

# PR #40735 完整报告

vllm-project/vllm

[MoE Refactor] Introduce RoutedExperts alias for FusedMoE and don't store SharedExperts in MK

合并时间: 2026-05-13 01:37

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/40735>

## 执行摘要

- 一句话: 引入 RoutedExperts 别名并解耦 SharedExperts 存储
- 推荐动作: 值得精读, 尤其关注 modular\_kernel.py 中如何将 SharedExperts 从内部状态改为方法参数, 这是典型的解耦模式。对于 MoE 重构系列, 理解此 PR 有助于后续理解 PR#38590 的更大变更。

## 功能与动机

根据 PR body, 这是为了 #38590 做的准备工作, 目的是引入 RoutedExperts (暂时作为 FusedMoE 的别名) 并传递 SharedExperts 作为参数而不是存储在 MK 中, 以便将来 RoutedExperts 能成为一个拥有量化方法和层数据的独立类。

## 实现拆解

1. 定义 RoutedExperts 别名并统一导出: 在 vllm/model\_executor/layers/fused\_moe/\_\_init\_\_.py 中定义 RoutedExperts = FusedMoE, 并额外导出 SharedExperts、FusedMoEQuantConfig 等常用类型, 使各量化方法可以从统一入口导入, 减少散落依赖。
2. 替换量化方法中的类型注解: 在 gptq\_marlin.py、modelopt.py、fp8.py、quark\_moe.py、mxfp4.py、awq\_marlin.py、humming.py、gguf.py、moe\_wna16.py 等文件中, 将所有 isinstance(layer, FusedMoE) 替换为 isinstance(layer, RoutedExperts), 对应的方法签名如 create\_weights、process\_weights\_after\_loading、\_setup\_kernel 等的 layer 参数类型也从 torch.nn.Module 或 FusedMoE 改为 RoutedExperts。
3. 修改基类方法签名: 在 fused\_moe\_method\_base.py 和 unquantized\_fused\_moe\_method.py 中, 为 apply、apply\_monolithic 等方法增加 shared\_experts: SharedExperts | None 和 shared\_experts\_input: torch.Tensor | None 参数, 使得 MK 不需要再内部持有 SharedExperts。
4. 重构 MK 核心实现: 在 modular\_kernel.py 中, FusedMoEKernelModularImpl.\_\_init\_\_ 不再接受 shared\_experts 参数; \_maybe\_apply\_shared\_experts、\_finalize 和 apply 等方法改为从外部接收 shared\_experts 实例, 并在调用链中传递。去掉了 self.shared\_experts 属性。
5. 更新测试与辅助函数: 对应修改了 marlin\_utils.py、compressed\_tensors.py 等辅助模块中的类型注解和导入路径; MoE 相关测试文件同步调整了类型引用。

关键文件：

- `vllm/model_executor/layers/fused_moe/__init__.py` (模块 MoE 层; 类别 source; 类型 data-contract; 符号 `RoutedExperts`, `SharedExperts`, `FusedMoEQuantConfig`, `FusedMoEParallelConfig`) : 定义 `RoutedExperts = FusedMoE` 别名, 并统一导出常用类型, 是整个重构的基础。
- `vllm/model_executor/layers/fused_moe/modular_kernel.py` (模块 核心内核; 类别 source; 类型 core-logic; 符号 `FusedMoEKernelModularImpl.init`, `_maybe_apply_shared_experts`, `_finalize`, `apply`) : 核心实现: 将 `SharedExperts` 从 MK 内部移除, 改为通过方法参数传递, 是架构解耦的关键变更。
- `vllm/model_executor/layers/quantization/gptq_marlin.py` (模块 量化器; 类别 source; 类型 data-contract; 符号 `get_moe_quant_method`, `process_weights_after_loading`, `_setup_kernel`, `get_fused_moe_quant_config`) : 代表性量化方法文件, 展示了类型注解替换和接口调整的典型模式。
- `vllm/model_executor/layers/fused_moe/fused_moe_method_base.py` (模块 MoE 层; 类别 source; 类型 data-contract; 符号 `mk_owns_shared_expert`, `mk_can_overlap_shared_experts`, `apply`, `apply_monolithic`) : MoE 量化方法的基类, 新增 `shared_experts` 参数改变了所有具体量化方法的契约。
- `vllm/model_executor/layers/quantization/modelopt.py` (模块 量化器; 类别 source; 类型 data-contract; 符号 `process_weights_after_loading`, `get_fused_moe_quant_config`, `create_weights`) : 另一个关键量化方法, 涉及的改动与 `gptq_marlin` 类似, 且包含突变→断言的变更。

关键符号: `RoutedExperts`, `FusedMoEKernelModularImpl.init`,  
`FusedMoEKernelModularImpl.apply`, `FusedMoEMethodBase.apply`,  
`process_weights_after_loading`, `_setup_kernel`, `get_fused_moe_quant_config`

## 关键源码片段

### `vllm/model_executor/layers/fused_moe/modular_kernel.py`

核心实现: 将 `SharedExperts` 从 MK 内部移除, 改为通过方法参数传递, 是架构解耦的关键变更。

```
# vllm/model_executor/layers/fused_moe/modular_kernel.py ( 关键变更摘要 )
```

```
class FusedMoEKernelModularImpl:
    def __init__(
        self,
        prepare_finalize: FusedMoEPrepareAndFinalizeModular,
        fused_experts: FusedMoEExpertsModular,
        # 不再接受 shared_experts 参数
        inplace: bool = False,
    ):
        self.prepare_finalize = prepare_finalize
        self.fused_experts = fused_experts
        # 不再存储 self.shared_experts
```

```

...

def _maybe_apply_shared_experts(
    self,
    shared_experts: SharedExperts | None, # 外部传入
    shared_experts_input: torch.Tensor | None,
):
    if shared_experts is not None:
        shared_experts.apply(
            shared_experts_input,
            SharedExpertsOrder.MK_INTERNAL_OVERLAPPED,
        )

def apply(
    self,
    ...
    shared_experts: SharedExperts | None = None, # 新增参数
    shared_experts_input: torch.Tensor | None = None,
) -> torch.Tensor:
    ...
    # 将 shared_experts 沿调用链传递
    self._maybe_apply_shared_experts(shared_experts, shared_experts_input)
    ...

```

## 评论区精华

- 占位符残留: gemini-code-assist[bot] 指出 `compressed_tensors.py` 和 `marlin_utils.py` 中存在未完成的注释和杂乱字符 (如 `# XXXXXXXX...` 和 `# ?`) , 要求清理。
- 导入路径错误: yzong-rh 指出 `fused_moe_method_base.py` 中类型检查导入的路径 `vllm.model_executor.layers.runner.shared_experts` 不正确, 应为 `vllm.model_executor.layers.fused_moe` 或 `...runner.shared_experts`; 作者接受并修复。
- `shared_experts` 传递遗漏: yzong-rh 在 `quark_moe.py` 和 `gptq_marlin.py` 中质疑 `shared_experts` 是否传递给 `moe_kernel.apply`, 引发作者补充参数传递。
- 旧导入路径问题: yzong-rh 指出许多地方仍从旧路径导入 `FusedMoEQuantConfig` 等类型; 作者回应将这些常用类型统一放到了 `__init__.py` 中以减少改动范围。
- 突变改为断言: 在 `modelopt.py` 中, `create_weights` 内的属性设置改为断言, 作者明确这是故意的, 因为这些属性应由 `layer` 预先设置。
- 遗留占位符 (other): 作者后续提交修复了这些占位符。
- 导入路径错误 (design): 作者承认并修复。
- `shared_experts` 是否传递到 `moe_kernel.apply (correctness)`: 作者补充了相应的参数传递。
- 旧导入路径大量残留 (design): 作者决定将这些类型添加到 `init.py` 的导出中, 以减少外部导入路径的变化。
- `create_weights` 中 `assign` 改为 `assert` (design): 作者解释这是故意的, 因为这些属性应由 `layer` 预先设置。

## 风险与影响

- 风险:

1. 大规模机械替换风险: 53 个文件中将 FusedMoE 替换为 RoutedExperts, 可能存在遗漏的引用, 尤其是动态类型检查或静态配置文件。但测试和类型提示可发现。
2. 接口签名变更: apply 等方法增加了 shared\_experts 和 shared\_experts\_input 参数, 若存在外部自定义 MoE 方法未同步更新, 将导致调用错误。
3. 旧导入路径残留: 虽然统一了导出, 但旧路径仍可能被第三方插件或尚未更新的内部模块使用, 造成导入失败。
4. 占位符未清理: 若提交前未完全清除 # XXXXXXXX... 等占位符, 将影响后续维护的代码质量。- 影响: 影响所有使用 MoE 层的量化方法 (FP8、GPTQ\_Marlin、ModelOpt、Quark、MXFP4 等) 以及 MoE 核心模块 (fused\_moe)。功能无变化, 但接口签名变更要求调用方传递 shared\_experts。开发者需了解 RoutedExperts 作为 FusedMoE 的新别名。团队应从旧的导入路径迁移到统一入口。- 风险标记: 核心路径变更, 大规模机械替换, 接口签名变更, 旧导入路径残留

## 关联脉络

- PR #38590 待定 (未在库中提供): 此 PR 明确为 #38590 做准备, 是 MoE 重构系列的关键前置步骤。
- PR #41055 [MoE Refactor] EPLB refactoring for FusedMoE: 同一重构系列, 简化 MoE 路由接口。
- PR #41046 [MoE Refactor] Move expert map related code into ExpertMapManager class: 同一重构系列, 抽取专家映射逻辑。
- PR #42334 [MoE Refactor] Move remaining experts classes to experts directory: 同一重构系列, 移动 experts 类, 与当前 PR 有文件交集 (fused\_moe/init.py)。