

# PR #26394 完整报告

sgl-project/sglang

[PD] Fix cross-rank queue divergence by gating metadata readiness before all-reduce

合并时间: 2026-05-26 21:59

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/26394>

## 执行摘要

- 一句话: 修复 PD 跨 rank 队列发散导致挂起的问题
- 推荐动作: 建议仔细阅读 `utils.py` 中 `_apply_metadata_gate` 的设计: 通过检查 `bootstrap_room` 而非额外 `all-reduce` 来同步元数据状态, 降低了通信开销, 是一个精巧的优化。同时注意 `_commit_transfer_to_req` 中 `conditional` 从重试转为直接 `abort` 的决策, 明确了前置条件。

## 功能与动机

在 PD 模式中, `poll_and_all_reduce` 返回 `KVPoll.Success` 只表示 KV 缓存传输完成, 但元数据通过独立的 RDMA 操作写入, 可能尚未在所有 TP rank 上可见。若某个 rank 因元数据就绪而将请求移出队列, 其他 rank 却未就绪, 则队列会发散, 后续 `all-reduce` 操作在不对齐的张量上运行, 导致挂起或静默损坏。

## 实现拆解

1. 在 `utils.py` 中添加元数据就绪检查: 新增 `_is_fake_transfer` 函数 (从 `decode.py` 迁移) 判断是否为模拟传输; 新增 `_apply_metadata_gate` 函数, 遍历 `poll` 结果, 对标记为 `Success` 的请求检查 `bootstrap_room` 字段, 若为 0 (元数据未写入) 则降级为 `Transferring`。该函数在 `all-reduce` 前被调用, 依赖分布式 `ReduceOp.MIN` 的语义自动约束。
2. 改造 `poll_and_all_reduce` 与 `poll_and_all_reduce_with_staging`: 两个函数新增 `decode_reqs`、`metadata_buffers`、`server_args` 可选参数, 当这些参数都非空时, 在 `all-reduce` 前调用 `_apply_metadata_gate`, 使得仅在所有 rank 的元数据就绪后状态才为 `Success`。
3. 在 `decode.py` 中移除本地定义、调整提交逻辑: 将 `_is_fake_transfer` 的本地定义迁移至 `utils.py`, 改从 `utils` 导入; 修改 `_commit_transfer_to_req` 方法: 原逻辑在 `actual_room == 0` 时返回 `False` 并等待下一轮 `poll`, 现改为记录 `error` 并直接中止请求 (因为 `gate` 已确保元数据就绪后才到达此路径), 从而消除队列不收敛的风险。

关键文件:

- `python/sglang/srt/disaggregation/utils.py` (模块 调度器; 类别 `source`; 类型 `core-logic`; 符号 `_apply_metadata_gate`, `_is_fake_transfer`, `poll_and_all_reduce`, `poll_and_all_reduce_with_staging`): 核心变更所在, 新增 `_apply_metadata_gate` 函数并

修改 `poll_and_all_reduce` 和 `poll_and_all_reduce_with_staging`, 在 all-reduce 前加入元数据就绪检查。

- `python/sglang/srt/disaggregation/decode.py` (模块 解码器; 类别 `source`; 类型 `core-logic`; 符号 `_commit_transfer_to_req`, `_poll_with_metadata_gate`): 移除本地 `_is_fake_transfer` 定义, 改为从 `utils` 导入; 调整 `_commit_transfer_to_req` 逻辑: 将 `metadata` 未就绪时的重试改为直接 `abort`, 消除队列发散残余可能性。

关键符号: `_apply_metadata_gate`, `_is_fake_transfer`, `poll_and_all_reduce`, `poll_and_all_reduce_with_staging`, `_commit_transfer_to_req`, `_poll_with_metadata_gate`

## 关键源码片段

### `python/sglang/srt/disaggregation/utils.py`

核心变更所在, 新增 `_apply_metadata_gate` 函数并修改 `poll_and_all_reduce` 和 `poll_and_all_reduce_with_staging`, 在 all-reduce 前加入元数据就绪检查。

```
# python/sglang/srt/disaggregation/utils.py
# 新增的元数据门控函数, 在 all-reduce 前检查 bootstrap_room
def _apply_metadata_gate(polls, decode_reqs, metadata_buffers, server_args) -> None:
    """Downgrade Success → Transferring for requests whose metadata hasn't landed.

    Mutates `polls` in-place. Called before all-reduce so that MIN across TP
    ranks naturally prevents any rank from committing before all ranks are ready.
    """
    for i, poll_val in enumerate(polls):
        if poll_val == int(KVPoll.Success):
            decode_req = decode_reqs[i]
            # 跳过模拟传输 (fake backend), 不施加门控
            if _is_fake_transfer(decode_req.req, server_args):
                continue
            # 读取当前 rank 的 metadata 区域: bootstrap_room 字段
            actual_room = metadata_buffers.bootstrap_room[
                decode_req.metadata_buffer_index, 0
            ].item()
            if actual_room == 0:
                # 元数据尚未写入, 将状态降级为 Transferring
                polls[i] = int(KVPoll.Transferring)

# 改造后的 poll_and_all_reduce, 加入可选 gate 参数
def poll_and_all_reduce(
    pollers,
    gloo_group: dist.ProcessGroup,
    decode_reqs=None,
    metadata_buffers: Optional[MetadataBuffers] = None,
    server_args: Optional[ServerArgs] = None,
):
    # ... 原有的 poll 逻辑 ...
    polls = [int(poller.poll()) for poller in pollers]
```

```

# 仅当三个可选参数同时提供时才应用 gate
if (
    decode_reqs is not None
    and metadata_buffers is not None
    and server_args is not None
):
    _apply_metadata_gate(polls, decode_reqs, metadata_buffers, server_args)

# 所有 rank 的 poll 结果做 ReduceOp.MIN
tensor_to_reduce = torch.tensor(polls, dtype=torch.uint8, device="cpu")
dist.all_reduce(tensor_to_reduce, op=dist.ReduceOp.MIN, group=gloo_group)
return tensor_to_reduce.tolist()

```

### python/sglang/srt/disaggregation/decode.py

移除本地 `_is_fake_transfer` 定义，改为从 `utils` 导入；调整 `_commit_transfer_to_req` 逻辑：将 `metadata` 未就绪时的重试改为直接 `abort`，消除队列发散残余可能性。

```

# python/sglang/srt/disaggregation/decode.py
# 调整后的 _commit_transfer_to_req, 不再返回 bool, 直接 commit 或 abort
def _commit_transfer_to_req(self, decode_req: DecodeRequest):
    idx = decode_req.metadata_buffer_index
    (
        output_id,
        ...
    ) = ...

    expected_room = decode_req.req.bootstrap_room
    actual_room = ...

    if _is_fake_transfer(decode_req.req, self.scheduler.server_args):
        pass
    elif actual_room == 0:
        # 此前 _poll_with_metadata_gate 已确保所有 rank 的 metadata 就绪
        # 若此处仍为 0, 说明存在严重不一致, 直接 abort
        logger.error(
            f"Metadata unexpectedly not ready after readiness gate: "
            f"request {decode_req.req.rid}, bootstrap_room={expected_room}, "
            f"metadata_buffer_index={idx}"
        )
        prepare_abort(
            decode_req.req,
            "Metadata unexpectedly not ready after readiness gate "
            "(bootstrap_room=0)",
            status_code=HTTPStatus.INTERNAL_SERVER_ERROR,
        )
        decode_req.kv_receiver.clear()
        decode_req.kv_receiver = None
        return
    elif actual_room != expected_room:

```

```
# 真实损坏检测
error_msg = (...)
prepare_abort(decode_req.req, error_msg, status_code=...)
decode_req.kv_receiver.clear()
decode_req.kv_receiver = None
return

# Success - commit the transfer
decode_req.req.output_ids.append(output_id[0].item())
...
```

## 评论区精华

无 review 讨论。作者直接提交并通过 CI 验证，所有解聚测试均通过。

- 暂无高价值评论线程

## 风险与影响

- 风险：

1. gate 逻辑绕过风险：若 `_apply_metadata_gate` 未正确调用（如参数传错或遗漏），仍可能出现原发散问题，但函数签名强制要求三个可选参数同时非空才会生效，减少误用。
2. 性能回归：新增的 `bootstrap_room` 读取与比较增加了 CPU 开销，但仅为  $O(N)$  的 Python 循环，且仅在 `poll` 路径执行，影响极小。
3. staging 路径一致性：`poll_and_all_reduce_with_staging` 同步改造后，若 `staging handler` 行为与 `gate` 逻辑冲突，可能导致状态不一致，但两者均降级为 `Transferring`，方向一致。
4. 缺少测试覆盖：无新增单元测试验证 `gate` 降级逻辑及跨 rank 一致性场景，依赖现有集成测试。
  - 影响：影响 PD 模式下所有使用 `poll_and_all_reduce` 及 `poll_and_all_reduce_with_staging` 的解码节点，修复了因元数据延迟导致的服务挂起或结果损坏问题，提升稳定性。不涉及 API 或模型输出变更，用户无感知。团队应关注 `gate` 逻辑的正确维护。
  - 风险标记：核心路径变更，缺少测试覆盖，`gate` 参数误用风险

## 关联脉络

- PR #26299 [PD] Fix top logprobs crash in prefill path: 同为 PD 模块的 bugfix，涉及 `disaggregation` 核心逻辑，共享相似测试集。
- PR #25834 [PD] Decode pre-allocation refactor: 该 PR 对 `decode` 队列预处理逻辑进行了重构，与本 PR 的队列管理间接相关。