

PR #27330 完整报告

sgl-project/sglang

[UnifiedTree]: Fix CP Reduce

合并时间: 2026-06-05 14:03

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/27330>

执行摘要

- 一句话: 修复 CP 模式下 `all_reduce/barrier` 的通信组
- 推荐动作: 此 PR 是必要的修复, 设计清晰, 推荐合入。可作为统一通信组调用的好例子。

功能与动机

PR title 明确为 "Fix CP Reduce", 表明问题存在于 CP (Context Parallel) 模式。在 CP 模式下, attention 通信组 (attention group) 可能比全局 `tp_group` 更小, 原代码硬编码使用 `self.tp_group` 进行 `all_reduce` 和 `barrier` 会导致通信错误。此次修复将通信操作统一委托给已经存在的 `_all_reduce_attn_groups` 和 `_barrier_attn_groups` 方法, 这些方法能正确识别并使用正确的通信组。

实现拆解

1. 统一 `all_reduce` 调用: 在 `can_terminate_prefetch` 和 `check_prefetch_progress` 中, 将原本条件判断 `if self.tp_world_size > 1:` 内的 `torch.distributed.all_reduce(states/..., group=self.tp_group)` 替换为直接调用 `self._all_reduce_attn_groups(states, torch.distributed.ReduceOp.MAX/MIN)`, 无需再手动检查 `tp_world_size`。
2. 统一 `barrier` 调用: 在 `release_aborted_request` 中, 将条件内的 `torch.distributed.barrier(group=self.tp_group)` 替换为 `self._barrier_attn_groups()`。
3. 统一 `drain` 中的 `all_reduce`: 在 `drain_storage_control_queues` 中同样替换了 `all_reduce` 调用。
4. 删除冗余条件: 移除了 `if self.tp_world_size > 1` 的条件包裹, 因为 `_all_reduce_attn_groups` 方法内部已处理单卡情况。

关键文件:

- `python/sglang/srt/mem_cache/unified_radix_cache.py` (模块 缓存层; 类别 `source`; 类型 `core-logic`; 符号 `can_terminate_prefetch`, `check_prefetch_progress`, `release_aborted_request`, `drain_storage_control_queues`): 唯一修改的文件, 修复了 `can_terminate_prefetch`、`check_prefetch_progress`、`release_aborted_request`、`drain_storage_control_queues` 四个方法中错误的通信组使用。

关键符号: `can_terminate_prefetch`, `check_prefetch_progress`, `release_aborted_request`, `drain_storage_control_queues`

关键源码片段

python/sglang/srt/mem_cache/unified_radix_cache.py

唯一修改的文件，修复了 `can_terminate_prefetch`、`check_prefetch_progress`、`release_aborted_request`、`drain_storage_control_queues` 四个方法中错误的通信组使用。

```
def can_terminate_prefetch(self, operation: PrefetchOperation) -> bool:
    # ... 前面的逻辑不变
    states = torch.tensor(
        [1 - int(can_terminate), int(operation_terminated)],
        dtype=torch.int,
    )
    # 之前是 :
    # if self.tp_world_size > 1:
    # torch.distributed.all_reduce(states, op=MAX, group=self.tp_group)
    # 现在使用统一方法，内部会处理 CP 通信组
    self._all_reduce_attn_groups(states, torch.distributed.ReduceOp.MAX)
    can_terminate = states[0].item() == 0
    operation_terminated = states[1].item() == 1
    return can_terminate or operation_terminated

def check_prefetch_progress(self, req_id: str) -> bool:
    # ... 前面的逻辑不变
    if self.tp_world_size > 1:
        # ... 构造 packed tensor
        # 之前是 :
        # torch.distributed.all_reduce(packed, op=MIN, group=self.tp_group)
        self._all_reduce_attn_groups(packed, torch.distributed.ReduceOp.MIN)
        min_completed_tokens = int(packed[0].item())
        # ...

def release_aborted_request(self, rid: str) -> None:
    # ... 前面的逻辑不变
    completed_tokens, _ = self.cache_controller.terminate_prefetch(operation)
    # 之前是 :
    # if self.tp_world_size > 1:
    # torch.distributed.barrier(group=self.tp_group)
    self._barrier_attn_groups() # 统一 barrier
    self.dec_host_lock_ref(last_host_node, anchor_lock_params)
    # ...

def drain_storage_control_queues(self) -> None:
    # ... 前面的逻辑不变
    # 之前是 :
    # if self.tp_world_size > 1:
    # torch.distributed.all_reduce(qsizes, op=MIN, group=self.tp_group)
    self._all_reduce_attn_groups(qsizes, torch.distributed.ReduceOp.MIN)
    qsize_list = list(map(int, qsizes.tolist()))
    # ...
```

评论区精华

该 PR 未产生 review 讨论。单次 commit 直接由作者合并，可能为小范围修复。

- 暂无高价值评论线程

风险与影响

- 风险：低风险，因为替换的方法 `_all_reduce_attn_groups` 和 `_barrier_attn_groups` 已在其他位置使用且经过测试，内部会处理通信组的选择和单 / 多卡的判断。变更仅移除条件检查和替换调用，不改变逻辑语义。但需确认所有调用路径中 `self._all_reduce_attn_groups` 和 `self._barrier_attn_groups` 的正确初始化。
- 影响：影响范围局限于 `unified_radix_cache.py` 中的 4 个方法，主要影响 CP 模式下统一缓存层的通信行为。非 CP 模式下行为不变（因为方法内会 fallback 到 `tp_group`）。对用户无直接感知，但对 CP 模式的正确性至关重要。
- 风险标记：无测试覆盖，核心路径变更

关联脉络

- PR #27264 [UnifiedTree]: Sync sidecar component hits across TP ranks and make SWA prefetch all-or-nothing: 同属 UnifiedTree 缓存层改造，涉及跨 TP 同步和通信组调整。