

PR #43898 完整报告

vllm-project/vllm

[ROCm][DSv4] Remove device pipeline stall in sparse attention

合并时间: 2026-05-29 15:42

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/43898>

执行摘要

- 一句话: 消除稀疏注意力 GPU 气泡
- 推荐动作: 建议合入。此 PR 很好地展示了如何通过消除 GPU 微气泡来提升性能, 是 ROCm 上 DSv4 推理链路中的一次精细优化。值得关注的设计点: 用 `torch.zeros` 合并赋值操作减少 kernel launch、用已知 host 值替代 D2H 同步获取 `indptr[-1]`。

功能与动机

在 ROCm 的 DSv4 路径上, 每次 `build_ragged_indices_from_dense` 调用因 `indptr[0] = 0` 的 host 开销和 `indptr[-1].item()` 的 D2H 拷贝同步导致明显的 GPU 气泡。PR 旨在通过消除同步来提升 GPU 利用率和整体吞吐。

实现拆解

1. 消除 host-device 同步: 将 `indptr` 的创建从 `torch.empty + indptr[0] = 0` (需 host 执行赋值, 造成同步) 改为 `torch.zeros(...)`, 利用 GPU 上的 zero kernel 一次性完成初始化, 避免独立的 host 触发 kernel。
2. 替换动态形状计算: 将 flat 缓冲区的尺寸从 `int(indptr[-1].item())` (需要 D2H 拷贝 `indptr[-1]` 的值, 触发同步) 改为静态计算 `indices.shape[0] * max_width` (即 `tokens * max_width`), 该值在 host 侧已知, 无需等待 device 结果。
3. 兼容性验证: PR 声明静态缓冲区大小与所有下游消费者 (`_sparse_attn_prefill_ragged_kernel`、`_sparse_attn_decode_ragged_kernel`、`_copy_ragged_to_graph_buffers`) 兼容, 因为这些消费者通过 `indptr[i]:indptr[i+1]` 访问, 不会读取超出 `indptr[-1]` 的范围。

关键文件:

- `vllm/v1/attention/ops/rocm_aiter_mla_sparse.py` (模块 注意力; 类别 source; 类型 core-logic; 符号 `build_ragged_indices_from_dense`): 唯一被修改的文件, 包含稀疏注意力构造函数 `build_ragged_indices_from_dense`, 通过两处关键优化消除 GPU 气泡。

关键符号: `build_ragged_indices_from_dense`

关键源码片段

`vllm/v1/attention/ops/rocm_aiter_mla_sparse.py`

唯一被修改的文件，包含稀疏注意力构造函数 `build_ragged_indices_from_dense`，通过两处关键优化消除 GPU 气泡。

```
# vllm/v1/attention/ops/rocm_aiter_mla_sparse.py 中 build_ragged_indices_from_dense 函数片段
# 优化：消除 GPU 气泡，减少 host-device 同步
```

```
def build_ragged_indices_from_dense(indices, lengths):
    max_width = indices.shape[1] if indices.ndim == 2 else 0
    lengths = lengths.clamp(min=0, max=max_width).contiguous()

    # 原写法: torch.empty + indptr[0] = 0 导致 host 参与赋值，产生同步
    # 新写法: torch.zeros 一次性在 device 上清零，仅一次 kernel launch
    indptr = torch.zeros(
        indices.shape[0] + 1, dtype=torch.int32, device=indices.device
    )
    torch.cumsum(lengths, dim=0, out=indptr[1:])

    if indices.numel() == 0:
        flat = torch.empty(0, dtype=torch.int32, device=indices.device)
    else:
        # 原写法: int(indptr[-1].item()) 需要 D2H 拷贝 indptr[-1]，触发同步
        # 新写法: 使用 host 侧已知的 tokens * max_width，无需等待 device
        flat = torch.empty(
            indices.shape[0] * max_width,
            dtype=torch.int32,
            device=indices.device,
        )
    # 后续填充和数据拷贝操作，下游消费者通过 indptr[i]:indptr[i+1] 访问，不会越界
    if flat.numel() > 0:
        # ... 填充逻辑
```

评论区精华

Reviewer [AndreasKaratzas](#) 询问是否可以使用 `torch.empty` 代替 `torch.zeros` 来进一步提升性能 ([Did you try empty here as well? It's a bit faster](#))。作者 [kliuae](#) 回复经验证两者性能相同，选择 `zeros` 是为了只产生一次 host dispatch 和一次 kernel enqueue，避免额外的 fill kernel launch。讨论了 GPU 调度的微观优化权衡，最终达成一致。

- 使用 `torch.zeros` 还是 `torch.empty` 创建 `indptr` (performance): 保持使用 `torch.zeros`，因为性能等价且 kernel 调度次数更少。

风险与影响

- 风险：主要风险在 `flat` 缓冲区静态尺寸增大可能的内存浪费：`indices.shape[0] * max_width` 可能远大于实际需要的 `indptr[-1]`，但作者通过分析三个下游消费者的索引行为（均以 `indptr` 为界）确认不会越界，且 `_copy_ragged_to_graph_buffers` 的预分配缓冲区 `max_entries_per_row == max_width` 确保兼容。内存增加量有限（`tokens * max_width` 通常接近实际使用量），风险低。

- 影响：影响范围：仅影响 ROCm 平台上 DeepSeek-V4 使用稀疏注意力的推理路径。性能提升：根据 benchmark, Output Token Throughput 提升 2.8%, TTFT 下降 6.07%, ITL 下降 1.48%。lm_eval gsm8k 准确率: 0.9515 (与 baseline 持平), 无精度回归。影响程度：中等, 针对特定硬件和模型优化, 代码改动量极小 (+4/-3), 但性能收益显著。
- 风险标记：仅影响 ROCm 平台, 未新增测试覆盖, 存在潜在内存浪费

关联脉络

- PR #43905 [DSv4] Move mHC tilelang kernels & Don't use CustomOP in dsv4/nvidia: 同为 DSv4 模型的性能优化 / 重构 PR, 涉及不同的底层实现 (tilelang kernel 迁移 vs 稀疏注意力同步消除), 属于同一模型家族的不同优化方向。
- PR #43891 [Model Refactoring] Remove unnecessary torch op registration for DSv4: 同为 DSv4 模型的清理 / 优化 PR, 移除不必要的 torch op 注册, 与当前 PR 共同提升 DSv4 在 ROCm 上的推理效率。