

PR #41034 完整报告

vllm-project/vllm

[BugFix][CPU] fix error on CPU runner shutdown

合并时间: 2026-04-29 12:08

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/41034>

执行摘要

- 一句话: 修复 CPU 后端 shutdown 时 accelerator 调用崩溃
- 推荐动作: 该 PR 修复明确, 方案简洁 (+11 行), 值得立即合并。建议后续为 CPU 后端增加对应的 shutdown 控制流测试, 确保 monkey-patch 不遗漏未来新增的 accelerator API。开发者可参考该模式处理其他跨后端的兼容性问题。

功能与动机

Issue #41033 报告 CPU 后端执行 vllm bench throughput 后 worker 进程退出时崩溃, 报错 `RuntimeError: Cannot access accelerator device when none is available`。根本原因是 `CPUModelRunner` 继承了 `GPUModelRunner.shutdown()`, 其中调用了 `torch.accelerator.synchronize()`, CPU 环境下无 accelerator 设备。

实现拆解

1. 新增猴补丁函数: 在 `vllm/v1/worker/cpu_model_runner.py` 末尾定义 `_set_torch_accelerator_to_noop()`, 将 `torch.accelerator.synchronize` 和 `torch.accelerator.empty_cache` 替换为空的 `lambda` 或函数 `noop`。
2. 初始化时调用: 在 `CPUModelRunner.__init__` 开头、调用 `super().__init__` 之前执行 `_set_torch_accelerator_to_noop()`, 以确保父类初始化及后续 shutdown 中所有对这两个 accelerator API 的调用都变成空操作。
3. 无需重写 shutdown: 由于 monkey-patch 覆盖了所有调用点, `CPUModelRunner` 无需重写 shutdown 方法, 完全复用 `GPUModelRunner.shutdown` 中除 accelerator 调用外的其他清理逻辑 (如释放模型权重、清空 KV cache 等)。注意 `gemini-code-assist` 建议的显式 shutdown 重写已在最终提交中被更优雅的 monkey-patch 方案替代, 由 review 讨论驱动。

关键文件:

- `vllm/v1/worker/cpu_model_runner.py` (模块 模型运行器; 类别 source; 类型 core-logic; 符号 `_set_torch_accelerator_to_noop`, `noop`): 核心变更文件: 新增 monkey-patch 函数并在初始化时调用, 修复 CPU 后端 shutdown 崩溃。

关键符号: `_set_torch_accelerator_to_noop`, `noop`

关键源码片段

vllm/v1/worker/cpu_model_runner.py

核心变更文件：新增 monkey-patch 函数并在初始化时调用，修复 CPU 后端 shutdown 崩溃。

```
# vllm/v1/worker/cpu_model_runner.py

# 在文件末尾新增 monkey-patch 函数
def _set_torch_accelerator_to_noop() -> None:
    """将 torch.accelerator 的 API 替换为空操作，
    避免 CPU 环境因无 accelerator 设备导致 RuntimeError。"""
    def noop(*args: Any, **kwargs: Any) -> None:
        # 空函数，忽略所有参数
        pass

    # 这两个方法在 GPUModelRunner.shutdown 中被调用；
    # CPU 上不需要 GPU 同步或缓存清空。
    torch.accelerator.synchronize = noop
    torch.accelerator.empty_cache = noop

class CPUModelRunner(GPUModelRunner):
    def __init__(self, vllm_config: VllmConfig, device: torch.device):
        # 在调用 super().__init__ 前 patch，
        # 确保父类的 shutdown 等方法不会因 accelerator 调用而崩溃。
        _set_torch_accelerator_to_noop()

        with _torch_cuda_wrapper():
            super().__init__(vllm_config, device)
        # ... 后续初始化代码不变
```

评论区精华

核心讨论在 PR 作者 @fadara01 与 reviewer @bigPYJ1151 之间展开：

- 初始方案：PR 最初尝试在 CPUModelRunner 中重写 shutdown 为 pass（空实现），但这样会丢失 GPUModelRunner.shutdown 中的有用清理逻辑（如释放模型权重、清空 KV cache），导致内存泄漏（gemini-code-assist 指出）。
- 改进方案：@bigPYJ1151 建议改为 monkey-patch torch.accelerator.synchronize 和 torch.accelerator.empty_cache，使其在 CPU 环境下变为 no-op，从而最大程度复用 GPU 代码。此建议被作者采纳并重新提交。
- 结论：最终方案以 monkey-patch 方式实现，既修复了崩溃，又保留了父类的清理逻辑，得到 reviewer 批准合并。gemini-code-assist 的初始评论（关于 shutdown 空实现的潜在泄漏）因方案变更而过时。
 - monkey-patch vs shutdown 空实现 (design): 作者采纳建议，改为在 __init__ 时对 torch.accelerator.synchronize 和 empty_cache 进行 monkey-patch。
 - 潜在内存泄漏风险 (correctness): 该评论针对弃用的空实现方案，最终 monkey-patch 方案保留了父类 shutdown 的清理逻辑，该问题不再存在。

风险与影响

- 风险：
 - 回归风险低：改动仅影响 CPU 后端，且仅修改了 `torch.accelerator` 的两个方法，不影响 GPU 或其他后端行为。
 - 遗漏其他 `accelerator` 调用：该方法只 `patch` 了已知的 `synchronize` 和 `empty_cache`，若未来 `GPUModelRunner` 新增其他 `torch.accelerator` 调用（如 `torch.accelerator.current_device`），CPU 端可能再次崩溃。建议考虑更通用的防御性包装（如 `context manager` 或 `try-except`）。
 - 无测试覆盖：本次变更未添加对应的单元测试或集成测试，无法自动验证 `monkey-patch` 生效且不破坏现有功能。
- 影响：
 - 用户 / 运维：CPU 后端用户将不再遇到 `worker` 退出时崩溃，`vllm bench throughput` 等命令可正常退出。
 - 系统 / 团队：该修复保持了 GPU 代码的高复用性，避免为 CPU 重复实现 `shutdown` 逻辑。改动极小（+11 行），维护成本低。
 - 影响范围：仅限 `vLLM v1` 架构、CPU 后端，其他后端（GPU、XPU、Intel GPU）不受影响。
 - 风险标记：缺少测试覆盖，核心路径变更

关联脉络

- 暂无明显关联 PR