



Amazon Elastic Container
Service (Amazon ECS)



SERVERLESS CONTAINERIZED WEATHER APP ON AWS ECS

Harvey Aland
March 28th, 2025

Introduction

In today's fast-paced technology landscape, the ability to build, deploy, and manage applications efficiently is crucial. This project demonstrates the end-to-end process of developing a containerized weather application using **Docker**, deploying it to **Amazon Elastic Container Registry (ECR)**, and running it on **Amazon Elastic Container Service (ECS)**. The infrastructure is provisioned using Terraform, showcasing the power of Infrastructure as Code (IaC).

This document provides a detailed walkthrough of the project, including the code, architecture, and deployment process. It also highlights the relevance of this project to both my personal growth and the broader industry.

Overview of the Project

The project involves building a simple weather application using Python and Flask, containerizing it with Docker, and deploying it to AWS using ECS and ECR. The infrastructure is provisioned using Terraform, which automates the creation of resources such as VPC, subnets, ECS clusters, and an Application Load Balancer (ALB).

Key Components

1. Application:

- A Python-based weather app that retrieves and displays weather data.
- Built using Flask for the backend and containerized using Docker.

2. Infrastructure:

- VPC: A Virtual Private Cloud with public and private subnets for secure networking.
- ECS: Elastic Container Service for running the containerized app.
- ECR: Elastic Container Registry for storing Docker images.
- ALB: Application Load Balancer for routing traffic to the ECS tasks.

3. Tools and Technologies:

- Docker: For containerizing the application.
- Terraform: For provisioning and managing AWS infrastructure.
- AWS CLI: For interacting with AWS services.
- GitHub: For version control and collaboration
- Python: For basic flask app
- HTML: Written front end content for app

Why This Project is Relevant

1. Skill Development:

- This project allowed me to gain hands-on experience with Docker, AWS ECR, ECS, and Terraform.
- I learned how to containerize an application and manage cloud infrastructure.

2. Portfolio Building:

- This project serves as a tangible demonstration of my ability to design, build, and deploy a modern cloud-native application.

3. Containerization:

- Containerization is a cornerstone of modern software development, enabling applications to run consistently across different environments. This project showcases my ability to leverage Docker for containerization.
-

4. Infrastructure as Code (IaC):

- Terraform is widely used for automating infrastructure provisioning. This project highlights my ability to use Terraform to manage cloud resources efficiently.

Flask application build

While researching what application I could create as a beginner to showcase containerizing an application in docker I found myself looking at different small apps like a simple message response or a shopping list. However, neither of which I felt had any real-world value.

I wanted to build something that I thought for the most part was cool and offered some form of practicality. So, I settled on a weather app since I rely on them daily, especially living in the unpredictable Northern English climate.

To get weather data, I used OpenWeatherMap.com, which provides a free API for region-specific forecasts. I went with London for my region of choice.

Firstly, I had to confirm that the API is working and outputting the data before I could build the app.

Checking API URL



Building the app using Python Flask

what is flask? - Flask is a lightweight and flexible web framework for python, it allowed me to build my web application quickly without needing a lot of boilerplate code. This was a viable option for me as I have experience using python and it was relatively straight forward for what I was building.

How Does Flask Work? - Flask runs a webserver that listens for requests from users, you can define routes (URLs) in python to determine what should happen when someone visits a specific page. For my project I chose flask to render a HTML page

Writing the Python Code

The Python script (app.py) does the following:

1. Initializes a Flask application.
2. Fetches real-time weather data from OpenWeatherMap for a specific city.
3. Defines routes for the application:
 - `/` → Displays the weather on the webpage.
 - `/health` → Used for health checks by AWS ALB.
4. Runs the app on port 5000 so it can be accessed inside the container.

The Python Flask code is annotated with notes which tells you what each line does. This script took a while to get right as it was my first time creating a flask app, I had to reference YouTube tutorials and a little help from AI too to get it right.

```

weather_app.py > health
1 import requests # This is like telling Python, "I want to use a tool to ask for information from the internet."
2 from flask import Flask, render_template # Telling Python that we are going to make a web app with Flask and sh
3
4 app = Flask(__name__) # This is creating a "Flask" app. Imagine it as a box where our web app lives.
5
6 API_KEY = 'd552c10073f57e64af138603e24430e7' # This is like your secret password to access the weather data.
7 CITY = 'London' # We want the weather for London.
8 BASE_URL = f'http://api.openweathermap.org/data/2.5/weather?q={CITY}&appid={API_KEY}&units=metric' # This is the
9
10 #Health check
11 @app.route('/health')
12 def health():
13     return "OK", 200 # ALB will receive a 200 status, meaning the app is healthy.
14
15
16 @app.route('/') # This means "When someone goes to our website's home page, do the following."
17 def index(): # This is the function that runs when someone visits the page.
18     response = requests.get(BASE_URL) # We ask the weather website for the data.
19     data = response.json() # We take the data the website gives us and change it into a format we can use.
20
21     # Check if everything went fine with the data
22     if data["cod"] == 200: # If the website says it's all good (code 200 means everything worked)
23         temperature = data["main"]["temp"] # We get the temperature from the data.
24         description = data["weather"][0]["description"] # We get the weather description (like "sunny" or "cloud
25         city_name = data["name"] # Extract the city name from the API response
26     else:
27         temperature = None # If something went wrong, we say there's no temperature.
28         description = "Error fetching data" # We tell the user there was an error.
29
30     return render_template('index.html', temperature=temperature, description=description, city=city_name) # Sh
31
32 app.run(host="0.0.0.0", port=5000, debug=True)

```

After the flask code was done, I had to create the requirements.txt file since the flask app depends on certain python libraries.

```

requirements.txt
1 flask # Runs web application
2 requests # Allows your app to fetch weather data from the OpenWeatherMap API
3 gunicorn # Provides a more robust server for running Flask in Docker (instead of the built-in development server).

```

HTML template for the app frontend

Next step was to create a frontend page that users of the app would see when accessing the app.

Seeing as I've used HTML before this part was relatively easy, upon researching how HTML and Python flask can be used together I found a template that would be best suited for my project and did some tweaks. Since I only needed a simple layout, this template worked well with minimal modifications.

The index.html file is a Jinja template used by Flask to dynamically display weather data fetched from OpenWeatherMap. When a user visits the app, Flask retrieves weather information and passes it to the template using the render_template() function. The placeholders inside {{ }} are replaced with real values, such as the city name, temperature, and weather condition. This ensures that each time the page loads, it reflects the latest weather data. The template also includes simple CSS styling to format the output neatly, making the app visually appealing and easy to read.

Flask App Code that links with Jinja HTML template

```
# Check if everything went fine with the data
if data["cod"] == 200: # If the website says it's all good (code 200 means everything worked)
    temperature = data["main"]["temp"] # We get the temperature from the data.
    description = data["weather"][0]["description"] # We get the weather description (like "sunny" or "cloudy").
    city_name = data["name"] # Extract the city name from the API response
else:
    temperature = None # If something went wrong, we say there's no temperature.
    description = "Error fetching data" # We tell the user there was an error.

return render_template('index.html', temperature=temperature, description=description, city=city_name) # Show the webpage with the temperature
```

HTML Code that links with the Flask App

```
<body>
  <h1>Weather in {{ city }}</h1>
  <!-- Displays the city name from the Flask app -->

  <div class="weather-box">
    <!-- This box contains the weather details -->

    <p><strong>Temperature:</strong> {{ temperature }}°C</p>
    <!-- Displays the temperature from the API -->

    <p><strong>Condition:</strong> {{ description }}</p>
    <!-- Displays the weather condition (e.g., "Cloudy") -->
  </div>
</body>

</html>
```

Later, when I build and test the app, you'll see that the Python Flask and HTML code work as expected with Jinja ({{ }}) replacing the placeholders with the weather data.

Building the flask app

Now the code and files were in order it was time to build the app and see if it was working.

1st Step: Installing required Python Packages

```
ModuleNotFoundError: No module named 'requests'
PS C:\Users\Harve\desktop\ECS> pip install requests
Collecting requests
  Downloading requests-2.32.3-py3-none-any.whl.metadata (4.6 kB)
Collecting charset-normalizer<4,>=2 (from requests)
  Downloading charset_normalizer-3.4.1-cp313-cp313-win_amd64.whl.metadata (36 kB)
Collecting idna<4,>=2.5 (from requests)
  Downloading idna-3.10-py3-none-any.whl.metadata (10 kB)
Collecting urllib3<3,>=1.21.1 (from requests)
  Downloading urllib3-2.3.0-py3-none-any.whl.metadata (6.5 kB)
Collecting certifi>=2017.4.17 (from requests)
  Downloading certifi-2025.1.31-py3-none-any.whl.metadata (2.5 kB)
Downloading requests-2.32.3-py3-none-any.whl (64 kB)
Downloading certifi-2025.1.31-py3-none-any.whl (166 kB)
Downloading charset_normalizer-3.4.1-cp313-cp313-win_amd64.whl (102 kB)
Downloading idna-3.10-py3-none-any.whl (70 kB)
Downloading urllib3-2.3.0-py3-none-any.whl (128 kB)
Installing collected packages: urllib3, idna, charset-normalizer, certifi, requests
Successfully installed certifi-2025.1.31 charset-normalizer-3.4.1 idna-3.10 requests-2.32.3 urllib3-2.3.0
[notice] A new release of pip is available: 24.3.1 -> 25.0.1
[notice] To update, run: python.exe -m pip install --upgrade pip
PS C:\Users\Harve\desktop\ECS>
```

2nd Step: Running the app locally by executing the python script to start the flask web server

Ran by using the following PS command: `python app.py`

3rd Step: Checking the app was running on port 5000

Using <http://127.0.0.1:5000> (Localhost)

127.0.0.1:5000

Weather in London

Temperature: 13.39°C

Condition: overcast clouds

Confirmed now that the app is running locally, the data is being retrieved correctly, and the HTML/CSS styles are applied as expected, ensuring the app's interface is visually consistent.

Containerize the Flask App using Docker

After a lengthy process of installing docker I was finally ready to start the main part of this project. The first step in doing so was to create a DockerFile.

What is a DockerFile?

- A **Dockerfile** is a text file script that contains instructions for Docker to build a container image for an application.
- It defines the base image, copies necessary files, installs dependencies, and specifies the command to run the app inside the container.
- For example, a Dockerfile for a Flask app installs Flask, copies files, and runs the app on port 5000.

After researching and watching Docker tutorials, I created a simple Dockerfile to define the environment and instruct Docker to build the base image. This file specifies the base image, copies the necessary application files, installs dependencies, and sets the command to run the app inside the container.

```
Dockerfile > ...
1  # Use an official Python runtime as a base image
2  FROM python:3.8-slim
3
4  # Set the working directory inside the container
5  WORKDIR /app
6
7  # Copy the current directory contents (your Flask app) into the container's working directory
8  COPY . /app
9
10 # Install the required dependencies inside the container
11 RUN pip install -r requirements.txt
12
13 # Expose port 5000 for the Flask app
14 EXPOSE 5000
15
16 # Command to run your app
17 CMD ["python", "weather_app.py"]
```


Building the Docker container

After creating the Dockerfile, I ran the `docker build -t flask-weather-app .` command in PowerShell to build the Docker image. This command instructed Docker to process the Dockerfile and create an image named flask-weather-app, which could then be used to run containers.

```
PS C:\Users\Harve\Desktop\ECS> docker build -t flask-weather-app .
[+] Building 12.3s (10/10) FINISHED
=> [internal] load build definition from Dockerfile
=> => transferring dockerfile: 498B
=> [internal] load metadata for docker.io/library/python:3.8-slim
=> [auth] library/python:pull token for registry-1.docker.io
=> [internal] load .dockerignore
=> => transferring context: 2B
=> [1/4] FROM docker.io/library/python:3.8-slim@sha256:1d52838af602b4b5a831beb13a0e4d073280665ea7be7f69ce2382f29c5a613f
=> => resolve docker.io/library/python:3.8-slim@sha256:1d52838af602b4b5a831beb13a0e4d073280665ea7be7f69ce2382f29c5a613f
=> => sha256:3971691a363796c39467aae4cdce6ef773273fe6bfc67154d01e1b589befb912 248B / 248B
=> => sha256:a3f1dfe736c5f959143f23d75ab522a60be2da902efac236f4fb2a153cc14a5d 14.53MB / 14.53MB
=> => sha256:030d7bdc20a63e3d22192b292d006a69fa3333949f536d62865d1bd0506685cc 3.51MB / 3.51MB
=> => sha256:302e3ee498053a7b5332ac79e8efebec16e900289fc1ecd1c754ce8fa047fcab 29.13MB / 29.13MB
=> => extracting sha256:302e3ee498053a7b5332ac79e8efebec16e900289fc1ecd1c754ce8fa047fcab
=> => extracting sha256:030d7bdc20a63e3d22192b292d006a69fa3333949f536d62865d1bd0506685cc
=> => extracting sha256:a3f1dfe736c5f959143f23d75ab522a60be2da902efac236f4fb2a153cc14a5d
=> => extracting sha256:3971691a363796c39467aae4cdce6ef773273fe6bfc67154d01e1b589befb912
=> [internal] load build context
=> => transferring context: 6.71kB
=> [2/4] WORKDIR /app
=> [3/4] COPY . /app
=> [4/4] RUN pip install -r requirements.txt
=> exporting to image
=> => exporting layers
=> => exporting manifest sha256:aa32824b21ca68b6259a98b9eaae359b46f320b4e57a4f66ec03be868cf999aa
=> => exporting config sha256:d82fb486b28c9f4e334f9047af3b9da89ce1c3629fe2694d9ac757ece1c161c2
=> => exporting attestation manifest sha256:73d3e845f85093e5ac463e0cd2c1a385689e795d2d9c13ed536b3132e436bc8c
=> => exporting manifest list sha256:5a856f1adc6fd37413609bc6ab0cc707686904d4f8b4b2f8e67f95f6erd38149
```

Running the Flask App inside the container

I used the `docker run` command to start the container and the `docker ps` command to verify that it was running. This confirmed that the application was successfully deployed inside the Docker container.

```
PS C:\Users\Harve\Desktop\ECS> docker run -d -p 5000:5000 --name weather-app flask-weather-app
529a003ea31ec07556e290fd447465737522aa60ad7da138b06ae301ddc41e8b
PS C:\Users\Harve\Desktop\ECS> docker ps
CONTAINER ID   IMAGE          COMMAND                  CREATED        STATUS        PORTS                               NAMES
529a003ea31e   flask-weather-app   "python weather_app...."   22 seconds ago   Up 21 seconds   0.0.0.0:5000->5000/tcp        weather-app
PS C:\Users\Harve\Desktop\ECS>
```

Image in Docker Desktop

Images [Give feedback](#)
View and manage your local and Docker Hub images. [Learn more](#)

Local

Hub repositories

161.91 MB / 224.01 MB in use

1 Images

Last refresh: 2 hours ago

Q Search

<input type="checkbox"/>	Name	Tag	Image ID	Created	Size	Actions
<input type="checkbox"/>	<div><div></div>flask-weather-app</div>	latest	c28fa1c8b646	2 hours ago	215.86 MB	<div><div></div><div></div><div></div></div>

Verifying app is running inside container

Once I verified that the container was running and the image was created, I checked the URL again to confirm that the app was successfully running inside the container and functioning as expected.

localhost:5000

Weather in London

Temperature: 6.04°C

Condition: few clouds

AWS Infrastructure deployed and managed via Terraform IAC

After successfully containerizing the Flask app with Docker, I used Terraform to define and deploy the necessary cloud infrastructure. The Terraform code automated the provisioning of resources.

AWS Infrastructure Breakdown

1. Networking

- Amazon VPC (Virtual Private Cloud):
 - A virtual network (aws_vpc) is created with a CIDR block of 10.0.0.0/16.
- Subnets:
 - Public and private subnets are created across two availability zones (eu-west-2a and eu-west-2b).
- Internet Gateway (aws_internet_gateway):
 - Attached to the VPC to enable internet access for public subnets.
- Route Tables:
 - Public and private route tables are configured to manage traffic routing.
- NAT Gateway (aws_nat_gateway):
 - Deployed in a public subnet to allow private subnets to access the internet securely.
- Elastic IP (EIP) (aws_eip):
 - Allocated and associated with the NAT Gateway to provide a public IP for outbound internet access.

2. Security

- Security Groups (aws_security_group):
 - ALB Security Group:
 - Allows inbound HTTP (port 80) and HTTPS (port 443) traffic.
 - Allows all outbound traffic.
 - ECS Security Group:
 - Restricts inbound traffic to only allow communication from the ALB on port 5000.
 - Allows all outbound traffic.

3. Load Balancing

- Application Load Balancer (ALB) (aws_lb):
 - Public-facing ALB deployed in public subnets.
 - Listens on port 80 (aws_lb_listener) and forwards traffic to the ECS tasks.
- Target Group (aws_lb_target_group):
 - Routes traffic to ECS tasks running on port 5000.
 - ALB performs health checks on the ECS tasks to ensure they are healthy and able to handle traffic.

4. Container Services

- Amazon ECS (Elastic Container Service):
 - ECS Cluster (aws_ecs_cluster):
 - Logical grouping of ECS tasks and services.
 - ECS Task Definition (aws_ecs_task_definition):
 - Defines the container image, CPU, memory, and networking requirements.
 - Uses AWS Fargate for serverless container execution.
 - ECS Service (aws_ecs_service):
 - Manages the desired number of ECS tasks (2 in this case).
 - Deploys tasks in private subnets with no public IPs.
 - Associates tasks with the ALB target group.

5. Container Registry

- Amazon ECR (Elastic Container Registry) (aws_ecr_repository):
 - Private Docker image repository for storing the application container image.
 - Image scanning is enabled to detect vulnerabilities on push.

6. IAM (Identity and Access Management)

- IAM Roles:
 - ECS Task Execution Role:
 - Grants permissions to pull images from ECR and write logs to CloudWatch.
 - ECS Task Role:
 - Grants permissions for ECS tasks to access other AWS services (e.g., S3, CloudWatch).
 - ECR Push Role:
 - Grants permissions to push Docker images to the ECR repository.
- IAM Policies:
 - Custom policies are attached to roles to define specific permissions (e.g., pushing images to ECR).
 -

I have included an un-modularized version of the terraform code which has been annotated for the sake of reviewing, please see attached file in the GIT repository.

Pushing the Docker Image to ECR

After defining and deploying the TF code I checked to see if ECS and ECR were created and running in the console, after verifying I then pushed the Docker Image to ECR.

The whole process of pushing the Docker image initially stumped me, as each use case can be different. I came across information online suggesting that shell scripts are often used to automate the Docker login process, because the login session typically expires after 12 hours. This means that if you need to make any updates to the app, you will have to go through the entire login process again.

However, since my app was intended as a one-time deployment for the purpose of showcasing my working experience, I decided to manually push the image. The infrastructure was destroyed after the project was completed, so it didn't make sense to automate a process for an app that would only be deployed once. The need for

automating the login and push process was unnecessary in this case, as I would not be frequently pushing new images or updating the app.

Logging in to docker via AWS CLI

```
PS C:\Windows\system32> aws ecr get-login-password --region eu-west-2 | docker login --username AWS --password-stdin [REDACTED].dkr.ecr.eu-west-2.amazonaws.com
Login Succeeded
```

Following the commands instructed in the AWS console I pushed the image to ECR

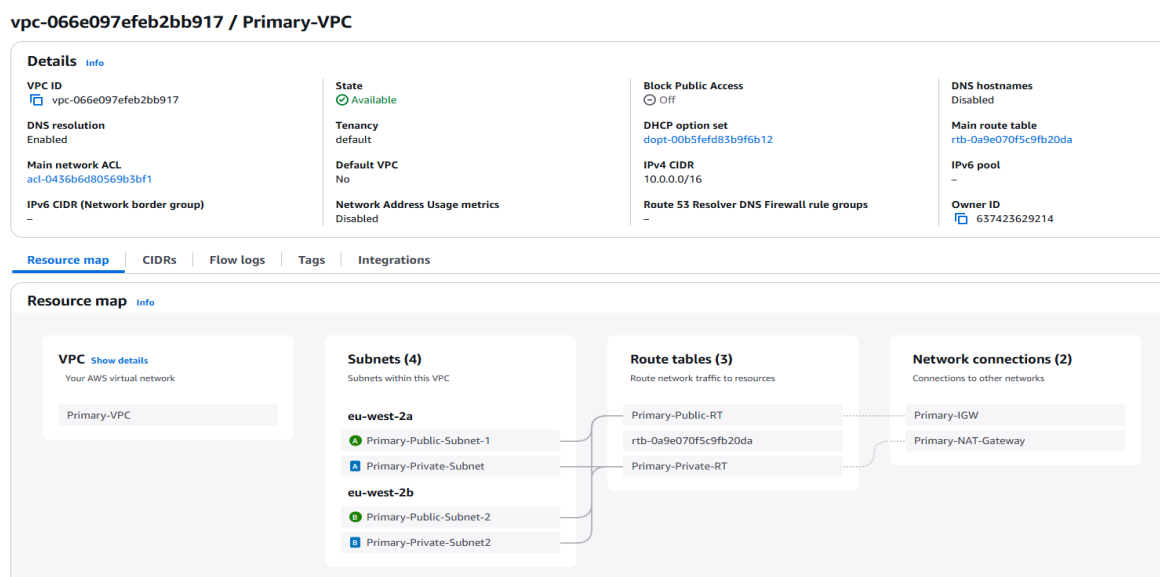
```
PS C:\Windows\system32> docker tag flask-weather-app:latest [REDACTED].dkr.ecr.eu-west-2.amazonaws.com/weather-repo:latest
PS C:\Windows\system32> docker push [REDACTED].dkr.ecr.eu-west-2.amazonaws.com/weather-repo:latest
The push refers to repository [REDACTED].dkr.ecr.eu-west-2.amazonaws.com/weather-repo]
a3f1dfe736c5: Pushed
6897804f8b41: Pushed
3971691a3637: Pushed
302e3ee49805: Pushed
030d7bdc20a6: Pushed
9eec49d4cddc: Pushed
dc76dc40244a: Pushed
db8a87431f0b: Pushed
latest: digest: sha256:c28fa1c8b646da8135d42618adadc6a02cc838ebfce3041d319d07da8e3d1b01 size: 856
PS C:\Windows\system32>
```

Deployment checks in AWS console

After pushing the docker image to ECR the only task left to do for this project was to do a final test of the app and see if all the infrastructure had been created as expected.

Networking

Here we have the **VPC** map showcasing the network topology and how the resources link together.



Load Balancing

Confirming the ALB listens on port 80 for HTTP traffic and routes it to ECS tasks running on port 5000 via a target group. The target group performs health checks on the ECS tasks to ensure only healthy tasks receive traffic.

ecs-target-group

Details

arn:aws:elasticloadbalancing:eu-west-2:637423629214:targetgroup/ecs-target-group/4fe6ba0daab3b69f

Target type IP	Protocol : Port HTTP: 5000	Protocol version HTTP1	VPC vpc-066e097efeb2bb917
IP address type IPv4	Load balancer ecs-alb		

2 Total targets	2 Healthy 0 Anomalous	0 Unhealthy	0 Unused	0 Initial	0 Draining
--------------------	-----------------------------	----------------	-------------	--------------	---------------

Distribution of targets by Availability Zone (AZ)

Select values in this table to see corresponding filters applied to the Registered targets table below.

Targets | Monitoring | Health checks | Attributes | Tags

Registered targets (2) [Info](#)

Anomaly mitigation: **Not applicable** [Deregister](#) [Register](#)

Target groups route requests to individual registered targets using the protocol and port number specified. Health checks are performed on all registered targets according to the target group's health check settings. Anomaly detection is automatically applied to HTTP/HTTPS target groups with at least 3 healthy targets.

Filter targets

<input type="checkbox"/>	IP address	Port	Zone	Health status	Health status details	Administrative override	Overri...	Anomaly detec
<input type="checkbox"/>	10.0.3.67	5000	eu-west-2b (eu...)	Healthy	-	No override	No overri...	Normal
<input type="checkbox"/>	10.0.1.25	5000	eu-west-2a (eu...)	Healthy	-	No override	No overri...	Normal

Checking ALB DNS

To confirm the ALB is functioning correctly and balancing the load for the application, the ALB DNS name is used for verification. Since this project does not have a real domain name, the ALB DNS name acts as a substitute to validate that the ALB is routing traffic as expected.

ecs-alb-1082124794.eu-west-2.elb.amazonaws.com

Weather in London

Temperature: 4.57°C

Condition: clear sky

ECS Cluster, ECS tasks provisioned & FARGATE Launch Type.

ecs-cluster

Last updated
March 15, 2025 at 20:09 (UTC)

Update clusterDelete cluster

Cluster overview

ARN
am:aws:ecs:eu-west-2:637423629214:cluster/ecs-cluster

Status
Active

CloudWatch monitoring
Default

Registered container instances
-

Services
Draining
-

Active
1

Tasks
Pending
-

Running
2

Services

Tasks

Infrastructure

Metrics

Scheduled tasks

Configuration

Tags

Services (1) Info

Manage tagsUpdateDelete serviceCreate

Filter services by value

Filter launch type
Any launch type

Filter service type
Any service type

< 1 >

Service nameARNStatusService...Deployments and tasksLast deploy...Task de...Launch type

ecs-service

am:aws:ecs:eu-west-2:637423629214:service

Active

REPLICA

2/2 Tasks running

CompletedView

ecs-task:1

FARGATE

ECR Repo

Private repositories (1)

View push commandsDeleteActionsCreate repository

Search by repository substring

Repository name

URI

Created at

Tag immutability

Encryption type

weather-repo

637423629214.dkr.ecr.eu-west-2.amazonaws.com/weather-repo

March 14, 2025, 17:42:46 (UTC-00)

Mutable

AES-256

Confirming Docker image has been pushed to ECR Repo

Image scan overview, status, and full vulnerabilities has moved to the image detail page. To access, click an image tag.

Images (3)

DeleteDetailsScanView push commands

Search artifacts

< 1 >

Image tag

Artifact type

Pushed at

Size (MB)

Image URI

Digest

Last recorded pull time

latest

Image Index

March 15, 2025, 14:56:22 (UTC-00)

53.94

Copy URI

sha256:c28fa1c8b646da8135d42618adadc6...

March 15, 2025, 15:06:07 (UTC-00)

-

Image

March 15, 2025, 14:56:21 (UTC-00)

53.94

Copy URI

sha256:3943af4d49cb5a642344eb04676a7f...

March 15, 2025, 14:56:22 (UTC-00)

-

Image

March 15, 2025, 14:56:21 (UTC-00)

0.00

Copy URI

sha256:17b468c59d0773db24e98c687b081...

March 15, 2025, 14:56:28 (UTC-00)

Health check set up for target group to use ECS running tasks

In the original deployment of this infrastructure, health checks were not configured for the ECS tasks, meaning the load balancer had no way to verify the health or availability of the tasks. Leveraging the skills I gained from earning the AWS Solutions Architect Associate (SAA) certification, I implemented health checks to add redundancy and improve reliability. This ensures the load balancer routes traffic only to healthy tasks, enhancing the overall resilience of the application.

Adding in the ALB Target Group Health Check

Before

```
# Target Group for ECS Service
resource "aws_lb_target_group" "alb_target_group" {
  name        = "ecs-target-group"
  port        = 80
  protocol    = "HTTP"
  target_type = "ip" # For Fargate, use IP target type

  vpc_id = aws_vpc.ecs_vpc.id
}
```

After

```
# Target Group for ECS Service
resource "aws_lb_target_group" "alb_target_group" {
  name        = "ecs-target-group"
  port        = 5000
  protocol    = "HTTP"
  target_type = "ip" # For Fargate, use IP target type
  vpc_id      = aws_vpc.Primary_vpc.id

  health_check {
    protocol          = "HTTP"
    path              = "/health" # Health check URL for ECS tasks
    interval          = 30
    timeout            = 5
    healthy_threshold = 3
    unhealthy_threshold = 3
  }

  tags = {
    Name        = "ECS Target Group"
    Environment = "production"
  }
}
```

The health check ensures ECS tasks are healthy by sending HTTP requests to the **/health** endpoint every 30 seconds. A task is marked healthy after 3 successful responses (200 OK) or unhealthy after 3 failures. The target group port is 5000 because the ECS tasks run the application on this port, ensuring traffic is forwarded to the correct endpoint.

Adding in the /health route in the flask app

Once the ALB target group (TG) code was updated to include health checks, I added a /health route to the Flask app. This route allows the ALB to send HTTP requests to the /health endpoint and receive a 200 OK response when the application is healthy. This ensures the ALB can properly monitor the ECS tasks and route traffic only to those that are functioning correctly.

```
#Health check
@app.route('/health')
def health():
    return "OK", 200 # ALB will receive a 200 status, meaning the app is healthy.
```

Modifying the ECS Task Definition (to reflect the health check path)

Exposes port 5000 where my Flask App (with /health endpoint) runs.

```
#ECS task Definition
resource "aws_ecs_task_definition" "ecs_task" {
  family              = "ecs-task"
  requires_compatibilities = ["FARGATE"]
  cpu                 = "256"
  memory              = "512"
  network_mode        = "awsvpc"

  container_definitions = jsonencode([
    {
      name       = "my-container"
      image      = "${aws_ecr_repository.weather-app-repo.repository_url}:latest"
      memory     = 512
      cpu        = 256
      essential  = true
      portMappings = [
        {
          containerPort = 5000
          hostPort       = 5000
          protocol       = "tcp"
        }
      ]
    }
  ])
  execution_role_arn = aws_iam_role.ecs_task_execution_role.arn
}
```

ECS service integration

Binds the ALB to my containers port with automatic registration/deregistration of tasks.

```
#ECS using Fargate
resource "aws_ecs_service" "ecs_service" {
  name           = "ecs-service"
  cluster        = aws_ecs_cluster.ecs_cluster.id
  task_definition = aws_ecs_task_definition.ecs_task.arn
  desired_count  = 2 # desired count of running tasks

  launch_type = "FARGATE"

  network_configuration {
    subnets      = [aws_subnet.Primary_sub_priv.id, aws_subnet.Primary_sub_priv2.id]
    assign_public_ip = false # Fargate tasks don't require public IPs
    security_groups = [aws_security_group.ecs_sg.id]
  }

  load_balancer {
    target_group_arn = aws_lb_target_group.alb_target_group.arn
    container_name   = "my-container"
    container_port   = 5000
  }

  depends_on = [
    aws_lb_listener.ecs_listener
  ]
}
```

Security group update to allow ALB to ECS communication over port 5000

```
# Security Group for ECS tasks (restricted to only allow traffic from ALB)
resource "aws_security_group" "ecs_sg" {
  name           = "ecs-sg"
  description    = "Allow traffic from ALB only"
  vpc_id         = aws_vpc.Primary_vpc.id

  # Allow inbound traffic from the ALB security group on HTTP/HTTPS ports
  ingress {
    from_port = 5000
    to_port   = 5000
    protocol  = "tcp"
    security_groups = [aws_security_group.alb_sg.id] # Only allow traffic from ALB
  }

  # Allow outbound traffic (usually all traffic)
  egress {
    from_port = 0
    to_port   = 0
    protocol  = "-1" # All traffic
    cidr_blocks = ["0.0.0.0/0"]
  }
}
```

The ECS security group only allows inbound traffic on port 5000 from the ALB's security group, ensuring tasks can't be accessed directly from the internet. Outbound traffic is unrestricted to let tasks pull images and send logs. This creates a secure one-way flow:

internet traffic must pass through the ALB first, protecting the containers while maintaining functionality.

Full Workflow:

1. The ALB sends health checks (GET /health) to each ECS task's private IP on port 5000. Your Flask app returns 200 OK response
2. ALB marks task healthy after 3 good responses
3. Unhealthy tasks get terminated by ECS
4. Service launches replacement tasks automatically

Considerations and improvements

HTTPS Listener

One aspect not implemented in this deployment but essential for a fully functioning, production-ready application is an HTTPS listener for the Application Load Balancer (ALB).

For this project, HTTP (port 80) was intentionally used instead of HTTPS to simplify development and testing. Since this is a learning exercise without production traffic or sensitive data, HTTP reduces complexity by:

1. Avoiding certificate management – No need to provision and validate ACM certificates.
2. Simplifying debugging – HTTP allows easier inspection of traffic and logs without encryption overhead.
3. Reducing costs – No expenses for SSL/TLS certificates or domain registration.

While HTTPS is critical for real-world application ensuring security, SEO benefits, and compliance with browser requirement. This project prioritizes core AWS infrastructure components, such as ALB, ECS, and health checks.

Future HTTPS Implementation

To implement HTTPS in a production environment, the following steps would be required:

1. Request an ACM certificate – Use AWS Certificate Manager (ACM) to generate an SSL/TLS certificate for the domain.

2. Add an HTTPS listener – Configure the ALB to listen on port 443 and associate it with the ACM certificate.
3. Redirect HTTP to HTTPS – Enforce secure traffic by setting up automatic redirection from HTTP to HTTPS.
4. The `alb_sg` Security group is already configured for HTTP & HTTPS, so no changes needed here.

Open Egress security Group Rules

1. ALB Security Group (`alb_sg`)
 - The ALB does not *initiate* outbound connections; it only receives incoming traffic and forwards it to ECS tasks.
 - Open egress (0.0.0.0/0) ensures that the ALB can send health check requests to backend services or interact with other AWS resources if needed (e.g., logging to CloudWatch, redirecting traffic).
2. ECS Security Group (`ecs_sg`)
 - The ECS tasks (containers) must be able to communicate externally for various reasons, such as:
 - Fetching dependencies or updates from the internet (if required).
 - Communicating with other AWS services (S3, DynamoDB, RDS, or third-party APIs).
 - Sending responses back to clients via the ALB.
 - If outbound traffic were restricted, the ECS tasks might not be able to access necessary services (e.g., an RDS database, external APIs, or even basic networking like DNS resolution).

NAT Gateway redundancy

My thought process was if both NAT Gateways are live and running, they will be incurring a lot of extra unnecessary cost which for me building this project is a main concern.

If I were to implement a dual NAT GATEWAY architecture, there would be two main points to consider. Cost saving and RPO recovery time.

How it could be done?

Having built a failover-based architecture in a previous project. I began thinking about how utilizing cloud watch alarms and lambda functions could achieve this NAT Gateway failover solution.

Cost-Saving Approach (Single NAT with Failover)

Using one active NAT Gateway with CloudWatch and Lambda failover minimizes costs by ensuring only one NAT processes traffic at a time. If it fails, CloudWatch triggers a Lambda function to update route tables, redirecting traffic to the backup NAT. This approach reduces hourly and data transfer costs while maintaining availability, making it ideal for budget-conscious environments.

High Availability Approach (Dual NAT with Static Routing)

A dual NAT setup with static routing ensures instant failover by distributing private subnets across two NATs. If one fails, the other continues handling traffic without delays. While this setup incurs higher costs, it eliminates downtime and is ideal for mission-critical workloads that require uninterrupted internet access.

I suppose it depends on what is more important for a company standpoint, I can see that the best probable solution would be to accept the extra cost of having two active NAT Gateways running and focus on providing the best possible service having minimum downtime which in turn saves cost and probably outweighs the cost of the two NAT Gateways. The golden question here would be how much is 2-3 minutes of your application being down going to cost you? It's a trade off to consider.

Project Summary

This project was a great experience for me in terms of building a wider range of working knowledge with non-AWS services like Docker. The idea of containerization was always interesting to me and having an opportunity to implement it into something I have built was rewarding, especially sitting back and watching it all come together during the testing period. This project has further motivated me to look outside the realm of AWS cloud computing and look into more diverse deployments in the future.

Through this project, I gained the confidence and practical experience needed to manage Docker deployments effectively. I also developed a deeper understanding of how Docker integrates with AWS services like ECS and ECR, and how the overall container workflow operates. This experience has equipped me with the skills and knowledge to tackle similar Docker-related projects confidently in the future.