# Hosting your Web Application on Azure

Azure Machine Learning provides a powerful platform for running and deploying large language models (LLMs). This guide will walk you through the steps to deploy your LLM and create an endpoint for your LLM in Azure Machine Learning.

## Prerequisites

- An active Azure account  - Setting Up a free Azure Account - NLP Community - Confluence (iabg.de)
- An API-accessible LLM, either from closed-source providers (GPT, Claude, Mistral Large etc.) OR a self-deployed LLM  - Deploying Large Language Model (LLM) on Azure

## Steps to Host your Web Application

- **Sign in to Azure Portal:**

    - Open your web browser and navigate to the Azure Portal: https://portal.azure.com.
    - Sign in with your Azure account credentials

- **From the home page, select "Create a resource" and then "Web App":**



- **Fill in the initial configuration for your Web App:**

For this tutorial, we are deploying a Python web-based application directly from a code repository. Other options (e.g. Docker deployment, other programming languages) can also be specified here.

Under "Basics", select your free trial subscription and a resource group (if you have created one). Specify a name for your web app (which will also be used as the final URL), select "Code" as publish mode, your prefered Python version and Region (e.g. "Germany West Central").

Basics    Database    Deployment    Networking    Monitoring    Tags    Review + create

App Service Web Apps lets you quickly build, deploy, and scale enterprise-grade web, mobile, and API apps running on any platform. Meet rigorous performance, scalability, security and compliance requirements while using a fully managed platform to perform infrastructure maintenance.  Learn more

**Project Details**

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

Subscription *  ⓘ                         | Free trial                                    ⌄ |

        Resource Group *  ⓘ               | TVL                                           ⌄ |
                                          Create new

**Instance Details**

Name *                                    | my-test-iabg-webapp                          ✓ |
                                                                            .azurewebsites.net

Publish *          ● Code   ○ Docker Container   ○ Static Web App

Runtime stack *                           | Python 3.12                                   ⌄ |

Operating System *     ● Linux    ○ Windows

Region *                                  | Germany West Central                          ⌄ |

                    ⓘ Not finding your App Service Plan? Try a different region or select your App
                       Service Environment.

**Pricing plans**

App Service plan pricing tier determines the location, features, cost and compute resources associated with your app.
Learn more ⧉

Linux Plan (Germany West Central) *  ⓘ    | ASP-TVL-b586 (F1)                             ⌄ |
                                          Create new

Pricing plan                              **Free F1** (Shared infrastructure)

**Zone redundancy**

An App Service plan can be deployed as a zone redundant service in the regions that support it. This is a deployment

| Review + create |    | < Previous |    | Next : Database > |

You can also customize other aspects of your applications in the other tabs. For example, if you want your app to be publicly accessible, ensure that "Enable public access" is turned on under "Networking":

Basics    Database    Deployment    **Networking**    Monitoring    Tags    Review + create

Web Apps can be provisioned with the inbound address being public to the internet or isolated to an Azure virtual network. Web Apps can also be provisioned with outbound traffic able to reach endpoints in a virtual network, be governed by network security groups or affected by virtual network routes. By default, your app is open to the internet and cannot reach into a virtual network. These aspects can also be changed after the app is provisioned.  Learn more ⧉

Enable public access *  ⓘ          ● On    ○ Off

⚠ Network injection is only available in Basic, Standard, Premium, Premium V2, and Premium V3 Dedicated App Service plans.

Enable network injection            ○ On    ● Off

In general, the default options for the advanced tabs are sufficient.
When you are done, click on "Review + create" in the lower part of the screen.

- **Wait for the App to be deployed**

when deployment is complete, click on "Go to resource". This page gives you an overview on everything related to your current WebApp deployment (usage, logs, monitoring, accessibility, etc..)

- **Connect a repository to your Web App**

    Under "Deployment Center" in the left pane, you can specify the source of your app code. Fill in the Source provider (in our example "GitHub"), repository name, branch and other details.
    This code repository will be copied in the app container. You can access it later on via SSH (see Debugging).



- **Set up an entrypoint for your Web App**

    Under "Settings" in the left pane, select the "General Settings" tab. Then, specify your entrypoint options.

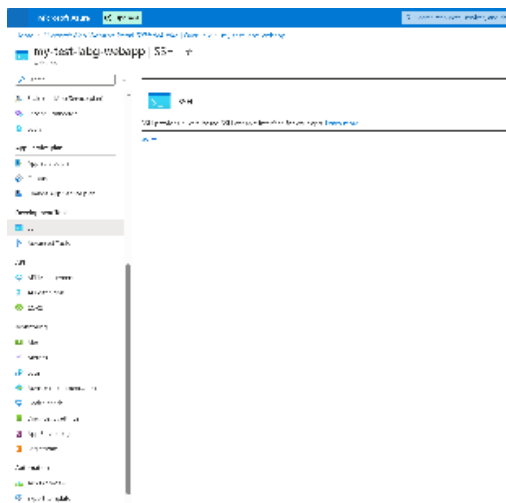In the example below, we specified a Python start-up script with CLI arguments (LLM API endpoint URL, LLM API key).



# Debugging your Web App

## SSH Connection

You can use the SSH connection option from the left pane to connect inside your app container. There you have all your application code and computing resources. You can also set up a debugger via SSH tunnel (e.g. in VSCode) or try running tests directly.

## Streaming Logs

Another possibility is to turn on logs in the Azure WebApp configs. In this way you can access the console stream of your program and gather diagnostic information.

To do so, select "App Service logs" from the left pane, and turn on Application logging in Filesystem mode (specify retention period and storage quota).



After this step, you can go to "Log stream" and access the application logs:

# ai-microclass-cloud-tutorial | Log stream ★ ⋯
Web App

🔔 Alerts
📊 Metrics
📄 Logs
💡 Advisor recommendations
💚 Health check
📋 Diagnostic settings
📲 App Service logs
📡 Log stream

**Automation**

👥 Tasks (preview)
📄 Export template

**Support + troubleshooting**

🩺 Resource health
❓ Support + Troubleshooting

⟳ Reconnect   📋 Copy   ❚❚ Pause   ✕ Clear

```
2024-03-17T00:41:42.330717612Z /tmp/8dc4619c92e24ca/antenv/lib/python3.11/site-
packages/langchain_core/_api/deprecation.py:117: LangChainDeprecationWarning:
deprecated in LangChain 0.1.0 and will be removed in 0.2.0. Use Use new agent
create_react_agent, create_json_agent, create_structured_chat_agent, etc. inste
2024-03-17T00:41:42.341581186Z   warn_deprecated(
/home/LogFiles/2024_03_17_lw1sdlwk0005YL_docker.log (https://ai-microclass-clo
tutorial.scm.azurewebsites.net/api/vfs/LogFiles/2024_03_17_lw1sdlwk0005YL_dock
2024-03-17T00:40:13.434Z INFO  - Starting container for site
2024-03-17T00:40:13.434Z INFO  - docker run -d --expose=8000 --name ai-microcla
WEBSITE_USE_DIAGNOSTIC_SERVER=false -e WEBSITE_SITE_NAME=ai-microclass-cloud-t
WEBSITE_AUTH_ENABLED=False -e WEBSITE_ROLE_INSTANCE_ID=0 -e WEBSITE_HOSTNAME=a
tutorial.azurewebsites.net -e WEBSITE_INSTANCE_ID=30020627b4609efb836376960cae
-e HTTP_LOGGING_ENABLED=1 appsvc/python:3.11_20240207.3.tuxprod python3 -m src
2024-03-17T00:40:17.142Z INFO  - Initiating warmup request to container ai-mic
for site ai-microclass-cloud-tutorial
2024-03-17T00:40:50.332z INFO  - Waiting for response to warmup request for con
tutorial_6_26cf8afa. Elapsed time = 33.1904272 sec
2024-03-17T00:41:08.185z INFO  - Waiting for response to warmup request for con
tutorial_6_26cf8afa. Elapsed time = 51.0439689 sec
2024-03-17T00:41:25.443Z INFO  - Waiting for response to warmup request for con
tutorial_6_26cf8afa. Elapsed time = 68.3017171 sec
2024-03-17T00:41:41.898Z INFO  - Waiting for response to warmup request for con
tutorial_6_26cf8afa. Elapsed time = 84.7564906 sec
2024-03-17T00:41:44.242Z INFO  - Container ai-microclass-cloud-tutorial_6_26cf
tutorial initialized successfully and is ready to serve requests.
2024-03-17T06:27:39.510Z INFO  -
Ending Log Tail of existing logs ---Starting Live Log Stream ---
```