

PR #25684 完整报告

sgl-project/sglang

[CI] Enable weight prefetch for 8-gpu-h200 basic tests

合并时间: 2026-05-19 05:12

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/25684>

执行摘要

- 一句话: 为 8-GPU H200 CI 测试启用权重预取
- 推荐动作: 该 PR 技术难度低、风险小, 但效益显著, 建议快速合并。其中值得关注的设计决策是在多线程加载的同时启用预取——两者叠加带来了 5.9x 的加载加速, 表明两者协作良好。后续可将此标志扩展到其他 8-GPU 测试, 以进一步加速 CI 套件。

功能与动机

8-GPU H200 CI 套件中有两个每次提交必运行的测试 (`TestMiniMaxM25Basic` 和 `TestUnifiedMambaHiCache`) 完全被检查点权重加载所主导。如 PR body 所述: "Without prefetch, the loader streams shards from disk before staging them to GPU, which inflates server boot time and slows the per-commit gate." 目标是通过已存在的 `--weight-loader-prefetch-checkpoints` 标志让磁盘读取与 H2D 拷贝重叠, 从而加速 CI 流水线。

实现拆解

1. 在 Minimax M2.5 测试中启用预取: 在 `test/registered/8-gpu-models/test_minimax_m25_basic.py` 的 `setUpClass` 中, 将 `--weight-loader-prefetch-checkpoints` 添加到 `popen_launch_server` 的 `other_args` 列表。该标志与现有的多线程加载配置 (`enable_multithread_load: true, num_threads: 64`) 协同工作。
2. 在 Mamba HiCache 测试中启用预取: 在 `test/registered/radix_cache/test_unified_radix_cache_kl_hicache.py` 的 `TestUnifiedMambaHiCache.setUpClass` 中, 将 `--weight-loader-prefetch-checkpoints` 添加到 `other_args` 列表末尾。另一个测试类 `TestUnifiedDeepSeekV4FlashHiCache` 保持不变。
3. 更新预估运行时间: 根据实测运行时间更新 `register_cuda_ci(est_time=...)` 调用的时间参数。Minimax 测试从 290s 降至 250s (实际 246s), Mamba 测试从 768s 降至 745s (实际 747s), 使调度器能更准确分配时间窗口。

关键文件:

- `test/registered/8-gpu-models/test_minimax_m25_basic.py` (模块 Minimax 测试; 类别 test; 类型 test-coverage): 为该 Minimax-M2.5 TP=8 测试添加 `--weight-loader-prefetch-checkpoints` 标志, 并将 `est_time` 从 290s 更新为 250s。是 CI 时间节省的主要贡献者 (权重加载从 120.88s→20.32s)。

- test/registered/radix_cache/test_unified_radix_cache_kl_hicache.py (模块 HiCache 测试; 类别 test; 类型 test-coverage) : 为 TestUnifiedMambaHiCache 测试添加 `--weight-loader-prefetch-checkpoints` 标志, 并将 `est_time` 从 768s 更新为 745s。权重加载从 57.05s→23.82s (2.4x 加速)。另一测试类 TestUnifiedDeepSeekV4FlashHiCache 未受影响。

关键符号: 未识别

关键源码片段

test/registered/8-gpu-models/test_minimax_m25_basic.py

为该 MiniMax-M2.5 TP=8 测试添加 `--weight-loader-prefetch-checkpoints` 标志, 并将 `est_time` 从 290s 更新为 250s。是 CI 时间节省的主要贡献者 (权重加载从 120.88s→20.32s)。

```
# test/registered/8-gpu-models/test_minimax_m25_basic.py
# 改动: 新增预取标志并更新 est_time
```

```
register_cuda_ci(est_time=250, stage="base-c", runner_config="8-gpu-h200")
```

```
class TestMiniMaxM25Basic(CustomTestCase):
    @classmethod
    def setUpClass(cls):
        cls.model = MINIMAX_M25_MODEL_PATH
        cls.base_url = DEFAULT_URL_FOR_TEST
        other_args = [
            "--trust-remote-code",
            "--tp", "8",
            "--ep-size", "8",
            "--mem-fraction-static", "0.85",
            "--reasoning-parser", "minimax-append-think",
            "--model-loader-extra-config",
            '{"enable_multithread_load": true, "num_threads": 64}',
            "--weight-loader-prefetch-checkpoints", # 新增: 启用权重预取以重叠磁盘读取与 H2D 拷贝
        ]
        cls.process = popen_launch_server(
            cls.model,
            cls.base_url,
            timeout=DEFAULT_TIMEOUT_FOR_SERVER_LAUNCH,
            other_args=other_args,
        )
```

test/registered/radix_cache/test_unified_radix_cache_kl_hicache.py

为 TestUnifiedMambaHiCache 测试添加 `--weight-loader-prefetch-checkpoints` 标志, 并将 `est_time` 从 768s 更新为 745s。权重加载从 57.05s→23.82s (2.4x 加速)。另一测试类 TestUnifiedDeepSeekV4FlashHiCache 未受影响。

```
# test/registered/radix_cache/test_unified_radix_cache_kl_hicache.py
# 改动: 仅在 TestUnifiedMambaHiCache 中启用预取, 另一测试类不变
```

```
register_cuda_ci(est_time=745, stage="base-c", runner_config="8-gpu-h200")
```

```
class TestUnifiedMambaHiCache(UnifiedRadixTreeTestMixin, CustomTestCase):
    # ...
    @classmethod
    def setUpClass(cls):
        cls.model = MAMBA_MODEL
        cls.base_url = DEFAULT_URL_FOR_TEST
        cls.process = popen_launch_server(
            cls.model,
            cls.base_url,
            timeout=DEFAULT_TIMEOUT_FOR_SERVER_LAUNCH,
            other_args=[
                "--tp-size", "4",
                "--chunked-prefill-size", "2048",
                # ... 其他标志 ...
                "--max-running-requests", "4",
                "--weight-loader-prefetch-checkpoints", # 新增: 启用权重预取
            ],
            env={"SGLANG_ENABLE_UNIFIED_RADIX_TREE": "1"},
        )
        cls.input_ids = get_input_ids(cls.model, num_samples=18)
```

评论区精华

无 review 评论或讨论。作者通过 `/rerun-test` 命令手动触发了两次 CI 运行（`test_minimax_m25_basic.py` 和 `test_unified_radix_cache_kl_hicache.py`），均通过。

- 暂无高价值评论线程

风险与影响

- 风险：风险极低。变更仅限于 CI 测试文件的配置参数（新增启动标志、更新时间估算），不修改任何模型代码、公共 API 或推理行为。该标志已是加载器中的现有功能，仅影响加载阶段的磁盘 I/O 和 H2D 拷贝调度。测试断言（GSM8K 准确率 >0.900、KL 散度检查、多轮缓存命中检查）保持不变，在手动触发的 CI 运行中均通过。
- 影响：影响范围：仅限于 base-c-test-8-gpu-h200 CI 套件中的两个测试文件。对用户：无直接影响。对系统：CI 流水线中这两项测试的总运行时间减少约 120 秒（原约 1112 秒，现约 993 秒），权重加载时间减少约 134 秒，加速了每次提交的 CI 门控。对团队：开发者等待 CI 结果的时间缩短，开发迭代更加顺畅。
- 风险标记：仅 CI 配置变更，低风险

关联脉络

- 暂无明显关联 PR