

PR #22405 完整报告

sgl-project/sglang

[CICD] [prefill-only] Consolidate prefill-only model E2E tests

合并时间: 2026-04-09 15:54

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/22405>

执行摘要

本 PR 对 sglang 仓库中的 prefill-only 模型端到端测试进行了大规模重构，将原本分散在四个无关目录的测试文件统一迁移至新的 `test/registered/prefill_only/` 目录，并清晰分离了输入嵌入功能测试。变更涉及 16 个文件的移动、重命名和新增，更新了 CODEOWNERS 以匹配新结构。作者已本地验证所有测试通过，旨在提升代码维护性、CI 覆盖率和团队协作效率，对用户无直接影响，但为后续测试扩展奠定基础。

功能与动机

为什么做：根据 PR body，prefill-only 模型测试此前“散落在四个无关目录中，没有统一的家”，导致维护困难、CI 覆盖不全（例如 `test/srt/test_multi_item_scoring.py` 未注册在 CI 中）。目标是通过“统一 prefill-only 模型端到端测试”来解决这些问题，提供一致的测试结构和更好的可管理性。

实现拆解

实现按模块分层进行：

- 新目录创建：建立 `test/registered/prefill_only/` 目录，集中存放所有 prefill-only 相关测试，包括从 `core/`、`models/`、`openai_server/basic/` 和 `embedding/` 移动来的文件。
- 概念分离：将原 `test/registered/embedding/` 重命名为 `test/registered/input_embedding/`，仅保留输入嵌入功能测试，避免与 prefill-only 测试混淆。
- 单元测试重组：移动 `test_pooler_score_and_pool.py` 从 `unit/` 根目录到 `unit/layers/`，以匹配源码模块 `sglang.srt.layers.pooler` 的层级。
- 计分测试优化：用两个新文件替换旧计分测试：
 - `test_score_engine.py`：专注于引擎 API 层的正确性测试，覆盖 CausalLM 和 SequenceClassification 模型的单点和批量评分。
 - `test_score_api.py`：专注于 HTTP `/v1/score` 端点的集成测试，验证响应结构和 CLI 参数传递。
- 基础设施更新：修改 `.github/CODEOWNERS`，将所有权从旧文件路径扩展到整个 `prefill_only` 目录。

评论区精华

review 中无实质性讨论，仅审核者 hnyls2002 批准。PR body 中作者单方面阐述了变更理由和验证结果，例如：“All new tests were verified locally before merging”，并提供了测试运行时间和结果。没有外部争议或深度技术交锋。

风险与影响

风险分析：

1. 路径变更风险：文件移动可能导致 CI 脚本或外部工具引用失效，但作者已更新 CODEOWNERS 并本地验证，降低了风险。
2. 测试覆盖验证：移除旧文件并新增测试需确保功能等价，作者通过对比 HuggingFace 参考和本地测试已覆盖。
3. 团队适应成本：新目录结构需开发者调整 workflow，但长期看维护效率提升。

影响分析：

- 用户影响：无直接影响，仅内部测试代码变更。
- 系统影响：测试更结构化，可能加速 CI 执行和问题调试。
- 团队影响：代码所有权更清晰，但需短期学习新布局。

关联脉络

从近期历史 PR 看，本 PR 与多个测试和 CI 优化 PR 相关联：

- PR #22418 将 Runai 模型加载测试移至夜间套件，与本 PR 的测试重组理念相似，都旨在优化 CI 管理。
- PR #22353 新增 Torch Profiler 分析 workflow，涉及测试代码扩展，与本 PR 的测试组织改进互补。
- PR #22400 为 CI 添加快速失败机制，与本 PR 的 run-ci 标签和测试效率提升主题一致。这些 PR 共同反映了仓库在强化测试基础设施和 CI 管道方面的持续演进趋势。