

# PR #37813 完整报告

vllm-project/vllm

[Perf] fuse kernels in gdn

合并时间: 2026-04-02 19:52

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/37813>

## 执行摘要

- 一句话: 融合 GDN 层的后卷积操作内核, 提升 Qwen 模型推理性能。
- 推荐动作: 建议技术管理者和工程师精读此 PR, 关注内核融合的设计决策 (如 Triton 内核的网格划分、内存布局优化) 和性能权衡。特别留意 review 中讨论的数值稳定性问题, 可作为未来内核开发的借鉴。

## 功能与动机

PR body 中的基准测试结果显示, 在相同配置下, 使用融合内核后输出 token 吞吐量从 351.64 tok/s 提升到 357.31 tok/s (约 1.6% 提升), 平均每输出 token 时间从 406.68 ms 减少到 400.48 ms。这证明了性能优化的需求, 旨在通过内核融合减少操作链开销, 提升 Qwen 模型的推理效率。

## 实现拆解

实现主要包括三个部分:

1. 新增 Triton 融合内核文件 `fused_gdn_prefill_post_conv.py`, 实现 `fused_post_conv_prep` 函数, 将 `split`、`rearrange`、`contiguous`、`l2norm` 和 `gating` 操作融合为单次内核执行。
2. 修改 `gdn_linear_attn.py` 中的 `_forward_core` 方法, 在预填充路径中使用新融合内核替换原有的离散操作序列, 并调整 L2 归一化设置。
3. 新增测试文件 `test_fused_gdn_post_conv.py`, 提供参考实现和全面测试, 验证内核正确性; 同时更新 `__init__.py` 导出新函数。

关键文件:

- `tests/kernels/test_fused_gdn_post_conv.py` (模块 测试): 新增的全面测试文件, 验证融合内核与参考实现的正确性, 覆盖多种配置 (如 Qwen3.5-35B/397B 配置、不同序列长度、L2 归一化开关), 是确保变更无回归的关键。
- `vllm/model_executor/layers/fla/ops/fused_gdn_prefill_post_conv.py` (模块 内核层): 新增的融合 Triton 内核实现, 核心性能优化点, 将多个操作融合为单次内核执行, 直接决定性能提升效果。
- `vllm/model_executor/layers/mamba/gdn_linear_attn.py` (模块 模型层): 修改的核心模型文件, 集成新融合内核到 `GDNLinearAttention` 层的预填充路径, 影响模型推理逻辑。

- vllm/model\_executor/layers/fla/ops/\_\_init\_\_.py (模块 内核层) : 更新导出列表, 添加新函数 fused\_post\_conv\_prep, 确保其他模块可访问。

关键符号: fused\_post\_conv\_prep, reference\_post\_conv, \_forward\_core

## 评论区精华

review 中主要讨论点:

- gemini-code-assist[bot] 指出融合内核中 softplus 实现  $\text{tl.log}(1.0 + \text{tl.exp}(x))$  可能存在数值稳定性问题, 建议使用更稳健的实现 (如  $\text{tl.where}(x > 0, x + \text{tl.log}(1.0 + \text{tl.exp}(-x)), \dots)$ ), 以避免大正值下的精度问题。
- 同一 review 还建议在 gdn\_linear\_attn.py 的断言中添加描述性错误信息, 以增强调试能力。
- vadiklyutiy 在评论中询问了关于非 torch.compile 边界缩短以自动融合的可能性, 作者 ZJY0516 回应表示已尝试且当前是最小边界。最终 vadiklyutiy 批准 PR, 并提到运行了 Qwen3.5 评估。
- softplus 数值稳定性问题 (correctness): 未在 PR 中直接解决, 但 review 提供了改进建议; 作者可能后续考虑, 当前依赖阈值设置。
- 断言错误信息增强 (style): PR 已采纳建议, 在最终代码中添加了错误消息。

## 风险与影响

- 风险: 风险点包括:
  - 数值稳定性风险: 融合内核中的 softplus 实现可能在大输入值时出现精度问题, 尽管当前阈值设置为 20.0, 但未采用更稳健方法, 可能影响模型输出一致性。
  - 正确性风险: 新内核替换了多个离散操作, 需确保与原始数学等价; 测试覆盖了多种配置 (如不同头数、序列长度、L2 归一化开关), 但可能未覆盖所有边缘情况。
  - 集成风险: 修改 gdn\_linear\_attn.py 可能影响其他路径 (如解码阶段), 需验证 L2 归一化设置 (use\_qk\_l2norm\_in\_kernel=False) 的正确性。
- 影响: 影响范围:
  - 用户影响: 对使用 Qwen 模型 (特别是 GDN 层) 的用户透明, 可能体验轻微性能提升, 但需确保模型输出无回归。
  - 系统影响: 减少预填充阶段的计算和内存开销, 提升推理吞吐量, 基准测试显示 token 吞吐量提升约 1.6%, 可能在高并发场景下放大效益。
  - 团队影响: 引入新的 Triton 内核, 增加了代码维护复杂度, 但提供了内核融合的设计示例, 可供其他优化参考。
- 风险标记: 数值稳定性风险, 集成逻辑变更

## 关联脉络

- PR #32996 Feature/silu block quant fusion v1: 类似的内核融合优化, 将 SiLU 乘法与分块 FP8 量化融合, 提升推理性能, 与本 PR 同属性能优化类别。

- PR #38684 [Perf] DSV3.2 Indexer Fused Weights Projection: 融合 DeepSeek V3.2 索引器中的投影层, 优化推理性能, 与本 PR 共享性能改进和模型特定优化主题。
- PR #38086 [ROCm] Enable VLLM triton FP8 moe for gfx1201, tuned for Qwen3-30B-A3B-FP8 tp=2 and Qwen/Qwen3.5-35B-A3B-FP8 tp=2: 针对 Qwen 模型的 Triton 内核优化, 涉及 FP8 MoE 后端, 与本 PR 在 Qwen 模型和内核调优方面有交集。