

PR #25173 完整报告

sgl-project/sglang

Refactor NIXL hicache. Add O_DIRECT support

合并时间: 2026-06-01 23:28

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/25173>

执行摘要

- 一句话: 重写 NIXL HiCache, 增加 O_DIRECT 支持及 mmap 分配器
- 推荐动作: 该 PR 值得仔细审查, 特别是 `nixl_registry.py` 中的上下文管理器设计模式和 `mmap_allocator.py` 中 `HugePage` 的支持实现。建议团队阅读其异常处理策略和对齐检查逻辑, 可作为高可靠 I/O 后端的参考。合入前需确认 AMD/NPU 环境 CI 通过, 不过当前 CI 结果均为绿色 (除不相关的失败)。

功能与动机

PR 描述指出: 'This is a rewrite of the HiCache NIXL connector. It fixes a number of bugs (thread synchronization, object store). It also adds an option of using O_DIRECT for file handles, which can improve I/O performance significantly.' 根本动机是修复现有实现中的线程安全性和对象存储后端缺陷, 同时通过 O_DIRECT 释放文件后端的 I/O 性能潜力。PR body 中提供了详细的性能测试数据, 证明 O_DIRECT 可大幅降低 PD 分解场景的 TTFT。

实现拆解

1. 重构注册管理器 (NixlRegistry): 新建 `nixl_registry.py`, 将原来分散在 `nixl_utils.py` 和 `hicache_nixl.py` 中的 NIXL 内存注册 / 注销逻辑统一为上下文管理器 `NixlRegistry`。提供 `storage()` 方法, 在进入时依次打开文件 (FILE 后端)、注册内存、构建 `xfer_descs`, 退出时自动 `deregister` 和关闭 `fd`, 确保异常安全。同时将设备 ID 分配线程化以避免 OBJ 后端的竞态。
2. 预注册主机内存: 在 `register_mem_pool_host()` 中对主机端 `kv_buffer` 或 `bounce buffer` 进行一次 NIXL 注册, 替代每次传输时重复注册。零拷贝模式注册整个 `kv_buffer`; 非零拷贝模式为 `set/get` 各注册一个 `bounce buffer`。这大幅减少了每笔传输的附加开销。
3. 新增 O_DIRECT 支持与对齐校验: 通过 `SGLANG_HICACHE_NIXL_USE_DIRECT_IO` 环境变量或配置中的 `use_direct_io` 键开启 O_DIRECT (默认开启)。当 O_DIRECT 激活且使用文件后端时, `HiCacheNixl` 标记 `needs_page_alignment=True`。在注册 `host` 内存时, 若为零拷贝模式则通过 `HostKVCache.is_stride_page_aligned()` 检查每个 `page` 指针是否 OS 页对齐; 若不对齐则回退到 `bounce buffer` 拷贝模式。对齐要求硬编码为 4KB (4096) 作为安全值。
4. 实现 mmap 后端的宿主分配器: 新增 `mmap_allocator.py`, 通过匿名 `mmap` (`MAP_SHARED | MAP_ANONYMOUS | MAP_POPULATE`) 分配主机内存, 保证基地址 OS 页对齐。可选支持 `HugePage` (`SGLANG_HUGEPAGE_SIZE=2MB|1GB`), 通过

`ctypes` 调用 `libc.mmap` 额外传递 `MAP_HUGETLB` 等标志。`HostTensorAllocator` 使用 `alloc_mmap` 替代 `torch.empty()`，使得 `kv_buffer.data_ptr()` 天然满足 `O_DIRECT` 对齐要求。

5. 重构测试套件：将原来的测试文件 `python/sglang/srt/mem_cache/storage/nixl/test_hicache_nixl_storage.py` 删除，替换为位于 `test/registered/unit/mem_cache/test_hicache_nixl_storage.py` 的新测试。新测试使用 `MockMemPoolHost` 模拟主机缓存池，覆盖零拷贝与非零拷贝布局、`O_DIRECT` 对齐与回退、`NIXL` API 正确性（未注册内存进行传输应该失败）、多线程并发、以及对象存储集成测试（需要 `MinIO` 二进制，但 `CI` 中会跳过）。同时注册了 `CI` 运行。

关键文件：

- `python/sglang/srt/mem_cache/storage/nixl/nixl_registry.py`（模块 `NIXL` 注册器；类别 `source`；类型 `dependency-wiring`；符号 `_buffer_sizes`, `NixlRegistry`, `init`, `_open_files`）：核心新增文件，定义了 `NixlRegistry` 上下文管理器，负责 `NIXL` 存储端注册 / 注销的生命周期管理。替换了原有的分散注册逻辑，是本次重构的核心。
- `python/sglang/srt/mem_cache/storage/nixl/hicache_nixl.py`（模块 `NIXL` 存储后端；类别 `source`；类型 `dependency-wiring`；符号 `register_buffers`, `register_files`, `register_objects`, `_execute_transfer`）：主修改文件，集成 `NixlRegistry`、预注册主机内存、`O_DIRECT` 对齐检查，是 `NIXL` 后端对外接口。
- `python/sglang/srt/mem_cache/storage/nixl/nixl_utils.py`（模块 `NIXL` 工具函数；类别 `source`；类型 `dependency-wiring`；符号 `get_use_direct_io`, `set_bucket`, `NixlRegistration`, `init`）：移除 `NixlRegistration` 类，将配置类简化并增强 `get_use_direct_io` 方法，减小工具模块责任。
- `python/sglang/srt/mem_cache/mmap_allocator.py`（模块 内存分配器；类别 `source`；类型 `dependency-wiring`；符号 `_alloc_hugepage`, `alloc_mmap`）：新增文件，实现基于 `mmap` 的主机内存分配器，确保页对齐以支持 `O_DIRECT` 零拷贝，并可选大页。是 `O_DIRECT` 功能的基础设施。
- `test/registered/unit/mem_cache/test_hicache_nixl_storage.py`（模块 测试套件；类别 `test`；类型 `test-coverage`；符号 `MockMemPoolHost`, `get_page_buffer_meta`, `get_dummy_flat_data_page`, `get_data_page`）：新增综合测试套件，覆盖零拷贝 / 非零拷贝、`O_DIRECT` 对齐、多线程和 `OBJ` 后端，提升 `CI` 中的测试覆盖。
- `python/sglang/srt/mem_cache/memory_pool_host.py`（模块 主机缓存池；类别 `source`；类型 `dependency-wiring`；符号 `HostTensorAllocator`, `is_stride_page_aligned`）：为支持 `O_DIRECT` 对齐检查，新增 `is_stride_page_aligned` 方法；`HostTensorAllocator` 改用 `mmap` 分配器。

关键符号：`NixlRegistry.storage`, `NixlRegistry._registered`, `NixlRegistry._open_files`, `HiCacheNixl.init`, `HiCacheNixl.register_mem_pool_host`, `HiCacheNixl._execute_transfer`, `HiCacheNixl._xfer_and_wait`, `HiCacheNixl._xfer_pre_registered`, `alloc_mmap`, `_alloc_hugepage`, `HostKVCache.is_stride_page_aligned`, `NixlBackendConfig.get_use_direct_io`

关键源码片段

python/sclang/srt/mem_cache/storage/nixl/nixl_registry.py

核心新增文件，定义了 NixlRegistry 上下文管理器，负责 NIXL 存储端注册 / 注销的生命周期管理。替换了原有的分散注册逻辑，是本次重构的核心。

```
"""NixlRegistry 上下文管理器，封装 NIXL 存储端的注册/注销生命周期。"""
```

```
class NixlRegistry:
    def __init__(self, agent, mem_type: str, file_manager=None):
        self.agent = agent
        self.mem_type = mem_type
        self.file_manager = file_manager
        # OBJ 插件使用进程内唯一 devId 映射对象，并发时必须加锁
        self._obj_devid_lock = threading.Lock()
        self._obj_devid_next = 1

    @contextmanager
    def _open_files(self, paths, create: bool):
        """打开所有指定路径的文件，yield fd 列表；退出时统一关闭。"""
        fds = []
        try:
            for path in paths:
                fd = self.file_manager.open_file(path, create=create)
                if fd is None: # 任一打开失败则整体失败
                    yield None
                    return
                fds.append(fd)
            yield fds
        finally:
            for fd in fds:
                self.file_manager.close_file(fd)

    @contextmanager
    def _registered(self, items, mem_type: str):
        """注册 items 到 NIXL，yield 注册句柄；退出时自动注销。"""
        reg = None
        if items:
            reg_descs = self.agent.get_reg_descs(items, mem_type)
            if reg_descs is not None:
                try:
                    reg = self.agent.register_memory(reg_descs)
                except Exception as e:
                    logger.error(f"Failed to register memory of type {mem_type}: {e}")
        try:
            yield reg
        finally:
            if reg is not None:
                try:
                    self.agent.deregister_memory(reg)
```

```

except Exception as e:
    logger.debug("deregister_memory skipped: %s", e)

@contextmanager
def storage(self, buffers, keys, direction):
    """完整流程: 打开文件 -> 注册 -> yield xfer_descs -> 注销 -> 关闭文件。"""
    sizes = _buffer_sizes(buffers)
    if sizes is None:
        yield None
        return

    if self.mem_type == "FILE":
        with self._open_files(keys, create=(direction == "WRITE")) as fds:
            if fds is None:
                yield None
                return
            tuples = [(0, sizes[i], fds[i], keys[i]) for i in range(len(keys))]
            with self._registered(tuples, "FILE") as reg:
                if reg is None:
                    yield None
                    return
                yield self.agent.get_xfer_descs(
                    [(0, sizes[i], fds[i]) for i in range(len(fds))], "FILE"
                )
    else: # OBJ
        # ... 类似处理
        pass

```

python/sglang/srt/mem_cache/mmap_allocator.py

新增文件，实现基于 mmap 的主机内存分配器，确保页对齐以支持 O_DIRECT 零拷贝，并可选大页。是 O_DIRECT 功能的基础设施。

```

"""mmap 后端主机分配器，保证页对齐，可选 HugePage 支持。"""

import ctypes, ctypes.util, mmap, os, weakref
import torch

_libc_name = ctypes.util.find_library("c") or "libc.so.6"
_libc = ctypes.CDLL(_libc_name, use_errno=True)
# ... 设置 mmap/munmap 签名
_MAP_POPULATE = getattr(mmap, "MAP_POPULATE", 0x08000)
_MAP_HUGETLB = 0x40000
_MAP_HUGE_2MB = 21 << 26
_MAP_HUGE_1GB = 30 << 26

def alloc_mmap(dims: tuple, dtype: torch.dtype) -> torch.Tensor:
    """分配一个 mmap 后端的 tensor，基地址 OS 页对齐。"""
    n_bytes = math.prod(dims) * torch.empty([], dtype=dtype).element_size()
    hugepage = (envs.SGLANG_HUGEPAGE_SIZE.get() or "").strip().upper()

```

```

if hugepage == "":
    page_size, extra_flags = mmap.PAGESIZE, 0
elif hugepage == "2MB":
    page_size, extra_flags = 2 * 1024 * 1024, _MAP_HUGETLB | _MAP_HUGE_2MB
elif hugepage == "1GB":
    page_size, extra_flags = 1024 * 1024 * 1024, _MAP_HUGETLB | _MAP_HUGE_1GB
else:
    # 未知值, 回退普通页
    page_size, extra_flags = mmap.PAGESIZE, 0

alloc_bytes = math.ceil(n_bytes / page_size) * page_size
if extra_flags and _libc:
    try:
        ptr = _libc.mmap(None, alloc_bytes,
                        mmap.PROT_READ | mmap.PROT_WRITE,
                        mmap.MAP_SHARED | mmap.MAP_ANONYMOUS | _MAP_POPULATE |
                        extra_flags,
                        -1, 0)
        if ptr == ctypes.c_void_p(-1).value:
            raise OSError
        array = (ctypes.c_uint8 * n_bytes).from_address(ptr)
        weakref.finalize(array, _libc.munmap, ptr, alloc_bytes)
        return torch.frombuffer(array, dtype=dtype, count=math.prod(dims)).reshape(dims)
    except OSError:
        # 大页失败, 降级为普通 mmap (alloc_bytes 重新计算)
# 普通 mmap 路径
mm = mmap.mmap(-1, alloc_bytes, flags=mmap.MAP_SHARED | mmap.MAP_ANONYMOUS |
MAP_POPULATE,
                prot=mmap.PROT_READ | mmap.PROT_WRITE)
return torch.frombuffer(mm, dtype=dtype, count=math.prod(dims)).reshape(dims)

```

评论区精华

1. MinIO 容器依赖: ishandhanani 和 xiezhq-hermann 认为不应该将 MinIO 加入生产容器, 作者回应后将其移除, 仅保留测试中可选执行(跳过)。
 2. libc 可移植性: gemini-code-assist 指出硬编码 libc.so.6 不兼容 musl, 建议使用 ctypes.util.find_library('c'), 作者采纳并修改代码。
 3. 传输异常处理: AI 审核指出 _xfer_and_wait 中若 agent.transfer 抛出异常, release_xfer_handle 不会被调用; 作者说明当前行为可接受, 未来计划使用 NIXL 异步接口改进。
 4. 硬编码 4096: ishandhanani 询问对齐检查中的 4096 含义, 作者解释为 O_DIRECT 的最小块大小安全值, 并补充了注释。
 5. 废弃方法处理: ishandhanani 建议对不再使用的旧方法(如旧版 register_buffers 的某些调用)抛出异常而非保留空实现, 作者最终采纳并抛出 NotImplementedError。
- MinIO 二进制是否应加入 Docker 镜像 (design): 作者移除了 Dockerfile 中 MinIO 的安装, 对象存储测试在缺少 MinIO 时跳过。

- libc 路径硬编码的可移植性 (correctness): 作者采纳建议, 修改为动态查找。
- `_xfer_and_wait` 异常路径资源泄漏 (correctness): 未修改代码, 但添加注释说明 (如果后续加入了注释); 当前保持原样。
- 对齐检查中硬编码 4096 的含义 (question): 作者添加了注释解释该值。
- 废弃方法是否应直接抛出异常 (design): 作者改为抛出 `NotImplementedError`。

风险与影响

- 风险:
 1. `O_DIRECT` 对齐风险: 若文件系统 / 硬件要求更大对齐尺寸 (如某些 NVMe 512 字节也可), 4KB 硬编码可能不够, 但业界普遍安全。若用户遇到对齐问题, 代码会 fallback 到 bounce buffer 而非失败, 但性能下降。
 2. mmap 分配器替换 `torch.empty()` 的影响: `HostTensorAllocator` 现在使用 mmap 而非 PyTorch 默认分配器。这改变了内存的 NUMA 策略和 fault-in 行为。`MAP_POPULATE` 立即分配物理页, 可能增加启动时间。CUDA 注册要求物理页固定, 而 mmap 的 `MAP_SHARED` 行为可能与传统方式不同 (但代码已注释要求)。整体迁移经过单元测试验证, 但若与其他内存模式 (如 IOMMU 或大页交错) 交互可能产生意外。
 3. 线程安全: 新代码引入了 `_obj_devid_lock` 处理 OBJ 后端的并发注册, 但未覆盖所有可能竞态 (如 NIXL agent 自身的线程安全)。测试中包含压力测试 (默认不运行), 未在 CI 中持续验证。
 4. 测试覆盖限制: 所有测试仅在 NVIDIA GPU 上通过, 未对 AMD、Intel 或 NPU 硬件验证, 该后端可能在这些平台上引入回归。
 5. 配置更改: 新增环境变量 `SGLANG_HICACHE_NIXL_USE_DIRECT_IO` 默认开启, 可能改变现有用户的默认行为 (如果之前未使用 `O_DIRECT`), 但预期性能提升是正向的。
 - 影响: 用户: 使用 NIXL 后端的用户将自动获得 `O_DIRECT` 带来的 I/O 性能提升 (特别是在 SSD RAID 上), 无需手动配置; 但需确保文件系统支持 `O_DIRECT` (Linux 默认支持)。使用 OBJ 后端的用户不受影响。系统: 新代码要求 Linux 内核 $\geq 2.6.23$ (用于 `MAP_POPULATE`), 大页需要 HugeTLB 支持。CI 中新增了单元测试, 但压力测试和 MinIO 测试默认跳过, 未显著增加 CI 时长。团队: 代码结构更清晰, 上下文管理器模式降低了资源泄漏风险。新分配的 mmap 路径与 PyTorch 内存管理解耦, 便于未来专用优化。
- 风险标记: 核心路径变更, 底层分配器替换, 缺少非 NVIDIA 测试, `O_DIRECT` 对齐风险, 默认行为变化

关联脉络

- PR #26615 [sgl] Window-aware LRU refresh for SWA prefix cache in unified cache: 同为 HiCache 相关的性能优化 PR, 修改了统一缓存的 LRU 策略。虽不直接重叠代码, 但属同一条功能线 (HiCache 性能改进)。
- PR #26883 [PP][Bugfix] Handle input_ids assignment in prepare_for_extend: 修改了 `cache_controller` 文件 (与本 PR 有交集), 但属于不同模块。