

PR #24954 完整报告

sgl-project/sglang

[Mamba] Fix extra_buffer overlap schedule races

合并时间: 2026-05-19 12:13

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/24954>

执行摘要

- 一句话: 将 Mamba 状态操作迁移到 forward stream 消除调度竞争
- 推荐动作: 该 PR 值得所有关注并发调度和 Mamba 模型的开发者精读, 其“捐赠模式”和“延迟操作到 forward stream”是处理调度器与前向流之间竞争的有效模式。建议合并后关注 HiCache 兼容性修复和 review 中提到的 GPU→CPU 同步优化。

功能与动机

关联 Issue #24221 描述了在 overlap scheduler 中, Mamba radix cache 可能在前向 pass 尚未完成写入时快照 Mamba 临时状态, 导致 radix cache 条目损坏。具体路径是 chunked prefill 的 stash 操作没有等待 copy_done 同步。本 PR 通过将所有 Mamba 状态操作移到 forward stream, 从根本上消除竞争窗口。

实现拆解

1. 将 MambaPool.alloc 中的立即清零操作拆分为独立的 clear_slots 方法, 原 alloc 不再执行清零, 让调度器可以分配插槽而不触发 GPU 操作。
2. 在 cache_unfinished_req 中, 将来自 req 的 Mamba 状态索引“捐赠”给 radix cache (即直接转移所有权), 然后为请求分配新的空白插槽; 替换原来的 fork_from 复制模式, 避免在读缓存时与 forward stream 产生数据竞争。
3. 在 prefix match 的 COW 路径中, 不在 scheduler stream 上立即复制, 而是将源索引记录到 req.mamba_cow_src_index 和 req.mamba_needs_clear, 实际复制 / 清零延迟到 forward stream 的 init_forward_metadata 阶段。
4. 在 hybrid_linear_attn_backend (以及 Mamba2 后端) 的 init_forward_metadata 中调用新方法 _execute_deferred_mamba_cow_and_clear, 执行之前收集到的清除和复制操作; 同时利用 is_draft_worker 跳过推测解码的 draft 阶段, 防止重复执行。
5. 提取公共辅助函数 set_mamba_track_indices_from_reqs 以减少 eagle_info_v2 等位置对 Mamba 追踪索引的重复计算; 统一 copy_from 接口支持批量索引; 删除不再使用的 fork_from 方法。

关键文件:

- python/sglang/srt/managers/schedule_batch.py (模块 调度器; 类别 source; 类型 core-logic; 符号 set_mamba_track_indices_from_reqs, _collect_deferred_mamba_cow_and_clear, prepare_for_split_prefill) : 核心调度批次类

，添加延迟 COW/clear 字段，新增收集和提取辅助函数，是 PR 主要逻辑入口

- python/sglang/srt/mem_cache/memory_pool.py (模块 缓存池; 类别 source; 类型 core-logic; 符号 clear_slots, copy_from, fork_from) : 将 alloc 中的清零独立成 clear_slots, 删除 fork_from, 修改 copy_from 支持批量操作, 为延迟执行奠定基础
- python/sglang/srt/mem_cache/mamba_radix_cache.py (模块 缓存树; 类别 source; 类型 core-logic; 符号 _alloc_mamba_slot) : cache_unfinished_req 采用捐赠模式替代 fork 复制, 避免与 forward stream 竞争
- 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 _alloc_mamba_slot) : 适配统一缓存的 Mamba 组件, 采用捐赠模式并添加 _alloc_mamba_slot
- python/sglang/srt/layers/attention/hybrid_linear_attn_backend.py (模块 注意力后端; 类别 source; 类型 core-logic; 符号 _execute_deferred_mamba_cow_and_clear) : 添加 _execute_deferred_mamba_cow_and_clear 在 forward stream 上执行延迟操作
- python/sglang/srt/model_executor/forward_batch_info.py (模块 前向元数据; 类别 source; 类型 data-contract) : 为 ForwardBatch 添加延迟操作字段, 连接调度和 forward stream
- python/sglang/srt/speculative/eagle_info_v2.py (模块 推测解码; 类别 source; 类型 dependency-wiring) : 使用公共辅助函数 set_mamba_track_indices_from_reqs, 消除重复代码
- python/sglang/srt/mem_cache/hi_mamba_radix_cache.py (模块 缓存树; 类别 source; 类型 core-logic) : 适配 HiCache 的捐赠模式变更

关键符号: set_mamba_track_indices_from_reqs,
_collect_deferred_mamba_cow_and_clear, clear_slots, copy_from, _alloc_mamba_slot,
_execute_deferred_mamba_cow_and_clear, prepare_for_split_prefill

关键源码片段

python/sglang/srt/managers/schedule_batch.py

核心调度批次类, 添加延迟 COW/clear 字段, 新增收集和提取辅助函数, 是 PR 主要逻辑入口

```
def set_mamba_track_indices_from_reqs(batch):
    """从请求对象构建 mamba_track_indices (权威来源)。
    避免之前在 eagle_info 中的重复实现, 并确保索引来源一致。
    """
    req_to_token_pool = batch.req_to_token_pool
    # 获取所有请求的 ping-pong 映射缓冲区, 形状 (bs, ping_pong_size)
    all_buffers = req_to_token_pool.req_index_to_mamba_ping_pong_track_buffer_mapping[
        batch.req_pool_indices
    ] # shape: (bs, ping_pong_size), int64, on device
    # 从每个请求的 mamba_next_track_idx 构建列索引
    idx = (
        torch.tensor(
            [req.mamba_next_track_idx for req in batch.reqs],
            dtype=torch.int64,
```

```

        pin_memory=True,
    )
    .unsqueeze(1) # shape: (bs, 1)
    .to(device=all_buffers.device, non_blocking=True)
)
# 通过 gather 取出正确的跟踪索引, 结果形状 (bs,)
batch.mamba_track_indices = (
    torch.gather(all_buffers, 1, idx).squeeze(1).to(torch.int64)
)

```

python/sglang/srt/mem_cache/memory_pool.py

将 alloc 中的清零独立成 clear_slots, 删除 fork_from, 修改 copy_from 支持批量操作, 为延迟执行奠定基础

```

def clear_slots(self, indices: torch.Tensor):
    """Zero out mamba state at the given pool indices. 必须在 forward stream 上执行。"""
    need_size = len(indices)
    # 清除 conv state
    for i in range(len(self.mamba_cache.conv)):
        t = self.mamba_cache.conv[i]
        # 扩展零张量以适应多个插槽, 避免 CPU-GPU 同步
        z = torch.zeros(1, dtype=t.dtype, device=t.device).expand(
            t.shape[0], need_size, *t.shape[2:]
        )
        t[:, indices] = z
    # 清除 temporal state
    t = self.mamba_cache.temporal
    z = torch.zeros(1, dtype=t.dtype, device=t.device).expand(
        t.shape[0], need_size, *t.shape[2:]
    )
    t[:, indices] = z

def copy_from(self, src_indices: torch.Tensor, dst_indices: torch.Tensor):
    """从源索引复制 mamba 状态到目标索引。在 forward stream 上执行。"""
    for i in range(len(self.mamba_cache.conv)):
        self.mamba_cache.conv[i][:, dst_indices] = self.mamba_cache.conv[i][
            :, src_indices
        ]
    self.mamba_cache.temporal[:, dst_indices] = self.mamba_cache.temporal[
        :, src_indices
    ]

```

评论区精华

- GPU→CPU 同步性能问题: gemini-code-assist 指出 `_collect_mamba_init_info` 中的 `all_pool_idx.tolist()` 和 `req.mamba_cow_src_index.item()` 引入 GPU→CPU 同步, 建议使用 GPU 原生方法。作者未回应, 该问题未解决。

- COW 执行位置: ispobock 询问是否可将 COW 移出 attention backend, hanming-lu 认为放入 init_forward_metadata (forward stream) 是合适的, 因为只有 linear backend 涉及 mamba。双方达成一致保留在 attention backend。
- disable extra_buffer 时的 KL 退化: opherlieber 报告在 extra_buffer 禁用时 chunked prefill 测试 KL 从 0.0006 退化到 0.06, 根源是延迟的 COW/clear 未正确传递。hanming-lu 通过 commit a3a2b4b 修复, 后续测试通过。
- HiCache 兼容性: hzh0425 提醒与 HiCache 的 Mamba 加载逻辑冲突, hanming-lu 承认并等待后续兼容方案。
 - GPU→CPU 同步性能问题 (performance): 作者未直接回应, 该性能问题在 PR 中未解决, 可能留待后续优化。
 - COW 执行位置讨论 (design): 达成一致保留在 attention backend, 因为该处是 forward stream 入口。
 - disable extra_buffer 时的 KL 退化 (correctness): hanming-lu 通过 commit a3a2b4b (fix no_buffer) 修复, 后续测试通过。

风险与影响

- 风险:
 1. 性能风险: 新的 _collect_mamba_init_info 方法引入 GPU→CPU 同步点 (tolist/item), 在调度器高负载时可能成为瓶颈。
 2. HiCache 兼容性: hzh0425 指出本 PR 与 HiCache 的 Mamba 加载逻辑冲突, 可能导致 HiCache 测试失败。虽然 PR 已合并, 但兼容性修复仍在待办。
 3. extra_buffer 禁用路径: 早期版本禁用 extra_buffer 时出现 KL 退化, 虽已修复但此路径的测试覆盖相对薄弱。
 4. MTP 推测解码: draft forward 可能重复执行延迟操作, 已通过清空 mamba_clear_indices 等字段修复, 但若未来引入新的 forward 路径需注意。- 影响: 用户: 使用 Mamba 模型时, overlap scheduler 下不再出现随机数据损坏导致的生成错误或悬挂。系统: 调度器与 forward stream 的隐式同步减少, 缓存操作更清晰。负面影响是 forward stream 增加了额外的清零 / 复制操作, 但只在 extend batch 时执行, 量级可控。团队: 代码可维护性提升, 重复逻辑被提取为公共函数。后续 HiCache 适配需要额外工作。
- 风险标记: 核心路径变更, HiCache 兼容性风险, extra_buffer 禁用路径退化, MTP draft 重复执行风险

关联脉络

- 暂无明显关联 PR