

PR #24934 完整报告

sgl-project/sglang

DeepSeek V4 MTP Support CP

合并时间: 2026-05-20 07:51

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/24934>

执行摘要

- 一句话: DeepSeek V4 MTP 支持 Context Parallel
- 推荐动作: 值得精读, 尤其是在 CP 与 speculative decoding 集成方面的实现模式。关注点: CP 元数据的准备时机、数据切分后如何恢复顺序, 以及 CP 与 dp_attention 的兼容性。设计权衡: 复用 NSA 已有的 CP 工具函数, 避免重复逻辑, 但引入了对 NSA 后端的隐式依赖。

功能与动机

在 DeepSeek V4 上同时使用 MTP (多 token 预测) 和 Context Parallel 可显著降低长序列 prefill 延迟, 但对 MTP 的 forward 路径缺乏 CP 适配。本 PR 填补这一缺口, 使 EAGLE 推测解码能在 CP 模式下正确运行。

实现拆解

1. 导入 CP 工具函数: 在 `deepseek_v4_nextn.py` 中添加来自 `sglang.srt.layers.attention.nsa.utils` 和 `sglang.srt.layers.utils.cp_utils` 的 CP 相关函数, 包括 `nsa_use_prefill_cp`、`cp_split_and_rebuild_data`、`cp_split_and_rebuild_position`、`cp_all_gather_rerange_output`、`prepare_context_parallel_metadata` 等。这些函数构成 CP 拆解 / 聚合的基础操作。
2. DeepseekV4ModelNextN 适配 CP: 在 `__init__` 中新增 `nsa_enable_prefill_cp` 标志和 `cp_size` 属性; 在 `forward` 中, 在调用 `self.decoder` 之前插入条件判断: 如果 `nsa_use_prefill_cp(forward_batch)` 为真, 则调用 `cp_split_and_rebuild_data` 和 `cp_split_and_rebuild_position` 将当前设备的数据按 CP 策略切分; decoder 输出后调用 `cp_all_gather_rerange_output` 将各分片结果聚集并恢复顺序, 保证后续 `hc_head` 处理正确。
3. DeepseekV4ForCausalLM 适配 CP: 在 `__init__` 中新增 `cp_rank` 和 `cp_size` 属性; 在 `forward` 中增加 CP 元数据准备逻辑: 通过 `can_nsa_cp_split` 判断是否需要切分, 若是则调用 `prepare_context_parallel_metadata` 生成 `attn_cp_metadata` 附加到 `forward_batch`, 并根据 round-robin-split 模式调整相关索引偏移。
4. 配套测试: 在 `test/registered/dsv4/test_deepseek_v4_flash_fp4_b200.py` 中新增 `TestDSV4FlashFP4B200Balanced_CP` 测试类, 启动服务时加入 `--attn-cp-size 4`、`--enable-nsa-prefill-context-parallel` 和 `--nsa-prefill-cp-mode round-robin-split`, 通过 GSM8K 检查验证准确率, 确保 CP+MTP 组合的正确性。

关键文件:

- `python/sglang/srt/models/deepseek_v4_nextn.py` (模块 MTP 解码; 类别 source; 类型 core-logic): 核心实现文件, 通过添加 NSA CP 的导入和条件分支, 使 MTP 模型支持 prefill 阶段的上下文并行; 修改了 `DeepseekV4ModelNextN` 和 `DeepseekV4ForCausalLM` 两个类的 `__init__` 和 `forward`, 涉及数据切分、位置重建和 `all-gather` 聚集等关键逻辑。
- `test/registered/dsv4/test_deepseek_v4_flash_fp4_b200.py` (模块 DSV4 测试; 类别 test; 类型 test-coverage; 符号 `TestDSV4FlashFP4B200Balanced_CP`, `setUpClass`, `tearDownClass`, `test_gsm8k`): 新增测试类 `TestDSV4FlashFP4B200Balanced_CP`, 验证在 `TP=4+DP=4+CP=4` 配置下 CP+MTP 的 GSM8K 准确率; 测试覆盖了新的命令行参数组合, 是 PR 功能正确性的关键验证。

关键符号: `DeepseekV4ModelNextN.init`, `DeepseekV4ModelNextN.forward`, `DeepseekV4ForCausalLM.init`, `DeepseekV4ForCausalLM.forward`, `TestDSV4FlashFP4B200Balanced_CP.setUpClass`, `TestDSV4FlashFP4B200Balanced_CP.tearDownClass`, `TestDSV4FlashFP4B200Balanced_CP.test_gsm8k`

关键源码片段

`python/sglang/srt/models/deepseek_v4_nextn.py`

核心实现文件, 通过添加 NSA CP 的导入和条件分支, 使 MTP 模型支持 prefill 阶段的上下文并行; 修改了 `DeepseekV4ModelNextN` 和 `DeepseekV4ForCausalLM` 两个类的 `__init__` 和 `forward`, 涉及数据切分、位置重建和 `all-gather` 聚集等关键逻辑。

```
def forward(self, input_ids, positions, forward_batch, input_embeds=None):
    # ... 省略 embedding 处理 ...

    # 如果启用 NSA prefill CP, 则在 decoder 前对 hidden states 和 positions 进行拆分
    if nsa_use_prefill_cp(forward_batch):
        hidden_states = cp_split_and_rebuild_data(forward_batch, hidden_states)
        positions = cp_split_and_rebuild_position(forward_batch, positions)

    # 实际 decoder 调用 (每个设备处理自己的部分)
    hidden_states = self.decoder(
        positions=positions,
        hidden_states=hidden_states,
        forward_batch=forward_batch,
        input_ids_global=input_ids_global,
    )

    # 如果启用了 CP, 则对所有设备的输出进行 all-gather 并恢复原始顺序
    if nsa_use_prefill_cp(forward_batch):
        hidden_states = cp_all_gather_rerange_output(
            hidden_states,
            self.cp_size,
            forward_batch,
```

```

        torch.cuda.current_stream(),
    )

    # 后续的 hc_head 处理 (不变)
    pre_hc_head = hidden_states.flatten(1)
    hidden_states = self.hc_head(...)

    return logits

```

test/registered/dsv4/test_deepseek_v4_flash_fp4_b200.py

新增测试类 TestDSV4FlashFP4B200Balanced_CP，验证在 TP=4+DP=4+CP=4 配置下 CP+MTP 的 GSM8K 准确率；测试覆盖了新的命令行参数组合，是 PR 功能正确性的关键验证。

```

class TestDSV4FlashFP4B200Balanced_CP(ServerSanityMixin, CustomTestCase):
    '''验证 CP + MTP 组合的集成测试。使用 TP=4, DP=4, CP=4, 并启用 NSA prefill CP。'''

    @classmethod
    def setUpClass(cls):
        cls.model = try_cached_model(MODEL)
        cls.base_url = DEFAULT_URL_FOR_TEST
        cls.process = popen_launch_server(
            cls.model, cls.base_url, timeout=SERVER_LAUNCH_TIMEOUT,
            other_args=[
                '--trust-remote-code',
                '--tp', '4',
                '--attn-cp-size', '4', # 设置 Context Parallel 度为 4
                '--enable-dp-attention',
                '--moe-a2a-backend', 'deepep',
                '--speculative-algorithm', 'EAGLE', # 启用 EAGLE 推测解码
                '--speculative-num-steps', '1',
                '--speculative-eagle-topk', '1',
                '--speculative-num-draft-tokens', '2',
                '--enable-nsa-prefill-context-parallel', # 启用 NSA 的 prefill CP
                '--nsa-prefill-cp-mode', 'round-robin-split', # 使用 round-robin split 模式
                '--deepep-config', DEEPEP_CONFIG,
            ],
            env=_DEEPEP_ENV,
        )

    @classmethod
    def tearDownClass(cls):
        if hasattr(cls, 'process') and cls.process:
            kill_process_tree(cls.process.pid)

    def test_gsm8k(self):
        _gsm8k_check(self) # 使用 GSM8K 数据集验证推理正确性

```

评论区精华

- Fridge003 指出 `get_attention_tp_rank()` 等函数已废弃，应改用 `get_attention_cp_size()` 和 `get_attention_cp_rank()`，已修复。
- 另一个废弃标志 `SGLANG_DEBUG_HACK_CP_CHECK_RANK_CONSISTENCY` 也被指出需要在 `main` 上移除，最终版本已无该标志。
- Fridge003 请求增加 CP+MTP 的集成测试，Paiiiiiiiiiiiiiiii 同意并实现，后续经过多轮 CI 调试通过。
- Deprecated flag `SGLANG_DEBUG_HACK_CP_CHECK_RANK_CONSISTENCY (correctness)`: 已移除，最终版本无此 flag。
- Wrong CP size/rank function (correctness): 已修正。
- Request for CP+MTP test (testing): 测试类 `TestDSV4FlashFP4B200Balanced_CP` 已添加并通过 CI。

风险与影响

- 风险:
 - 核心路径变更: MTP forward 增加了 CP 分支，可能在非 CP 配置下引入回归；缺少对无 CP 场景的针对性测试。
 - 依赖 NSA 的 CP 机制: 本实现紧密绑定于 NSA 的 prefill CP，若未来 NSA 后端重构，需同步维护。
 - 性能开销: CP 拆分和 all-gather 引入额外通信和重排，若 CP 配置不合理可能导致性能反降。
 - 测试覆盖有限: 仅覆盖 $TP=4+DP=4+CP=4$ 的 balanced 场景，未覆盖其他 CP 大小或非 round-robin 模式。
- 影响:
 - 用户: 启用 `--enable-nsa-prefill-context-parallel` 和 `--attn-cp-size` 后，DeepSeek V4 推理可使用 CP+MTP。若无此 PR，在 CP 模式下无法正确使用 speculative decoding。用户需同步设置 `--nsa-prefill-cp-mode`。
 - 系统: 增加了约 60 行核心模型逻辑，以及新的测试类。代码维护者需留意 CP 与 speculative 的兼容性变化。
 - 团队: 需确保后续对 NSA CP 的修改不破坏此功能，最好建立 CI 覆盖。
 - 风险标记: 核心路径变更，依赖 NSA CP，测试覆盖有限

关联脉络

- PR #25396 fix: fix deepseek v4 CP error: 修复了 DeepSeek V4 Context Parallel 的先前错误，为本 PR 提供基础。
- PR #25729 fix(dsv4): upgrade forward metadata on main stream for large PP size: 涉及 DeepSeek V4 forward 元数据处理，与 CP 模式下 forward 流程间接关联。
- PR #25465 verify_done: wait not synchronize: 优化 speculative decoding v2 的同步机制，本 PR 也涉及 speculative decoding 与 CP 的配合。