

PR #26108 完整报告

sgl-project/sglang

FutureMap: debug-assert that gather sees a stashed value

合并时间: 2026-05-23 04:10

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/26108>

执行摘要

- 一句话: FutureMap 增加 consume 前值非负断言
- 推荐动作: 值得精读。该 PR 展示了在 GPU 异步执行环境中使用编译断言进行不变量检查的实践, 对理解 FutureMap 的语义契约和 CI 诊断能力有帮助。

功能与动机

FutureMap 的 `output_tokens_buf` 存在 WRWR 不变量被破坏的可能 (如两次连续写入导致的 WW 模式), 之前静默容错会导致第一个写入的值被静默丢弃。该 PR 意图在 CI 中尽早发现此类 bug。

实现拆解

1. 模块级环境变量检查: 在 `overlap_utils.py` 模块顶部读取 `SGLANG_IS_IN_CI` 环境变量, 计算 `_DEBUG_ASSERT` 标志, 决定是否启用断言逻辑。
2. 新增编译态断言函数: 定义 `_assert_nonneg_and_invalidate` 函数, 使用 `@torch.compile(dynamic=True)` 编译, 内部调用 `torch._assert_async` 断言所有 `gather` 值 `>= 0`, 然后将被读取的 `buf` 位置写回 `-1`。一次 kernel launch 完成断言和置零。
3. 初始化与 `gather` 处集成: `__init__` 时根据 `_DEBUG_ASSERT` 将 `output_tokens_buf` 初始化为全 `-1` (否则保持原 `torch.empty`); 在 `resolve_future` 和 `_resolve_spec_extras` 的 `gather` 操作后分别插入断言调用, 将 `gathered` 值 (`batch.input_ids` 或 `draft_input.bonus_tokens`) 与对应的 `indices` 传入 `_assert_nonneg_and_invalidate`。

关键文件:

- `python/sglang/srt/managers/overlap_utils.py` (模块 调度器; 类别 `source`; 类型 `core-logic`; 符号 `_assert_nonneg_and_invalidate`): 唯一修改的文件; 新增了 `_assert_nonneg_and_invalidate` 函数, 修改了 `__init__`、`resolve_future` 和 `_resolve_spec_extras` 以在 CI 中插入断言。

关键符号: `_assert_nonneg_and_invalidate`

关键源码片段

python/sglang/srt/managers/overlap_utils.py

唯一修改的文件；新增了 `_assert_nonneg_and_invalidate` 函数，修改了 `__init__`、`resolve_future` 和 `_resolve_spec_extras` 以在 CI 中插入断言。

```
# 模块顶部，定义编译态断言函数：gather 后检查非负并置 -1
@torch.compile(dynamic=True)
def _assert_nonneg_and_invalidate(
    values: torch.Tensor, buf: torch.Tensor, indices: torch.Tensor
) -> None:
    """Fused: assert all `values >= 0` and scatter -1 into `buf[indices]`.
    Compiled so the reduction + assert + scatter run as one kernel launch."""
    torch._assert_async((values >= 0).all())
    buf[indices] = -1

# FutureMap 的 __init__ 修改：CI 模式下将 output_tokens_buf 初始化为全 -1
self.output_tokens_buf = (
    torch.full((self.req_pool_size,), -1, dtype=torch.int64, device=self.device)
    if _DEBUG_ASSERT
    else torch.empty(
        (self.req_pool_size,), dtype=torch.int64, device=self.device
    )
)

# resolve_future 中，在非 spec 路径的 gather 后插入断言
if self.spec_algo.is_none():
    _resolve_future_token_ids(batch.input_ids, self.output_tokens_buf)
    if _DEBUG_ASSERT:
        _assert_nonneg_and_invalidate(
            batch.input_ids, self.output_tokens_buf, batch.req_pool_indices
        )
else:
    self._resolve_spec_extras(batch)

# _resolve_spec_extras 中，在 gather bonus_tokens 后插入断言
draft_input.bonus_tokens = self.output_tokens_buf[indices]
if _DEBUG_ASSERT:
    _assert_nonneg_and_invalidate(
        draft_input.bonus_tokens, self.output_tokens_buf, indices
    )
```

评论区精华

无 review 讨论。PR 作者独自提交，无外部评论。

- 暂无高价值评论线程

风险与影响

- 风险：

- 性能风险：断言仅在 CI 启用，生产环境无影响。torch._assert_async 在 GPU 上异步，不阻塞流水线，影响可忽略。
- 回归风险：断言仅在 gather 后执行，不影响逻辑行为。但若 CI 中先前未发现的 WWRR 模式因断言暴露，可能导致 CI 失败，需配合修复 upstream 或其他 PR。
- 兼容性风险：torch._assert_async 需要较新 PyTorch 版本，但仓库已广泛使用，风险低。
- 影响：
 - 影响范围：仅影响 overlap_utils.py 一个文件。
 - 影响程度：中等。修复了静默数据错误暴露的渠道，帮助功能开发者早期发现跨迭代 relay 逻辑的 bug。对普通用户无感知，仅在 CI 中增加早期故障检测。
 - 风险标记：核心路径变更，缺少测试覆盖

关联脉络

- PR #26085 drop FutureIndices wrapper class: 修改了同一文件 overlap_utils.py 中的 FutureMap 类，涉及 indices 的使用方式，与本 PR 的断言逻辑密切相关。