

apiVersion: apps/v1

kind: Deployment

metadata:

name: qwenvl-deployment

spec:

replicas: 2

selector:

matchLabels:

app: qwenvl

template:

metadata:

labels:

app: qwenvl

spec:

containers:

- name: qwenvl

image: public.ecr.aws/d6u1k1m2/qwenvlm:latest

resources:

requests:

memory: "20Gi"

limits:

memory: "25Gi"