

PR #24565 完整报告

sgl-project/sglang

Expand support matrix for pypi wheel release

合并时间: 2026-05-07 08:39

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/24565>

执行摘要

- 一句话: 扩展 PyPI 发布矩阵至多 Python 版本和双架构
- 推荐动作: 该 PR 是基础设施改进, 值得运维和 CI 团队关注。矩阵构建和并行发布的设计模式可复用至其他 Python 包发布流程。

功能与动机

为 sglang 提供更广泛的 PyPI 支持矩阵, 覆盖 Python 3.10-3.13 以及 x86_64 和 aarch64 架构, 同时确保在自托管构建节点上正确清理残留文件、集成 `setuptools-rust` 以构建原生 gRPC 扩展。

实现拆解

1. 拆分构建与发布为两个独立 Job: 将原 `publish` 拆分为 `build` 和 `publish`, 前者负责在所有配置下构建 wheel, 后者统一收集并发布到 PyPI。
2. 引入策略矩阵: 在 `build job` 中设置 `strategy.matrix`, 组合 `python-version: ["3.10", "3.11", "3.12", "3.13"]` 和 `arch: [x86_64, aarch64]`, 并利用 `include` 映射到对应自托管 runner (`x64-docker-build-node / arm-docker-build-node`), 实现多平台并行构建。
3. 添加 Rust 和 `protoc` 依赖: 新增 `dtolnay/rust-toolchain@stable` 安装 Rust 工具链, 并保留 `install_protoc.sh` 步骤, 确保 `setuptools-rust` 能编译 `sglang-grpc` 原生扩展。
4. 清理自托管节点残留文件: 在 `checkout` 前使用 Alpine 容器删除先前构建遗留的 `root` 所有权文件, 避免 `actions/checkout` `EACCES` 错误。
5. 制品上传与下载: `build job` 末尾使用 `upload-artifact@v4` 按矩阵组合命名上传 wheel 文件; `publish job` 通过 `download-artifact@v4` 下载所有制品, 合并后一次性发布到 PyPI。

关键文件:

- `.github/workflows/release-pypi.yml` (模块 CI/CD; 类别 `infra`; 类型 `infrastructure`): 核心变更文件, 完整重构了 PyPI 发布工作流, 引入矩阵构建、多架构支持、Rust 编译依赖和制品管理。

关键符号: 未识别

评论区精华

无 review 讨论。

- 暂无高价值评论线程

风险与影响

- 风险：
 - 构建环境依赖：aarch64 构建依赖自托管 arm-docker-build-node runner 可用性，若 runner 不可用或配置不当将阻塞发布。
 - 版本一致性：矩阵中各构建可能因缓存或依赖版本差异产生不一致的 wheel，需确保依赖锁定。
 - 安全问题：自托管 runner 上的 Docker 清理操作可能影响其他作业，但已设置 `ll true` 避免严格失败。
- 影响：
 - 用户：用户可从 PyPI 直接安装适用于 Python 3.10-3.13 和 arm64 架构的 sglang 版本，无需自行编译，降低使用门槛。
 - 系统：发布流程并发构建增多，GitHub Actions 分钟消耗增加，但总体发布效率提升。
 - 团队：发布矩阵扩展减少手动干预，降低出错概率。
 - 风险标记：自托管 runner 可用性，依赖版本一致性，Docker 清理副作用

关联脉络

- 暂无明显关联 PR