

[Event dashboard](#) > [AgentCore Runtime](#) > Hosting MCP server

## Hosting MCP server

### Overview

[Amazon Bedrock AgentCore Runtime](#) lets you deploy Model Context Protocol (MCP) servers without managing infrastructure. This guide covers how to host MCP tools using the [Amazon Bedrock AgentCore Python SDK](#).

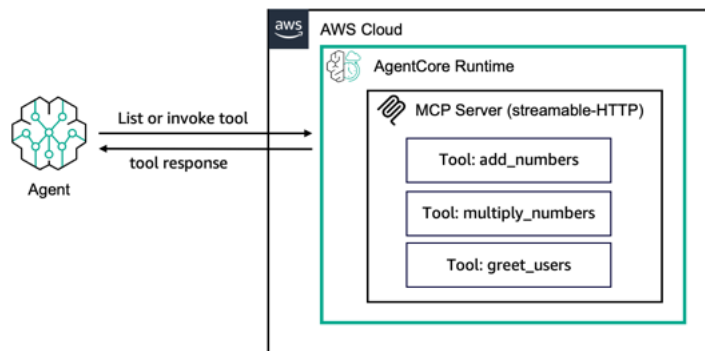
The SDK wraps your agent functions as MCP-compatible servers, handling protocol details so you can focus on logic. It supports both HTTP and MCP protocols for tool and agent interaction.

For tool hosting, the SDK uses **Stateless Streamable HTTP** with the `Mcp-Session-Id` header for [session isolation](#). Servers must run on `0.0.0.0:8000` and expose `/mcp`.

AgentCore automatically adds `Mcp-Session-Id` to support stateless, session-aware communication. The `InvokeAgentRuntime` API passes through all payloads, enabling seamless MCP message proxying.

### High Level Architecture

We will use a very simple MCP server with 3 tools: `add_numbers`, `multiply_numbers` and `greet_users`



- Agent Initialization:** The agent starts and is configured to interact with the configured tools via the MCP Server hosted on AWS using the AgentCore Runtime.
- Tool Discovery or Invocation:** The agent sends a request to the MCP Server (using streamable-HTTP) to either list available tools or invoke a specific tool (e.g., `add_numbers`, `multiply_numbers`, or `greet_users`).
- Tool Execution:** The MCP Server, running inside the AgentCore Runtime, processes the request and routes it to the appropriate tool implementation.
- Tool Response Returned:** Once the selected tool completes its logic (e.g., performs a computation or generates a greeting), the MCP Server returns the response to the agent.
- Agent Processes Result:** The agent receives the result from the MCP Server and uses it to continue its reasoning or return a final response to the user or system.

### Getting Started

#### Environment Requirements

Before starting this lab, ensure you have:

- **Python 3.10+** installed on your system
- **AWS credentials** configured with appropriate permissions
- **Docker** running for containerization
- **Amazon Bedrock AgentCore SDK** installed
- **Strands Agents** framework installed

## Create MCP Server

### Key MCP Server Components Explained:

- **FastMCP**: Creates an MCP server that can host your tools
- **@mcp.tool()**: Decorator that turns your Python functions into MCP tools
- **stateless\_http=True**: Required for AgentCore Runtime compatibility
- **Tools**: Three simple tools demonstrating different types of operations

## Create Local Testing Client

### ► Code Snippet

## Testing Locally

To test your MCP server locally:

### 1. Terminal 1: Start the MCP server

```
1 python amazon-bedrock-agentcore-workshop/01-AgentCore-runtime/02-hosting-MCP-server/mcp_server.py
```

### 2. Terminal 2: Run the test client

```
1 python amazon-bedrock-agentcore-workshop/01-AgentCore-runtime/02-hosting-MCP-server/my_mcp_client.py
```

The screenshot shows a code editor with two terminals. Terminal 1 is running the MCP server, and Terminal 2 is running the test client. The output in Terminal 2 shows the client making requests to the server, including a request for available tools. The server responds with a list of tools: add\_numbers, multiply\_numbers, and greet\_user.

```

jossai@842f577d057d 02-hosting-MCP-server % python my_mcp_client.py
INFO:httpx:HTTP Request: POST http://localhost:8000/mcp/ "HTTP/1.1 307 Temporary Redirect"
INFO:httpx:HTTP Request: POST http://localhost:8000/mcp/ "HTTP/1.1 200 OK"
INFO:mcp.client.streamable_http:Negotiated protocol version: 2025-06-18
INFO:httpx:HTTP Request: POST http://localhost:8000/mcp/ "HTTP/1.1 307 Temporary Redirect"
INFO:httpx:HTTP Request: POST http://localhost:8000/mcp/ "HTTP/1.1 202 Accepted"
INFO:httpx:HTTP Request: POST http://localhost:8000/mcp/ "HTTP/1.1 307 Temporary Redirect"
INFO:httpx:HTTP Request: POST http://localhost:8000/mcp/ "HTTP/1.1 200 OK"

Available tools:
- add_numbers: Add two numbers together
- multiply_numbers: Multiply two numbers together
- greet_user: Greet a user by name
jossai@842f577d057d 02-hosting-MCP-server %

```

You should see your three tools listed in the output.

## Setting up Amazon Cognito for Authentication

AgentCore Runtime requires authentication. We'll use Amazon Cognito to provide JWT tokens for accessing our deployed MCP server. We also modify the Python path so that modules in the parent

directory can be imported. It imports utility functions used to set up the Cognito user pool and IAM role

The code below is the helper function to create and configure a Cognito user pool. The result is stored in `cognito_config`, which contains values such as `user_pool_id`, `client_id`, and `discovery_url`.

#### ▼ Code Snippet

```
1 import sys
2 import os
3
4 # Get the current notebook's directory
5 current_dir = os.path.dirname(os.path.abspath('__file__' if '__file__' in globals() else '.'))
6
7 utils_dir = os.path.join(current_dir, '..')
8 utils_dir = os.path.abspath(utils_dir)
9
10 # Add to sys.path
11 sys.path.insert(0, utils_dir)
12
13 from utils import create_agentcore_role, setup_cognito_user_pool
```



```
1 cognito_config = setup_cognito_user_pool()
```



## Create IAM Execution role for the AgentCore Runtime

This code snippet creates an IAM role specifically for the MCP server. This role is required for deploying to AgentCore Runtime and is tied to the `tool_name`.

#### ▼ Code Snippet

```
1 tool_name = "mcp_server_ac"
2 agentcore_iam_role = create_agentcore_role(agent_name=tool_name)
```



## Configuring AgentCore Runtime Deployment

This sets up the AWS SDK session, determines the current region, and checks for the required files (`mcp_server.py`, `requirements.txt`). If files are missing, it raises an error.

#### ▼ Code Snippet

```
1 from bedrock_agentcore_starter_toolkit import Runtime
2 from boto3.session import Session
3 import time
4 import os
5
6 boto_session = Session()
7 region = boto_session.region_name
8
9 required_files = ['mcp_server.py', 'requirements.txt']
10 for file in required_files:
11     if not os.path.exists(file):
12         raise FileNotFoundError(f"Required file {file} not found")
```



Next, we initialize the `Runtime` instance and set up authentication configuration using Cognito. This ensures only valid tokens from the user pool can invoke the MCP server.

#### ▼ Code Snippet



```

1  agentcore_runtime = Runtime()
2
3  auth_config = {
4      "customJWTAuthorizer": {
5          "allowedClients": [
6              cognito_config['client_id']
7          ],
8          "discoveryUrl": cognito_config['discovery_url'],
9      }
10 }

```

Here, we configure the MCP runtime using this script as the entry point. It creates an ECR repository automatically and prepares the runtime to deploy the MCP server using the defined IAM role and authentication.

#### ▼ Code Snippet



```

1  response = agentcore_runtime.configure(
2      entrypoint="mcp_server.py",
3      execution_role=agentcore_iam_role['Role']['Arn'],
4      auto_create_ecr=True,
5      requirements_file="requirements.txt",
6      region=region,
7      authorizer_configuration=auth_config,
8      protocol="MCP",
9      agent_name=tool_name
10 )

```

#### 📘 Documentation Note

We use the provided [starter toolkit](#) to configure the AgentCore Runtime with an entrypoint, execution role, and requirements file. It will also auto-create the Amazon ECR repository during setup. A Dockerfile will be generated from your application code as part of this step.

## Launch MCP Server to AgentCore Runtime

This next script initializes a Bedrock AgentCore Runtime session using the Runtime class. It sets the AWS region, checks for the required source files, and defines a custom JWT authentication configuration using Amazon Cognito credentials. The `configure()` method sets up the runtime deployment with all necessary parameters—such as IAM role, Python entrypoint, and required libraries—preparing it for launch.

#### ▼ Code Snippet



```

1  from bedrock_agentcore_starter_toolkit import Runtime
2  from boto3.session import Session
3  import time
4  import os
5
6  boto_session = Session()
7  region = boto_session.region_name
8
9  required_files = ['mcp_server.py', 'requirements.txt']
10 for file in required_files:
11     if not os.path.exists(file):
12         raise FileNotFoundError(f"Required file {file} not found")
13
14 agentcore_runtime = Runtime()
15
16 auth_config = {
17     "customJWTAuthorizer": {
18         "allowedClients": [
19             cognito_config['client_id']
20         ],

```

```

21         "discoveryUrl": cognito_config['discovery_url'],
22     }
23 }
24
25 response = agentcore_runtime.configure(
26     entrypoint="mcp_server.py",
27     execution_role=agentcore_iam_role['Role']['Arn'],
28     auto_create_ecr=True,
29     requirements_file="requirements.txt",
30     region=region,
31     authorizer_configuration=auth_config,
32     protocol="MCP",
33     agent_name=tool_name
34 )

```

Now, we launch the configured runtime, deploying the MCP server to Amazon Bedrock AgentCore Runtime. The result includes identifiers like `agent_arn` and `agent_id`.

```
1 launch_result = agentcore_runtime.launch()
```



With the Dockerfile ready, we deploy the MCP server to AgentCore Runtime. This step creates both the Amazon ECR repository and the AgentCore Runtime environment. Let's wait for about 20 seconds for the AgentCore runtime status to be ready before we store the configurations for remote access.

Before we can invoke our deployed MCP server, let's store the Agent ARN and Cognito configuration in AWS Systems Manager Parameter Store and AWS Secrets Manager for easy retrieval:

#### ▼ Code Snippet

```

1 import boto3
2 import json
3
4 ssm_client = boto3.client('ssm', region_name=region)
5 secrets_client = boto3.client('secretsmanager', region_name=region)
6
7 try:
8     cognito_credentials_response = secrets_client.create_secret(
9         Name='mcp_server/cognito/credentials',
10        Description='Cognito credentials for MCP server',
11        SecretString=json.dumps(cognito_config)
12    )
13 except secrets_client.exceptions.ResourceExistsException:
14     secrets_client.update_secret(
15         SecretId='mcp_server/cognito/credentials',
16         SecretString=json.dumps(cognito_config)
17     )
18
19 agent_arn_response = ssm_client.put_parameter(
20     Name='/mcp_server/runtime/agent_arn',
21     Value=launch_result.agent_arn,
22     Type='String',
23     Description='Agent ARN for MCP server',
24     Overwrite=True
25 )

```



## Creating Remote Testing Client

This snippet will create a test client that will retrieve the necessary credentials from AWS and connect to the deployed server.

#### ▼ Code Snippet

```

1 %%writefile my_mcp_client_remote.py
2 import asyncio

```



```

3 import boto3
4 import json
5 import sys
6 from boto3.session import Session
7
8 from mcp import ClientSession
9 from mcp.client.streamable_http import streamablehttp_client
10
11 async def main():
12     boto_session = Session()
13     region = boto_session.region_name
14
15     print(f"Using AWS region: {region}")
16
17     try:
18         ssm_client = boto3.client('ssm', region_name=region)
19         agent_arn_response = ssm_client.get_parameter(Name='/mcp_server/runtime/agent_arn')
20         agent_arn = agent_arn_response['Parameter']['Value']
21         print(f"Retrieved Agent ARN: {agent_arn}")
22
23         secrets_client = boto3.client('secretsmanager', region_name=region)
24         response = secrets_client.get_secret_value(SecretId='mcp_server/cognito/credentials')
25         secret_value = response['SecretString']
26         parsed_secret = json.loads(secret_value)
27         bearer_token = parsed_secret['bearer_token']
28         print("✓ Retrieved bearer token from Secrets Manager")
29
30     except Exception as e:
31         print(f"Error retrieving credentials: {e}")
32         sys.exit(1)
33
34     if not agent_arn or not bearer_token:
35         print("Error: AGENT_ARN or BEARER_TOKEN not retrieved properly")
36         sys.exit(1)
37
38     encoded_arn = agent_arn.replace(':', '%3A').replace('/', '%2F')
39     mcp_url = f"https://bedrock-agentcore.{region}.amazonaws.com/runtimes/{encoded_arn}/invocation"
40     headers = {
41         "authorization": f"Bearer {bearer_token}",
42         "Content-Type": "application/json"
43     }
44
45     print(f"\nConnecting to: {mcp_url}")
46     print("Headers configured ✓")
47
48     try:
49         async with streamablehttp_client(mcp_url, headers, timeout=120, terminate_on_close=False) as client:
50             read_stream,
51             write_stream,
52             _ =
53         ):
54             async with ClientSession(read_stream, write_stream) as session:
55                 print("\n🔧 Initializing MCP session...")
56                 await session.initialize()
57                 print("✓ MCP session initialized")
58
59                 print("\n🔧 Listing available tools...")
60                 tool_result = await session.list_tools()
61
62                 print("\n📋 Available MCP Tools:")
63                 print("=" * 50)
64                 for tool in tool_result.tools:
65                     print(f"🔧 {tool.name}")
66                     print(f"   Description: {tool.description}")
67                     if hasattr(tool, 'inputSchema') and tool.inputSchema:
68                         properties = tool.inputSchema.get('properties', {})
69                         if properties:
70                             print(f"   Parameters: {list(properties.keys())}")
71                     print()
72
73                 print(f"✅ Successfully connected to MCP server!")
74                 print(f"Found {len(tool_result.tools)} tools available.")
75

```

```

76     except Exception as e:
77         print(f"✖ Error connecting to MCP server: {e}")
78         sys.exit(1)
79
80 if __name__ == "__main__":
81     asyncio.run(main())

```

## Testing

Here are a few ways to test your config:

### 1. Testing Your Deployed MCP Server

```

1  print("Testing deployed MCP server...")
2  print("=" * 50)
3  !python my_mcp_client_remote.py

```



### 2. Invoking MCP Tools Remotely

This is an advanced client that goes beyond listing tools. It will also invokes them to showcase the full capabilities of the MCP system.

#### ▼ Code Snippet

```

1  %%writefile invoke_mcp_tools.py
2  import asyncio
3  import boto3
4  import json
5  import sys
6  from boto3.session import Session
7
8  from mcp import ClientSession
9  from mcp.client.streamable_http import streamablehttp_client
10
11 async def main():
12     boto_session = Session()
13     region = boto_session.region_name
14
15     print(f"Using AWS region: {region}")
16
17     try:
18         ssm_client = boto3.client('ssm', region_name=region)
19         agent_arn_response = ssm_client.get_parameter(Name='/mcp_server/runtime/agent_arn')
20         agent_arn = agent_arn_response['Parameter']['Value']
21         print(f"Retrieved Agent ARN: {agent_arn}")
22
23         secrets_client = boto3.client('secretsmanager', region_name=region)
24         response = secrets_client.get_secret_value(SecretId='mcp_server/cognito/credentials')
25         secret_value = response['SecretString']
26         parsed_secret = json.loads(secret_value)
27         bearer_token = parsed_secret['bearer_token']
28         print("✓ Retrieved bearer token from Secrets Manager")
29
30     except Exception as e:
31         print(f"Error retrieving credentials: {e}")
32         sys.exit(1)
33
34     encoded_arn = agent_arn.replace(':', '%3A').replace('/', '%2F')
35     mcp_url = f"https://bedrock-agentcore.{region}.amazonaws.com/runtimes/{encoded_arn}/invocat
36     headers = {
37         "authorization": f"Bearer {bearer_token}",
38         "Content-Type": "application/json"
39     }
40
41     print(f"\nConnecting to: {mcp_url}")
42

```



```

43     try:
44         async with streamablehttp_client(mcp_url, headers, timeout=120, terminate_on_close=False):
45             read_stream,
46             write_stream,
47             _ =
48         ):
49         async with ClientSession(read_stream, write_stream) as session:
50             print("\n🔧 Initializing MCP session...")
51             await session.initialize()
52             print("✓ MCP session initialized")
53
54             print("\n🔧 Listing available tools...")
55             tool_result = await session.list_tools()
56
57             print("\n📋 Available MCP Tools:")
58             print("=" * 50)
59             for tool in tool_result.tools:
60                 print(f"🔧 {tool.name}: {tool.description}")
61
62             print("\n🧪 Testing MCP Tools:")
63             print("=" * 50)
64
65             try:
66                 print("\n➕ Testing add_numbers(5, 3)...")
67                 add_result = await session.call_tool(
68                     name="add_numbers",
69                     arguments={"a": 5, "b": 3}
70                 )
71                 print(f"📄 Result: {add_result.content[0].text}")
72             except Exception as e:
73                 print(f"🚫 Error: {e}")
74
75             try:
76                 print("\n✖ Testing multiply_numbers(4, 7)...")
77                 multiply_result = await session.call_tool(
78                     name="multiply_numbers",
79                     arguments={"a": 4, "b": 7}
80                 )
81                 print(f"📄 Result: {multiply_result.content[0].text}")
82             except Exception as e:
83                 print(f"🚫 Error: {e}")
84
85             try:
86                 print("\n👋 Testing greet_user('Alice')...")
87                 greet_result = await session.call_tool(
88                     name="greet_user",
89                     arguments={"name": "Alice"}
90                 )
91                 print(f"📄 Result: {greet_result.content[0].text}")
92             except Exception as e:
93                 print(f"🚫 Error: {e}")
94
95             print("\n✅ MCP tool testing completed!")
96
97         except Exception as e:
98             print(f"🚫 Error connecting to MCP server: {e}")
99             sys.exit(1)
100
101 if __name__ == "__main__":
102     asyncio.run(main())

```


## ➡ Try it out

### At an AWS Event

If you are following the workshop via workshop studio, go to JupyterLab in SageMaker Studio. In the JupyterLab UI navigate to [01-AgentCore-runtime/02-hosting-MCP-server](#).



## Self-paced

Here's a notebook that lets you try out the above: [hosting\\_mcp\\_server.ipynb](#) .