

PR #21392 完整报告

sgl-project/sglang

Refactor: decouple segment tracking from comm registration

合并时间: 2026-05-06 17:07

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/21392>

执行摘要

- 一句话: 解耦 NCCL 注册与分配, 支持跨组内存复用
- 推荐动作: 值得精读, 尤其关注 C++ 跟踪 + Python 延迟注册的边界设计, 以及如何在避免数据复制的同时保持正确性。适合理解对称内存分配器的演进。

功能与动机

When multiple communication groups share a single global MemPool, memory blocks released by one group's comm may be reused by another group's comm. However, symmetric memory requires buffers to be registered with a specific ncclComm via ncclCommWindowRegister. Reusing memory across groups causes the registration to be associated with the wrong communicator. So redesign symmetric memory allocator to defer NCCL window registration from allocation-time to context exit-time. (from PR body)

实现拆解

1. C++ 层段跟踪: 在 `pynccl_allocator.py` 的内嵌 C++ 代码中, 新增 `Segment` 结构体, 全局 `g_segments` 向量和 `g_comm_registration_index` 映射。`nccl_alloc_plug` 调用 `track_segment` 只插入段信息, 不再读取环境变量调用 `ncclCommWindowRegister`。`nccl_free_plug` 在内存池销毁时清空所有跟踪数据 (`g_segments` 和 `g_comm_registration_index`)。新增 `nccl_allocator_register_segments_with_comm(comm_ptr)` 函数, 通过 `g_comm_registration_index` 记录每个 communicator 已注册的索引, 仅注册新添加的段。
2. Python 层延迟注册: `SymmetricMemoryContext` 现在所有组共享一个 MemPool (`_shared_mem_pool`)。`__init__` 不再通过环境变量传递 comm 指针, 而是保存当前 comm 句柄。`__exit__` 时调用新增的 `_register_segments_for_comm` 方法, 通过 ctypes 调用 C++ 函数 `nccl_allocator_register_segments_with_comm`。这确保即使是池中重用的内存, 只要未在当前 communicator 上注册过, 就会被正确注册。
3. 新增 benchmark 脚本: `benchmark/bench_pynccl_allocator/bench_segment_tracking.py` 对比两种段跟踪方法的 CPU 时间: 新的 C++ 跟踪与旧 PyTorch 的 `mem_pool.snapshot()`。结果显示新方法快约 25 倍。同时移除了旧的环境变量 `SGLANG_TMP_NCCL_COMM_VALUE` 传递机制和分配时注册路径。

关键文件:

- `python/sglang/srt/distributed/device_communicators/pynccl_allocator.py` (模块 对称内存; 类别 `source`; 类型 `core-logic`; 符号 `get_nccl_mem_pool`, `_register_segments_for_comm`): 核心文件, 将 NCCL 窗口注册从分配时移至上下文退出时, 并实现 C++ 段跟踪与索引注册机制
- `benchmark/bench_pynccl_allocator/bench_segment_tracking.py` (模块 性能基准; 类别 `source`; 类型 `benchmark`; 符号 `setup_segments`, `bench_register_segments_with_comm`, `bench_mempool_snapshot`, `bench_with_various_segment_counts`): 新增 `benchmark`, 量化新跟踪方法的 CPU 优势, 验证设计

关键符号: `nccl_alloc_plug`, `nccl_allocator_register_segments_with_comm`, `get_nccl_mem_pool`, `_register_segments_for_comm`, `SymmetricMemoryContext.exit`

关键源码片段

`python/sglang/srt/distributed/device_communicators/pynccl_allocator.py`

核心文件, 将 NCCL 窗口注册从分配时移至上下文退出时, 并实现 C++ 段跟踪与索引注册机制

```
// C++ 源码片段: 段跟踪与索引注册
// 全局段列表以插入顺序维护 (FIFO)
static std::vector<Segment> g_segments;
static std::mutex g_segment_mutex;

// 每个 communicator 的下一个待注册索引
static std::unordered_map<uintptr_t, size_t> g_comm_registration_index;

// 分配时只跟踪段, 不注册
void* nccl_alloc_plug(size_t size, int device, void* stream) {
    void* ptr;
    NCCLCHECK(ncclMemAlloc(&ptr, size));
    // 仅记录段信息, 推迟注册到上下文退出
    track_segment(ptr, size);
    return ptr;
}

// 注册所有未注册段到指定 communicator
int nccl_allocator_register_segments_with_comm(uintptr_t comm_ptr) {
    std::lock_guard<std::mutex> lock(g_segment_mutex);
    // 获取该 comm 的下一个待注册索引 (默认 0)
    size_t& start_idx = g_comm_registration_index[comm_ptr];
    // 仅注册自上次以来新增的段
    for (size_t i = start_idx; i < g_segments.size(); ++i) {
        auto& seg = g_segments[i];
        ncclComm_t comm = (ncclComm_t)(comm_ptr);
        ncclWindow_t win;
```

```

    NCCLCHECK(ncclCommWindowRegister(
        comm, seg.ptr, seg.size, &win, NCCL_WIN_COLL_SYMMETRIC));
}
start_idx = g_segments.size();
return 0;
}

# Python 端: __exit__ 时延迟注册
class SymmetricMemoryContext:
    def __enter__(self):
        # ... 初始化共享 MemPool, 分配内存
        return self

    def __exit__(self, exc_type, exc_val, exc_tb):
        if not self.exited:
            # 将跟踪段注册到当前 communicator
            self._register_segments_for_comm()
            self.exited = True

    def _register_segments_for_comm(self):
        """
        通过 ctypes 调用 C++ 注册函数, 只注册本 comm 的未注册段
        """
        global _register_func
        # _register_func 在模块加载时已绑定 nccl_allocator_register_segments_with_comm
        ret = _register_func(self._comm_ptr)
        assert ret == 0, f"Register segments failed with error code {ret}"

```

评论区精华

1. 将段注册逻辑完全移到 C++ 内部: nvcastet 建议不要将段列表从 C++ 复制到 Python, 而是让 C++ 内部维护注册索引, 避免数据复制和同步开销。作者采纳, 实现了 `g_comm_registration_index` 映射。
2. `_ptr_to_registered_comms` 潜在泄漏: gemini-code-assist[bot] 指出字典只添加不删除。作者解释 PyTorch mempool 在池存活时不会释放内存给分配器, 因此字典不会无限增长。
3. `untrack_segment` 线性扫描性能: gemini-code-assist[bot] 指出使用线性扫描 $O(N)$ 。作者回应 N 不大, 无需复杂数据结构; nvcastet 也认为当前假设下不会频繁调用 `untrack`。
4. 验证注册返回码: nvcastet 要求 C++ 注册函数返回非零时 `assert`。作者添加了断言。
 - 将段注册逻辑完全移到 C++ 内部 (design): 作者采纳, 实现了 `g_comm_registration_index` 映射。
 - `_ptr_to_registered_comms` 内存泄漏 (correctness): 作者解释 PyTorch mempool 在池存活时不会释放内存给分配器, 因此字典不会无限增长。
 - `untrack_segment` 线性扫描性能 (performance): 作者回应 N 不大, 无需复杂数据结构。nvcastet 也认为当前假设下不会频繁调用 `untrack`。
 - 验证注册返回码 (correctness): 作者添加了 `assert` 判断。

风险与影响

- 风险：
 - 假设局限：C++ 端假设内存池销毁前不会有单个段释放，若未来支持部分释放，跟踪逻辑将需要调整。
 - 注册顺序依赖：C++ 层依赖 `g_segments` 的 FIFO 顺序，多线程并发分配可能导致索引漂移，但当前分配是单线程。
 - 兼容性：仅测试了 FP8 模型，FP4（如 DeepSeek-R1 FP4）未验证，可能存在问题。
 - 缺少单元测试：虽有 benchmark，但无针对新逻辑的单元测试，回归风险依赖集成测试。
- 影响：
 - 用户：无直接接口变化，但对称内存性能提升，TP8 场景 e2e 提升 6.85%，DP8 提升 2%。
 - 系统：减少 CPU 开销（移除 snapshot），降低内存碎片（共享 MemPool），提升跨组内存复用正确性。
 - 团队：代码结构更清晰，注册逻辑集中在 C++ 内部，易于维护。
 - 风险标记：核心路径变更，假设依赖，缺少单元测试，FP4 未验证

关联脉络

- PR #19329 Related PR: PR body 引用 #19329 issue comment 提供上下文。
- PR #20153 Related PR: PR body 引用 #20153，可能也是相关修改。