

Task8

Task 8: AI Model Deployment & MLOps

We have used ECS to host the Docker task and added an ALB (Application Load Balancer) to the publicly exposed Docker service. This is all achieved using a CloudFormation template and ECS Fargate.

The Docker service takes a voice sound file and converts it into text.

The famous Docker image used for this is `onerahmet/openai-whisper-asr-webservice:latest`.

Dockerfile & Kubernetes YAML Files

CloudFormation File

```
AWSTemplateFormatVersion: '2010-09-09'

Description: Deploy Whisper ASR API to ECS Fargate
Parameters:
  WhisperModel:
    Type: String
    Default: tiny
    AllowedValues: [tiny, base, small, medium, large]
    Description: Whisper model to use

Resources:
  WhisperVPC:
    Type: AWS::EC2::VPC
    Properties:
      CidrBlock: 10.0.0.0/16
      EnableDnsSupport: true
      EnableDnsHostnames: true
      Tags: [{ Key: Name, Value: WhisperVPC }]

  WhisperSubnet1:
```

Type: AWS::EC2::Subnet

Properties:

VpcId: !Ref WhisperVPC

CidrBlock: 10.0.1.0/24

AvailabilityZone: !Select [0, !GetAZs '']

MapPublicIpOnLaunch: true

WhisperSubnet2:

Type: AWS::EC2::Subnet

Properties:

VpcId: !Ref WhisperVPC

CidrBlock: 10.0.2.0/24

AvailabilityZone: !Select [1, !GetAZs '']

MapPublicIpOnLaunch: true

WhisperInternetGateway:

Type: AWS::EC2::InternetGateway

WhisperAttachGateway:

Type: AWS::EC2::VPCGatewayAttachment

Properties:

VpcId: !Ref WhisperVPC

InternetGatewayId: !Ref WhisperInternetGateway

WhisperRouteTable:

Type: AWS::EC2::RouteTable

Properties:

VpcId: !Ref WhisperVPC

WhisperRoute:

Type: AWS::EC2::Route

DependsOn: WhisperAttachGateway

Properties:

RouteTableId: !Ref WhisperRouteTable

DestinationCidrBlock: 0.0.0.0/0

GatewayId: !Ref WhisperInternetGateway

WhisperSubnetRouteTableAssoc1:

Type: AWS::EC2::SubnetRouteTableAssociation

Properties:

SubnetId: !Ref WhisperSubnet1
RouteTableId: !Ref WhisperRouteTable

WhisperSubnetRouteTableAssoc2:

Type: AWS::EC2::SubnetRouteTableAssociation

Properties:

SubnetId: !Ref WhisperSubnet2
RouteTableId: !Ref WhisperRouteTable

WhisperSecurityGroup:

Type: AWS::EC2::SecurityGroup

Properties:

GroupDescription: Allow HTTP access
VpcId: !Ref WhisperVPC
SecurityGroupIngress:

- IpProtocol: tcp
FromPort: 9000
ToPort: 9000
CidrIp: 0.0.0.0/0

WhisperCluster:

Type: AWS::ECS::Cluster

WhisperTaskExecutionRole:

Type: AWS::IAM::Role

Properties:

AssumeRolePolicyDocument:

Statement:

- Effect: Allow
Principal:
Service: ecs-tasks.amazonaws.com
Action: sts:AssumeRole

ManagedPolicyArns:

- arn:aws:iam::aws:policy/service-role/AmazonECSTaskExecutionRolePolicy

WhisperTaskDefinition:

Type: AWS::ECS::TaskDefinition

Properties:

Family: whisper-task

```
RequiresCompatibilities: [FARGATE]
Cpu: 512
Memory: 1024
NetworkMode: awsvpc
ExecutionRoleArn: !GetAtt WhisperTaskExecutionRole.Arn
ContainerDefinitions:
```

- Name: whisper
Image: onerahmet/openai-whisper-asr-

```
webservice:latest
```

```
PortMappings:
```

- ContainerPort: 9000

```
Environment:
```

- Name: ASR_MODEL
Value: !Ref WhisperModel

```
WhisperService:
```

```
Type: AWS::ECS::Service
```

```
DependsOn: WhisperALBListener
```

```
Properties:
```

```
Cluster: !Ref WhisperCluster
```

```
LaunchType: FARGATE
```

```
DesiredCount: 1
```

```
NetworkConfiguration:
```

```
AwsvpcConfiguration:
```

```
AssignPublicIp: ENABLED
```

```
SecurityGroups: [!Ref WhisperSecurityGroup
```

```
Subnets: [!Ref WhisperSubnet1, !Ref WhisperSubnet2]
```

```
TaskDefinition: !Ref WhisperTaskDefinition
```

```
LoadBalancers:
```

- ContainerName: whisper
ContainerPort: 9000
TargetGroupArn: !Ref WhisperTargetGroup

```
WhisperALB:
```

```
Type: AWS::ElasticLoadBalancingV2::LoadBalancer
```

```
Properties:
```

```
Name: whisper-alb
```

```
Subnets: [!Ref WhisperSubnet1, !Ref WhisperSubnet2]
```

```
SecurityGroups: [!Ref WhisperSecurityGroup]
```

Scheme: internet-facing

Type: application

WhisperTargetGroup:

Type: AWS::ElasticLoadBalancingV2::TargetGroup

Properties:

Port: 9000

Protocol: HTTP

VpcId: !Ref WhisperVPC

TargetType: ip

HealthCheckPath: /docs

WhisperALBListener:

Type: AWS::ElasticLoadBalancingV2::Listener

Properties:

LoadBalancerArn: !Ref WhisperALB

Port: 9000

Protocol: HTTP

DefaultActions:

- **Type:** forward

TargetGroupArn: !Ref WhisperTargetGroup

Outputs:

WhisperAPIURL:

Description: Whisper REST API URL

Value: !Join ["", ["http://", !GetAtt WhisperALB.DNSName, ":9000"]]

Steps to Deploy the Model

Deploy this using the AWS CLI deploy command:

```
#!/bin/bash
```

```
aws cloudformation deploy --region ap-south-1 \  
  --template-file ./main.yaml \  
  --stack-name ecsaimodel \  
  --tags madeFromCLI=yes anotherTagForAllStackResources=okay \  
  --capabilities CAPABILITY_NAMED_IAM
```

--no-execute-changeset

Screenshot of the model running on ECS

Whisper Asr WebService

Service health | Elastic C... | Load balancer details | E... | whisper-alb-765361359

Not secure whisper-alb-765361359.ap-south-1.elb.amazonaws.com:9000/docs#/Endpoints/asr_post

Name	Description
encode	Encode audio first through ffmpeg
boolean (query)	Default value : true
task	Available values : transcribe, translate
string (query)	Default value : transcribe
language	Available values : af, am, ar, as, az, ba, be, bg, bn, bo, br, bs, ca, cs, cy, da, de, el, en, es, et, eu, fa, fi, fo, fr, gl, gu, ha, haw, he, hi, hr, ht, hu, hy, id, is, it, ja, jw, ka, kk, km, kn, ko, la, lb, ln, lo, lt, lv, mg, mi, mk, ml, mn, mr, ms, mt, my, ne, nl, nn, no, oc, pa, pl, ps, pt, ro, ru, sa, sd, si, sk, sl, sn, so, sq, sr, su, sv, sw, ta, te, tg, th, tk, tl, tr, tt, uk, ur, uz, vi, yi, yo, yue, zh
initial_prompt	initial_prompt
output	Available values : txt, vtt, srt, tsv, json
string (query)	Default value : txt

Request body required

audio_file * required

string (binary)

```
ooumua@fedora:~/Downloads/test1$ curl -X POST http://whisper-alb-765361359.ap-south-1.elb.amazonaws.com:9000/asr \
-F audio_file=@rec1.flac
Hello world, this is Aryan Pandey. Thank you.
[ooumua@fedora test1]$
```

Whisper Asr WebService ^{1.8.2} ^{OAS 3.1}

[openapi.json](#)

Whisper ASR WebService is a general-purpose speech recognition webservice.

[MIT License](#)

Endpoints

POST /asr Asr

Parameters

Name	Description
encode	Encode audio first through ffmpeg
boolean (query)	Default value : true
task	Available values : transcribe, translate
string (query)	Default value : transcribe
language	Available values : af, am, ar, as, az, ba, be, bg, bn, bo, br, bs, ca, cs, cy, da, de, el, en, es, et, eu, fa, fi, fo, fr, gl, gu, ha, haw, he, hi, hr, ht, hu, hy, id, is, it, ja, jw, ka, kk, km, kn, ko, la, lb, ln, lo, lt, lv, mg, mi, mk, ml, mn, mr, ms, mt, my, ne, nl, nn, no, oc, pa, pl, ps, pt, ro, ru, sa, sd, si, sk, sl, sn, so, sq, sr, su, sv, sw, ta, te, tg, th, tk, tl, tr, tt, uk, ur, uz, vi, yi, yo, yue, zh
initial_prompt	

Try it out

```
ooumua@fedora:~/Downloads/test1$ curl -X POST http://whisper-alb-765361359.ap-south-1.elb.amazonaws.com:9000/asr \
-F audio_file=@rec1.flac
Hello world, this is Aryan Pandey. Thank you.
[ooumua@fedora test1]$
```

Amazon Elastic Container Service

Clusters > ecsaimodel-WhisperCluster-6QbnxBNgKJXh > Services > ecsaimodel-WhisperService-nM1UjeswHwLW > Health

ecsaimodel-WhisperService-nM1UjeswHwLW

Last updated April 15, 2025 at 03:03 (UTC+5:30) [Update service](#) [Delete service](#)

Service overview

Status: Active Tasks (1 Desired): 0 Pending | 1 Running Task definition: [revision whisper-task:1](#) Deployment status: Success

Health and metrics | Tasks | Logs | Deployments | Events | Configuration and networking | Service auto scaling | Tags

Status

Service name: [ecsaimodel-WhisperService-nM1UjeswHwLW](#) Service ARN: [arn:aws:ecs:ap-south-1:222634371739:service/ecsaimodel-WhisperCluster-6QbnxBNgKJXh/ecsaimodel-WhisperService-nM1UjeswHwLW](#) Deployments current state: 1 Completed task Created at: [April 15, 2025 at 02:52 \(UTC+5:30\)](#)

Health check grace period: 0 seconds

Load balancer health

Load balancer	Load balancer type	Listeners	Target group	Targets
whisper-alb	Application Load Balancer	HTTP:9000	ecsaim-Whisp-UFGBACU1I3XX Details	1 Healthy 0 Unhealthy

© 2025, Amazon Web Services, Inc. or its affiliates. [Privacy](#) [Terms](#) [Cookie preferences](#)

EC2 > Load balancers > whisper-alb

whisper-alb

[Actions](#)

Details

Load balancer type: Application Status: Active VPC: [vpc-039de96d703ff2ba8](#) Load balancer IP address type: IPv4

Scheme: Internet-facing Hosted zone: [ZP97RAFLXTNZK](#) Availability Zones: [subnet-07c5088d0eba5ecee](#) ap-south-1a (aps1-az1) [subnet-0753968c972545143](#) ap-south-1b (aps1-az3) Date created: April 15, 2025, 02:49 (UTC+05:30)

Load balancer ARN: [arn:aws:elasticloadbalancing:ap-south-1:222634371739:loadbalancer/app/whisper-alb/ac8abdccb0a6e001](#) DNS name: [whisper-alb-765361359.ap-south-1.elb.amazonaws.com \(A Record\)](#)

Listeners and rules (1) | Network mapping | Resource map | Security | Monitoring | Integrations | Attributes | Capacity | Tags

[Manage rules](#) [Manage listener](#) [Add listener](#)

A listener checks for connection requests on its configured protocol and port. Traffic received by the listener is routed according to the default action and any additional rules.

< 1 > [Settings](#)

© 2025, Amazon Web Services, Inc. or its affiliates. [Privacy](#) [Terms](#) [Cookie preferences](#)

CloudFormation console showing the 'Stacks' view for the 'ecsaiproject' stack. The stack is in the 'Active' state and contains 5 resources. The 'Resources' tab is selected, displaying a table of 17 resources.

Logical ID	Physical ID	Type	Status
WhisperALB	arn:aws:elasticloadbalancing:ap-south-1:222634371739:loadbalancer/app/whisper-alb/ac8abdcdb0a6e001	AWS::ElasticLoadBalancingV2::LoadBalancer	CREATE_COMPLETE
WhisperALBListener	arn:aws:elasticloadbalancing:ap-south-1:222634371739:listener/app/whisper-alb/ac8abdcdb0a6e001/67f2e7c236028d33	AWS::ElasticLoadBalancingV2::Listener	CREATE_COMPLETE
WhisperAttachGateway	IGWVjpc-039de96d703ff2ba8	AWS::EC2::VPCGatewayAttachment	CREATE_COMPLETE
WhisperCluster	ecsaiproject-WhisperCluster-6QbnxBNgKJXh	AWS::ECS::Cluster	CREATE_COMPLETE

CloudFormation console showing the 'Stacks' view for the 'ecsaiproject' stack. The stack is in the 'Active' state and contains 5 resources. The 'Events' tab is selected, displaying a timeline of events for the stack.

Event Name	Event Type	Timestamp
WhisperService	CREATE_COMPLETE	2025-04-15 02:49:00 UTC+0530
WhisperALBListener	CREATE_COMPLETE	2025-04-15 02:49:30 UTC+0530
WhisperRoute	CREATE_COMPLETE	2025-04-15 02:50:00 UTC+0530
WhisperALB	CREATE_COMPLETE	2025-04-15 02:50:30 UTC+0530
WhisperTaskDefinition	CREATE_COMPLETE	2025-04-15 02:51:00 UTC+0530
WhisperSubnetRouteTableAs...	CREATE_COMPLETE	2025-04-15 02:51:30 UTC+0530
WhisperSubnetRoute TableAs...	CREATE_COMPLETE	2025-04-15 02:52:00 UTC+0530
WhisperAttachGateway	CREATE_COMPLETE	2025-04-15 02:52:30 UTC+0530
WhisperSecurityGroup	CREATE_COMPLETE	2025-04-15 02:53:00 UTC+0530
WhisperSubnet1	CREATE_COMPLETE	2025-04-15 02:53:30 UTC+0530
WhisperTargetGroup	CREATE_COMPLETE	2025-04-15 02:54:00 UTC+0530
WhisperSubnet2	CREATE_COMPLETE	2025-04-15 02:54:30 UTC+0530
WhisperRoute Table	CREATE_COMPLETE	2025-04-15 02:55:00 UTC+0530
WhisperVPC	CREATE_COMPLETE	2025-04-15 02:55:30 UTC+0530
WhisperInternetGateway	CREATE_COMPLETE	2025-04-15 02:56:00 UTC+0530
WhisperCluster	CREATE_COMPLETE	2025-04-15 02:56:30 UTC+0530
WhisperTaskExecutionRole	CREATE_COMPLETE	2025-04-15 02:57:00 UTC+0530