

PR #22187 完整报告

sgl-project/sglang

[HiSparse]: Add benchmark for hisparse kernel

合并时间: 2026-04-13 12:49

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/22187>

执行摘要

本 PR 为 hisparse JIT 内核新增基准测试脚本 `bench_hispase.py`，通过模拟缓存加载场景，测量不同配置下的内核执行延迟。此变更属于测试基础设施扩展，不影响核心功能，旨在为性能优化提供数据支撑，已集成到 CI 中。

功能与动机

动机源于需要系统化评估 hisparse JIT 内核的性能，PR body 中的性能数据表格显示了在不同 batch size、hot buffer size 和 miss rate 下的延迟指标。结合历史 PR 22155（为 hisparse 添加 CI 测试），此 PR 是持续改进测试覆盖的一部分，旨在标准化性能基准，助力内核优化决策。

实现拆解

实现集中在单个文件 `python/sglang/jit_kernel/benchmark/bench_hispase.py`，关键改动点包括：

- 配置定义：设置 batch sizes (1、10、100)、hot buffer sizes (4096、8192) 和 miss rates (0.2、0.001)，覆盖典型场景。
- 输入构建：使用 torch 创建模拟张量，如 top_k_tokens、host_cache 和 device_buffer，模拟缓存命中与缺失。
- 内核调用：调用 `load_cache_to_device_buffer_mla` 函数执行 hisparse 内核。
- 时间测量：通过循环计时测量延迟，并注册到 CI 的 `stage-b-kernel-benchmark-1-gpu-large` 套件中。

代码示例片段（基于 patch_excerpt）：

```
CONFIGS = [  
    (batch_size, hot_buffer_size, miss_rate, batch_size * round(TOP_K * miss_rate))  
    for batch_size, hot_buffer_size, miss_rate in itertools.product(  
        BATCH_SIZES, HOT_BUFFER_SIZES, MISS_RATES  
    )  
]
```

评论区精华

gemini-code-assist[bot] 在 review 中提出了多项优化建议，核心讨论如下：

“增加 WARMUP_ROUNDS 和 ROUNDS 以获得更稳定结果。”“使用 `host_cache.normal_()` 代替 `copy_` 以减少内存分配。”“建议使用 `triton.testing.do_bench` 进行更鲁棒的统计测量。”“修正函数返回值，避免重复相同值。”

讨论聚焦于提升基准测试的准确性和效率，无重大争议，作者通过提交修复了 lint 问题，可能已采纳部分建议。

风险与影响

风险分析：

- 测量准确性依赖输入构建和计时逻辑的正确性，如有误可能导致性能数据偏差。
- 外部库（如 triton）版本兼容性可能影响结果可复现性。
- 代码风格问题已在提交中修复，降低了维护风险。

影响分析：

- 影响范围限于测试基础设施，为 hisparse 内核提供性能基准，支持团队监控和优化工作。
- 集成到 CI 中，有助于防止性能回归，但不影响用户功能或系统核心路径。

关联脉络

从历史 PR 看，此 PR 与 PR 22155（为 hisparse 添加 CI 测试）直接相关，两者共同构建 hisparse 内核的测试生态。近期 PR 中涉及 jit-kernel 和基准测试的较多（如 PR 22631、22649），表明团队在持续强化内核性能验证。此 PR 是这一趋势的具体体现，为后续优化提供了数据基础。