

PR #38450 完整报告

vllm-project/vllm

[ROCm][CI] Fix cross-attention dispatch for encoder-decoder models

合并时间: 2026-03-29 13:08

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/38450>

执行摘要

本 PR 修复了 ROCm 平台下编码器 - 解码器模型 (如 Whisper、BART) 交叉注意力层的后端调度错误, 通过从 ROCM_ATT_N 和 ROCM_AITER_FA 后端移除 ENCODER_DECODER 支持, 确保调度选择正确工作的后端 (如 ROCM_AITER_UNIFIED_ATT_N 或 TRITON_ATT_N), 从而解决波束搜索结果不匹配问题, 提升了模型可靠性和测试覆盖。

功能与动机

本 PR 旨在解决测试 `test_whisper_beam_search_single_beam` 在 ROCm 上的失败问题。PR body 中解释: 在编码器 - 解码器模型中, 交叉注意力层 (AttentionType.ENCODER_DECODER) 在 ROCM_ATT_N 和 ROCM_AITER_FA 后端上计算错误, 当 `max_query_len > 1` (如波束搜索) 时导致输出不匹配, 而贪婪解码因 `max_query_len=1` 跳过错误路径不受影响。Issue 评论提到此修复由 PR #38321 的讨论引发, 凸显了跨后端调度的重要性。

实现拆解

关键改动按模块拆解如下:

- 注意力后端模块: 在 `rocm_attn.py` 和 `rocm_aiter_fa.py` 中, 修改 `supports_attn_type` 方法, 移除 ENCODER_DECODER 支持, 并添加详细注释解释技术原因 (如 ROCM_ATT_N 的 prefill kernel 错误处理缓存语义, ROCM_AITER_FA 的序列边界和 causal mask 问题)。
- 平台调度模块: 在 `rocm.py` 中, 改进 `get_attn_backend_cls` 函数的日志记录, 当后端因不兼容被跳过时, 添加日志显示跳过原因和替代后端, 提升调试能力。
- 测试模块: 在 `test_transcription_validation_whisper.py` 中, 扩展测试参数化以覆盖多个 ROCm 后端 (如默认、TRITON_ATT_N、ROCM_AITER_UNIFIED_ATT_N), 并禁用 prefix caching 以减少非确定性, 确保修复验证。
- 文档与工具模块: 更新 `attention_backends.md` 文档表格, 并修改 `generate_attention_backend_docs.py` 生成器, 使其自动计算支持类型, 防止未来不一致。

评论区精华

review 讨论中只有一个关键线程: `gemini-code-assist[bot]` 在文档文件上指出 ROCM_ATT_N 条目仍显示支持所有类型, 与代码变更矛盾。作者 `AndreasKaratzas` 回应:

"This was a statically set value in the auto generator. I modified this such that it collects the number of possible attention types." 这导致修改文档生成器，确保支持类型从代码自动派生，解决了文档一致性问题。讨论无争议，结论清晰。

风险与影响

技术风险：1) 文档不一致风险已通过自动化生成器缓解；2) 移除 ENCODER_DECODER 支持可能导致用户配置意外回退到其他后端，但日志改进有助于监控；3) 测试扩展可能增加 CI 执行时间，但提高了覆盖质量。总体风险低，因为变更针对特定后端和场景。

影响分析：对用户，ROCm 平台上的编码器 - 解码器模型波束搜索将产生正确结果，提升用户体验；对系统，后端调度更精准，避免计算错误，增强稳定性；对团队，测试覆盖扩展和日志改进有助于未来调试和 CI 维护。

关联脉络

本 PR 与历史 PR #38321 关联（如 Issue 评论所述），后者可能涉及类似交叉注意力问题，表明 vLLM 社区在持续优化 ROCm 后端支持。从近期历史 PR 看，多个 ROCm 相关 PR（如 #38415、#38413、#38252）关注 CI 和基础设施改进，本 PR 是其中一环，专注于注意力后端调度和测试修复，反映了 vLLM 在 ROCm 平台上的成熟度提升趋势。