

# PR #24768 完整报告

sgl-project/sglang

[PrefillDelayer] support NCCL all-gather for cross-DP info sync

合并时间: 2026-05-10 12:20

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/24768>

## 执行摘要

- 一句话: PrefillDelayer 支持 NCCL all-gather 避免 GPU↔CPU 同步
- 推荐动作: 该 PR 值得查看, 因为它修复了一个潜在的性能问题, 并且设计清晰。特别关注 PrefillDelayer 中条件选择 gather 组和设备的方式, 以及调度器中简化的传递逻辑, 可作为模块间依赖注入的范例。

## 功能与动机

当设置 `SGLANG_NCCL_ALL_GATHER_IN_OVERLAP_SCHEDULER_SYNC_BATCH=1` 或 `disable_overlap_schedule=True` 时, 调度器其他部分 (如 `scheduler_dp_attn_mixin`) 已使用 NCCL 设备组, 但 PrefillDelayer 仍使用 gloo CPU 组, 这迫使每次调度迭代进行一次额外的 GPU↔CPU 同步。此外, 原代码在 `disable_overlap_schedule=True` 时, `_global_info_buffer` 创建在 GPU 设备上但 `local_info` 在 CPU 上并由 gloo 组 gather, 存在不协调。本 PR 旨在消除该同步开销并修复潜在的不一致。

## 实现拆解

1. PrefillDelayer 构造函数增加 `device_group` 参数 ( `python/sglang/srt/managers/prefill_delayer.py` ) : 新增可选参数 `device_group=None`, 并在初始化逻辑中根据条件决定使用 NCCL 设备组还是 CPU 组。
2. 条件选择 gather 组和设备: 根据 `server_args.disable_overlap_schedule` 或环境变量 `SGLANG_NCCL_ALL_GATHER_IN_OVERLAP_SCHEDULER_SYNC_BATCH` 判断, 若满足任一条件则使用 `device_group` 和 `device`, 否则使用 `cpu_group` 和 `"cpu"`。同时将 `_global_info_buffer` 的创建设备改为 `self._gather_device`, 确保与 gather 组一致。
3. 更新 `_gather_info` 方法: 将 `local_info` 的创建设备从固定 `"cpu"` 改为 `self._gather_device`, 且 all-gather 使用的 `group` 从 `self._cpu_group` 改为 `self._gather_group`, 确保实际使用的组与设备匹配。
4. 调度器中传递 `device_group` (`python/sglang/srt/managers/scheduler.py`) : 在 `init_schedule_policy` 方法中创建 PrefillDelayer 时, 新增传入 `device_group=self.tp_group.device_group`, 并将 `device` 参数简化为无条件传递 `self.tp_group.device`, 因为选择逻辑已移至 delayer 内部。同时移除了之前按条件选择 `device` 的三元表达式。

5. 导入调整：在 `prefill_delayer.py` 中新增 `from sglang.srt.environ import envs` 以读取环境变量。

关键文件：

- `python/sglang/srt/managers/prefill_delayer.py` (模块 调度器；类别 `source`；类型 `dependency-wiring`；符号 `PrefillDelayer.init`, `PrefillDelayer._gather_info`)：核心修改文件，新增 `device_group` 参数并实现了条件选择 `gather` 组和设备，修复了 `all-gather` 路径的同步问题。
- `python/sglang/srt/managers/scheduler.py` (模块 调度器；类别 `source`；类型 `core-logic`；符号 `init_schedule_policy`)：作为调用方，传递 `device_group` 并简化 `device` 参数的传递，是本 PR 的另一关键改动。

关键符号：`PrefillDelayer.init`, `PrefillDelayer._gather_info`, `init_schedule_policy`

## 关键源码片段

### `python/sglang/srt/managers/prefill_delayer.py`

核心修改文件，新增 `device_group` 参数并实现了条件选择 `gather` 组和设备，修复了 `all-gather` 路径的同步问题。

```
# python/sglang/srt/managers/prefill_delayer.py
```

```
class PrefillDelayer:
    def __init__(
        self,
        dp_size: int,
        attn_tp_size: int,
        cpu_group,
        server_args,
        max_delay_passes: int,
        token_usage_low_watermark: Optional[float],
        metrics_collector: Optional["SchedulerMetricsCollector"] = None,
        device: Optional["torch.device"] = "cpu",
        device_group=None, # <- 新增参数，用于 NCCL all-gather
    ):
        # ... 其他初始化 ...
        # 判别是否使用 NCCL 设备组：
        # 当 disable_overlap_schedule = True 或环境变量 SGLANG_NCCL_ALL_GATHER_IN_OVERLAP_SCHEDULER_SYNC_BATCH 被设置时使用
        use_nccl = (
            server_args.disable_overlap_schedule
            or envs.SGLANG_NCCL_ALL_GATHER_IN_OVERLAP_SCHEDULER_SYNC_BATCH.get()
        )
        if use_nccl:
            assert device_group is not None, (
                "device_group is required when using NCCL for PrefillDelayer all-gather"
            )
            self._gather_group = device_group # 使用 NCCL 组
```

```

        self._gather_device = device # 使用 GPU 设备 (tp_group.device)
    else:
        self._gather_group = cpu_group # 保持原有 gloo CPU 组
        self._gather_device = "cpu" # 设备为 CPU

# 全局 buffer 现在创建在选定的 gather 设备上, 保证一致性
self._global_info_buffer = torch.empty(
    (dp_size_dim, attn_tp_size, 5),
    dtype=torch.int64,
    device=self._gather_device,
)
# 不再需要 self._cpu_group 字段

def _gather_info(self, ...):
    local_info = torch.tensor(
        [...],
        device=self._gather_device, # 改为动态设备, 而非固定 "cpu"
        dtype=torch.int64,
    )
    torch.distributed.all_gather_into_tensor(
        self._global_info_buffer.flatten(),
        local_info,
        group=self._gather_group, # 使用选定的 group
    )

```

### python/sglang/srt/managers/scheduler.py

作为调用方, 传递 `device_group` 并简化 `device` 参数的传递, 是本 PR 的另一关键改动。

```

# python/sglang/srt/managers/scheduler.py

def init_schedule_policy(self):
    # ...
    self.prefill_delayer = PrefillDelayer(
        dp_size=self.dp_size,
        attn_tp_size=self.attn_tp_size,
        cpu_group=self.tp_cpu_group,
        device_group=self.tp_group.device_group, # <-- 新增参数, 传入 NCCL 组
        server_args=self.server_args,
        metrics_collector=...,
        max_delay_passes=...,
        token_usage_low_watermark=...,
        device=self.tp_group.device, # 简化: 无条件传入设备
    )

```

## 评论区精华

无 reviewer 评论。该 PR 由 ispobock 直接批准。

- 暂无高价值评论线程

## 风险与影响

- 风险：本 PR 改动量小 (+25/-9)，且条件分支清晰，风险较低。主要风险在于：若用户设置了 `SGLANG_NCCL_ALL_GATHER_IN_OVERLAP_SCHEDULER_SYNC_BATCH` 但未提供 `device_group`（例如调度器未正确初始化），断言会触发报错，但这属于合理的防御性编程。此外，默认行为（`disable_overlap_schedule=False` 且未设置环境变量）保持不变，不会影响现有用户。
- 影响：影响范围限定在启用了 `PrefillDelayer` 且满足 NCCL all-gather 条件的场景。对于这些场景，可消除每次调度迭代中额外的 GPU↔CPU 同步，提升性能。对于其他用户无行为变化。对系统其他模块无负面影响。
- 风险标记：核心路径变更，缺少测试覆盖

## 关联脉络

- PR #24735 [Spec] Move `accept_tokens` off `EagleDraftInput`; pass via method arg: 同属调度器模块的依赖重构，关注模块间参数传递的简化与一致性
- PR #24097 Restrict `fa_skip_kv_cache` to non-MLA backends: 同为 attention 相关性能修复，关注条件分支对性能的影响