

PR #26994 完整报告

sgl-project/sglang

jit_kernel tests: bump multiprocessing_test timeout 90s -> 240s (cold JIT cache)

合并时间: 2026-06-03 05:33

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/26994>

执行摘要

- 一句话: 调高 JIT kernel 测试超时阈值
- 推荐动作: PR 变更简单, 值得关注的是其根因分析思路 (对比不同参数化测试耗时、推断冷 JIT 缓存)。建议阅读 PR body 中的“smoking gun”分析, 理解如何从 CI 日志中定位非死锁类超时问题。长期方案 (固定 JIT 缓存路径) 值得跟进。

功能与动机

`test_tp_qknorm[2]` 在 CI 中持续超时 (>90s), 而同一 job 中后续参数化测试 (4 进程、8 进程) 仅需 ~45-52s 即可通过。根因是首次调用时 triton + cutlass JIT 编译开销大 (实测 60-180s), 后续参数化测试复用已 warm 的缓存。PR body 明确指出: “This is the minimum to stop the flake bleeding.”

实现拆解

变更仅涉及一个文件 `python/sglang/jit_kernel/tests/utils.py`, 核心修改是:

1. 修改默认超时值: 将 `multiprocess_test` 函数签名中的 `timeout: int = 90` 改为 `timeout: int = 240`。
2. 更新 docstring: 补充默认值设定依据, 说明首次调用的 JIT 编译开销可达 60-180 秒, 后续参数化测试因缓存预热可降低至 ~60 秒。
3. 不需要改动调用方: `test_custom_all_reduce.py` 和 `test_tp_qknorm.py` 都继承新默认值。

关键文件:

- `python/sglang/jit_kernel/tests/utils.py` (模块 测试工具; 类别 test; 类型 test-coverage; 符号 `multiprocess_test`): 修改默认超时值并更新 docstring, 是 PR 唯一变更的文件。

关键符号: `multiprocess_test`

评论区精华

两位审核者 (DarkSharpness、hnyls2002) 均直接批准, 无 review 评论。讨论主要为 PR body 中的根因分析: 通过对比 `test_tp_qknorm[2]` 超时 (>90s) 与 `[4]` (~45s)、`[8]` (~52s) 通过的现象, 结合 worker stdout 仅输出 `OMP_NUM_THREADS` 警告, 推断首次调用需编译大量 kernel 导致超时。

- 暂无高价值评论线程

风险与影响

- 风险：风险极低：
 - 超时从 90s 增加到 240s 仅影响测试用例的等待时间，不会改变测试逻辑。
 - 如果测试本身存在死锁或无限循环，240s 后仍会失败。
 - 不会影响生产代码。
 - 长期来看，增加超时可能掩盖真正性能回归，但 PR body 已明确此为临时措施，并计划通过固定 JIT 缓存路径彻底解决。
 - 影响：影响范围仅限于 jit-kernel 测试套件中的 multiprocess_test 函数调用者，主要影响 test_custom_all_reduce 和 test_tp_qknorm 两个测试文件。CI 中预期不再因冷 JIT 编译超时而失败。对用户无影响。
- 风险标记：临时缓解措施，仅测试配套

关联脉络

- PR #24757 Optimize ngram decode id computation: 同属 jit-kernel 模块，引入了新的 JIT kernel 测试，可能类似地受冷 JIT 缓存影响。