
 - Type "Add AI Endpoint" and you will be prompted for the following information:
 - Endpoint type, select "completion"
 - Completion label, the default of "Completion" is fine
 - Completion AI Service, select AzureOpenAI
 - Completion deployment or model id e.g., text-davinci-003
 - Completion endpoint URI e.g., <https://contoso-openai.azure.com/>
 - Completion endpoint API key (this will be stored in VS Code secure storage)
3. If you have the details of an Azure OpenAI endpoint that you want to use
 - Type "Add AI Endpoint" and you will be prompted for the following information:
 - Endpoint type, select "completion"
 - Completion label, the default of "Completion" is fine
 - Completion AI Service, select OpenAI
 - Completion deployment or model id e.g., text-davinci-003
 - Completion endpoint API key (this will be stored in VS Code secure storage)
4. If you have the details of an OpenAI endpoint that you want to use
 - Type "Add AI Endpoint" and you will be prompted for the following information:
 - Endpoint type, select "completion"
 - Completion label, the default of "Completion" is fine
 - Completion AI Service, select OpenAI
 - Completion deployment or model id e.g., text-davinci-003
 - Completion endpoint API key (this will be stored in VS Code secure storage)

Once you have a AI endpoint configured proceed as follows:

1. Select the semantic function you want to execute
2. Select the "Run Function" icon which is shown in the Functions view
3. You will be prompted to enter any arguments the semantic function requires
4. The response will be displayed in the Output view in the "Semantic Kernel" section

Create a semantic function

1. Once you have installed the Semantic Kernel Tools extension you will see a new Semantic Kernel option in the activity bar

- We recommend you clone the semantic-kernel repository and open this in your VS Code workspace
2. Click the Semantic Kernel icon to open Semantic Kernel Functions view
 3. Click the "Add Semantic Skill" icon in the Semantic Kernel Functions view title bar
 4. You will be prompted to select a folder
- This will be the location of the plugin which will contain your new Semantic Function
 - Create a new folder called MyPlugin1 in this directory <location of your clone>\semantic-kernel\samples\skills
 - Select this new folder as your Plugin folder
5. You will be prompted for a function name, enter MyFunction1
 6. You will be prompted for a function description2
 7. A new prompt text file will be automatically created for your new function
 8. You can now enter your prompt.

Troubleshooting

- Enabling Trace Level Logs
- You can enable trace level logging for the Semantic Kernel using the following steps:
 1. Open settings (Ctrl + ,)
 2. Type "Semantic Kernel"
 3. Select Semantic Kernel Tools -> Configuration
 4. Change the log level to "Trace"
 5. Repeat the steps to execute a semantic function and this time you should see trace level debugging of the semantic kernel execution

Below is a list of possible errors you might receive and details on how to address them.

Errors creating a Semantic Function

- Unable to create function prompt file for <name>
- An error occurred creating the skprompt.txt file for a new semantic function. Check you can create new folders and files in the location specified for the semantic function.
- Function <name> already exists. Found function prompt file: <file name>
- A skprompt.txt file already exists for the semantic function you are trying to create. Switch to the explorer view to find the conflicting file.
- Unable to create function configuration file for <file name>

- An error occurred creating the config.json file for a new semantic function. Check you can create new folders and files in the location specified for the semantic function.
- Configuration file for <file> already exists. Found function config file: <file name>
- A config.json file already exists for the semantic function you are trying to create. Switch to the explorer view to find the conflicting file.

Errors configuring an AI Endpoint

- Unable to find any subscriptions. Please log in with a user account that has access to a subscription where OpenAI resources have been deployed.
- The user account you specified to use when logging in to Microsoft does not have access to any subscriptions. Please try again with a different account.
- Unable to find any resource groups. Please log in with a user account that has access to a subscription where OpenAI resources have been deployed.
- The user account you specified to use when logging in to Microsoft does not have access to any resource groups in the subscription you selected. Please try again with a different account or a different subscription.
- Unable to find any OpenAI resources. Please log in with a user account that has access to a subscription where OpenAI resources have been deployed.
- The user account you specified to use when logging in to Microsoft does not have access to any Azure OpenAI resources in the resource group you selected. Please try again with a different account or a different resource group.
- Unable to find any OpenAI model deployments. Please log in with a user account that has access to a subscription where OpenAI model deployments have been deployed.
- The user account you specified to use when logging in to Microsoft does not have access to any deployment models in the Azure OpenAI resource you selected. Please try again with a different account or a different Azure OpenAI resource.
- Unable to access the Azure OpenAI account. Please log in with a user account that has access to an Azure OpenAI account.
- The user account you specified to use when logging in to Microsoft does not have access to the Azure OpenAI account in the Azure OpenAI resource you selected. Please try again with a different account or a different Azure OpenAI resource.
- Unable to access the Azure OpenAI account keys. Please log in with a user account that has access to an Azure OpenAI account.
- The user account you specified to use when logging in to Microsoft does not have access to the Azure OpenAI account keys in the Azure OpenAI resource you selected. Please try again with a different account or a different Azure OpenAI resource.

- Settings does not contain a valid AI completion endpoint configuration. Please run the "Add Azure OpenAI Endpoint" or "Add AI Endpoint" command to configure a valid endpoint.
- You have not configured an AI endpoint. Please refer to the first part of the Execute a Semantic Function section above.

Errors executing a Semantic Function

- ModelNotAvailable – unable to fetch the list of model deployments from Azure (Unauthorized)
- This failure comes when calling the Azure OpenAI REST API. Check the AzureOpenAI resource you are using is correctly configured.

Take the next step

Now you can start building your own Semantic Functions

It all starts with an ask

Support for Semantic Kernel

Article • 04/06/2023



👋 Welcome! There are a variety of ways to get supported in the Semantic Kernel (SK) world.

Your preference	What's available
Read the docs	This learning site ↗ is the home of the latest information for developers
Visit the repo	Our open-source GitHub repository ↗ is available for perusal and suggestions
Realtime chat	Visit our Discord channel ↗ to get supported quickly with our CoC actively enforced
Realtime video	We will be hosting regular office hours that will be announced in our Discord channel ↗

More support information

- [Frequently Asked Questions \(FAQs\)](#)
- [Hackathon Materials](#)
- [Code of Conduct](#)

Next step

[Run the samples](#)

Semantic Kernel FAQ's

Article • 05/23/2023



Why is the Kernel only in C# and Python?

Both C# and Python are popular coding languages and we're actively adding additional languages based on community feedback. Both [Java](#) and [TypeScript](#) are on our roadmap and being actively developed in experimental branches.

Where are the sample plugins?

We have [sample apps](#) and plugins you can try out so you can quickly learn the concepts of Semantic Kernel.

How do I get help or provide feedback?

There are a variety of [support options available!](#)

Is something up with my OpenAI or Azure OpenAI key?

Depending upon the model you are trying to access, there may be times when your key may not work because of high demand. Or, because your access to the model is limited by the plan you're currently signed up for — so-called "throttling". In general, however, your key will work according to the plan agreement with your LLM AI provider.

Why aren't my Jupyter notebooks coming up in my VSCode or Visual Studio?

First of all, you'll need to be running locally on your own machine to interact with the Jupyter notebooks. If you've already cleared that hurdle, then all you need to do is to install the [Polyglot Extension](#) which requires .NET 7 to be installed. For complete information on the latest release of Polyglot Extension you can learn more [here](#).

Next step

[Get more support](#)

Contributor Covenant Code of Conduct

Article • 05/23/2023



Our Pledge

In the interest of fostering an open and welcoming environment, we as owners, contributors and maintainers pledge to making participation in our project and our community a harassment-free experience for everyone, regardless of age, body size, disability, ethnicity, gender identity and expression, level of experience, nationality, personal appearance, race, religion, or sexual identity and orientation.

Our Standards

Examples of behavior that contributes to creating a positive environment include:

- Using welcoming and inclusive language
- Being respectful of differing viewpoints and experiences
- Gracefully accepting constructive criticism
- Focusing on what is best for the community
- Showing empathy towards other community members

Examples of unacceptable behavior by participants include:

- The use of sexualized language or imagery and unwelcome sexual attention or advances
- Trolling, insulting/derogatory comments, and personal or political attacks
- Public or private harassment
- Publishing others' private information, such as a physical or electronic address, without explicit permission
- Other conduct which could reasonably be considered inappropriate in a professional setting

Our Responsibilities

Project maintainers are responsible for clarifying the standards of acceptable behavior and are expected to take appropriate and fair corrective action in response to any

instances of unacceptable behavior.

Project maintainers have the right and responsibility to remove, edit, or reject comments, commits, code, wiki edits, issues, and other contributions that are not aligned to this Code of Conduct, or to ban temporarily or permanently any contributor for other behaviors that they deem inappropriate, threatening, offensive, or harmful.

Scope

This Code of Conduct applies both within project spaces and in public spaces when an individual is representing the project or its community. Examples of representing a project or community include using an official project e-mail address, posting via an official social media account, or acting as an appointed representative at an online or offline event. Representation of a project may be further defined and clarified by project maintainers.

Enforcement

Instances of abusive, harassing, or otherwise unacceptable behavior may be reported by contacting the project team by using #moderation in the [Discord community](#). The project team will review and investigate all complaints, and will respond in a way that it deems appropriate to the circumstances. The project team is obligated to maintain confidentiality with regard to the reporter of an incident. Further details of specific enforcement policies may be posted separately.

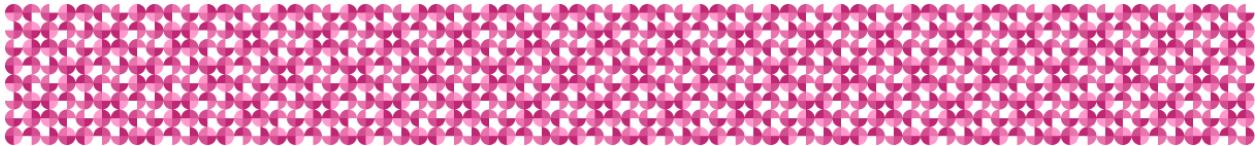
Project maintainers who do not follow or enforce the Code of Conduct in good faith may face temporary or permanent repercussions as determined by other members of the project's leadership.

Attribution

This Code of Conduct is adapted from the [Contributor Covenant]
[<https://www.contributor-covenant.org/>], version 1.4, available [here](#).

Schillace Laws of Semantic AI

Article • 05/23/2023



Consider the future of this decidedly "semantic" AI

The "Schillace Laws" were formulated after working with a variety of Large Language Model (LLM) AI systems to date. Knowing them will accelerate your journey into this exciting space of reimagining the future of software engineering. Welcome!

- 1. Don't write code if the model can do it; the model will get better, but the code won't.** The overall goal of the system is to build very high leverage programs using the LLM's capacity to plan and understand intent. It's very easy to slide back into a more imperative mode of thinking and write code for aspects of a program. Resist this temptation – to the degree that you can get the model to do something reliably now, it will be that much better and more robust as the model develops.
- 2. Trade leverage for precision; use interaction to mitigate.** Related to the above, the right mindset when coding with an LLM is not "let's see what we can get the dancing bear to do," it's to get as much leverage from the system as possible. For example, it's possible to build very general patterns, like "build a report from a database" or "teach a year of a subject" that can be parameterized with plain text prompts to produce enormously valuable and differentiated results easily.
- 3. Code is for syntax and process; models are for semantics and intent.** There are lots of different ways to say this, but fundamentally, the models are stronger when they are being asked to reason about meaning and goals, and weaker when they are being asked to perform specific calculations and processes. For example, it's easy for advanced models to write code to solve a sudoku generally, but hard for them to solve a sudoku themselves. Each kind of code has different strengths and it's important to use the right kind of code for the right kind of problem. The boundaries between syntax and semantics are the hard parts of these programs.
- 4. The system will be as brittle as its most brittle part.** This goes for either kind of code. Because we are striving for flexibility and high leverage, it's important to not

hard code anything unnecessarily. Put as much reasoning and flexibility into the prompts and use imperative code minimally to enable the LLM.

5. **Ask Smart to Get Smart.** Emerging LLM AI models are incredibly capable and "well educated" but they lack context and initiative. If you ask them a simple or open-ended question, you will get a simple or generic answer back. If you want more detail and refinement, the question has to be more intelligent. This is an echo of "Garbage in, Garbage out" for the AI age.
6. **Uncertainty is an exception throw.** Because we are trading precision for leverage, we need to lean on interaction with the user when the model is uncertain about intent. Thus, when we have a nested set of prompts in a program, and one of them is uncertain in its result ("One possible way...") the correct thing to do is the equivalent of an "exception throw" - propagate that uncertainty up the stack until a level that can either clarify or interact with the user.
7. **Text is the universal wire protocol.** Since the LLMs are adept at parsing natural language and intent as well as semantics, text is a natural format for passing instructions between prompts, modules and LLM based services. Natural language is less precise for some uses, and it is possible to use structured language like XML sparingly, but generally speaking, passing natural language between prompts works very well, and is less fragile than more structured language for most uses. Over time, as these model-based programs proliferate, this is a natural "future proofing" that will make disparate prompts able to understand each other, the same way humans do.
8. **Hard for you is hard for the model.** One common pattern when giving the model a challenging task is that it needs to "reason out loud." This is fun to watch and very interesting, but it's problematic when using a prompt as part of a program, where all that is needed is the result of the reasoning. However, using a "meta" prompt that is given the question and the verbose answer and asked to extract just the answer works quite well. This is a cognitive task that would be easier for a person (it's easy to imagine being able to give someone the general task of "read this and pull out whatever the answer is" and have that work across many domains where the user had no expertise, just because natural language is so powerful). *So, when writing programs, remember that something that would be hard for a person is likely to be hard for the model, and breaking patterns down into easier steps often gives a more stable result.*
9. **Beware "pareidolia of consciousness"; the model can be used against itself.** It is very easy to imagine a "mind" inside an LLM. But there are meaningful differences between human thinking and the model. An important one that can be exploited is

that the models currently don't remember interactions from one minute to the next. So, while we would never ask a human to look for bugs or malicious code in something they had just personally written, we can do that for the model. It might make the same kind of mistake in both places, but it's not capable of "lying" to us because it doesn't know where the code came from to begin with. _This means we can "use the model against itself" in some places – it can be used as a safety monitor for code, a component of the testing strategy, a content filter on generated content, etc. _

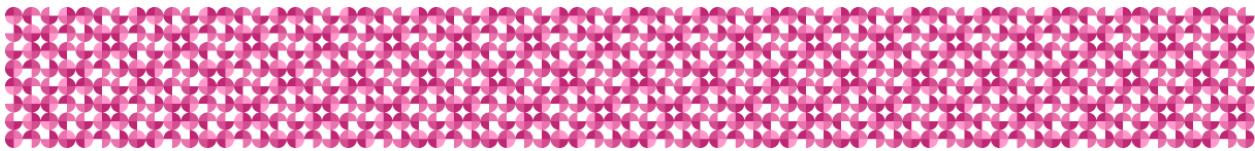
Take the next step

If you're interested in LLM AI models and feel inspired by the Schillace Laws, be sure to visit the Semantic Kernel GitHub repository and add a star to show your support!

[Go to the SK GitHub repository](#)

Responsible AI and Semantic Kernel

Article • 05/23/2023



What is a Transparency Note?

An AI system includes not only the technology, but also the people who will use it, the people who will be affected by it, and the environment in which it is deployed. Creating a system that is fit for its intended purpose requires an understanding of how the technology works, what its capabilities and limitations are, and how to achieve the best performance. Microsoft's Transparency Notes are intended to help you understand how our AI technology works, the choices system owners can make that influence system performance and behavior, and the importance of thinking about the whole system, including the technology, the people, and the environment. You can use Transparency Notes when developing or deploying your own system, or share them with the people who will use or be affected by your system.

Microsoft's Transparency Notes are part of a broader effort at Microsoft to put our AI Principles into practice. To find out more, see the [Microsoft AI principles](#).

Introduction to Semantic Kernel

Semantic Kernel (SK) is a lightweight SDK that lets you easily mix conventional programming languages with the latest in Large Language Model (LLM) AI "prompts" with templating, chaining, and planning capabilities out-of-the-box.

The basics of Semantic Kernel

Semantic Kernel (SK) builds upon the following five concepts:

Concept	Short Description
Kernel	The kernel orchestrates a user's ASK expressed as a goal
Planner	Planner breaks it down into steps based upon resources that are available
Plugins	Plugins are customizable resources built from LLM AI prompts and native code

Concept	Short Description
Memories	Memories are customizable resources that manage contextual information
Connectors	Connectors are customizable resources that enable external data access

Use cases for LLM AI

Intended uses

The general intended uses include:

- Chat and conversation interaction: Users can interact with a conversational agent that responds with responses drawn from trusted documents such as internal company documentation or tech support documentation; conversations must be limited to answering scoped questions.
- Chat and conversation creation: Users can create a conversational agent that responds with responses drawn from trusted documents such as internal company documentation or tech support documentation; conversations must be limited to answering scoped questions.
- Code generation or transformation scenarios: For example, converting one programming language to another, generating docstrings for functions, converting natural language to SQL.
- Journalistic content: For use to create new journalistic content or to rewrite journalistic content submitted by the user as a writing aid for pre-defined topics. Users cannot use the application as a general content creation tool for all topics. May not be used to generate content for political campaigns.
- Question-answering: Users can ask questions and receive answers from trusted source documents such as internal company documentation. The application does not generate answers ungrounded in trusted source documentation.
- Reason over structured and unstructured data: Users can analyze inputs using classification, sentiment analysis of text, or entity extraction. Examples include analyzing product feedback sentiment, analyzing support calls and transcripts, and refining text-based search with embeddings.
- Search: Users can search trusted source documents such as internal company documentation. The application does not generate results ungrounded in trusted source documentation.
- Summarization: Users can submit content to be summarized for pre-defined topics built into the application and cannot use the application as an open-ended

summarizer. Examples include summarization of internal company documentation, call center transcripts, technical reports, and product reviews.

- Writing assistance on specific topics: Users can create new content or rewrite content submitted by the user as a writing aid for business content or pre-defined topics. Users can only rewrite or create content for specific business purposes or pre-defined topics and cannot use the application as a general content creation tool for all topics. Examples of business content include proposals and reports. For journalistic use, see above Journalistic content use case.

Considerations when choosing a use case for LLM AI

There are some considerations:

- Not suitable for open-ended, unconstrained content generation. Scenarios where users can generate content on any topic are more likely to produce offensive or harmful text. The same is true of longer generations.
- Not suitable for scenarios where up-to-date, factually accurate information is crucial unless you have human reviewers or are using the models to search your own documents and have verified suitability for your scenario. The service does not have information about events that occur after its training date, likely has missing knowledge about some topics, and may not always produce factually accurate information.
- Avoid scenarios where use or misuse of the system could result in significant physical or psychological injury to an individual. For example, scenarios that diagnose patients or prescribe medications have the potential to cause significant harm.
- Avoid scenarios where use or misuse of the system could have a consequential impact on life opportunities or legal status. Examples include scenarios where the AI system could affect an individual's legal status, legal rights, or their access to credit, education, employment, healthcare, housing, insurance, social welfare benefits, services, opportunities, or the terms on which they are provided.
- Avoid high stakes scenarios that could lead to harm. Each LLM AI model reflects certain societal views, biases and other undesirable content present in the training data or the examples provided in the prompt. As a result, we caution against using the models in high-stakes scenarios where unfair, unreliable, or offensive behavior might be extremely costly or lead to harm.
- Carefully consider use cases in high stakes domains or industry: Examples include but are not limited to healthcare, medicine, finance or legal.

- Carefully consider well-scoped chatbot scenarios. Limiting the use of the service in chatbots to a narrow domain reduces the risk of generating unintended or undesirable responses.
- Carefully consider all generative use cases. Content generation scenarios may be more likely to produce unintended outputs and these scenarios require careful consideration and mitigations.

Characteristics and limitations of LLM AI

When it LLM AI models, there are particular fairness and responsible AI issues to consider. People use language to describe the world and to express their beliefs, assumptions, attitudes, and values. As a result, publicly available text data typically used to train large-scale natural language processing models contains societal biases relating to race, gender, religion, age, and other groups of people, as well as other undesirable content. These societal biases are reflected in the distributions of words, phrases, and syntactic structures.

Evaluating and integrating Semantic Kernel for your use

When getting ready to deploy any AI-powered products or features, the following activities help to set you up for success:

- Understand what it can do: Fully assess the capabilities of any AI system you are using to understand its capabilities and limitations. Understand how it will perform in your particular scenario and context by thoroughly testing it with real life conditions and data.
- Respect an individual's right to privacy: Only collect data and information from individuals for lawful and justifiable purposes. Only use data and information that you have consent to use for this purpose.
- Legal review: Obtain appropriate legal advice to review your solution, particularly if you will use it in sensitive or high-risk applications. Understand what restrictions you might need to work within and your responsibility to resolve any issues that might come up in the future. Do not provide any legal advice or guidance.
- Human-in-the-loop: Keep a human-in-the-loop and include human oversight as a consistent pattern area to explore. This means ensuring constant human oversight of the AI-powered product or feature, and maintaining the role of humans in decision making. Ensure you can have real-time human intervention in the solution

to prevent harm. This enables you to manage situations when the AI model does not perform as required.

- Security: Ensure your solution is secure and has adequate controls to preserve the integrity of your content and prevent unauthorized access.
- Build trust with affected stakeholders: Communicate the expected benefits and potential risks to affected stakeholders. Help people understand why the data is needed and how the use of the data will lead to their benefit. Describe data handling in an understandable way.
- Customer feedback loop: Provide a feedback channel that allows users and individuals to report issues with the service once it's been deployed. Once you've deployed an AI-powered product or feature it requires ongoing monitoring and improvement -- be ready to implement any feedback and suggestions for improvement. Establish channels to collect questions and concerns from affected stakeholders (people who may be directly or indirectly impacted by the system, including employees, visitors, and the general public). Examples of such channels are:
 - Feedback features built into app experiences, An easy-to-remember email address for feedback, Anonymous feedback boxes placed in semi-private spaces, and Knowledgeable representatives on site. Feedback: Seek out feedback from a diverse sampling of the community during the development and evaluation process (for example, historically marginalized groups, people with disabilities, and service workers). See: [Community Jury](#).
 - User Study: Any consent or disclosure recommendations should be framed in a user study. Evaluate the first and continuous-use experience with a representative sample of the community to validate that the design choices lead to effective disclosure. Conduct user research with 10-20 community members (affected stakeholders) to evaluate their comprehension of the information and to determine if their expectations are met.

Learn more about responsible AI

- Microsoft responsible AI resources
- [Microsoft Azure Learning course on responsible AI](#)