

PR #19548 完整报告

sgl-project/sglang

fix: support PP2+CP8+TP8 (PP with context parallelism)

合并时间: 2026-03-17 00:51

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/19548>

执行摘要

本次 PR 修复了 SGLang 调度器中 Pipeline Parallelism (PP) 与 Context Parallelism (CP) 结合使用的通信问题, 支持 PP2+CP8+TP8 等配置, 通过修改调度器发送 / 接收逻辑并放宽配置断言, 确保在 H20 等硬件上并行部署正常工作。

功能与动机

动机源于 H20 GPU 集群上 PP 与 CP 结合使用时出现的问题, 导致无法生成输出。PR body 引用 PR #19504 的评论, 明确需修复此配置下的 bug。Issue 评论中 [yiakwy-xpu-ml-framework-team](#) 验证修改后 PP2+CP 工作正常, 强调了修复的紧迫性。

实现拆解

主要改动点按模块拆解:

- 调度器模块 (`scheduler_pp_mixin.py`) :
 - 修改 `_pp_send_pyobj_to_next_stage` 和 `_pp_recv_pyobj_from_prev_stage` 函数, 添加条件 `self.attn_tp_rank == 0 and self.attn_cp_rank == 0`, 确保只有 TP 和 CP rank 为零的进程进行进程间通信。
 - 在接收后添加 CP 广播代码块: `python if self.attn_cp_size > 1: data = broadcast_pyobj(data, self.attn_cp_group.rank, self.attn_cp_cpu_group, src=self.attn_cp_group.ranks[0],)`
- 配置模块 (`server_args.py`) :
 - 在 `_handle_context_parallelism` 方法中, 将硬性断言 `assert self.pp_size == 1` 改为条件性检查, 仅当 `enable_nsa_prefill_context_parallel` 为假时触发, 允许 PP 与 CP 在 NSA 预填充上下文并行启用时共存。
 - 设置 `attn_cp_size = tp_size` 以适配配置路径。

评论区精华

Review 讨论中关键交锋包括:

- ShangmingCai 建议移除断言, 但 whybeyoung 反驳:

"The assert blocks the case: context parallelism on (`attn_cp_size > 1`) with pipeline parallelism on (`pp_size > 1`) while NSA prefill context parallel is off (`enable_nsa_prefill_context_parallel=False`). That combination is not

implemented/tested; removing the check would allow an invalid config to start." 最终采纳条件性断言，平衡灵活性与安全性。

- [yiakwy-xpu-ml-framework-team](#) 提出混淆点：关于 `attn_cp_size` 定义和 TP 约束，指出 TP=2 时也应支持 CP，但此问题被标记为未解决，需后续澄清。

风险与影响

风险：

1. 核心调度通信路径变更可能引入死锁或数据不一致，特别是在分布式环境中依赖特定 rank 假设。
2. 配置检查放宽可能导致用户误用未测试组合，引发运行时错误。
3. 缺少针对 PP+CP 组合的专门测试，依赖 CI 覆盖可能不足。

影响：

- 积极影响：扩展并行配置选项，支持更高效的大规模模型部署，提升 H20 等硬件资源利用率。
- 潜在影响：需用户确保正确设置 `enable_nsa_prefill_context_parallel` 等变量，团队需跟进未解决的 `attn_cp_size` 讨论以维护代码质量。

关联脉络

从历史 PR 看，本次 PR 与近期调度器优化（如 PR #22577 修复空闲检测、PR #22453 修复 HiSparse 解码侧）共享调度器模块的演进趋势，反映团队在分布式并行性和性能调优上的持续投入。虽然未直接关联到特定功能线，但共同支撑 SGLang 在高负载环境下的稳定性和扩展性。