

PR #37029 完整报告

vllm-project/vllm

[Hardware][XPU] Align memory usage with cuda on xpu

合并时间: 2026-03-25 18:14

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/37029>

执行摘要

本 PR 在 XPU 平台添加了 `empty_cache` 调用以对齐 CUDA 内存分析, 但 review 指出核心测量逻辑有误, 可能返回不准确的内存使用数据, 影响性能监控。变更仅涉及单个文件, ROCm 部分被回滚, 风险中等, 建议关注后续修复。

功能与动机

动机基于 issue #19033, CUDA 平台已添加 `empty_cache` 优化内存分析, 为确保跨平台一致性, XPU 和 ROCm 平台需同步此更改。PR body 中明确表述: "xpu/rocm platform should add this to keep memory profiling part align with cuda." 目的是统一内存分析行为, 避免平台差异导致监控偏差。

实现拆解

修改仅限一个文件, 按模块拆解如下:

- 模块: `platforms/xpu`
- 文件: `vllm/platforms/xpu.py`
- 关键改动: 在 `get_current_memory_usage` 方法中, 添加一行代码 `torch.xpu.empty_cache()`, 位于 `torch.xpu.reset_peak_memory_stats(device)` 之前。代码片段如下:

```
def get_current_memory_usage(
    cls,
    device: torch.types.Device | None = None
) -> float:
    torch.xpu.empty_cache()
    torch.xpu.reset_peak_memory_stats(device)
    return torch.xpu.max_memory_allocated(device)
```

- 注意: ROCm 平台的类似更改在提交 `fe14502` 中被回滚, 仅保留 XPU 部分。

评论区精华

review 讨论聚焦于测量逻辑正确性, `gemini-code-assist[bot]` 提出关键质疑:

"The logic for `get_current_memory_usage` on XPU appears incorrect. Calling `torch.xpu.max_memory_allocated(device)` after `reset_peak_memory_stats` will likely

return 0, not the current memory usage. To align with the ROCm and CUDA platforms, this function should use `torch.xpu.mem_get_info()`..." 此评论建议使用 `mem_get_info` 计算当前内存使用，但 PR 未采纳，仅添加了 `empty_cache`。两个人类 reviewer (yma11 和 bigPYJ1151) 批准了 PR，未进一步讨论，导致逻辑问题悬而未决。

风险与影响

风险：

- 测量不准确风险： `get_current_memory_usage` 方法在 XPU 上使用 `max_memory_allocated`，在重置统计后可能返回 0，而非真实内存使用，导致内存分析数据错误。
- 兼容性风险：未真正对齐 CUDA 和 ROCm 的实现（后者可能使用不同方法），跨平台比较可能产生不一致结果。

影响：

- 用户影响：对终端用户透明，但开发者在进行性能调优时可能依赖错误的内存数据。
- 系统影响：仅影响 XPU 平台的内存监控模块，不涉及核心推理或调度逻辑，影响范围有限。
- 团队影响：需后续跟进测量逻辑修复，以避免长期技术债务。

关联脉络

从历史 PR 分析中，未发现直接关联的 PR（如修改相同文件或功能线）。本 PR 基于 issue #19033，但材料中未提供该 issue 细节，可能涉及更早的 CUDA 内存优化 PR。近期历史 PR 中多有 ROCm 相关更改（如 38413、38415），但本 PR 的 ROCm 部分被回滚，暗示可能存在平台特定考量或测试问题。整体看，这反映了 vLLM 项目在跨硬件平台内存管理对齐上的持续努力，但当前实现存在缺陷，需后续优化。