

# PR #39083 完整报告

vllm-project/vllm

[FEAT] [Perf] [Gemma4] Fused Gemma4 Routing Function Triton

合并时间: 2026-04-19 17:57

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/39083>

## 执行摘要

- 一句话: 为 Gemma4 模型添加 Triton 融合路由函数, 显著提升 MoE 推理性能。
- 推荐动作: 建议工程师精读此 PR, 特别是 Triton 内核设计部分, 展示了如何通过向量化排序和减少内存操作优化 MoE 路由。关注性能权衡 (如 num\_warps 设置) 和数值稳定性处理 (如硬编码常数)。对于从事内核优化或模型特定加速的开发者, 这是一个有价值的案例。

## 功能与动机

根据 PR body, 目的是提高 Gemma4 模型的性能, 因为原自定义路由函数引入许多同步点和全局内存读写, 且不被 torch.compile 捕获, 导致性能瓶颈。新 Triton 内核旨在减少这些开销, 提升推理效率。

## 实现拆解

1. 新增 Triton 路由内核: 在 vllm/model\_executor/models/gemma4.py 中定义 @triton.jit 装饰的 \_gemma4\_routing\_kernel 函数, 实现向量化排序和 top-K 选择, 通过整数位转换和掩码操作减少同步和全局内存访问。
2. 提供 Python 包装函数: 同一文件中添加 gemma4\_fused\_routing\_kernel\_triton 函数, 封装内核调用, 处理张量连续性和内存分配, 默认 num\_warps=1。
3. 更新模型路由逻辑: 在 Gemma4MoE 类中, 将自定义路由函数替换为新 Triton 内核, 通过 gemma4\_fused\_routing\_kernel\_triton 调用, 确保与原有 PyTorch 参考实现 gemma4\_routing\_function\_torch 兼容。
4. 添加测试覆盖: 创建 tests/kernels/moe/test\_gemma4router.py, 定义 test\_gemma4\_routing\_kernel\_triton 测试函数, 使用参数化测试验证 Triton 内核与参考实现的数值等价性, 包括边缘情况如最大上下文长度 250K。
5. 微调配置: 在 pyproject.toml 中添加 VALU 常量, 可能用于后续 Triton 内核开发或文档。

关键文件:

- vllm/model\_executor/models/gemma4.py (模块 模型执行; 类别 source; 类型 core-logic; 符号 \_gemma4\_routing\_kernel, gemma4\_fused\_routing\_kernel\_triton, gemma4\_routing\_function\_torch): 核心实现文件, 包含新增的 Triton JIT 路由内核、Python 包装函数及模型集成逻辑, 直接决定 Gemma4 MoE 路由性能。
- tests/kernels/moe/test\_gemma4router.py (模块 测试; 类别 test; 类型 test-coverage; 符号 sort\_by\_id, test\_gemma4\_routing\_kernel\_triton): 新增测试文件, 验证 Triton 路

由内核与 PyTorch 参考实现的数值等价性，覆盖多种 token 数量、数据类型和边缘情况，确保功能正确性。

- pyproject.toml (模块 配置; 类别 config; 类型 configuration) : 配置文件微调, 添加 VALU 常量, 可能用于 Triton 内核开发或文档, 影响较小但显示了对工具链的更新。

关键符号: `_gemma4_routing_kernel`, `gemma4_fused_routing_kernel_triton`, `gemma4_routing_function_torch`, `sort_by_id`, `test_gemma4_routing_kernel_triton`

## 关键源码片段

### vllm/model\_executor/models/gemma4.py

核心实现文件, 包含新增的 Triton JIT 路由内核、Python 包装函数及模型集成逻辑, 直接决定 Gemma4 MoE 路由性能。

```
@triton.jit
def _gemma4_routing_kernel(
    gating_ptr,
    per_expert_scale_ptr,
    topk_weights_ptr,
    topk_ids_ptr,
    E: tl.constexpr,
    K: tl.constexpr,
    BLOCK_E: tl.constexpr,
):
    pid = tl.program_id(0) # 每个程序处理一个 token 序列
    offs_e = tl.arange(0, BLOCK_E)
    valid = offs_e < E # 掩码, 处理实际专家数

    # 加载门控输出并转换为 float32 以确保数值稳定性
    logits = tl.load(
        gating_ptr + pid * E + offs_e,
        mask=valid,
        other=-float("inf"),
    ).to(tl.float32)

    max_l = tl.max(logits, axis=0) # 计算最大值用于稳定 softmax

    # 将 float32 转换为可排序的整数键, 避免浮点排序开销
    MIN32 = -2147483648
    logit_bits = logits.to(tl.int32, bitcast=True)
    sign_b = logit_bits >> 31
    key = tl.where(sign_b == 0, logit_bits ^ -1, logit_bits ^ MIN32)
    key = tl.where(valid, key, 0x7FFFFFFF) # 无效位置设为最大值以排到末尾
    sk64 = key.to(tl.int64) & 0x00000000FFFFFFFF
    packed = (sk64 << 32) | offs_e.to(tl.int64) # 打包键和索引
    sorted_p = tl.sort(packed, descending=False) # 向量化排序

    # 从排序结果中提取所有键和专家 ID
```

```

all_keys = ((sorted_p >> 32) & 0x00000000FFFFFFFF).to(tl.int32)
all_ids = (sorted_p & 0x00000000FFFFFFFF).to(tl.int32)

# 逆转换恢复原始 logit 位
sign_k = all_keys >> 31
all_bits = tl.where(sign_k < 0, all_keys ^ -1, all_keys ^ MIN32)
all_logits = all_bits.to(tl.float32, bitcast=True)

# 计算所有元素的 raw_exp (使用 exp2 和硬编码常数近似 log2(e))
all_raw_exp = tl.math.exp2((all_logits - max_l) * 1.4426950408889634)

# 仅对 top-K 元素求和并进行重归一化
top_mask = offs_e < K
renorm_raw = tl.sum(tl.where(top_mask, all_raw_exp, 0.0), axis=0)
renorm_raw = tl.where(renorm_raw > 0.0, renorm_raw, 1.0)
inv_renorm = 1.0 / renorm_raw

# 加载 top-K 专家的缩放因子
all_scales = tl.load(
    per_expert_scale_ptr + all_ids.to(tl.int64),
    mask=top_mask,
    other=1.0,
).to(tl.float32)

# 计算最终权重: 归一化后乘以缩放因子
all_weights = (all_raw_exp * inv_renorm * all_scales).to(tl.float32)

# 存储结果: top-K 权重和 ID
base_off = pid * K + offs_e
tl.store(topk_ids_ptr + base_off, all_ids, mask=top_mask)
tl.store(topk_weights_ptr + base_off, all_weights, mask=top_mask)

```

## tests/kernels/moe/test\_gemma4router.py

新增测试文件，验证 Triton 路由内核与 PyTorch 参考实现的数值等价性，覆盖多种 token 数量、数据类型和边缘情况，确保功能正确性。

```

def sort_by_id(w, ids):
    order = ids.argsort(dim=-1) # 按专家 ID 排序以消除并列差异
    return w.gather(1, order), ids.gather(1, order)

# Gemma4 MoE 模型支持最大上下文长度 250K，测试覆盖边界值
@pytest.mark.parametrize("num_tokens", [1, 2, 2048, 250000])
@pytest.mark.parametrize("num_experts", [128]) # Gemma4 专家数
@pytest.mark.parametrize("topk", [8]) # Gemma4 top-k 值
@pytest.mark.parametrize("dtype", [torch.bfloat16, torch.half, torch.float32])
def test_gemma4_routing_kernel_triton(
    num_tokens: int,
    num_experts: int,
    topk: int,

```

```

dtype: torch.dtype,
):
torch.manual_seed(0) # 固定随机种子以确保可重复性
gating = torch.randn(num_tokens, num_experts, dtype=dtype, device="cuda")
scales = torch.rand(num_experts, dtype=torch.float32, device="cuda")

# 分别获取 PyTorch 参考实现和 Triton 内核的结果
ref_w, ref_ids = gemma4_routing_function_torch(gating, topk, scales)
tri_w, tri_ids = gemma4_fused_routing_kernel_triton(gating, topk, scales)

# 排序以处理并列情况
ref_ws, ref_is = sort_by_id(ref_w, ref_ids)
tri_ws, tri_is = sort_by_id(tri_w, tri_ids)

ids_match = (ref_is == tri_is).all().item() # 检查专家 ID 是否一致
weights_match = torch.allclose(ref_ws, tri_ws, atol=1e-2, rtol=1e-2) # 检查权重差异在容差内
all_match = ids_match and weights_match
if not all_match:
    bad = (ref_is != tri_is).any(dim=-1).nonzero(as_tuple=True)[0]
    if len(bad):
        r = bad[0].item()
        print(f"第一行错误 {r}: ref_ids={ref_ids[r].tolist()} tri_ids={tri_ids[r].tolist()}")
    assert all_match # 断言确保所有测试通过

```

## 评论区精华

review 中核心讨论点:

- 使用 LOG2E 常量: gemini-code-assist[bot] 建议使用 `tl.LOG2E` 提高可读性, 但作者 `tjtanaa` 指出 Triton 内核不支持非 `constexpr` 全局变量, 因此保持硬编码常数。
- `num_warps` 默认值: gemini-code-assist[bot] 建议 `num_warps=4` 以更好利用 GPU, 作者基于内部基准坚持使用 `num_warps=1`, 认为这是最优设置。
- `docstring` 修正: gemini-code-assist[bot] 指出 `docstring` 中的复制粘贴痕迹, 作者已更新以准确描述内核功能。
- 与 `torch.compile` 比较: ProExpertProg 询问是否可直接应用 `torch.compile`, 作者提供微基准测试数据, 显示 Triton 内核在 A100、H100、MI300X 等 GPU 上仍优于编译后的 PyTorch 版本。
- 测试容差设置: ZJY0516 质疑测试中 `atol=1e-2` 是否太大, 作者解释 `bfloat16` 有较大数值发散, 并引用现有 `fused_topk` 测试阈值以确保一致性。
  - 使用 LOG2E 常量建议 (`correctness`): 作者拒绝建议, 保持硬编码常数 `1.4426950408889634`, 确保内核可执行。
  - `num_warps` 默认值设置 (`performance`): 保持 `num_warps=1` 作为默认值, 避免不必要的性能开销。
  - 与 `torch.compile` 性能比较 (`performance`): Triton 内核性能优势明显, 因此选择实现新内核而非依赖编译优化。

- 测试容差设置质疑 (testing): 保持现有容差设置, 认为足以覆盖数值差异且与项目标准一致。

## 风险与影响

- 风险: 技术风险包括:
  - 数值精度: 在 bfloat16 等低精度数据类型下, Triton 内核可能与 PyTorch 参考实现有微小差异 (测试中最大误差约  $3.89e-03$ ), 虽然测试容差已覆盖, 但需监控实际部署中的数值稳定性。
  - 兼容性: 优化仅针对 Gemma4 模型 (128 专家、top-k=8), 不通用; 且依赖 Triton 和 CUDA 平台, 可能影响非 CUDA 后端或未来模型扩展。
  - 维护性: 新增 Triton 内核增加了代码复杂性, 未来更改需确保内核优化逻辑保持, 并可能引入版本兼容性问题。
- 影响: 影响分析:
  - 用户: Gemma4 模型用户将获得显著的推理性能提升, 端到端吞吐量在 A100、H100、MI300X、B60 等 GPU 上提升 4-32%, 降低延迟。
  - 系统: 减少路由计算的开销, 降低同步和内存访问, 提高 GPU 资源利用率, 可能为类似模型优化提供参考。
  - 团队: 为模型特定性能优化树立模板, 但增加了 Triton 内核维护负担, 需确保测试覆盖充分。
  - 风险标记: 数值精度差异, 模型特定优化, Triton 依赖

## 关联脉络

- PR #40273 Fix MoE backend selection for LoRA (unquantized MoE): 同样涉及 MoE 层优化, 但针对后端选择而非路由函数, 可参考对 MoE 组件的维护模式。
- PR #40143 [Core] Reduce mm scheduler, get\_num\_embed overhead: 性能优化主题相似, 通过缓存减少开销, 展示了内核级性能改进的常见手法。
- PR #39909 Added general ND x ND matmul and unit test for it: 涉及通用计算函数重构和测试, 与本 PR 的模型特定优化形成对比, 但都关注数值正确性和性能。