

PR #24125 完整报告

sgl-project/sglang

[AMD] Skip redundant CatArrayBatchedCopy in GLM-5 NSA TileLang decode

合并时间: 2026-05-13 17:55

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/24125>

执行摘要

- 一句话: 跳过 GLM-5 NSA TileLang 解码中冗余的 CatArrayBatchedCopy
- 推荐动作: 值得精读。该 PR 展示了如何通过分析数据流中的冗余操作实现零成本优化, 并通过精密的条件控制确保向后兼容。设计决策 (零拷贝视图、gated cat-skip、仅限 HIP) 可作为同类优化工程的范例。重点关注 forward_mla.py 中 forward_absorb_core 的 decode 分支和 nsa_backend.py 中 forward_decode 的 q_all 传递逻辑。

功能与动机

在 GLM-5 NSA TileLang decode on ROCm 上, fused-rope 路径调用了冗余的 `CatArrayBatchedCopy<OpaqueType<1u>, ...>` 内核, 每次 decode step 每层执行一次, 重建一个已经存在的张量 (与 `q_cat` 字节相同)。这是一种纯开销, 需要消除。PR body 明确描述: “The cat is pure overhead — same data, same layout, fresh allocation + copy.” 关联 Issue #2879 提供了 preshuffled layout 的支持背景。

实现拆解

1. 调整调用方数据契约 (forward_mla.py: forward_absorb_core) : 在 `_skip_rope_for_nsa_tilelang_fused()` 路径中, 当处于 decode/idle 模式时, 不再将 `q_cat` 切片为 `q_nope_fused/q_pe_fused` 再分别传入 `attn_mqa`, 而是直接传递 `q_cat` 作为 `q` 参数, 并将 `q_rope` 设为 None, 同时将 llama_4_scaling 的乘法改为对 `q_cat` 前 `kv_lora_rank` 维度的原地操作。Prefill 路径保持原有 split 形式, 因为 forward_extend 要求 `q_rope` 非空。
2. 实现解码器零拷贝视图 (nsa_backend.py: forward_decode) : 在 `q_rope is not None` 分支外新增 else 分支——当 `q_rope` 为 None 时, 认为调用方已传入拼接好的 `q` (即 `q_cat`), 直接通过 `q.contiguous().view(-1, tp_q_head_num, head_dim)` 建立零拷贝视图作为 `q_all`, 并据此计算 `q_nope/q_rope` 视图。同时标注 `q_all` 非空, 供后续 impl 分支使用。
3. 条件跳过 cat 操作 (nsa_backend.py) : 在 tilelang 和 aiter impl 分支中, 将原 `if q_rope is not None` 的条件改为 `if q_all is None or not _is_hip`。这样, 当 HIP 后端且 `q_all` 已由零拷贝视图提供时, 跳过内部的 `concat_mla_absorb_q_general` 或 `torch.cat`; 非 HIP 后端始终保持原有 re-cat 行为, 字节级一致。flashmla_sparse 和 flashmla_kv 分支保持原状, 因为它们依赖于 CUDA 驱动, 在 HIP 上不可达。
4. 精度与性能验证: 在 MI355X TP=8 上提供 GSM8K 精度数据 (0.941) 和端到端基准 (吞吐 +1.4%, TPOT -1.0%), 并确认 CI 中失败任务均未涉及本 PR 代码路径。

关键文件:

- `python/sglang/srt/models/deepseek_common/attention_forward_methods/forward_mla.py` (模块 注意力层; 类别 `source`; 类型 `data-contract`; 符号 `forward_absorb_core`): 修改了核心注意力前向方法 `forward_absorb_core`, 在 `decode` 路径下直接传递 `q_cat` 并设 `q_rope=None`, 改变了与 `attn_mqa` 的数据契约, 是 `cat-skip` 的入口端。
- `python/sglang/srt/layers/attention/nsa_backend.py` (模块 注意力层; 类别 `source`; 类型 `core-logic`; 符号 `forward_decode`): 修改了 `forward_decode` 方法, 添加零拷贝视图分支并调整 `tilelang/aiter` 的 `cat` 逻辑, 是 `cat-skip` 的接收端。

关键符号: `forward_absorb_core`, `forward_decode`

关键源码片段

`python/sglang/srt/models/deepseek_common/attention_forward_methods/forward_mla.py`

修改了核心注意力前向方法 `forward_absorb_core`, 在 `decode` 路径下直接传递 `q_cat` 并设 `q_rope=None`, 改变了与 `attn_mqa` 的数据契约, 是 `cat-skip` 的入口端。

```
def forward_absorb_core(...):
    # ... 省略前置代码 ...
    q_cat, _, k_pe_fused, _ = fused_qk_rope_cat_and_cache_mla(...)
    save_kv_cache = False
    # On decode, pass q_cat directly to attn_mqa with q_rope=None so
    # nsa_backend.forward_decode reuses q_cat as a zero-copy view
    # (`q.contiguous().view(...)` fast-path) instead of running the
    # redundant `concat_mla_absorb_q_general(q_nope_fused, q_pe_fused)`
    # that would otherwise rebuild a tensor byte-identical to q_cat.
    if forward_batch.forward_mode.is_decode_or_idle():
        if llama_4_scaling is not None:
            # llama_4_scaling applies only to the q_nope portion;
            # mutate in place via the slice view of q_cat.
            q_cat[..., :self.kv_lora_rank] *= llama_4_scaling
        attn_output = self.attn_mqa(
            q_cat, # pass full q_cat directly
            None,
            None,
            forward_batch,
            q_rope=None, # signal: already concatenated
            k_rope=k_pe_fused,
            save_kv_cache=save_kv_cache,
            **(dict(topk_indices=topk_indices) if topk_indices is not None else {}),
        )
    else:
        # Prefill keeps split form because forward_extend asserts q_rope is not None
        q_nope_fused = q_cat[..., :self.kv_lora_rank]
        q_pe_fused = q_cat[..., self.kv_lora_rank:]
        if llama_4_scaling is not None:
```

```

    q_nope_fused *= llama_4_scaling
    attn_output = self.attn_mqa(
        q_nope_fused,
        None, None, forward_batch,
        q_rope=q_pe_fused,
        k_rope=k_pe_fused,
        save_kv_cache=save_kv_cache,
        **(dict(topk_indices=topk_indices) if topk_indices is not None else {}),
    )

```

python/sglang/srt/layers/attention/nsa_backend.py

修改了 `forward_decode` 方法，添加零拷贝视图分支并调整 `tilelang/aiter` 的 `cat` 逻辑，是 `cat-skip` 的接收端。

```

def forward_decode(self, ..., q_rope, ...):
    # ... 省略前置代码 ...
    if q_rope is not None:
        q_nope = q.view(-1, layer.tp_q_head_num, layer.v_head_dim)
        q_rope = q_rope.view(-1, layer.tp_q_head_num, layer.head_dim - layer.v_head_dim)
        q_all = None # signal: need concat in impl block
    else:
        # Caller passed already-concatenated q (q_all = q). Reuse it directly
        # via a zero-copy view; the impl-specific blocks below will skip the
        # otherwise redundant concat_mla_absorb_q_general call.
        q_all = q.contiguous().view(-1, layer.tp_q_head_num, layer.head_dim)
        q_nope = q_all[:, :, :layer.v_head_dim]
        q_rope = q_all[:, :, layer.v_head_dim:]

    # ... page_table setup ...

    if self.nsa_decode_impl == "tilelang":
        # Cat-skip (HIP-only): when caller passes q_rope=None on HIP, q_all
        # has already been set to a zero-copy view; the `not_is_hip` clause
        # keeps CUDA / MUSA paths byte-identical by always re-cat.
        if q_all is None or not _is_hip:
            q_all = concat_mla_absorb_q_general(q_nope, q_rope)
        return self._forward_tilelang(q_all=q_all, ...)
    elif self.nsa_decode_impl == "aiter":
        if q_all is None or not _is_hip:
            q_all = torch.cat([q_nope, q_rope], dim=-1)
        return self._forward_aiter(q_all=q_all, ...)
    # ... other impls unchanged ...

```

评论区精华

Review 中核心讨论围绕修改范围的限制：

- 1am9trash指出：“I think we can only make the change in amd side (e.g. tilelang/aiter backend). This change may be never reached in nv code path.” 建议将 `cat-skip` 严格

限定在 AMD 后端。

- Jacob0226 采纳并回复: “Good catch, thanks! Done in 1f2b7c48c.” 随即提交修正: 恢复 flashmla_sparse 和 flashmla_kv 的原 if q_ropes is not None 逻辑, 仅保留 tilelang 和 aiter 的 cat-skip, 并明确注释为 HIP-only。
- 最终获得 1am9trash 和 kkHuang-amd 的 Approval。
- 将 cat-skip 限制在 AMD 后端 (design): 采纳 reviewer 建议, 仅 tilelang 和 aiter 分支保留 cat-skip, flashmla 分支保持原状。

风险与影响

- 风险:

1. 数据流契约变更: 调用方 (forward_absorb_core) 改变了 attn_mqa 的参数约定 (q_ropes=None), 可能影响其他未预期的后端。但通过严格条件 forward_mode.is_decode_or_idle() 和仅在 _skip_ropes_for_nsa_tilelang_fused() 路径下执行, 风险可控。
2. 非 HIP 后端行为保持: 通过 if q_all is None or not _is_hip 和保留原 flashmla 分支, 确保 CUDA/MUSA/XPU 等后端的行为字节级不变, 降低回归风险。
3. 原地内存操作: q_cat[..., :self.kv_lora_rank] *= llama_4_scaling 对 q_cat 进行原地乘法, 可能影响后续 q_cat 的其他引用。但 q_cat 在此处是 fused_qk_ropes_cat_and_cache_mla 的局部输出, 后续仅用于 attn_mqa 调用, 因此安全。
4. 缺少单元测试: PR 未添加新的单元测试, 仅依赖端到端精度和性能测试, 可能遗漏边缘情况 (如 batch size 极端值、与 topk_indices 交互等)。- 影响: 用户: 仅影响使用 GLM-5 模型且在 AMD GPU (MI355X) 上启用 NSA TileLang decode 的用户, 将获得约 1.4% 吞吐提升和 1% TPOT 降低。其他模型和后端无影响。系统: 无新增依赖、配置或环境变量更改。团队: AMD 相关 CI 任务 (pr-test-amd) 会覆盖该路径; 维护者需注意未来调整 forward_absorb_core 或 nsa_backend 时保持此数据契约的一致性。

- 风险标记: 核心路径数据流变更, 仅限 AMD/ROCM, 缺少测试覆盖

关联脉络

- PR #23562 [AMD] Preshuffled paged MQA + page_size=64 for GLM-5 NSA TileLang decode: 被 PR body 标记为 baseline, 本 PR 的优化基于该 PR 提供的 preshuffled 布局和 aiter 兼容路径。
- PR #2879 Support preshuffled layout in indexer_k_quant_and_cache / cp_gather_indexer_k_quant_cache: 关联的外部 Issue (ROCM/aiter#2879), 为本 PR 提供必要的 preshuffled 底层支持, 与本 PR 共同构成完整的性能优化栈。