

PR #24856 完整报告

sgl-project/sglang

Fix TRTLLM MHA routing for draft extend

合并时间: 2026-05-13 06:48

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/24856>

执行摘要

- 一句话: 修复 draft extend 路由到 decode kernel 导致的非法内存访问
- 推荐动作: 该 PR 作为紧急 bugfix 值得精读, 尤其是涉及注意力后端的路由逻辑的开发者。建议在后续优化中评估是否可在特定条件下 (如单 batch 且无 IMA 风险) 对 draft extend 也使用 decode kernel 以恢复性能。

功能与动机

在 Qwen3.5-397B-A17B-FP8 模型使用 TRTLLM MHA 后端和 NEXTN 推测解码时, DRAFT_EXTEND_V2 被错误路由到 decode kernel, 导致非法内存访问 (CUDA error: an illegal memory access was encountered)。此 PR 旨在修复该路由逻辑, 确保 draft extend 使用 context kernel。

实现拆解

1. 修改路由条件: 在 `python/sglang/srt/layers/attention/trtllm_mha_backend.py` 的 `forward_extend` 方法中, 将原先判断 `is_target_verify()` or `is_draft_extend_v2()` 的复合条件简化为仅判断 `is_target_verify()`。
2. 保留 decode kernel 给 TARGET_VERIFY: 当 `forward_mode.is_target_verify()` 为真时, 依然使用 `flashinfer.decode.trtllm_batch_decode_with_kv_cache`; 其余情况 (包括 DRAFT_EXTEND_V2) 使用 `flashinfer.prefill.trtllm_batch_context_with_kv_cache`。
3. 无其他文件变动: 仅此一个文件, 改动量极小 (1 行新增, 4 行删除)。

关键文件:

- `python/sglang/srt/layers/attention/trtllm_mha_backend.py` (模块 注意力; 类别 source; 类型 core-logic): 核心注意力后端, 修改了 `forward_extend` 方法中的 kernel 路由逻辑, 是 PR 的唯一变动文件。

关键符号: 未识别

关键源码片段

`python/sglang/srt/layers/attention/trtllm_mha_backend.py`

核心注意力后端, 修改了 `forward_extend` 方法中的 kernel 路由逻辑, 是 PR 的唯一变动文件。

第 851 行附近: `forward_extend` 方法中的 kernel 路由选择

```

# 本 PR 将 DRAFT_EXTEND_V2 从 decode 路径中移除, 避免非法内存访问
if forward_batch.forward_mode.is_target_verify():
    # 只有 TARGET_VERIFY 才使用 decode kernel
    o = flashinfer.decode.trtllm_batch_decode_with_kv_cache(
        query=q,
        kv_cache=kv_cache,
        workspace_buffer=self.workspace_buffer,
        block_tables=page_table,
        seq_lens=self.forward_metadata.cache_seq_lens_int32,
        max_seq_len=self.max_context_len,
        bmm1_scale=bmm1_scale,
        bmm2_scale=bmm2_scale,
        window_left=layer.sliding_window_size,
        sinks=attention_sink,
        skip_softmax_threshold_scale_factor=(
            envs.SGLANG_SKIP_SOFTMAX_DECODE_THRESHOLD_SCALE_FACTOR.get()
        ),
        out_dtype=self.q_data_type,
        q_len_per_req=self.forward_metadata.max_seq_len_q,
    )
else:
    # 所有其他模式 (包括 DRAFT_EXTEND_V2) 走 context kernel
    o = flashinfer.prefill.trtllm_batch_context_with_kv_cache(
        query=q,
        kv_cache=kv_cache,
        workspace_buffer=self.workspace_buffer,
        block_tables=page_table,
        seq_lens=self.forward_metadata.cache_seq_lens_int32,
        max_q_len=self.forward_metadata.max_seq_len_q,
        max_kv_len=self.max_context_len,
        bmm1_scale=bmm1_scale,
        bmm2_scale=bmm2_scale,
        batch_size=self.forward_metadata.cu_seq_lens_q.shape[0] - 1,
        cum_seq_lens_q=self.forward_metadata.cu_seq_lens_q,
        cum_seq_lens_kv=self.forward_metadata.cu_seq_lens_k,
        window_left=layer.sliding_window_size,
        sinks=attention_sink,
        # 省略部分参数, 延续原有代码风格
    )

```

评论区精华

审核人 Fridge003 直接批准, 无 review 评论。但后续 yhyang201 在 issue 中提供了另一个 config 的数据: 在 Qwen3.5-NVFP4、TP=4、bs=1、accept_len 固定为 4 的场景下, decode 路径不会触发 IMA 且性能提升 4.7% (530→555 tok/s), 推测原因是 context 路径使用低占用的 8-CTA PersistentContext 内核, 而 decode 路径使用 KV-split MultiCtasKv 内核更好地填充了 SM。

- 暂无高价值评论线程

风险与影响

- 风险：本 PR 只涉及单行路由条件变更，风险较低。主要风险在于：
 - 对 DRAFT_EXTEND_V2 使用 context kernel 可能带来性能回归（如 yhyang201 观察到的场景），但这是保证正确性的必要妥协。
 - 如果未来有其他模式（如 DRAFT_EXTEND_V1）也需要 decode kernel，需额外调整。
 - 影响：影响范围：仅影响使用 TRTLLM MHA 后端并启用推测解码（NEXTN/MTP）的用户，特别是 Qwen3.5 等大模型。修复后，原本因 CUDA 非法内存访问而失败的推理任务可以正常运行，但 draft extend 阶段可能略有性能下降。影响程度：对受影响用户是关键 bugfix，优先级高；对其他用户无影响。
- 风险标记：核心路径变更

关联脉络

- 暂无明显关联 PR