Back to Blog    ‹ Newer Article    Older Article ›

···

# Create an Azure OpenAI, LangChain, ChromaDB, and Chainlit Chat App in Container Apps using Terraform

By  Paolo Salvatori

Published Jul 27 2023 06:47 AM              17.8K Views

This article shows how to quickly build chat applications using Python and leveraging powerful technologies such as OpenAI ChatGPT models, Embedding models, LangChain framework, ChromaDB vector database, and Chainlit, an open-source Python package that is specifically designed to create user interfaces (UIs) for AI applications. These applications are hosted on Azure Container Apps, a fully managed environment that enables you to run microservices and containerized applications on a serverless platform.

- **Simple Chat**: This simple chat application utilizes OpenAI's language models to generate real-time completion responses.
- **Documents QA Chat**: This chat application goes beyond simple conversations. Users can upload up to 10 `.pdf` and `.docx` documents, which are then processed to create vector embeddings. These embeddings are stored in ChromaDB for efficient retrieval. Users can pose questions about the uploaded documents and view the Chain of Thought, enabling easy exploration of the reasoning process. The completion message contains links to the text chunks in the documents that were used as a source for the response.

Both applications use a user-defined managed identity to authenticate and authorize against Azure OpenAI Service (AOAI) and Azure Container Registry (ACR) and use Azure Private Endpoints to connect privately and securely to these services. The chat UIs are built using Chainlit, an open-source Python package designed

explicitly for creating AI applications. Chainlit seamlessly integrates with <u>LangChain</u>, <u>LlamaIndex</u>, and <u>LangFlow</u>, making it a powerful tool for easily developing ChatGPT-like applications.

By following our example, you can quickly create sophisticated chat applications that utilize cutting-edge technologies, empowering users with intelligent conversational capabilities.

You can find the code and Visio diagrams in the companion <u>GitHub</u> repository. Also, check the following articles:
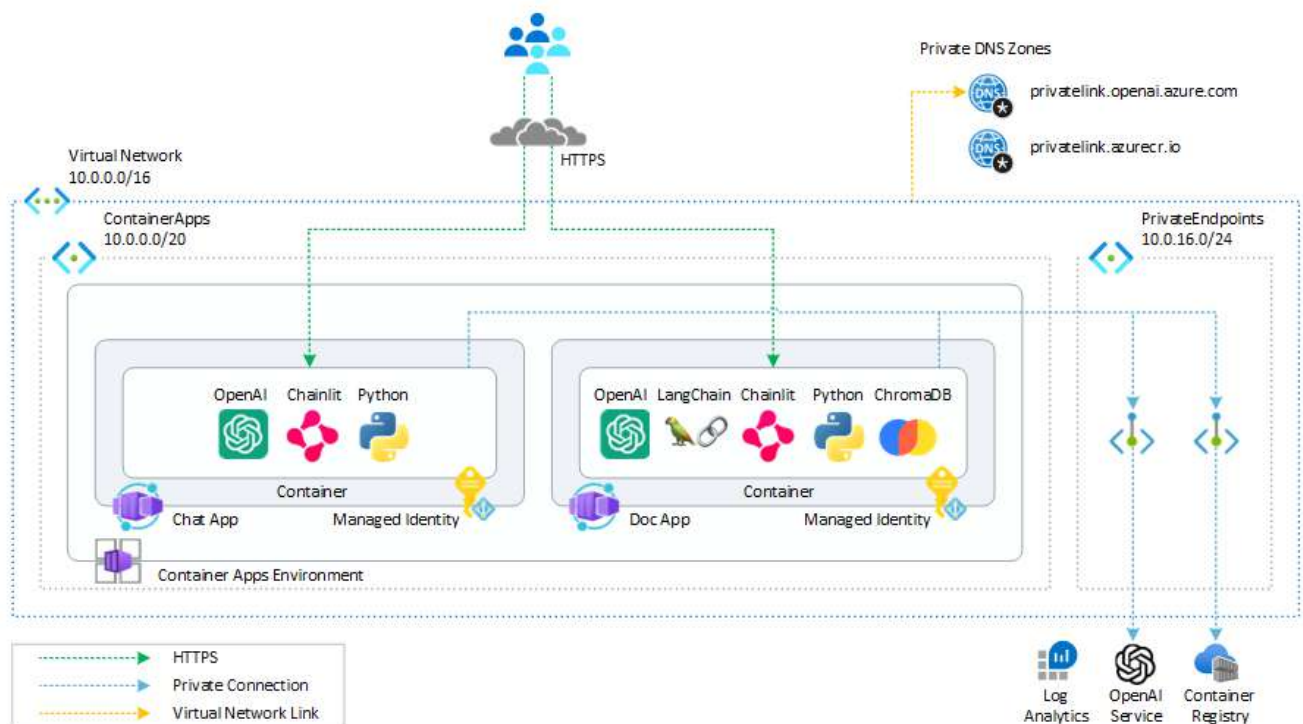
- <u>Deploy and run an Azure OpenAI ChatGPT application on AKS via Bicep</u>
- <u>Deploy and run an Azure OpenAI ChatGPT application on AKS via Terraform</u>

## Prerequisites

- An active <u>Azure subscription</u>. If you don't have one, create a <u>free Azure account</u> before you begin.
- <u>Visual Studio Code</u> installed on one of the <u>supported platforms</u> along with the <u>HashiCorp Terraform</u>.
- Azure CLI version 2.49.0 or later installed. To install or upgrade, see <u>Install Azure CLI</u>.
- `aks-preview` Azure CLI extension of version 0.5.140 or later installed
- <u>Terraform v1.5.2 or later</u>.

## Architecture

The following diagram shows the architecture and network topology of the sample:



This sample provides two sets of Terraform modules to deploy the infrastructure and the chat applications.

## Infrastructure Terraform Modules

You can use the Terraform modules in the `terraform/infra` folder to deploy the infrastructure used by the sample, including the <u>Azure Container Apps Environment</u>, <u>Azure OpenAI Service (AOAI)</u>, and <u>Azure Container Registry (ACR)</u>, but not the <u>Azure Container Apps (ACA)</u>. The Terraform modules in the `terraform/infra` folder

deploy the following resources:

- azurerm_virtual_network: an Azure Virtual Network with two subnets:
    - `ContainerApps` : this subnet hosts the Azure Container Apps Environment.
    - `PrivateEndpoints` : this subnet contains the Azure Private Endpoints to the Azure OpenAI Service (AOAI) and Azure Container Registry (ACR) resources.
- azurerm_container_app_environment: the Azure Container Apps Environment hosting the Azure Container Apps.
- azurerm_cognitive_account: an Azure OpenAI Service (AOAI) with a GPT-3.5 model used by the chatbot applications. Azure OpenAI Service gives customers advanced language AI with OpenAI GPT-4, GPT-3, Codex, and DALL-E models with Azure's security and enterprise promise. Azure OpenAI co-develops the APIs with OpenAI, ensuring compatibility and a smooth transition from one to the other. The Terraform modules create the following models:
    - GPT-35: a `gpt-35-turbo-16k` model is used to generate human-like and engaging conversational responses.
    - Embeddings model: the `text-embedding-ada-002` model is to transform input documents into meaningful and compact numerical representations called embeddings. Embeddings capture the semantic or contextual information of the input data in a lower-dimensional space, making it easier for machine learning algorithms to process and analyze the data effectively. Embeddings can be stored in a vector database, such as ChromaDB or Facebook AI Similarity Search, explicitly designed for efficient storage, indexing, and retrieval of vector embeddings.
- azurerm_user_assigned_identity: a user-defined managed identity used by the chatbot applications to acquire a security token to call the Chat Completion API of the ChatGPT model provided by the Azure OpenAI Service and to call the Embedding model.
- azurerm_container_registry: an Azure Container Registry (ACR) to build, store, and manage container images and artifacts in a private registry for all container deployments. In this sample, the registry stores the container images of the two chat applications.
- azurerm_private_endpoint: an Azure Private Endpoint is created for each of the following resources:
    - Azure OpenAI Service (AOAI)
    - Azure Container Registry (ACR)
- azurerm_private_dns_zone: an Azure Private DNS Zone is created for each of the following resources:
    - Azure OpenAI Service (AOAI)
    - Azure Container Registry (ACR)
- azurerm_log_analytics_workspace: a centralized Azure Log Analytics workspace is used to collect the diagnostics logs and metrics from all the Azure resources:
    - Azure OpenAI Service (AOAI)
    - Azure Container Registry (ACR)
    - Azure Container Apps (ACA)

## Application Terraform Modules

You can use these Terraform modules in the `terraform/apps` To deploy the Azure Container Apps (ACA) using the Docker container images stored in the Azure Container Registry you deployed in the previous step.

- azurerm_container_app: this sample deploys the following applications:
    - **chatapp**: this simple chat application utilizes OpenAI's language models to generate real-time completion responses.

- **docapp**: This chat application goes beyond conversations. Users can upload up to 10 `.pdf` and `.docx` documents, which are then processed to create vector embeddings. These embeddings are stored in ChromaDB for efficient retrieval. Users can pose questions about the uploaded documents and view the Chain of Thought, enabling easy exploration of the reasoning process. The completion message contains links to the text chunks in the files that were used as a source for the response.

## Azure Container Apps

Azure Container Apps (ACA) is a serverless compute service provided by Microsoft Azure that allows developers to easily deploy and manage containerized applications without the need to manage the underlying infrastructure. It provides a simplified and scalable solution for running applications in containers, leveraging the power and flexibility of the Azure ecosystem.

With Azure Container Apps, developers can package their applications into containers using popular containerization technologies such as Docker. These containers encapsulate the application and its dependencies, ensuring consistent execution across different environments.

Powered by Kubernetes and open-source technologies like Dapr, KEDA, and envoy, the service abstracts away the complexities of managing the infrastructure, including provisioning, scaling, and monitoring, allowing developers to focus solely on building and deploying their applications. Azure Container Apps handles automatic scaling, and load balancing, and natively integrates with other Azure services, such as Azure Monitor and Azure Container Registry (ACR), to provide a comprehensive and secure application deployment experience.

Azure Container Apps offers benefits such as rapid deployment, easy scalability, cost-efficiency, and seamless integration with other Azure services, making it an attractive choice for modern application development and deployment scenarios.

## Azure OpenAI Service

The Azure OpenAI Service is a platform offered by Microsoft Azure that provides cognitive services powered by OpenAI models. One of the models available through this service is the ChatGPT model, which is designed for interactive conversational tasks. It allows developers to integrate natural language understanding and generation capabilities into their applications.

Azure OpenAI Service provides REST API access to OpenAI's powerful language models including the GPT-3, Codex and Embeddings model series. In addition, the new GPT-4 and ChatGPT model series have now reached general availability. These models can be easily adapted to your specific task, including but not limited to content generation, summarization, semantic search, and natural language-to-code translation. Users can access the service through REST APIs, Python SDK, or our web-based interface in the Azure OpenAI Studio.

You can use Embeddings model to transform raw data or inputs into meaningful and compact numerical representations called embeddings. Embeddings capture the semantic or contextual information of the input data in a lower-dimensional space, making it easier for machine learning algorithms to process and analyze the data effectively. Embeddings can be stored in a vector database, such as ChromaDB or Facebook AI Similarity Search (FAISS), explicitly designed for efficient storage, indexing, and retrieval of vector embeddings.

The Chat Completion API, which is part of the Azure OpenAI Service, provides a dedicated interface for interacting with the ChatGPT and GPT-4 models. This API is currently in preview and is the preferred method for accessing these models. The GPT-4 models can only be accessed through this API.

GPT-3, GPT-3.5, and GPT-4 models from OpenAI are prompt-based. With prompt-based models, the user interacts with the model by entering a text prompt, to which the model responds with a text completion. This completion is the model's continuation of the input text. While these models are compelling, their behavior is also very sensitive to the prompt. This makes prompt construction a critical skill to develop. For more information, see Introduction to prompt engineering.

Prompt construction can be complex. In practice, the prompt acts to configure the model weights to complete the desired task, but it's more of an art than a science, often requiring experience and intuition to craft a successful prompt. The goal of this article is to help get you started with this learning process. It attempts to capture general concepts and patterns that apply to all GPT models. However, it's essential to understand that each model behaves differently, so the learnings may not apply equally to all models.

Prompt engineering refers to the process of creating instructions called prompts for Large Language Models (LLMs), such as OpenAI's ChatGPT. With the immense potential of LLMs to solve a wide range of tasks, leveraging prompt engineering can empower us to save significant time and facilitate the development of impressive applications. It holds the key to unleashing the full capabilities of these huge models, transforming how we interact and benefit from them. For more information, see Prompt engineering techniques.

## Vector Databases

A vector database is a specialized database that goes beyond traditional storage by organizing information to simplify the search for similar items. Instead of merely storing words or numbers, it leverages vector embeddings - unique numerical representations of data. These embeddings capture meaning, context, and relationships. For instance, words are represented as vectors, whereas similar words have similar vector values.

The applications of vector databases are numerous and powerful. In language processing, they facilitate the discovery of related documents or sentences. By comparing the vector embeddings of different texts, finding similar or related information becomes faster and more efficient. This capability benefits search engines and recommendation systems, which can suggest relevant articles or products based on user interests.

In the realm of image analysis, vector databases excel in finding visually similar images. By representing images as vectors, a simple comparison of vector values can identify visually similar images. This capability is precious for tasks like reverse image search or content-based image retrieval.

Additionally, vector databases find applications in fraud detection, anomaly detection, and clustering. By comparing vector embeddings of data points, unusual patterns can be detected, and similar items can be grouped together, aiding in effective data analysis and decision-making.

Here is a list of the most popular vector databases:

- ChromaDB is a powerful database solution that stores and retrieves vector embeddings efficiently. It is commonly used in AI applications, including chatbots and document analysis systems. By storing embeddings in ChromaDB, users can easily search and retrieve similar vectors, enabling faster and more accurate matching or recommendation processes. ChromaDB offers excellent scalability high performance, and supports various indexing techniques to optimize search operations. It is a versatile tool that enhances the functionality and efficiency of AI applications that rely on vector embeddings.

- Facebook AI Similarity Search (FAISS) is another widely used vector database. Facebook AI Research develops it and offers highly optimized algorithms for similarity search and clustering of vector embeddings. FAISS is known for its speed and scalability, making it suitable for large-scale applications. It offers different indexing methods like flat, IVF (Inverted File System), and HNSW (Hierarchical Navigable Small World) to organize and search vector data efficiently.

- SingleStore: SingleStore aims to deliver the world's fastest distributed SQL database for data-intensive applications: SingleStoreDB, which combines transactional + analytical workloads in a single platform.

- Astra DB: DataStax Astra DB is a cloud-native, multi-cloud, fully managed database-as-a-service based on Apache Cassandra, which aims to accelerate application development and reduce deployment time for applications from weeks to minutes.
- Milvus: Milvus is an open source vector database built to power embedding similarity search and AI applications. Milvus makes unstructured data search more accessible and provides a consistent user experience regardless of the deployment environment. Milvus 2.0 is a cloud-native vector database with storage and computation separated by design. All components in this refactored version of Milvus are stateless to enhance elasticity and flexibility.
- Qdrant: Qdrant is a vector similarity search engine and database for AI applications. Along with open-source, Qdrant is also available in the cloud. It provides a production-ready service with an API to store, search, and manage points—vectors with an additional payload. Qdrant is tailored to extended filtering support. It makes it useful for all sorts of neural network or semantic-based matching, faceted search, and other applications.
- Pinecone: Pinecone is a fully managed vector database that makes adding vector search to production applications accessible. It combines state-of-the-art vector search libraries, advanced features such as filtering, and distributed infrastructure to provide high performance and reliability at any scale.
- Vespa: Vespa is a platform for applications combining data and AI online. Building such applications on Vespa helps users avoid integration work to get features, and it can scale to support any amount of traffic and data. To deliver that, Vespa provides a broad range of query capabilities, a computation engine with support for modern machine-learned models, hands-off operability, data management, and application development support. It is free and open source to use under the Apache 2.0 license.
- Zilliz: Milvus is an open-source vector database, with over 18,409 stars on GitHub and 3.4 million+ downloads. Milvus supports billion-scale vector search and has over 1,000 enterprise users. Zilliz Cloud provides a fully-managed Milvus service made by the creators of Milvus. This helps to simplify the process of deploying and scaling vector search applications by eliminating the need to create and maintain complex data infrastructure. As a DBaaS, Zilliz simplifies the process of deploying and scaling vector search applications by eliminating the need to create and maintain complex data infrastructure.
- Weaviate: Weaviate is an open-source vector database used to store data objects and vector embeddings from ML-models, and scale into billions of data objects from the same name company in Amsterdam. Users can index billions of data objects to search through and combine multiple search techniques, such as keyword-based and vector search, to provide search experiences.

This sample makes of ChromaDB vector database, but you can easily modify the code to use another vector database. You can even use Azure Cache for Redis Enterprise to store the vector embeddings and compute vector similarity with high performance and low latency. For more information, see Vector Similarity Search with Azure Cache for Redis Enterprise

# LangChain

LangChain is a software framework designed to streamline the development of applications using large language models (LLMs). It serves as a language model integration framework, facilitating various applications like document analysis and summarization, chatbots, and code analysis.

LangChain's integrations cover an extensive range of systems, tools, and services, making it a comprehensive solution for language model-based applications. LangChain integrates with the major cloud platforms such as Microsoft Azure, Amazon AWS, and Google, and with API wrappers for various purposes like news, movie information, and weather, as well as support for Bash, web scraping, and more. It also supports multiple

language models, including those from OpenAI, Anthropic, and Hugging Face. Moreover, LangChain offers various functionalities for document handling, code generation, analysis, debugging, and interaction with databases and other data sources.

## Chainlit

Chainlit is an open-source Python package that is specifically designed to create user interfaces (UIs) for AI applications. It simplifies the process of building interactive chats and interfaces, making developing AI-powered applications faster and more efficient. While Streamlit is a general-purpose UI library, Chainlit is purpose-built for AI applications and seamlessly integrates with other AI technologies such as LangChain, LlamaIndex, and LangFlow.

With Chainlit, developers can easily create intuitive UIs for their AI models, including ChatGPT-like applications. It provides a user-friendly interface for users to interact with AI models, enabling conversational experiences and information retrieval. Chainlit also offers unique features, such as displaying the Chain of Thought, which allows users to explore the reasoning process directly within the UI. This feature enhances transparency and enables users to understand how the AI arrives at its responses or recommendations.

For more information, see the following resources:

- Documentation
- Examples
- API Reference
- Cookbook

## Deploy the Infrastructure

Before deploying the Terraform modules in the `terraform/infra` folder, specify a value for the following variables in the terraform.tfvars variable definitions file.

```
1  name_prefix = "Blue"
2  location    = "EastUS"
```

This is the definition of each variable:

- `prefix` : specifies a prefix for all the Azure resources.
- `location` : specifies the region (e.g., EastUS) where deploying the Azure resources.

**NOTE**: Make sure to select a region where Azure OpenAI Service (AOAI) supports both GPT-3.5/GPT-4 models like `gpt-35-turbo-16k` and Embeddings models like `text-embedding-ada-002` .

## OpenAI Module

The following table contains the code from the `terraform/infra/modules/openai/main.tf` Terraform module used to deploy the Azure OpenAI Service.

```
1   resource "azurerm_cognitive_account" "openai" {
2     name                          = var.name
3     location                      = var.location
4     resource_group_name           = var.resource_group_name
5     kind                          = "OpenAI"
6     custom_subdomain_name         = var.custom_subdomain_name
7     sku_name                      = var.sku_name
8     public_network_access_enabled = var.public_network_access_enabled
9     tags                          = var.tags
10
11    identity {
12      type = "SystemAssigned"
13    }
14
15    lifecycle {
16      ignore_changes = [
17        tags
18      ]
19    }
20  }
21
22  resource "azurerm_cognitive_deployment" "deployment" {
23    for_each             = {for deployment in var.deployments: deployment.name => deployment}
24
25    name                 = each.key
26    cognitive_account_id = azurerm_cognitive_account.openai.id
27
28    model {
29      format  = "OpenAI"
30      name    = each.value.model.name
31      version = each.value.model.version
32    }
33
34    scale {
35      type = "Standard"
36    }
37  }
38
39  resource "azurerm_monitor_diagnostic_setting" "settings" {
40    name                       = "DiagnosticsSettings"
41    target_resource_id         = azurerm_cognitive_account.openai.id
42    log_analytics_workspace_id = var.log_analytics_workspace_id
43
44    enabled_log {
45      category = "Audit"
46
47      retention_policy {
48        enabled = true
49        days    = var.log_analytics_retention_days
50      }
51    }
52
53    enabled_log {
54      category = "RequestResponse"
55
56      retention_policy {
57        enabled = true
58        days    = var.log_analytics_retention_days
59      }
60
```

```
61      }
62
63      enabled_log {
64        category = "Trace"
65
66        retention_policy {
67          enabled = true
68          days     = var.log_analytics_retention_days
69        }
70      }
71
72      metric {
73        category = "AllMetrics"
74
75        retention_policy {
76          enabled = true
77          days     = var.log_analytics_retention_days
78        }
79      }
    }
```

Azure Cognitive Services uses custom subdomain names for each resource created through the Azure portal, Azure Cloud Shell, Azure CLI, Bicep, Azure Resource Manager (ARM), or Terraform. Unlike regional endpoints, which were common for all customers in a specific Azure region, custom subdomain names are unique to the resource. Custom subdomain names are required to enable authentication features like Azure Active Directory (Azure AD). We need to specify a custom subdomain for our Azure OpenAI Service, as our chatbot applications will use an Azure AD security token to access it. By default, the `terraform/infra/modules/openai/main.tf` module sets the value of the `custom_subdomain_name` parameter to the lowercase name of the Azure OpenAI resource. For more information on custom subdomains, see Custom subdomain names for Cognitive Services.

This Terraform module allows you to pass an array containing the definition of one or more model deployments in the `deployments` variable. For more information on model deployments, see Create a resource and deploy a model using Azure OpenAI. The `openai_deployments` variable in the `terraform/infra/variables.tf` file defines the structure and the default models deployed by the sample:

```
1    variable "openai_deployments" {
2      description = "(Optional) Specifies the deployments of the Azure OpenAI Service"
3      type = list(object({
4        name = string
5        model = object({
6          name = string
7          version = string
8        })
9        rai_policy_name = string
10     }))
11     default = [
12       {
13         name = "gpt-35-turbo-16k"
14         model = {
15           name = "gpt-35-turbo-16k"
16           version = "0613"
17         }
18         rai_policy_name = ""
19       },
20       {
21         name = "text-embedding-ada-002"
22         model = {
23           name = "text-embedding-ada-002"
24           version = "2"
25         }
26         rai_policy_name = ""
27       }
28     ]
29   }
```

Alternatively, you can use the Terraform module for deploying Azure OpenAI Service. to deploy Azure OpenAI Service.

## Private Endpoint Module

The `terraform/infra/main.tf` the module creates Azure Private Endpoints and Azure Private DNDS Zones for each of the following resources:

- Azure OpenAI Service (AOAI)
- Azure Container Registry (ACR)

In particular, it creates an Azure Private Endpoint and Azure Private DNDS Zone to the Azure OpenAI Service as shown in the following code snippet:

```
1   module "openai_private_dns_zone" {
2     source                   = "./modules/private_dns_zone"
3     name                     = "privatelink.openai.azure.com"
4     resource_group_name      = azurerm_resource_group.rg.name
5     tags                     = var.tags
6     virtual_networks_to_link = {
7       (module.virtual_network.name) = {
8         subscription_id = data.azurerm_client_config.current.subscription_id
9         resource_group_name = azurerm_resource_group.rg.name
10      }
11    }
12  }
13
14  module "openai_private_endpoint" {
15    source                       = "./modules/private_endpoint"
16    name                         = "${module.openai.name}PrivateEndpoint"
17    location                     = var.location
18    resource_group_name          = azurerm_resource_group.rg.name
19    subnet_id                    = module.virtual_network.subnet_ids[var.vm_subnet_name]
20    tags                         = var.tags
21    private_connection_resource_id = module.openai.id
22    is_manual_connection         = false
23    subresource_name             = "account"
24    private_dns_zone_group_name  = "AcrPrivateDnsZoneGroup"
25    private_dns_zone_group_ids   = [module.openai_private_dns_zone.id]
26  }
```

Below you can read the code of the `terraform/infra/modules/private_endpoint/main.tf` module, which is used to create Azure Private Endpoints:

```
1   resource "azurerm_private_endpoint" "private_endpoint" {
2     name                = var.name
3     location            = var.location
4     resource_group_name = var.resource_group_name
5     subnet_id           = var.subnet_id
6     tags                = var.tags
7
8     private_service_connection {
9       name                           = "${var.name}Connection"
10      private_connection_resource_id = var.private_connection_resource_id
11      is_manual_connection           = var.is_manual_connection
12      subresource_names              = try([var.subresource_name], null)
13      request_message                = try(var.request_message, null)
14    }
15
16    private_dns_zone_group {
17      name                 = var.private_dns_zone_group_name
18      private_dns_zone_ids = var.private_dns_zone_group_ids
19    }
20
21    lifecycle {
22      ignore_changes = [
23        tags
24      ]
25    }
26  }
```

## *Private DNS Zone Module*

In the following box, you can read the code of
the `terraform/infra/modules/private_dns_zone/main.tf` module, which is utilized to create the Azure Private
DNS Zones.

```
1   resource "azurerm_private_dns_zone" "private_dns_zone" {
2     name                = var.name
3     resource_group_name = var.resource_group_name
4     tags                = var.tags
5
6     lifecycle {
7       ignore_changes = [
8         tags
9       ]
10    }
11  }
12
13  resource "azurerm_private_dns_zone_virtual_network_link" "link" {
14    for_each = var.virtual_networks_to_link
15
16    name                  = "link_to_${lower(basename(each.key))}"
17    resource_group_name   = var.resource_group_name
18    private_dns_zone_name = azurerm_private_dns_zone.private_dns_zone.name
19    virtual_network_id    = "/subscriptions/${each.value.subscription_id}/resourceGroups/${each.val
20
21    lifecycle {
22      ignore_changes = [
23        tags
24      ]
25    }
26  }
```

## Workload Managed Identity Module

Below you can read the code of the `terraform/infra/modules/managed_identity/main.tf` module, which is used to create the Azure Managed Identity used by the Azure Container Apps to pull container images from the Azure Container Registry, and by the chat applications to connect to the Azure OpenAI Service. You can use a system-assigned or user-assigned managed identity from Azure Active Directory (Azure AD) to let Azure Container Apps access any Azure AD-protected resource. For more information, see Managed identities in Azure Container Apps. You can pull container images from private repositories in an Azure Container Registry using user-assigned or user-assigned managed identities for authentication to avoid using administrative credentials. For more information, see Azure Container Apps image pull with managed identity. This user-defined managed identity is assigned the Cognitive Services User role on the Azure OpenAI Service namespace and ACRPull role on the Azure Container Registry (ACR). By assigning the above roles, you grant the user-defined managed identity access to these resources.

```
1   resource "azurerm_user_assigned_identity" "workload_user_assigned_identity" {
2     name                 = var.name
3     resource_group_name  = var.resource_group_name
4     location             = var.location
5     tags                 = var.tags
6
7     lifecycle {
8       ignore_changes = [
9         tags
10      ]
11    }
12  }
13
14  resource "azurerm_role_assignment" "cognitive_services_user_assignment" {
15    scope                 = var.openai_id
16    role_definition_name  = "Cognitive Services User"
17    principal_id          = azurerm_user_assigned_identity.workload_user_assigned_identity.principal_
18    skip_service_principal_aad_check = true
19  }
20
21  resource "azurerm_role_assignment" "acr_pull_assignment" {
22    scope                 = var.acr_id
23    role_definition_name  = "AcrPull"
24    principal_id          = azurerm_user_assigned_identity.workload_user_assigned_identity.principal_
25    skip_service_principal_aad_check = true
26  }
```

## Deploy the Applications

Before deploying the Terraform modules in the `terraform/apps` folder, specify a value for the following variables in the Terraform.tfvars variable definitions file.

```
1   resource_group_name             = "BlueRG"
2   container_app_environment_name  = "BlueEnvironment"
3   container_registry_name         = "BlueRegistry"
4   workload_managed_identity_name  = "BlueWorkloadIdentity"
5   container_apps                  = [
6     {
7       name                        = "chatapp"
8       revision_mode               = "Single"
9       ingress                     = {
10        allow_insecure_connections  = true
11        external_enabled            = true
12        target_port                 = 8000
13        transport                   = "http"
14        traffic_weight              = {
15          label                     = "default"
16          latest_revision           = true
17          revision_suffix           = "default"
18          percentage                = 100
19        }
20      }
21      template                    = {
22        containers                = [
23          {
24            name                    = "chat"
25            image                   = "chat:v1"
26            cpu                     = 0.5
27            memory                  = "1Gi"
28            env                     = [
29              {
30                name                = "TEMPERATURE"
31                value               = 0.9
32              },
33              {
34                name                = "AZURE_OPENAI_BASE"
35                value               = "https://blueopenai.openai.azure.com/"
36              },
37              {
38                name                = "AZURE_OPENAI_KEY"
39                value               = ""
40              },
41              {
42                name                = "AZURE_OPENAI_TYPE"
43                value               = "azure_ad"
44              },
45              {
46                name                = "AZURE_OPENAI_VERSION"
47                value               = "2023-06-01-preview"
48              },
49              {
50                name                = "AZURE_OPENAI_DEPLOYMENT"
51                value               = "gpt-35-turbo-16k"
52              },
53              {
54                name                = "AZURE_OPENAI_MODEL"
55                value               = "gpt-35-turbo-16k"
56              },
57              {
58                name                = "AZURE_OPENAI_SYSTEM_MESSAGE"
59                value               = "You are a helpful assistant."
60
```

```
        },
        {
          name                    = "MAX_RETRIES"
          value                   = 5
        },
        {
          name                    = "BACKOFF_IN_SECONDS"
          value                   = "1"
        },
        {
          name                    = "TOKEN_REFRESH_INTERVAL"
          value                   = 2700
        }
      ]
      liveness_probe            = {
        failure_count_threshold = 3
        initial_delay           = 30
        interval_seconds        = 60
        path                    = "/"
        port                    = 8000
        timeout                 = 30
        transport               = "HTTP"
      }
      readiness_probe = {
        failure_count_threshold = 3
        interval_seconds        = 60
        path                    = "/"
        port                    = 8000
        success_count_threshold = 3
        timeout                 = 30
        transport               = "HTTP"
      }
      startup_probe = {
        failure_count_threshold = 3
        interval_seconds        = 60
        path                    = "/"
        port                    = 8000
        timeout                 = 30
        transport               = "HTTP"
      }
    }
  ]
  min_replicas              = 1
  max_replicas              = 3
  }
},
{
  name                      = "docapp"
  revision_mode             = "Single"
  ingress                   = {
    allow_insecure_connections = true
    external_enabled        = true
    target_port             = 8000
    transport               = "http"
    traffic_weight          = {
      label                 = "default"
      latest_revision       = true
      revision_suffix       = "default"
      percentage            = 100
    }
```

```
121          }
122        template                      = {
123          containers                  = [
124            {
125              name                    = "doc"
126              image                   = "doc:v1"
127              cpu                     = 0.5
128              memory                  = "1Gi"
129              env                     = [
130                {
131                  name                    = "TEMPERATURE"
132                  value                   = 0.9
133                },
134                {
135                  name                    = "AZURE_OPENAI_BASE"
136                  value                   = "https://blueopenai.openai.azure.com/"
137                },
138                {
139                  name                    = "AZURE_OPENAI_KEY"
140                  value                   = ""
141                },
142                {
143                  name                    = "AZURE_OPENAI_TYPE"
144                  value                   = "azure_ad"
145                },
146                {
147                  name                    = "AZURE_OPENAI_VERSION"
148                  value                   = "2023-06-01-preview"
149                },
150                {
151                  name                    = "AZURE_OPENAI_DEPLOYMENT"
152                  value                   = "gpt-35-turbo-16k"
153                },
154                {
155                  name                    = "AZURE_OPENAI_MODEL"
156                  value                   = "gpt-35-turbo-16k"
157                },
158                {
159                  name                    = "AZURE_OPENAI_ADA_DEPLOYMENT"
160                  value                   = "text-embedding-ada-002"
161                },
162                {
163                  name                    = "AZURE_OPENAI_SYSTEM_MESSAGE"
164                  value                   = "You are a helpful assistant."
165                },
166                {
167                  name                    = "MAX_RETRIES"
168                  value                   = 5
169                },
170                {
171                  name                    = "CHAINLIT_MAX_FILES"
172                  value                   = 10
173                },
174                {
175                  name                    = "TEXT_SPLITTER_CHUNK_SIZE"
176                  value                   = 1000
177                },
178                {
179                  name                    = "TEXT_SPLITTER_CHUNK_OVERLAP"
180                  value                   Skip to form content
181
```

```
182              },
183              {
184                 name                    = "EMBEDDINGS_CHUNK_SIZE"
185                 value                   = 16
186              },
187              {
188                 name                    = "BACKOFF_IN_SECONDS"
189                 value                   = "1"
190              },
191              {
192                 name                    = "CHAINLIT_MAX_SIZE_MB"
193                 value                   = 100
194              },
195              {
196                 name                    = "TOKEN_REFRESH_INTERVAL"
197                 value                   = 2700
198              }
199           ]
200           liveness_probe = {
201              failure_count_threshold = 3
202              initial_delay           = 30
203              interval_seconds        = 60
204              path                    = "/"
205              port                    = 8000
206              timeout                 = 30
207              transport               = "HTTP"
208           }
209           readiness_probe = {
210              failure_count_threshold = 3
211              interval_seconds        = 60
212              path                    = "/"
213              port                    = 8000
214              success_count_threshold = 3
215              timeout                 = 30
216              transport               = "HTTP"
217           }
218           startup_probe = {
219              failure_count_threshold = 3
220              interval_seconds        = 60
221              path                    = "/"
222              port                    = 8000
223              timeout                 = 30
224              transport               = "HTTP"
225           }
226        }
227     ]
228     min_replicas                = 1
229     max_replicas                = 3
       }
     }]
```

This is the definition of each variable:

- `resource_group_name` : specifies the name of the resource group that contains the infrastructure resources: Azure OpenAI Service, Azure Container Registry, Azure Container Apps Environment, Azure Log Analytics, and user-defined managed identity.
- `container_app_environment_name` : the name of the Azure Container Apps Environment in which to deploy the chat applications.

- `container_registry_name` : the name of <u>Azure Container Registry</u> used to hold the container images of the chat applications.
- `workload_managed_identity_name` : the name of the <u>user-defined managed identity</u> used by the chat applications to authenticate with <u>Azure OpenAI Service</u> and <u>Azure Container Registry</u>.
- `container_apps` : the definition of the two chat applications. The application configuration does not specify the following data because the `container_app` module later defines this information:
  - `image` : This field contains the name and tag of the container image but not the login server of the <u>Azure Container Registry</u>.
  - `identity` : The identity of the container app.
  - `registry` : The registry hosting the container image for the application.
  - `AZURE_CLIENT_ID` : The client ID of the user-defined managed identity used by the application to authenticate with <u>Azure OpenAI Service</u> and <u>Azure Container Registry</u>.
  - `AZURE_OPENAI_TYPE` : This environment variable specifies the authentication type with <u>Azure OpenAI Service</u>: if you set the value of the `AZURE_OPENAI_TYPE` environment variable to `azure` , you need to specify the OpenAI key as a value for the `AZURE_OPENAI_KEY` environment variable. Instead, if you set the value to `azure_ad` in the application code, assign an Azure AD security token to the `openai_api_key` property. For more information, see <u>How to switch between OpenAI and Azure OpenAI endpoints with Python</u>.

## Container App Module

The `terraform/apps/modules/container_app/main.tf` module is utilized to create the <u>Azure Container Apps</u>. The module defines and uses the following <u>data source</u> for the <u>Azure Container Registry</u>, <u>Azure Container Apps Environment</u>, and <u>user-defined managed identity</u> created when deploying the infrastructure. These data sources are used to access the properties of these Azure resources.

```
1   data "azurerm_container_app_environment" "container_app_environment" {
2     name                 = var.container_app_environment_name
3     resource_group_name  = var.resource_group_name
4   }
5
6   data "azurerm_container_registry" "container_registry" {
7     name                 = var.container_registry_name
8     resource_group_name  = var.resource_group_name
9   }
10
11  data "azurerm_user_assigned_identity" "workload_user_assigned_identity" {
12    name                = var.workload_managed_identity_name
13    resource_group_name = var.resource_group_name
14  }
```

The module creates and utilizes the following local variables:

```
1   locals {
2     identity = {
3       type        = "UserAssigned"
4       identity_ids = [data.azurerm_user_assigned_identity.workload_user_assigned_identity.id]
5     }
6     identity_env = {
7       name        = "AZURE_CLIENT_ID"
8       secret_name = null
9       value       = data.azurerm_user_assigned_identity.workload_user_assigned_identity.client_id
10    }
11    registry = {
12      server      = data.azurerm_container_registry.container_registry.login_server
13      identity    = data.azurerm_user_assigned_identity.workload_user_assigned_identity.id
14    }
15  }
```

This is the explanation of each local variable:

- `identity` : uses the resource ID of the <u>user-defined managed identity</u> to define the `identity` block for each container app deployed by the module.
- `identity_env` : uses the client ID of the <u>user-defined managed identity</u> to define the value of the `AZURE_CLIENT_ID` environment variable that is appended to the list of environment variables of each container app deployed by the module.
- `registry` : uses the login server of the <u>Azure Container Registry</u> to define the `registry` block for each container app deployed by the module.

Here is the complete Terraform code of the module:

```hcl
data "azurerm_container_app_environment" "container_app_environment" {
  name                = var.container_app_environment_name
  resource_group_name = var.resource_group_name
}

data "azurerm_container_registry" "container_registry" {
  name                = var.container_registry_name
  resource_group_name = var.resource_group_name
}

data "azurerm_user_assigned_identity" "workload_user_assigned_identity" {
  name                = var.workload_managed_identity_name
  resource_group_name = var.resource_group_name
}

locals {
  identity = {
    type         = "UserAssigned"
    identity_ids = [data.azurerm_user_assigned_identity.workload_user_assigned_identity.id]
  }
  identity_env = {
    name        = "AZURE_CLIENT_ID"
    secret_name = null
    value       = data.azurerm_user_assigned_identity.workload_user_assigned_identity.client_id
  }
  registry = {
    server      = data.azurerm_container_registry.container_registry.login_server
    identity    = data.azurerm_user_assigned_identity.workload_user_assigned_identity.id
  }
}

resource "azurerm_container_app" "container_app" {
  for_each                  = {for app in var.container_apps: app.name => app}

  container_app_environment_id = data.azurerm_container_app_environment.container_app_environment
  name                      = each.key
  resource_group_name       = var.resource_group_name
  revision_mode             = each.value.revision_mode
  tags                      = each.value.tags

  template {
    max_replicas    = each.value.template.max_replicas
    min_replicas    = each.value.template.min_replicas
    revision_suffix = each.value.template.revision_suffix

    dynamic "container" {
      for_each = each.value.template.containers

      content {
        cpu     = container.value.cpu
        image   = "${data.azurerm_container_registry.container_registry.login_server}/${container
        memory  = container.value.memory
        name    = container.value.name
        args    = container.value.args
        command = container.value.command

        dynamic "env" {
          for_each = container.value.env == null ? [local.identity_env] : concat(container.value.
```

```
    content {
      name        = env.value.name
      secret_name = env.value.secret_name
      value       = env.value.value
    }
  }

  dynamic "liveness_probe" {
    for_each = container.value.liveness_probe == null ? [] : [container.value.liveness_prob

    content {
      port                    = liveness_probe.value.port
      transport               = liveness_probe.value.transport
      failure_count_threshold = liveness_probe.value.failure_count_threshold
      host                    = liveness_probe.value.host
      initial_delay           = liveness_probe.value.initial_delay
      interval_seconds        = liveness_probe.value.interval_seconds
      path                    = liveness_probe.value.path
      timeout                 = liveness_probe.value.timeout

      dynamic "header" {
        for_each = liveness_probe.value.header == null ? [] : [liveness_probe.value.header]

        content {
          name  = header.value.name
          value = header.value.value
        }
      }
    }
  }

  dynamic "readiness_probe" {
    for_each = container.value.readiness_probe == null ? [] : [container.value.readiness_pr

    content {
      port                    = readiness_probe.value.port
      transport               = readiness_probe.value.transport
      failure_count_threshold = readiness_probe.value.failure_count_threshold
      host                    = readiness_probe.value.host
      interval_seconds        = readiness_probe.value.interval_seconds
      path                    = readiness_probe.value.path
      success_count_threshold = readiness_probe.value.success_count_threshold
      timeout                 = readiness_probe.value.timeout

      dynamic "header" {
        for_each = readiness_probe.value.header == null ? [] : [readiness_probe.value.heade

        content {
          name  = header.value.name
          value = header.value.value
        }
      }
    }
  }

  dynamic "startup_probe" {
    for_each = container.value.startup_probe == null ? [] : [container.value.startup_probe]

    content {
      port                              = startup_probe.value.port
```

```
              transport               = startup_probe.value.transport
              failure_count_threshold = startup_probe.value.failure_count_threshold
              host                    = startup_probe.value.host
              interval_seconds        = startup_probe.value.interval_seconds
              path                    = startup_probe.value.path
              timeout                 = startup_probe.value.timeout

              dynamic "header" {
                for_each = startup_probe.value.header == null ? [] : [startup_probe.value.header]

                content {
                  name  = header.value.name
                  value = header.value.name
                }
              }
            }
          }

          dynamic "volume_mounts" {
            for_each = container.value.volume_mounts == null ? [] : [container.value.volume_mounts]

            content {
              name = volume_mounts.value.name
              path = volume_mounts.value.path
            }
          }
        }
      }

      dynamic "volume" {
        for_each = each.value.template.volume == null ? [] : each.value.template.volume

        content {
          name         = volume.value.name
          storage_name = volume.value.storage_name
          storage_type = volume.value.storage_type
        }
      }
    }
  }

  dynamic "dapr" {
    for_each = each.value.dapr == null ? [] : [each.value.dapr]

    content {
      app_id       = dapr.value.app_id
      app_port     = dapr.value.app_port
      app_protocol = dapr.value.app_protocol
    }
  }

  dynamic "identity" {
    for_each = each.value.identity == null ? [local.identity] : [each.value.identity]

    content {
      type         = identity.value.type
      identity_ids = identity.value.identity_ids
    }
  }

  dynamic "ingress" {
```

```
182        for_each = each.value.ingress == null ? [] : [each.value.ingress]
183
184        content {
185          target_port                = ingress.value.target_port
186          allow_insecure_connections = ingress.value.allow_insecure_connections
187          external_enabled           = ingress.value.external_enabled
188          transport                  = ingress.value.transport
189
190          dynamic "traffic_weight" {
191            for_each = ingress.value.traffic_weight == null ? [] : [ingress.value.traffic_weight]
192
193            content {
194              percentage      = traffic_weight.value.percentage
195              label           = traffic_weight.value.label
196              latest_revision = traffic_weight.value.latest_revision
197              revision_suffix = traffic_weight.value.revision_suffix
198            }
199          }
200        }
201      }
202
203    dynamic "registry" {
204      for_each = each.value.registry == null ? [local.registry] : concat(each.value.registry, [loca
205
206      content {
207        server   = registry.value.server
208        identity = registry.value.identity
209      }
210    }
211
212    dynamic "secret" {
213      for_each = nonsensitive(toset([for pair in lookup(var.container_app_secrets, each.key, []) :
214
215      content {
216        name  = secret.key
217        value = local.container_app_secrets[each.key][secret.key]
218      }
219    }
220  }
221 }
```

As you can notice, the module uses the login server of the Azure Container Registry to create the fully qualified name of the container image of the current container app.

## Managed identities in Azure Container Apps

Each chat application makes use of a DefaultAzureCredential object to acquire a security token from Azure Active Directory and authenticate and authorize with Azure OpenAI Service (AOAI) and Azure Container Registry (ACR) using the credentials of the user-defined managed identity associated with the container app. You can use a managed identity in a running container app to authenticate and authorize with any service that supports Azure AD authentication. With managed identities:

- Container apps and applications connect to resources with the managed identity. You don't need to manage credentials in your container apps.
- You can use role-based access control to grant specific permissions to a managed identity.

- System-assigned identities are automatically created and managed. They are deleted when your container app or container app is deleted.
- You can add and delete user-assigned identities and assign them to multiple resources. They are independent of your container app or the container app's lifecycle.
- You can use managed identity to authenticate with a private Azure Container Registry without a username and password to pull containers for your Container App.
- You can use managed identity to create connections for Dapr-enabled applications via Dapr components

For more information, see Managed identities in Azure Container Apps. The workloads running in a container app can use the Azure Identity client libraries to acquire a security token from the Azure Active Directory. You can choose one of the following approaches inside your code:

- Use `DefaultAzureCredential`, which will attempt to use the `WorkloadIdentityCredential`.
- Create a `ChainedTokenCredential` instance that includes `WorkloadIdentityCredential`.
- Use `WorkloadIdentityCredential` directly.

The following table provides the minimum package version required for each language's client library.

| Language | Library | Minimum Version | Example |
|----------|---------|-----------------|---------|
| .NET | Azure.Identity | 1.9.0 | Link |
| Go | azidentity | 1.3.0 | Link |
| Java | azure-identity | 1.9.0 | Link |
| JavaScript | @azure/identity | 3.2.0 | Link |
| Python | azure-identity | 1.13.0 | Link |

**NOTE**: When using Azure Identity client library with Azure Container Apps, the client ID of the managed identity must be specified. When using the `DefaultAzureCredential`, you can explicitly specify the client ID of the container app managed identity in the `AZURE_CLIENT_ID` environment variable.

## Simple Chat Application
The Simple Chat Application is a large language model-based chatbot that allows users to submit general-purpose questions to a GPT model, generating and streaming back human-like and engaging conversational responses. The following picture shows the welcome screen of the chat application.

## Funny Chat 🤖

Hey there, curious minds and merry souls! 🌟 Looking for a burst of laughter and a dash of wit? You've landed in the right spot! Welcome to our funny chat, where the doors to imagination and humor are wide open! 🚪

### Useful Links 🔗

This chat is built using Chainlit:

- **Documentation:** Get started with Chainlit Documentation 📘
- **Discord Community:** Join Chainlit Discord to ask questions, share projects, and connect with other developers! 💬

> Type your message here...

Built with Chainlit

You can modify the welcome screen in markdown by editing the `chainlit.md` file at the project's root. If you do not want a welcome screen, leave the file empty. The following picture shows what happens when a user submits a new message in the chat.

Chainlit can render messages in markdown format and provides classes to support the following elements:

- Audio: The `Audio` class allows you to display an audio player for a specific audio file in the chatbot user interface. You must provide either a URL or a path or content bytes.
- Avatar: The `Avatar` class allows you to display an avatar image next to a message instead of the author's name. You need to send the element once. Next,, if an avatar's name matches an author's name, the avatar will be automatically displayed. You must provide either a URL or a path or content bytes.
- File: The `File` class allows you to display a button that lets users download the file's content. You must provide either a URL or a path or content bytes.
- Image: The `Image` class is designed to create and handle image elements to be sent and displayed in the chatbot user interface. You must provide either a URL or a path or content bytes.
- Pdf: The `Pdf` class allows you to display a PDF hosted remotely or locally in the chatbot UI. This class either takes a URL of a PDF hosted online or the path of a local PDF.
- Pyplot: The `Pyplot` class allows you to display a Matplotlib pyplot chart in the chatbot UI. This class takes a pyplot figure.
- TaskList: The `TaskList` class allows you to display a task list next to the chatbot UI.
- Text: The `Text` class allows you to display a text element in the chatbot UI. This class takes a string and creates a text element that can be sent to the UI. It supports the markdown syntax for formatting text. You must provide either a URL or a path or content bytes.

Chainlit provides three <u>display options</u> that determine how an element is rendered in the context of its use. The ElementDisplay type represents these options. The following display options are available:

- `Side` : this option displays the element on a sidebar. The sidebar is hidden by default and opened upon element reference click.
- `Page` : this option displays the element on a separate page. The user is redirected to the page upon an element reference click.
- `Inline` : this option displays the element below the message. If the element is <u>global</u>, it is displayed if it is explicitly mentioned in the message. If the element is <u>scoped</u>, it is displayed regardless of whether it is expressly mentioned in the message.

You can click the user icon on the UI to access the chat settings and choose, for example, between the light and dark themes.

The application is built in Python. Let's take a look at the individual parts of the application code. The Python code starts by importing the necessary packages/modules in the following section.

```
1   # Import packages
2   import os
3   import sys
4   import time
5   import openai
6   import random
7   import logging
8   import chainlit as cl
9   from azure.identity import DefaultAzureCredential
10  from dotenv import load_dotenv
11  from dotenv import dotenv_values
12
13  # Load environment variables from .env file
14  if os.path.exists(".env"):
15      load_dotenv(override=True)
16      config = dotenv_values(".env")
```

These are the libraries used by the chat application:

1. `os` : This module provides a way of interacting with the operating system, enabling the code to access environment variables, file paths, etc.
2. `sys` : This module provides access to some variables used or maintained by the interpreter and functions that interact with the interpreter.
3. `time` : This module provides various time-related time manipulation and measurement functions.
4. `openai` : The OpenAI Python library provides convenient access to the OpenAI API from applications written in Python. It includes a pre-defined set of classes for API resources that initialize themselves dynamically from API responses, making it compatible with a wide range of versions of the OpenAI API. You can find usage examples for the OpenAI Python library in our <u>API reference</u> and the <u>OpenAI Cookbook</u>.
5. `random` : This module provides functions to generate random numbers.
6. `logging` : This module provides flexible logging of messages.
7. `chainlit as cl` : This imports the <u>Chainlit</u> library and aliases it as `cl` . Chainlit is used to create the UI of the application.

8. `DefaultAzureCredential` from `azure.identity`: when the `openai_type` property value is `azure_ad,` a `DefaultAzureCredential` object from the [Azure Identity client library for Python - version 1.13.0(https://learn.microsoft.com/en-us/python/api/overview/azure/identity-readme?view=azure-python) is used to acquire security token from the Azure Active Directory using the credentials of the user-defined managed identity, whose client ID is defined in the `AZURE_CLIENT_ID` environment variable.

9. `load_dotenv` and `dotenv_values` from `dotenv`: Python-dotenv reads key-value pairs from a `.env` file and can set them as environment variables. It helps in the development of applications following the 12-factor principles.

The `requirements.txt` file under the `src` folder contains the list of packages used by the chat applications. You can restore these packages in your environment using the following command:

```
1   pip install -r requirements.txt --upgrade
```

Next, the code reads environment variables and configures the OpenAI settings.

```
1    # Read environment variables
2    temperature = float(os.environ.get("TEMPERATURE", 0.9))
3    api_base = os.getenv("AZURE_OPENAI_BASE")
4    api_key = os.getenv("AZURE_OPENAI_KEY")
5    api_type = os.environ.get("AZURE_OPENAI_TYPE", "azure")
6    api_version = os.environ.get("AZURE_OPENAI_VERSION", "2023-06-01-preview")
7    engine = os.getenv("AZURE_OPENAI_DEPLOYMENT")
8    model = os.getenv("AZURE_OPENAI_MODEL")
9    system_content = os.getenv("AZURE_OPENAI_SYSTEM_MESSAGE", "You are a helpful assistant.")
10   max_retries = int(os.getenv("MAX_RETRIES", 5))
11   backoff_in_seconds = float(os.getenv("BACKOFF_IN_SECONDS", 1))
12
13   # Configure OpenAI
14   openai.api_type = api_type
15   openai.api_version = api_version
16   openai.api_base = api_base
17   openai.api_key = api_key
```

Here's a brief explanation of each variable and related environment variable:

1. `temperature`: A float value representing the temperature for Create chat completion method of the OpenAI API. It is fetched from the environment variables with a default value of 0.9.

2. `api_base`: The base URL for the OpenAI API.

3. `api_key`: The API key for the OpenAI API.

4. `api_type`: A string representing the type of the OpenAI API.

5. `api_version`: A string representing the version of the OpenAI API.

6. `engine`: The engine used for OpenAI API calls.

7. `model`: The model used for OpenAI API calls.

8. `system_content`: The content of the system message used for OpenAI API calls.

9. `max_retries`: The maximum number of retries for OpenAI API calls.

10. `backoff_in_seconds`: The backoff time in seconds for retries in case of failures.

In the next section, the code sets the default Azure credential based on the `api_type` and configures a logger for logging purposes.

```
1   # Set default Azure credential
2   default_credential = DefaultAzureCredential() if openai.api_type == "azure_ad" else None
3
4   # Configure a logger
5   logging.basicConfig(
6       stream=sys.stdout,
7       format='[%(asctime)s] {%(filename)s:%(lineno)d} %(levelname)s - %(message)s',
8       level=logging.INFO
9   )
10  logger = logging.getLogger(__name__)
```

Here's a brief explanation:

1. `default_credential` : It sets the default Azure credential to `DefaultAzureCredential()` if the `api_type` is "azure_ad"; otherwise, it is set to `None` .
2. `logging.basicConfig()` : This function configures the logging system with specific settings.
    ◦ `stream` : The output stream where log messages will be written. Here, it is set to `sys.stdout` for writing log messages to the standard output.
    ◦ `format` : The format string for log messages. It includes the timestamp, filename, line number, log level, and the actual log message.
    ◦ `level` : The logging level. It is set to `logging.INFO` , meaning only messages with the level `INFO` and above will be logged.
3. `logger` : This creates a logger instance named after the current module ( `__name__` ). The logger will be used to log messages throughout the code.

Next, the code defines a helper function `backoff` that takes an integer `attempt` and returns a float value representing the backoff time for exponential retries in case of API call failures.

```
1   def backoff(attempt: int) -> float:
2       return backoff_in_seconds * 2 ** attempt + random.uniform(0, 1)
```

The backoff time is calculated using the `backoff_in_seconds` and `attempt` variables. It follows the formula `backoff_in_seconds * 2 ** attempt + random.uniform(0, 1)` . This formula increases the backoff time exponentially with each attempt and adds a random value between 0 and 1 to avoid synchronized retries. Then, the application defines a function called `refresh_openai_token()` to refresh the OpenAI security token if needed.

```
1   def refresh_openai_token():
2       token = cl.user_session.get('openai_token')
3       if token is None or token.expires_on < int(time.time()) - 1800:
4           cl.user_session.set('openai_token', default_credential.get_token("https://cognitiveservic
5           openai.api_key = cl.user_session.get('openai_token').token
```

The function follows these steps:

1. It fetches the current token from `cl.user_session` (which seems to be a part of the `chainlit` library) using the key `'openai_token'` . The user session is a dictionary that stores the user's session data. The id and env keys are reserved for the session ID and environment variables. Other keys can be used to store arbitrary data in the user's session.

2. It checks if the token is `None` or if its expiration time (`expires_on`) is less than the current time minus 1800 seconds (30 minutes).

Next, the code defines a function called `start_chat` that is used to initialize the when the user connects to the application or clicks the `New Chat` button.

```
1   .on_chat_start
2   async def start_chat():
3       # Sending Avatars for Chat Participants
4       await cl.Avatar(
5           name="Chatbot",
6           url="https://cdn-icons-png.flaticon.com/512/8649/8649595.png"
7       ).send()
8       await cl.Avatar(
9           name="Error",
10          url="https://cdn-icons-png.flaticon.com/512/8649/8649595.png"
11      ).send()
12      await cl.Avatar(
13          name="User",
14          url="https://media.architecturaldigest.com/photos/5f241de2c850b2a36b415024/master/w_1600%
15      ).send()
16
17      # Initializing message_history in user session
18      system_content = "Welcome to the chat!"
19      cl.user_session.set("message_history", [{"role": "system", "content": system_content}])
```

Here is a brief explanation of the function steps:

- `@cl.on_chat_start` : The on chat start decorator registers a callback function `start_chat()` to be called when the Chainlit chat starts. It is used to set up the chat and send avatars for the Chatbot, Error, and User participants in the chat.
- `cl.Avatar()` : the Avatar class allows you to display an avatar image next to a message instead of the author name. You need to send the element once. Next, if an avatar's name matches an author's name, the avatar will be automatically displayed. You must provide either a URL or a path or content bytes.
- `cl.user_session.set()` : This API call sets a value in the user session dictionary. In this case, it initializes the `message_history` in the user's session with a system content message, indicating the chat's start.

Finally, the application defines the method called whenever the user sends a new message in the chat.

```
1    .on_message
2    async def main(message: str):
3        # Fetching message history from user session
4        message_history = cl.user_session.get("message_history")
5
6        # Appending user's message to message history
7        message_history.append({"role": "user", "content": message})
8
9        # Creating an empty Chainlit response message
10       msg = cl.Message(content="")
11
12       # Retry the OpenAI API call if it fails
13       for attempt in range(max_retries):
14           try:
15               # Refresh the OpenAI security token if using Azure AD
16               if openai.api_type == "azure_ad":
17                   refresh_openai_token()
18
19               # Sending the message to OpenAI and streaming the response
20               async for stream_resp in await openai.ChatCompletion.acreate(
21                   engine=engine,
22                   model=model,
23                   messages=message_history,
24                   temperature=temperature,
25                   stream=True
26               ):
27                   if stream_resp and len(stream_resp.choices) > 0:
28                       token = stream_resp.choices[0]["delta"].get("content", "")
29                       await msg.stream_token(token)
30               break
31           # Exception handling for different types of errors during the API call (Timeout, APIError
```

Here is a detailed explanation of the function steps:

- `@cl.on_message` : The on message decorator registers a callback function `main(message: str)` to be called when the user submits a new message in the chat. It is the main function responsible for handling the chat logic.

- `cl.user_session.get()` : This API call retrieves a value from the user's session data stored in the user session dictionary. In this case, it fetches the `message_history` from the user's session to maintain the chat history.

- `message_history.append()` : This API call appends a new message to the `message_history` list. It is used to add the user's message and the assistant's response to the chat history.

- `cl.Message()` : This API call creates a Chainlit Message object. The `Message` class is designed to send, stream, edit, or remove messages in the chatbot user interface. In this sample, the `Message` object is used to stream the OpenAI response in the chat.

- `msg.stream_token()` : The stream token method of the Message class streams a token to the response message. It is used to send the response from the OpenAI Chat API in chunks to ensure real-time streaming in the chat.

- `await openai.ChatCompletion.acreate()` : This API call sends a message to the OpenAI Chat API in an asynchronous mode and streams the response. It uses the provided `message_history` as context for generating the assistant's response.

- The section also includes an exception handling block that retries the OpenAI API call in case of specific errors like timeouts, API errors, connection errors, invalid requests, service unavailability, and other non-

retriable errors. You can replace this code with a general-purpose retrying library for Python like <u>Tenacity</u>.

Below, you can read the complete code of the application.

```python
# Import packages
import os
import sys
import time
import openai
import random
import logging
import chainlit as cl
from azure.identity import DefaultAzureCredential
from dotenv import load_dotenv
from dotenv import dotenv_values

# Load environment variables from .env file
if os.path.exists(".env"):
    load_dotenv(override = True)
    config = dotenv_values(".env")


# Read environment variables
temperature = float(os.environ.get("TEMPERATURE", 0.9))
api_base = os.getenv("AZURE_OPENAI_BASE")
api_key = os.getenv("AZURE_OPENAI_KEY")
api_type = os.environ.get("AZURE_OPENAI_TYPE", "azure")
api_version = os.environ.get("AZURE_OPENAI_VERSION", "2023-06-01-preview")
engine = os.getenv("AZURE_OPENAI_DEPLOYMENT")
model = os.getenv("AZURE_OPENAI_MODEL")
system_content = os.getenv("AZURE_OPENAI_SYSTEM_MESSAGE", "You are a helpful assistant.")
max_retries = int(os.getenv("MAX_RETRIES", 5))
backoff_in_seconds = float(os.getenv("BACKOFF_IN_SECONDS", 1))
token_refresh_interval = int(os.getenv("TOKEN_REFRESH_INTERVAL", 1800))


# Configure OpenAI
openai.api_type = api_type
openai.api_version = api_version
openai.api_base = api_base
openai.api_key = api_key


# Set default Azure credential
default_credential = DefaultAzureCredential(
) if openai.api_type ==  "azure_ad" else None


# Configure a logger
logging.basicConfig(stream = sys.stdout,
                    format = '[%(asctime)s] {%(filename)s:%(lineno)d} %(levelname)s - %(message)s
                    level = logging.INFO)
logger = logging.getLogger(__name__)

def backoff(attempt : int) -> float:
    return backoff_in_seconds * 2**attempt + random.uniform(0, 1)


# Refresh the OpenAI security token every 45 minutes
def refresh_openai_token():
    token = cl.user_session.get('openai_token')
    if token ==  None or token.expires_on < int(time.time()) - token_refresh_interval:
        cl.user_session.set('openai_token', default_credential.get_token(
            "https://cognitiveservices.azure.com/.default"))
        openai.api_key = cl.user_session.get('openai_token').token


@cl.on_chat_start
async def start_chat():
```

```python
    await cl.Avatar(
        name = "Chatbot",
        url = "https://cdn-icons-png.flaticon.com/512/8649/8649595.png"
    ).send()
    await cl.Avatar(
        name = "Error",
        url = "https://cdn-icons-png.flaticon.com/512/8649/8649595.png"
    ).send()
    await cl.Avatar(
        name = "User",
        url = "https://media.architecturaldigest.com/photos/5f241de2c850b2a36b415024/master/w_160
    ).send()
    cl.user_session.set(
        "message_history",
        [{"role": "system", "content": system_content}],
    )


@cl.on_message
async def main(message: str):
    message_history = cl.user_session.get("message_history")
    message_history.append({"role": "user", "content": message})

    # Create the Chainlit response message
    msg = cl.Message(content = "")

    # Retry the OpenAI API call if it fails
    for attempt in range(max_retries):
        try:
            # Refresh the OpenAI security token if using Azure AD
            if openai.api_type ==  "azure_ad":
                refresh_openai_token()

            # Send the message to OpenAI in an asynchronous mode and stream the response
            async for stream_resp in await openai.ChatCompletion.acreate(
                engine = engine,
                model = model,
                messages = message_history,
                temperature = temperature,
                stream = True
            ):
                if stream_resp and len(stream_resp.choices) > 0:
                    token = stream_resp.choices[0]["delta"].get("content", "")
                    await msg.stream_token(token)
            break
        except openai.error.Timeout:
            # Implement exponential backoff
            wait_time = backoff(attempt)
            logger.exception(f"OpenAI API timeout occurred. Waiting {wait_time} seconds and tryin
            time.sleep(wait_time)
        except openai.error.APIError:
            # Implement exponential backoff
            wait_time = backoff(attempt)
            logger.exception(f"OpenAI API error occurred. Waiting {wait_time} seconds and trying
            time.sleep(wait_time)
        except openai.error.APIConnectionError:
            # Implement exponential backoff
            wait_time = backoff(attempt)
            logger.exception(f"OpenAI API connection error occurred. Check your network settings,
            time.sleep(wait_time)
        except openai.error.Inval
```

```
121        # Implement exponential backoff
122        wait_time = backoff(attempt)
123        logger.exception(f"OpenAI API invalid request. Check the documentation for the specif
124        time.sleep(wait_time)
125    except openai.error.ServiceUnavailableError:
126        # Implement exponential backoff
127        wait_time = backoff(attempt)
128        logger.exception(f"OpenAI API service unavailable. Waiting {wait_time} seconds and tr
129        time.sleep(wait_time)
130    except Exception as e:
131        logger.exception(f"A non retriable error occurred. {e}")
132        break
133
134
    message_history.append({"role": "assistant", "content": msg.content})
    await msg.send()
```

You can run the application locally using the following command. The `-w` flag` indicates auto-reload whenever we make changes live in our application code.

```
1    chainlit run app.py -w
```

## Documents QA Chat

The Documents QA Chat application allows users to submit up to 10 `.pdf` and `.docx` documents. The application processes the uploaded documents to create vector embeddings. These embeddings are stored in ChromaDB vector database for efficient retrieval. Users can pose questions about the uploaded documents and view the Chain of Thought, enabling easy exploration of the reasoning process. The completion message contains links to the text chunks in the files that were used as a source for the response. The following picture shows the chat application interface. As you can see, you can click the `Browse` button and choose up to 10 `.pdf` and `.docx` documents to upload. Alternatively, you can drag and drop the files over the control area.

Please upload up to 10 `.pdf` or `.docx` files to begin.

11:33:36 AM

Drag and drop files here
Limit 100mb.

Browse Files

Type your message here...

Built with Chainlit

After uploading the documents, the application creates and stores embeddings to [ChromaDB](#) vector database. During the phase, the UI shows a message `Processing <file-1>, <file-2>...`, as shown in the following picture:

Please upload up to 10 `.pdf` or `.docx` files to begin.

Please upload up to 10 `.pdf` or `.docx` files to begin.

11:37:22 AM

11:37:26 AM    Processing `Building_Microservices.pdf` ...

X Stop task

Type your message here...

Built with Chainlit

When the code finished creating embeddings, the UI is ready to receive user's questions:

Understanding the individual steps for generating a specific answer can become challenging as your chat application grows in complexity. To solve this issue, Chainlit allows you to easily explore the reasoning process right from the user interface using the Chain of Thought. If you are using the LangChain integration, every intermediary step is automatically sent and displayed in the Chainlit UI just clicking and expanding the steps, as shown in the following picture:

To see the text chunks that the large language model used to originate the response, you can click the sources links, as shown in the following picture:

In the Chain of Thought, below each message, you can find an `edit` button, as a pencil icon, if a prompt generated that message. Clicking on it opens the Prompt Playground dialog, allowing you to modify and iterate on the prompt as needed.

**Prompt playground** ⑦                                          ✕

> What Are Microservices? | 3Key Benefits
> The benefits of microservices are many and varied. Many of these benefits can be laid
> a
> t the door of any distributed system. Microservices, however, tend to achieve these
> benefits to a greater degree primarily due to how far they take the concepts behind
> distributed systems and service-oriented architecture.
> Technology Heterogeneity
> With a system composed of multiple, collaborating services, we can decide to use dif-
> ferent technologies inside each one. This allows us to pick the right tool for each job,
> rather than having to select a more standardized, one-size-fits-all approach that often
> ends up being the lowest common denominator.
> If one part of our system needs to improve its performance, we might decide to use a
> different technology stack that is better able to achieve the performance levels
> Source: 69-pl
>
> Content: We'll begin with an introduction to microservices, including the key benefits as
> well as some of the downsides.
> xiv | PrefaceChapter 2, The Evolutionary Architect
> This chapter discusses the difficulties we face in terms of making trade-offs as
> architects, and covers specifically just how many things we need to think about
> with microservices.
> Chapter 3, How to Model Services
> Here we'll start to define the boundary of microservices, using techniques from
> domain-driven design to help focus our thinking.
> Chapter 4, Integration
> This is where we start getting a bit deeper into specific technology implications,
> as we discuss what sorts of service collaboration techniques will help us most.
> We'll also delve into the topic of user interfaces and integrating with legacy and
> commercial off-the-shelf (COTS) products.
> Chapter 5, Splitting the Monolith
> Many people get interested in microservices as an antidote to large, hard-to-
> Source: 54-pl
> What is a microservice? What are the main characteristics of a microservice?
> A microservice is a small, autonomous service that works together with other microservices. The main characteristics of a
> microservice are that it is small and focused on doing one thing well. It is also autonomous, meaning it can independently
> operate and make decisions. Microservices also allow for technology heterogeneity, meaning different technologies can be
> used within each microservice. Additionally, microservices promote resilience and scaling, and they are easy to deploy.
> SOURCES: 7-pl, 63-pl, 69-pl

**Model**
text-davinci-003 ▾

Temperature                    0

Stop sequences
⊗                          ✕

Top P                          1

Frequency penalty              0

Presence penalty               0

[Submit]  ↺

Let's take a look at the individual parts of the application code. The Python code starts by importing the necessary packages/modules in the following section.

```python
1    # Import packages
2    import os
3    import io
4    import sys
5    import time
6    import openai
7    import random
8    import logging
9    import chainlit as cl
10   from pypdf import PdfReader
11   from docx import Document
12   from azure.identity import DefaultAzureCredential
13   from dotenv import load_dotenv
14   from dotenv import dotenv_values
15   from langchain.embeddings.openai import OpenAIEmbeddings
16   from langchain.text_splitter import RecursiveCharacterTextSplitter
17   from langchain.vectorstores import Chroma
18   from langchain.chains import RetrievalQAWithSourcesChain
19   from langchain.chat_models import AzureChatOpenAI
20   from langchain.prompts.chat import (
21       ChatPromptTemplate,
22       SystemMessagePromptTemplate,
23       HumanMessagePromptTemplate,
24   )
25
26   # These three lines swap the stdlib sqlite3 lib with the pysqlite3 package
27   __import__('pysqlite3')
28   sys.modules['sqlite3'] = sys.modules.pop('pysqlite3')
29
30   # Load environment variables from .env file
31   if os.path.exists(".env"):
32       load_dotenv(override=True)
33       config = dotenv_values(".env")
```

These are the libraries used by the chat application:

1. `os` : This module provides a way of interacting with the operating system, enabling the code to access environment variables, file paths, etc.

2. `sys` : This module provides access to some variables used or maintained by the interpreter and functions that interact with the interpreter.

3. `time` : This module provides various time-related functions for time manipulation and measurement.

4. `openai` : The OpenAI Python library provides convenient access to the OpenAI API from applications written in Python. It includes a pre-defined set of classes for API resources that initialize themselves dynamically from API responses, which makes it compatible with a wide range of versions of the OpenAI API. You can find usage examples for the OpenAI Python library in our API reference and the OpenAI Cookbook.

5. `random` : This module provides functions to generate random numbers.

6. `logging` : This module provides flexible logging of messages.

7. `chainlit as cl` : This imports the Chainlit library and aliases it as `cl`. Chainlit is used to create the UI of the application.

8. `DefaultAzureCredential` from `azure.identity` : when the `openai_type` property value is `azure_ad`, a `DefaultAzureCredential` object from the Azure Identity client library for Python - version 1.13.0 is

used to acquire security token from the Azure Active Directory using the credentials of the user-defined managed identity, whose client ID is defined in the `AZURE_CLIENT_ID` environment variable.

9. `load_dotenv` and `dotenv_values` from `dotenv` : Python-dotenv reads key-value pairs from a `.env` file and can set them as environment variables. It helps in the development of applications following the 12-factor principles.

10. `langchain` : Large language models (LLMs) are emerging as a transformative technology, enabling developers to build applications that they previously could not. However, using these LLMs in isolation is often insufficient for creating a truly powerful app - the real power comes when you can combine them with other sources of computation or knowledge. LangChain library aims to assist in the development of those types of applications.

The `requirements.txt` file under the `src` folder contains the list of packages used by the chat applications. You can restore these packages in your environment using the following command:

```
pip install -r requirements.txt --upgrade
```

Next, the code reads environment variables and configures the OpenAI settings.

```python
# Read environment variables
temperature = float(os.environ.get("TEMPERATURE", 0.9))
api_base = os.getenv("AZURE_OPENAI_BASE")
api_key = os.getenv("AZURE_OPENAI_KEY")
api_type = os.environ.get("AZURE_OPENAI_TYPE", "azure")
api_version = os.environ.get("AZURE_OPENAI_VERSION", "2023-06-01-preview")
chat_completion_deployment = os.getenv("AZURE_OPENAI_DEPLOYMENT")
embeddings_deployment = os.getenv("AZURE_OPENAI_ADA_DEPLOYMENT")
model = os.getenv("AZURE_OPENAI_MODEL")
max_size_mb = int(os.getenv("CHAINLIT_MAX_SIZE_MB", 100))
max_files = int(os.getenv("CHAINLIT_MAX_FILES", 10))
text_splitter_chunk_size = int(os.getenv("TEXT_SPLITTER_CHUNK_SIZE", 1000))
text_splitter_chunk_overlap = int(os.getenv("TEXT_SPLITTER_CHUNK_OVERLAP", 10))
embeddings_chunk_size = int(os.getenv("EMBEDDINGS_CHUNK_SIZE", 16))
max_retries = int(os.getenv("MAX_RETRIES", 5))
backoff_in_seconds = float(os.getenv("BACKOFF_IN_SECONDS", 1))
token_refresh_interval = int(os.getenv("TOKEN_REFRESH_INTERVAL", 1800))

# Configure system prompt
system_template = """Use the following pieces of context to answer the users question.
If you don't know the answer, just say that you don't know, don't try to make up an answer.
ALWAYS return a "SOURCES" part in your answer.
The "SOURCES" part should be a reference to the source of the document from which you got your an

Example of your response should be:

---

The answer is foo
SOURCES: xyz

---

Begin!
----------------
{summaries}"""
messages = [
    SystemMessagePromptTemplate.from_template(system_template),
    HumanMessagePromptTemplate.from_template("{question}"),
]
prompt = ChatPromptTemplate.from_messages(messages)
chain_type_kwargs = {"prompt": prompt}

# Configure OpenAI
openai.api_type = api_type
openai.api_version = api_version
openai.api_base = api_base
openai.api_key = api_key
```

Here's a brief explanation of each variable and related environment variable:

1. `temperature` : A float value representing the temperature for Create chat completion method of the OpenAI API. It is fetched from the environment variables with a default value of 0.9.

2. `api_base` : The base URL for the OpenAI API.

3. `api_key` : The API key for the OpenAI API.

4. `api_type` : A string representing the type of the OpenAI API.

5. `api_version` : A string representing the version of the OpenAI API.

6. `chat_completion_deployment` : the name of the Azure OpenAI GPT model for chat completion.
7. `embeddings_deployment` : the name of the Azure OpenAI deployment for embeddings.
8. `model` : The model used for chat completion calls (e.g, `gpt-35-turbo-16k` ).
9. `max_size_mb` : the maximum size for the uploaded documents.
10. `max_files` : the maximum number of documents that can be uploaded.
11. `text_splitter_chunk_size` : the maximum chunk size used by
    the `RecursiveCharacterTextSplitter` object.
12. `text_splitter_chunk_overlap` : the maximum chunk overlap used by
    the `RecursiveCharacterTextSplitter` object.
13. `embeddings_chunk_size` : the maximum chunk size used by the `OpenAIEmbeddings` object.
14. `max_retries` : The maximum number of retries for OpenAI API calls.
15. `backoff_in_seconds` : The backoff time in seconds for retries in case of failures.
16. `system_template` : The content of the system message used for OpenAI API calls.

In the next section, the code sets the default Azure credential based on the `api_type` and configures a logger for logging purposes.

```
1   # Set default Azure credential
2   default_credential = DefaultAzureCredential() if openai.api_type == "azure_ad" else None
3
4   # Configure a logger
5   logging.basicConfig(
6       stream=sys.stdout,
7       format='[%(asctime)s] {%(filename)s:%(lineno)d} %(levelname)s - %(message)s',
8       level=logging.INFO
9   )
10  logger = logging.getLogger(__name__)
```

Here's a brief explanation:

1. `default_credential` : It sets the default Azure credential to `DefaultAzureCredential()` if
   the `api_type` is "azure_ad"; otherwise, it is set to `None` .
2. `logging.basicConfig()` : This function configures the logging system with specific settings.
   ○ `stream` : The output stream where log messages will be written. Here, it is set to `sys.stdout` for
     writing log messages to the standard output.
   ○ `format` : The format string for log messages. It includes the timestamp, filename, line number, log
     level, and the actual log message.
   ○ `level` : The logging level. It is set to `logging.INFO` , meaning only messages with the
     level `INFO` and above will be logged.
3. `logger` : This creates a logger instance named after the current module ( `__name__` ). The logger will be
   used to log messages throughout the code.

Next, the code defines a helper function `backoff` that takes an integer `attempt` and returns a float value representing the backoff time for exponential retries in case of API call failures.

```
1   def backoff(attempt: int) -> float:
2       return backoff_in_seconds * 2 ** attempt + random.uniform(0, 1)
```

The backoff time is calculated using the `backoff_in_seconds` and `attempt` variables. It follows the formula `backoff_in_seconds * 2 ** attempt + random.uniform(0, 1)` . This formula increases the backoff time exponentially with each attempt and adds a random number between 0 and 1 to avoid synchronized retries.

Then, the application defines a function called `refresh_openai_token()` to refresh the OpenAI security token if needed.

```
1  def refresh_openai_token():
2      token = cl.user_session.get('openai_token')
3      if token is None or token.expires_on < int(time.time()) - token_refresh_interval:
4          cl.user_session.set('openai_token', default_credential.get_token("https://cognitiveservic
5          openai.api_key = cl.user_session.get('openai_token').token
```

The function follows these steps:

1. It fetches the current token from `cl.user_session` (which seems to be a part of the `chainlit` library) using the key `'openai_token'`. The <u>user session</u> is a dictionary that stores the user's session data. The id and env keys are reserved for the session ID and environment variables. Other keys can be used to store arbitrary data in the user's session.
2. It checks if the token is `None` or if its expiration time (`expires_on`) is less than the current time minus 1800 seconds (30 minutes).

Next, the code defines a function called `start_chat` that is used to initialize the when the user connects to the application or clicks the `New Chat` button.

```
1   .on_chat_start
2   async def start_chat():
3       # Sending Avatars for Chat Participants
4       await cl.Avatar(
5           name="Chatbot",
6           url="https://cdn-icons-png.flaticon.com/512/8649/8649595.png"
7       ).send()
8       await cl.Avatar(
9           name="Error",
10          url="https://cdn-icons-png.flaticon.com/512/8649/8649595.png"
11      ).send()
12      await cl.Avatar(
13          name="User",
14          url="https://media.architecturaldigest.com/photos/5f241de2c850b2a36b415024/master/w_1600%
15      ).send()
```

Here is a brief explanation of the function steps:

- `@cl.on_chat_start` : The <u>on chat start</u> decorator registers a callback function `start_chat()` to be called when the Chainlit chat starts. It is used to set up the chat and send <u>avatars</u> for the Chatbot, Error, and User participants in the chat.
- `cl.Avatar()` : the <u>Avatar</u> class allows you to display an avatar image next to a message instead of the author's name. You need to send the element once. Next, if an avatar's name matches an author's name, the avatar will be automatically displayed. You must provide either a URL or a path or content bytes.

The following code is used to initialize the large language model (LLM) chain used to reply to questions on the content of the uploaded documents.

```
1   # Initialize the file list to None
2      files = None
3
4      # Wait for the user to upload a file
5      while files is None:
6          files = await cl.AskFileMessage(
7              content=f"Please upload up to {max_files} `.pdf` or `.docx` files to begin.",
8              accept=["application/pdf", "application/vnd.openxmlformats-officedocument.wordprocess
9              max_size_mb=max_size_mb,
10             max_files=max_files,
11             timeout=86400,
12             raise_on_timeout=False
13         ).send()
```

The AskFileMessage API call prompts the user to upload up to a specified number of `.pdf` or `.docx` files. The uploaded files are stored in the `files` variable. The process continues until the user uploads files. For more information, see AskFileMessage.

The following code processes each uploaded file by extracting its content.

1. The text content of each file is stored in the list `all_texts`.
2. This code performs text processing and chunking. It checks the file extension to read the file content accordingly, depending on if it's a `.pdf` or a `.docx` document.
3. The text content is split into smaller chunks using the RecursiveCharacterTextSplitter LangChain object.
4. Metadata is created for each chunk and stored in the `metadatas` list.
5. If `openai.api_type == "azure_ad"`, the code invokes the `refresh_openai_token()` that gets a security token from Azure AD to communicate with the Azure OpenAI Service.

```
1    # Create a message to inform the user that the files are being processed
2        content = ''
3        if (len(files)  ==  1):
4            content = f"Processing `{files[0].name}`..."
5        else:
6            files_names = [f"`{f.name}`" for f in files]
7            content = f"Processing {', '.join(files_names)}..."
8        msg = cl.Message(content = content, author = "Chatbot")
9        await msg.send()
10
11       # Create a list to store the texts of each file
12       all_texts = []
13
14       # Process each file uploaded by the user
15       for file in files:
16
17           # Create an in-memory buffer from the file content
18           bytes = io.BytesIO(file.content)
19
20           # Get file extension
21           extension = file.name.split('.')[-1]
22
23           # Initialize the text variable
24           text = ''
25
26           # Read the file
27           if extension == "pdf":
28               # ...
29           elif extension == "docx":
30               # ...
31
32           # Split the text into chunks
33           text_splitter = RecursiveCharacterTextSplitter(
34               chunk_size=text_splitter_chunk_size,
35               chunk_overlap=text_splitter_chunk_overlap
36           )
37           texts = text_splitter.split_text(text)
38
39           # Add the chunks and metadata to the list
40           all_texts.extend(texts)
41
42       # Create a metadata for each chunk
43       metadatas = [{"source": f"{i}-pl"} for i in range(len(all_texts))]
44
45       #  Refresh the OpenAI security token if using Azure AD
46       if openai.api_type == "azure_ad":
47           refresh_openai_token()
```

The next piece of code performs the following steps:

1. It creates an <u>OpenAIEmbeddings</u> configured to use the embeddings model in the Azure OpenAI Service to create embeddings from text chunks.

2. It creates a <u>ChromaDB</u> vector database using the `OpenAIEmbeddings` object, the text chunks list, and the metadata list.

3. It creates an <u>AzureChatOpenAI</u> LangChain object based on the GPR model hosted in Azure OpenAI Service.

4. It creates a chain using the [RetrievalQAWithSourcesChain.from_chain_type](#) API call uses previously created models and stores them as retrievers.

5. It stores the metadata and text chunks in the user session using the `cl.user_session.set()` API call.

6. It creates a message to inform the user that the files are ready for queries, and finally returns the `chain`.

7. The `cl.user_session.set("chain", chain)` call stores the LLM chain in the [user session](#) dictionary for later use.

```
1   #  Refresh the OpenAI security token if using Azure AD
2      if openai.api_type  ==  "azure_ad":
3          refresh_openai_token()
4
5      # Create a Chroma vector store
6      embeddings = OpenAIEmbeddings(
7          deployment = embeddings_deployment,
8          openai_api_key = openai.api_key,
9          openai_api_base = openai.api_base,
10         openai_api_version = openai.api_version,
11         openai_api_type = openai.api_type,
12         chunk_size = embeddings_chunk_size)
13
14     # Create a Chroma vector store
15     db = await cl.make_async(Chroma.from_texts)(
16         all_texts, embeddings, metadatas = metadatas
17     )
18
19     # Create an AzureChatOpenAI llm
20     llm = AzureChatOpenAI(
21         temperature = temperature,
22         openai_api_key = openai.api_key,
23         openai_api_base = openai.api_base,
24         openai_api_version = openai.api_version,
25         openai_api_type = openai.api_type,
26         deployment_name = chat_completion_deployment)
27
28     # Create a chain that uses the Chroma vector store
29     chain = RetrievalQAWithSourcesChain.from_chain_type(
30         llm = llm,
31         chain_type = "stuff",
32         retriever = db.as_retriever(),
33         return_source_documents = True,
34         chain_type_kwargs = chain_type_kwargs
35     )
36
37     # Save the metadata and texts in the user session
38     cl.user_session.set("metadatas", metadatas)
39     cl.user_session.set("texts", all_texts)
40
41     # Create a message to inform the user that the files are ready for queries
42     content = ''
43     if (len(files)  ==  1):
44         content = f"`{files[0].name}` processed. You can now ask questions!"
45     else:
46         files_names = [f"`{f.name}`" for f in files]
47         content = f"{', '.join(files_names)} processed. You can now ask questions."
48     msg.content = content
49     msg.author = "Chatbot"
50     await msg.update()
51
52     # Store the chain in the user session
53     cl.user_session.set("chain", chain)
```

The following code handles the communication with the OpenAI API and incorporates retrying logic in case the API calls fail due to specific errors.

- @cl.on_message : The on message decorator registers a callback function main(message: str) to be called when the user submits a new message in the chat. This is the main function responsible for handling

the chat logic.

- `cl.user_session.get("chain")` : this call retrieves the LLM chain from the <u>user session</u> dictionary.
- The `for` loop allows multiple attempts, up to `max_retries` , to communicate with the chat completion API and handles different types of API errors, such as timeout, connection error, invalid request, and service unavailability.
- `await chain.acall` : The asynchronous call to the <u>RetrievalQAWithSourcesChain.acall</u> executes the LLM chain with the user message as an input.

```
1   .on_message
2   async def run(message: str):
3       # Retrieve the chain from the user session
4       chain = cl.user_session.get("chain")
5
6       # Initialize the response
7       response =  None
8
9       # Retry the OpenAI API call if it fails
10      for attempt in range(max_retries):
11          try:
12              # Refresh the OpenAI security token if using Azure AD
13              if openai.api_type == "azure_ad":
14                  refresh_openai_token()
15
16              # Ask the question to the chain
17              response = await chain.acall(message, callbacks=[cl.AsyncLangchainCallbackHandler()])
18              break
19          except openai.error.Timeout:
20              # Exception handling for timeout error
21              # Implement exponential backoff
22              wait_time = backoff(attempt)
23              logger.exception(f"OpenAI API timeout occurred. Waiting {wait_time} seconds and tryin
24              time.sleep(wait_time)
25          except openai.error.APIError:
26              # Exception handling for API error
27              # Implement exponential backoff
28              wait_time = backoff(attempt)
29              logger.exception(f"OpenAI API error occurred. Waiting {wait_time} seconds and trying
30              time.sleep(wait_time)
31          except openai.error.APIConnectionError:
32              # Exception handling for API connection error
33              # Implement exponential backoff
34              wait_time = backoff(attempt)
35              logger.exception(f"OpenAI API connection error occurred. Check your network settings,
36              time.sleep(wait_time)
37          except openai.error.InvalidRequestError:
38              # Exception handling for invalid request error
39              # Implement exponential backoff
40              wait_time = backoff(attempt)
41              logger.exception(f"OpenAI API invalid request. Check the documentation for the specif
42              time.sleep(wait_time)
43          except openai.error.ServiceUnavailableError:
44              # Exception handling for service unavailable error
45              # Implement exponential backoff
46              wait_time = backoff(attempt)
47              logger.exception(f"OpenAI API service unavailable. Waiting {wait_time} seconds and tr
48              time.sleep(wait_time)
49          except Exception as e:
50              # Exception handling for non-retriable errors
51              logger.exception(f"A non-retriable error occurred. {e}")
52              break
```

The code below extracts the answers and sources from the API response and formats them to be sent as a message.

- The answer and sources are obtained from the response dictionary.

- The sources are then processed to find corresponding texts in the user session metadata ( `metadatas` ) and create `source_elements` using `cl.Text()` .
- `cl.Message().send()` : the <u>Message</u> API creates and displays a message containing the answer and sources, if available.

```
1   # Get the answer and sources from the response
2       answer = response["answer"]
3       sources = response["sources"].strip()
4       source_elements = []
5
6       # Get the metadata and texts from the user session
7       metadatas = cl.user_session.get("metadatas")
8       all_sources = [m["source"] for m in metadatas]
9       texts = cl.user_session.get("texts")
10
11      if sources:
12          found_sources = []
13
14          # Add the sources to the message
15          for source in sources.split(","):
16              source_name = source.strip().replace(".", "")
17              # Get the index of the source
18              try:
19                  index = all_sources.index(source_name)
20              except ValueError:
21                  continue
22              text = texts[index]
23              found_sources.append(source_name)
24              # Create the text element referenced in the message
25              source_elements.append(cl.Text(content=text, name=source_name))
26
27          if found_sources:
28              answer += f"\nSources: {', '.join(found_sources)}"
29          else:
30              answer += "\nNo sources found"
31
32      await cl.Message(content=answer, elements=source_elements).send()
```

Below, you can read the complete code of the application.

```python
# Import packages
import os
import io
import sys
import time
import openai
import random
import logging
import chainlit as cl
from pypdf import PdfReader
from docx import Document
from azure.identity import DefaultAzureCredential
from dotenv import load_dotenv
from dotenv import dotenv_values
from langchain.embeddings.openai import OpenAIEmbeddings
from langchain.text_splitter import RecursiveCharacterTextSplitter
from langchain.vectorstores import Chroma
from langchain.chains import RetrievalQAWithSourcesChain
from langchain.chat_models import AzureChatOpenAI
from langchain.prompts.chat import (
    ChatPromptTemplate,
    SystemMessagePromptTemplate,
    HumanMessagePromptTemplate,
)

# These three lines swap the stdlib sqlite3 lib with the pysqlite3 package
__import__('pysqlite3')
sys.modules['sqlite3'] = sys.modules.pop('pysqlite3')

# Load environment variables from .env file
if os.path.exists(".env"):
    load_dotenv(override = True)
    config = dotenv_values(".env")

# Read environment variables
temperature = float(os.environ.get("TEMPERATURE", 0.9))
api_base = os.getenv("AZURE_OPENAI_BASE")
api_key = os.getenv("AZURE_OPENAI_KEY")
api_type = os.environ.get("AZURE_OPENAI_TYPE", "azure")
api_version = os.environ.get("AZURE_OPENAI_VERSION", "2023-06-01-preview")
chat_completion_deployment = os.getenv("AZURE_OPENAI_DEPLOYMENT")
embeddings_deployment = os.getenv("AZURE_OPENAI_ADA_DEPLOYMENT")
model = os.getenv("AZURE_OPENAI_MODEL")
max_size_mb = int(os.getenv("CHAINLIT_MAX_SIZE_MB", 100))
max_files = int(os.getenv("CHAINLIT_MAX_FILES", 10))
text_splitter_chunk_size = int(os.getenv("TEXT_SPLITTER_CHUNK_SIZE", 1000))
text_splitter_chunk_overlap = int(os.getenv("TEXT_SPLITTER_CHUNK_OVERLAP", 10))
embeddings_chunk_size = int(os.getenv("EMBEDDINGS_CHUNK_SIZE", 16))
max_retries = int(os.getenv("MAX_RETRIES", 5))
backoff_in_seconds = float(os.getenv("BACKOFF_IN_SECONDS", 1))

# Configure system prompt
system_template = """Use the following pieces of context to answer the users question.
If you don't know the answer, just say that you don't know, don't try to make up an answer.
ALWAYS return a "SOURCES" part in your answer.
The "SOURCES" part should be a reference to the source of the document from which you got your an

Example of your response should be:
```

```
---

The answer is foo
SOURCES: xyz

---

Begin!
----------------
{summaries}"""
messages = [
    SystemMessagePromptTemplate.from_template(system_template),
    HumanMessagePromptTemplate.from_template("{question}"),
]
prompt = ChatPromptTemplate.from_messages(messages)
chain_type_kwargs = {"prompt": prompt}

# Configure OpenAI
openai.api_type = api_type
openai.api_version = api_version
openai.api_base = api_base
openai.api_key = api_key

# Set default Azure credential
default_credential = DefaultAzureCredential(
) if openai.api_type  ==  "azure_ad" else None

# Configure a logger
logging.basicConfig(stream = sys.stdout,
                    format = '[%(asctime)s] {%(filename)s:%(lineno)d} %(levelname)s - %(message)s
                    level = logging.INFO)
logger = logging.getLogger(__name__)

# Refresh the OpenAI security token every 45 minutes
def refresh_openai_token():
    token = cl.user_session.get('openai_token')
    if token  ==  None or token.expires_on < int(time.time()) - 1800:
        cl.user_session.set('openai_token', default_credential.get_token(
            "https://cognitiveservices.azure.com/.default"))
        openai.api_key = cl.user_session.get('openai_token').token

def backoff(attempt : int) -> float:
    return backoff_in_seconds * 2**attempt + random.uniform(0, 1)

@cl.on_chat_start
async def start():
    await cl.Avatar(
        name = "Chatbot",
        url = "https://cdn-icons-png.flaticon.com/512/8649/8649595.png"
    ).send()
    await cl.Avatar(
        name = "Error",
        url = "https://cdn-icons-png.flaticon.com/512/8649/8649595.png"
    ).send()
    await cl.Avatar(
        name = "User",
        url = "https://media.architecturaldigest.com/photos/5f241de2c850b2a36b415024/master/w_160
    ).send()

    # Initialize the file list to None for container
```

```python
    files = None

    # Wait for the user to upload a file
    while files  ==  None:
        files = await cl.AskFileMessage(
            content = f"Please upload up to {max_files} `.pdf` or `.docx` files to begin.",
            accept = ["application/pdf",
                    "application/vnd.openxmlformats-officedocument.wordprocessingml.document"],
            max_size_mb = max_size_mb,
            max_files = max_files,
            timeout = 86400,
            raise_on_timeout = False
        ).send()

    # Create a message to inform the user that the files are being processed
    content = ''
    if (len(files)  ==  1):
        content = f"Processing `{files[0].name}`..."
    else:
        files_names = [f"`{f.name}`" for f in files]
        content = f"Processing {', '.join(files_names)}..."
    msg = cl.Message(content = content, author = "Chatbot")
    await msg.send()

    # Create a list to store the texts of each file
    all_texts = []

    # Process each file uplodaded by the user
    for file in files:

        # Create an in-memory buffer from the file content
        bytes = io.BytesIO(file.content)

        # Get file extension
        extension = file.name.split('.')[-1]

        # Initialize the text variable
        text = ''

        # Read the file
        if extension  ==  "pdf":
            reader = PdfReader(bytes)
            for i in range(len(reader.pages)):
                text +=  reader.pages[i].extract_text()
        elif extension  ==  "docx":
            doc = Document(bytes)
            paragraph_list = []
            for paragraph in doc.paragraphs:
                paragraph_list.append(paragraph.text)
            text = '\n'.join(paragraph_list)

        # Split the text into chunks
        text_splitter = RecursiveCharacterTextSplitter(
            chunk_size = text_splitter_chunk_size,
            chunk_overlap = text_splitter_chunk_overlap)
        texts = text_splitter.split_text(text)

        # Add the chunks and metadata to the list
        all_texts.extend(texts)
```

```python
    # Create a metadata for each chunk
    metadatas = [{"source": f"{i}-pl"} for i in range(len(all_texts))]

    #  Refresh the OpenAI security token if using Azure AD
    if openai.api_type  ==  "azure_ad":
        refresh_openai_token()

    # Create a Chroma vector store
    embeddings = OpenAIEmbeddings(
        deployment = embeddings_deployment,
        openai_api_key = openai.api_key,
        openai_api_base = openai.api_base,
        openai_api_version = openai.api_version,
        openai_api_type = openai.api_type,
        chunk_size = embeddings_chunk_size)

    # Create a Chroma vector store
    db = await cl.make_async(Chroma.from_texts)(
        all_texts, embeddings, metadatas = metadatas
    )

    # Create an AzureChatOpenAI llm
    llm = AzureChatOpenAI(
        temperature = temperature,
        openai_api_key = openai.api_key,
        openai_api_base = openai.api_base,
        openai_api_version = openai.api_version,
        openai_api_type = openai.api_type,
        deployment_name = chat_completion_deployment)

    # Create a chain that uses the Chroma vector store
    chain = RetrievalQAWithSourcesChain.from_chain_type(
        llm = llm,
        chain_type = "stuff",
        retriever = db.as_retriever(),
        return_source_documents = True,
        chain_type_kwargs = chain_type_kwargs
    )

    # Save the metadata and texts in the user session
    cl.user_session.set("metadatas", metadatas)
    cl.user_session.set("texts", all_texts)

    # Create a message to inform the user that the files are ready for queries
    content = ''
    if (len(files)  ==  1):
        content = f"`{files[0].name}` processed. You can now ask questions!"
    else:
        files_names = [f"`{f.name}`" for f in files]
        content = f"{', '.join(files_names)} processed. You can now ask questions."
    msg.content = content
    msg.author = "Chatbot"
    await msg.update()

     # Store the chain in the user session
    cl.user_session.set("chain", chain)

@cl.on_message
async def run(message: str):
    # Retrieve the chain from the user session
```

```python
chain = cl.user_session.get("chain")

# Initialize the response
response =  None

# Retry the OpenAI API call if it fails
for attempt in range(max_retries):
    try:
        # Refresh the OpenAI security token if using Azure AD
        if openai.api_type  ==  "azure_ad":
            refresh_openai_token()

        # Ask the question to the chain
        response = await chain.acall(message, callbacks = [cl.AsyncLangchainCallbackHandler()
        break
    except openai.error.Timeout:
        # Implement exponential backoff
        wait_time = backoff(attempt)
        logger.exception(f"OpenAI API timeout occurred. Waiting {wait_time} seconds and tryin
        time.sleep(wait_time)
    except openai.error.APIError:
        # Implement exponential backoff
        wait_time = backoff(attempt)
        logger.exception(f"OpenAI API error occurred. Waiting {wait_time} seconds and trying
        time.sleep(wait_time)
    except openai.error.APIConnectionError:
        # Implement exponential backoff
        wait_time = backoff(attempt)
        logger.exception(f"OpenAI API connection error occurred. Check your network settings,
        time.sleep(wait_time)
    except openai.error.InvalidRequestError:
        # Implement exponential backoff
        wait_time = backoff(attempt)
        logger.exception(f"OpenAI API invalid request. Check the documentation for the specif
        time.sleep(wait_time)
    except openai.error.ServiceUnavailableError:
        # Implement exponential backoff
        wait_time = backoff(attempt)
        logger.exception(f"OpenAI API service unavailable. Waiting {wait_time} seconds and tr
        time.sleep(wait_time)
    except Exception as e:
        logger.exception(f"A non retriable error occurred. {e}")
        break

# Get the answer and sources from the response
answer = response["answer"]
sources = response["sources"].strip()
source_elements = []

# Get the metadata and texts from the user session
metadatas = cl.user_session.get("metadatas")
all_sources = [m["source"] for m in metadatas]
texts = cl.user_session.get("texts")

if sources:
    found_sources = []

    # Add the sources to the message
    for source in sources.split(","):
        source_name = source.Skip to main content (button)
```

```
303        # Get the index of the source
304        try:
305            index = all_sources.index(source_name)
306        except ValueError:
307            continue
308        text = texts[index]
309        found_sources.append(source_name)
310        # Create the text element referenced in the message
311        source_elements.append(cl.Text(content = text, name = source_name))
312
313    if found_sources:
314        answer +=  f"\nSources: {', '.join(found_sources)}"
315    else:
        answer +=  "\nNo sources found"

    await cl.Message(content = answer, elements = source_elements).send()
```

You can run the application locally using the following command. The `-w` flag` indicates auto-reload whenever we make changes live in our application code.

```
1 | chainlit run app.py -w
```

## Build Docker Images

You can use the `src/01-build-docker-images.sh` Bash script to build the Docker container image for each container app.

```
1   #!/bin/bash
2
3   # Variables
4   source ./00-variables.sh
5
6   # Use a for loop to build the docker images using the array index
7   for index in ${!images[@]}; do
8     # Build the docker image
9     docker build -t ${images[$index]}:$tag -f Dockerfile --build-arg FILENAME=${filenames[$index]}
10  done
```

Before running any script in the `src` folder, make sure to customize the value of the variables inside the `00-variables.sh` file located in the same folder. This file is embedded in all the scripts and contains the following variables:

```
1   # Variables
2
3   # Azure Container Registry
4   prefix="Blue"
5   acrName="${prefix}Registry"
6   acrResourceGrougName="${prefix}RG"
7   location="EastUS"
8
9   # Python Files
10  docAppFile="doc.py"
11  chatAppFile="chat.py"
12
13  # Docker Images
14  docImageName="doc"
15  chatImageName="chat"
16  tag="v1"
17  port="8000"
18
19  # Arrays
20  images=($docImageName $chatImageName)
21  filenames=($docAppFile $chatAppFile)
```

The `Dockerfile` under the `src` folder is parametric and can be used to build the container images for both chat applications.

```
# app/Dockerfile

# # Stage 1 - Install build dependencies

# A Dockerfile must start with a FROM instruction that sets the base image for the container.
# The Python images come in many flavors, each designed for a specific use case.
# The python:3.11-slim image is a good base image for most applications.
# It is a minimal image built on top of Debian Linux and includes only the necessary packages to
# The slim image is a good choice because it is small and contains only the packages needed to ru
# For more information, see:
# * https://hub.docker.com/_/python
# * https://docs.streamlit.io/knowledge-base/tutorials/deploy/docker
FROM python:3.11-slim AS builder

# The WORKDIR instruction sets the working directory for any RUN, CMD, ENTRYPOINT, COPY and ADD i
# If the WORKDIR doesn't exist, it will be created even if it's not used in any subsequent Docker
# For more information, see: https://docs.docker.com/engine/reference/builder/#workdir
WORKDIR /app

# Set environment variables.
# The ENV instruction sets the environment variable <key> to the value <value>.
# This value will be in the environment of all "descendant" Dockerfile commands and can be replac
# For more information, see: https://docs.docker.com/engine/reference/builder/#env
ENV PYTHONDONTWRITEBYTECODE 1
ENV PYTHONUNBUFFERED 1

# Install git so that we can clone the app code from a remote repo using the RUN instruction.
# The RUN comand has 2 forms:
# * RUN <command> (shell form, the command is run in a shell, which by default is /bin/sh -c on L
# * RUN ["executable", "param1", "param2"] (exec form)
# The RUN instruction will execute any commands in a new layer on top of the current image and co
# The resulting committed image will be used for the next step in the Dockerfile.
# For more information, see: https://docs.docker.com/engine/reference/builder/#run
RUN apt-get update && apt-get install -y \
  build-essential \
  curl \
  software-properties-common \
  git \
  && rm -rf /var/lib/apt/lists/*

# Create a virtualenv to keep dependencies together
RUN python -m venv /opt/venv
ENV PATH="/opt/venv/bin:$PATH"

# Clone the requirements.txt which contains dependencies to WORKDIR
# COPY has two forms:
# * COPY <src> <dest> (this copies the files from the local machine to the container's own filesy
# * COPY ["<src>",... "<dest>"] (this form is required for paths containing whitespace)
# For more information, see: https://docs.docker.com/engine/reference/builder/#copy
COPY requirements.txt .

# Install the Python dependencies
RUN pip install --no-cache-dir --no-deps -r requirements.txt

# Stage 2 - Copy only necessary files to the runner stage

# The FROM instruction initializes a new build stage for the application
FROM python:3.11-slim
```

```
61   # Define the filename to copy as an argument
62   ARG FILENAME
63
64   # Deefine the port to run the application on as an argument
65   ARG PORT=8000
66
67   # Set an environment variable
68   ENV FILENAME=${FILENAME}
69
70   # Sets the working directory to /app
71   WORKDIR /app
72
73   # Copy the virtual environment from the builder stage
74   COPY --from=builder /opt/venv /opt/venv
75
76   # Set environment variables
77   ENV PATH="/opt/venv/bin:$PATH"
78
79   # Clone the $FILENAME containing the application code
80   COPY $FILENAME .
81
82   # Copy the chainlit.md file to the working directory
83   COPY chainlit.md .
84
85   # Copy the .chainlit folder to the working directory
86   COPY ./.chainlit ./.chainlit
87
88   # The EXPOSE instruction informs Docker that the container listens on the specified network ports
89   # For more information, see: https://docs.docker.com/engine/reference/builder/#expose
90   EXPOSE $PORT
91
92   # The ENTRYPOINT instruction has two forms:
93   # * ENTRYPOINT ["executable", "param1", "param2"] (exec form, preferred)
94   # * ENTRYPOINT command param1 param2 (shell form)
95   # The ENTRYPOINT instruction allows you to configure a container that will run as an executable.
96   # For more information, see: https://docs.docker.com/engine/reference/builder/#entrypoint
     CMD chainlit run $FILENAME --port=$PORT
```

## Test applications locally

You can use the `src/02-run-docker-container.sh` Bash script to test the containers for
the `sender`, `processor`, and `receiver` applications.

```bash
#!/bin/bash

# Variables
source ./00-variables.sh

# Print the menu
echo "================================="
echo "Run Docker Container (1-3): "
echo "================================="
options=(
  "Doc"
  "Chat"
)
name=""
# Select an option
COLUMNS=0
select option in "${options[@]}"; do
  case $option in
    "Doc")
      docker run -it \
      --rm \
      -p $port:$port \
      -e AZURE_OPENAI_BASE=$AZURE_OPENAI_BASE \
      -e AZURE_OPENAI_KEY=$AZURE_OPENAI_KEY \
      -e AZURE_OPENAI_MODEL=$AZURE_OPENAI_MODEL \
      -e AZURE_OPENAI_DEPLOYMENT=$AZURE_OPENAI_DEPLOYMENT \
      -e AZURE_OPENAI_ADA_DEPLOYMENT=$AZURE_OPENAI_ADA_DEPLOYMENT \
      -e AZURE_OPENAI_VERSION=$AZURE_OPENAI_VERSION \
      -e AZURE_OPENAI_TYPE=$AZURE_OPENAI_TYPE \
      -e TEMPERATURE=$TEMPERATURE \
      --name $docImageName \
      $docImageName:$tag
      break
    ;;
    "Chat")
      docker run -it \
      --rm \
      -p $port:$port \
      -e AZURE_OPENAI_BASE=$AZURE_OPENAI_BASE \
      -e AZURE_OPENAI_KEY=$AZURE_OPENAI_KEY \
      -e AZURE_OPENAI_MODEL=$AZURE_OPENAI_MODEL \
      -e AZURE_OPENAI_DEPLOYMENT=$AZURE_OPENAI_DEPLOYMENT \
      -e AZURE_OPENAI_VERSION=$AZURE_OPENAI_VERSION \
      -e AZURE_OPENAI_TYPE=$AZURE_OPENAI_TYPE \
      -e TEMPERATURE=$TEMPERATURE \
      --name $chatImageName \
      $chatImageName:$tag
      break
    ;;
    "Quit")
      exit
    ;;
    *) echo "invalid option $REPLY" ;;
  esac
done
```

# Push Docker containers to the Azure Container Registry

You can use the `src/03-push-docker-image.sh` Bash script to push the Docker container images for the `sender`, `processor`, and `receiver` applications to the Azure Container Registry (ACR)

```bash
#!/bin/bash

# Variables
source ./00-variables.sh

# Login to ACR
echo "Logging in to [${acrName,,}] container registry..."
az acr login --name ${acrName,,}

# Retrieve ACR login server. Each container image needs to be tagged with the loginServer name of
echo "Retrieving login server for the [${acrName,,}] container registry..."
loginServer=$(az acr show --name ${acrName,,} --query loginServer --output tsv)

# Use a for loop to tag and push the local docker images to the Azure Container Registry
for index in ${!images[@]}; do
  # Tag the local sender image with the loginServer of ACR
  docker tag ${images[$index],,}:$tag $loginServer/${images[$index],,}:$tag

  # Push the container image to ACR
  docker push $loginServer/${images[$index],,}:$tag
done
```

## Monitoring

Azure Container Apps provides several built-in observability features that together give you a holistic view of your container app's health throughout its application lifecycle. These features help you monitor and diagnose the state of your app to improve performance and respond to trends and critical problems.

You can use the `Log Stream` panel on the Azure Portal to see the logs generated by a container app, as shown in the following screenshot.

Alternatively, you can click open the `Logs` panel, as shown in the following screenshot, and use a Kusto Query Language (KQL) query to filter, project, and retrieve only the desired data.

## Review deployed resources

You can use the Azure portal to list the deployed resources in the resource group, as shown in the following picture:

| | | | |
|---|---|---|---|
| ∨ Azure OpenAI | | | |
| ☐ 🟢 BlueOpenAI | Azure OpenAI | East US | ••• |
| ∨ Container App | | | |
| ☐ 🔷 chatapp | Container App | East US | ••• |
| ☐ 🔷 docapp | Container App | East US | ••• |
| ∨ Container Apps Environment | | | |
| ☐ 🔷 BlueEnvironment | Container Apps Environment | East US | ••• |
| ∨ Container registry | | | |
| ☐ ☁ BlueRegistry | Container registry | East US | ••• |
| ∨ Log Analytics workspace | | | |
| ☐ 📊 BlueWorkspace | Log Analytics workspace | East US | ••• |
| ∨ Managed Identity | | | |
| ☐ 🔑 BlueRegistryIdentity | Managed Identity | East US | ••• |
| ☐ 🔑 BlueWorkloadIdentity | Managed Identity | East US | ••• |
| ∨ Network Interface | | | |
| ☐ 🔷 BlueEnvironment | Container Apps Environment | East US | ••• |
| ∨ Container registry | | | |
| ☐ ☁ BlueRegistry | Container registry | East US | ••• |
| ∨ Log Analytics workspace | | | |
| ☐ 📊 BlueWorkspace | Log Analytics workspace | East US | ••• |
| ∨ Managed Identity | | | |
| ☐ 🔑 BlueRegistryIdentity | Managed Identity | East US | ••• |
| ☐ 🔑 BlueWorkloadIdentity | Managed Identity | East US | ••• |
| ∨ Network Interface | | | |
| ☐ 🖥 BlueOpenAIPrivateEndpoint.nic.7cde7fdf-5590-4db7-9e7b-91··· | Network Interface | East US | ••• |
| ☐ 🖥 BlueRegistryPrivateEndpoint.nic.5767ff52-4ee8-484a-832d-d4··· | Network Interface | East US | ••• |
| ∨ Private DNS zone | | | |
| ☐ 🔵 privatelink.azurecr.io | Private DNS zone | Global | ••• |
| ☐ 🔵 privatelink.openai.azure.com | Private DNS zone | Global | ••• |
| ∨ Private endpoint | | | |
| ☐ ‹I› BlueOpenAIPrivateEndpoint | Private endpoint | East US | ••• |
| ☐ ‹I› BlueRegistryPrivateEndpoint | Private endpoint | East US | ••• |
| ∨ Virtual network | | | |
| ☐ ‹··› BlueVNet | Virtual network | East US | ••• |

You can also use Azure CLI to list the deployed resources in the resource group:

```
1  az resource list --resource-group <resource-group-name>
```

You can also use the following PowerShell cmdlet to list the deployed resources in the resource group:

```
1  Get-AzResource -ResourceGroupName <resource-group-name>
```

## Clean up resources

You can delete the resource group using the following Azure CLI command when you no longer need the resources you created. This will remove all the Azure resources.

```
1  az group delete --name <resource-group-name>
```

Alternatively, you can use the following PowerShell cmdlet to delete the resource group and all the Azure resources.

```
1  Remove-AzResourceGroup -Name <resource-group-name>
```

👍 3 Likes

## 6 Comments

---

**cicorias** Microsoft

Jul 27 2023 08:40 A

Thanks again for such relevant and timely guidance Paolo!!

👍 1 Like

---

**paolosalvatori** Microsoft

Jul 27 2023 09:01 A

Thanks @cicorias, I thought to create an Azure Container Apps + Azure OpenAI sample after my articles on AKS + Azure OpenAI.

👍 1 Like

---

**merveguel** Copper Contributor

Aug 15 2023 02:32

Hi @paolosalvatori, thanks for such detailed post, I was wondering if this article can be implemented using the free-tier version of the prerequisite tools?
Best, Merve

👍 0 Likes

---

**paolosalvatori** Microsoft

Aug 15 2023 03:17

Hi @merveguel, all the pre-requisites are free of charge. I think you can use a free Azure account to test the entire architecture, but I'm not sure about Azure OpenAI Service. I should check, but I'm OOF now. If you don't have an Azure account already, the best way is to open one. Alternatively, check the pricing page for the various services and the Azure Free Account FAQ: https://azure.microsoft.com/en-us/free/free-account-faq

👍 1 Like

Aug 29 2023 01:02

**Heman_k** Copper Contributor

Hi @paolosalvatori , Thanks for sharing this very detailed article. Covers a lot of ground!!

On a related topic, can you please comment on the vector search capabilities in Azure Search? How does it compare it with vector search in other similar products?

Also, I am specifically interested in the support for HNSW algorithm in in Azure Search. Does Azure Search support HNSW? I have seen sample code where there is a module called HnswVectorSearchAlgorithmConfiguration in the azure.search.documents.indexes.models package, but I am having problems importing this module, at least in Python 3.9.  Any info or suggestions?

👍 0 Likes

Aug 29 2023 11:59

**paolosalvatori** Microsoft

Thanks @Heman_k

I'm not an Azure Cognitive Search subject-matter expert, but I can surely affirm that Azure Cognitive Search has strong vector search capabilities that allow you to search and retrieve similar vectors based on their similarity scores. While I do not have direct comparisons to other similar products or other vector databases, such as Chroma or FAISS, Azure Cognitive Search's vector search capabilities are designed to provide efficient and accurate results.
Regarding the HNSW algorithm, Azure Search does support it. The HNSW algorithm is a commonly used algorithm for approximate nearest neighbor search, and it is available in Azure Cognitive Search for vector search scenarios.

Please check Add vector search - Azure Cognitive Search | Microsoft Learn.
As for the specific issue you mentioned regarding importing the HnswVectorSearchAlgorithmConfiguration module in Python 3.9, it's possible that the module may not be available or accessible in that specific version. It is recommended to check the documentation or reach out to Microsoft Azure support for further assistance or alternative approaches.
Overall, Azure Search offers robust vector search capabilities, including support for the HNSW algorithm, which can be valuable for creating efficient and accurate search experiences.

P.S. If you found my article and sample interesting and helpful, please like the article and star the GitHub project, thanks :)

👍 1 Like

You must be a registered user to add a comment. If you've already registered, sign in. Otherwise, register and sign in.

✕

Comment

## Co-Authors

paolosalvatori

## Version history

**Last update:**   Jul 27 2023 07:40 AM
**Updated by:**     paolosalvatori

## Labels

| | |
|---|---|
| App | **101** |
| Cloud Native Apps | **70** |
| Data & AI | **110** |

‹ Previous    Next ›

## Share

**What's new**

Surface Pro 9

Surface Laptop 5

Surface Studio 2+

Surface Laptop Go 2

Surface Laptop Studio

Surface Duo 2

Microsoft 365

Windows 11 apps

**Microsoft Store**

Account profile

Download Center

Microsoft Store support

Returns

Order tracking

Virtual workshops and training

Microsoft Store Promise

Flexible Payments

**Education**

Microsoft in education

Devices for education

Microsoft Teams for Education

Skip to main content

Microsoft 365 Education

Education consultation appointment

Educator training and development

Deals for students and parents

Azure for students

## Business

Microsoft Cloud

Microsoft Security

Dynamics 365

Microsoft 365

Microsoft Power Platform

Microsoft Teams

Microsoft Industry

Small Business

## Developer & IT

Azure

Developer Center

Documentation

Microsoft Learn

Microsoft Tech Community

Azure Marketplace

AppSource

Visual Studio

## Company

Careers

About Microsoft

Company news

Privacy at Microsoft

Investors

Diversity and inclusion

Accessibility

Sustainability