

PR #27038 完整报告

sgl-project/sglang

[sglang] Fix Mamba COW over-releasing SWA locks (cascade-evict assert crash)

合并时间: 2026-06-03 01:42

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/27038>

执行摘要

- 一句话: 修复 Mamba COW 路径 SWA 锁误释放
- 推荐动作: 强烈建议阅读此 PR, 它展示了一个易忽略的锁上下文传递问题, 属于典型的并发 bug 模式。设计上, lock_ref 的 inc/dec 需要严格配对且携带边界信息, 对理解 SGLang 缓存系统的锁模型很有帮助。建议作者补充单元测试, 验证 CoW 路径下锁计数的正确性。

功能与动机

Mamba 组件的 `finalize_match_result` 方法在 CoW 路径中, 为分配 Mamba 池槽位需要临时驱逐缓存节点。在调用 `inc_lock_ref` 后再调用 `dec_lock_ref` 时, 未传递锁上下文 (`swa_uuid_for_lock`), 导致 SWA 锁释放逻辑遍历到根节点, 错误地减少了其他请求的锁计数。这破坏了锁计数一致性, 最终触发断言失败和服务器崩溃。PR body 明确描述了该问题: “A missing uuid makes it walk all the way to root and over-decrement SWA locks”。

实现拆解

1. 在 `python/sglang/srt/mem_cache/unified_cache_components/mamba_component.py` 的 `finalize_match_result` 方法中, 将原先丢失返回值的 `self.cache.inc_lock_ref(last_node)` 调用改为捕获返回值 `lock_result`。
2. 将对应的 `self.cache.dec_lock_ref(last_node)` 调用改为 `self.cache.dec_lock_ref(last_node, lock_result.to_dec_params())`, 传递锁操作参数, 确保 SWA 释放时使用正确的 `swa_uuid_for_lock` 边界 token。
3. 新增的注释清晰解释了 bug 根因和作用: 如果缺失 `uuid`, SWA 释放会越过当前请求窗口边界一直走到根节点, 错误减少祖先节点上其他驻留请求的 SWA 锁计数。

关键文件:

- `python/sglang/srt/mem_cache/unified_cache_components/mamba_component.py` (模块缓存层; 类别 `source`; 类型 `core-logic`; 符号 `finalize_match_result`): 本次变更的唯一文件, 核心修复位置。在 `finalize_match_result` 方法的 CoW 路径中, 将 `inc_lock_ref` 返回值传递给 `dec_lock_ref`, 确保 SWA 锁释放正确。

关键符号: `finalize_match_result`

关键源码片段

python/sglang/srt/mem_cache/unified_cache_components/mamba_component.py

本次变更的唯一文件，核心修复位置。在 `finalize_match_result` 方法的 CoW 路径中，将 `inc_lock_ref` 返回值传递给 `dec_lock_ref`，确保 SWA 锁释放正确。

```
# python/sglang/srt/mem_cache/unified_cache_components/mamba_component.py
# 第 92-104 行，修复后的 CoW 路径关键片段
```

```
def finalize_match_result(self, result, params, value_chunks, best_value_len):
    # ... 省略前面的匹配逻辑 ...
    mamba_value = last_node.component_data[self.component_type].value
    if cow_mamba and mamba_value is not None:
        assert req is not None
        if req.mamba_pool_idx is None:
            dst_index = self.cache.req_to_token_pool.mamba_pool.alloc(1)
            if dst_index is None:
                # 修复：捕获 inc 结果，携带 swa_uuid_for_lock 边界信息给 dec
                # 如果没有 uuid，SWA 释放会越过窗口走到根节点，
                # 错误减少其他请求在祖先节点上的 SWA 锁计数
                lock_result = self.cache.inc_lock_ref(last_node)
                self.cache.evict(EvictParams(num_tokens=0, mamba_num=1))
                dst_index = self.cache.req_to_token_pool.mamba_pool.alloc(1)
                # 传递 to_dec_params() 确保释放停在窗口边界
                self.cache.dec_lock_ref(last_node, lock_result.to_dec_params())
                assert dst_index is not None, "Can not alloc mamba cache"
                req.mamba_pool_idx = dst_index[0]
            req.mamba_cow_src_index = mamba_value
            req.mamba_needs_clear = False
```

评论区精华

Review 中 hzh0425 提问：“In what scenarios would these two components be used together? It seems they're always used separately.” 作者 bixue2010 回复：“yea, in our model, we are using them together.” 这表明 Mamba 和 SWA 同时使用的场景主要来自作者的模型，可能非主流配置，但修复是普遍必需的。Hzh0425 和 ispobock 均批准了 PR。

- Mamba 和 SWA 同时使用的场景 (question): 作者 bixue2010 回复：在他们的模型中，两者是同时使用的。

风险与影响

- 风险：变更仅修改了 CoW 路径中的锁释放调用，且传递的参数来自刚刚的 `inc_lock_ref` 返回值，逻辑上正确且安全。但增加参数传递可能对 `dec_lock_ref` 内部的锁释放行为产生影响，需要保证 `to_dec_params()` 返回的对象与 `dec_lock_ref` 期望的接口完全兼容。由于代码库中其他所有锁释放调用都已正确传递 `uuid`，此修复只是补齐了遗漏路径，回归风险低。未增加单元测试，建议后续补充针对 CoW 路径的锁计数一致性测试。

- 影响：直接影响：修复了 Mamba + SWA 混合使用场景下因锁计数不一致导致的服务器崩溃（cascade-evict assert crash）。影响范围仅限于使用 Mamba 缓存且同时启用 SWA 的模型（如作者的模型）。对于不触发 CoW 路径的场景，无影响。由于 Mamba 和 SWA 通常不在同一模型中同时使用，大多数用户不会遇到此问题。
- 风险标记：缺少测试覆盖

关联脉络

- PR #27064 Fix stale import after kl_nightly rename: 涉及相同的 Mamba 缓存测试（test_unified_radix_cache_kl_mamba.py），表明 Mamba 缓存模块处于活跃维护中。
- PR #26948 Improve type annotations in unified radix cache: 与 unified radix cache 相关，本次修改的 mamba_component 属于其组件。
- PR #27070 Relax mamba unified cache kl threshold: 修改了 Mamba 缓存测试的 KL 阈值，涉及同一缓存组件。