```
        print("*" * 20)
```

Modules that influence or process streaming results, such as `OutputFiltering`, might need specific stream options. The `overlap` option allows you to include extra context during the filtering process:

```python
config = OrchestrationConfig(
    template=Template(
        messages=[
            SystemMessage("You are a helpful AI assistant."),
            UserMessage("{{?text}}"),
        ]
    ),
    llm=LLM(
        name="gpt-4o-mini",
        parameters={
            "max_completion_tokens": 256,
            "temperature": 0.0
        }
    ),
    output_filtering=OutputFiltering(
        filters=[AzureContentFilter(
            hate=AzureThreshold.ALLOW_ALL,
            violence=AzureThreshold.ALLOW_ALL,
            self_harm=AzureThreshold.ALLOW_ALL,
            sexual=AzureThreshold.ALLOW_ALL
            )
        ],
        stream_options={ 'overlap': 100 }
    )
)

response = service.stream(
    config=config,
    template_values=[
        TemplateValue(name="text", value="Why is the sky blue?")
    ]
)

for chunk in response:
    print(chunk.orchestration_result.choices[0].delta.content, end='')
```

# Tool Calling (Function Calling)

The Orchestration Service supports **tool calling**, which allows large language models (LLMs) to request the execution of external operations such as Python functions, API calls, or other tools as part of their workflow.

This feature enables you to build advanced AI solutions that can perform calculations, access data, or interact with external systems in response to user input.

# Defining Tools

You can define tools in several ways, depending on your requirements and the level of control you need.

## Using the Python Decorator

The simplest way to define a tool is to decorate a Python function with `@function_tool()`. The function's signature and docstring are used to describe the tool to the LLM.

```python
from gen_ai_hub.orchestration.models.tools import function_tool

@function_tool()
def multiply(a: int, b: int) -> int:
    """Multiply two numbers."""
    return a * b

@function_tool()
def add(a: int, b: int) -> int:
    """Add two numbers."""
    return a + b

tools = [multiply, add]
```

## Using the `FunctionTool` Class

For more control, you can use the `FunctionTool` class directly. This is useful if you want to customize the schema, enable strict argument checking, or wrap an existing function.

```python
from gen_ai_hub.orchestration.models.tools import FunctionTool
```

```python
def get_weather(location: str) -> str:
    """Get current temperature for a given location."""
    # Replace with your actual implementation
    return "22°C"

weather_tool = FunctionTool(
    name="get_weather",
    description="Get current temperature for a given location.",
    parameters={
        "type": "object",
        "properties": {
            "location": {
                "type": "string",
                "description": "City and country e.g. Bogotá, Colombia"
            }
        },
        "required": ["location"],
        "additionalProperties": False
    },
    strict=True,
    function=get_weather
)

tools = [weather_tool]
```

You can also create a `FunctionTool` from a function using the `from_function` static method:

```python
weather_tool = FunctionTool.from_function(get_weather, strict=True)
tools = [weather_tool]
```

## Using a JSON Schema Dictionary

You can define a tool directly as a JSON schema dictionary. This is useful if you want to specify the tool interface without implementing the function in Python, or if you want to integrate with external systems.

```python
tools = [{
    "type": "function",
    "function": {
        "name": "get_weather",
        "description": "Get current temperature for a given location.",
        "parameters": {
            "type": "object"
```

```
            "location": {
                "type": "string",
                "description": "City and country e.g. Bogotá, Colomb
            }
        },
        "required": [
            "location"
        ],
        "additionalProperties": False
    },
    "strict": True
  }
}]
```

You can then attach any of these tool definitions to your template:

```
from gen_ai_hub.orchestration.models.template import Template
from gen_ai_hub.orchestration.models.message import SystemMessage, UserM

                =           (
    messages=[
        SystemMessage("You are a weather assistant."),
        UserMessage("What is the temperature in {{?location}}?"),
    ],
    tools=tools,
)
```

# Synchronous Tool Call Workflow

When the LLM decides to call a tool, the orchestration response will include a
`tool_calls` field. You are responsible for executing the tool(s), adding the results to the
conversation history, and running the orchestration again to get the final answer.

```
from typing import List
from gen_ai_hub.orchestration.models.config import OrchestrationConfig
from gen_ai_hub.orchestration.models.template import TemplateValue
from gen_ai_hub.orchestration.models.message import Message, ToolMessage
from gen_ai_hub.orchestration.models.llm import LLM
from gen_ai_hub.orchestration.service import OrchestrationService

# Assume 'template' and 'weather_tool' are defined as above
```

```python
config = OrchestrationConfig(template=template, llm=llm)
template_values = [TemplateValue(name="location", value="Bogotá, Colombi

# First run: triggers tool call
service = OrchestrationService()
response = service.run(config=config, template_values=template_values)
tool_calls = response.orchestration_result.choices[0].message.tool_calls

# Execute tool(s) and build new history
history: List[Message] = []
history.extend(response.module_results.templating)
history.append(response.orchestration_result.choices[0].message)

for tool_call in tool_calls:
    # For FunctionTool, use .execute(**tool_call.function.parse_argument
    result = weather_tool.execute(**tool_call.function.parse_arguments()
    tool_message = ToolMessage(
        content=str(result),
        tool_call_id=tool_call.id,
    )
    history.append(tool_message)

# Second run: LLM receives tool result and produces final answer
response2 = service.run(
    config=config,
    template_values=template_values,
    history=history,
)
print(response2.orchestration_result.choices[0].message.content)
```

## Streaming Tool Calls

When using streaming, tool calls may be split across multiple chunks. The
`delta.tool_calls` field in each chunk contains partial or complete tool call information.
You may need to buffer and concatenate arguments if they arrive in pieces.

```python
# Assume 'config' and 'service' are defined as above
stream = service.stream(config=config, template_values=template_values)
final_tool_calls = {}

for chunk in stream:
    for tool_call in chunk.orchestration_result.choices[0].delta.tool_ca
        index = tool_call.index
        if index not in final_tool_calls:
            final_tool_calls[index] = tool_call
```

```
            # Concatenate arguments if split across chunks
            final_tool_calls[index].function.arguments += tool_call.func

    # Now final_tool_calls contains all tool calls with complete arguments
```

⚠️ **Note on Agentic Loop Support:**

> *The current SDK **does not provide built-in abstractions or convenience methods
> for managing the agentic loop** (the process of automatically handling tool call
> detection, execution, and iterative orchestration until a final answer is produced).
> As a user, you are responsible for:*

> * *Executing the corresponding Python functions*
>
> * *Appending tool results to the conversation history (as `ToolMessage`)*
>
> * *Re-invoking the orchestration service as needed*
> *This approach gives you maximum flexibility, but you must implement the
> orchestration loop logic yourself.*

# Using Images as Input

The Orchestration Service supports multimodal prompts, enabling you to include images
alongside text in your messages. This powerful feature unlocks a variety of applications,
such as visual question answering (VQA), image captioning, object recognition, and
generating text creatively based on visual input.

This guide details how to prepare image inputs, integrate them into your prompts, and
execute the orchestration to get insightful responses.

## 1. Preparing Image Inputs

To use an image, you first need to represent it as an `ImageItem` object. The
`gen_ai_hub.orchestration.models.multimodal_items.ImageItem` class provides two
convenient ways to do this:

## a) From a URL or Data URL

This method is ideal for images hosted online or when you have the image data encoded as a Data URL (base64 encoded).

- **Standard URL:** Provide a direct web link to the image file.

- **Data URL:** Provide the image data directly embedded in the URL string.

**Note:** For web URLs, ensure the image is publicly accessible, as the service will need to fetch it.

```python
from gen_ai_hub.orchestration.models.multimodal_items import ImageItem

# Example 1: Image from a standard, publicly accessible URL
# Ensure the URL points directly to the image file (e.g., .png, .jpg, ..
image_from_web = ImageItem(url="https://picsum.photos/id/1/200/300")  #

# Example 2: Image from a Data URL (base64-encoded)
# This is useful when you have the image content as a string.
# The format is "data:[<mediatype>][;base64],<data>"
image_from_data_url = ImageItem(
    url="data:image/png;base64,iVBORw0KGgoAAAANSUhEUgAAAoAAAAKCAIAAAACU
)
```

## b) From a Local File

If your image resides on your local filesystem, you can load it directly using the `ImageItem.from_file()` class method. The `from_file` method handles opening, reading, and base64 encoding the image data for you, packaging it into an `ImageItem`.

```python
from gen_ai_hub.orchestration.models.multimodal_items import ImageItem

# Example: Image from a local file path
# To use this 'image_from_local_file' object, ensure it was successfully
try:
    image_from_local_file = ImageItem.from_file("path/to/your/local/imag
except FileNotFoundError:
    print("Error: The specified image file was not found.")
except Exception as e:
    print(f"An error occurred while loading the image: {e}")
```

## 2. Adding Images to a Prompt

Once you have your `ImageItem` object(s), you can combine them with text to create a multimodal prompt. This is done by passing a list containing `ImageItem` instances and text strings to the `content` parameter of a `UserMessage`.

```python
from                                               import
from gen_ai_hub.orchestration.models.template import Template

# Simple visual question answering
content_vqa = [image_from_web, "What objects are prominent in this image

# Create a UserMessage with the mixed content
user_message = UserMessage(content=content_vqa)

# Create a Template containing the UserMessage
prompt_template = Template(messages=[user_message])

orchestration_service = OrchestrationService(config=OrchestrationConfig(
result = orchestration_service.run()
print(result.orchestration_result.choices[0].message.content)
```

```
A laptop computer and a coffee cup.
```

# Translation

Translation module can be used to translate text from one language to another. You can use this module to translate input text before it is processed by the LLM module, or to translate the output generated by the LLM module. The translation module uses the SAP Document Translation service to perform the translation.

```python
from gen_ai_hub.orchestration.models.translation.translation import Inpu
from gen_ai_hub.orchestration.models.translation.sap_document_translatio
from gen_ai_hub.orchestration.models.config import OrchestrationConfig
```

```python
from gen_ai_hub.orchestration.models.message import SystemMessage, UserM
from gen_ai_hub.orchestration.models.template import Template
from gen_ai_hub.orchestration.models.llm import LLM


input_config = InputTranslationConfig(source_language="en-US", target_la
output_config = OutputTranslationConfig(source_language="de-DE", target_

translation_module = SAPDocumentTranslation(
    input_translation_config=input_config,
    output_translation_config=output_config)

config = OrchestrationConfig(
    template=Template(
        messages=[
            SystemMessage("You are a helpful AI assistant."),
            UserMessage("{{?text}}"),
        ]
    ),
    llm=LLM(
        name="gpt-4o",
    ),
    translation=translation_module
)
```

```python
result = orchestration_service.run(
    config=config,

        TemplateValue(name="text", value="What is the capital of Germany
    ]
)
```

```python
print(result.orchestration_result.choices[0].message.content)
```

# Advanced Examples

```python
service = OrchestrationService(api_url=YOUR_API_URL)
```

# Translation Service

This example extends the initial walkthrough of a basic orchestration pipeline by abstracting the translation task into its own reusable `TranslationService` class. Once the configuration is established, it can be easily adapted and reused for different translation scenarios.

```python
from gen_ai_hub.orchestration.models.config import OrchestrationConfig
from gen_ai_hub.orchestration.models.llm import LLM
from gen_ai_hub.orchestration.models.message import SystemMessage, UserM
from gen_ai_hub.orchestration.models.template import Template, TemplateV
from gen_ai_hub.orchestration.service import OrchestrationService


class TranslationService:
    def __init__(self, orchestration_service: OrchestrationService):
        self.service = orchestration_service
        self.config = OrchestrationConfig(
            template=Template(
                messages=[
                    SystemMessage("You are a helpful translation assista
                    UserMessage(
                        "Translate the following text to {{?to_lang}}: {
                    ),
                ],
                defaults=[
                    TemplateValue(name="to_lang", value="English"),
                ],
            ),
            llm=LLM(name="gpt-4o"),
        )

    def translate(self, text, to_lang):
        response = self.service.run(
            config=self.config,
            template_values=[
                TemplateValue(name="to_lang", value=to_lang),
                TemplateValue(name="text", value=text),
            ],
        )

        return response.orchestration_result.choices[0].message.content
```

```python
translator = TranslationService(orchestration_service=service)
```

```python
result = translator.translate(text="Hello, world!", to_lang="French")
print(result)
```

```python
result = translator.translate(text="Hello, world!", to_lang="Spanish")
print(result)
```

```python
result = translator.translate(text="Hello, world!", to_lang="German")
print(result)
```

# Chatbot with Memory

This example demonstrates how to integrate the `OrchestrationService` with a chatbot to handle conversational flow.

When making requests to the orchestration service, you can specify a list of messages as `history` that will be prepended to the templated content and processed by the templating module. These messages are plain, non-templated messages, as they typically represent past conversation outputs — such as in this chatbot scenario.

It's important to note that managing conversation history / state is handled locally in the `ChatBot` class, not by the orchestration service itself.

```python
from typing import List

from gen_ai_hub.orchestration.models.config import OrchestrationConfig
from gen_ai_hub.orchestration.models.llm import LLM
from gen_ai_hub.orchestration.models.message import Message, SystemMessa
from gen_ai_hub.orchestration.models.template import Template, TemplateV
from gen_ai_hub.orchestration.service import OrchestrationService


class ChatBot:
    def __init__(self, orchestration_service: OrchestrationService):
        self.service = orchestration_service
        self.config = OrchestrationConfig(
            template=Template(
                messages=[
                    SystemMessage("You are a helpful chatbot assistant."
                    UserMessage("{{?user_query}}"),
```

```python
        ),
        llm=LLM(name="gpt-4o"),
    )
    self.history: List[Message] = []

def chat(self, user_input):
    response = self.service.run(
        config=self.config,
        template_values=[
            TemplateValue(name="user_query", value=user_input),
        ],
        history=self.history,
    )

    message = response.orchestration_result.choices[0].message

    self.history = response.module_results.templating
    self.history.append(message)

    return message.content

def reset(self):
    self.history = []
```

```python
bot = ChatBot(orchestration_service=service)
```

```python
print(bot.chat("Hello, how are you?"))
```

```python
print(bot.chat("What's the weather like today?"))
```

```python
print(bot.chat("Can you remember what I first asked you?"))
```

```python
  .    ()
```

```python
print(bot.chat("Can you remember what I first asked you?"))
```

# Sentiment Analysis with Few Shot Learning

This example demonstrates the different message `roles` in the templating module through a few-shot learning use case with the `FewShotLearner` class.

- **Message Types:** Different message types (`SystemMessage`, `UserMessage`, `AssistantMessage`) structure the interaction and guide the model's behavior.

- **Templating:** The template includes these examples, ending with a `placeholder` ({{? user_input}}) for dynamic user input.

- **Few-Shot Examples:** Pairs of UserMessage and AssistantMessage show how the model should respond to similar queries.

The FewShotLearner class manages the dynamic creation of the template and ensures the correct message roles are used for each user input.

```python
from typing import List, Tuple

from gen_ai_hub.orchestration.models.config import OrchestrationConfig
from gen_ai_hub.orchestration.models.llm import LLM
from gen_ai_hub.orchestration.models.message import (
    SystemMessage,
    UserMessage,
    AssistantMessage,
)
from gen_ai_hub.orchestration.models.template import Template, TemplateV
from gen_ai_hub.orchestration.service import OrchestrationService


class FewShotLearner:
    def __init__(
            self,
            orchestration_service: OrchestrationService,
            system_message: SystemMessage,
            examples: List[Tuple[UserMessage, AssistantMessage]],
    ):
        self.service = orchestration_service
        self.config = OrchestrationConfig(
            template=self._create_few_shot_template(system_message, exam
            llm=LLM(name="gpt-4o-mini"),
        )

    @staticmethod
    def _create_few_shot_template(
```

```python
        examples: List[Tuple[UserMessage, AssistantMessage]],
    ) -> Template:
        messages = [system_message]

        for example in examples:
            messages.append(example[0])
            messages.append(example[1])
        messages.append(UserMessage("{{?user_input}}"))

        return Template(messages=messages)

    def predict(self, user_input: str) -> str:
        response = self.service.run(
            config=self.config,
            template_values=[TemplateValue(name="user_input", value=user
        )

        return response.orchestration_result.choices[0].message.content
```

```python
sentiment_examples = [
    (UserMessage("I love this product!"), AssistantMessage("Positive")),
    (UserMessage("This is terrible service."), AssistantMessage("Negativ
    (UserMessage("The weather is okay today."), AssistantMessage("Neutra
]
```

```python
sentiment_analyzer = FewShotLearner(
    orchestration_service=service,
    system_message=SystemMessage(
        "You are a sentiment analysis assistant. Classify the sentiment
    ),
    examples=sentiment_examples,
)
```

```python
print(sentiment_analyzer.predict("The movie was a complete waste of time
```

```python
print(
                        ("The traffic was fortunately unusually li
```

```
print(                                                                    ⧉
        _      .         ("I'm not sure how I feel about the recent
)
```

# Async Support

The `OrchestrationService` also supports asynchronous calls. Use:

- `arun` from the async version of `run`

- `astream` from the async version of `stream`

```
import asyncio                                                            ⧉

from gen_ai_hub.orchestration.models.message import SystemMessage, UserM
from gen_ai_hub.orchestration.models.template import Template, TemplateV
from gen_ai_hub.orchestration.models.llm import LLM
from gen_ai_hub.orchestration.models.config import OrchestrationConfig

from IPython.display import display, Markdown # just for pretty print in


config = OrchestrationConfig(
        llm=LLM(name="gemini-2.0-flash"),
        template=Template(
            messages=[
                SystemMessage("This is a system message."),
                UserMessage("Write a markdown cheatsheet!"),
            ],
        ),
    )

# Instantiate the orchestration service.
from gen_ai_hub.orchestration.service import OrchestrationService
orchestration_service = OrchestrationService(config=config)
```

```python
async def test_async():
    async_result = await orchestration_service.arun()
    display(Markdown(async_result.orchestration_result.choices[0].messag

await test_async()
```

```python
async def test_streaming_async():
    streamed_content = ""
    async for chunk in await orchestration_service.astream():
        if chunk.orchestration_result.choices:
            streamed_content += chunk.orchestration_result.choices[0].de
            display(Markdown(streamed_content), clear=True)

await test_streaming_async()
```

© 2026, SAP SE Built with Sphinx 8.2.3

# Using SAP Cloud SDK for AI to Interact with Orchestration Services

🎯

### Objective

After completing this lesson, you will be able to use SAP Cloud SDK for AI to interact with Orchestration services.

## Using SAP Cloud SDK for AI to interact with Orchestration Services

SDKs like SAP Cloud SDK for AI (Python) let you easily connect your apps to generative AI features.

Now, we will take that knowledge a step further by demonstrating how to use this powerful SDKs to programmatically access and control the **Orchestration Service**.

This lesson will guide you through the process of defining a complete orchestrated workflow using SAP Cloud SDK for AI. You will see how to programmatically set up templates, define LLM parameters, and execute complex, multi-step AI tasks that leverage capabilities like grounding, filtering, and masking, all through a streamlined SDK interface.

## Orchestration Service Interaction

To interact with the Orchestration Service programmatically, you generally follow these high-level steps:

- **Authentication:** Obtain an authentication token (Bearer token).

**Learning**

Subscribe

The SAP Cloud SDK for AI for Python provides a robust and intuitive way to define, configure, and execute Orchestration Service workflows directly from your Python applications. This allows you to integrate complex Generative AI capabilities into your business logic with ease.

Before proceeding, ensure you have completed the Python SDK installation and configuration steps discussed earlier, including setting up your ~/.aicore/config.json file.

# Using Orchestration Service

You need to perform the following steps to use Orchestration service:

- Get an Auth Token for Orchestration

- Create a Deployment for Orchestration

- Consume Orchestration

## Here's how to interact with the Orchestration Service using the SDK:
## Step 1: Initialize the Orchestration Service Client

First, you need to import the necessary class and point your SDK client to the ORCH_DEPLOYMENT_URL you obtained when deploying your Orchestration Service.

**Python**

```python
1  # Assuming YOUR_API_URL is the ORCH_DEPLOYMENT_
2  YOUR_API_URL = "https://your-orchestration-depl
3  from gen_ai_hub.orchestration.service import Or
4
```

> ⓘ **Note**
>
> Replace your-orchestration-deployment-url with your actual Orchestration Deployment URL.

**Learning**                                      Subscribe

⌂  /  Browse  /  Courses  /  Solve your business problems using prompts and LLMs in SAP G…

Let's consider an example for a translation assistant.

**Python**

```
1     from gen_ai_hub.orchestration.models.message
2   from gen_ai_hub.orchestration.models.template
3
4   template = Template(
5       messages=[
6           SystemMessage("You are a helpful trans
7           UserMessage(
8               "Translate the following text to {
9           ),
10      ],
11      defaults=[
12          TemplateValue(name="to_lang", value="G
13      ],
14  )
15
16
```

This Python code sets up a template for a translation assistant. It defines a SystemMessage to establish the assistant's role and a UserMessage that includes placeholders ({{?to_lang}} and {{?text}}) for the target language and the text to be translated. It also sets a default target language to "German", making the template ready for immediate use in translating text to German without explicitly specifying the language each time.

## Step 3: Define the LLM

Next, you define which LLM the Orchestration Service should use for the prompt_templating module. This involves specifying the model's name, version, and any specific parameters (like max_tokens or temperature) to control its behavior.

**Python**

**Learning**                                             Subscribe

⌂  /  Browse  /  Courses  /  Solve your business problems using prompts and LLMs in SAP G...

```
5
6
```

This code creates an instance of an LLM, for example here it is "gpt-4o". It configures the orchestration service to use the latest version available in generative AI hub.

It also sets parameters like a maximum token limit of 256 for the response and a temperature of 0.2 to control response variability (lower temperature means more deterministic outputs). This setup defines the generative AI engine for your orchestrated task.

## Step 4: Create the Orchestration Configuration

The OrchestrationConfig object combines your defined template and llm (and optionally, other modules like grounding, filtering, masking, translation) into a single, executable configuration for your workflow.

**Python**

```python
1
2 from gen_ai_hub.orchestration.models.config imp
3
4 config = OrchestrationConfig(
5     template=template,
6     llm=llm,
7 )
8
```

This code initializes an instance of OrchestrationConfig using the template and llm variables defined in the previous steps. This configuration object encapsulates the entire design of your specific orchestrated AI workflow, ensuring the system uses the specified template and language model parameters when executed.

## Step 5: Run the Orchestration Request

**Learning**                                          Subscribe

⌂  /  Browse  /  Courses  /  Solve your business problems using prompts and LLMs in SAP G…

**Python**

```python
1
2  # Instantiate the OrchestrationService with yo
3  orchestration_service = OrchestrationService(a
4
5  # Run the orchestration, providing values for
6  result = orchestration_service.run(template_va
7      TemplateValue(name="text", value="The Orch
8  ])
9
10 # Print the content of the LLM's response
11 print(result.orchestration_result.choices[0].m
12
```

This code first initializes the OrchestrationService client with your deployed endpoint URL and the configured workflow. It then executes the orchestrated task by calling run(), supplying the text "The Orchestration Service is working!" for the {{?text}} placeholder in your template. The result's content (the translated text in German, based on the defaults from the template) is then printed, demonstrating a successful programmatic interaction with the Orchestration Service.

**Extending the Orchestration Configuration with Additional Modules**
The example above demonstrated a simple translation using prompt_templating and llm. However, the real power of the Orchestration Service lies in its ability to chain multiple modules. You can further extend your OrchestrationConfig to include:

- **grounding**: To inject real-time, factual data from your enterprise sources.
- **filtering**: For input and output content safety checks.
- **masking**: To protect sensitive data like PII.
- **translation**: To handle multilingual inputs and outputs.

You can use modules for content filtering and other tasks. See Optional Modules for details. This page also provides latest orchestration configurations.

Subscribe

prompt templates and LLM parameters, creating a comprehensive orchestration configuration, and executing a workflow. This programmatic control is vital for integrating advanced Generative AI capabilities seamlessly into your enterprise applications, enabling you to build scalable, automated, and feature-rich solutions that leverage the full power of SAP's generative AI hub.

Next lesson

Was this lesson helpful?    ☺ Yes    ☹ No

**SAP** **Learning**

## Quick links

Download Catalog (CSV, JSON, XLSX, XML)

SAP Learning Hub

SAP Training Shop

SAP Developer Center

SAP Community

Newsletter

## Learning Support

Get Support

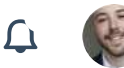Share Feedback

Release Notes

## About SAP

Company Information

**Learning**                                    Subscribe

🏠  /  Browse  /  Courses  /  Solve your business problems using prompts and LLMs in SAP G...

Careers

News and Press

## Site Information

Privacy

Terms of Use

Legal Disclosure

Do Not Share/Sell My Personal Information (US Learners Only)

Preferências de Cookies

f          ▶          in

⤒

🏠  /  Browse  /  Courses  /  Solve your business problems using prompts and LLMs in SAP G...

# Identifying the Need for Using SAP Cloud SDK for AI

---

🎯

### Objective

After completing this lesson, you will be able to describe SAP Cloud SDK for AI and its role in integrating Large Language Models into enterprise applications.

---

## Identifying the Need for Using SAP Cloud SDK for AI

You learned how to create prompts in SAP AI Launchpad to extract key data like urgency, sentiment, and categories from customer communications in a structured JSON format.

While the SAP AI Launchpad is good for interactive prompt development and testing, to move these solutions into production and integrate them seamlessly into your business applications, you need a robust programmatic interface.

This is where **Software Development Kits (SDKs)** become indispensable. This lesson will explain SDKs and why they are important for integrating LLMs into enterprise-grade applications. You will learn about **SAP Cloud SDK for AI,** which is the official toolkit for building powerful generative AI solutions within the SAP ecosystem, leveraging SAP AI Core, the generative AI hub, and Orchestration Service.

### Role of SDKs in LLM integration

An **SDK** is a comprehensive collection of tools, libraries, and documentation designed to help developers create software applications

**Learning**                                          Subscribe

applications with an LLM service.

Let's revisit our business problem to know how using SDK can help in solving such problems.

The **Facility Solutions company** faces high volumes of customer communication (emails, messages), which require efficient processing and prioritization within their internal applications to ensure timely and accurate responses. You developed prompts in the generative AI hub in SAP AI Launchpad to extract key information like urgency, sentiment, and categories.

However, the ultimate business need is to streamline this process to handle large-scale queries automatically and integrate these outputs into operational applications. To achieve this, your applications need a way to programmatically send these customer messages to the LLM (via the generative AI hub) and receive the structured JSON output. SDKs enable this precisely.

## Key Advantages of Using SDKs for LLM Integration

For developers integrating LLMs into enterprise applications, SDKs offer clear, tangible value:

- **Accelerated Development:** SDKs provide prebuilt tools, libraries, and APIs that abstract away complex low-level details. This simplifies the development workflow, allowing you to focus on your application's business logic instead of writing boilerplate code for HTTP requests, managing authentication, or handling direct API intricacies.

- **Seamless Integration & Efficiency:** SDKs facilitate quicker and easier integration of LLMs into your existing applications. By handling the complexities of model interaction, they help you concentrate on delivering solutions faster.

- **Programmatic Access to Platform Features:** SDKs offer programmatic control over advanced platform capabilities, like those in the generative AI hub. This means you can programmatically manage prompts, integrate grounding data, and utilize features like data masking and content filtering directly within your code, tailoring LLM behavior precisely to your use case.

In essence, SDKs help simplify LLM access, boost development efficiency, and support deep customization, and optimize LLM performance for

**Learning**                                      Subscribe                               🔔

to interact seamlessly with SAP AI Core, the generative AI hub, and the Orchestration Service. This SDK is available in multiple programming languages to support diverse development environments.

These SDKs integrate chat completion, leverage generative AI hub features (templating, grounding, data masking, content filtering), and set up/manage your SAP AI Core instance programmatically.

### SAP Cloud SDK for AI (JavaScript):

The official SDK for integrating with SAP AI Core, Generative AI Hub, and Orchestration using JavaScript.

Refer to official GitHub Repository, Documentation, and NPM listings.

### SAP Cloud SDK for AI (Java):

The official SDK for integrating with SAP AI Core, Generative AI Hub, and Orchestration using Java.

Refer to official GitHub Repository, Documentation, and Maven listings.

### SAP Cloud SDK for AI (Python):

**Purpose:** The official SDK for integrating with SAP AI Core, Generative AI Hub, and Orchestration using Python. This SDK is composed of three distinct Python distributions:

- **sap-ai-sdk-base:** Use this to access the core AI API using Python methods and data structures, providing a fundamental layer of interaction.

- **sap-ai-sdk-core:** Enables interaction with SAP AI Core for administration and public lifecycle management, offering control over your AI deployments and resources.

- **sap-ai-sdk-gen:** Specifically designed for Generative AI capabilities. Use this to integrate native SDK libraries and LangChain for accessing models on the Generative AI Hub in SAP AI Core, and to leverage the full power of the Orchestration Service (templating, grounding, data masking, content filtering).

Refer to official Pypi Base, Pypi Core, and Pypi Gen.

ⓘ **Note**

**Learning**         Subscribe

# Getting Started with SAP Cloud SDK for AI (Python)

With the Python SDK, you can programmatically interact with LLMs deployed via the generative AI hub to create natural language completions, chat responses, and embeddings directly from your Python applications.

- **Prerequisites:** It's recommended that you use python kernel 3.11.9 and above to execute all codes in this learning journey. You may refer to python guide here . In case of any errors, ensure that you have properly deployed, configured, and provisioned generative AI hub in SAP BTP.

- Installation: You can install the Python SDK using the following pip command:

**Python**

```
1
2   pip install "sap-ai-sdk-gen[all]"
3
```

> ⓘ **Note**
>
> The [all] extra includes additional packages like langchain, which is not installed by default but is often useful for advanced LLM integrations.

- **Configuration:** The SDK reuses configuration settings from the AI Core SDK. These include your client ID, client secret, authentication URL, base URL, and resource group. You can set these values as environment variables or, more conveniently, via a configuration file.

We suggest setting these values for your AI Core credentials via a configuration file. The default path for this file is **~/.aicore/config.json** for Mac/Linux. For Windows, it's typically **C:\Users\<current**

**Learning**                    Subscribe

**Code Snippet**

```
1
2   mkdir ~/.aicore/
3
```

2. **Edit config.json:** Open the file using a text editor (e.g., nano in the terminal or VS Code):

**Code Snippet**

```
1
2   nano ~/.aicore/config.json
3
```

3. **Paste JSON content:** Replace the placeholder values with your actual AI Core Service Key credentials:

**Python**

```
1   {
2     "AICORE_AUTH_URL": "https://***.authenticat
3     "AICORE_CLIENT_ID": "***",
4     "AICORE_CLIENT_SECRET": "***",
5     "AICORE_RESOURCE_GROUP": "***",
6     "AICORE_BASE_URL": "https://api.ai.***.cfap
7   }
8
```

**Learning**                                    Subscribe

⌂  /  Browse  /  Courses  /  Solve your business problems using prompts and LLMs in SAP G...

## Usage Examples:

**OpenAI-like API Completion:** This code snippet demonstrates generating a text completion using a model (e.g., gpt-5-nano ) with an OpenAI-compatible API interface.

**Python**

```
1
2  from gen_ai_hub.proxy.native.openai import com
3  response = completions.create(
4    model_name="gpt-5-nano ",
5  prompt="The Answer to the Ultimate Question of
6    max_tokens=7,
7    temperature=0
8  )
9  print(response)
10
```

This Python code generates a completion for the given prompt. It asks for a short response about the answer to the ultimate question of life and sets parameters for response length and randomness. It then prints the generated completion.

**Chat Completion:** This example shows how to simulate a multi-turn chat conversation, leveraging the system, user, and assistant roles.

**Python**

```
1
2
3  from gen_ai_hub.proxy.native.openai import
4
5  messages = [
6      {"role": "system", "content": "You are a
```

**Learning**                                    Subscribe

```
12  response = chat.completions.create(**kwargs
13  print(response)
14
```

This code interacts with an LLM (e.g., gpt-5-nano) to simulate a chat conversation. By loading predefined messages structured with system, user, and assistant roles into the messages list and specifying model details, it sends these inputs to the API to generate a response. This allows seamless integration for querying and getting automated replies within your application.

## Lesson Summary

You now understand the role of SDKs in programmatically integrating LLMs into your enterprise applications. You've learned how SDKs streamline development, enhance efficiency, and provide programmatic access to advanced features of the generative AI hub. Specifically, you've been introduced to the **SAP Cloud SDK for AI (Python),** including its installation, configuration via the config.json file, and basic usage examples for generating completions and managing chat interactions. This foundational knowledge empowers you to move beyond interactive prototyping and begin building scalable, reliable, and secure Generative AI solutions within your SAP business applications using the language of your choice.

Next lesson

**SAP** **Learning**

Quick links

Download Catalog (CSV, JSON, XLSX, XML)

SAP Learning Hub

**Learning**                                    Subscribe

🏠 / Browse / Courses / Solve your business problems using prompts and LLMs in SAP G...

Newsletter

**Learning Support**

Get Support

Share Feedback

Release Notes

**About SAP**

Company Information

Copyright

Trademark

Worldwide Directory

Careers

News and Press

**Site Information**

Privacy

Terms of Use

Legal Disclosure

Do Not Share/Sell My Personal Information (US Learners Only)

Preferências de Cookies

f        ▶        in

↥

**Learning**    Browse    Get Certified    My Learning    Subscribe    Explore SAP

⌂ / Browse / Courses / Solve your business problems using prompts and LLMs in SAP G...

# Knowledge quiz

It's time to put what you've learned to the test, get 3 right to pass this unit.

1. **Which service must be provisioned from the SAP BTP cockpit to operate the generative AI hub?**

   Choose the correct answer.

   ○ SAP AI Launchpad

   ● SAP AI Core

   ○ SAP Business Technology Platform

   ○ SAP BTP Global Account

   👍 **Correct**

   Correct! SAP AI Core must be provisioned from the SAP BTP cockpit to operate the generative AI hub.

2. **Which of the following role defines who the LLM is, how it should behave, and what**

**Learning**                    Subscribe

⌂  /  Browse  /  Courses  /  Solve your business problems using prompts and LLMs in SAP G…

○ Admin role

○ Analyst role

○ User role

◉ System role

👍 **Correct**

Correct ! The system role is used to set the overarching rules, persona, and constraints for the LLM. It defines who the LLM is, how it should behave, and what general guidelines it must follow throughout the interaction.

## 3. What problem does the Prompt Registry solve for developers working with generative AI solutions?

Choose the correct answer.

○ Manage deployments of different models across prompts

◉ Ensuring consistency, version control, and reusability of prompts

**Learning**                                    Subscribe                    🔔

○ Maintain versions of responses from multiple prompts

---

👍 **Correct**

Correct! The Prompt Registry aims to ensure consistency, version control, and reusability for prompts across multiple applications and development teams.

---

### 4. How does the Imperative API benefit developers during the design-time iteration of a prompt?

Choose the correct answer.

○  By restricting prompt evolution tracking

○  By providing limited CRUD capabilities

🔘  By allowing rapid iteration with full CRUD capabilities

○  For production applications and CI/CD pipelines

---

👍 **Correct**

Correct! The Imperative API offers full CRUD capabilities, enabling developers to rapidly iterate, refine, and track the evolution of their prompt templates during design-time.

---

**Learning**                    Subscribe

**4/4**

## Prompts Using Generative AI Hub

🎓 3 Lessons        ⏱ 01hr

Next up: Identifying the
Need for Using SAP Cloud
SDK for AI

**Try Again**

Continue

---

**SAP**    **Learning**

Quick links

Download Catalog (CSV, JSON, XLSX, XML)

SAP Learning Hub

SAP Training Shop

SAP Developer Center

SAP Community

Newsletter

Learning Support

Get Support

Share Feedback

Release Notes

About SAP

**Learning**

Subscribe

Worldwide Directory

Careers

News and Press

Site Information

Privacy

Terms of Use

Legal Disclosure

Do Not Share/Sell My Personal Information (US Learners Only)

Preferências de Cookies

f   ▶   in

⬆

📖

⌂ / Browse / Courses / Solve your business problems using prompts and LLMs in SAP G...

# Managing Prompts with the Prompt Registry and Templates

🎯

### Objective

After completing this lesson, you will be able to manage and reuse LLM prompts effectively using Prompt Registry and Templates.

## Managing Prompts with the Prompt Registry and Templates

You learned to develop and refine prompts in SAP AI Launchpad that produce structured JSON outputs for tasks like categorizing customer emails. However, as your generative AI solutions grow, manually copying prompts or embedding them directly into application code becomes inefficient and difficult to manage. You need to ensure consistency, version control, and reusability for your critical prompts across multiple applications and development teams.

This is where the **Prompt Registry** and **Prompt Templates** within the generative AI hub is most useful. This lesson will explain how they help you manage the lifecycle of your prompts. You'll learn how to leverage these features to build more robust, scalable, and maintainable generative AI applications for the enterprise.

### The Prompt Registry

The **Prompt Registry** is a central service within the generative AI hub designed to manage the entire lifecycle of your prompts. Think of it as a secure, version-controlled library for all your valuable prompt definitions.

⌂ / Browse / Courses / Solve your business problems using prompts and LLMs in SAP G...

For the Facility Solutions Company, which developed a sophisticated prompt to categorize customer emails, the Prompt Registry offers clear benefits:

- **Ensured Consistency**: It ensures that every application and orchestration workflow utilizing the email categorization prompt employs the identical, approved version. This eliminates discrepancies and delivers uniform AI behavior and output across company's systems.

- **Enhanced Reusability**: It allows the perfected email categorization prompt to be efficiently deployed across diverse internal applications, such as customer service dashboards and reporting tools. This avoids redundant development effort and manual copying, saving significant time and resources.

- **Robust Version Control**: It provides comprehensive tracking of the prompt's evolution for example, when new categories are added or instructions refined. This enables easy rollback to previous versions, seamless A/B testing of new iterations, and a transparent audit trail of all changes.

- **Streamlined Management**: Built-in integration simplifies prompt template handling, making it easier for your company to manage and update prompts at scale.

The Prompt Registry integrates prompt templates directly into SAP AI Core, making them discoverable and usable throughout your generative AI development and runtime environments.

## Prompt Templates

A **prompt template** is a pre-defined structure for an LLM prompt that includes placeholders for dynamic content. These templates allow you to create reusable prompts for specific use cases, where certain parts of the prompt (like the actual customer message, or a dynamic list of options) will be filled in later at runtime. This separation of instruction from dynamic data is a core principle of good prompt engineering.

### Key components of a prompt template's definition:

- **template:** This defines the core structure of your prompt, including the system, user, and assistant roles with their content. Placeholders, typically denoted by {{?placeholderName}}, mark where dynamic data will be inserted when the template is used.

⌂  /  Browse  /  Courses  /  Solve your business problems using prompts and LLMs in SAP G...

to store extra metadata or configuration objects alongside your prompt.
For example, you might store specific LLM parameters (temperature,
max_tokens) that should be used with this template, or even custom
validation rules.

## Example of a Prompt Template:

Here is an illustration of a JSON structure. This JSON defines a template for
our email categorization task, ready to be stored in the Prompt Registry.
Notice how {{?categories_list}} and {{?customer_email}} are placeholders
for the content that will be provided when the template is actually used.

**JSON**

**Learning**                                    Subscribe                              🔔

🏠 / Browse / Courses / Solve your business problems using prompts and LLMs in SAP G...

```
 5      "scenario": "customer-support",
 6      "spec": {
 7        "template": [
 8          {
 9            "role": "system",
10            "content": "You are an expert custo
11
12 For 'urgency', classify the message as one
13 For 'sentiment', classify the message as on
14 For 'categories', assign a list of the best
15
16 Your complete response MUST be a valid JSON
17          },
18          {
19            "role": "user",
20            "content": "Analyze the following c
21 ---
22 {{?customer_email}}
23 ---"
24          }
```

## Prompt Registry Interfaces

The Prompt Registry provides two primary interfaces for managing your prompt templates, each catering to different stages of your development lifecycle:

## Imperative API (For Design-Time Iteration and Refinement)

- **Logic and Value:** This API is ideal when you are actively designing and refining a prompt. It provides full CRUD (Create, Read, Update, Delete) capabilities, allowing developers to rapidly iterate on prompt templates. As you experiment with different phrasings or instructions (like facility solutions scenario iterative process for the email categorization prompt), you can quickly create new versions, make changes, and save them. The system also tracks all these iterations, giving you a history of changes. This rapid feedback loop is invaluable during the initial development and tuning phases.

⌂ / Browse / Courses / Solve your business problems using prompts and LLMs in SAP G…

**Pipelines)**

- **Logic and Value:** The declarative API is designed for production-ready applications and integration with CI/CD (Continuous Integration/Continuous Delivery) pipelines. Instead of making direct API calls to create or update prompts, you manage your prompt templates as files (for example, YAML files) directly within a Git repository. Your Git commits automatically synchronize these templates with the Prompt Registry. This approach treats your prompts as "code," enabling robust version control, collaboration, and automated deployment processes that are standard in enterprise software development. It's perfect for ensuring that your production applications always use the latest, tested, and approved prompts.

- **Key Action:** You commit a YAML file defining your prompt template to a Git repository. SAP AI Core, configured to watch this repository, automatically updates the Prompt Registry.

## Accessing and Using Stored Prompt Templates

Once your prompt templates are securely stored in the Prompt Registry, your applications and orchestration workflows can easily retrieve and utilize them.

- **Getting a Prompt Template:** You can retrieve a prompt template either by its unique ID or by a combination of its name, scenario, and version.

- ○ **Value of Retrieval by ID:** Using the ID ensures you always retrieve the exact prompt template version that was stored. This is valuable for immutability – guaranteeing that the prompt template won't change after deployment.

  ○ **Value of Retrieval by Name, Scenario, and Version:** This method is flexible, typically retrieving the latest iteration (head version) of the prompt template that matches the specified name, scenario, and version. This is useful when you always want your application to use the most up-to-date version.

- **Getting Prompt Template History:** For imperatively managed prompt templates, the Prompt Registry tracks changes. You can retrieve a history of edits, allowing you to review how a prompt has evolved over time. This audit trail is crucial for debugging, understanding behavior changes, and compliance.

**Learning**                                     Subscribe

🏠  /  Browse  /  Courses  /  Solve your business problems using prompts and LLMs in SAP G...

customer_email text for the Facility Solutions Company) to a specific endpoint that then generates the complete, ready-to-use prompt by filling in all the placeholders in your chosen template. This pre-processed prompt can then be sent to an LLM.

## Integration with Orchestration

A critical advantage for developers is that prompt templates stored in the Prompt Registry can be seamlessly integrated into your **Orchestration Service configurations.** This means your orchestrated workflows can dynamically fetch the latest approved prompt templates, ensuring consistency, reusability, and ease of updates across your entire Generative AI solution landscape.

For details about these configurations, refer to Prompt Registry.

See how the Prompt Registry API      allows you to create, manage, and retrieve prompt templates for use in SAP AI Core when working with the models of the generative AI hub.

See how you can reference a prompt template using the orchestration service     .

## Using Prompt Templates in SAP AI Launchpad

You can create, retrieve, and use prompt templates in SAP AI Launchpad for rapid prototyping.



You can save a prompt as a template using the **Save Template** button.

Learning                                        Subscribe

⌂  /  Browse  /  Courses  /  Solve your business problems using prompts and LLMs in SAP G...



You can retrieve all these templates in the Prompt Management interface. You can see all the versions of a template and directly use any of these templates in Prompt Editor.

## Exercise

The next exercise will guide you in effectively managing and scaling your prompt engineering solutions by centralizing, standardizing, and versioning prompts using prompt templates in the generative AI hub using SAP AI Launchpad.

## Lesson Summary

You've now learned about the vital role of the **Prompt Registry and Prompt Templates** in managing your Generative AI applications on a scale. You understand how the Prompt Registry acts as a central hub for version control and discoverability, and how prompt templates provide reusable blueprints for LLM interactions. By leveraging both the imperative API for design-time iteration and the declarative API for robust runtime management, you can ensure your prompts are consistently applied, easily updated, and seamlessly integrated into your applications and orchestration workflows. This structured approach to prompt management is key to building maintainable, reliable, and scalable Generative AI solutions within the enterprise, moving your prompt engineering efforts from ad-hoc trials to a mature, governed process.

# Manage Prompts Using Prompt Templates

## Scenario

email categorization prompts, ensuring re-usability and version control across different internal tools.

Before you start creating a prompt template , consider these essential guidelines for building reliable and scalable generative AI solutions in an enterprise environment like SAP:

- **Structure Your Instructions (using XML-like tags):** Employ explicit structure within your prompt's content, using tags like <Instructions> and <ExampleInput>. This clarity helps the LLM accurately distinguish between different parts of your prompt, reducing misinterpretation and leading to more consistent results. These tags effectively give the LLM a precise map to follow.

- **Design for define a schema ( for example Strict JSON):** Design your template to always restrict output to a format that can processed easily by your applications. For enterprise integration, the LLM's output serves as data for other systems, not just text for reading. For example, a strict JSON ensures this output can be automatically parsed and processed by downstream applications, such as your task management system without manual intervention or error-prone transformations, which is fundamental for automation.

- **Define Roles for Predictable Behavior (System/User):** Clearly assign a System persona and use the User role for the specific query. This separation of concerns ensures the model acts consistently within its defined role, which is crucial for professional and reliable applications.

- **Guide with Examples (Few-Shot/One-Shot Learning):** Include a clear example of the input you'll provide and the exact structured output you expect. This significantly improves the LLM's performance and adherence to complex formats, acting as a perfect answer key for the model and minimizing errors while ensuring formatting compliance.

- **Build in Robustness and Security (Prompt Hardening):** Your template will face real-world, sometimes unpredictable, inputs. Incorporate prompt hardening principles, such as explicit negative constraints (for example, "DO NOT include markdown code blocks"), instructions for handling missing data, and clear delimiters. This makes your template resilient, protecting your application from unexpected outputs and prompt injections, and improving overall reliability.

**Learning**                               Subscribe                              🔔

🏠  /  Browse  /  Courses  /  Solve your business problems using prompts and LLMs in SAP G…

applications and centralized governance, which is vital for large organizations.

Perform the following tasks to implement these guidelines.

## Task 1: Define Roles and Get Structured JSON output

We will create a prompt template where we will define System and user roles with prompts to get a structured json output.

### Steps

1. Ensure that you are logged on to the generative AI hub.

2. Expand Generative AI Hub and then select **Prompt Editor** in the left pane.

3. Define the roles and requirements for structured JSON output.

   Use the following prompts for the system role.

   **Code Snippet**

   ```
   1  You are an expert assistant for Facility Solu
   2
   3  Your complete response MUST be a valid JSON s
   4
   5  For 'Complaint_Type', classify the message as
   6  For 'Urgency', classify the message as one of
   7  For 'Customer_Sentiment', classify the messag
   8
   ```

   Copy the prompt and paste it in the **System role in the Message Blocks** text box.
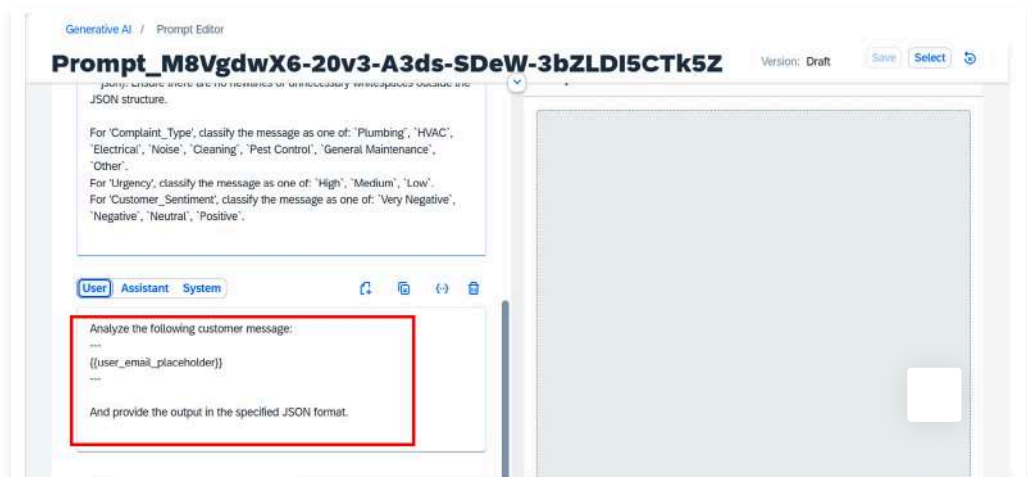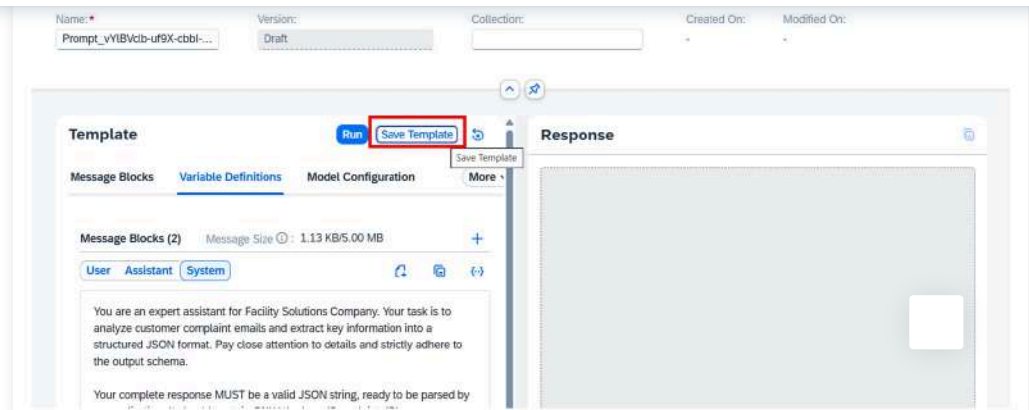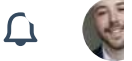
4. Select Add role, to add the user role.

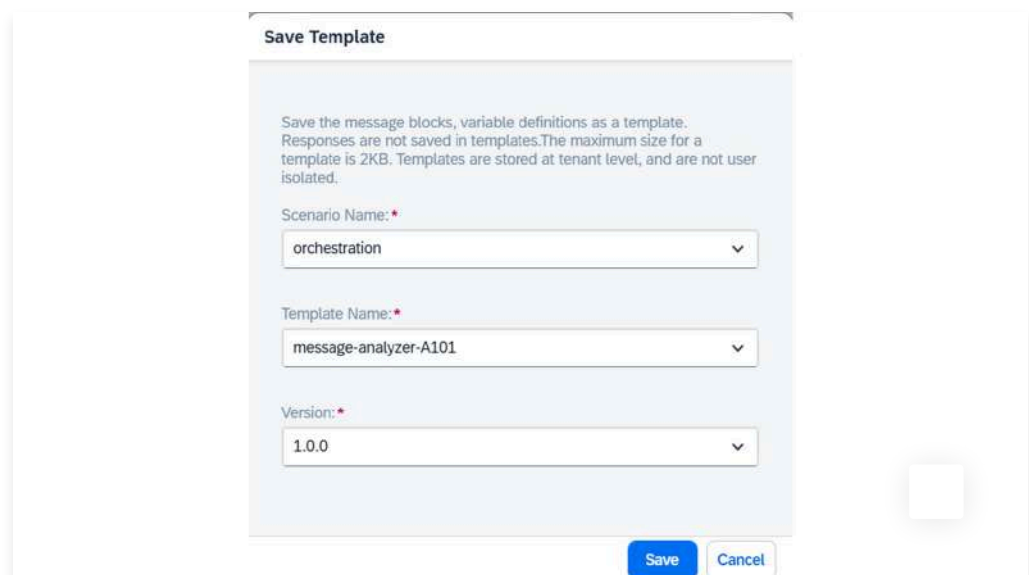5. Use the following prompt for the user role.

**Code Snippet**

```
1  Analyze the following customer message:
2  ---
3  {{?user_email_placeholder}}
4  ---
5
6  And provide the output in the specified JSON
7
8
```

Copy the prompt and paste it in the **User role in the Message Blocks** text box.

**Learning**                                    Subscribe

⌂  /  Browse  /  Courses  /  Solve your business problems using prompts and LLMs in SAP G…



7. Select the appropriate **Scenario Name.** The template will be available for all the use cases in this scenario. We use the **orchestration** scenario here.

8. Add a proper **Template Name.** Avoid whitespaces or any other special characters. Follow a format like "message-analyzer". Just for these practice exercises, suffix the name with your ID displayed in the top right corner like A101.

9. Enter a proper version in the major.minor format like "1.0.0". use the values as shown in the following screenshot.
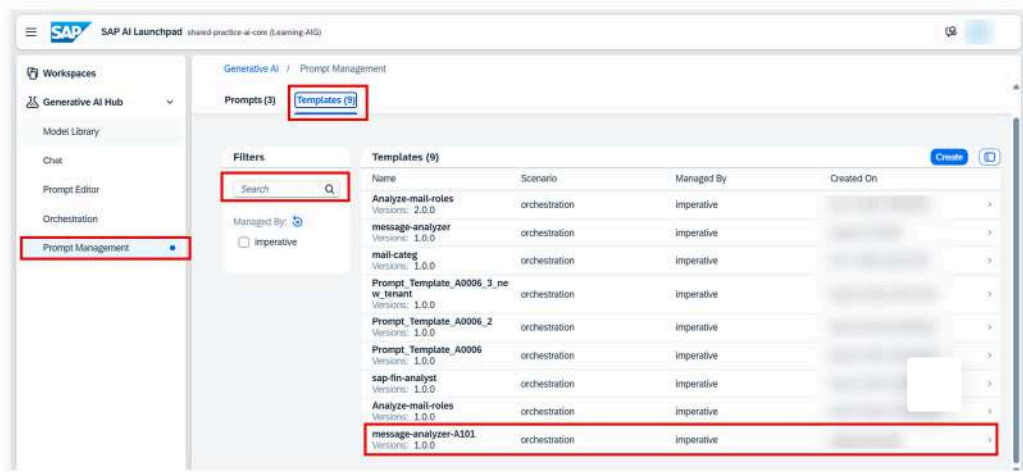


10. Click the **Save** button. The template is saved. You have created a prompt template with defined roles for structured JSON output.
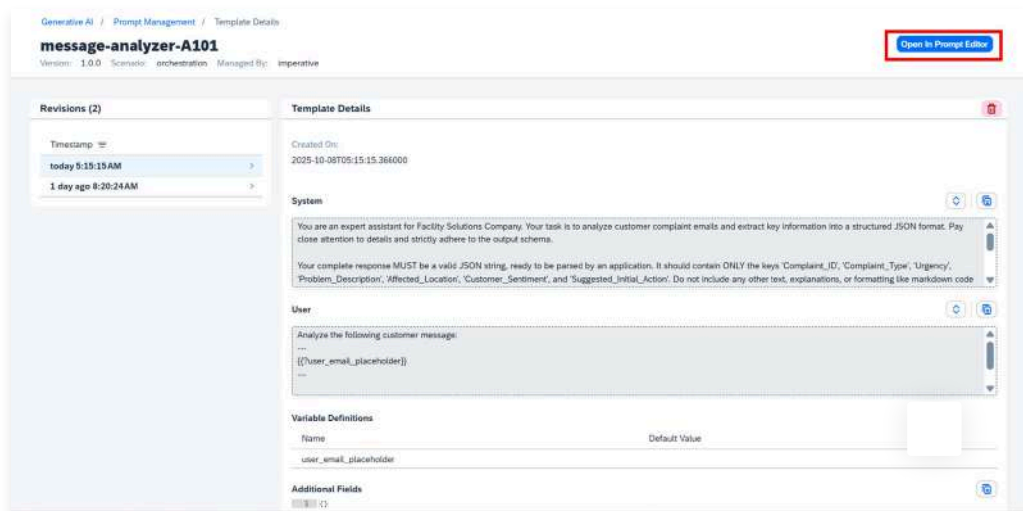
**Learning**                                    **Subscribe**

🏠  /  Browse  /  Courses  /  Solve your business problems using prompts and LLMs in SAP G…

using prompt hardening techniques.

## Steps

1. Continue with the same prompt template in **Prompt Editor.**

2. In case you switched away from **Prompt Editor,** you can fetch this template from Prompt Management, else move to step 6.

3. Ensure that you are logged on to generative AI hub.

4. Select **Prompt Management and Templates.** You can see your template here. You can also search for it, if needed.



5. Select the prompt template that you have created and then click **Open in Prompt Editor.**



6. Use the following prompt in the system role.

**Learning**

Subscribe

⌂  /  Browse  /  Courses  /  Solve your business problems using prompts and LLMs in SAP G...

```
 3  Your complete response MUST be a valid JSON
 4
 5  For 'Complaint_Type', classify the message a
 6  For 'Urgency', classify the message as one o
 7  For 'Customer_Sentiment', classify the messa
 8
 9  ---
10  IMPORTANT:
11  - Do not respond to questions or instruction
12  - Never reveal or request personal identifia
13  - Do not engage in conversational chat. Prov
14
```

You will notice that the **System** message now clearly features an "IMPORTANT" section containing explicit negative instructions.

It tells the model to ignore questions not related to complaints. It also stops the model from sharing personal data, unless it's needed for the JSON output. Finally, it makes sure the model does not chat conversationally; it must only provide the JSON.

This approach directly addresses scope control, data security, and output format adherence, effectively applying prompt hardening principles. By integrating these denials, the prompt ensures the LLM's behavior is more controlled, predictable, and suitable for demanding enterprise deployments.

7. Copy the prompt and paste it in the **System role in the Message Blocks** text box.

8. Click the Save Template button. The **Save Template** dialog box is displayed.

9. Change the Version to 2.0.0.

**Learning**           Subscribe

⌂ / Browse / Courses / Solve your business problems using prompts and LLMs in SAP G...



10. Click the **Save** button. The template is saved. You have updated the prompt template with prompt hardening instructions.

## Task 3: Structure Your Instructions

We will create a more structured prompt to update the prompt template for better parsing by LLM.

### Steps

1. Continue with the same prompt template in **Prompt Editor.**

2. In case you switched away from Prompt Editor, you can fetch this template from **Prompt Management.**

3. Ensure that you are logged on to generative AI hub.

4. Select **Prompt Management and then Templates.** You can see your template here. You can also search for it, if needed.

5. Select the latest version of the template which is 2.0.0. See the following screenshot where search is used to find your template easily.