

# PR #23031 完整报告

sgl-project/sglang

Revert "feat: Support MXFP4 quantized dense models on AMD CDNA2/CDNA3 GPUs (#19143)"

合并时间: 2026-04-17 12:53

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/23031>

## PR #23031 分析报告

### 执行摘要

PR #23031 回退了 PR #19143 引入的 AMD CDNA2/CDNA3 GPU 上 MXFP4 量化模型支持，以修复因 petit\_kernel 0.0.3 版本缺少 CP310 wheel 导致的 AMD CI 安装阶段失败。此次变更通过删除相关配置文件和调整依赖，确保了 CI 构建稳定性，但暂时移除了 MXFP4 功能，影响特定用户群体。

### 功能与动机

原始 PR #19143 添加了对 MXFP4 量化模型的支持，但新版本 petit\_kernel 0.0.3 只提供 CP312 wheels，而 AMD CI 环境需要 CP310 wheels，导致依赖安装阶段崩溃。根据 PR body 中的表述：“breaks AMD CI in install dependency stage”，revert 决策旨在快速恢复 CI 通过性，避免持续构建失败影响开发流程。

### 实现拆解

本次变更主要通过文件删除和重构完成，具体步骤如下：

1. 删除 MXFP4 相关文件：移除 python/sglang/srt/layers/quantization/petit\_mxfp4.py 和 python/sglang/srt/layers/quantization/petit\_nvfp4.py，这两个文件定义了 PetitMxfp4Config 和 PetitNvFp4Config 类，彻底撤销 MXFP4 量化支持。
2. 重构 petit.py 文件：将 python/sglang/srt/layers/quantization/petit.py 从简单的导入 shim 重写为直接包含 PetitNvFp4Config 实现，但移除了 MXFP4 相关代码，仅保留 NVFP4 逻辑。关键变更示例如下：

```
class PetitNvFp4Config(QuantizationConfig):
    """Config class for Petit FP4."""

    def __init__(
        self,
        is_checkpoint_nvfp4_serialized: bool = False,
        kv_cache_quant_algo: str = None,
        group_size: int = None,
        exclude_modules: List[str] = None,
    ) -> None:
        self.is_checkpoint_nvfp4_serialized = is_checkpoint_nvfp4_serialized
        if is_checkpoint_nvfp4_serialized:
```

```

        logger.warning(
            "Detected nvfp4 checkpoint. Please note that the "
            "format is experimental and subject to change."
        ) # 警告用户检查点格式可能变动
self.group_size = group_size
self.kv_cache_quant_algo = kv_cache_quant_algo
self.exclude_modules = exclude_modules

@classmethod
def get_name(cls) -> str:
    return "petit_nvfp4" # 返回配置标识符, 用于量化方法选择

@classmethod
def from_config(cls, config: Dict[str, Any]) -> "PetitNvFp4Config":
    quant_config = cls.get_from_keys(config, ["quantization"])
    quant_method = quant_config["quant_algo"]
    group_size = quant_config.get("group_size", None)
    verify_petit_nvfp4_supported(quant_method, group_size) # 调用工具函数验证支持性
    is_checkpoint_nvfp4_serialized = "NVFP4" in quant_method
    kv_cache_quant_algo = quant_config["kv_cache_quant_algo"]
    if not kv_cache_quant_algo:
        kv_cache_quant_algo = "auto" # 设置默认值以防配置缺失
    exclude_modules = quant_config.get("exclude_modules", None)
    if not (group_size and kv_cache_quant_algo and (exclude_modules is not None)):
        raise ValueError(
            "NVFP4 quantization requires group size and "
            "kv_cache_quant_algo specified in hf_quant_config.json"
        ) # 严格校验必要参数, 避免运行时错误
    return cls(
        is_checkpoint_nvfp4_serialized,
        kv_cache_quant_algo,
        group_size,
        exclude_modules,
    )

```

1. 更新工具函数: 修改 `python/sglang/srt/layers/quantization/petit_utils.py`, 删除 `_check_petit_mxfp4_supported` 等 MXP4 辅助函数, 简化错误处理, 确保未安装 `petit_kernel` 时能优雅抛出错误。
2. 调整导入和配置: 修改 `python/sglang/srt/layers/quantization/__init__.py`、`python/sglang/srt/configs/model_config.py` 和 `python/sglang/srt/server_args.py`, 移除对 `PetitMxfp4Config` 的引用, 并更新 `quantization` 配置映射, 防止系统加载已删除的功能。
3. 依赖配置更新: 修改 `python/pyproject_other.toml`, 调整 `petit_kernel` 版本或相关配置, 以匹配回退后的依赖要求, 确保 CI 安装阶段兼容性。

### `python/sglang/srt/layers/quantization/petit.py`

重构文件内容, 从导入 `shim` 改为直接实现 `PetitNvFp4Config`, 但移除 MXP4 相关代码, 简化依赖并集中 NVFP4 逻辑。

```

class PetitNvFp4Config(QuantizationConfig):
    """Config class for Petit FP4."""

    def __init__(
        self,
        is_checkpoint_nvfp4_serialized: bool = False,
        kv_cache_quant_algo: str = None,
        group_size: int = None,
        exclude_modules: List[str] = None,
    ) -> None:
        self.is_checkpoint_nvfp4_serialized = is_checkpoint_nvfp4_serialized
        if is_checkpoint_nvfp4_serialized:
            logger.warning(
                "Detected nvfp4 checkpoint. Please note that the "
                "format is experimental and subject to change."
            )
        self.group_size = group_size
        self.kv_cache_quant_algo = kv_cache_quant_algo
        self.exclude_modules = exclude_modules

    @classmethod
    def get_name(cls) -> str:
        return "petit_nvfp4" # 返回配置名称, 用于识别量化方法

    @classmethod
    def get_supported_act_dtypes(cls) -> List[torch.dtype]:
        return [torch.bfloat16, torch.half] # 支持BF16和FP16激活数据类型

    @classmethod
    def from_config(cls, config: Dict[str, Any]) -> "PetitNvFp4Config":
        quant_config = cls.get_from_keys(config, ["quantization"]) # 从配置中提取量化部分
        quant_method = quant_config["quant_algo"]
        group_size = quant_config.get("group_size", None)
        verify_petit_nvfp4_supported(quant_method, group_size) # 验证NVFP4支持性
        is_checkpoint_nvfp4_serialized = "NVFP4" in quant_method
        kv_cache_quant_algo = quant_config["kv_cache_quant_algo"]
        if not kv_cache_quant_algo:
            kv_cache_quant_algo = "auto" # 默认自动选择KV缓存量化算法
        exclude_modules = quant_config.get("exclude_modules", None)
        if not (group_size and kv_cache_quant_algo and (exclude_modules is not None)):
            raise ValueError(
                "NVFP4 quantization requires group size and "
                "kv_cache_quant_algo specified in hf_quant_config.json"
            ) # 校验必要配置项
        return cls(
            is_checkpoint_nvfp4_serialized,
            kv_cache_quant_algo,
            group_size,
            exclude_modules,

```

)

## 评论区精华

review 中仅有 HaiShaw 的 approval，无具体评论。这表明团队基于 CI 失败事实快速决策 revert，未引发技术争议或深入讨论，反映了对构建稳定性的优先考虑。

## 风险与影响

技术风险：

- 功能移除风险：MXFP4 量化支持被彻底删除，依赖此功能的 AMD GPU 用户无法使用相关模型，可能导致生产中断。
- 依赖回退风险：回退到 petit\_kernel 0.0.2 可能引入旧版本的安全漏洞或性能退化，需评估长期影响。
- 兼容性问题：代码中残留的 NVFP4 支持可能未充分测试，在特定配置下引发运行时错误。
- 配置断裂风险：model\_config.py 中的 quantization 验证列表移除 petit\_mxfp4，若用户配置仍引用该键，可能导致解析失败。

影响评估：

- 用户影响：AMD CDNA2/CDNA3 GPU 用户暂时无法使用 MXFP4 量化模型，需等待后续修复或替代方案。
- 系统影响：CI 安装阶段稳定性得到修复，减少构建失败，提升开发效率。
- 团队影响：凸显依赖版本管理的重要性，推动团队加强 CI 环境锁定和测试覆盖；量化模块开发者需重新评估依赖兼容性策略。

## 关联脉络

直接关联 PR #19143，该 PR 引入了 MXFP4 量化支持，本 revert 揭示了一个典型问题：依赖版本不匹配如何导致功能回退。结合近期历史 PR，如 #22888（修复 NPU 环境检查）和 #22994（统一环境变量读取），可见团队在持续优化基础设施和跨平台兼容性。此次变更提醒开发者在引入新依赖时需严格验证多版本支持，以避免类似 CI 中断。