

envs:

MODEL\_NAME: xtuner/llava-llama-3-8b-v1\_1

MODEL\_ARCH: llava

HF\_HUB\_ENABLE\_HF\_TRANSFER: True

resources:

# accelerators: {L4:4, A100:4, A100:8, A100-80GB:2, A100-80GB:4, A100-80GB:8} ## Large

models

accelerators: [L4, A10g, A100, A100, A100-80GB, T4, M60] ## Small models

# cpus: 32+

memory: 32+

# use\_spot: True

# disk\_size: 512 # Ensure model checkpoints (~246GB) can fit.

# disk\_tier: best

ports: 8080 # Expose to internet traffic.

service:

readiness\_probe:

path: /v1/chat/completions

post\_data:

model: \$MODEL\_NAME

messages:

- role: user

content: Hello! What is your name?

max\_tokens: 1

readiness\_probe: /v1/models

replica\_policy:

min\_replicas: 1

max\_replicas: 10

target\_qps\_per\_replica: 2.5

upscale\_delay\_seconds: 300

downscale\_delay\_seconds: 1200

setup: |

pip install hf\_transfer

docker pull openmmlab/lmdeploy:latest

run: |

# Serve With Docker

docker run -d --runtime nvidia --gpus all \

-v ~/.cache/huggingface:/root/.cache/huggingface \

--env "HUGGING\_FACE\_HUB\_TOKEN=hf\_ksMHvhGLTINtdSHXBihthxFFjfbWlszaaM" \

-p 8080:8080 \

--ipc=host \

openmmlab/lmdeploy:latest \

lmdeploy serve api\_server \$MODEL\_NAME --model-name \$MODEL\_ARCH --server-port 8080