

PR #43100 完整报告

vllm-project/vllm

[BugFix] Fix Humming MoE deploy error

合并时间: 2026-06-03 00:32

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/43100>

执行摘要

- 一句话: 修复 Humming MoE 部署时 quant config 与 schema 初始化遗漏
- 推荐动作: 建议合并。该 PR 修复了明确的部署阻塞 bug, 改动量小且经过本地验证。建议后续为该路径补充测试, 防止回归。

功能与动机

部署 ISTA-DASLab/Qwen3.6-35B-A3B-2Bit-GSQ 时遭遇断言错误 `assert self.moe_quant_config is not None`, 以及 `layer.input_schemas` 为空。分析发现 `moe_quant_config` 从未被初始化, 而 `input_schemas` 和 `weight_schemas` 仅在非 Humming 格式的检查点分支中被赋值, 导致 Humming 格式检查点下 schema 未初始化。

实现拆解

实现拆解:

1. Schema 初始化修复: 在 `humming.py` 的 `process_weights_after_loading` 方法中, 将 `layer.weight_schemas[sublayer_name]` 和 `layer.input_schemas[sublayer_name]` 的赋值语句从 `if force_requant` 条件块内移出 (减少缩进一级), 确保无论是否进行强制重量化, `schema` 都会被正确设置。
2. quant config 初始化修复: 在 MoE 专家创建之前, 添加 `layer.ensure_moe_quant_config_init()` 调用, 该方法会调用 `get_fused_moe_quant_config` 并将结果赋值给 `self.moe_quant_config`, 从而消除后续的断言错误。
3. 测试: 未添加自动化测试, 但通过手动运行指定模型验证了修复的正确性。

关键文件:

- `vllm/model_executor/layers/quantization/humming.py` (模块 量化模块; 类别 `source`; 类型 `data-contract`): 修复了导致 Humming MoE 部署失败的两个关键 bug: `moe_quant_config` 未初始化导致断言失败, 以及 `schema` 分配未覆盖所有分支导致 `input_schemas` 为空。

关键符号: `process_weights_after_loading`

关键源码片段

vllm/model_executor/layers/quantization/humming.py

修复了导致 Humming MoE 部署失败的两个关键 bug: `moe_quant_config` 未初始化导致断言失败, 以及 `schema` 分配未覆盖所有分支导致 `input_schemas` 为空。

```
# -*- 文件 : vllm/model_executor/layers/quantization/humming.py -*-

# 旧版: schema 赋值缩进在 force_requant 块内, 导致 Humming 格式 checkpoint
# 下不会被赋值
# 新版: 始终赋值 schema
layer.weight_schemas[sublayer_name] = weight_schema
layer.input_schemas[sublayer_name] = input_schema

# force requant (origin quant setting -> fp16/bf16 -> new_quant setting)
assert isinstance(weight_schema, HummingWeightSchema)
force_requant = self.force_weight_schema is not None
if force_requant and weight_schema != self.force_weight_schema:
    # ... 重量化逻辑 (不变) ...

# prepare layer config from humming kernel
HummingMethod.prepare_layer_meta(
    layer=layer,
    shape_n=configs["shape_n"],
    shape_k=configs["shape_k"],
    pad_n_to_multiple=256,
    pad_k_to_multiple=128,
    input_schema=input_schema,
    weight_schema=weight_schema,
    has_bias=self.moe.has_bias,
    num_experts=layer.num_experts,
    torch_dtype=layer.param_dtype,
    sublayer_name=sublayer_name,
)

# preprocess weight for inference
HummingMethod.transform_humming_layer(layer, sublayer_name=sublayer_name)

# use moe modular
experts: HummingIndexedExperts | HummingGroupedExperts
# 新增: 确保 moe_quant_config 已初始化
layer.ensure_moe_quant_config_init()
assert self.moe_quant_config is not None
if get_humming_moe_gemm_type() == "indexed":
    experts = HummingIndexedExperts(layer, self.moe, self.moe_quant_config)
else:
    experts = HummingGroupedExperts(layer, self.moe, self.moe_quant_config)
self.experts = experts
```

评论区精华

- gemini-code-assist指出初始修改中 `self.get_fused_moe_quant_config(layer)` 的返回值未被赋值，导致 `self.moe_quant_config` 仍为 `None`，并建议直接赋值。
- jinzhen-lin建议调用 `layer.ensure_moe_quant_config_init()` 方法，这是一个更简洁的封装。PR 作者接受了该建议并修改了代码。
- jinzhen-lin最终批准了 PR。
- `get_fused_moe_quant_config` 返回值未赋值 (correctness): PR 作者随后采纳了 jinzhen-lin 的建议，改用 `layer.ensure_moe_quant_config_init()`，内部包含了赋值。
- 建议使用 `layer.ensure_moe_quant_config_init()` (correctness): PR 作者提交了包含该改动的版本。

风险与影响

- 风险：风险较低。修改范围仅限于 Humming 量化层中 MoE 权重加载流程，且改动逻辑清晰（减少缩进、调用已有方法）。但缺少自动化测试覆盖，如果未来代码重构或 `ensure_moe_quant_config_init` 实现发生变化，可能再次引入类似问题。
- 影响：修复后，使用 Humming packed quantization 的 MoE 模型（如 Qwen3.6-35B-A3B-2Bit-GSQ）可以正常部署和推理。该修复不影响其他量化方案或非 MoE 模型。
- 风险标记：缺少测试覆盖

关联脉络

- 暂无明显关联 PR