

PR #21907 完整报告

sgl-project/sglang

[Fix] Add `_MOE_TP` to `graph_capture` for MoE models with `ep>1`

合并时间: 2026-04-03 17:33

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/21907>

执行摘要

本 PR 修复了 MoE 模型在专家并行 (ep) 大于 1 时, 由于 CUDA 图捕获中缺失 `_MOE_TP` 通信组而导致的段错误 (退出代码 -9)。通过重构 `graph_capture()` 函数确保所有相关通信组被正确捕获, 并添加了上下文并行 (CP) 的 CUDA 图禁用条件, 恢复了受影响配置 (如 `Qwen3-235B-FP8 --tp=8 --ep=2`) 的正常运行。这是一个关键 bugfix, 提升了分布式 MoE 场景下 CUDA 图的稳定性。

功能与动机

问题根源: PR #18233 将 MoE allreduce 从通用的 `_TP` 组切换到专用的 `_MOE_TP` 组, 但未同步更新 `graph_capture()` 函数。这导致在 CUDA 图回放时, 自定义 allreduce 内核 (`cross_device_reduce_1stage`) 尝试访问未注册的 IPC 句柄, 引发非法内存访问和进程崩溃。

影响范围: 所有 MoE 模型配置满足 $1 < ep < tp$ 且 `moe_tp_size > 1` 时受影响, 例如 `Qwen3-235B-FP8 --tp=8 --ep=2`。CI 证据显示在 PR #18233 合并后出现段错误, 而之前正常。

验证结果: 在 8xH200 上测试, 修复后 `Qwen3-235B-FP8` 配置通过测试 (gsm8k 96.4%, 3124 tok/s)。

实现拆解

主要改动集中在两个文件:

1. `python/sglang/srt/distributed/parallel_state.py`:

- 重构 `graph_capture()` 函数, 将原有的条件分支替换为通用逻辑。
- 使用 `contextlib.ExitStack` 和 `seen` 集合管理多个通信组的捕获, 确保 `_TP`、`_MOE_EP` 和 `_MOE_TP` 组 (如果存在且唯一) 都被纳入 CUDA 图。
- 关键代码片段:

2. `python/sglang/srt/server_args.py`:

- 在 `_handle_pieewise_cuda_graph()` 方法中添加条件, 当 `attn_cp_size > 1` 时禁用分段 CUDA 图。
- 这反映了上下文并行 (CP) 目前不支持 CUDA 图捕获的现状。

评论区精华

由于本 PR 没有正式的 review 评论，讨论亮点主要从提交历史中推断：

- Fridge003 的提交：添加了 'disable cp for pcg' 提交，表明在合并前对 CP 的 CUDA 图支持进行了调整，可能基于内部讨论或测试发现 CP 与 CUDA 图不兼容。
- 决策结论：采用保守策略，在 CP 启用时禁用 CUDA 图，避免潜在问题。

风险与影响

技术风险：

- `graph_capture()` 重构依赖 `id(group)` 唯一性，如果通信组对象被意外复制或重用，可能导致捕获遗漏。
- CP 禁用条件可能影响依赖 CUDA 图优化的 CP 配置性能，但这是权衡后的安全选择。
- 修改涉及分布式通信核心路径，需在多样 MoE 配置下充分测试（如不同 `tp`、`ep`、`moe_tp_size` 组合）。

影响分析：

- 用户影响：修复了受影响的 MoE 模型用户的段错误问题，恢复模型功能；对不使用 MoE 或 `ep=1` 的用户无影响。
- 系统影响：提升了 CUDA 图在复杂分布式配置下的稳定性和可靠性。
- 团队影响：强调了通信组变更时需同步更新所有相关上下文（如 CUDA 图捕获）的协作规范。

关联脉络

- 直接关联：PR #18233 是本问题的根源，它引入了 `_MOE_TP` 组但未更新 `graph_capture()`，导致本 PR 的修复需求。
- 演进趋势：从近期历史 PR 看，仓库持续优化分布式性能（如 #21511 启用 FP8 KV 缓存、#21591 添加 HiSparse 缓存传输），本 PR 是这一趋势中的稳定性修复，确保 MoE 等高级特性在 CUDA 图优化下可靠工作。
- 跨 PR 模式：反映了基础设施变更（如新通信组）需全面测试的教训，未来类似改动应更注重影响面分析。