

PR #21873 完整报告

sgl-project/sglang

[Misc] Add network timeout to eval dataset downloads

合并时间: 2026-04-02 04:16

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/21873>

执行摘要

- 一句话: 为评估数据集下载添加网络超时, 避免 CI 因网络问题无限挂起。
- 推荐动作: 该 PR 变更直接, 适合快速浏览以了解 CI 优化措施。关注点: 超时值 30 秒的合理性, 以及是否需要在其他类似场景中推广此模式。

功能与动机

根据 PR body 描述, `pandas.read_csv(url)` 默认无超时设置, 当网络到 Azure Blob Storage 缓慢或不可达时, CI 作业会无限挂起。具体案例包括: <https://github.com/sgl-project/sglang/actions/runs/23846257201/job/69514908983> (MMLU 下载挂起 15 分钟) 和 <https://github.com/sgl-project/sglang/actions/runs/23846257201/job/69579769522> (挂起 11 分钟后按文件超时)。PR #21800 修复了 `simple_eval_common.py` 和 `simple_eval_mgsm.py`, 但遗漏了 `pandas.read_csv()` 路径, 本 PR 旨在补全这些遗漏。

实现拆解

实现分为两部分: 1) 在 `python/sglang/test/` 目录下的 `simple_eval_mmlu.py`、`simple_eval_gpqa.py`、`simple_eval_math.py` 中, 修改 `pandas.read_csv()` 调用, 通过检查文件名是否包含 `://"` 来区分本地文件和远程 URL, 对远程 URL 添加 `storage_options={"timeout": 30}` 参数; 2) 在 `python/sglang/test/simple_eval_common.py` 和 `sgl-model-gateway/e2e_test/infra/simple_eval_common.py` 中, 将 `requests.get()` 调用的超时从 120 秒调整为 30 秒, 以保持一致性。

关键文件:

- `python/sglang/test/simple_eval_mmlu.py` (模块 `test`): 修改了 MMLU 评估数据集的下载逻辑, 添加超时以避免 CI 挂起, 是核心变更文件之一。
- `python/sglang/test/simple_eval_gpqa.py` (模块 `test`): 修改了 GPQA 评估数据集的下载逻辑, 添加超时以避免 CI 挂起, 是核心变更文件之一。
- `python/sglang/test/simple_eval_math.py` (模块 `test`): 修改了 MATH 评估数据集的下载逻辑, 添加超时以避免 CI 挂起, 是核心变更文件之一。
- `sgl-model-gateway/e2e_test/infra/simple_eval_common.py` (模块 `model-gateway`): 调整了模型网关中数据集下载的超时设置, 从无超时改为 30 秒, 确保一致性。

关键符号: `pandas.read_csv`, `download_dataset`

评论区精华

无 review 评论，但 PR body 中提到了与 PR #21800 的关联，表明这是对先前修复的补充。提交信息显示作者将超时从 120 秒调整为 30 秒，但未提供调整理由的讨论。

- 暂无高价值评论线程

风险与影响

- 风险：技术风险较低：1) 超时设置可能过短（30 秒），在正常网络波动下导致不必要的下载失败，影响 CI 稳定性；2) 修改涉及多个文件，但逻辑简单，回归风险小；3) 对本地文件路径（不含 "://"）的 `pandas.read_csv()` 调用保持不变，兼容性良好。
- 影响：影响范围限于 CI 测试环境：1) 对用户无直接影响，仅影响内部测试流程；2) 系统层面，避免了 CI 无限挂起，提高了资源利用率和测试可靠性；3) 团队层面，减少了因网络问题导致的 CI 阻塞，提升了开发效率。
- 风险标记：超时设置可能过短

关联脉络

- PR #21800 [Misc] Add network timeout to eval dataset downloads: 本 PR 是 #21800 的补充，修复了 #21800 中遗漏的 `pandas.read_csv()` 路径，两者共同解决 CI 因网络问题挂起的问题。