

PR #24332 完整报告

sgl-project/sglang

[Codex] Diffusion handle non-contiguous CFG communication

合并时间: 2026-05-06 17:27

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/24332>

执行摘要

- 一句话: 修复扩散模型 CFG 并行中非连续张量通信崩溃
- 推荐动作: 建议合并此 PR, 因为它修复了 CPG 并行在 JoyAI 等模型上的功能性崩溃, 并带来了显著的性能提升。但在合并前, 应评估 review 中提出的 in-place 语义问题——如果调用者依赖原始张量更新, 需复制回结果 (如 `input_.copy_(contiguous_result)`); 若当前所有调用者都不依赖, 则可忽略。建议补充一个单元测试用例, 覆盖非连续张量输入场景。

功能与动机

修复 `jdopensee/JoyAI-Image-Edit-Diffusers` 模型在 `--num-gpus=2` 自动启用 CFG 并行时, 去噪阶段因非连续张量调用 `torch.distributed.all_reduce` 和 `broadcast` 而崩溃的问题。PR body 中描述了错误信息 `ValueError: Tensors must be contiguous` 和具体触发点。

实现拆解

1. 在 `cfg_model_parallel_all_reduce` 中添加连续化检查: 在 `python/sglang/multimodal_gen/runtime/distributed/communication_op.py` 的 `cfg_model_parallel_all_reduce` 函数中, 在调用底层 `all_reduce` 前检查 `input_.is_contiguous()`, 若非连续则通过 `input_.contiguous()` 创建连续副本, 确保底层分布式调用成功。
2. 在 `broadcast` 中添加连续化检查: 在 `python/sglang/multimodal_gen/runtime/distributed/group_coordinator.py` 的 `broadcast` 方法中, 同样在调用 `torch.distributed.broadcast` 前对非连续张量进行 `contiguous()` 转换。
3. 保持快速路径不变: 仅当张量非连续时才执行额外复制, 不影响已有连续张量的性能。

关键文件:

- `python/sglang/multimodal_gen/runtime/distributed/communication_op.py` (模块 分布式通信; 类别 `source`; 类型 `core-logic`; 符号 `cfg_model_parallel_all_reduce`): 修改 `cfg_model_parallel_all_reduce`, 在调用底层 `all_reduce` 前对非连续张量执行 `contiguous()`, 修复 CFG 并行中的崩溃。review 评论对此处的 in-place 语义提出风险。
- `python/sglang/multimodal_gen/runtime/distributed/group_coordinator.py` (模块 分布式通信; 类别 `source`; 类型 `core-logic`; 符号 `GroupCoordinator.broadcast`): 修改 `broadcast` 方法, 添加与 `all_reduce` 相同的 `contiguous()` 保护, 修复 JoyAI 模型在 `broadcast` 上的崩溃。

关键符号: `cfg_model_parallel_all_reduce`, `GroupCoordinator.broadcast`

关键源码片段

[python/sglang/multimodal_gen/runtime/distributed/communication_op.py](#)

修改 `cfg_model_parallel_all_reduce`, 在调用底层 `all_reduce` 前对非连续张量执行 `contiguous()`, 修复 CFG 并行中的崩溃。review 评论对此处的 `in-place` 语义提出风险。

```
# python/sglang/multimodal_gen/runtime/distributed/communication_op.py

def cfg_model_parallel_all_reduce(
    input_: torch.Tensor,
    op: torch._C._distributed_c10d.ReduceOp = torch._C._distributed_c10d.ReduceOp.SUM,
) -> torch.Tensor:
    """All-reduce the input tensor across CFG parallel group."""
    # 修复: 底层 all_reduce 要求连续张量, 非连续时显式转为连续
    # 注意: contiguous() 可能创建副本, 破坏 in-place 语义
    if not input_.is_contiguous():
        input_ = input_.contiguous()
    return get_cfg_group().all_reduce(input_, op=op)
```

[python/sglang/multimodal_gen/runtime/distributed/group_coordinator.py](#)

修改 `broadcast` 方法, 添加与 `all_reduce` 相同的 `contiguous()` 保护, 修复 JoyAI 模型在 `broadcast` 上的崩溃。

```
# python/sglang/multimodal_gen/runtime/distributed/group_coordinator.py

class GroupCoordinator:
    # ...
    def broadcast(self, input_: torch.Tensor, src: int = 0, async_op: bool = False):
        """Broadcast the input tensor.
        NOTE: `src` is the local rank of the source rank.
        """
        assert src < self.world_size, f"Invalid src rank ({src})"

        # Bypass the function if we are using only 1 GPU.
        if self.world_size == 1:
            return input_
        # Broadcast.
        # 修复: 底层 broadcast 要求连续张量, 非连续时显式转为连续
        if not input_.is_contiguous():
            input_ = input_.contiguous()
        torch.distributed.broadcast(
            input_,
            src=self.ranks[src],
            group=self.device_group,
            async_op=async_op,
        )
        return input_
```

评论区精华

review 评论来自 `gemini-code-assist[bot]`，指出在 `cfg_model_parallel_all_reduce` 中对非连续输入调用 `contiguous()` 会创建副本，破坏了对原始张量原地修改的语义（`all_reduce` 通常期望 in-place），可能导致调用者继续使用未规约的原始张量。建议将结果复制回原始张量以保持行为一致。该评论尚未被作者回复或解决。

- `all_reduce` 非连续张量 contiguity 破坏 in-place 语义 (correctness): 未解决。PR 作者未回复或修改代码。

风险与影响

• 风险:

1. 语义破坏风险: `all_reduce` 文档和习惯用法支持原地修改，但当前实现仅在非连续时返回新张量，调用者若依赖原始张量更新将得到未规约结果。这可能影响依赖 in-place 行为的代码路径（如后续张量共享引用）。
2. 性能影响: `contiguous()` 调用会触发内存复制和重新排列，对非连续大张量有额外开销。但当非连续时触发，且 CFG 并行本身较少见非连续场景，风险可控。
3. 测试覆盖不足: 源码变更未附带新测试用例，仅依赖现有单元测试和手工验证。
 - 影响: 用户影响: 修复了 JoyAI 等扩散模型在 CFG 并行模式下的崩溃，使这些模型能正常使用自动 CFG 并行加速。根据 PR 数据，对 JoyAI 可带来 24-27% 的端到端延迟改善。
 - 系统影响: 改动集中在两个通信函数，影响范围仅限于扩散模型的 CFG 并行路径，不涉及序列并行或其他通信原语。
 - 团队影响: 小型改动（4 行增加），易于审查和合并。

- 风险标记: 语义破坏风险，缺少测试覆盖

关联脉络

- 暂无明显关联 PR