

# PR #33773 完整报告

vllm-project/vllm

[ROCm][FEAT] Integrate aiter gemm w8a8 ptpc

合并时间: 2026-04-16 09:55

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/33773>

## 执行摘要

- 一句话: 在 AMD ROCm 平台集成 aiter GEMM 内核, 优化 FP8 推理性能。
- 推荐动作: 建议工程师精读此 PR, 重点关注内核选择逻辑 (如 `can_implement` 方法如何实现条件分发) 以及权重处理流程。这对于理解 ROCm 平台性能优化和量化内核集成有重要参考价值。

## 功能与动机

根据 PR 描述, 目的是利用 aiter 内核优化 AMD ROCm 平台上的 FP8 推理性能, 相比标准 `torch.scaled_mm` 基线, 能为主流大语言模型带来显著的性能提升, 同时保持精度在合理范围内。

## 实现拆解

1. 新增内核类: 在 `vllm/model_executor/kernels/linear/scaled_mm/aiter.py` 中, 添加了 `AiterPreshuffledPerTokenFp8ScaledMMLinearKernel` 和 `AiterPerTokenFp8ScaledMMLinearKernel` 类, 实现了 `is_supported`、`can_implement` 等方法, 用于检查平台支持和配置兼容性 (如每令牌激活量化和每通道权重量化)。这为核心逻辑提供了性能优化的入口。
2. 更新底层操作: 修改 `vllm/_aiter_ops.py`, 增加了 `_load_gemm_tuned_configs` 和 `_check_kernel_tuned` 函数来加载和检查调优配置, 并注册了新的 GEMM 操作如 `_rocm_aiter_preshuffled_per_token_w8a8_gemm_impl`, 确保依赖库的集成和错误处理。
3. 集成到内核选择: 在 `vllm/model_executor/kernels/linear/__init__.py` 中, 导入新内核并将其添加到 ROCm 平台的 FP8 内核优先级列表中, 更新了 `init_fp8_linear_kernel` 函数的参数顺序, 以支持权重形状传递。
4. 调整量化方案: 更新多个量化方案文件 (如 `compressed_tensors_w8a8_fp8.py`、`fp8.py` 等), 修正量化键映射和权重处理逻辑, 确保与新内核的兼容性。
5. 测试配套: 在 `tests/utils.py` 中有细微调整, 但测试覆盖主要依赖于现有框架, 未添加大规模测试文件。

关键文件:

- `vllm/model_executor/kernels/linear/scaled_mm/aiter.py` (模块 内核实现; 类别 `source`; 类型 `core-logic`; 符号 `AiterPreshuffledPerTokenFp8ScaledMMLinearKernel`, `is_supported`, `can_implement`, `process_weights_after_loading`): 核心实现文件, 新增

了 AiterPreshuffledPerTokenFp8ScaledMMLinearKernel 等内核类，定义了平台支持和配置检查逻辑。

- vllm/\_aiter\_ops.py (模块 底层操作; 类别 source; 类型 dependency-wiring; 符号 \_load\_gemm\_tuned\_configs, \_check\_kernel\_tuned, \_rocm\_aiter\_preshuffled\_per\_token\_w8a8\_gemm\_impl, \_rocm\_aiter\_w8a8\_gemm\_impl) : 底层操作文件, 新增了调优配置加载函数和新的 GEMM 实现, 支持 aiter 库的集成。
- vllm/model\_executor/kernels/linear/\_\_init\_\_.py (模块 内核选择; 类别 source; 类型 entrypoint) : 内核选择入口文件, 更新了导入和内核优先级列表, 确保新内核被正确集成。
- vllm/model\_executor/layers/quantization/compressed\_tensors/schemes/compressed\_tensors\_w8a8\_fp8.py (模块 量化方案; 类别 source; 类型 data-contract) : 量化方案文件, 更新了量化键映射和权重形状处理, 以兼容新内核。

关键符号: AiterPreshuffledPerTokenFp8ScaledMMLinearKernel.is\_supported, AiterPreshuffledPerTokenFp8ScaledMMLinearKernel.can\_implement, \_rocm\_aiter\_preshuffled\_per\_token\_w8a8\_gemm\_impl, \_load\_gemm\_tuned\_configs

## 关键源码片段

### vllm/model\_executor/kernels/linear/scaled\_mm/aiter.py

核心实现文件, 新增了 AiterPreshuffledPerTokenFp8ScaledMMLinearKernel 等内核类, 定义了平台支持和配置检查逻辑。

```
class AiterPreshuffledPerTokenFp8ScaledMMLinearKernel(FP8ScaledMMLinearKernel):
    @classmethod
    def is_supported(
        cls, compute_capability: int | None = None
    ) -> tuple[bool, str | None]:
        # 检查是否在ROCm平台, 且aiter线性FP8功能已启用
        if not current_platform.is_rocm():
            return False, "requires ROCm."
        if not rocm_aiter_ops.is_linear_fp8_enabled():
            return (
                False,
                "requires setting `VLLM_ROCM_USE_AITER=1` "
                "and `VLLM_ROCM_USE_AITER_LINEAR=1`. "
                "`VLLM_ROCM_USE_AITER_LINEAR` default is True.",
            )
        try:
            import aiter # noqa: F401
        except Exception:
            return False, "requires aiter library to be installed."
        return True, None

    @classmethod
    def can_implement(cls, c: FP8ScaledMMLinearLayerConfig) -> tuple[bool, str | None]:
        # 检查配置是否支持每令牌激活量化和每通道权重量化 (PTPC)
        is_ptpc = (
```

```

        c.activation_quant_key.scale.group_shape.is_per_token()
        and c.weight_quant_key.scale.group_shape.is_per_channel()
    )
    if c.weight_shape is None:
        return False, "weight_shape is required for Aiter kernels"
    N, K = c.weight_shape
    fp8_dtype = current_platform.fp8_dtype()

    if c.out_dtype is not torch.bfloat16:
        return False, "requires bfloat16 output dtype."

    if not is_ptpc:
        return (
            False,
            "requires per token activation scales and per channel weight scales.",
        )

    # 确保维度对齐, 以利用aiter优化
    if not (N % 16 == 0 and K % 16 == 0):
        return (
            False,
            f"requires N and K dimensions divisible by 16, received "
            f"N={N} and K={K}.",
        )

    # 仅在调优配置可用时使用shuffled gemm, 否则回退到其他内核
    if not rocm_aiter_ops.is_shuffled_per_token_w8a8_gemm_tuned(N, K, fp8_dtype):
        return (
            False,
            f"requires a tuned configuration for N: {N} and K: {K} "
            f"and fp8 dtype {fp8_dtype}.",
        )

    return True, None

```

## vllm/\_aiter\_ops.py

底层操作文件, 新增了调优配置加载函数和新的 GEMM 实现, 支持 aiter 库的集成。

```

def _load_gemm_tuned_configs(
    q_dtype_w: torch.dtype, csv_path: str
) -> set[tuple[int, int, int]]:
    # 从CSV文件加载调优配置, 筛选指定数据类型, 缓存结果以提高性能
    try:
        df = pd.read_csv(csv_path).drop_duplicates()
        df = df[df["q_dtype_w"] == str(q_dtype_w)]
        return set(zip(df["N"].astype(int), df["K"].astype(int), df["M"].astype(int)))
    except Exception:
        # 文件读取失败时返回空集, 避免崩溃
        return set()

```

```

def _rocm_aiter_preshuffled_per_token_w8a8_gemm_impl(
    A: torch.Tensor,
    B: torch.Tensor,
    As: torch.Tensor,
    Bs: torch.Tensor,
    bias: torch.Tensor | None = None,
    output_dtype: torch.dtype = torch.float16,
) -> torch.Tensor:
    # 调用aiter库的bpreshuffle GEMM操作，实现高性能w8a8矩阵乘法
    from aiter import gemm_a8w8_bpreshuffle
    output = gemm_a8w8_bpreshuffle(A, B, As, Bs, None, output_dtype)
    if bias is not None:
        # 当前aiter不支持bias融合，因此单独添加
        output.add_(bias)
    return output

```

## 评论区精华

Review 中核心讨论包括：

- Bias 融合性能问题：Gemini-bot 指出在 `_rocm_aiter_preshuffled_per_token_w8a8_gemm_impl` 中，`bias` 应融合到 GEMM 操作中以提升性能，但 `vllm` 回应称当前 `aiter` 实现不支持此功能，因此保持分离操作。
- 类型提示和命名一致性：tjtanaa 建议移除冗余代码（如 `mutates_args=[]`）并保持命名语义（如使用 `'preshuffled'` 而非 `'shuffled'`），`vllm` 已相应更新。
- 无关变更处理：BadrBasowid 和 tjtanaa 质疑了量化文件中添加的 `ptpc_fp8` 等条目，确认为合并冲突导致的不必要变更，已被移除。结论：`bias` 融合暂缓，代码清理和命名优化已完成，确保 PR 专注于核心功能。
- Bias 融合在 GEMM 操作中的性能优化 (performance): 由于 `aiter` 库限制，`bias` 融合暂缓实施，保持分离操作。
- 命名一致性及代码清理 (design): 命名已更新为 `'preshuffled'`，冗余代码已移除。
- 量化文件中无关变更的处理 (correctness): 无关变更已回滚，确保 PR 专注于核心功能。

## 风险与影响

- 风险：技术风险包括：
  - 性能依赖风险：新内核的性能高度依赖 `aiter` 库的调优配置，若配置缺失或未优化，可能无法达到预期吞吐量提升。
  - 精度退化风险：FP8 量化本身可能引入精度损失，PR 中提及在 LMeval 评估中有轻微波动，需用户监控模型输出。
  - 兼容性与复杂性：新增内核选择路径增加了系统复杂性，可能影响其他量化方案或平台（如 CUDA）的稳定性，需确保回退机制可靠。
- 影响：对用户影响：在 AMD ROCm 平台上运行 FP8 量化模型（如 Qwen、DeepSeek）时，可获得显著的吞吐量提升，但需设置环境变量（如 `VLLM_ROCM_USE_AITER=1`）并确保 `aiter` 库已安装。对系统影响：扩展了线性层内核选择机制，为 ROCm 平台提供了更多优

化选项，可能轻微增加初始化开销。对团队影响：工程师需熟悉 aiter 集成细节和 FP8 量化配置，以便进一步调优和问题排查。

- 风险标记：外部依赖性能影响，精度退化风险，配置复杂性增加

## 关联脉络

- PR #36247 [ 需查询完整标题，但根据讨论推测与 DeepSeek 模型修复相关 ]: 在 review 中被提及，涉及量化方案映射的修复，与本 PR 中压缩张量文件的调整相关。