

PR #24799 完整报告

sgl-project/sglang

[AMD] Fix DeepSeek import cascade by supporting both pre- and post-#2958 `aiter`fused_qk_rmsnorm` APIs`

合并时间: 2026-05-11 14:41

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/24799>

执行摘要

- 一句话: 兼容新旧 `aiter fused_qk_rmsnorm` API 修复 AMD DeepSeek 崩溃
- 推荐动作: 值得精读。此 PR 展示了如何通过特性探测优雅地处理上游库 API 变更, 避免破坏性升级。设计模式值得借鉴: 保持调用侧接口不变, 使用适配器模式兼容新旧 API。特别推荐给从事硬件后端集成或内核库依赖管理的工程师。

功能与动机

ROCm/aiter#2958 将 `aiter.ops.fused_qk_norm_rope_cache_quant` 中的公共 `fused_qk_rmsnorm` 重命名为私有 `_fused_qk_rmsnorm`, 并引入新的统一入口点 (in-place, kwargs-only, no-return 签名)。由于 `forward_mla.py` 在模块加载时导入此符号, 重命名导致 `ImportError`, 级联影响 28 个 SGLang 模型模块 (`deepseek_v2`, `deepseek_nextn` 等), 最终导致 `DeepseekV3ForCausalLM` 被认为不支持, 回退到 Transformers 后端, 并在 `modeling_deepseek.py` 中崩溃。该问题导致 MI35x 上所有 8 个 TP rank 崩溃。

实现拆解

本 PR 仅修改一个文件 `forward_mla.py`, 通过特性探测实现双向兼容:

步骤 1: 移除直接导入, 替换为 `try/except` 探测

- 文件: `python/sglang/srt/models/deepseek_common/attention_forward_methods/forward_mla.py`
- 原代码直接 `from aiter.ops.fused_qk_norm_rope_cache_quant import fused_qk_rmsnorm as fused_qk_rmsnorm_bf16`
- 改为先尝试导入新 API, 若失败再回退到旧 API。

步骤 2: 实现新 API 的适配包装

- `try` 块中导入 `aiter.ops.enum.QuantType` 和 `aiter.ops.fused_qk_rmsnorm_group_quant.fused_qk_rmsnorm` (新统一接口)。
- 定义 `fused_qk_rmsnorm_bf16(q, q_weight, q_eps, k, k_weight, k_eps)` 包装函数, 内部预分配 `q_out` 和 `k_out` 为 `torch.empty_like`, 然后调用新 API 并将 `quant_type` 设为 `_AiterQuantType.No` (因为旧 API 不执行量化)。

- 包装函数的签名与旧 API 完全一致（6 个位置参数，返回 (q_out, k_out) 元组），保证了调用侧 forward_mla.py 中的 _use_aiter 分支代码无需任何修改。

步骤 3: 保留旧导入作为 fallback

- except ImportError: 块中保留原始导入方式，确保仍在使用旧版 aiter 的环境不受影响。
- 注意 batched_gemm_a8w8... 的导入位于 if _use_aiter: 块的末尾，即在 try/except 之后，确保在所有情况下都能正确导入。

关键文件:

- python/sglang/srt/models/deepseek_common/attention_forward_methods/forward_mla.py (模块 模型核心; 类别 source; 类型 data-contract; 符号 fused_qk_rmsnorm_bf16) : 唯一修改的文件。核心变更: 将硬编码的旧 API 导入替换为 try/except 特性探测, 为新 API 提供适配包装以保持调用侧接口不变。

关键符号: fused_qk_rmsnorm_bf16

关键源码片段

python/sglang/srt/models/deepseek_common/attention_forward_methods/forward_mla.py

唯一修改的文件。核心变更: 将硬编码的旧 API 导入替换为 try/except 特性探测, 为新 API 提供适配包装以保持调用侧接口不变。

```
if _use_aiter:
    # 尝试导入新版 aiter 的 unified fused_qk_rmsnorm 入口
    # 该入口在 ROCm/aiter#2958 中被引入, in-place、仅 kwargs、无返回值
    # 使用 try/except 兼容新旧版本, 无需原子性升级 aiter
    try:
        from aiter.ops.enum import QuantType as _AiterQuantType
        from aiter.ops.fused_qk_rmsnorm_group_quant import (
            fused_qk_rmsnorm as _aiter_fused_qk_rmsnorm_unified,
        )

        # 包装函数: 保持与旧 API 一致的 6 位置参数 + (q, k) 元组返回签名
        # 确保调用侧 (forward_mla.py 中的 _use_aiter 分支) 无需任何修改
        def fused_qk_rmsnorm_bf16(q, q_weight, q_eps, k, k_weight, k_eps):
            # 预分配输出缓冲区 - 旧 API 内部自行分配并返回, 新 API 要求调用方传入
            q_out = torch.empty_like(q)
            k_out = torch.empty_like(k)
            # 调用新 unified API, 传入所有参数, 量化类型设为 No (因为旧 API 不量化)
            _aiter_fused_qk_rmsnorm_unified(
                q_out_quantized=q_out,
                k_out=k_out,
                q=q,
                q_weight=q_weight,
                q_epsilon=q_eps,
                k=k,
                k_weight=k_weight,
```

```

        k_epsilon=k_eps,
        quant_type=_AiterQuantType.No,
    )
    return q_out, k_out

except ImportError:
    # 若新 API 不存在, 回退到旧版 aiter 的导入方式 (pre-#2958)
    from aiter.ops.fused_qk_norm_rope_cache_quant import (
        fused_qk_rmsnorm as fused_qk_rmsnorm_bf16,
    )

# 这部分导入不受 API 变更影响, 在 try/except 之后确保所有路径都能执行
from aiter.ops.triton.batched_gemm_a8w8_a_per_token_group_prequant_w_per_batched_
tensor_quant import (
    batched_gemm_a8w8_a_per_token_group_prequant_w_per_batched_tensor_quant,
)

```

评论区精华

1. 特性探测 vs 统一分支策略: HaiShaw 建议使用 amd/aiter-ci 分支进行 CI 集成, 避免各种 API 处理。但 1am9trash 指出本 PR 的 try/except 机制可同时支持新旧 API, 无需等待 aiter 版本升级即可合并, 确保 aiter 集成 CI 和 sglang CI 都能通过。
 2. 没有独立测试文件: review 中未要求添加测试, 但 PR 描述中建议了验证步骤 (模块导入冒烟测试、端到端回归、数值一致性检查)。
- 兼容策略选择 (design): 采用 try/except 双向兼容策略被接受并合入; HaiShaw 建议单独使用 amd/aiter-ci 分支进行 aiter CI 集成。

风险与影响

- 风险:
 1. 兼容性风险: 对新 API 的包装依赖 QuantType.No 枚举值, 若后续 aiter 修改枚举名或删除该值可能导致导入失败。但 fallback 机制能保证回退到旧 API, 风险可控。
 2. 性能风险: 新路径额外引入两次 torch.empty_like 调用, 但正如 PR 描述中分析的, 这是 Sub-microsecond 级开销, 在逐层逐 step 的 MLA 前向中可忽略。旧路径完全无包装。
 3. 缺少测试覆盖: 当前没有针对此适配的单元测试, 依赖手动验证。后续增加测试可防止回归。
- 影响:
 1. 用户端: 修复了 AMD GPU (MI35x) 上 DeepSeek-R1-MXFP4 系列模型的崩溃问题, 使得这些模型可以正常使用。对 NVIDIA GPU 无任何影响。
 2. 系统端: 模块加载时的 ImportError 不再级联吞没, ModelRegistry 能正确注册 DeepseekV3ForCausalLM 等架构, 避免回退到 Transformers 后端。
 3. 团队端: 减少了 SGLang 对 aiter 版本原子升级的依赖, 降低了双仓库联合维护的协调成本。 - 风险标记: 核心路径变更, 缺少测试覆盖

关联脉络

- PR #2958 Refactor to add TypeBasedDispatcher to simplify dispatching: aiter 端的同等 issue, 正是此 PR 修复的上游 API 变更。