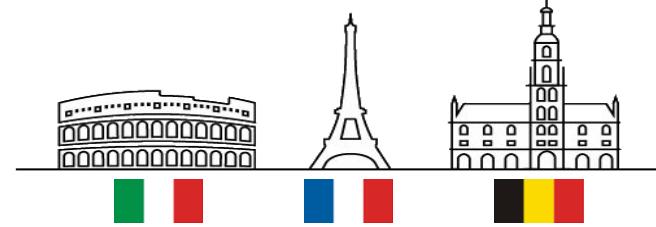


APIM ❤️ OpenAI

exploring the generative AI Gateway pattern

Massimo Crippa

Lead Architect at **codit**
proximus NXT



Generative AI

has ignited a remarkable range of possibilities

All industry sectors

are embracing AI advancements





Explore and Understand

run and accelerate experimentation

Unique challenges

could slow down or even completely block adoption

Strategy

ensure the effective use of AI services



Azure model catalog - flexible deployment options



Managed compute

- Deployed to managed Azure VMs
- One click deployment
- Pay per GPU billing
- 100+ open models



Serverless APIs and
Model as a Service

- Ready to use model APIs
- Pay per token
- Pay per provisioned capacity
- 30+ flagship models

Azure OpenAI Service

GPT-4

GPT-4-Turbo

GPT-3.5-Turbo

GPT-4 for Vision

GPT-4o

DALL·E 3

Generative Text Models, with varying capabilities and uses

Generative
Image Model



Tokens
& Quotas

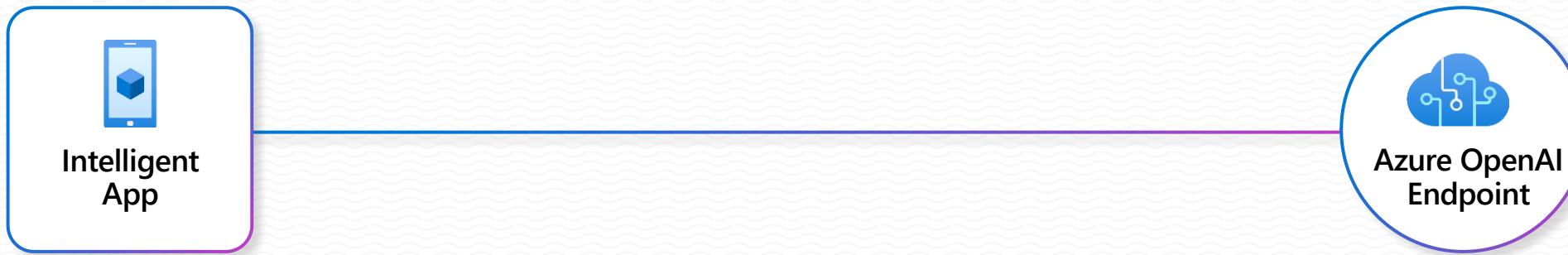


Pay-as-you-go
& Provisioned
Throughput Units (PTUs)

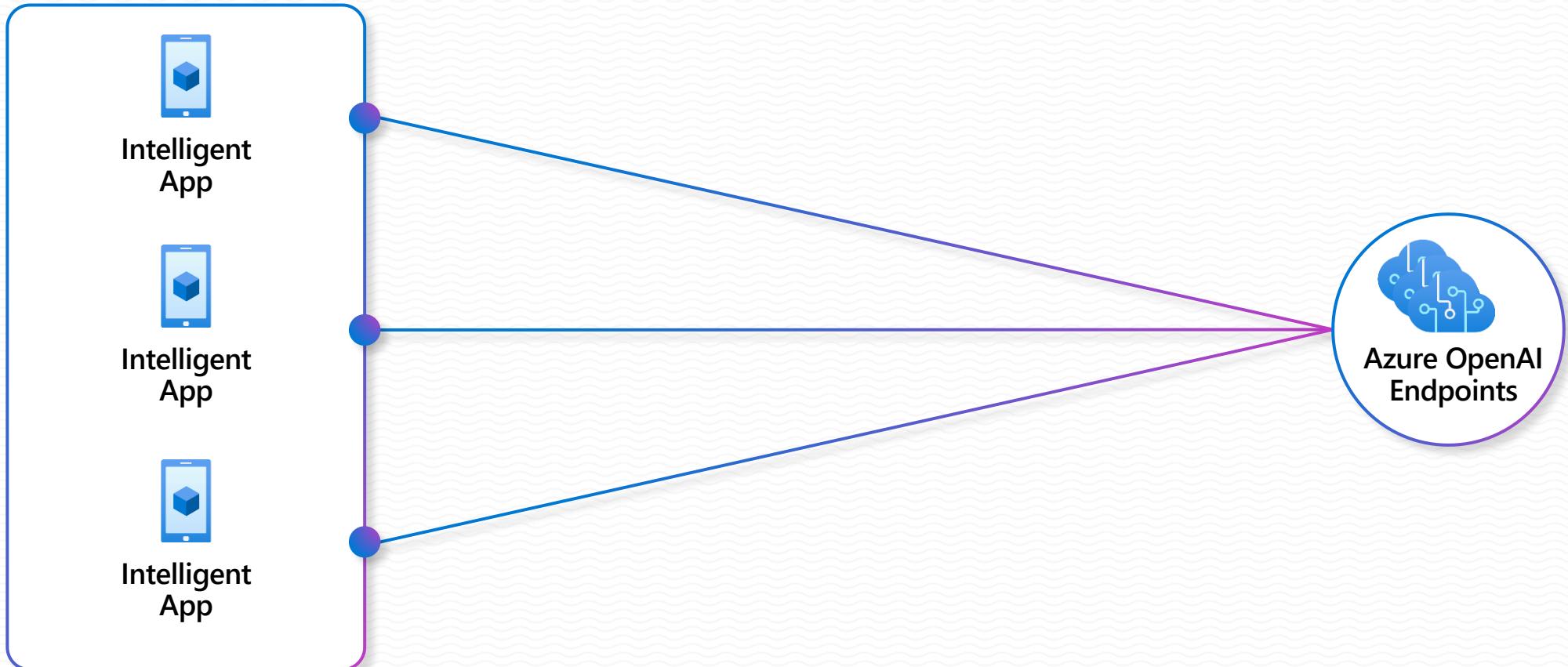


API & SDK

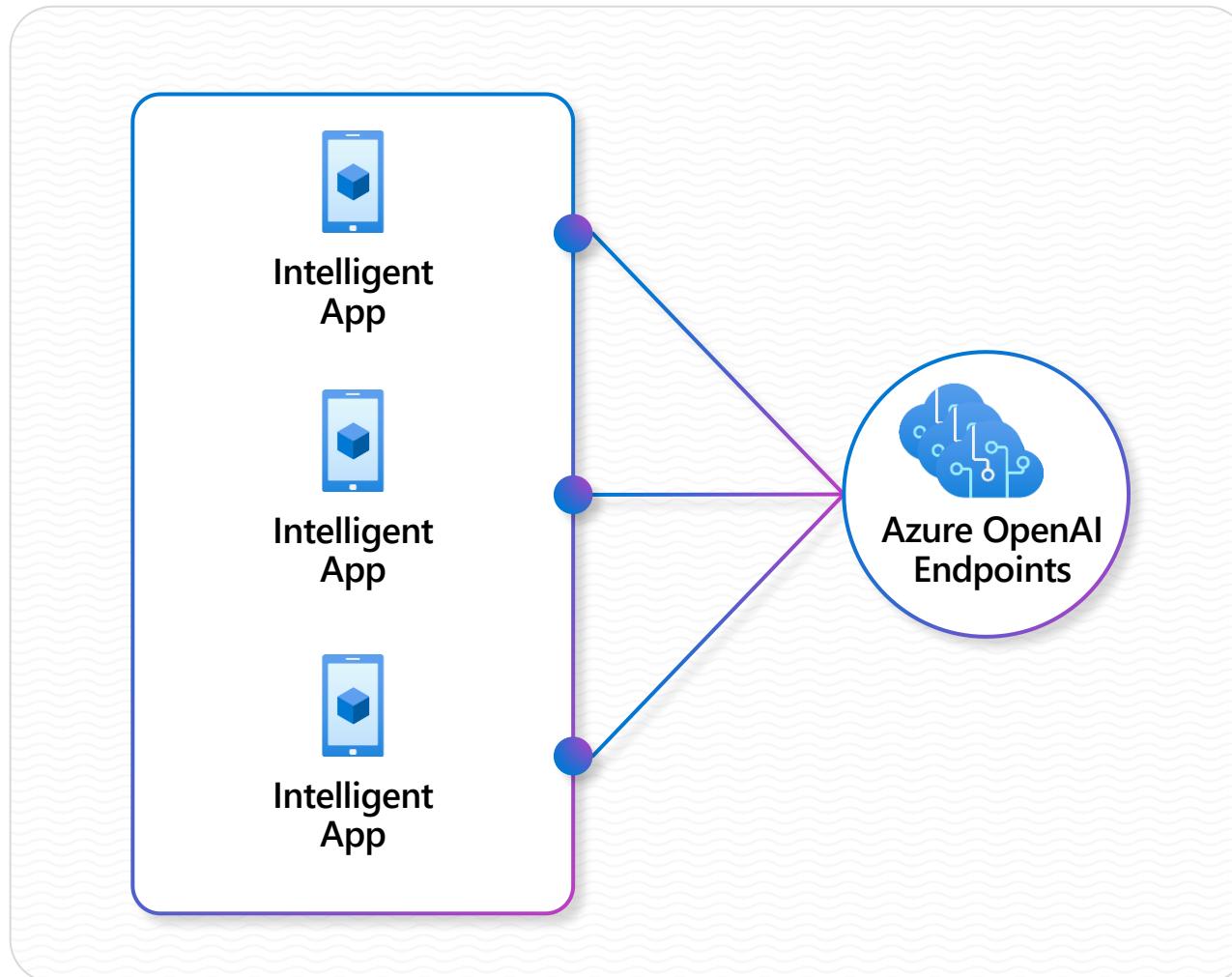
Connecting an Intelligent App with OpenAI



Scaling Up: Multiple Apps, Multiple OpenAI Endpoints



Scaling Up: Multiple Apps, Multiple OpenAI Endpoints



Scaling Challenges

Track token usage



Multiple OpenAI endpoints



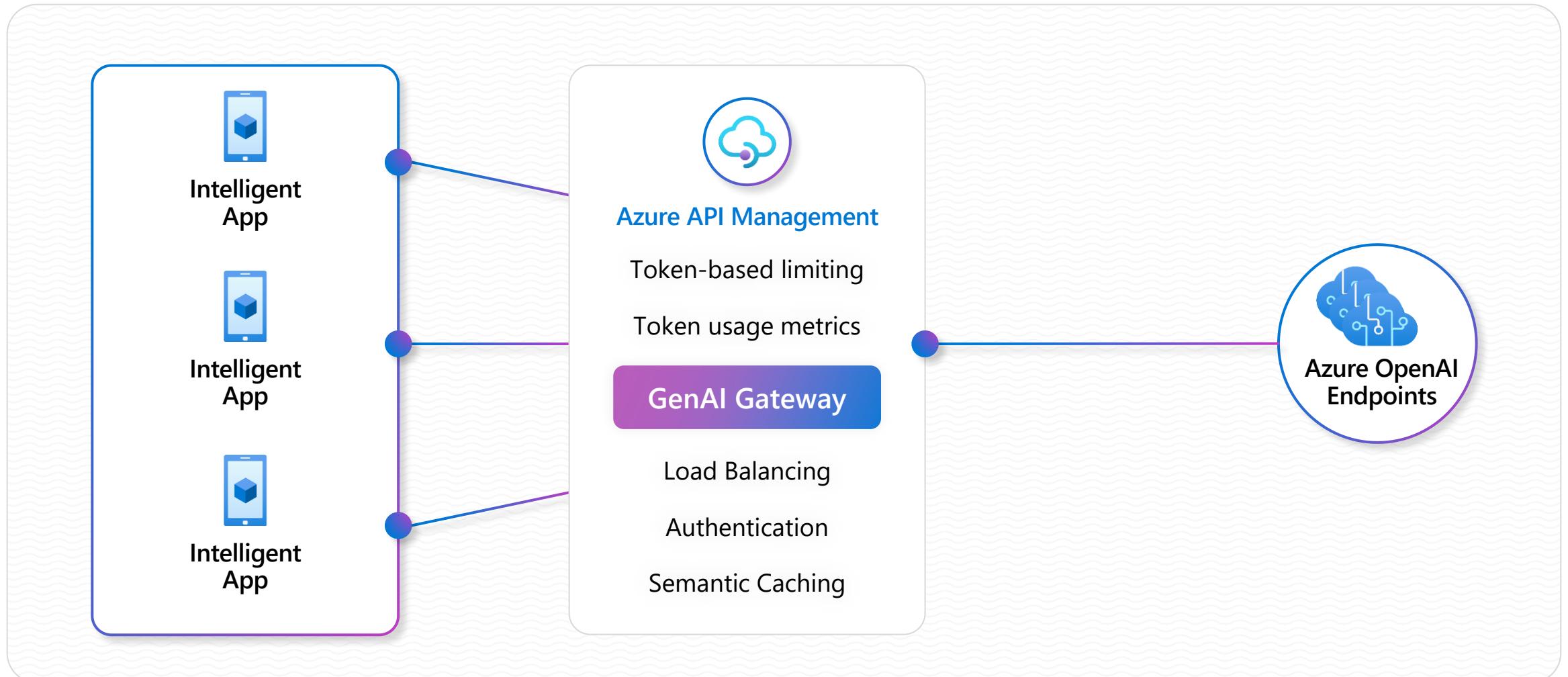
Authentication and authorization



Assign token-based limits



Delegate security, scalability and observability to the GenAI Gateway



Explore GenAI Gateway – where to start?



GenAI Gateway Labs

<https://aka.ms/apim/genai/labs>
(<https://aka.ms/ai-gateway>)

→ experimentation-driven approach
that pave the road to success

README Code of conduct MIT license

APIM ❤️ OpenAI - 🚀 Labs for the [GenAI Gateway](#) capabilities of Azure API Management

Open Source

What's new ⚡

- the [Content filtering](#) and [Prompt shielding](#) labs.
- the [Model routing](#) lab with OpenAI model based routing.
- the [Prompt flow](#) lab to try the [Azure AI Studio Prompt Flow](#) with Azure API Management.
- priority and weight parameters to the [Backend pool load balancing](#) lab.
- the [Streaming](#) tool to test OpenAI streaming with Azure API Management.
- the [Tracing](#) tool to debug and troubleshoot OpenAI APIs using [Azure API Management tracing capability](#).
- image processing to the [GPT-4o inferencing](#) lab.
- the [Function calling](#) lab with a sample API on Azure Functions.



Security

Security



Specialized security policies : OAuth, API key/token validation, schema validation, and threat detection tailored to your API's structure.

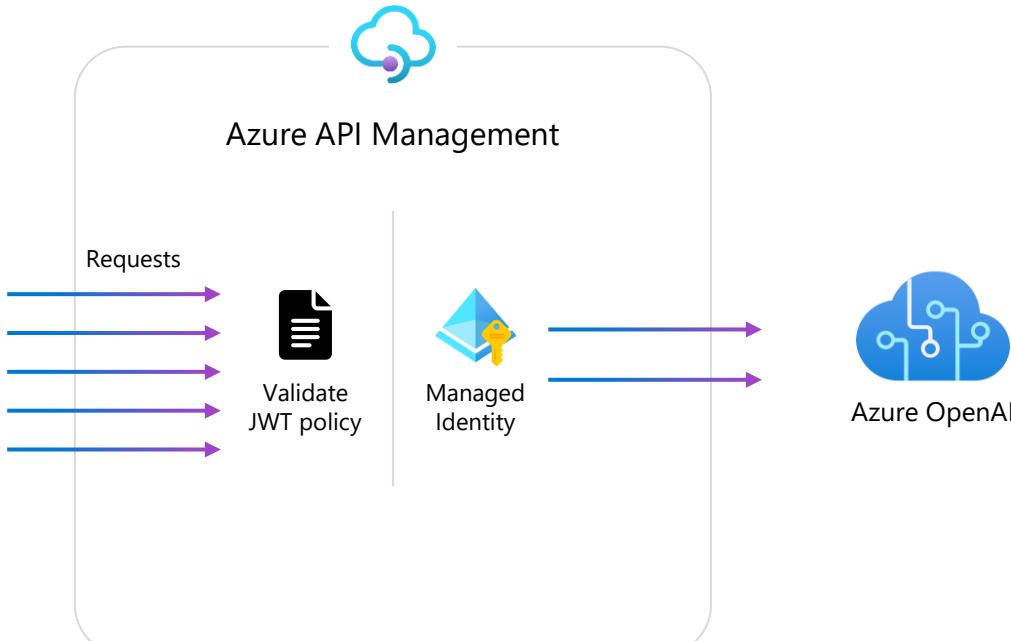


Positive Security Model : explicitly defines what is allowed rather than just relying on detecting abnormalities.



Traffic Control : Rate limits and quotas protect your AI resources from abuse and uncontrolled usage spikes.

Authentication and Authorization



Configure managed identity authentication

Validate claims in JWT to manage access to OpenAI endpoints

Authenticate API consumers using subscription keys

request-forwarding

- policy.xml
- README.MD
- request-forwarding.ipynb

[6] Assign a role to enable APIM to access OpenAI API

This lab uses a zero trust security strategy with a key less approach using an [Azure Managed Identity](#). The following script assigns the `Cognitive Services OpenAI User` role to the OpenAI API.

[Generate](#) [+ Code](#) [+ Markdown](#)

```
if mock_disabled:
    openai_resource_name = openai_resources[0].get("name")
    openai_resource_stdout = ! az cognitiveservices account show --name {openai_resource_name} --resource-group {resource_group}
    openai_resource = json.loads(openai_resource_stdout.n)
    openai_resource_id = openai_resource.get("id")
    role_assignment_stdout = ! az role assignment create --assignee {apim_managed_identity} \
        --role "Cognitive Services OpenAI User" \
        --scope {openai_resource_id}
```

OpenAI > All operations > Policies

HTTP GenAI-Gateway / 01-OpenAI-AuthWith-Apikey

POST https://mcx-aigw-apim.azure-api.net/openai/deployments/gpt-35-turbo/chat/completions?api-version=2024-02-01

Params • Authorization Headers (10) Body • Scripts Tests Settings

Headers • 8 hidden

Key	Value
<input checked="" type="checkbox"/> api-key	3b3643559e344baeb5ec0253bd4a222d
<input checked="" type="checkbox"/> Content-Type	application/json
Key	Value

```
1 <polices>
2   <inbound>
3     <base />
4     <validate-azure-ad-token tenant-id="7517bc42-bcf8-4916-a677-b5753051f846">
5       <client-application-ids>
6         <application-id>b3c82555-3a08-4e66-8e4a-001314c78e74</application-id>
7       </client-application-ids>
8     </validate-azure-ad-token>
9     <authentication-managed-identity resource="https://cognitiveservices.azure.com" output-token-va
10    <set-header name="Authorization" exists-action="override">
11      <value>@("Bearer " + (string)context.Variables["managed-id-access-token"])</value>
12    </set-header>
13    <set-backend-service backend-id="openai-backend-pool" />
14  </inbound>
15  <backend>
```

Body Cookies Headers (19) Test Results

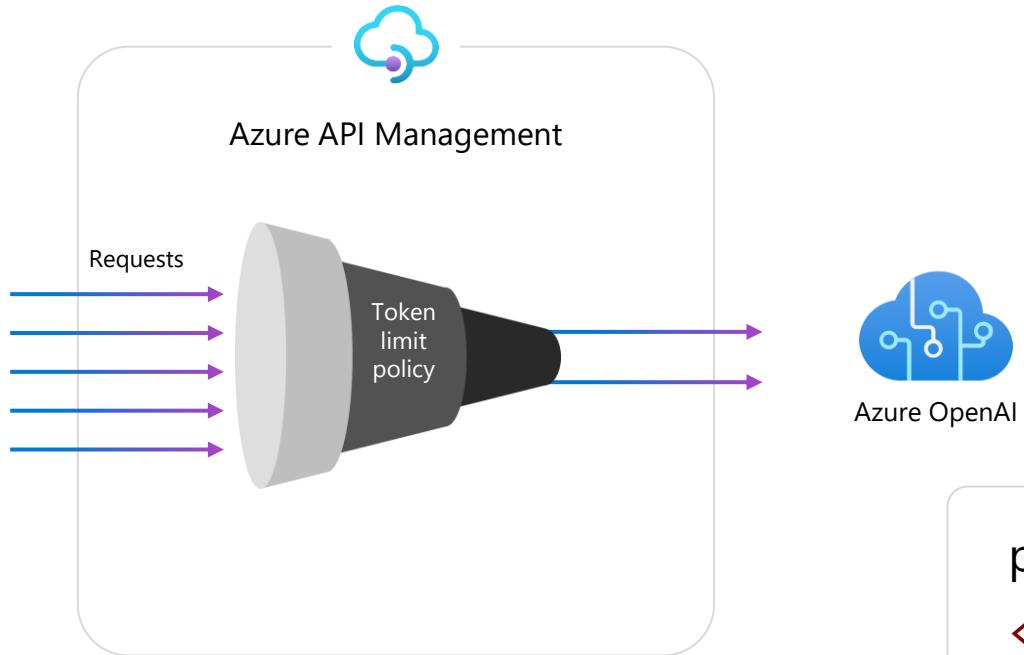
Pretty Raw Preview Visualize JSON

```
19   },
20   "severity": "safe"
21 },
22 "finish_reason": "stop",
23 "index": 0,
24 "message": {
25   "content": "Oh sure, let me just consult my extensive collection of antique sundials an
26   "role": "assistant"
27 }
```

Token count

timestamp	model	region	promptTokens	completionTokens	totalTokens	remainingTokens	remainingRequests
10/18/2024, 10:17:09.534 AM		France Central	32	76	108	19952	18
10/18/2024, 10:17:06.337 AM		Sweden Central	32	31	63	19952	18
10/18/2024, 10:17:01.652 AM		France Central	32	37	69	19968	18
10/18/2024, 10:16:58.415 AM		Sweden Central	32	85	117	19968	19
10/18/2024, 10:16:56.518 AM		France Central	32	19	51	19984	19

[GA] Azure OpenAI Token Limit policy



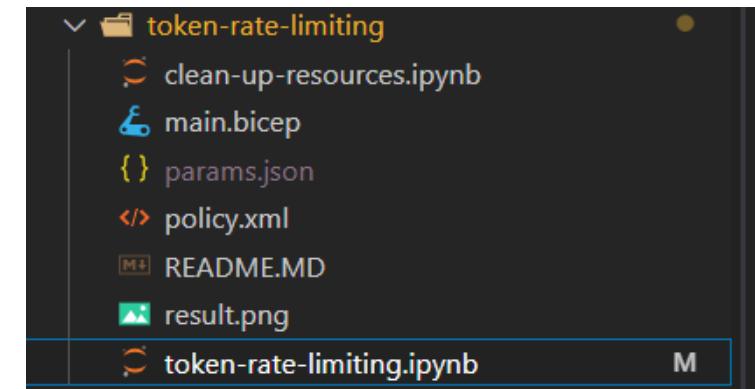
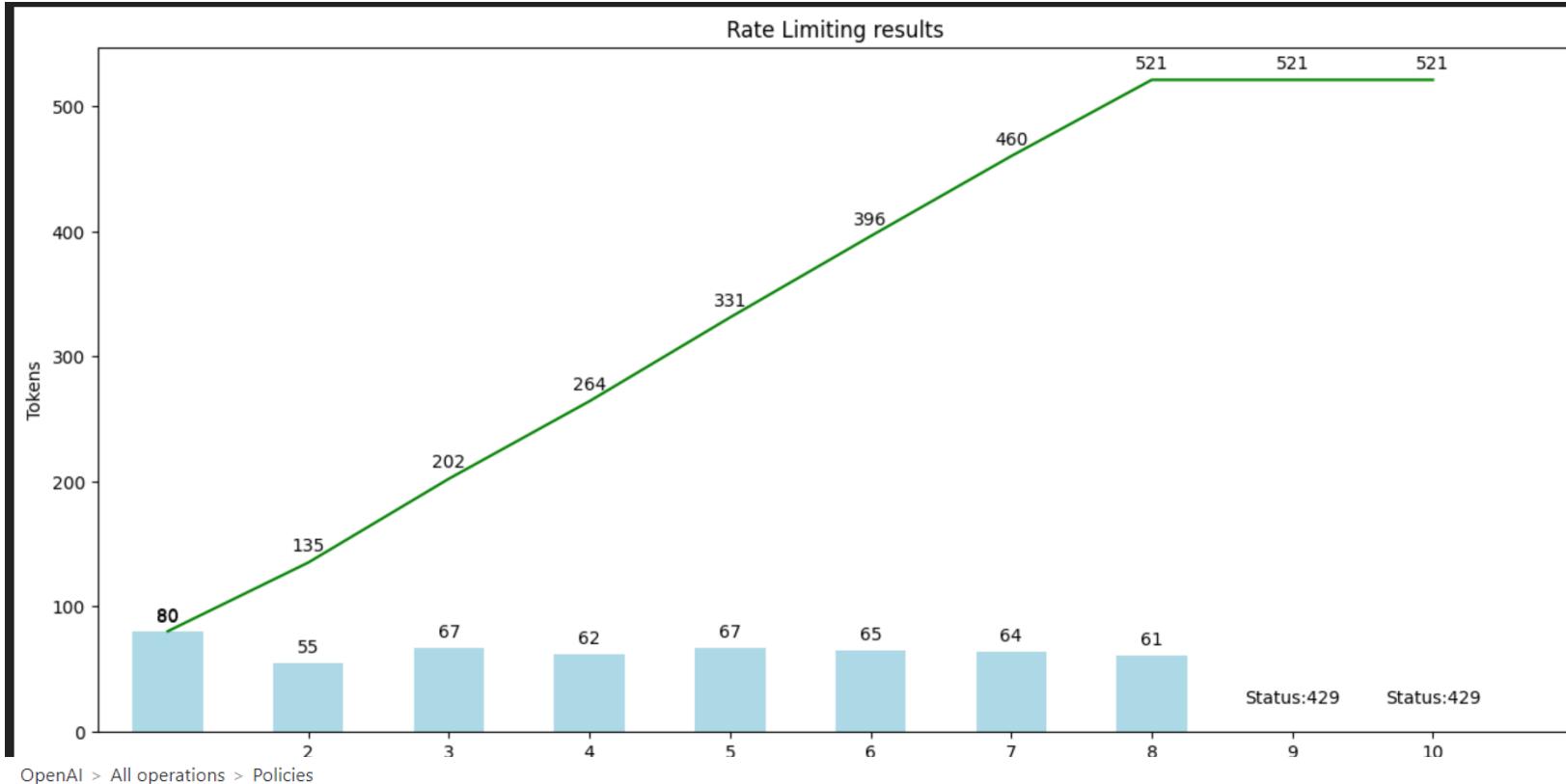
Configure tokens per minute (TPM) limits based on counter keys



Define policy behavior for throttling

policy.xml

```
<azure-openai-token-limit  
    counter-key="@({context.Subscription.Id})"  
    tokens-per-minute="1000"  
    estimate-prompt-tokens="false" />
```



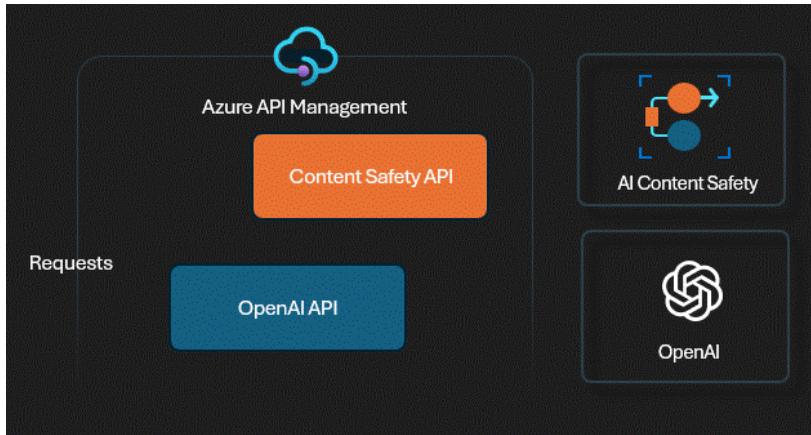
```

1 <policies>
2   <inbound>
3     <base />
4     <authentication-managed-identity resource="https://cognitiveservices.azure.com" output-token-variable-name="managed-id-access-token" ignore-error="false" />
5     <set-header name="Authorization" exists-action="override">
6       <value>@(Bearer " + (string)context.Variables["managed-id-access-token"])</value>
7     </set-header>
8     <set-backend-service backend-id="openai1" />
9     <azure-openai-token-limit counter-key="@({context.RequestIpAddress})" tokens-per-minute="500" estimate-prompt-tokens="false" remaining-tokens-variable-name="remainingTokens" />
10   </inbound>
11   <backend>
12     <base />
13   </backend>
14   <outbound>
15     <base />
16   </outbound>
17   <on-error>
18     <base />
19   </on-error>
20 </policies>

```

Deployment type	Fine-tune	Capacity
Standard		20K TPM

Prompt guarding, prompt shielding



```
<policies>
<inbound>
    <!-- Get the request body and store it in a context variable -->
    <base />
    <set-variable name="requestBody" value="@((context.Request.Body.As<string>(preserveContent: true) ?? ""))" />
    <!-- Define a list of prohibited words or phrases -->
    <set-variable name="prohibitedWords" value="Zeus,Apollo" />
    <!-- Check if the request body contains any prohibited content -->
    <choose>
        <when condition="@{
            var prohibitedWords = context.Variables.GetValueOrDefault<string>("prohibitedWords").Split(',');
            var requestBody = context.Variables.GetValueOrDefault<string>("requestBody");
            return prohibitedWords.Any(word => requestBody.IndexOf(word, StringComparison.OrdinalIgnoreCase) >= 0);
        }">
            <!-- If prohibited content is found, block the request and return an error response -->
            <return-response>
                <set-status code="400" reason="Bad Request" />
                <set-body>@{
                    var error = new {
                        error = new {
                            code = 400,
                            message = "Bad Request",
                            description = "Prompt Guard - Ancient Greek mythology not allowed.",
                            details = new {
                                errorType = "InvalidInput",
                                inputPrompt = "Zeus, Apollo and similar"
                            }
                        };
                    return Newtonsoft.Json.JsonConvert.SerializeObject(error);
                }</set-body>
            </return-response>
        </when>
    </choose>
</inbound>
```



Check against a list of allowed or denied expressions



Analyzes LLM inputs and detects User Prompt attacks and Document attacks



Leverage the APIM policies to build additional Generative AI capabilities.



Load Balancing

Load Balancing



Intelligent Load Balancing : Distribute traffic intelligently (location, weight, priority, ..) across different AI endpoints deployed in multiple regions.

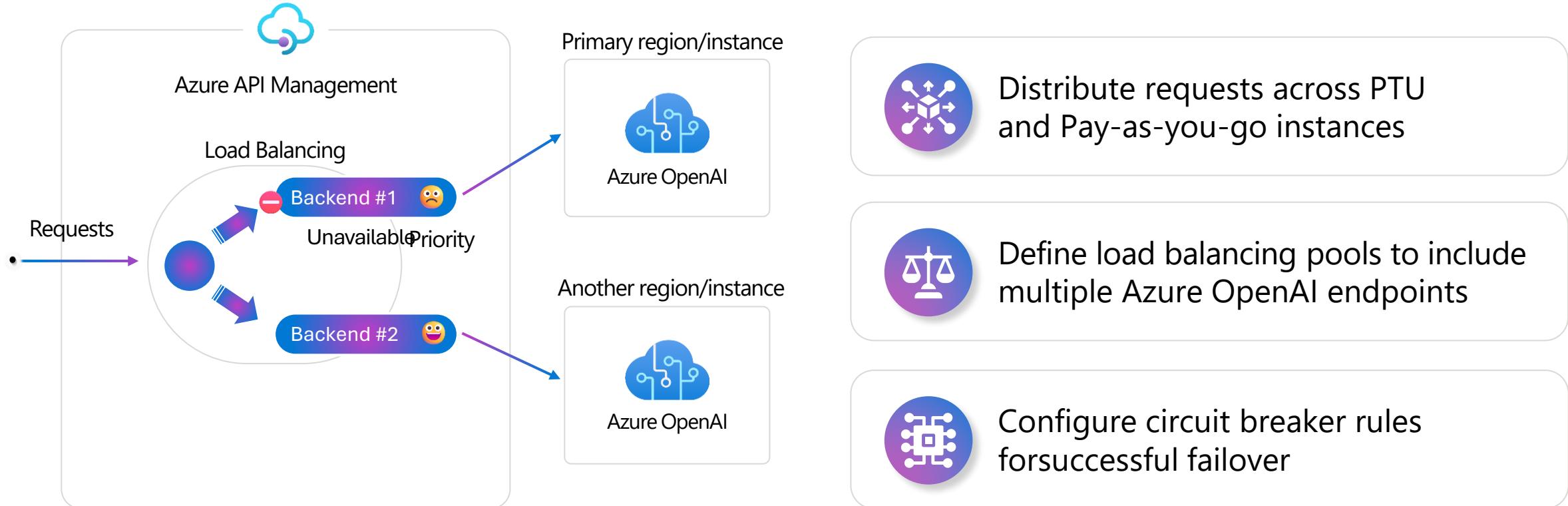


Failover and Resilience : detect failures and automatically redirect traffic to healthy endpoints, improving your application's reliability.



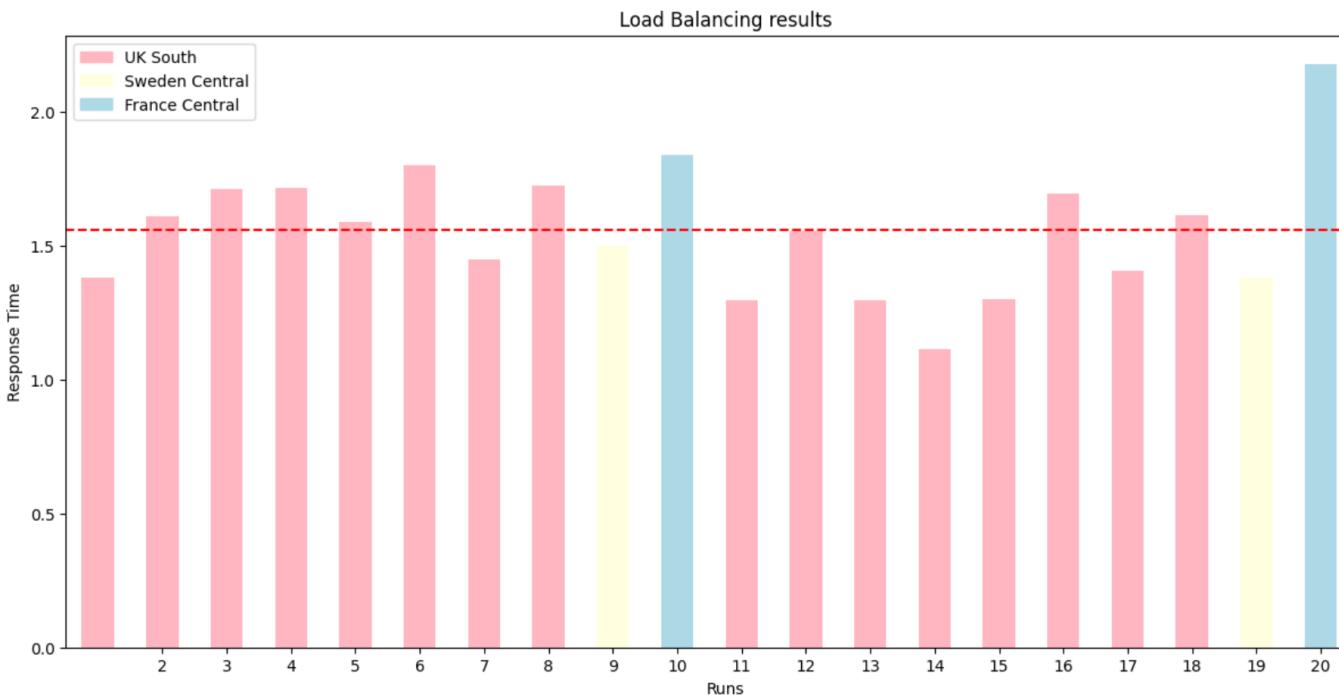
A/B Testing : seamlessly transition across multiple versions of deployed models without disrupting existing applications.

[GA] Load Balancer and Circuit Breaker



backend-pool-load-balancing

- backend-pool-load-balancing.ipynb
- clean-up-resources.ipynb
- main.bicep
- {} params.json
- <> policy.xml
- README.MD
- result.png
- > built-in-logging
- > content-filtering
- > developer-tooling
- > function-calling
- > GPT-4o-inferencing
- > message-storing
- > model-routing
- > prompt-flow
- > request-forwarding
 - <> policy.xml
 - README.MD
 - request-forwarding.ipynb
- > response-streaming
- > semantic-caching
- > slim-self-hosting
- > token-metrics-emitting
- > token-rate-limiting
- > vector-searching
- > tools



```
resource backendOpenAI 'Microsoft.ApiManagement/service/backends@2023-09-01-preview' = [for (config, i) in openAIConfig: if(length(openAIConfig) > 0) {
  name: config.name
  parent: apimService
  properties: {
    description: 'backend description'
    url: '${cognitiveServices[i].properties.endpoint}/openai'
    protocol: 'http'
    circuitBreaker: {
      rules: [
        {
          failureCondition: [
            count: 1
            errorReasons: [
              'Server errors'
            ]
            interval: 'PT5M'
            statusCodeRanges: [
              {
                min: 429
                max: 429
              }
            ]
          ]
          name: 'openAIBreakerRule'
          tripDuration: 'PT1M'
          acceptRetryAfter: true // respects the Retry-After header
        }
      ]
    }
  }
}
```

```
resource backendPoolOpenAI 'Microsoft.ApiManagement/service/backends@2023-09-01-preview' = if(length(openAIConfig) > 1) {
  name: openAIBackendPoolName
  parent: apimService
  properties: {
    description: openAIBackendPoolDescription
    type: 'Pool'
    // protocol: 'http' // the protocol is not needed in the Pool type
    // url: '${cognitiveServices[0].properties.endpoint}/openai' // the url is not needed in the Pool type
    pool: {
      services: [for (config, i) in openAIConfig: {
        id: '/backends/${backendOpenAI[i].name}'
        priority: config.priority
        weight: config.weight
      }]
    }
}
```

Yellow arrows point to the 'statusCodeRanges' field in the first code block and the 'pool' field in the second code block, highlighting specific configuration details.



Observability

Observability



Visibility into the AI usage patterns and quotas. By proxying AI calls, we gain a unified view of usage, performance, cost, health and more.

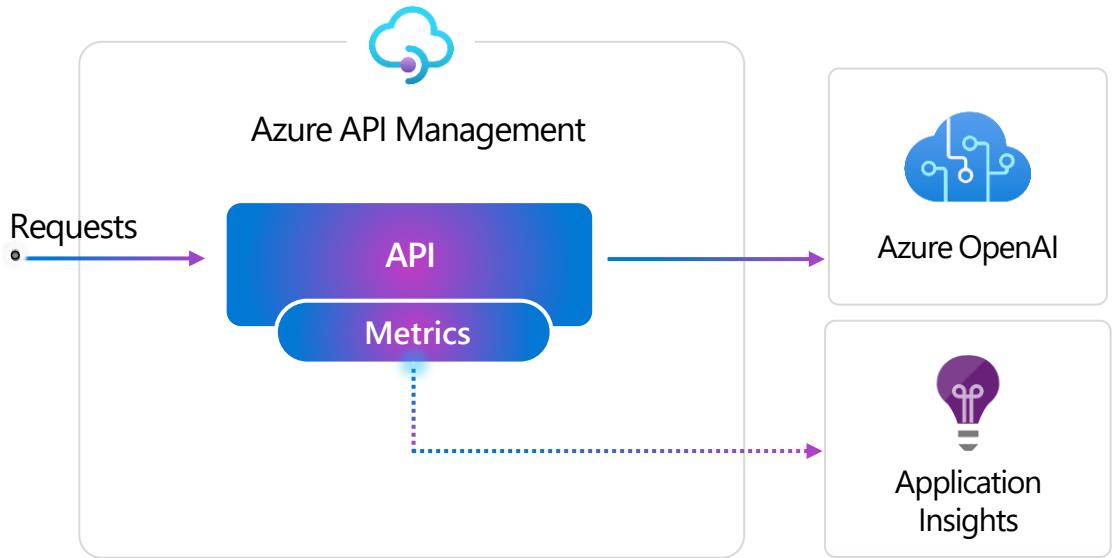


Get data for debugging, troubleshooting, and the metrics necessary to ensure reliability and availability.



End to end traceability for your Intelligent applications.

[GA] Azure OpenAI Token Metric policy

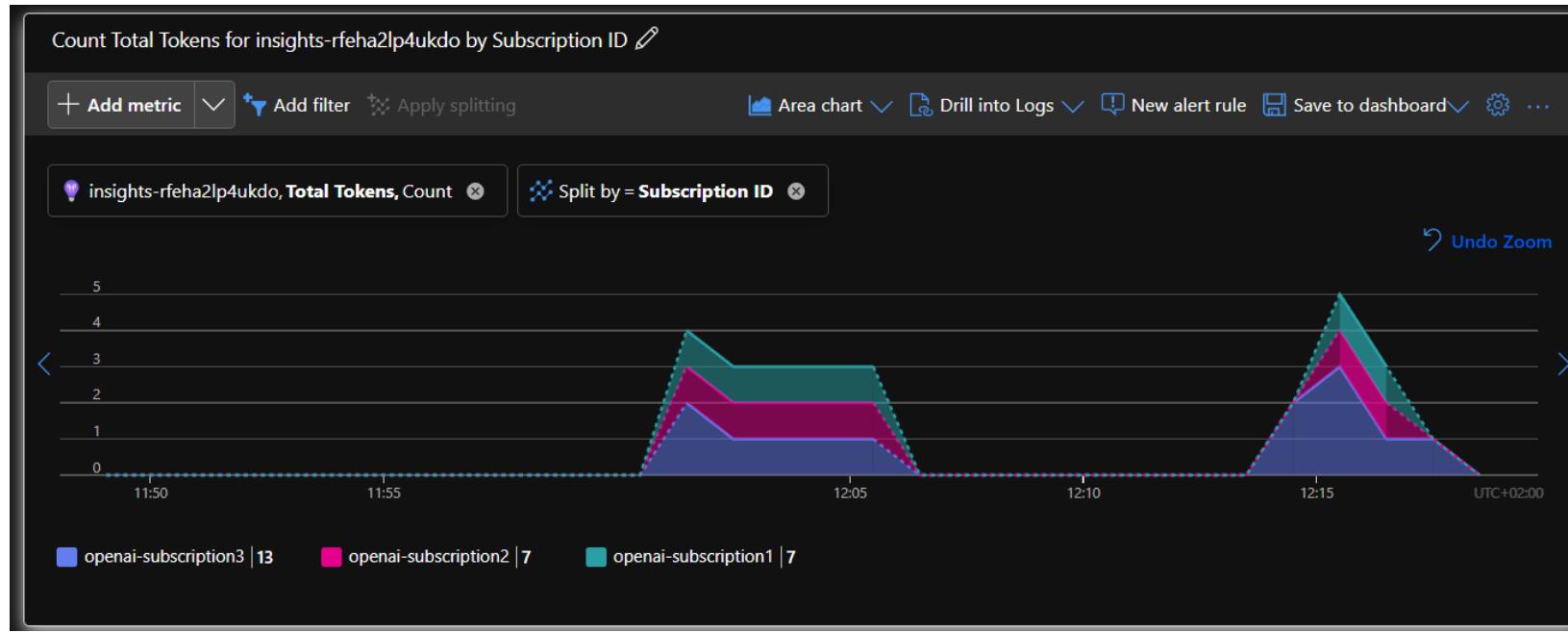


Facilitate accurate cross-charging
based on token consumption

Collect token usage data

policy.xml

```
<azure-openai-emit-token-metric  
namespace="AzureOpenAI">  
    <dimension name="User ID" />  
    <dimension name="Subscription ID"  
/>  
</azure-openai-emit-token-metric>
```

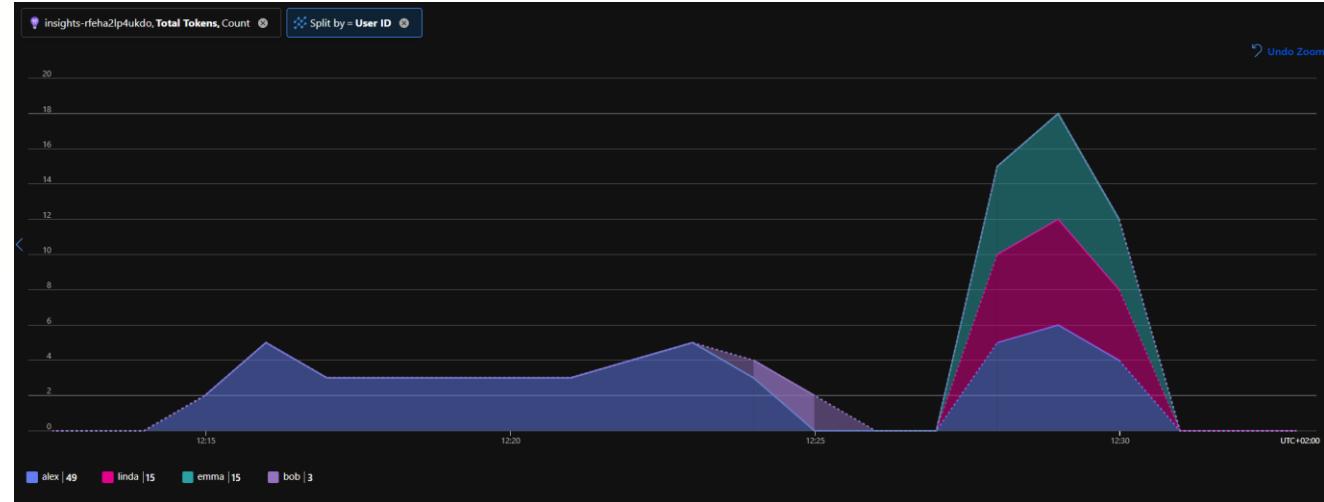


OpenAI > All operations > Policies

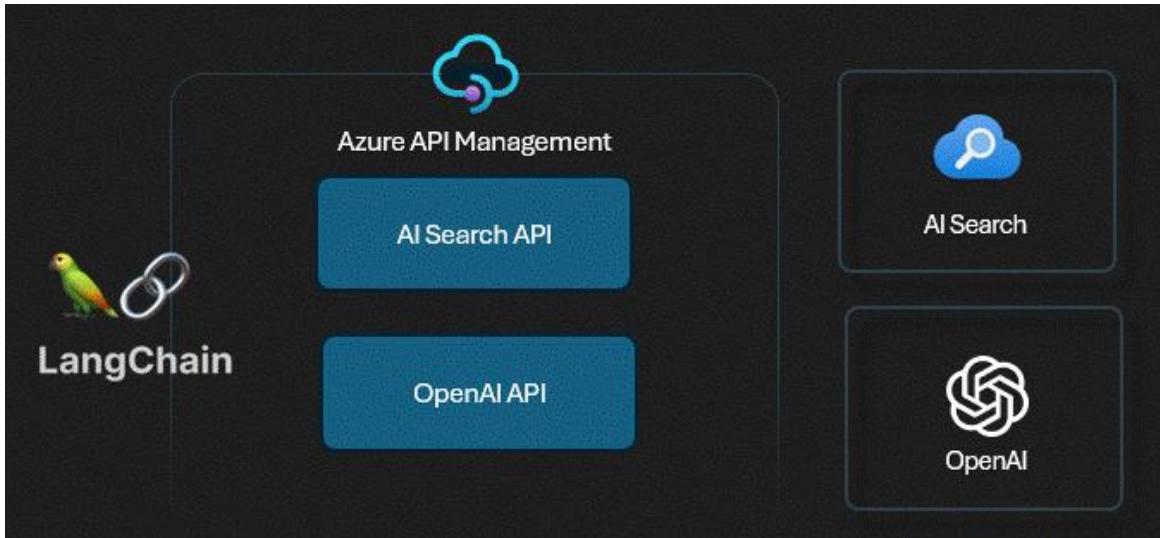
```

1 <policies>
2   <inbound>
3     <base />
4     <authentication-managed-identity resource="https://cognitiveservices.azure.com" output-token-variable-name="managed-id-access-token" ignore-error="false" />
5     <set-header name="Authorization" exists-action="override">
6       <value>@("Bearer " + (string)context.Variables["managed-id-access-token"])</value>
7     </set-header>
8     <set-backend-service backend-id="openai1" />
9     <azure-openai-emit-token-metric namespace="openai">
10       <dimension name="Subscription ID" value="@({context.Subscription.Id})" />
11       <dimension name="Client IP" value="@({context.RequestIpAddress})" />
12       <dimension name="API ID" value="@({context.Api.Id})" />
13       <dimension name="User ID" value="@({context.Request.Headers.GetValueOrDefault("x-user-id", "N/A")})" />
14     </azure-openai-emit-token-metric>
15   </inbound>
16   <backend>
17     <base />
18   </backend>
19   <outbound>
20     <base />

```



Explore the RAG pattern : built-in logging



The requests are logged into Application Insights and metrics available in Azure Monitor.

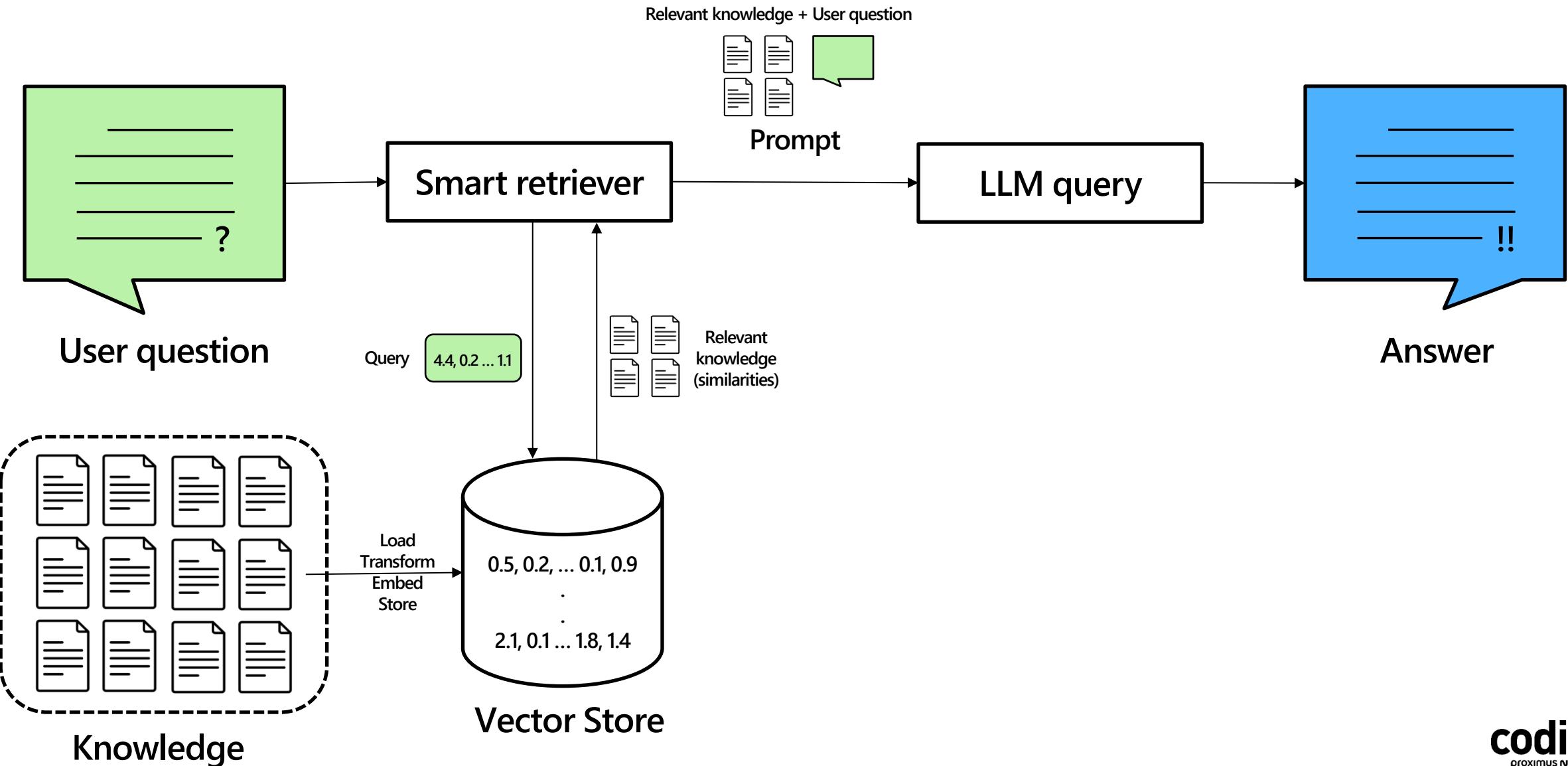


Enables tracking request/response details and token usage with the provided notebook

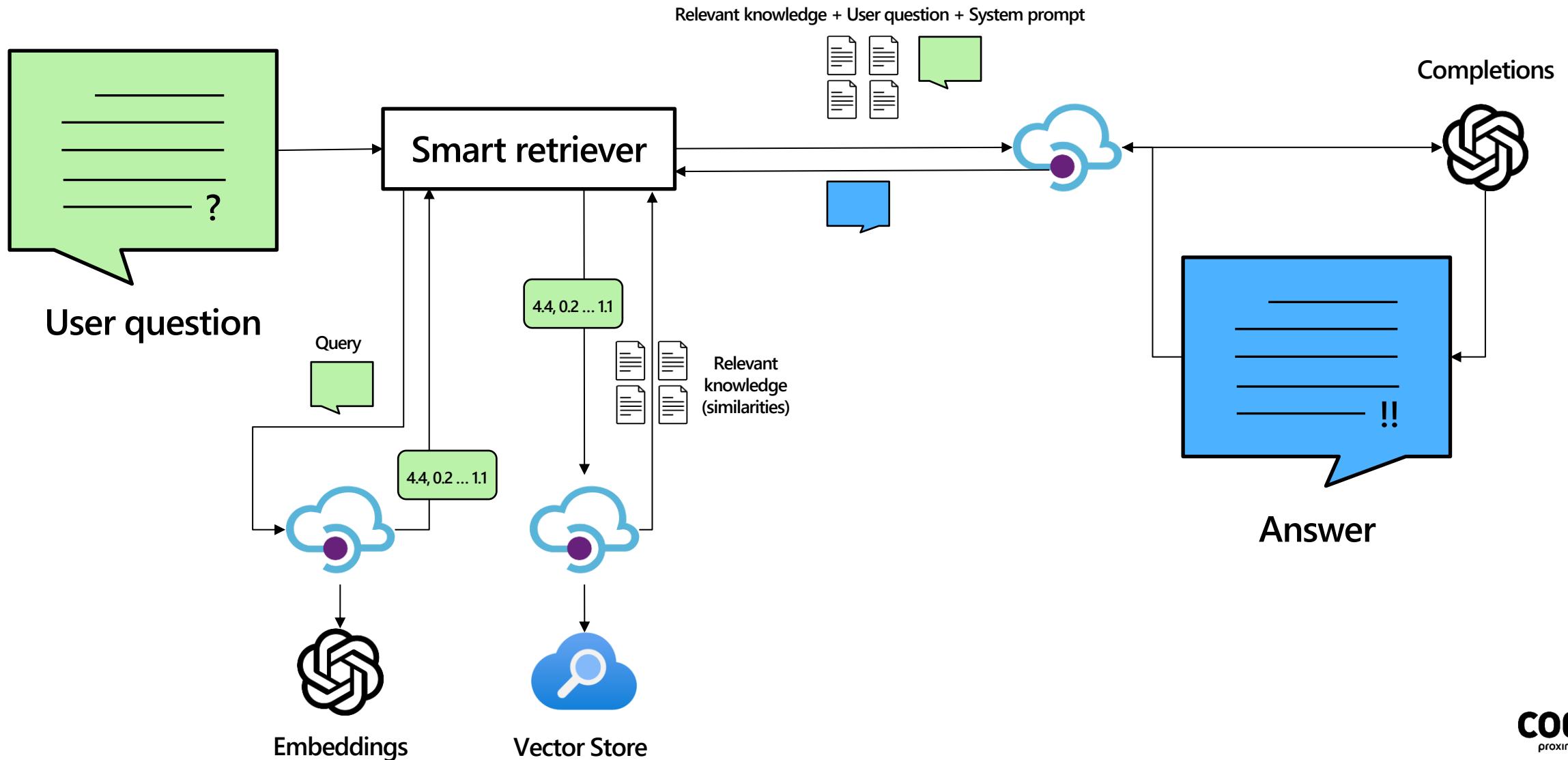


Enables the creation of Azure dashboards for a single pane of glass monitoring approach.

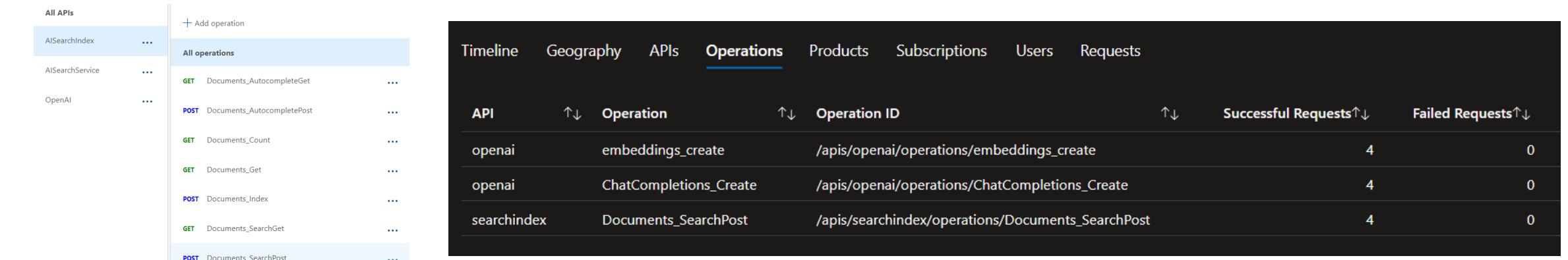
Explore the RAG pattern



AI Gateway : RAG pattern

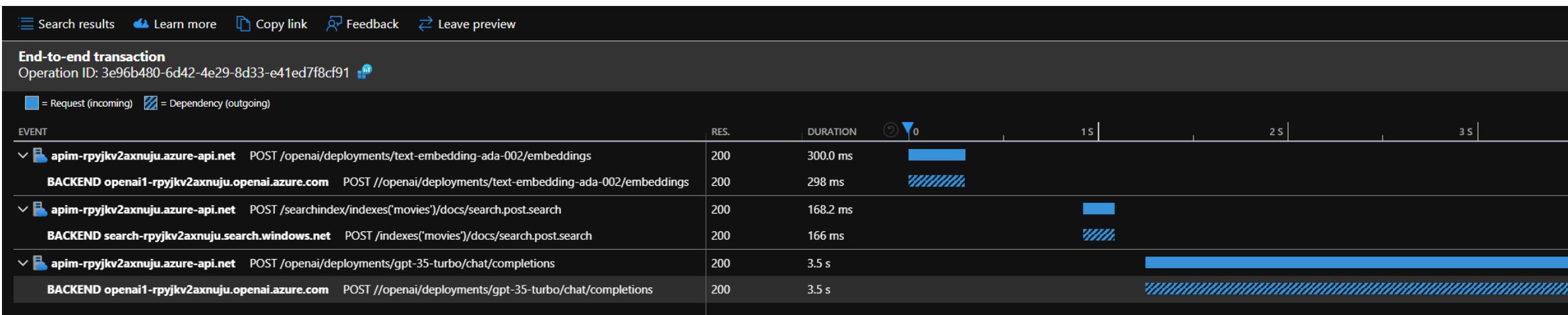


Explore the RAG pattern : built-in logging



The screenshot shows the Azure API Management Operations blade. On the left, a sidebar lists "All APIs" with sections for AISearchIndex, AISystemService, and OpenAI. Under OpenAI, several operations are listed with their HTTP methods and URLs. On the right, the main area has tabs for Timeline, Geography, APIs, Operations (which is selected), Products, Subscriptions, Users, and Requests. The Operations table displays the following data:

API	Operation	Operation ID	Successful Requests	Failed Requests
openai	embeddings_create	/apis/openai/operations/embeddings_create	4	0
openai	ChatCompletions_Create	/apis/openai/operations/ChatCompletions_Create	4	0
searchindex	Documents_SearchPost	/apis/searchindex/operations/Documents_SearchPost	4	0

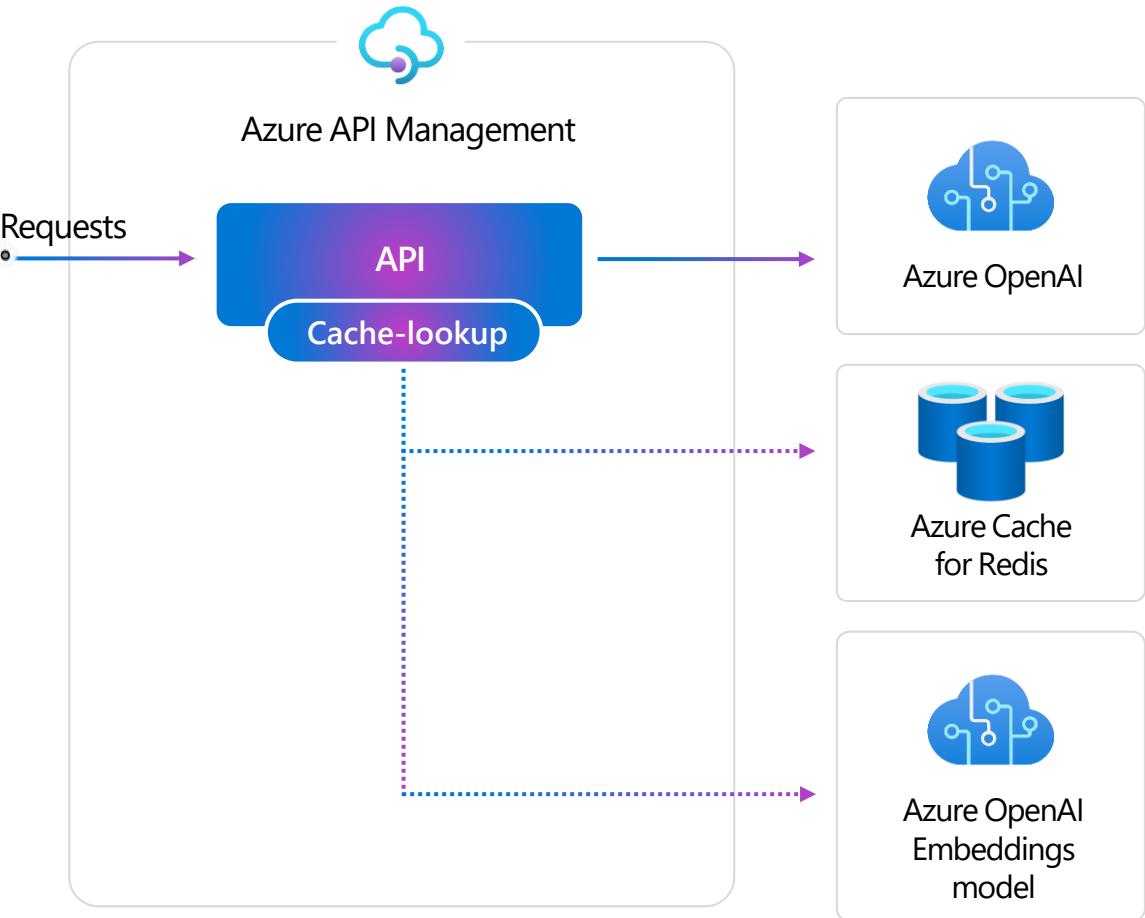


The screenshot shows the Azure Application Insights Log Analytics blade. At the top, there are navigation links: Search results, Learn more, Copy link, Feedback, and Leave preview. Below that, a section titled "End-to-end transaction" shows an operation ID: 3e96b480-6d42-4e29-8d33-e41ed7f8cf91. A legend indicates that blue squares represent "Request (incoming)" and blue diagonal stripes represent "Dependency (outgoing)". The main area displays a table of events with columns for EVENT, RES., DURATION, and a timeline bar. The events show requests from "apim-rpyjkv2axnuju.azure-api.net" and dependencies to "BACKEND" services like "openai1-rpyjkv2axnuju.openai.azure.com" and "search-rpyjkv2axnuju.search.windows.net".



... many more

[Preview] Azure OpenAI Semantic Caching policy



Configure semantic caching for all API consumers

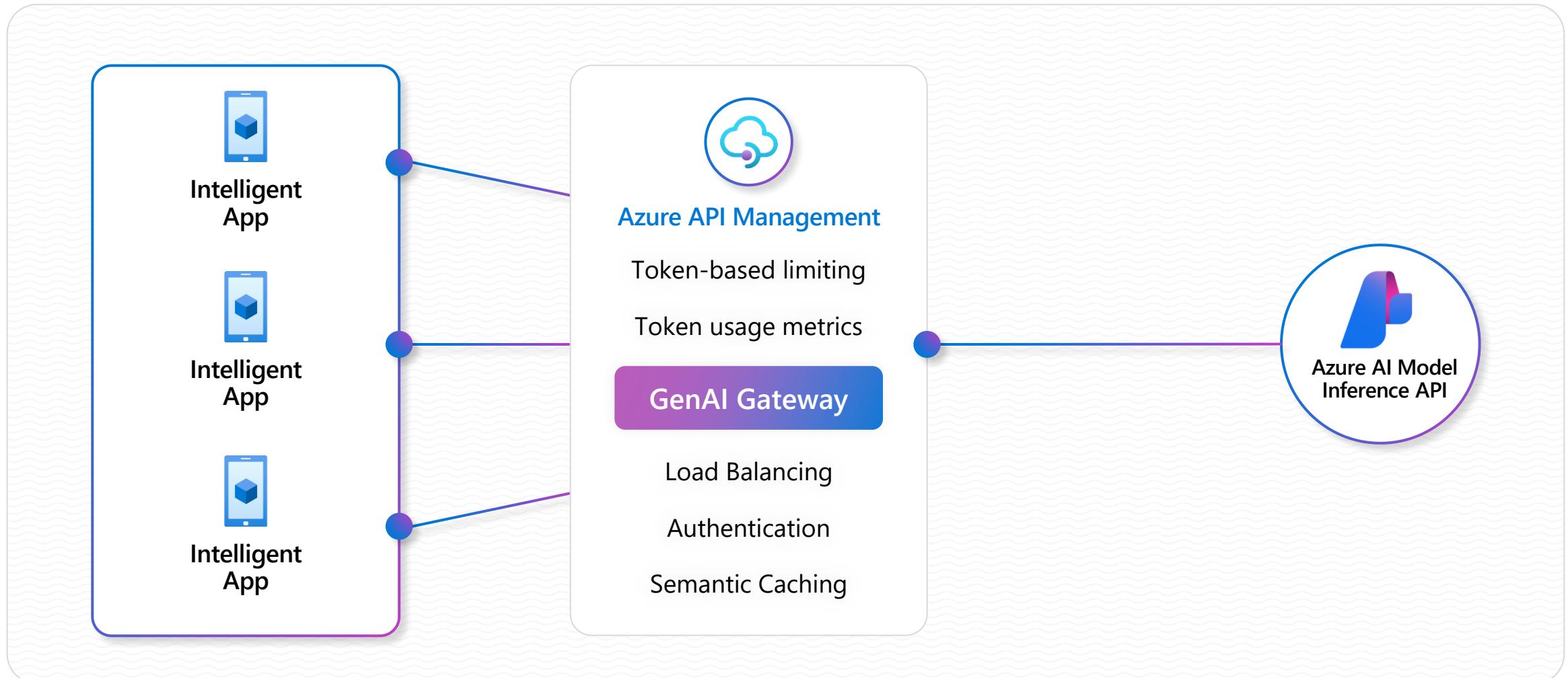


Define similarity score threshold for caching

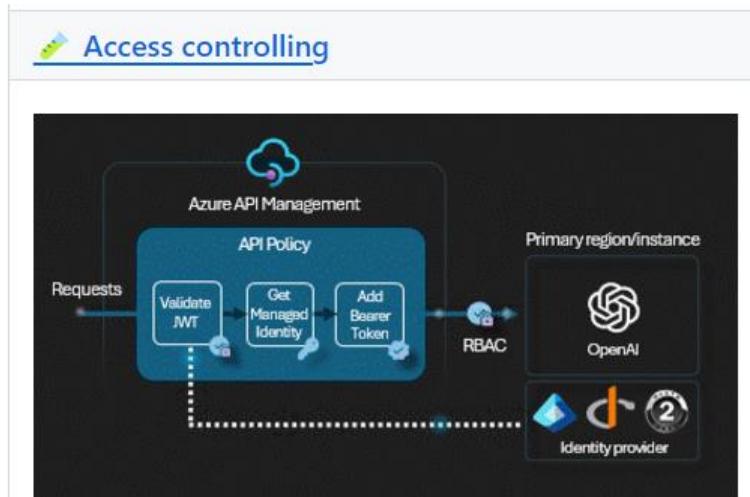
policy.xml

```
<azure-openai-semantic-cache-lookup  
    score-threshold="0.05"  
    embeddings-backend-id="azure-openai-backend">  
    <vary-by>@(context.Subscription.Id)"</vary-by>  
</azure-openai-semantic-cache-lookup>
```

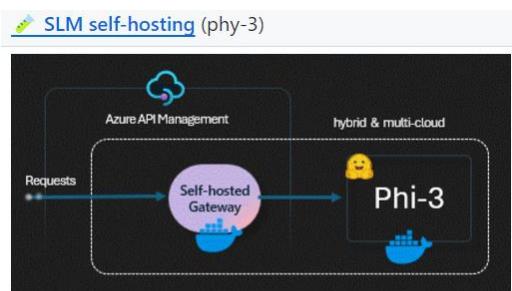
[Preview] Azure AI Model Inference API Support



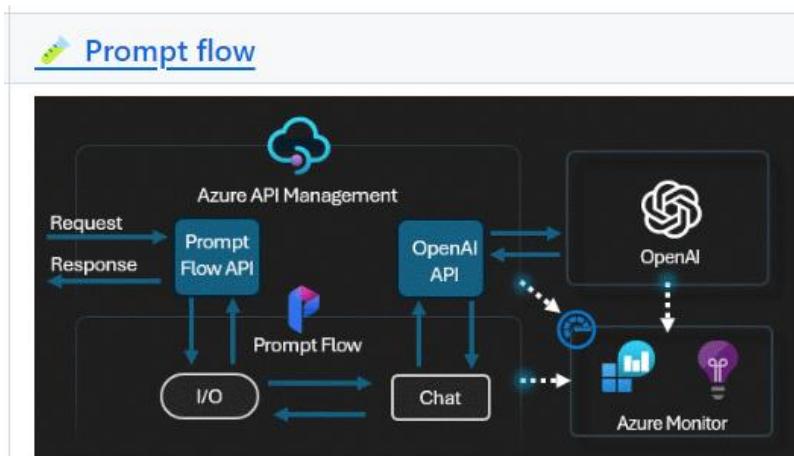
... many more



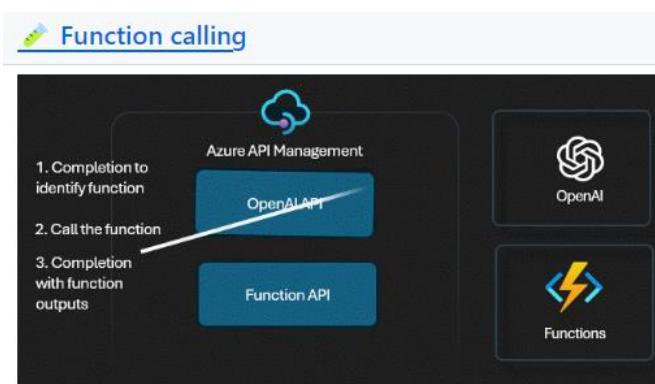
Playground to try the [OAuth 2.0 authorization feature](#) using identity provider to enable more fine-grained access to OpenAPI APIs by particular users or client.



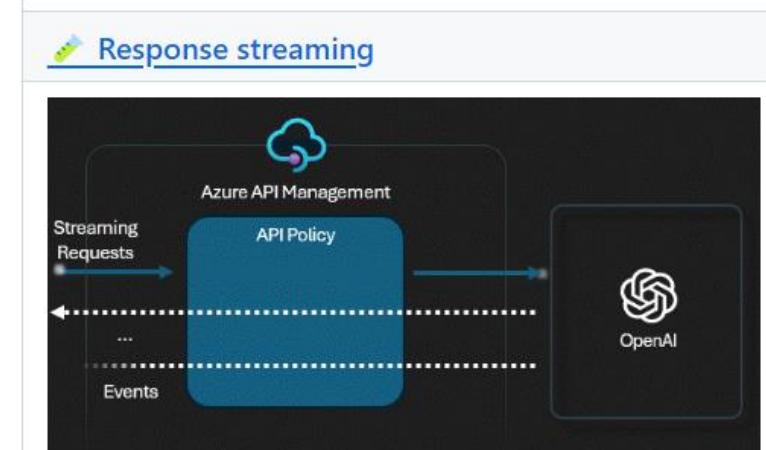
Playground to try the self-hosted [phy-3 Small Language Model \(SLM\)](#) through the [Azure API Management self-hosted gateway](#) with OpenAI API compatibility.



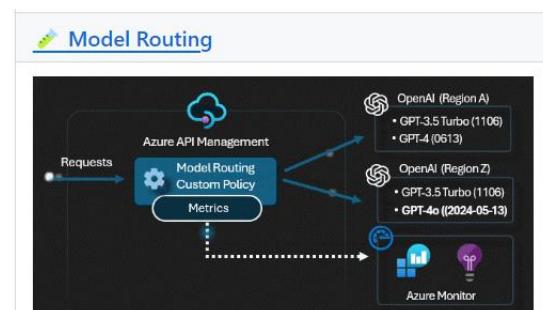
Playground to try the [Azure AI Studio Prompt Flow](#) with Azure API Management.



Playground to try the OpenAI [function calling](#) feature with an Azure Functions API that is also managed by Azure API Management.



Playground to try response streaming with Azure API Management and Azure OpenAI endpoints to explore the advantages and shortcomings associated with [streaming](#).



Playground to try routing to a backend based on Azure OpenAI model and version.

Get Started with Azure API Management and OpenAI



Documentation

<https://aka.ms/apim/openai-docs>



AI Hub Gateway Accelerator

<https://aka.ms/ai-hub-gateway>



GenAI Gateway Labs

<https://aka.ms/apim/genai/labs>



GenAI Gateway Accelerator

<https://aka.ms/apim-genai-lza>

Thank you