容器服务 AI 模型部署





【版权声明】

©2013-2025 腾讯云版权所有

本文档(含所有文字、数据、图片等内容)完整的著作权归腾讯云计算(北京)有限责任公司单独所有,未经腾讯云 事先明确书面许可,任何主体不得以任何形式复制、修改、使用、抄袭、传播本文档全部或部分内容。前述行为构成 对腾讯云著作权的侵犯,腾讯云将依法采取措施追究法律责任。

【商标声明】



腾讯云

及其它腾讯云服务相关的商标均为腾讯云计算(北京)有限责任公司及其关联公司所有。本文档涉及的第三方主体的 商标,依法由权利人所有。未经腾讯云及有关权利人书面许可,任何主体不得以任何方式对前述商标进行使用、复 制、修改、传播、抄录等行为,否则将构成对腾讯云及有关权利人商标权的侵犯,腾讯云将依法采取措施追究法律责 任。

【服务声明】

本文档意在向您介绍腾讯云全部或部分产品、服务的当时的相关概况,部分产品、服务的内容可能不时有所调整。 您所购买的腾讯云产品、服务的种类、服务标准等应由您与腾讯云之间的商业合同约定,除非双方另有约定,否则, 腾讯云对本文档内容不做任何明示或默示的承诺或保证。

【联系我们】

我们致力于为您提供个性化的售前购买咨询服务,及相应的技术售后服务,任何问题请联系 4009100100或 95716。



文档目录

AI 模型部署 Embedding 模型部署 AI 网关部署



AI 模型部署 Embedding 模型部署

最近更新时间: 2025-10-13 14:56:12

本文档介绍如何在 TKE 上部署 Embedding 模型。

简介

Embedding 模型是将文本、图像等非结构化数据转化为低维稠密向量的 AI 模型,这些向量能精准捕捉数据的语义或特征关联,让机器可通过向量计算理解数据间的相似性,是语义检索、推荐系统、聚类分析等场景的核心基础,在文档检索中,可通过对比查询与文档的 Embedding 向量快速匹配相关内容。

tke-ai-playbook 是 TKE 团队开源的 AI 大模型相关脚本,包含模型下载、部署推理服务、性能测试等模块,可在 TKE 一站式体验 AI 相关功能,点此获取 开源链接。

环境准备

步骤1: 创建 TKE 标准集群

- 1. 登录 腾讯云容器服务控制台,选择左侧导航栏中的集群。
- 2. 单击集群列表上方的新建。
- 3. 在集群类型中,选择标准集群。
- 4. 在**创建集群 > 网络配置**中,选择 VPC-CNI,其余参数保持默认即可。



5. 在**创建集群 > 组件配置**中,存储组件勾选 **CFS**(用于持久化存储大模型权重文件;若创建时未勾选,后续可在 "集群 - 组件管理"中安装,操作详情请参见 通过组件管理页安装)。





6. 创建下载节点池。

选择 3 个节点(用于后续 3 个并发下载任务),推荐机型 SA5.LARGE8(性价比高,满足下载需求),系统 盘≥300GB,并分配 "免费公网 IP" 及带宽≥100Mbps(保障单节点下载速度,缩短整体耗时)。



7. 创建推理节点:需 1 个 GPU 节点(大模型推理依赖高显存),推荐选择 PNV5b.8XLARGE96机型(单卡48GB 显存,满足模型需求)。

步骤2: 创建存储组件

需创建 StorageClass (CFS 存储)和 PVC,确保模型下载后能被推理服务直接调用:



 进入集群的存储 > StorageClass 页面,单击新建,存储类型选 "文件存储 CFS",可用区与推理节点一致 (减少跨区访问延迟),其余参数默认;



2. 进入**存储 > PVC** 页面,单击**新建**,名称为 ai-model(与后续下载脚本参数一致,避免手动修改配置),存储 类选择上一步创建的 StorageClass 对象。



部署模型

本实践使用 vllm 推理框架,该框架支持的 embedding 模型主要有: BGE 系列模型和 E5-Mistral 系列模型,这两个系列的模型均可使用 vllm 框架进行模型推理,本文采用 E5-mistral-7b-instruct 模型完成验证。 E5-mistral-7b-instruct: 基于 Mistral-7B 架构微调的指令遵循模型,融合了 E5 系列在检索增强领域的优势与 Mistral 的高效推理能力;在中文与多语言理解、长文本处理和检索增强生成(RAG)场景中表现突出,并且拥有较长的上下文长度,支持 vllm 框架下的推理服务。

步骤1: 下载依赖

登录集群中的节点,节点镜像为 TencentOS 可以使用以下命令进行安装:

yum install -y jq



curl https://raw.githubusercontent.com/helm/helm/main/scripts/get-helm-3
| bash

步骤2: 下载模型

进入集群中,拉取 tke-ai-playbook,机器若提示安装 git,请先安装 git 后再拉取代码:

```
git clone https://github.com/tkestack/tke-ai-playbook.git
```

执行 playbook 中的下载脚本 tke-llm-downloader.sh,该脚本可以使用多个节点并发拉取模型,加快模型下载速度,PVC 名称需与前述一致,模型名称参数则选择目标模型名称,可前往 ModelScope 获取模型名称。
tke-llm-downloader.sh 是一个在 Kubernetes 集群中自动化下载大语言模型的工具脚本,主要功能包括:

- 模型下载管理: 支持从 ModelScope 或 Hugging Face 下载大语言模型到 Kubernetes PVC 存储;
- 并发下载优化: 通过 Kubernetes Job 的 completions 和 parallelism 参数实现多 Pod 并发下载,提升下载速度;
- 节点调度控制: 避免多个下载 Pod 调度到同一节点。

```
bash scripts/tke-llm-downloader.sh --pvc ai-model --completions 3 --
parallelism 3 --model intfloat/e5-mistral-7b-instruct
```

预期结果:三个 Pod 状态为 Completed 后,模型下载工作则执行完成。

tke-llm-downloader-6bcrn-0-dq2mv	0/1	Completed	0	
tke-llm-downloader-6bcrn-1-7nlfs	0/1	Completed	0	
tke-llm-downloader-6bcrn-2-5kmq7	0/1	Completed	0	

步骤3: 模型部署

vllm-inference-tke 是一个部署 vllm 推理框架并暴露服务的 Helm Chart 包,可支持单机单卡、单机多卡、 多机多卡等多种部署模式,修改自定义参数即可部署自建大模型。

进入 vllm-inference-tke 目录,运行 vllm-inference-tke 的 Chart 包,该 Chart 可以在 TKE 上部署基于 vLLM 的 OpenAI 兼容 API 服务,APIKey 需要自行设置,默认无需 APIKey,修改 model name 为自己想要部署的模型,并且 PVC 修改为自己集群的 PVC 名称,具体参数修改参考下述,目前默认参数支持单机单卡部署。

```
cd tke-ai-playbook/helm-charts/vllm-inference-tke
vim values.yaml
helm install vllm-service .
```

本文提供部署 E5-mistral-7b-instruct 模型的具体参数,具体的 values.yaml 文件内容如下:



```
model:
 pvc:
   enabled: true
   path: "intfloat/e5-mistral-7b-instruct"
 local:
    enabled: false
    path: "/data0/intfloat/e5-mistral-7b-instruct"
server:
 replicas: 1
 lwsGroupSize: 1
  image: "ccr.ccs.tencentyun.com/tke-ai-playbook/vllm-openai:v0.10.1-
  imagePullPolicy: IfNotPresent
 apiKey: ""
 resources:
    requests:
      nvidia.com/gpu: 1
    limits:
      nvidia.com/gpu: 1
 args:
   tpSize: 1
   ppSize: 1
   epEnabled: false
   maxModelLen: 2048
   maxBatchSize: 32
 extraArgs:
  - --disable-log-requests
  - --cuda-graph-sizes 1 2 4 8 16 24 32
 env:
  - name: VLLM_WORKER_MULTIPROC_METHOD
   value: "spawn"
 service:
   enabled: true
    type: LoadBalancer
    port: 60000
```



```
labels: {}

podAnnotations: {}
```

部署的模型会通过 service 暴露端点和端口,查看 Pod 日志看模型状态是否健康。

```
kubectl get svc | grep vllm-service
kubectl logs -f pod-name
```

预期结果:

```
api_server.py:1847] Starting vLLM API server 0 on http://0.0.0.0:60000
launcher.py:29] Available routes are:
APIServer pid=1)
                      INFO 09-25
                                     19:54:39
                      INFO 09-25
                                     19:54:39
                                                  [launcher.py:37]
APIServer pid=1)
                                                                       Route: /openapi.json, Methods: HEAD, GET
                                                                      Route: /docs, Methods: HEAD, GET
Route: /docs/oauth2-redirect, Methods: HEAD, GET
                      INFO 09-25
                                     19:54:39
                                                  launcher.py:37]
APIServer pid=1)
                      INFO 09-25
                                    19:54:39
(APIServer pid=1)
                                                  [launcher.py:37]
                                                                      Route: /redoc, Methods: HEAD, GET
Route: /health, Methods: GET
Route: /load, Methods: GET
Route: /ping, Methods: POST
APIServer pid=1)
                      INFO 09-25
                                     19:54:39
                                                  launcher.py:37]
                      INFO 09-25
                                     19:54:39
(APIServer pid=1)
                                                  [launcher.py:37]
                      INFO 09-25
                                     19:54:39
APIServer pid=1)
                                                  launcher.py:37]
                      INFO 09-25
                                    19:54:39
(APIServer pid=1)
                                                 [launcher.py:37]
                                                                       Route: /ping, Methods: GET
Route: /tokenize, Methods: POST
APIServer pid=1)
                      INFO 09-25
                                     19:54:39
                                                 [launcher.py:37]
                      INFO 09-25
                                     19:54:39
(APIServer pid=1)
                                                 [launcher.py:37]
                                                                      Route: /detokenize, Methods: POST
Route: /v1/models, Methods: GET
Route: /version, Methods: GET
                      INFO 09-25
                                     19:54:39
APIServer pid=1)
                                                 [launcher.py:37]
                      INFO 09-25
APIServer pid=1)
                                     19:54:39
                                                 [launcher.py:37]
APIServer pid=1)
                      INFO 09-25
                                                 [launcher.py:37]
```

步骤4:测试向量嵌入

```
curl -X POST http://xx.xx.xx.xx:60000/v1/embeddings \
  -H 'Content-Type: application/json' \
  -H 'Authorization: Bearer 3bxxxxxxxxxxxxx1a1d' \
  -d '{
    "model": "intfloat/e5-mistral-7b-instruct",
    "input": "用Python实现快速排序算法"
}'
```

预期结果:



常见问题

部署 vllm-service 后日志显示未找到模型?

检查模型是否已下载完成,Pod 状态应为 Completed。如果模型下载的 Pod 还是 Running 状态,请等待下载结束后重启服务。

部署 vllm-service 后日志显示 CudaOutofMemory?

说明 GPU 显存不足,请选用更大显存的 GPU,部署本文模型需单卡 L20 及以上机型。



AI 网关部署

最近更新时间: 2025-10-10 10:36:21

本文档介绍在 TKE 上如何部署 Envoy AI Gateway 并接入大模型。

简介

AI 阿关是面向大模型服务的专用流量治理组件,基于传统 API 阿关扩展而来,核心能实现多模型按需切换、租户鉴权与配额管控、Token 级限流、内容安全合规校验,还可在模型异常时自动切换保障稳定性。

Envoy Al Gateway 是一个开源的、基于 Kubernetes 原生架构的 Al 网关,专门用于管理和路由大模型服务流量。该项目构建在成熟的 Envoy Proxy 之上,为应用客户端与各种 Al 服务提供商之间提供了一个统一、安全且可扩展的接入层。

环境准备

- 1. 已创建并部署好 TKE 集群,建议选择**香港地域**集群,如果您还没有集群,请参见 创建集群。
- 2. 已创建节点池,并且节点池内有至少3个节点,推荐机型为 SA5.LARGE8。

环境验证

验证 kubectl:

kubectl version --client

预期结果:

Client Version: v1.32.2-tke.6 Kustomize Version: v5.5.0

验证 helm 安装:

helm version

预期结果:

version.BuildInfo{Version:"v3.19.0",

GitCommit: "3d8990f0836691f0229297773f3524598f46bda6",

GitTreeState:"clean", GoVersion:"go1.24.7"}



如果显示未找到 helm 命令,TencentOS 系统可以用下列命令一键安装:

```
curl https://raw.githubusercontent.com/helm/helm/main/scripts/get-helm-3
| bash
```

验证 curl 安装:

```
curl --version
```

```
curl 8.4.0 (x86_64-koji-linux-gnu) libcurl/8.4.0 OpenSSL/3.0.12 zlib/1.2.13 brotli/1.1.0 libidn2/2.3.4 libpsl/0.21.2 (+libidn2/2.3.4) libssh/0.10.5/openssl/zlib nghttp2/1.58.0 OpenLDAP/2.6.5 Release-Date: 2023-10-11 Protocols: dict file ftp ftps gopher gophers http https imap imaps ldap ldaps mqtt pop3 pop3s rtsp scp sftp smb smbs smtp smtps telnet tftp Features: alt-svc AsynchDNS brotli GSS-API HSTS HTTP2 HTTPS-proxy IDN IPv6 Kerberos Largefile libz NTLM PSL SPNEGO SSL threadsafe TLS-SRP UnixSocket
```

部署 AI 网关

第一步: 部署 Envoy Gateway

Envoy Al Gateway 建立在 Envoy Gateway 之上,使用 Helm 安装并等待部署准备就绪:

```
helm upgrade -i eg oci://docker.io/envoyproxy/gateway-helm \
--version v0.0.0-latest \
--namespace envoy-gateway-system \
--create-namespace

kubectl wait --timeout=2m -n envoy-gateway-system deployment/envoy-
gateway --for=condition=Available
```

第二步: 部署 Envoy Al Gateway

安装 AI Gateway 的 Helm chart,完成后等待部署准备就绪:

```
helm upgrade -i aieg oci://docker.io/envoyproxy/ai-gateway-helm \
--version v0.0.0-latest \
```



```
--namespace envoy-ai-gateway-system \
--create-namespace

kubectl wait --timeout=2m -n envoy-ai-gateway-system deployment/ai-
gateway-controller --for=condition=Available
```

第三步:配置 Envoy Al Gateway

安装 Envoy Al Gateway 后,将特定于 Al Gateway 的配置应用于 Envoy Gateway,重新启动部署,然后等待其准备就绪:

```
kubectl apply -f https://raw.githubusercontent.com/envoyproxy/ai-
gateway/main/manifests/envoy-gateway-config/redis.yaml
kubectl apply -f https://raw.githubusercontent.com/envoyproxy/ai-
gateway/main/manifests/envoy-gateway-config/config.yaml
kubectl apply -f https://raw.githubusercontent.com/envoyproxy/ai-
gateway/main/manifests/envoy-gateway-config/rbac.yaml
kubectl rollout restart -n envoy-gateway-system deployment/envoy-gateway
kubectl wait --timeout=2m -n envoy-gateway-system deployment/envoy-
gateway --for=condition=Available
```

第四步: 检查健康状态

检查 AI 网关 Pod:

```
kubectl get pods -n envoy-ai-gateway-system
```

检查 Envoy Gateway 容器:

```
kubectl get pods -n envoy-gateway-system
```

第五步: 部署测试后端

首先部署包含测试后端的基本 AI 网关设置:

```
kubectl apply -f https://raw.githubusercontent.com/envoyproxy/ai-
gateway/main/examples/basic/basic.yaml
```



等待网关容器准备就绪:

设置网关 URL:

```
export GATEWAY_URL=$(kubectl get gateway/envoy-ai-gateway-basic -o
jsonpath='{.status.addresses[0].value}')
```

验证 URL 是否可用:

```
echo $GATEWAY_URL
```

第六步: 测试 AI 网关

打开一个新终端,并使用以下命令向 AI 网关发送测试请求:

```
curl -H "Content-Type: application/json" \
-d '{ "model": "some-cool-self-hosted-model", "messages":
  [{"role":"system", "content": "Hi."}]}' \
$GATEWAY_URL/v1/chat/completions
```

预期结果:

```
{"choices":[{"message":{"role":"assistant", "content":"I am inevitable."}}]}
```

使用 AI 网关接入大模型

第一步: 下载配置模板

本文实践以 Kimi 大模型为例,先获取 AI 网关的默认模型配置模板,后续第二步修改模板中的 hostname 和 API_KEY,需要具备相关运营商的 API_KEY。

```
curl -O https://raw.githubusercontent.com/envoyproxy/ai-
gateway/main/examples/basic/openai.yaml
```



第二步:修改配置模板

编辑文件 openai.yaml 替换其中的 hostname 和 OPENAI_API_KEY,下面以 Kimi 为例,yaml 文件配置参考如下:

```
apiVersion: aigateway.envoyproxy.io/v1alpha1
kind: AIGatewayRoute
metadata:
  name: envoy-ai-gateway-basic-openai
  namespace: default
spec:
  parentRefs:
    - name: envoy-ai-gateway-basic
      kind: Gateway
      group: gateway.networking.k8s.io
  rules:
    - matches:
        - headers:
            - type: Exact
              name: x-ai-eg-model
              value: kimi-k2-0905-preview
      backendRefs:
        - name: envoy-ai-gateway-basic-openai
apiVersion: aigateway.envoyproxy.io/v1alpha1
kind: AIServiceBackend
metadata:
  name: envoy-ai-gateway-basic-openai
  namespace: default
spec:
  schema:
    name: OpenAI
  backendRef:
    name: envoy-ai-gateway-basic-openai
    kind: Backend
    group: gateway.envoyproxy.io
apiVersion: aigateway.envoyproxy.io/v1alpha1
kind: BackendSecurityPolicy
metadata:
```



```
name: envoy-ai-gateway-basic-openai-apikey
 namespace: default
spec:
 targetRefs:
   - group: aigateway.envoyproxy.io
      kind: AIServiceBackend
      name: envoy-ai-gateway-basic-openai
 type: APIKey
 apiKey:
    secretRef:
      name: envoy-ai-gateway-basic-openai-apikey
      namespace: default
apiVersion: gateway.envoyproxy.io/v1alpha1
kind: Backend
metadata:
 name: envoy-ai-gateway-basic-openai
 namespace: default
spec:
 endpoints:
    - fqdn:
        hostname: api.moonshot.cn
        port: 443
apiVersion: gateway.networking.k8s.io/v1alpha3
kind: BackendTLSPolicy
metadata:
 name: envoy-ai-gateway-basic-openai-tls
 namespace: default
spec:
 targetRefs:
    - group: 'gateway.envoyproxy.io'
      kind: Backend
      name: envoy-ai-gateway-basic-openai
 validation:
    wellKnownCACertificates: "System"
   hostname: api.moonshot.cn
apiVersion: v1
kind: Secret
```



```
metadata:
   name: envoy-ai-gateway-basic-openai-apikey
   namespace: default

type: Opaque
stringData:
   apiKey: sk-xxxxxxxxxxxx # Replace with your kimi API key.
```

第三步: 应用配置

应用更新的配置并等待网关容器准备就绪。

```
kubectl apply -f openai.yaml
kubectl wait pods --timeout=2m \
   -l gateway.envoyproxy.io/owning-gateway-name=envoy-ai-gateway-basic \
   -n envoy-gateway-system \
   --for=condition=Ready
```

第四步: 测试配置

预期结果:

```
{"id":"chatcmpl-
68d64bc7a5422d1970e278da","object":"chat.completion","created":175887456
7,"model":"kimi-k2-0905-preview","choices":[{"index":0,"message":
{"role":"assistant","content":"Hi there! How can I help you
```



```
today?"},"finish_reason":"stop"}],"usage":
{"prompt_tokens":9,"completion_tokens":11,"total_tokens":20}}
```

常见问题

1. 在执行 helm 之后,安装 chart 命令超时,无法安装 chart?

推荐集群选择香港地域,防止由于网络延迟等问题导致安装超时。

2. curl 命令测试之后,返回报错: No matching route found. It is likely that the model specified your request is not configured in the Gateway

检查 yaml 文件中的模型文件是否与 curl 命令中的一致。