

PR #22045 完整报告

sgl-project/sglang

[CI] Adjust CI server launch timeout

合并时间: 2026-04-03 22:38

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/22045>

执行摘要

本 PR 在测试框架中为模型配置添加了可自定义的服务器启动超时参数 `launch_timeout`, 允许针对特定测试 (如 8-GPU 模型测试) 设置更长超时, 以解决因模型加载时间过长导致的 CI 失败。变更涉及测试工具类和具体测试用例, 风险较低, 但参数文档未更新可能影响可维护性。

功能与动机

该 PR 旨在修复一个具体的 CI 失败 (链接: <https://github.com/sgl-project/sglang/actions/runs/23928333410/job/69789912643>)。从关联 Issue 评论中作者执行了重跑测试命令 `/rerun-test test_ring_2_5_1t.py`, 结合修改文件 `test_ring_2_5_1t.py`, 可推断失败原因是 8-GPU 模型测试中服务器启动时间超过默认超时, 导致 CI 中断。因此, PR 通过扩展测试配置能力, 为这类资源密集型测试提供更灵活的超时控制。

实现拆解

实现分为三个层次:

- 配置扩展: 在 `python/sglang/test/test_utils.py` 的 `ModelConfig` 类中新增可选参数 `launch_timeout`, 类型为 `Optional[float]`, 默认 `None`。
- 逻辑集成: 在 `python/sglang/test/accuracy_test_runner.py` 的两个函数 `_run_simple_eval` 和 `_run_nemo_skills_eval` 中, 将硬编码超时替换为条件回退逻辑:
- 用例适配: 在 `test/registered/8-gpu-models/test_ring_2_5_1t.py` 中, 为测试模型显式设置 `launch_timeout=1800` (30 分钟), 覆盖默认值。

评论区精华

Review 中仅有一条来自 `gemini-code-assist[bot]` 的评论, 聚焦于代码文档:

The `launch_timeout` parameter is added to the constructor but not documented in the class docstring or type hints for clarity. Adding a brief description would improve maintainability.

该评论指出新增参数缺乏文档描述, 可能影响开发者理解和使用, 但未引发进一步讨论或修改。

风险与影响

- 风险:

- 条件表达式 `model.launch_timeout or DEFAULT_TIMEOUT_FOR_SERVER_LAUNCH` 在 `launch_timeout` 为 0 时会错误地回退到默认值（因为 0 在 Python 中视为 False），但当前使用场景中超时值应大于 0，风险可控。
- 参数文档缺失可能增加维护成本，需后续补充。
- 影响：
 - 对 CI 稳定性：可减少特定测试因超时导致的失败，提升测试可靠性。
 - 对开发者：提供了自定义超时的灵活性，但需注意文档缺口。
 - 不影响现有测试默认行为，兼容性良好。

关联脉络

从近期历史 PR 看，本 PR 与多个 CI 基础设施调整相关（如 PR#22031 临时禁用准确性测试、PR#22001 修复工作流作业名称），共同反映了团队对 CI 稳定性和效率的持续优化。本 PR 采用“延长超时”而非“禁用测试”的策略，体现了在保证测试覆盖的前提下解决资源争用问题的思路。未来类似大型模型测试可能需要借鉴此模式，但需平衡超时设置与 CI 执行时间。