

# PR #41175 完整报告

vllm-project/vllm

[ROCm][Bugfix]: W4A4 MOE using emulation instead of AITER on MXFP4-supported hardware

合并时间: 2026-04-30 05:39

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/41175>

## 执行摘要

- 一句话: 修复 ROCm W4A4 MOE 错误使用仿真后端
- 推荐动作: 该 PR 值得精读以了解 PR#39801 引入的回归及修复方式, 展示了配置变更如何意外影响调度逻辑, 以及 review 中如何通过追问缩小修复范围。对于维护 ROCm/ 量化栈的工程师有参考价值。

## 功能与动机

PR#39801 意外关闭了所有夸克量化 W4A4 模型的 `rocm_aiter_fused_moe`, 迫使它们使用 `MXFP4.Emulation` 后端。正确的列表应包含 W4A16、W4A8 和 W4A4, 相关依据来自 AITER 调度代码 (参见链接)。

## 实现拆解

1. 更新 Quark MOE 初始化中的原生 OCP MX 方案列表: 在文件 `vllm/model_executor/layers/quantization/quark/quark_moe.py` 中, 将 `_AITER_NATIVE_OCP_MX_SCHEMES` 从 `("w_mxfp4",)` 扩展为 `("w_mxfp4", "w_mxfp4_a_mxfp4")`, 并简化对应的 TODO 注释。这使得 `w_mxfp4_a_mxfp4` (W4A4) 方案不再被认定需要仿真, 从而在满足硬件和 AITER 可用条件时使用原生 MOE 内核而非 EMULATION 后端。
2. 更新 ROCm AITER MOE 中的量化方法注释: 在文件 `vllm/model_executor/layers/fused_moe/rocm_aiter_fused_moe.py` 中, 调整注释以明确说明 `mxfp4` 包含 W4A4 和 W4A16 使用 `BLOCK_1X32`, 并注明 `MXFP6` 和 `MXFP8` 当前不支持 AITER, 改用仿真。
3. 移除 W4A8 支持: 在第二版提交中移除了对 `w4a8` 的添加, 因为 review 指出 AITER CK MOE 并不支持 W4A8。

关键文件:

- `vllm/model_executor/layers/quantization/quark/quark_moe.py` (模块 量化层; 类别 source; 类型 data-contract): 核心修复文件: 将 `_AITER_NATIVE_OCP_MX_SCHEMES` 从仅包含 `w_mxfp4` 扩展到包含 `w_mxfp4_a_mxfp4`, 使 W4A4 模型在支持 MXFP4 的硬件上正确使用原生 MOE 内核而非仿真。
- `vllm/model_executor/layers/fused_moe/rocm_aiter_fused_moe.py` (模块 MOE 层; 类别 source; 类型 data-contract): 辅助修改: 更新注释以清晰说明 MXFP4 方案 (W4A4、W4A16) 使用 `BLOCK_1X32`, 且 `MXFP6/MXFP8` 当前不支持 AITER 改用仿真。

关键符号: QuarkOCP\_MX\_MoEMethod.init, rocm\_aiter\_fused\_experts

## 关键源码片段

vllm/model\_executor/layers/quantization/quark/quark\_moe.py

核心修复文件: 将 `_AITER_NATIVE_OCP_MX_SCHEMES` 从仅包含 `w_mxfp4` 扩展到包含 `w_mxfp4_a_mxfp4`, 使 W4A4 模型在支持 MXFP4 的硬件上正确使用原生 MOE 内核而非仿真。

```
# 文件 : vllm/model_executor/layers/quantization/quark/quark_moe.py
# 上下文 : QuarkOCP_MX_MoEMethod.__init__ 中确定是否使用仿真

# TODO: Remove once all OCP MX schemes use the kernel abstraction
# 此元组明确列出 AITER CK MOE 原生支持哪些 OCP MX 方案。
# 只有方案在此列表中, 且平台支持 MX 且 AITER 启用时, 才会使用原生内核。
# 否则回退到 EMULATION 后端 (软件模拟, 性能较低)。
_AITER_NATIVE_OCP_MX_SCHEMES = ("w_mxfp4", "w_mxfp4_a_mxfp4")

# emulate 判断: 若平台不支持 MX、或方案不在原生列表中、
# 且 mxfp4_backend 为 NONE 或 AITER 未启用, 则标记为仿真。
self.emulate = (
    not current_platform.supports_mx()
    or self.ocp_mx_scheme not in _AITER_NATIVE_OCP_MX_SCHEMES
) and (
    self.mxfp4_backend is Mxfp4MoeBackend.NONE or not self.use_rocm_aiter_moe
)
```

vllm/model\_executor/layers/fused\_moe/rocm\_aiter\_fused\_moe.py

辅助修改: 更新注释以清晰说明 MXFP4 方案 (W4A4、W4A16) 使用 BLOCK\_1X32, 且 MXFP6/MXFP8 当前不支持 AITER 改用仿真。

```
# 文件 : vllm/model_executor/layers/fused_moe/rocm_aiter_fused_moe.py
# 上下文 : rocm_aiter_fused_experts 函数中的量化方法选择
else:
    quant_method = QuantMethod.NO.value
    # mxfp4 i.e. w4a4, w4a16 uses BLOCK_1X32
    # mxfp6 and mxfp8 are unsupported in AITER currently and use emulation instead
    # 此处明确: 只有 w4a4 和 w4a16 走 AITER 原生内核, 其他 MX 方案回退仿真。
    if quant_config.use_mxfp4_w4a4 or quant_config.use_mxfp4_w4a16:
        quant_method = QuantMethod.BLOCK_1X32.value
    # ... 后续处理 w8a8 等其他量化方案
```

## 评论区精华

Review 中 BowenBao 提出了两个关键问题:

1. W4A8 是否真正受支持: BowenBao 询问 `use_mxfp4_w4a8` 是否确实被 ROCm AITER fused MOE 支持, 指出它走的是 `aiter triton fmoe` 路径而非此处逻辑。开发者随后在第二版提交中移除了 W4A8 的支持。

2. w\_mxfp4\_a\_fp8 不应列入: BowenBao 指出 w\_mxfp4\_a\_fp8 也未被 CK MOE 支持, 不应放入该列表, 因为 aiter triton moe 在 QuarkOCP\_MX\_MoEMethod\_OSS 中单独处理。最终交付版本未包含该方案。

结论: 经讨论, 仅保留 W4A4 的修复, 移除了 W4A8 的添加。

- W4A8 是否真正受 AITER CK MOE 支持 (question): 开发者确认 W4A8 不受支持, 在第二版提交中移除了相关代码。
- w\_mxfp4\_a\_fp8 是否应列入原生列表 (question): 最终交付版本未包含 w\_mxfp4\_a\_fp8, 仅修复 w\_mxfp4\_a\_mxfp4。

## 风险与影响

### • 风险:

1. 回归风险: 修改了 Quark MOE 路由逻辑, 可能导致其他 OCP MX 方案 (如 MXFP6/MXFP8) 意外进入不支持的代码路径。但当前方案列表仅包含 w\_mxfp4 和 w\_mxfp4\_a\_mxfp4, 其余方案仍走仿真路径, 风险可控。
2. 硬件兼容性: 变更依赖 current\_platform.supports\_mx() 和 rocm\_aiter\_ops.is\_fused\_moe\_enabled(), 若这些函数返回错误结果可能导致不正确的路由。
3. 缺少测试覆盖: 本次变更未附带测试文件, 可能遗漏对 W4A4 native 路径的 CI 验证。

### • 影响:

1. 用户影响: 使用夸克量化 W4A4 模型的 ROCm 用户将获得性能提升 (从仿真切换为原生内核)。
2. 系统影响: 仅影响 ROCm 平台上的 MXFP4 量化 MOE 路径, 不涉及其他平台或量化方案。
3. 团队影响: 变更极小 (+4/-7), 无需额外维护负担。 - 风险标记: 核心路径变更, 缺少测试覆盖

## 关联脉络

- PR #39801 [假设标题] 导致 W4A4 错误使用仿真后端的 PR: 当前 PR 修复了 #39801 引入的回归问题 (将 W4A4 模型错误路由到仿真后端)。
- PR #41102 [假设标题] 解决量化 CI 问题的 PR: Issue 评论中 AndreasKaratzas 提到量化问题已在 #41102 中处理, 与当前修复相关联。