

PR #24226 完整报告

sgl-project/sglang

[BugFix] fix(hicache): fix two slot-reuse races in DecodeKVCacheOffloadManager

合并时间: 2026-05-21 09:20

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/24226>

执行摘要

- 一句话: 修复 HiCache 解码卸载管理器 slot 重用竞态条件
- 推荐动作: 建议合入。这是一个关键正确性修复, 解决了 DecodeKVCacheOffloadManager 中导致数据污染的两个竞态条件, 并重新启用了被禁用的端到端测试。实现中通过 offload_inflight 计数器将释放时刻与异步操作完成关联, 是处理带外操作序列化的良好设计模式, 值得参考。

功能与动机

PR 描述指出 `DecodeKVCacheOffloadManager` 存在两个相关的 slot 重用竞态: 第一个是 `offload_kv_cache` 在请求首次增量 offload 时调用 `token_to_kv_pool_allocator.free(token_indices[:state.prefill_len])`, 但解码请求仍通过 `req_to_token_pool.req_to_token` 访问那些 slot, 导致并发新请求重用后产生跨请求 KV 读取污染; 第二个是请求结束后 `cache_controller` 可能有未完成的异步 device→host 拷贝, 过早释放 slot 会破坏拷贝内容。表现症状为持续负载下约 5% 的响应会中途退化为来自其他请求提示词的文本片段循环。

实现拆解

1. 延迟 prefill 对齐 slot 释放: 在 `offload_kv_cache` 中删除对 `token_indices[:state.prefill_len]` 的提前释放调用, 改为在 `_release_finished_req` 中统一释放。释放顺序为先 prefill 部分 (`[0:prefill_len]`) 再 committed 部分 (`[prefill_len:committed_len]`)。
2. 跟踪异步 offload 飞行计数: 在 `__init__` 中新增 `offload_inflight` 字典。每次调用 `offload_kv_cache` 时通过 `_mark_offload_started` 递增计数; 每次 `_check_offload_progress` 处理一个 ack 时通过 `_mark_offload_finished` 递减计数。计数归零时清理字典条目。提供 `_has_inflight_offload` 查询方法。
3. 调整释放触发时机: 在 `_check_offload_progress` 中, 之前会根据请求是否已完成决定立即释放或仅释放增量槽位。现在改为统一调用 `_mark_offload_finished`, 然后检查 `req.finished()` and `not self._has_inflight_offload(req.rid)`, 满足时才调用 `_release_finished_req`, 并传递正确的 `start_offset` (prefill 长度)。
4. 增强防御性检查: 在 `_release_finished_req` 入口添加对 `req.req_pool_idx` 是否为 `None` 或 `-1` 的判断, 防止因 `ReqToTokenPool.free` 已将索引设为 `None` 导致非幂等操作 (如 `tree_cache.protected_size_` 重复递减、`host pool` 重复释放) 重复执行。

5. 配套测试更新：在 `test_specv2_kvcache_offloading.py` 中新增 6 个单元测试用例，覆盖 `prefill` 存在时释放、`prefill` 长度为 0 时跳过、`finalize_release` 创建 `state` 等场景。同时移除 `test_disaggregation_decode_offload.py` 中的 `disabled` 标记，重新启用之前因该竞态问题被禁用的端到端测试。

关键文件：

- `python/sglang/srt/disaggregation/decode_kvcache_offload_manager.py`（模块 卸载管理；类别 `source`；类型 `core-logic`；符号 `_mark_offload_started`, `_mark_offload_finished`, `_has_inflight_offload`, `init`）：核心源文件，修复两个 `slot` 重用竞态：删除提前释放 `prefill slot`、引入 `offload_inflight` 计数器跟踪异步操作、调整释放逻辑到请求完成后。
- `test/registered/disaggregation/test_specv2_kvcache_offloading.py`（模块 单元测试；类别 `test`；类型 `test-coverage`；符号 `_FinishedEvent`, `synchronize`, `test_release_finished_req_frees_prefill_when_state_present`, `test_release_finished_req_skips_prefill_free_when_prefill_len_zero`）：新增 6 个单元测试，覆盖 `prefill` 释放、`inflight offload` 延迟释放等关键场景，确保修复逻辑正确且不易退化。
- `test/registered/disaggregation/test_disaggregation_decode_offload.py`（模块 端到端测试；类别 `test`；类型 `test-coverage`）：移除 `disabled` 标记，重新启用端到端测试，确认 `fix` 通过集成验证。

关键符号：`init`, `offload_kv_cache`, `_mark_offload_started`, `_mark_offload_finished`, `_has_inflight_offload`, `_check_offload_progress`, `_release_finished_req`, `finalize_release_on_finish`

关键源码片段

`python/sglang/srt/disaggregation/decode_kvcache_offload_manager.py`

核心源文件，修复两个 `slot` 重用竞态：删除提前释放 `prefill slot`、引入 `offload_inflight` 计数器跟踪异步操作、调整释放逻辑到请求完成后。

```
# decode_kvcache_offload_manager.py (关键片段)
```

```
class DecodeKVCacheOffloadManager:
    def __init__(self, ...):
        # ... existing init ...
        self.offload_inflight = {} # 新增：跟踪每个请求未完成的异步 offload 数量

    def _mark_offload_started(self, rid):
        """每次发起异步 offload 时递增计数"""
        self.offload_inflight[rid] = self.offload_inflight.get(rid, 0) + 1

    def _mark_offload_finished(self, rid):
        """每次 ack 到达时递减计数，计数归零则删除条目"""
        count = self.offload_inflight.get(rid, 0)
        if count <= 1:
```

```

        self.offload_inflight.pop(rid, None)
    else:
        self.offload_inflight[rid] = count - 1

def _has_inflight_offload(self, rid):
    return self.offload_inflight.get(rid, 0) > 0

def offload_kv_cache(self, req) -> bool:
    # ... 早期释放 prefill slot 的代码被删除, 改为注释说明 ...
    # 在发起异步 offload 前, 先标记开始
    self._mark_offload_started(req.rid)
    self.ongoing_offload[ack_id] = (req, host_indices, incremental_tokens, time.time(), start,
    end)
    state.inc_len += incremental_aligned_len
    return True

def _check_offload_progress(self, finish_count):
    # ... 处理 ack ...
    self._mark_offload_finished(req.rid) # 每处理一个 ack 就减少计数
    # 原先的提前释放分支被删除, 改为统一的条件判断
    if req.finished() and not self._has_inflight_offload(req.rid):
        state = self.offloaded_state.get(req.rid)
        start_offset = state.prefill_len if state is not None else start
        self._release_finished_req(req, start_offset)

def _release_finished_req(self, req, start_offset):
    # 防御性检查: 如果 req_pool_idx 已经被清理则跳过
    if req.req_pool_idx is None or req.req_pool_idx == -1:
        return
    # 先释放 prefill 部分 [0:start_offset]
    if start_offset > 0:
        self.token_to_kv_pool_allocator.free(
            self.req_to_token_pool.req_to_token[req.req_pool_idx, :start_offset]
        )
    # 再释放 committed 部分 [start_offset:committed_len]
    committed_len = req.pop_committed_kv_cache()
    if committed_len > start_offset:
        self.token_to_kv_pool_allocator.free(
            self.req_to_token_pool.req_to_token[req.req_pool_idx, start_offset:committed_len]
        )
    req.req_pool_idx = None # 防止重复释放

```

评论区精华

- Code Review 建议添加 req_pool_idx 防御检查: gemini-code-assist[bot] 在 _check_offload_progress 和 finalize_release_on_finish 的 diff hunks 上建议在调用 _release_finished_req 前增加对 req.req_pool_idx 是否为 None 或 -1 的检查, 以预防重复释放导致的双重递减或崩溃。该建议在第三个 commit 中被采纳, 在

_release_finished_req 入口添加了防御性 guard。

- 合并者对模块未来的展望：ShangmingCai 在批准时提到团队正在开发解码侧 HiCache 支持，该 offload manager 可能在任务完成后被替换，暗示当前修复是短期解决方案。
 - req_pool_idx 防御检查建议 (correctness): 已采纳：在 _release_finished_req 入口添加了防御性 guard (第三个 commit)。
 - finalize_release_on_finish 中的类似检查 (correctness): 已采纳：在 _release_finished_req 加统一 guard 后，finalize_release_on_finish 的调用路径也被保护。

风险与影响

- 风险：
 - 性能退化：slot 保留更久可能增加 GPU 内存压力，尤其在并发请求量大时。但 PR 作者指出这是避免数据损坏的必要代价，且释放时机优化不会显著影响吞吐。
 - 防御检查遗漏：若其他代码路径将 req_pool_idx 置为其他非法值（如 -2），仍可能触发错误。但现有检查覆盖了 None 和 -1，已涵盖已知路径。
 - 测试覆盖有限：单元测试基于 mock，未覆盖实际多线程并发场景。端到端测试已重新启用，但可能仍有未暴露的时序问题。
- 影响：
 - 用户影响：使用 HiCache 解码 offload 的用户在持续负载下将消除约 5% 的响应异常（输出中混入其他请求的文本）。对于对输出一致性要求高的场景，变更影响显著。
 - 系统影响：slot 保留更久可能导致 GPU memory pool 使用率轻微上升，但通常可承受。
 - 团队影响：修复降低了 offload 机制的已知稳定性隐患，为后续解码侧 HiCache 重构积累了更可靠的基线。
 - 风险标记：竞态条件修复，异步操作序列化，slot 释放时机

关联脉络

- PR #20622 Disable flaky decode offload e2e test: 该 PR 禁用了因相同竞态条件而失败的端到端测试。本 PR 修复了根本原因并移除 disabled 标记以重新启用该测试。