

# PR #37283 完整报告

vllm-project/vllm

[Releases] [ROCm] Enable Nightly Docker Image and Wheel Releases for ROCm

合并时间: 2026-03-27 00:32

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/37283>

## 执行摘要

本 PR 通过重构 Buildkite 发布流水线，为 ROCm 平台启用了 nightly Docker 镜像和 Python wheel 的自动化发布。这一变更显著提升了 ROCm 用户的开发体验，通过缓存机制优化了长达 3 小时的构建时间，并经过多次迭代解决了配置管理和发布流程中的关键问题。

## 功能与动机

此 PR 的主要功能是启用 ROCm 的 nightly 构建发布，以解决 issue #36703 中用户对最新构建的需求。动机源于 ROCm 依赖构建时间过长（超过 3 小时），需通过缓存基础 Docker 镜像和 wheels 来改善效率。PR body 详细描述了用户将如何通过 Docker Hub 获取标签如 `vllm/vllm-openai-rocm:nightly` 的镜像，以及通过 S3 路径下载 wheels。

## 实现拆解

实现方案按模块拆解如下：

- CI 流水线重构：修改 `.buildkite/release-pipeline.yaml`，移除手动配置输入，整合 ROCm 基础镜像和 wheels 构建步骤，并添加 nightly 镜像推送任务。关键变更包括简化构建逻辑，使用缓存键基于 `Dockerfile.rocm_base` 哈希。
- 新增发布脚本：`.buildkite/scripts/push-nightly-builds-rocm.sh` 是新增脚本，负责从 ECR 拉取镜像并推送到 Docker Hub，支持标签如 `base-nightly` 和 `nightly-<commit>`。
- 缓存优化：修改 `.buildkite/scripts/cache-rocm-base-wheels.sh`，移除对 `PYTHON_VERSION` 和 `PYTORCH_ROCM_ARCH` 的依赖，缓存键仅基于 `Dockerfile` 内容，提升稳定性。
- 辅助脚本更新：更新 `.buildkite/scripts/annotate-release.sh` 和 `cleanup-nightly-builds.sh` 以支持 ROCm 特定仓库，确保发布注释和清理操作正确。

## 评论区精华

Review 讨论中突出了几个关键交锋：

- `S3_BUCKET` 变量变更：gemini-code-assist[bot] 指出变量从 `vllm-wheels` 改为 `vllm-wheels-dev` 可能影响发布目标，tjtanaa 回应是试用值并已更正。
- `DRY_RUN` 设置：评论指出 `DRY_RUN` 设置为 1 会阻止发布，tjtanaa 解释为试用运行，后续提交中移除。

- 配置提取的脆性：gshtras 评论从 Dockerfile.rocm\_base 提取配置易受未来变化影响，tjtanaa 随后移除这些提取，仅基于文件内容生成缓存键，决策结论是避免隐藏依赖，提升设计鲁棒性。

## 风险与影响

### 风险分析：

- 缓存键生成仅依赖 Dockerfile.rocm\_base 哈希，若环境变量未纳入，可能导致缓存无效或构建不一致。
- 发布脚本中的网络或权限错误可能阻止镜像推送到 Docker Hub，影响 nightly 发布可用性。
- 清理脚本 cleanup-nightly-builds.sh 修改后支持多仓库，但参数错误可能导致误删标签，需谨慎测试。

### 影响分析：

- 对用户：ROCm 用户现在可便捷获取 nightly 构建，加速开发和测试迭代。
- 对系统：CI 流水线增加 ROCm nightly 发布步骤，可能轻微增加资源消耗，但通过缓存优化整体效率。
- 对团队：简化发布流程，减少手动操作，促进更频繁的集成测试。

## 关联脉络

此 PR 与历史 PR #32264 密切相关，重用了其缓存逻辑为基础 Docker 镜像和 wheels 提供支持。从近期历史 PR 看，如 #38161 和 #38167，vLLM 仓库持续优化 ROCm 相关 CI 测试和性能，本 PR 是这一趋势的一部分，旨在扩展发布基础设施以支持更稳定的 ROCm 生态。关联 issue #36703 提供了直接需求背景，推动了此功能的实现。