

# PR #22577 完整报告

sgl-project/sglang

Add hisparse staging + decode offload guards to `is_fully_idle()`

合并时间: 2026-04-11 15:12

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/22577>

## 执行摘要

- 一句话: 修复调度器空闲检测逻辑, 增加 HiSparse 和 KV 卸载的临时状态检查。
- 推荐动作: 建议相关模块的工程师精读此 PR, 了解调度器空闲检测的完整状态机。特别关注 `is_fully_idle()` 函数中各种异步操作的检查逻辑, 这对理解系统在复杂状态下的行为很重要。同时建议检查相关测试是否充分覆盖新增的状态检查。

## 功能与动机

PR body 明确指出: `is_fully_idle()` 函数在 HiSparse (#20343) 和解码 KV 卸载功能添加之前编写, 这些功能将空闲检查放在 `self_check_during_idle` 或各自的事件循环中, 但从未传播到 `is_fully_idle()`。这意味着 `is_fully_idle()` 的其他使用者 (`flush_cache`、`hicache attach/detach`、`release_memory_occupation`) 可能在内存池处于临时状态时执行, 存在破坏性操作风险。

## 实现拆解

在 `scheduler.py` 文件的 `is_fully_idle()` 函数中, 在 `if not for_health_check:` 分支下添加了两个检查: 1. 如果存在 `decode_offload_manager`, 检查 `ongoing_offload` 队列是否为空; 2. 如果启用 HiSparse, 检查 `hisparse_coordinator` 是否有进行中的暂存操作。这两个检查遵循现有的 HiCache 异步操作检查模式。

关键文件:

- `python/sglang/srt/managers/scheduler.py` (模块 `scheduler`): 这是调度器核心文件, `is_fully_idle()` 函数是关键的空闲检测逻辑, 直接影响多个系统操作的安全性。

关键符号: `is_fully_idle`

## 评论区精华

Review 讨论较少, 只有 `ispobock` 的批准评论, 没有具体技术讨论。PR body 中作者详细说明了动机和实现思路, 强调遵循现有 HiCache 模式。测试计划包括运行 `test_scheduler_pause_generation.py` 测试和依赖现有 CI 覆盖。

- 空闲检测逻辑完整性 (`correctness`): 添加了 `decode_offload_manager.ongoing_offload` 和 `hisparse_coordinator.has_ongoing_staging()` 检查, 遵循现有 HiCache 模式。

## 风险与影响

- 风险：风险较低：1. 变更范围小（仅 6 行代码），逻辑清晰；2. 遵循现有 HiCache 检查模式，保持一致性；3. 健康检查路径不受影响（for\_health\_check=True 时跳过新增检查）；4. 潜在风险是新增检查可能引入误判，导致 is\_fully\_idle() 返回 false positive，但考虑到这些检查针对的是已知的临时状态，这种风险较小。
- 影响：影响范围中等：1. 直接影响调度器的空闲检测逻辑，影响所有依赖 is\_fully\_idle() 的操作（flush\_cache、hicache attach/detach、release\_memory\_occupation）；2. 确保在 HiSparse 暂存和 KV 卸载的临时状态下不会执行破坏性操作，提升系统稳定性；3. 对用户透明，但可能影响系统在特定状态下的行为；4. 需要确保测试覆盖充分，特别是 HiSparse 和 decode offload 的相关场景。
- 风险标记：核心路径变更，状态机完整性

## 关联脉络

- PR #22453 [HiSparse-pd] Add device-buffer budget and fix logical pool admission in decode side: 同样涉及 HiSparse 功能，修改了 decode.py 文件，与本 PR 的 HiSparse 检查相关。
- PR #22554 [mem] Introduce PoolStats dataclass; unify pool metrics and token\_usage: 同样修改了 scheduler.py 文件，涉及调度器内存管理相关逻辑。
- PR #22559 [metrics] Add PoolStats.update\_scheduler\_stats to deduplicate metrics assignment: 同样修改了调度器相关文件（scheduler\_runtime\_checker\_mixin.py），涉及调度器状态管理。