

# PR #39753 完整报告

vllm-project/vllm

[Model] Use mm\_features for Ernie-4.5 VL M-RoPE

合并时间: 2026-04-14 16:11

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/39753>

## 执行摘要

- 一句话: Ernie-4.5 VL 模型的 M-RoPE 位置计算从基于 token 扫描重构为使用 mm\_features 数据契约。
- 推荐动作: 建议工程师精读此 PR, 关注如何通过 iter\_mm\_grid\_thw 辅助方法将复杂 token 扫描逻辑抽象为清晰的数据迭代器, 以及多模态数据契约 (mm\_features) 在设计中的应用。对于涉及位置编码或模型重构的任务, 这是一个值得借鉴的设计决策示例。

## 功能与动机

根据 PR body, 目的是实现 #32656, 使 Ernie-4.5 VL 模型与其他多模态模型保持一致, 使用新的 mm\_features 驱动 M-RoPE 实现, 避免从 token ID 重建图像和视频区域, 提升代码维护性和一致性。

## 实现拆解

1. 修改主源码文件: 在 vllm/model\_executor/models/ernie45\_vl.py 中, 移除 get\_mrope\_input\_positions 方法的旧逻辑 (基于 input\_tokens 扫描和 itertools.groupby), 添加新的 iter\_mm\_grid\_thw 辅助方法。关键符号: get\_mrope\_input\_positions, iter\_mm\_grid\_thw。变更原因: 简化逻辑, 直接利用 mm\_features 中的偏移和网格数据, 避免复杂 token 解析。影响: 使位置计算更高效且与 vLLM 多模态数据契约对齐。
2. 新增辅助方法: iter\_mm\_grid\_thw 方法从 mm\_features 提取偏移和网格尺寸, 并根据模型配置应用空间和时间卷积尺寸, 返回迭代器。关键符号: iter\_mm\_grid\_thw。变更原因: 封装多模态数据处理, 提供清晰的接口。影响: 提高代码可读性和可测试性。
3. 更新位置计算逻辑: 在 get\_mrope\_input\_positions 中, 使用 iter\_mm\_grid\_thw 迭代处理多模态特征, 计算文本和多模态位置的 ID, 并合并为位置张量。关键符号: np.broadcast\_to, np.indices (但 review 建议使用 PyTorch)。变更原因: 从数据驱动角度重构, 减少状态管理。影响: 确保位置计算准确, 但依赖 NumPy 可能引入风格不一致。
4. 添加测试覆盖: 新增文件 tests/model\_executor/test\_ernie45\_vl\_mrope.py, 包含三个单元测试: test\_get\_mrope\_input\_positions\_text\_only, test\_get\_mrope\_input\_positions\_single\_image, test\_get\_mrope\_input\_positions\_interleaved\_image\_and\_video。关键符号: \_force\_cpu\_default\_device, DummyConfig, make\_model, make\_mm\_feature。变更原因: 验证重构后方法的正确性, 覆盖不同输入场景。影响: 降低回归风险, 提供可重复测试。

5. 配套调整：移除 import itertools，添加 Iterator 导入到 collections.abc，调整导入关系以支持新逻辑。变更原因：清理未使用的依赖，适应新代码结构。影响：最小化代码变更范围，保持模块整洁。

关键文件：

- vllm/model\_executor/models/ernie45\_vl.py（模块 模型执行器；类别 source；类型 core-logic；符号 get\_mrope\_input\_positions, iter\_mm\_grid\_thw）：主源码文件，重构了 M-RoPE 位置计算方法，从基于 token 扫描切换到使用 mm\_features 数据契约。
- tests/model\_executor/test\_ernie45\_vl\_mrope.py（模块 测试模块；类别 test；类型 test-coverage；符号 \_force\_cpu\_default\_device, DummyConfig, make\_model, make\_mm\_feature）：新增的单元测试文件，专门验证 Ernie-4.5 VL 模型 M-RoPE 位置计算的正确性，覆盖多种输入场景。

关键符号：get\_mrope\_input\_positions, iter\_mm\_grid\_thw

## 关键源码片段

### vllm/model\_executor/models/ernie45\_vl.py

主源码文件，重构了 M-RoPE 位置计算方法，从基于 token 扫描切换到使用 mm\_features 数据契约。

```
def iter_mm_grid_thw(
    self, mm_features: list[MultiModalFeatureSpec]
) -> Iterator[tuple[int, int, int, int]]:
    """
    迭代处理多模态特征，返回每个特征的偏移和调整后的网格尺寸。
    用于M-RoPE位置计算，从mm_features直接提取数据，避免token解析。
    """
    # 从模型配置获取卷积尺寸，用于缩放网格
    spatial_conv_size = self.config.spatial_conv_size
    temporal_conv_size = self.config.temporal_conv_size

    # 按偏移量排序特征，确保顺序处理
    for mm_feature in sorted(mm_features, key=lambda f: f.mm_position.offset):
        if mm_feature.data is None:
            # 数据缺失时抛出错误，保证计算可靠性
            raise ValueError("M-RoPE计算需要多模态特征数据")

        offset = mm_feature.mm_position.offset # 特征在输入序列中的起始位置
        if mm_feature.modality == "image":
            # 提取图像网格尺寸 (t, h, w) 并应用空间卷积缩放
            t, h, w = mm_feature.data["image_grid_thw"].data.tolist()
            yield offset, t, h // spatial_conv_size, w // spatial_conv_size
        elif mm_feature.modality == "video":
            # 提取视频网格尺寸并应用时间和空间卷积缩放
            t, h, w = mm_feature.data["video_grid_thw"].data.tolist()
            yield (
                offset,
```

```
        t // temporal_conv_size,
        h // spatial_conv_size,
        w // spatial_conv_size,
    )
```

## tests/model\_executor/test\_ernie45\_vl\_mrope.py

新增的单元测试文件，专门验证 Ernie-4.5 VL 模型 M-RoPE 位置计算的正确性，覆盖多种输入场景。

```
def test_get_mrope_input_positions_single_image():
    """
    测试单图像输入时的M-RoPE位置计算，验证位置矩阵和调整值delta。
    """
    # 创建虚拟模型实例，使用默认配置
    model = make_model(DummyConfig())

    # 构建单图像的多模态特征，模拟图像在偏移1处，网格尺寸为(1,4,4)
    mm_features = [
        make_mm_feature(
            modality="image",
            offset=1, # 图像在输入序列中的起始索引
            length=4, # 特征长度，用于占位
            grid_thw=(1, 4, 4), # 网格尺寸 (t, h, w)
        )
    ]

    # 调用重构后的方法，输入模拟token序列
    positions, delta = model.get_mrope_input_positions(
        input_tokens=[10, 20, 21, 22, 23, 30, 31], # 包含文本和图像token
        mm_features=mm_features,
    )

    # 预期位置矩阵，基于新逻辑计算
    expected = torch.tensor(
        [
            [0, 1, 1, 1, 1, 3, 4], # 第一维位置
            [0, 1, 1, 2, 2, 3, 4], # 第二维位置
            [0, 1, 2, 1, 2, 3, 4], # 第三维位置
        ]
    )

    # 断言验证计算正确性
    assert torch.equal(positions, expected)
    assert delta == -2 # 位置调整值，反映多模态跨度的影响
```

## 评论区精华

review 中，gemini-code-assist[bot] 提出了三个高优先级评论，建议将 `get_mrope_input_positions` 方法中的 NumPy 操作（如 `np.broadcast_to` 和 `np.indices`）替

换为 PyTorch 操作（如 `torch.arange` 和 `torch.meshgrid`），以保持与 vLLM 代码库的一致性并避免潜在的 64 位整数默认问题。DarkLight1337 批准了 PR 并提供了 LM-eval 结果，显示主分支与 PR 分支的性能无明显差异（0.623 vs 0.6197 exact\_match），表明变更未引入功能性回归。结论：代码风格建议未被明确采纳，但 PR 因功能正确而被批准，讨论聚焦于非关键的一致性优化。

- 代码风格一致性建议 (style): 建议未被明确采纳，但 PR 因功能正确和 LM-eval 结果相似而被 DarkLight1337 批准，表明风格问题被视为非关键优化。

## 风险与影响

- 风险：技术风险较低：
- 回归风险：变更涉及核心位置计算方法，但新增单元测试覆盖了文本、单图像、图像视频交错场景，验证了正确性，降低了错误计算风险。
- 性能影响：从基于 token 扫描转向数据驱动，可能微调计算开销，但 LM-eval 结果显示无显著性能下降，影响可忽略。
- 兼容性问题：新逻辑依赖 `mm_features` 数据结构，如果传入数据不符合预期（如缺失 `mm_position.offset` 或 `grid` 数据），可能引发 `ValueError` 或计算错误，但测试模拟了标准场景。
- 代码风格不一致：使用 NumPy 而非 PyTorch 可能在未来维护中引入混淆，但非功能性风险。
- 影响：影响范围中等：
- 用户影响：对最终用户透明，不改变 API 或使用方式，但确保了 Ernie-4.5 VL 模型在多模态推理中的位置编码准确性。
- 系统影响：仅限于 Ernie-4.5 VL 模型的 M-RoPE 计算模块，不影响其他模型或系统组件，提升了与 vLLM 多模态生态的一致性。
- 团队影响：为工程师提供了从传统 token 解析到现代数据驱动设计的案例，有助于统一多模态模型实现模式，简化后续维护。
- 风险标记：数据契约依赖，核心路径变更

## 关联脉络

- PR #39217 [Mistral Grammar] Fix tool and reasoning parsing: 同属多模态模型中的解析逻辑改进，共享使用数据契约（如 `mm_features`）处理多模态输入的模式，反映 vLLM 在多模态集成上的持续演进。