

# PR #38178 完整报告

vllm-project/vllm

[CI] Fix conch kernel crash on 3D input by reshaping to 2D before GEMM

合并时间: 2026-03-27 00:46

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/38178>

## 执行摘要

此 PR 修复了 conch 内核在使用 transformers backend 处理 3D 输入时的崩溃问题，通过将输入重塑为 2D 以适应 GEMM 操作，影响范围仅限于该内核，已通过测试验证，确保量化模型运行稳定。

## 功能与动机

当通过 transformers backend 运行量化模型（如 AWQ）时，conch 内核因接收到 3D 张量（batch, seq\_len, hidden\_dim）而崩溃，错误提示 "ValueError: too many values to unpack (expected 2)"。这是因为 `mixed_precision_gemm` 函数期望 2D 输入，但 backend 传递了 3D 张量，需要修复以支持该场景。

## 实现拆解

在文件 `vllm/model_executor/kernels/linear/mixed_precision/conch.py` 的 `apply_weights` 函数中，进行了以下更改：

```
x_2d = x.reshape(-1, x.shape[-1])
out_shape = x.shape[:-1] + (self.config.partition_weight_shape[1],)
output = mixed_precision_gemm(
    x=x_2d,
    ...
)
return output.reshape(out_shape)
```

重塑输入为 2D 以匹配 GEMM 要求，并在计算后恢复原始形状，参考了现有 `machete` 内核的实现模式。

## 评论区精华

review 中无实质性讨论；只有自动评论机器人 (`gemini-code-assist[bot]`) 指出更改，以及维护者 `Isotr0py` 的批准，无争议或技术交锋。

## 风险与影响

风险：重塑操作简单，但需确保输入形状在所有边界情况下正确处理，例如当张量维度不标准时可能引入错误或性能开销；测试已覆盖相关场景，降低了风险。

影响：修复了特定崩溃，使 transformers backend 支持量化模型运行，对系统性能影响微小（仅增加少量重塑操作），提升了 ROCm 平台上的兼容性。

## 关联脉络

与历史 PR 如 #38161（ROCm 平台量化测试修复）和 #38083（量化精度修复）相关，均涉及 ROCm 和量化功能的稳定性改进，表明团队在持续优化多平台支持和内核正确性，形成跨 PR 的量化内核演进趋势。