

Ollama推理框架使用指南

1 商品说明

Ollama 是一款跨平台推理框架客户端（MacOS、Windows、Linux），专为无缝部署大型语言模型（LLM）（如 Llama 2、Mistral、Llava 等）而设计。通过一键式设置，Ollama 可以在本地运行 LLM，将所有交互数据保存在自己的机器上，从而提高数据的私密性和安全性。

本商品在鲲鹏云的AArch64架构上Ubuntu24.04和HCE2.0系统中进行安装后以镜像提供给用户使用。

2 商品购买

您可以在云商店搜索“Ollama推理框架”。

其中，地域、规格、推荐配置使用默认，购买方式根据您的需求选择按需/按月/按年，短期使用推荐按需，长期使用推荐按月/按年，确认配置后点击“立即购买”。

2.1 商品支持自定义 ECS 购买，具体见章节 3.1.1

2.2 使用 RFS 模板直接部署



必填项填写后，点击 下一步

参数名称	值	类型	描述
EC2实例代码	<input type="text"/>	字符串	EC2实例的操作系统，系统默认值，或从影子实例或AMI中，或从EC2实例的AMI中，小写字母，数字和特殊字符（仅支持ASCII）。
实例大小	4x	number	设置实例大小（范围4x到4x10）。
实例大小	5x	number	设置实例的实例大小，基于实例的实例，可以设置为0，可设置实例的实例大小，默认值为0。
版本	最新	字符串	选择版本。
主机IPv4地址	192.168.0.18	字符串	取值范围：192.168.0.18 to 192.168.0.24, 172.16.0.1 to 172.16.0.24, or 192.168.0.18 to 192.168.0.24
子网IPv4地址	192.168.10.0/24	字符串	必须是CIDR格式，必须是VPC的默认子网，子网掩码不能超过24。
子网网关	192.168.10.1	字符串	子网网关，必须是子网内的有效IP地址。
行集类型（不兼容应用的人）	最新	字符串	rowSet类型，可以是rowSet或rowSet，或行集。
行集类型（不兼容应用的人）	month	字符串	设置行集TypePaaS的生活及与行集，取值范围：month, year。
行集类型（不兼容应用的人）	1	字符串	设置行集TypePaaS的生活及与行集，取值范围：periodType-month（周期类型为月），或：periodType-year（周期类型为年），或：periodType-day（周期类型为日）。

高级设置

启用高级设置。启用高级设置后，您可以自定义实例的实例大小，实例大小不能超过实例大小，实例大小不能超过实例大小。

安全防护

启用安全防护。启用安全防护后，您可以自定义实例的安全策略。

创建直接计划后，点击 确定

创建直接计划

通过执行计划，可以自定义实例的实例大小。

执行计划名称

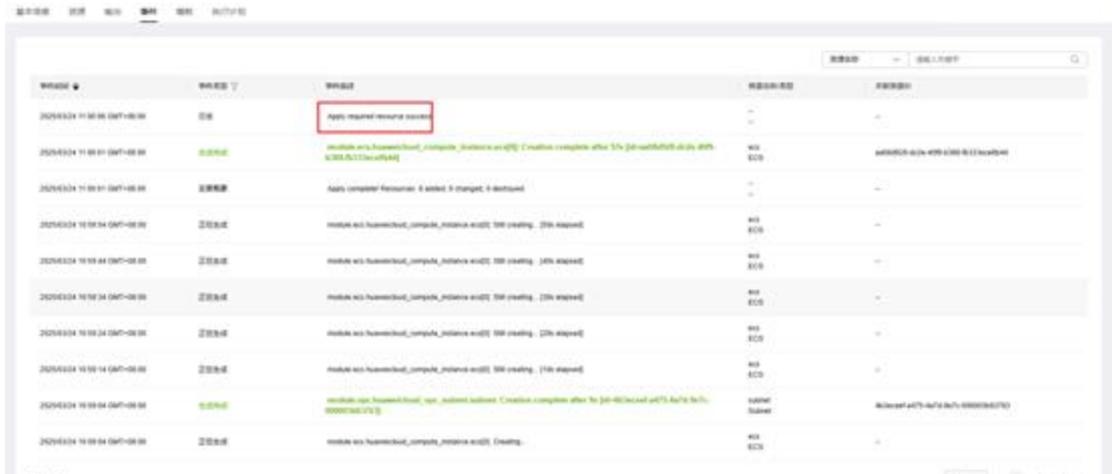
描述

确定 取消

点击 部署



如下图“Apply required resource success.”即为资源创建完成



3 商品资源配置

商品支持ECS控制台配置，下面对资源配置的方式进行介绍。

3.1 ECS 控制台配置

3.1.1 准备工作

在使用ECS控制台配置前，需要您提前配置好安全组规则。

安全组规则的配置如下：

- 入方向规则放通端口11434，源地址内必须包含您的客户端ip，否则无法访问
- 入方向规则放通CloudShell连接实例使用的端口22，以便在控制台登录调试。
- 出方向规则一键放通

3.1.2 创建 ECS

前提工作准备好后，选择ECS控制台配置跳转到购买ECS页面，ECS资源的配置如下图所示：

基础配置

计费模式

包年/包月

按需计费

竞价计费

按需计费实例不支持备案。 [了解备案限制](#)

区域

华北-北京四

推荐区域 华北-北京四

| 华南-广州

华东-上海一

华北-乌兰察布一

西南-贵阳一

云服务器创建后无法更改区域；不同区域之间内网互不相通，请就近选择靠近您业务的区域，减少网络时延。 [如何选择区域](#)

可用区

随机分配

可用区1

可用区2

可用区3

可用区7

 随机至多可用区

实例

规格类型选型

业务场景选型

CPU架构

x86计算

鲲鹏计算

实例筛选

--请选择vCPUs--

--请选择内存--

请输入规格名称模糊搜索

 隐藏售罄的规格**鲲鹏通用计算增强型**

鲲鹏内存优化型

鲲鹏超高I/O型

CSDN @p_xcn

操作系统

镜像

公共镜像

私有镜像

共享镜像

市场镜像

Huawei Cloud
EulerOS

CentOS



Ubuntu



EulerOS



Debian



KylinOS



UnionTechOS



openEuler

Huawei Cloud EulerOS 2.0 64bit for kAi2p with HDK 23.0.1 and CANN ...

存储与备份

系统盘

磁盘类型

系统盘大小(GiB)

通用型SSD

-

40

+

IOPS上限2,280, IOPS突发上限8,000 [高级设置](#)

增加一块数据盘

您还可以挂载 23 块磁盘 (云硬盘)

 开启备份

CSDN @p_xcn

值得注意的是：

- VPC您可以自行创建
- 安全组选择3.1.1章节中配置的安全组
- 弹性公网IP选择现在购买，推荐选择“按流量计费”，带宽大小可设置为5Mbit/s
- 高级配置需要在高级选项支持注入自定义数据，所以登录凭证不能选择“密码”，选择创建后设置
- 其余默认或按规则填写即可。

4 商品使用

4.1 登录服务器查看 Ollama 服务状态

通过命令 `systemctl status ollama` 查看服务状态。

```
root@image-hce-24u48g-rx1137 ~# systemctl status ollama
● ollama.service - Ollama Service
   Loaded: loaded (/etc/systemd/system/ollama.service; enabled; vendor preset: disabled)
   Active: active (running) since Tue 2025-06-03 15:39:51 CST; 3min 15s ago
     Main PID: 2283 (ollama)
        Tasks: 13 (limit: 297559)
       Memory: 38.5M
      CGroup: /system.slice/ollama.service
              └─ 2283 /usr/local/bin/ollama serve

Jun 03 15:39:54 image-hce-24u48g-rx1137 ollama[2283]: time=2025-06-03T15:39:54.235+08:00 level=INFO source=images.go:479 msg="t
Jun 03 15:39:54 image-hce-24u48g-rx1137 ollama[2283]: time=2025-06-03T15:39:54.236+08:00 level=INFO source=images.go:486 msg="t
Jun 03 15:39:54 image-hce-24u48g-rx1137 ollama[2283]: time=2025-06-03T15:39:54.236+08:00 level=INFO source=routes.go:1287 msg="p
Jun 03 15:39:54 image-hce-24u48g-rx1137 ollama[2283]: time=2025-06-03T15:39:54.977+08:00 level=INFO source=gpu.go:377 msg="look
Jun 03 15:39:55 image-hce-24u48g-rx1137 ollama[2283]: time=2025-06-03T15:39:55.298+08:00 level=INFO source=types.go:130 msg="inp
Jun 03 15:42:45 image-hce-24u48g-rx1137 ollama[2283]: [GIN] 2025/06/03 - 15:42:45 | 200 | 52.194µs | 127.0.0.1 | HEAD
Jun 03 15:42:45 image-hce-24u48g-rx1137 ollama[2283]: [GIN] 2025/06/03 - 15:42:45 | 200 | 116.479µs | 127.0.0.1 | GET
Jun 03 15:42:49 image-hce-24u48g-rx1137 ollama[2283]: [GIN] 2025/06/03 - 15:42:49 | 200 | 23.712µs | 127.0.0.1 | HEAD
Jun 03 15:42:49 image-hce-24u48g-rx1137 ollama[2283]: [GIN] 2025/06/03 - 15:42:49 | 200 | 778.457µs | 127.0.0.1 | GET
```



4.2 下载模型启动模型参考 Ollama 手册

完整操作参考 [Ollama手册](#)