

PR #39688 完整报告

vllm-project/vllm

[fix][MOE] Fix MOE experts `intermediate_size` dimension not being narrowed before weight loading

合并时间: 2026-04-14 17:35

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/39688>

执行摘要

该 PR 修复了 MOE 专家权重加载中的一个关键 bug: 当 `hidden_size` 和 `intermediate_size` 维度同时存在填充时, `_narrow_expert_data_for_padding` 方法仅裁剪了 `hidden_size` 维度, 导致后续张量复制出现形状不匹配错误。修复方案扩展了该方法以支持双维度裁剪, 并添加了测试用例。这是一个重要的 bugfix, 确保了 MOE 模型权重加载的可靠性。

功能与动机

问题根源: PR #37010 引入了 `_narrow_expert_data_for_padding` 方法来处理填充的 `hidden_size` 维度, 但该方法仅裁剪 `hidden_size` 维度, 而旧逻辑 (通过切片 `expert_data[:loaded_weight.shape[0], :loaded_weight.shape[1]]`) 同时处理了 `hidden_size` 和 `intermediate_size` 两个维度。这导致在最后一个 TP 分片存在填充时, `intermediate_size` 维度未正确裁剪, 引发 `RuntimeError: The size of tensor a (768) must match the size of tensor b (704) at non-singleton dimension 1`。

触发场景: 在 MXFP4 GEMM 等场景中, 当 backend (如 DeepEP) 对 `hidden_size` 和 `intermediate_size` 进行向上取整填充时, 权重参数会大于 checkpoint 中的原始权重, 需要裁剪后才能正确加载。

实现拆解

核心修改

1. `vllm/model_executor/layers/fused_moe/layer.py`:

- 修改 `_narrow_expert_data_for_padding` 方法:

```
python def _narrow_expert_data_for_padding( expert_data: torch.Tensor, loaded_weight: torch.Tensor, hidden_dim: int, shard_dim: int | None = None, # 从-1改为None ) -> torch.Tensor: dims = (hidden_dim,) if shard_dim is None else (hidden_dim, shard_dim) if loaded_weight.ndim > 0: for dim in dims: if (0 <= dim < expert_data.ndim and dim < loaded_weight.ndim and expert_data.shape[dim] > loaded_weight.shape[dim]): expert_data = expert_data.narrow(dim, 0, loaded_weight.shape[dim]) return expert_data
```
- 在 `_load_per_channel_weight_scale`、`_load_w13`、`_load_w2` 三个权重加载函数中调用该方法时传递 `shard_dim` 参数。

2. tests/kernels/moe/test_moe_weight_loading_padded.py:

- 新增 test_narrow_shard_dim 测试用例，模拟 w2 权重加载时两个维度同时填充的场景：
 - 设置填充的 hidden_size=3072（原始 2688）、填充的 intermediate_size=1024（原始 896）
 - 验证裁剪后张量数据正确复制到填充张量的左上角区域

调用链更新

函数	修改点
<code>_load_per_channel_weight_scale</code>	传递 <code>shard_dim</code> 参数
<code>_load_w13</code>	传递 <code>shard_dim</code> 参数
<code>_load_w2</code>	传递 <code>shard_dim</code> 参数

评论区精华

设计权衡：关于 `shard_dim` 参数默认值的讨论：

```
tomeras91: "in PyTorch, -1 is a valid dimension index (meaning 'last dimension'), so using it as a sentinel for 'skip' is a bit ambiguous. If someone ever passes shard_dim=-1 intending the last dim, the narrowing would silently be skipped. Consider using Optional[int] = None instead."
```

结论：作者采纳建议，将默认值从 `-1` 改为 `None`，避免了 API 歧义，体现了良好的接口设计实践。

风险与影响

技术风险

1. 回归风险：修复了 PR #37010 引入的回归，但需确保新逻辑在所有 MOE 权重加载路径（包括不同量化格式、TP 分片配置）下均正确工作。
2. 兼容性：`_narrow_expert_data_for_padding` 接口变更（新增 `shard_dim` 参数）保持了向后兼容（默认 `None`），不影响现有调用。
3. 测试覆盖：新增测试仅覆盖了 w2 场景的双维度填充，未覆盖 w1/w3 等其他路径，可能存在未发现的边缘情况。

影响范围

- 用户影响：修复了 MOE 模型（如 Qwen2-MoE）在特定配置下权重加载失败的问题，提升模型部署成功率。
- 系统影响：确保 MOE 专家权重正确加载，避免静默数据损坏或运行时崩溃。
- 团队影响：为权重加载逻辑中的多维度处理提供了参考模式，后续类似修复可借鉴。

关联脉络

历史 PR 关联

- PR #37010: 直接关联, 当前 PR 修复了该 PR 引入的回归问题。
- PR #39717、PR #39705: 同属量化相关 bugfix, 都涉及维度不匹配导致的静默数据损坏问题, 反映了团队对量化场景下数据完整性的持续关注。

演进趋势

从近期历史 PR 看, vLLM 团队在以下方向持续投入:

1. 量化优化: 频繁修复量化内核中的维度匹配和数据损坏问题 (如 #39717、#39705)。
2. MOE 支持: 当前 PR 是 MOE 专家权重加载的重要修复, 配合 #37010 等 PR, 显示团队正在完善 MOE 模型的端到端支持。
3. API 设计: review 中关于默认值设计的讨论体现了团队对接口清晰性和安全性的重视。

该 PR 是 MOE 模型支持链条中的关键一环, 确保了权重加载的可靠性, 为后续 MOE 性能优化和功能扩展奠定了基础。