

# PR #38396 完整报告

vllm-project/vllm

[AMD][CI] Update DeepEP branch

合并时间: 2026-04-18 03:30

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/38396>

## 执行摘要

- 一句话: 更新 ROCm 平台 DeepEP 版本并调整 CI 测试配置, 修复 gfx942/gfx950 编译问题。
- 推荐动作: 此 PR 主要涉及基础设施更新, 对于关注 ROCm 平台或 CI/CD 流程的工程师值得浏览, 特别是 Dockerfile 中构建参数的用法和 CI 测试迁移的决策。对于核心模型推理或性能优化工程师, 优先级较低。

## 功能与动机

根据 PR 描述和 Issue 评论, 此 PR 的目的是更新 DeepEP 分支到一个能正确为 gfx942 和 gfx950 架构进行提前编译的版本, 以部分解决 issue #37709。作者在评论中解释, 当前 ROCm DeepEP 仅支持 gfx942 和 gfx950 架构 (引用自 DeepEP 仓库的 setup.py), 而更新版本后, DeepEP 会在链接时将适当的内核打包到其二进制文件中。此外, 由于 CI 环境中目前没有 MI355 代理, 需要将测试用例迁移到 MI325 节点来验证变更。

## 实现拆解

1. 更新 DeepEP 依赖版本和构建参数: 在 docker/Dockerfile.rocm 中, 将 DeepEP 分支从 e84464ec 更新为 5d90af8b, 并新增构建参数 DEEPEP\_ROCM\_ARCH, 其值设为 gfx942;gfx950。同时, 将 ROC-SHMEM 的构建过程从直接调用 cmake 改为使用脚本 `bash ../scripts/build_configs/all_backends`, 并传递 `-DGPU_TARGETS="${DEEPEP_ROCM_ARCH}"` 参数, 以确保针对指定 GPU 架构编译。
2. 调整 CI 测试配置: 在 .buildkite/test-amd.yaml 中, 将 DeepEP 数据并行测试命令从 MI355 节点的测试步骤移动到 MI325 节点的测试步骤。具体来说, 在 MI325 节点的 Distributed Tests (2 GPUs)(H100-MI325) 步骤中添加命令 - `VLLM_LOGGING_LEVEL=DEBUG python3 examples/offline_inference/data_parallel.py --model=Qwen/Qwen1.5-MoE-A2.7B -tp=1 -dp=2 --max-model-len=2048 --all2all-backend=deepep_high_throughput`, 并从 MI355 节点的对应步骤中移除该命令。
3. 清理构建配置: 在提交历史中, 第三个提交移除了 ROC-SHMEM 构建中冗余的 `-DCMAKE_POSITION_INDEPENDENT_CODE=ON` 标志, 简化了 Dockerfile。

关键文件:

- docker/Dockerfile.rocm (模块 部署脚本; 类别 infra; 类型 infrastructure) : 更新了 DeepEP 版本和构建参数, 修复 gfx942/gfx950 架构的编译问题, 是 PR 的核心变更。

- `.buildkite/test-amd.yaml` (模块 CI 配置; 类别 `config`; 类型 `configuration`): 调整 CI 测试配置, 将 DeepEP 数据并行测试从 MI355 节点迁移到 MI325 节点, 以适配当前 CI 环境。

关键符号: 未识别

## 评论区精华

review 中主要有两个讨论点:

- GPU 架构硬编码问题: `gemini-code-assist[bot]` 指出 `GPU_TARGETS` 被硬编码, 但引入了 `DEEPEP_ROCM_ARCH` 构建参数, 建议使用该参数以提高可配置性。作者在实现中已采纳此建议, 通过 `-DGPU_TARGETS="${DEEPEP_ROCM_ARCH}"` 传递参数。
- 架构支持范围: `gshtras` 询问测试是否不应在 250s (可能指 `gfx250` 架构) 上运行, 作者回复澄清 ROCm DeepEP 目前仅支持 `gfx942` 和 `gfx950`, 并提供了 DeepEP 仓库的代码链接作为证据。整体讨论较少, `tjtanaa` 在确认修复链接和预期状态后批准了 PR。
- GPU 架构硬编码与可配置性 (`design`): 作者在实现中已采纳建议, 通过 `-DGPU_TARGETS="${DEEPEP_ROCM_ARCH}"` 传递参数, 使构建过程更灵活。
- DeepEP 支持的 GPU 架构范围 (`question`): 作者提供了 DeepEP 仓库的代码链接作为证据, 确认架构支持限制, 解释了测试迁移的原因。

## 风险与影响

- 风险: 技术风险较低, 主要涉及基础设施和 CI 配置:
- 兼容性风险: DeepEP 版本更新可能引入未知的二进制兼容性问题, 但 PR 描述中的测试结果显示数据并行测试通过, 降低了风险。
- CI 配置风险: 将测试从 MI355 移动到 MI325 节点, 如果两个节点的硬件或软件环境存在差异, 可能导致测试结果不一致。但 PR 描述表明测试在 MI325 上通过, 且 MI355 节点当前不可用, 因此风险可控。
- 构建过程风险: ROC-SHMEM 构建过程改为脚本驱动, 如果脚本行为不稳定或与未来版本不兼容, 可能影响 Docker 镜像构建。但变更较小, 且基于官方脚本, 风险有限。
- 影响: 影响范围主要限于 ROCm 平台的基础设施和 CI 流水线:
- 对用户的影响: 普通用户无直接影响, 除非他们使用基于更新后 Docker 镜像的自定义部署或依赖 DeepEP 功能。
- 对系统的影响: 修复了 `gfx942/gfx950` 架构的编译问题, 提升了 ROCm 平台上 DeepEP 的可用性和正确性。
- 对团队的影响: CI 测试配置调整确保了 DeepEP 相关测试能在可用节点上运行, 提高了测试覆盖率和可靠性, 有助于持续集成。
- 风险标记: 依赖版本更新, CI 配置迁移

## 关联脉络

- PR #39978 [ROCm][CI] Build fastsafetensors from source so it links against libamdhip64: 同样涉及 ROCm 平台的 CI 和 Dockerfile 更新, 聚焦于依赖构建和链接问题

，可对比基础设施改进模式。

- PR #39953 [ROCm] Fix TurboQuant on ROCm: backend routing, flash-attn compat, int64 overflow: 同为 ROCm 平台的 bugfix，涉及后端路由和兼容性修复，展示了跨平台支持的持续优化。