# STEPS TO SETUP APIM

Home > Resource groups >

# Create a resource group ...

**Basics**   Tags   Review + create

**Resource group** - A container that holds related resources for an Azure solution. The resource group can include all the resources for the solution, or only those resources that you want to manage as a group. You decide how you want to allocate resources to resource groups based on what makes the most sense for your organization. Learn more

Subscription * ⓘ                         | Free Trial                    ⌄ |

Resource group name * ⓘ                 |                                |

Region * ⓘ                               | (US) East US                  ⌄ |

---

☰  **Microsoft Azure**   ⊕ Upgrade  🔍 Search resources, services, and docs (G+/)     🟦 Copilot   ⊡ 🔔 ⚙ ❓ 🗗

Home >

# Resource groups 📌 ...
Relanto

+ Create   ⚙ Manage view ⌄   ↻ Refresh   ↓ Export to CSV   ⌁ Open query   | ⬚ Assign tags

⌞ Create a new resource group.

Filter for any field...   Subscription equals **all**   Location equals **all** ✕   ⛢ Add filter

Showing 1 to 1 of 1 records.                                              No grouping          ⌄

☐ Name ↑↓                                Subscription ↑↓                  Location ↑↓

☐ ⬡ apim_poc                             Free Trial                      East US

Home > Azure AI services | Azure OpenAI >

# Create Azure OpenAI ⋯

① **Basics**    ② Network    ③ Tags    ④ Review + submit

Azure OpenAI Service provides access to OpenAI's powerful language models, including all the latest OpenAI models. These models can be easily adapted to your specific tasks, including but not limited to content generation, summarization, image understanding, semantic search, and natural language to code translation. Top use cases include Call Centers, Virtual Assistants, Accessibility, Content Generation, and Code Development. The service also features the Assistants API, Fine Tuning capabilities and many ways to connect your data to the service for conversational experiences. The service can be scaled through Standard (tokens) and Provisioned (PTUs) deployment types.

Learn more

## Project Details

| | |
|---|---|
| Subscription * ⓘ | Free Trial ▾ |
| Resource group * ⓘ | ▾ |
| | Create new |

Previous    **Next**

---

Home > Azure AI services

⬧ **Azure AI services** | Azure OpenAI 📌 ⋯
Azure AI services

🔍 Search   ✕   «

+ Create   ✎ Manage deleted resources   ⚙ Manage view ∨   ↻ Refresh   ↓ Export to CSV   ⅋ Open query   ⋯

- ⬡ Overview
- ⬡ All Azure AI services
- ∨ Azure AI services
  - ⬧ Azure AI services
  - ⬧ **Azure OpenAI**
  - ⬧ AI Search
  - ⊙ Computer vision
  - ⬧ Face API
  - ⬧ Custom vision

Filter for any field...    Subscription equals **all**    Type equals **all**    ⁺⊽ Add filter

Showing 1 to 1 of 1 records.    No grouping ▾   ⊟ List

| ☐ Name ↑↓ | Kind ↑↓ | Location ↑↓ | Custom Domain ... ↑↓ | Pricing tier |
|---|---|---|---|---|
| ☐ ⬧ apim-poc-azure-ai-project | OpenAI | East US | ⬡ apim-poc-azure-a... | S0 |

1. Open the Azure open AI service that you created.
2. Go to the endpoint section and Copy the Service Endpoint.
3. Click on "Go to AI Foundry portal"



4. Go to the Deployments Section

5. Create a new Deployment Model and Deploy the model you need in that section by selecting the model and version.

1. Go to **API Management services** and create APIM**:**

# Create API Management service ...
API Management service

Basics | Monitor + secure | Virtual network | Managed identity | Tags | Review + install

## Project details

Select the subscription to manage deployed resources and costs. Use resource groups like folders to organize and manage all your resources.

| Subscription * ⓘ | Free Trial ▾ |
| Resource group * ⓘ | apim_poc ▾ |

Create new

## Instance details

| Region * ⓘ | (US) East US ▾ |
| Resource name * | |
| Organization name * ⓘ | Enter organization name |

[Review + create] [< Previous] [Next: Monitor + secure >]

Note: 1. Use standard tier to meet the requirements

2. Enable status of Managed identity section

# Create API Management service ...
API Management service

❌ Basics | Monitor + secure | Virtual network | **Managed identity** | Tags | Review + install

A system assigned managed identity enables Azure resources to authenticate to cloud services (e.g., Azure Key Va without storing credentials in code. Once enabled, all necessary permissions can be granted via Azure role-based-access control. The lifecycle of this type of managed identity is tied to the lifecycle of this resource. Additionally, ea resource (e.g., Virtual Machine) can only have one system assigned managed identity. Learn more

### System assigned managed identity
Enable system assigned identity to grant the resource access to other existing resources.

| Status | ☑ |

[Review + create] [< Previous] [Next: Tags >]

1. Now after creating APIM go to the Azure openAI service's Access control and click on Add role assignment



2. Add Cognitive Services OpenAI User(Built in) role:

**Note:**

(1) Assign access to Managed identity

(2) Select the APIM according to the image below

Download the json file from below link and upload it in the next step(This is for openai gpt-4o model)

Ref: https://github.com/HoussemDellai/ai-course/blob/main/300_apim_genai_openai/openapi-AzureOpenAI-2024-10-21-

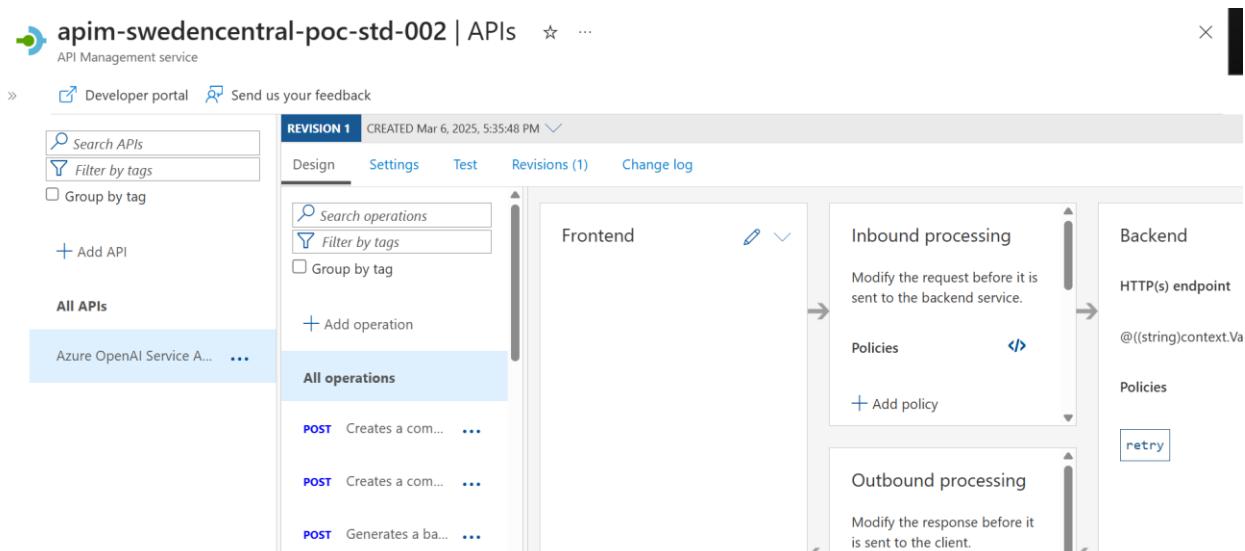Note: Add a suffix for your APIM URL endpoint.



## Create from OpenAPI specification

**Basic** | Full

* **OpenAPI specification**    `https://`    or    **Select a file**
                                                (maximum size 4 MiB)

**Include required query parameters in operation templates**    ☑

* **Display name**    GPT endpoints

* **Name**    gpt-endpoints

**API URL suffix**    openai

**Base URL**
`http(s)://apim-swedencentral-poc-std-002.azure-api.net/openai`

**Create**    **Cancel**

1. Go to the API which was just imported and go to the design section

## 2. Go to the policy in the inbound processing and add the below policy



Policy used (Replace your backend URL's in the JArray section):

```
<policies>

  <inbound>

    <base />

    <rate-limit-by-key calls="100" renewal-period="10" counter-
key="@(context.Request.IpAddress)" />

    <set-variable name="backendArray" value="@(

      new JArray(

        "https://ai-services-eastus-demo-prod.cognitiveservices.azure.com/openai",

        "https://ai-services-eastus-demo-prod-002.cognitiveservices.azure.com/openai",

        "https://apim-poc-azure-ai-project.openai.azure.com/openai"

      )

    )" />

    <set-variable name="backendIndex" value="@{

      var backends = (JArray)context.Variables["backendArray"];

      return (int)(DateTime.UtcNow.Ticks % backends.Count);
```

```xml
        }" />
    <set-variable name="backendSelection" value="@{

        var backends = (JArray)context.Variables["backendArray"];

        return (string)backends[(int)context.Variables["backendIndex"]];

        }" />

    <set-backend-service base-url="@((string)context.Variables["backendSelection"])" />

    <authentication-managed-identity resource="https://cognitiveservices.azure.com"
output-token-variable-name="msi-access-token" ignore-error="false" />

    <set-header name="Authorization" exists-action="override">

        <value>@("Bearer " + (string)context.Variables["msi-access-token"])</value>

    </set-header>

  </inbound>

  <backend>

    <retry condition="@(context.Response == null || context.Response.StatusCode == 500
|| context.Response.StatusCode == 429 || context.Response.StatusCode != 200)"
count="20" interval="1">

        <forward-request timeout="@((100 +
(context.Variables.ContainsKey(&quot;retryCount&quot;) ?
(int)context.Variables[&quot;retryCount&quot;] * 100 : 0)))" buffer-request-body="true" />

        <set-variable name="backendIndex" value="@{

            var backends = (JArray)context.Variables["backendArray"];

            return ((int)context.Variables["backendIndex"] + 1) % backends.Count;

        }" />

        <set-variable name="backendSelection" value="@{

            var backends = (JArray)context.Variables["backendArray"];

            return (string)backends[(int)context.Variables["backendIndex"]];

        }" />
```

```
      <set-backend-service base-url="@((string)context.Variables["backendSelection"])" />

    </retry>

  </backend>

  <outbound>

    <base />

    <set-header name="X-Selected-Backend" exists-action="override">

      <value>@((string)context.Variables["backendSelection"])</value>

    </set-header>

  </outbound>

  <on-error>

    <base />

  </on-error>

</policies>
```

3.  Specify the header and query parameter name and enable the checkbox(This will be used while making an API call)

1. Create a subscription key in the Subscription section of APIs under APIM



2. Retrieve the key and copy

This Apim's performance can be tested using Apache benchmark docker image

Run the below Docker command:
**docker run --rm -v <Path to your payload file>:/data jordi/ab -n 10 -c 10 -l -T "application/json" -H "api-key: <your-subscription-key>" -p /data/apim-body.json -v 2 -e /data/apim-errors.csv -g /data/apim-results.tsv [https://apim-swedencentral-poc-std-002.azure-api.net/openai/deployments/gpt-4o/chat/completions?api-version=2024-08-01-preview](https://apim-swedencentral-poc-std-002.azure-api.net/openai/deployments/gpt-4o/chat/completions?api-version=2024-08-01-preview)**

**Note:**
-n represents the number of requests
-c represents the concurrency
That means in every second c number of requests are sent parallelly.

Go to the Monitoring section of the APIM and query to get the logs.

**Query:**

1. **To get the Details:**

ApiManagementGatewayLogs

| where TimeGenerated > ago(1d)

| where IsRequestSuccess == true

| where ResponseCode == 200

| project TimeGenerated, BackendUrl| top 10 by TimeGenerated desc

**2. To get the Count of the Requests hit to particular Backend:**

```
ApiManagementGatewayLogs
| where TimeGenerated > ago(1d)
| where IsRequestSuccess == true
| where ResponseCode == 200
| project TimeGenerated, BackendUrl
| top 50 by TimeGenerated desc  // Fetch last 30 requests
| summarize HitCount = count() by BackendUrl
| order by HitCount desc
```