

PR #22003 完整报告

sgl-project/sglang

Support moe_dp_size = 1 for various attention_cp_size

合并时间: 2026-04-21 02:58

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/22003>

执行摘要

- 一句话: 支持 MoE 数据并行大小与注意力上下文并行大小解耦, 提升配置灵活性。
- 推荐动作: 该 PR 值得精读, 特别是 `communicator.py` 中散射模式的扩展和 `dp_attention.py` 中新通信函数的设计。关注 MOE_FULL 模式如何平衡数据完整性和通信开销。

功能与动机

从 PR body 中引用: 'Previously, we can only support `attention_cp_size == moe_dp_size` which is too restricted. In the real world case, we should let the MoE part unchanged and only apply the context parallel into attention layer.'

实现拆解

1. 扩展散射模式: 在 `python/sglang/srt/layers/communicator.py` 的 `ScatterMode` 枚举中添加 `MOE_FULL`, 用于当 `moe_dp_size < attn_cp_size` 时, 在 MoE 组内保持数据完整。
2. 引入通信辅助函数: 在 `python/sglang/srt/layers/dp_attention.py` 中添加 `get_moe_cp_group`、`is_enable_moe_cp_allgather` 等函数, 提供 MoE 上下文并行的组管理和条件判断。
3. 调整层通信逻辑: 在 `communicator.py` 中修改 `_compute_mlp_mode` 方法, 支持 `MOE_FULL` 模式, 并新增 `_gather_hidden_states_and_residual_moe` 和 `_scatter_hidden_states_moe` 方法来处理数据聚集和散射。
4. 更新并行组初始化: 在 `python/sglang/srt/distributed/parallel_state.py` 中, 当 `attn_cp_size > moe_dp_size` 时, 将 `_MOE_DP` 组设置为 `_ATTN_CP` 组的别名, 简化通信。
5. 适配模型代码: 修改 `qwen2_moe.py` 和 `qwen3_moe.py`, 更新导入和断言, 支持 `attn_cp_size` 整除 `moe_dp_size` 而非强制相等。
6. 补充测试覆盖: 在 `test/registered/4-gpu-models/test_qwen3_30b.py` 中添加新测试类 `TestQwen330BCP`, 验证配置 `--tp-size 4 --moe-dp-size 1 --attn-cp-size 2` 的正确性。

关键文件:

- `python/sglang/srt/layers/communicator.py` (模块 通信层; 类别 `source`; 类型 `core-logic`; 符号 `_gather_hidden_states_and_residual_moe`, `_scatter_hidden_states_moe`): 核心通信层逻辑, 新增 `MOE_FULL` 散射模式和 `gather/scatter` 方法, 影响数据流分发。

- `python/sglang/srt/layers/dp_attention.py` (模块 并行通信; 类别 source; 类型 core-logic; 符号 `get_moe_cp_group`, `get_moe_cp_rank`, `get_moe_cp_size`, `is_enable_moe_cp_allgather`) : 提供 MoE 上下文并行的组管理函数, 是通信调用的入口点。
- `python/sglang/srt/distributed/parallel_state.py` (模块 并行状态; 类别 source; 类型 core-logic) : 关键并行组初始化逻辑, 修改 MOE_DP 组以支持新配置, 影响整体通信拓扑。
- `test/registered/4-gpu-models/test_qwen3_30b.py` (模块 测试套件; 类别 test; 类型 test-coverage; 符号 `TestQwen330BCP`, `setUpClass`, `tearDownClass`, `test_gsm8k`) : 新增测试类验证新配置的正确性, 确保功能稳定。
- `python/sglang/srt/models/qwen2_moe.py` (模块 模型层; 类别 source; 类型 data-contract) : 模型层适配, 更新数据流转逻辑以支持新并行配置。

关键符号: `_gather_hidden_states_and_residual_moe`, `_scatter_hidden_states_moe`, `get_moe_cp_group`, `get_moe_cp_rank`, `get_moe_cp_size`, `is_enable_moe_cp_allgather`, `moe_cp_all_gather_into_tensor`

关键源码片段

`python/sglang/srt/layers/communicator.py`

核心通信层逻辑, 新增 MOE_FULL 散射模式和 `gather/scatter` 方法, 影响数据流分发。

```
@classmethod
def _compute_mlp_mode(cls, context: _LayerModeComputationContext):
    if context.is_layer_sparse:
        if (
            not get_moe_a2a_backend().is_none()
            or should_use_flashinfer_cutlass_moe_fp4_allgather()
        ):
            return ScatterMode.SCATTERED
        # NSA CP 目前不支持 MOE_FULL; 回退到 FULL
        if is_enable_moe_cp_allgather() and not is_nsa_enable_prefill_cp():
            return ScatterMode.MOE_FULL # 当启用 MoE CP allgather 时, 使用 MOE_FULL 模式
        return ScatterMode.FULL
    else:
        return (
            ScatterMode.SCATTERED
            if enable_moe_dense_fully_dp()
            else ScatterMode.FULL
        )
```

评论区精华

Review 中主要讨论了代码风格和设计细节。ch-wan 指出函数命名应遵循 `is_xxx_enabled` 规范, Shunkangz 已修正。关于 PCG (Piecewise CUDA Graph) 的注释被质疑为混淆, Shunkangz 解释已禁用相关路径。此外, 对 `_MOE_DP` 组初始化逻辑的讨论澄清了不同并行配置下组构成的差异。

- 函数命名规范 (style): 已采纳建议, 函数命名为 `is_enable_moe_cp_allgather`。
- PCG 相关性注释 (design): 注释被简化或移除, 确保清晰性。
- MOE_DP 组初始化逻辑 (correctness): 维持原逻辑, 确保并行组正确初始化。

风险与影响

- 风险: 技术风险包括:
 1. 形状不匹配: 在 `communicator.py` 的 `should_fuse_mlp_allreduce_with_next_layer` 中, 当 `MOE_FULL` 激活时, 需禁用融合以避免残留张量形状错误。
 2. 并行死锁: `parallel_state.py` 中 `_MOE_DP` 组可能错误初始化, 导致通信失败。
 3. 回归错误: 新散射模式可能影响现有配置, 需确保向后兼容。 - 影响: 用户影响: 工程师现在可以配置 `moe_dp_size=1` 并调整 `attn_cp_size` 以优化注意力层并行, 提升吞吐量 (PR 中显示从 941.495 token/s 到 1953.284 token/s)。系统影响: 扩展了并行策略的灵活性, 可能影响调度器和缓存管理。团队影响: 需更新相关文档和测试, 确保新配置稳定性。 - 风险标记: 形状不匹配风险, 并行组初始化风险

关联脉络

- PR #22914 [Refactor] Deduplicate NSA utils.py into cp_utils.py for context parallel: 同样涉及上下文并行逻辑重构, 共享并行策略调整主题。
- PR #23202 [core] Always-on StreamingSession in UnifiedRadixCache: 涉及调度和缓存优化, 与本 PR 的并行配置扩展相辅相成。