

# PR #25440 完整报告

sgl-project/sglang

Fix LoRA pool not appearing in /v1/loads

合并时间: 2026-05-16 09:21

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/25440>

## 执行摘要

- 一句话: 修复 /v1/loads 接口 LoRA 状态不显示 Bug
- 推荐动作: 该 PR 虽然代码量极小, 但修复了一个从 #16976 引入的长期 Bug, 对使用 LoRA 的部署具有实际意义。建议读者关注: 1) 幽灵属性 `lora_scheduler` 的来历 (或许可以通过仓库搜索确认其是否在其他地方被误用); 2) `hasattr` 防御式编程可能隐藏此类 Bug, 直接使用已知属性更安全。

## 功能与动机

`/v1/loads` 端点用于对外暴露负载指标, 但 LoRA 池相关信息始终缺失。原 PR body 明确指出 `lora_scheduler` 是一个不存在的幽灵名称 (phantom name), 由 #16976 引入, 导致 LoRA 被启用时端点仍然返回 LoRA 池不可用。

## 实现拆解

1. 定位问题: 在 `python/sglang/srt/observability/scheduler_metrics_mixin.py` 的 `get_loads` 方法中, LoRA 指标分支的守卫条件使用了 `hasattr(self, "lora_scheduler") and self.lora_scheduler is not None`。
2. 分析根因: `lora_scheduler` 从未在任何类上定义, `hasattr` 始终返回 `False`, 导致分支不可达。
3. 替换守卫: 将守卫条件改为 `if self.enable_lora:`。`self.enable_lora` 是 `Scheduler` 的真实属性, 准确反映 LoRA 是否启用, 且与 `Scheduler` 中其他 LoRA 相关逻辑的门控一致。
4. 验证: 由于是机械性修复, 且 `self.enable_lora` 已在多处使用, 风险极低。

关键文件:

- `python/sglang/srt/observability/scheduler_metrics_mixin.py` (模块 可观测性; 类别 source; 类型 core-logic; 符号 `get_loads`): 修复的核心文件, loRA 指标守卫条件从错误的 `hasattr(self, 'lora_scheduler')` 改为 `self.enable_lora`。

关键符号: `get_loads`

## 关键源码片段

`python/sglang/srt/observability/scheduler_metrics_mixin.py`

修复的核心文件，LoRA 指标守卫条件从错误的 `hasattr(self, 'lora_scheduler')` 改为 `self.enable_lora`。

```
# 路径 : python/sglang/srt/observability/scheduler_metrics_mixin.py
# 行 1043-1050: 修复 LoRA 指标守卫条件
```

```
lora = None
if include_all or "lora" in include:
    # 修复前 : hasattr(self, "lora_scheduler") and self.lora_scheduler is not None
    # lora_scheduler 是幽灵属性，从未被定义，导致此分支永远不执行
    # 修复后 : 使用 self.enable_lora，该属性由服务器参数真实设置
    if self.enable_lora:
        lora = LoRAMetrics(
            slots_used=self.stats.lora_pool_slots_used,
            slots_total=self.stats.lora_pool_slots_total,
            utilization=self.stats.lora_pool_utilization,
        )
```

## 评论区精华

无 review 讨论。PR 为自 merge（作者 merge），变更简单直接，未产生争议。

- 暂无高价值评论线程

## 风险与影响

- 风险：风险极低。修改仅涉及一行守卫条件，且 `self.enable_lora` 是早已存在的布尔属性，在 Scheduler 中多处使用。唯一潜在风险是如果存在某种奇怪场景下 `enable_lora` 为 True 但实际 LoRA 相关组件未正确初始化，不过这与守卫语义一致——既然 LoRA 已启用，就应该展示（或报错）而非静默跳过。
- 影响：影响范围：仅影响 `/v1/loads` 端点中 `lora` 字段的返回。此前 LoRA 池始终不展示（即使 LoRA 已启用），修复后正确返回 LoRA 利用率指标。用户影响：使 LoRA 用户能够通过 `/v1/loads` 获取真实池状态，便于监控和运维。系统影响：无性能或稳定性影响。回滚：可直接 revert 无副作用。
- 风险标记：缺少测试覆盖

## 关联脉络

- PR #16976 Add `/v1/loads` endpoint for load metrics: 引入原始 Bug 的 PR，使用不存在的 `lora_scheduler` 属性。