

PR #27459 完整报告

sgl-project/sglang

[core] Probe `set_kv_buffer` / `set_mla_kv_buffer` slot ids for OOB

合并时间: 2026-06-07 10:51

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/27459>

执行摘要

- 一句话: 在 KV 写入路径添加越界探测
- 推荐动作: 值得合并。这是一个低风险、高 ROI 的调试增强, 在 spec 测试中已证明有效。建议在更广泛的 CI 中逐步启用 SGLANG_ENABLE_ASYNC_ASSERT, 以最大化收益。

功能与动机

`set_kv_buffer` 是 KV 缓存的频繁写入路径, 但此前对 `loc` 参数没有任何值检查。如果 `loc` 含有过时的 slot id, 可能导致非法地址崩溃或静默破坏相邻 KV (问题远距离显现, 难以调试)。PR body 明确指出“`move_kv_cache` already probes its `tgt_loc` / `src_loc` the same way, but the much hotter `set_kv_buffer` write path had no value check on `loc`”。

实现拆解

1. 在 `MHATokenToKVPool.set_kv_buffer` 中添加探测: 在函数入口处插入 `maybe_detect_oob(loc, 0, self.size + self.page_size, "set_kv_buffer (MHA)")`, 在 KV 存储前捕获越界 slot id。
2. 在 `MHATokenToKVPoolFP4.set_kv_buffer` 中添加探测: 同样的逻辑, 标记为 "`set_kv_buffer (MHA-FP4)`"。
3. 在 `MLATokenToKVPool.set_kv_buffer` 和 `set_mla_kv_buffer` 中添加探测: 两个方法分别标记为 "`set_kv_buffer (MLA)`" 和 "`set_mla_kv_buffer (MLA)`", 覆盖 MLA 场景。
4. 在 `MLATokenToKVPoolFP4.set_kv_buffer` 和 `set_mla_kv_buffer` 中添加探测: 标记为 "`set_kv_buffer (MLA-FP4)`" 和 "`set_mla_kv_buffer (MLA-FP4)`", 覆盖 MLA + FP4 场景。

边界 `self.size + self.page_size` 与 `move_kv_cache` 一致, 表示物理缓冲区高度, 包含保留的填充尾部和 CUDA graph 的零填充行, 不会误报。

关键文件:

- `python/sglang/srt/mem_cache/memory_pool.py` (模块内存池; 类别 source; 类型 core-logic): 所有变更集中在此文件, 在 `set_kv_buffer` 和 `set_mla_kv_buffer` 的四个变体方法中插入越界探测。

关键符号: `set_kv_buffer`, `set_mla_kv_buffer`, `maybe_detect_oob`

关键源码片段

python/sglang/srt/mem_cache/memory_pool.py

所有变更集中在此文件，在 `set_kv_buffer` 和 `set_mla_kv_buffer` 的四个变体方法中插入越界探测。

```
# python/sglang/srt/mem_cache/memory_pool.py
# 在 MHATokenToKVPool.set_kv_buffer 入口处 (第 1177 行)
def set_kv_buffer(self, layer, loc, cache_k, cache_v, ...):
    # Catch stale slot ids here instead of as illegal-addr / silent KV
    # corruption in the store_kvcache write (gated on SGLANG_ENABLE_ASYNC_ASSERT).
    # 边界 self.size + self.page_size 与 move_kv_cache 使用的一致,
    # 包含 padding 和 cuda-graph 零填充行, 避免误报。
    maybe_detect_oob(loc, 0, self.size + self.page_size, "set_kv_buffer (MHA)")
    ... # 后续原有逻辑

# 同样在 MLATokenToKVPool.set_kv_buffer 入口处 (第 1867 行)
def set_kv_buffer(self, layer, loc, cache_k, cache_v, ...):
    maybe_detect_oob(loc, 0, self.size + self.page_size, "set_kv_buffer (MLA)")
    ...

# 其余两个 FP4 变体完全一致, 仅标签字符串不同
```

评论区精华

PR 讨论较少 (仅 3 条 comments, 其中 2 条为自动化 bot 回复), 无实质性 review 争议。作者在 PR body 中已清晰说明设计决策: 边界选择、门控条件、与 `move_kv_cache` 的一致性等等。

- 暂无高价值评论线程

风险与影响

- 风险: 风险极低。变更仅在 `SGLANG_ENABLE_ASYNC_ASSERT` 启用时生效 (默认关闭), 生产环境零开销。同时边界 `self.size + self.page_size` 与已有探测一致, 不会引入误报。唯一的潜在风险是如果未来修改了物理缓冲区布局但忘记更新此边界, 但概率低且会被后续测试覆盖。
- 影响:
 - 用户: 无影响, 默认不启用。
 - 系统: 在 spec CI 中默认启用后, 有助于调试阶段更快定位 KV 缓存越界问题, 减少排查时间。
 - 团队: 提高了 KV 缓存路径的可观测性, 降低调试难度。但需要确保 CSpec 之外的测试场景也适时启用该标志。
 - 风险标记: 门控默认关闭, 仅调试增强

关联脉络

- PR #27461 Enable async-assert invariant probes by default in CI: 本 PR 的探测依赖 SGLANG_ENABLE_ASYNC_ASSERT, 而 PR 27461 正好在 CI 中默认启用了该标志, 两者协同工作。
- PR #27364 [perf] reduce radix cache match overhead by changing the match algorithm: 同一仓库中对 KV 缓存路径的改动, 与本 PR 共享 memory_pool.py 文件, 且关注点类似 (避免非预期行为)。