# GPUStack- 单机版 用户手册

# 一、产品概述

## 1.1 产品介绍

GPUStack 云服务是基于开源 GPUStack 构建的托管式 AI 模型部署平台,让您无需管理基础设施,即可在高性能 GPU 集群上轻松部署和运行各类 AI 模型。

## 1.2 产品核心能力

**资源管理**:提供自动化 GPU 资源调度与集群管理,支持异构 GPU 设备统一纳管,实现资源利用率最大化与成本最优化;

**模型部署**:支持主流开源大模型一键部署,兼容 Hugging Face、ModelScope 等模型源,集成 vLLM、SGLang 和 TensorRT-LLM 等高性能推理引擎,满足不同场景性能需求;

**智能运维**:内置自动扩缩容、故障转移与负载均衡机制,提供实时性能监控与告警,确保服务高可用性与稳定性;

**安全管控**:提供完善的认证授权体系与网络隔离策略,支持私有化部署与数据安全保障,满足企业级安全合规要求。

## 1.3 产品优势

自动处理底层 GPU 资源调度、模型优化和扩展,让您专注于应用开发而非运维。

**零运维负担**:无需管理 GPU 驱动、CUDA 版本或集群配置;
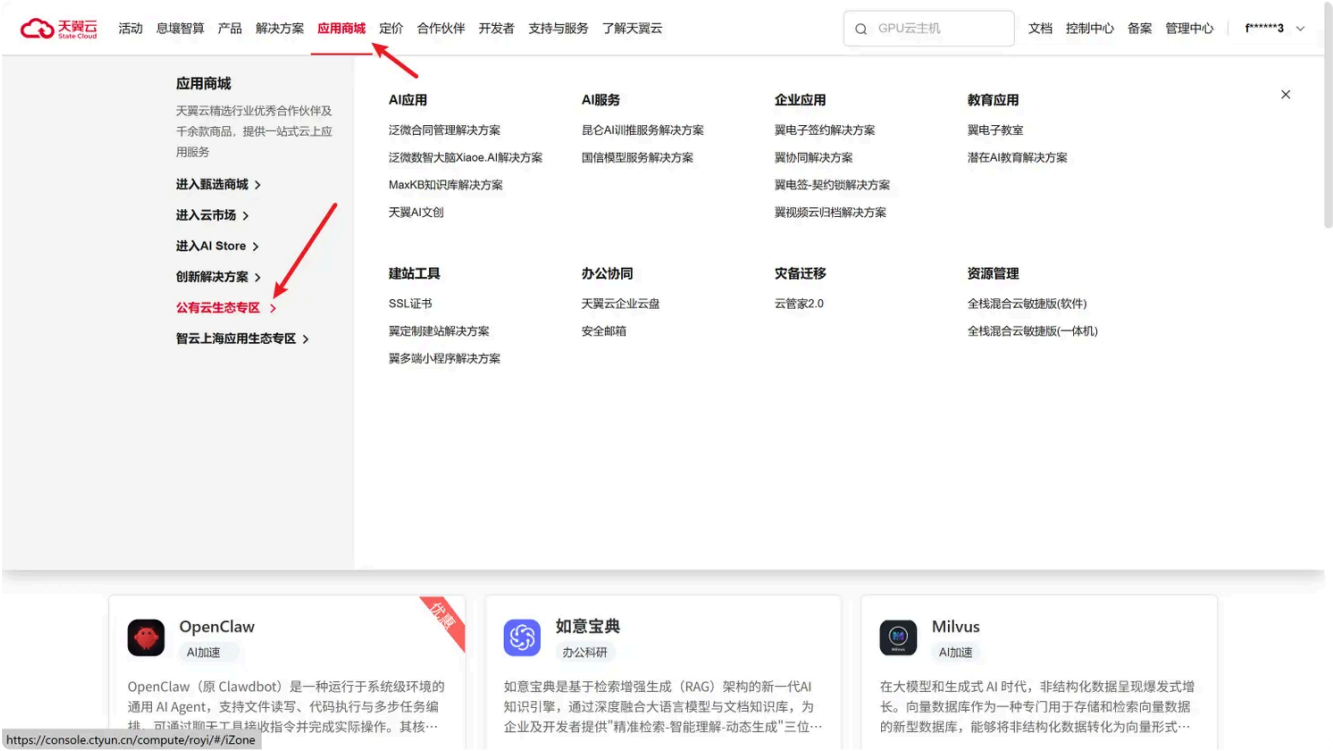
**开箱即用**:集成 vLLM、SGLang 和 TensorRT-LLM 等高性能推理引擎,支持自定义推理框架;

**一键部署**:支持从 Hugging Face、ModelScope、或本地直接部署,支持自动扩缩容、版本升降级;

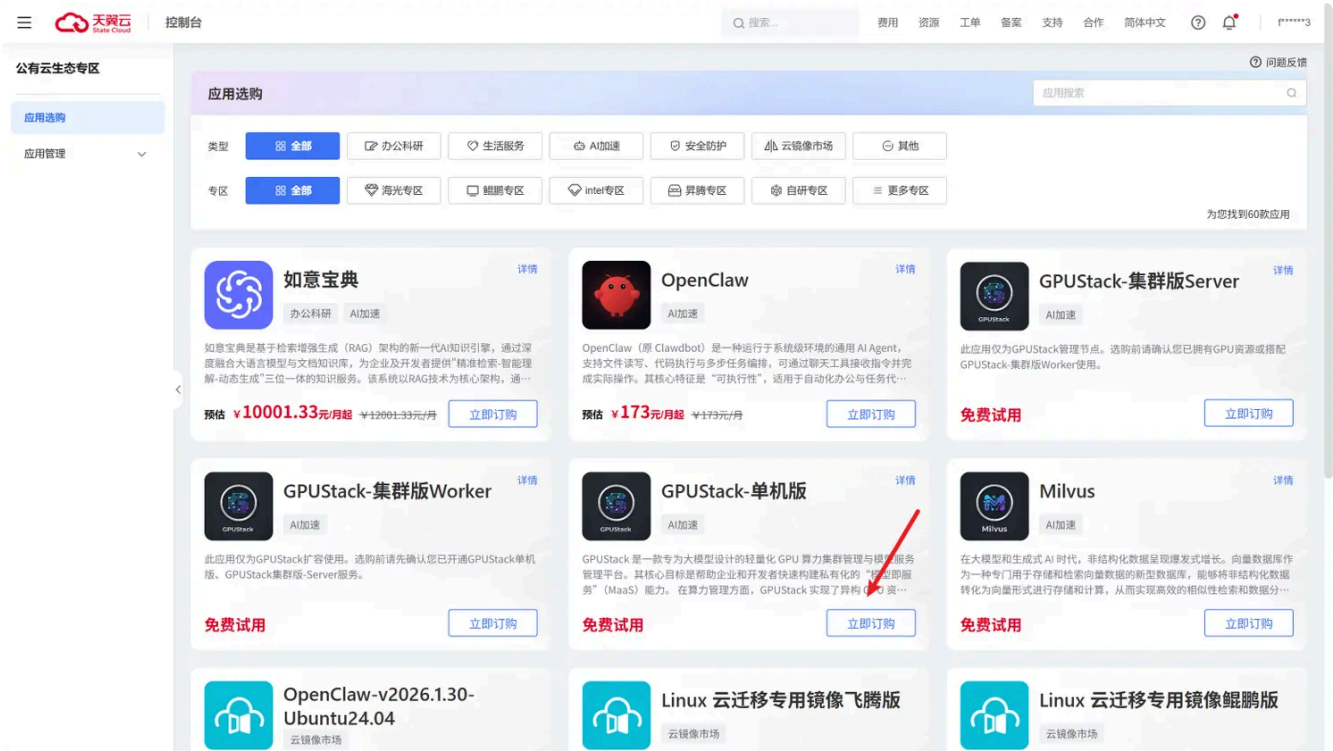**性能优化配置:**提供预调优模式,用于低延迟或高吞吐量;
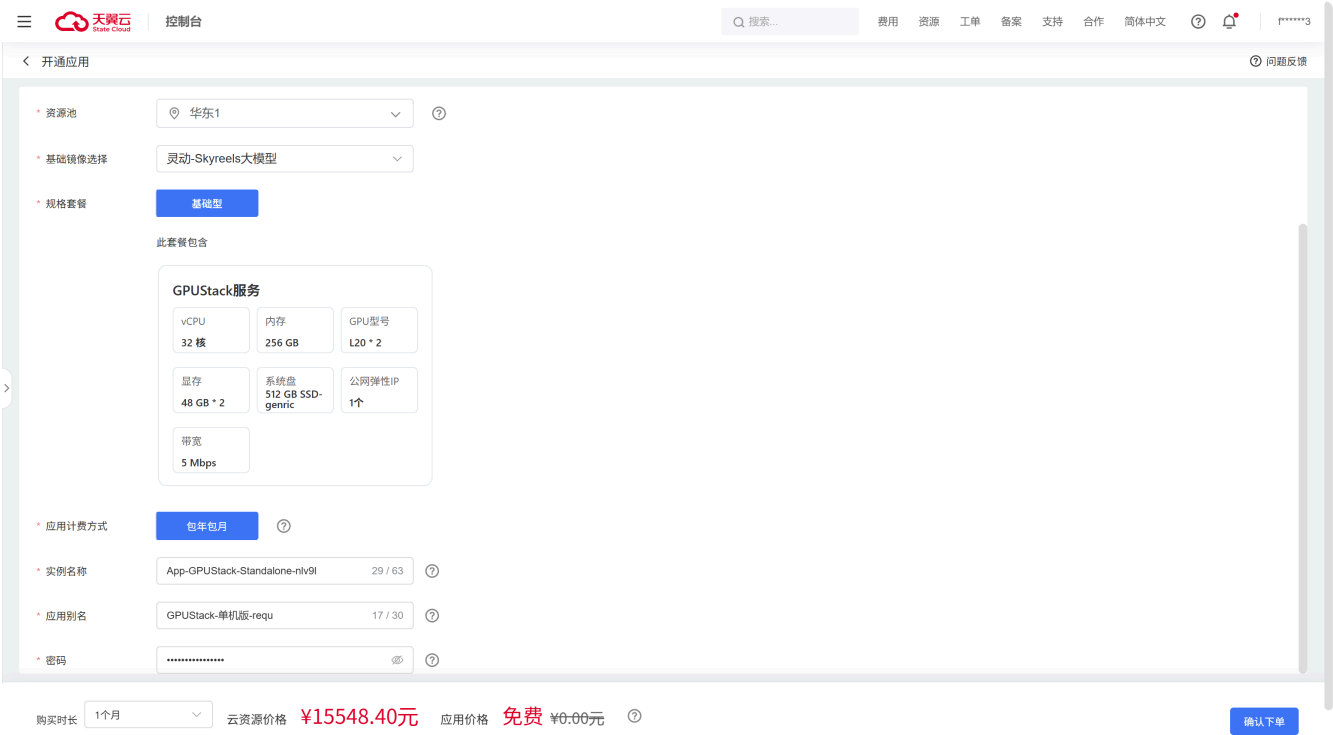
**运维能力:** 支持自动故障恢复、负载均衡、监控、认证和访问控制。

# 二、使用指南

## 2.1 登录天翼云官网，选择【应用商城】-【公有云生态专区】，点击立即选购，进入应用专区页



## 2.2 选择"GPUStack 单机版"，点击立即订购

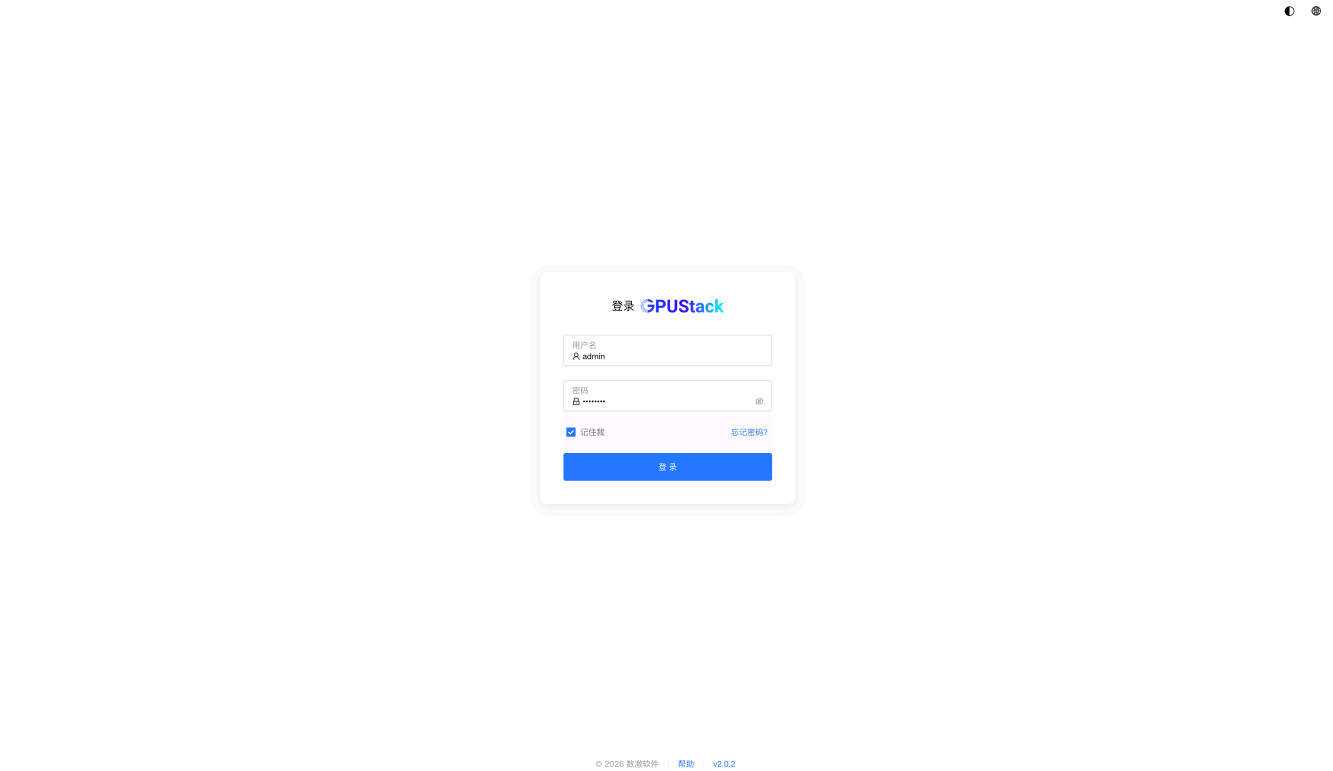## 2.3 按提示订购"GPUstack 单机版"，点击确认下单即可成功购买服务。



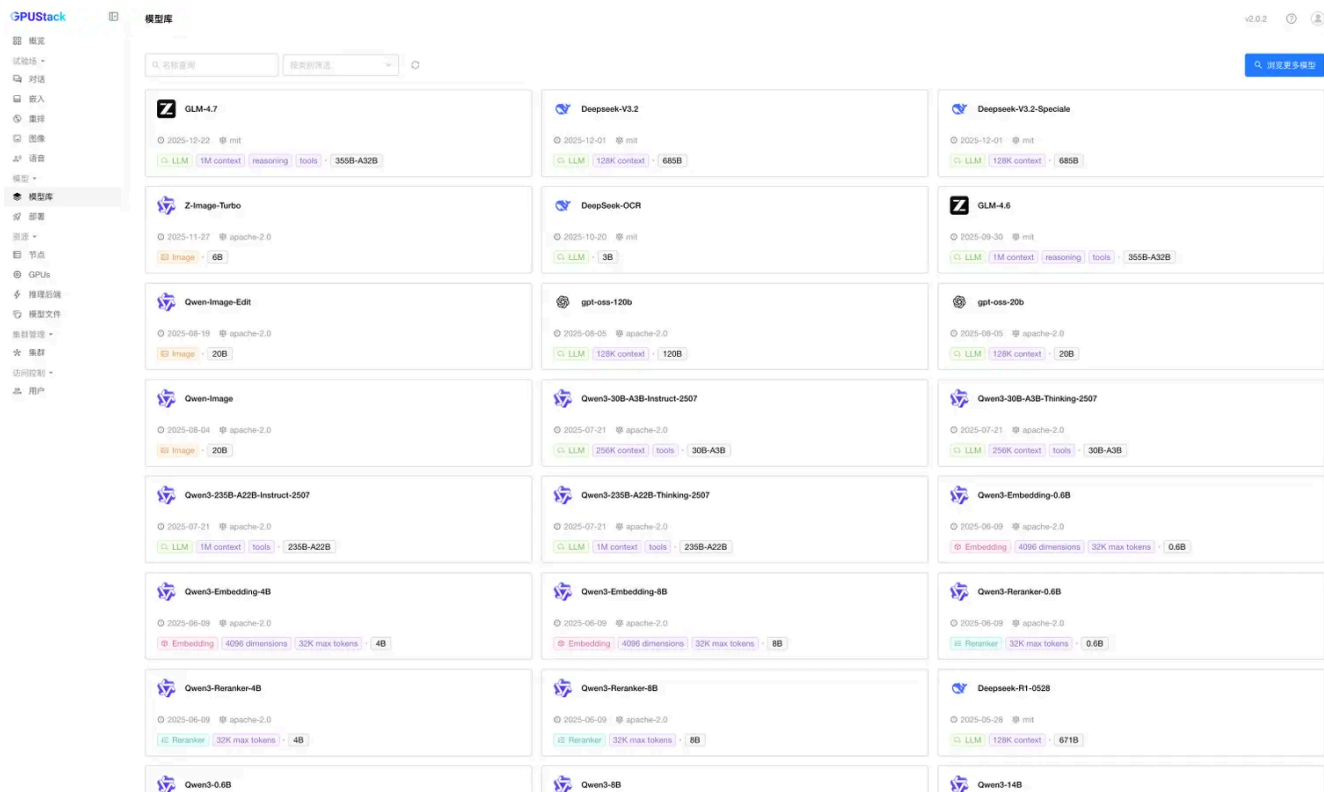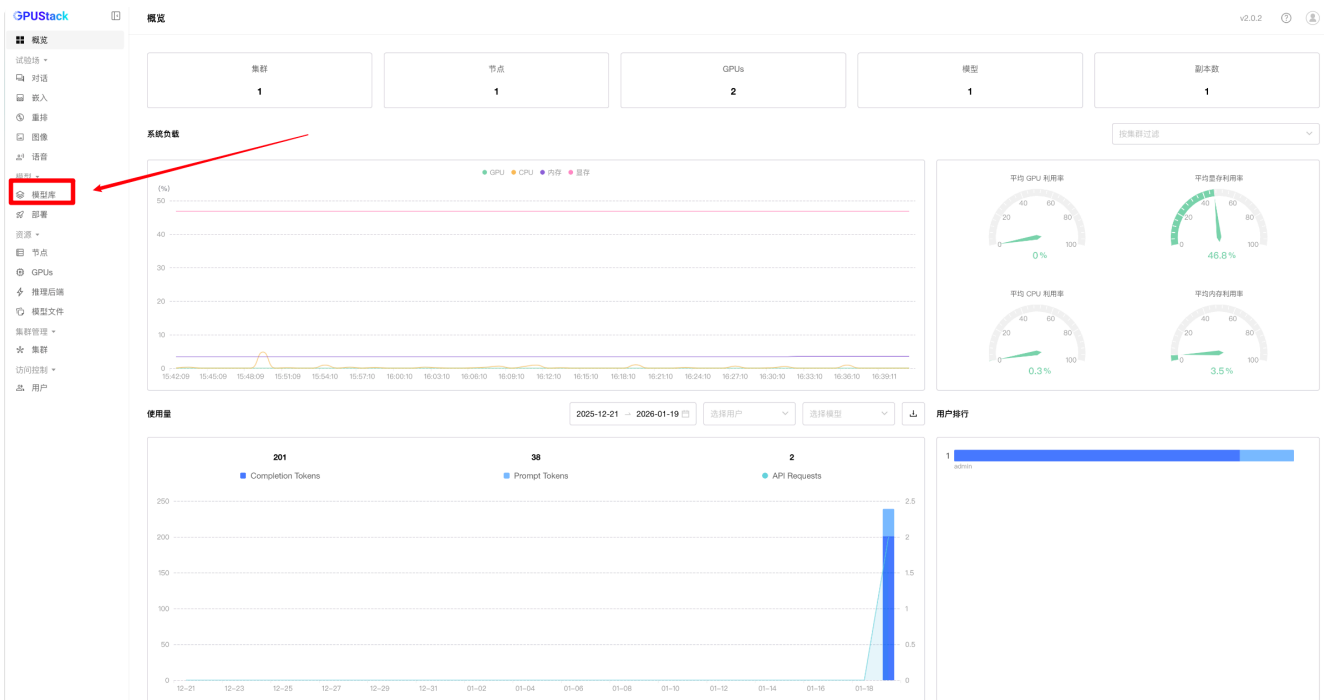## 2.4 在【我的应用】页，查看应用状态，当应用状态为运行中时，代表服务部署完成。
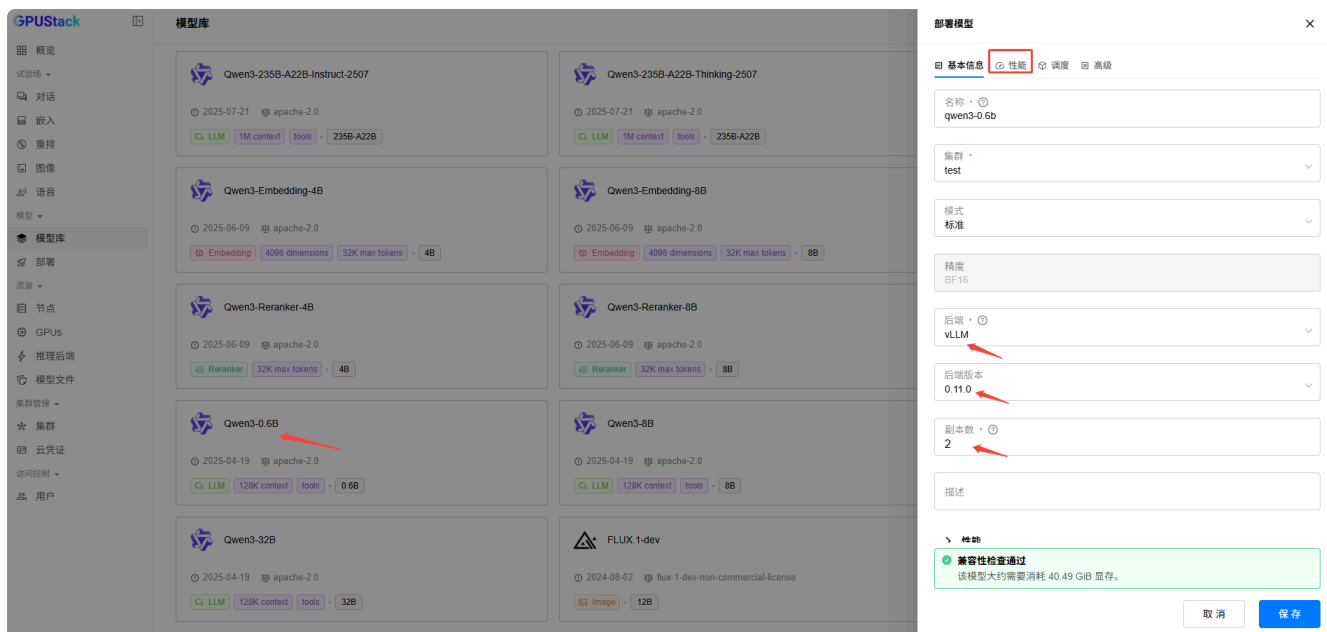
## 2.5 点击查看应用按钮，点击【立即使用】或在浏览器中输入复制的应用入口，即可访问 GPUstack
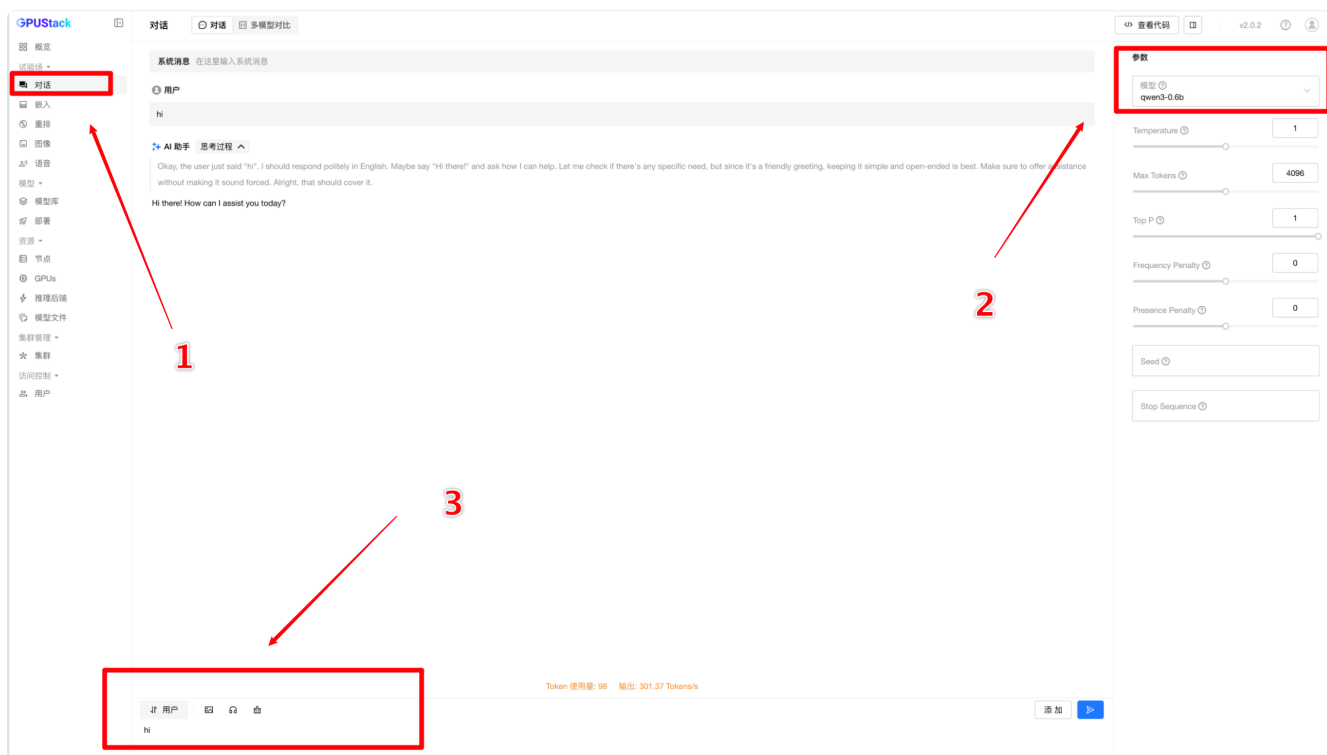


## 2.6 GPUstack 默认账号为 admin，密码为 gpustack。（可在应用中修改）



## 2.7 部署模型。进入模型库菜单，选择合适的模型开始部署；

# Screenshot 1 (概览)

GPUStack

- 概览
- 试验场
  - 对话
  - 嵌入
  - 重排
  - 图像
  - 语音
- 模型
  - 模型库
  - 部署
- 资源
  - 节点
  - GPUs
  - 推理后端
  - 模型文件
- 集群管理
  - 集群
- 访问控制
  - 用户

概览     v2.0.2

| 集群 | 节点 | GPUs | 模型 | 副本数 |
|---|---|---|---|---|
| 1 | 1 | 2 | 1 | 1 |

系统负载    按集群过滤

GPU CPU 内存 显存

平均 GPU 利用率 **0%**    平均显存利用率 **46.8%**

平均 CPU 利用率 **0.3%**    平均内存利用率 **3.5%**

使用量   2025-12-21 — 2026-01-19   选择用户   选择模型   ⬇    用户排行

| 201 | 38 | 2 |
|---|---|---|
| Completion Tokens | Prompt Tokens | API Requests |

admin

# Screenshot 2 (模型库)

模型库     v2.0.2

名称查询    搜索到筛选     浏览更多模型

**GLM-4.7**
2025-12-22   mit
LLM   1M context   reasoning   tools · 355B-A32B

**Deepseek-V3.2**
2025-12-01   mit
LLM   128K context · 685B

**Deepseek-V3.2-Speciale**
2025-12-01   mit
LLM   128K context · 685B

**Z-Image-Turbo**
2025-11-27   apache-2.0
Image · 6B

**DeepSeek-OCR**
2025-10-20   mit
LLM · 3B

**GLM-4.6**
2025-09-30   mit
LLM   1M context   reasoning   tools · 355B-A32B

**Qwen-Image-Edit**
2025-08-19   apache-2.0
Image · 20B

**gpt-oss-120b**
2025-08-05   apache-2.0
LLM   128K context · 120B

**gpt-oss-20b**
2025-08-05   apache-2.0
LLM   128K context · 20B

**Qwen-Image**
2025-08-04   apache-2.0
Image · 20B

**Qwen3-30B-A3B-Instruct-2507**
2025-07-21   apache-2.0
LLM   256K context   tools · 30B-A3B

**Qwen3-30B-A3B-Thinking-2507**
2025-07-21   apache-2.0
LLM   256K context   tools · 30B-A3B

**Qwen3-235B-A22B-Instruct-2507**
2025-07-21   apache-2.0
LLM   1M context   tools · 235B-A22B

**Qwen3-235B-A22B-Thinking-2507**
2025-07-21   apache-2.0
LLM   1M context   tools · 235B-A22B

**Qwen3-Embedding-0.6B**
2025-06-09   apache-2.0
Embedding   4096 dimensions   32K max tokens · 0.6B

**Qwen3-Embedding-4B**
2025-06-09   apache-2.0
Embedding   4096 dimensions   32K max tokens · 4B

**Qwen3-Embedding-8B**
2025-06-09   apache-2.0
Embedding   4096 dimensions   32K max tokens · 8B

**Qwen3-Reranker-0.6B**
2025-06-09   apache-2.0
Reranker   32K max tokens · 0.6B

**Qwen3-Reranker-4B**
2025-06-09   apache-2.0
Reranker   32K max tokens · 4B

**Qwen3-Reranker-8B**
2025-06-09   apache-2.0
Reranker   32K max tokens · 8B

**Deepseek-R1-0528**
2025-05-28   mit
LLM   128K context · 671B

**Qwen3-0.6B**

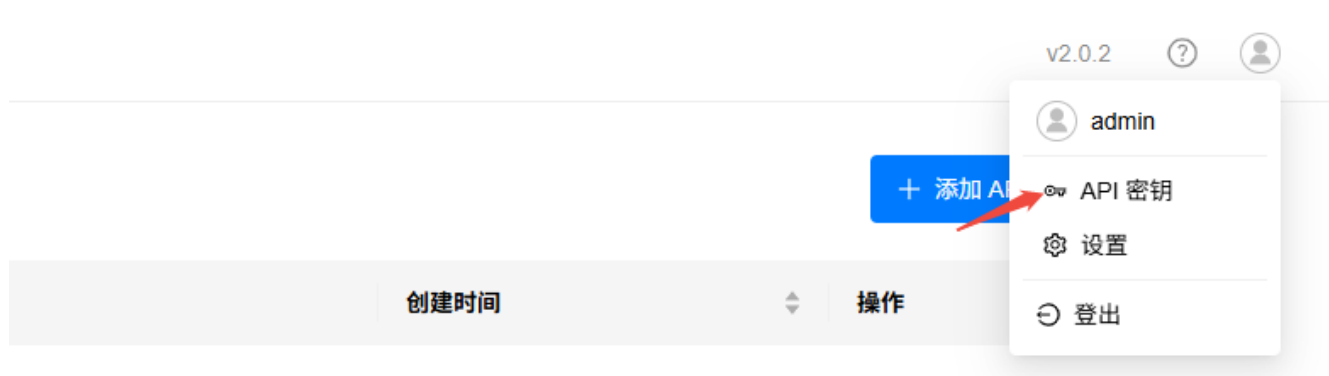**Qwen3-8B**

**Qwen3-14B**

## 2.8 等待模型部署成功后，便可进入试验场，进行验证；





## 2.9 API 密钥管理

1. 将鼠标悬停在用户头像上，选择" `API 密钥` "
2. 点击" `创建密钥` "生成新密钥
3. **安全提示**：该密钥仅在创建时可见一次，请妥善保存



# 三、常见问题

## Q1: 如何升级或降级已部署的模型版本?

**A**: GPUStack 支持无缝模型版本切换

1. 在模型详情页点击"编辑"按钮
2. 在"模型版本"下拉菜单中选择目标版本
3. 手动删除原来的示例，即可自动创建新的副本
4. 支持新旧副本共存，可以在用户无感知情况下切换版本

## Q2: 本地模型文件如何导入 GPUStack?

A：挂载共享存储，需要在集群创建之前完成挂载并存储，否则 worker 节点无法访问

1. 根据页面提示在创建集群时，进行缓存卷挂载，会自动生成对应的 worker 指令
2. 此时在部署模型时，可以选择本地模型进行部署

## 版权与声明