

apiVersion: apps/v1

kind: Deployment

metadata:

name: request-router-deployment

spec:

replicas: 2

selector:

matchLabels:

app: request-router

template:

metadata:

labels:

app: request-router

spec:

containers:

- name: request-router

image: public.ecr.aws/d6u1k1m2/lambda-router:latest

ports:

- containerPort: 8000

env:

- name: COGVLM\_ENDPOINT

value: "http://cogvlm-service.default.svc.cluster.local" # Fully qualified DNS name

- name: QWENVL\_ENDPOINT

value: "http://qwenvl-service.default.svc.cluster.local" # Fully qualified DNS name

---

apiVersion: v1

kind: Service

metadata:

name: request-router-service

spec:

type: LoadBalancer

selector:

app: request-router

ports:

- protocol: TCP

port: 80

targetPort: 8000