

# PR #40020 完整报告

vllm-project/vllm

[kv\_offload] Add multi-tier KV cache offloading framework

合并时间: 2026-05-13 22:21

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/40020>

## 执行摘要

- 一句话: 新增多级 KV 缓存卸载框架, 支持链式二级存储 / 网络
- 推荐动作: 值得精读, 特别是抽象接口设计和异步批处理模式。可关注层次化管理器的错误处理和生命周期管理。

## 功能与动机

现有 CPU 单级卸载容量有限, 无法满足大规模 KV 缓存或跨节点共享需求。本 PR 实现 RFC #38260 提出的 Multi-tier KV offloading 方案, 将卸载堆栈扩展为可链接的二级层 (如存储、网络), 并保持对上层模块透明的接口。

## 实现拆解

1. 抽象接口层: 在 base.py 中定义 SecondaryTierManager ABC, 包含异步 lookup/submit\_store/submit\_load/get\_finished 接口, 以及 JobMetadata/JobResult 数据类型。
2. CPU 主层包装: manager.py 中 CPUPrimaryTierOffloadingManager 继承 CPUOffloadingManager, 暴露 prepare\_read/complete\_read 等别名, 并提供 get\_kv\_memoryview() 零拷贝视图。
3. 多级编排器: TieringOffloadingManager 实现级联存储 (store 到所有次级)、分阶段提升 (load 从次级到主层后再给 GPU)、引用计数保护、get\_finished 限频 (每引擎步一次) 及 load 批处理。
4. 示例实现: example/\_\_init\_\_.py 中的 ExampleSecondaryTier 演示 LRU 驱逐和异步模拟。
5. 配置与工厂: spec.py 读取 secondary\_tiers 配置, factory.py 通过 get\_tier\_type() 静态方法创建对应类型实例。
6. 配套修改: SharedOffloadRegion 新增 create\_kv\_memoryview(); CPUOffloadingSpec.\_create\_handlers 扩展分支以兼容 TieringOffloadingSpec。
7. 测试: test\_tiering\_offloading.py 共 16 个测试用例覆盖基础、级联、提升、部分查找、无次级层等场景。

关键文件:

- vllm/v1/kv\_offload/tiering/manager.py (模块 卸载编排器; 类别 source; 类型 core-logic ; 符号 PendingPromotion, CPUPrimaryTierOffloadingManager, init,

get\_kv\_memoryview) : 核心编排器, 包含 CPUPrimaryTierOffloadingManager 和 TieringOffloadingManager, 实现级联存储、提升、引用计数保护等核心逻辑。

- vllm/v1/kv\_offload/tiering/example/\_\_init\_\_.py (模块 二次层示例; 类别 source; 类型 core-logic; 符号 \_JobMetadata, ExampleSecondaryTier, init, lookup) : 示例二级层实现, 演示如何实现 SecondaryTierManager 接口, 包含 LRU 驱逐和异步模拟。
- vllm/v1/kv\_offload/tiering/spec.py (模块 卸载规范; 类别 source; 类型 dependency-wiring; 符号 TieringOffloadingSpec, init, get\_manager, \_create\_handlers) : 入口 Spec, 负责解析配置并组装 TieringOffloadingManager 堆栈。
- vllm/v1/kv\_offload/tiering/base.py (模块 抽象接口; 类别 source; 类型 core-logic; 符号 JobMetadata, JobResult, SecondaryTierManager, init) : 定义了 SecondaryTierManager 抽象接口和 JobMetadata/JobResult 数据类型, 是框架扩展的契约。
- tests/v1/kv\_offload/test\_tiering\_offloading.py (模块 测试; 类别 test; 类型 test-coverage; 符号 \_mock\_mmap\_region, to\_keys, count\_hits, TestExampleSecondaryTier) : 涵盖 ExampleSecondaryTier 和 TieringOffloadingManager 的单元测试, 16 个测试用例验证核心行为。

关键符号: TieringOffloadingManager.prepare\_store, TieringOffloadingManager.lookup, TieringOffloadingManager.prepare\_load, CPUPrimaryTierOffloadingManager.get\_kv\_memoryview, ExampleSecondaryTier.submit\_store, ExampleSecondaryTier.submit\_load, TieringOffloadingSpec.get\_manager, create\_secondary\_tier

## 关键源码片段

### vllm/v1/kv\_offload/tiering/manager.py

核心编排器, 包含 CPUPrimaryTierOffloadingManager 和 TieringOffloadingManager, 实现级联存储、提升、引用计数保护等核心逻辑。

```
# 包装 CPUOffloadingManager 并添加读写别名, 方便二级层调用
class CPUPrimaryTierOffloadingManager(CPUOffloadingManager):
    def __init__(
        self,
        num_blocks: int,
        mmap_region: SharedOffloadRegion,
        cache_policy: str = 'lru',
        enable_events: bool = False,
    ):
        super().__init__(
            num_blocks=num_blocks,
            cache_policy=cache_policy,
            enable_events=enable_events,
        )
        self.mmap_region = mmap_region
        # 别名: read/write 用于 CPU <-> secondary, load/store 用于 CPU <-> GPU
        self.prepare_read = self.prepare_load
        self.complete_read = self.complete_load
```

```

self.prepare_write = self.prepare_store
self.complete_write = self.complete_store
# 持有一个内存视图供二级层零拷贝访问 CPU 缓存
self._kv_memoryview = mmap_region.create_kv_memoryview()

def get_kv_memoryview(self) -> memoryview:
    """返回形状为 (num_blocks, row_stride_bytes) 的 memoryview"""
    return self._kv_memoryview

def shutdown(self) -> None:
    super().shutdown()
    self._kv_memoryview.release()
    self._mmap_region.cleanup()

```

### vllm/v1/kv\_offload/tiering/example/\_\_init\_\_.py

示例二级层实现，演示如何实现 SecondaryTierManager 接口，包含 LRU 驱逐和异步模拟。

# 纯内存实现，用于测试和作为实现参考

```

class ExampleSecondaryTier(SecondaryTierManager):
    def __init__(
        self,
        vllm_config: 'VllmConfig',
        primary_kv_view: memoryview,
        max_blocks: int = 1000,
        simulate_async: bool = False,
    ):
        super().__init__(vllm_config, primary_kv_view)
        self.max_blocks = max_blocks
        self.simulate_async = simulate_async
        self.blocks: OrderedDict[OffloadKey, bool] = OrderedDict()
        self.completed_jobs: list[JobResult] = []
        self.pending_jobs: list[_JobMetadata] = []

    def submit_store(self, job_metadata: JobMetadata) -> None:
        job_id = job_metadata.job_id
        keys = job_metadata.keys
        block_ids = job_metadata.block_ids
        assert len(keys) == len(block_ids), 'keys / block_ids 长度不匹配'
        # 过滤已存在的块
        blocks_to_store = [k for k in keys if k not in self.blocks]
        if not blocks_to_store:
            return
        # LRU 驱逐
        num_evict = len(blocks_to_store) - (self.max_blocks - len(self.blocks))
        if num_evict > 0:
            protected = set(keys)
            evicted = []
            for key in self.blocks:
                if key not in protected:

```

```

        evicted.append(key)
        if len(evicted) == num_evict:
            break
    else:
        return # 无法腾出足够空间
    for key in evicted:
        del self.blocks[key]
# 记录传输作业
internal = _JobMetadata(job_id=job_id, keys=blocks_to_store, is_store=True)
if self.simulate_async:
    self.pending_jobs.append(internal)
else:
    self._complete_store_job(internal)

def lookup(self, key: OffloadKey, req_context: ReqContext) -> bool | None:
    return key in self.blocks

```

## vllm/v1/kv\_offload/tiering/spec.py

入口 Spec，负责解析配置并组装 TieringOffloadingManager 堆栈。

```

# 读取配置，创建主层和所有二级层
class TieringOffloadingSpec(CPUOffloadingSpec):
    def get_manager(self) -> OffloadingManager:
        if not self._manager:
            # 创建 scheduler 端 mmap
            world_size = self.vllm_config.parallel_config.world_size
            scheduler_mmap = SharedOffloadRegion(
                instance_id=self.vllm_config.instance_id,
                total_size_bytes=self.cpu_page_size_per_worker * world_size * self.num_blocks,
                num_blocks=self.num_blocks,
                rank=None,
                num_workers=world_size,
                cpu_page_size=self.cpu_page_size_per_worker,
            )
            self._scheduler_mmap = scheduler_mmap
            # 创建 CPU 主层
            primary_tier = CPUPrimaryTierOffloadingManager(
                num_blocks=self.num_blocks,
                cache_policy=self.eviction_policy,
                enable_events=enable_events,
                mmap_region=scheduler_mmap,
            )
            # 为所有二级层提供同一份内存视图（零拷贝共享）
            primary_kv_view = primary_tier.get_kv_memoryview()
            secondary_tiers = []
            for cfg in self.secondary_tier_configs:
                tier = create_secondary_tier(cfg, primary_kv_view, self.vllm_config)
                secondary_tiers.append(tier)
            self._manager = TieringOffloadingManager(

```

```
        primary_tier=primary_tier,  
        secondary_tiers=secondary_tiers,  
        enable_events=enable_events,  
    )  
    return self._manager
```

## 评论区精华

- @orozery 建议 `get_finished` 限制到每引擎步一次，作者在 `take_events` 中实现去重。
- @liranschour 和 @orozery 认为单块 load 效率低，作者引入 `PendingPromotion` 按 (tier, req) 分组延迟批量提交。
- @gemini-code-assist 指出 `memoryview` 可能因底层张量重分配不安全，作者添加 `assert` 确认零拷贝并计划后续跟踪。
- @orozery 建议移除 `store_threshold` 并用静态 `get_tier_type()` 替代 `tier_name`，作者采纳。
- @orozery 最终批准，非关键问题留后续 PR。
- 文件结构重组为 `tiering/` 子目录 (design): 按建议重组，将文件从 `vllm/v1/kv_offload/abstract.py` 等地迁移到 `tiering/`。
- `get_finished` 调用频率优化 (performance): 在 `take_events` 中批量处理已完成任务，减少轮询频率。
- 批处理 load 请求以提升性能 (performance): 引入 `PendingPromotion`，在引擎步末批量提交 load 请求。
- 内存视图零拷贝安全性 (correctness): 添加 `assert` 确认 `np` 和 `torch` 共享存储，留待后续彻底解决。
- 移除 `store_threshold` 和 `tier_name` 配置 (design): 移除 `store_threshold`，添加 `get_tier_type` 静态方法，配置不再需要 `tier_name`。

## 风险与影响

- 风险:
  - 零拷贝安全: `memoryview` 引用的缓冲区需严格生命周期管理，虽有 `assert` 但生产环境需额外保障。
  - 性能风险: 异步轮询和引用计数增加 `Scheduler` 负担，但限频和批处理已缓解。
  - 配置兼容性: `secondary_tiers` 格式错误将引发异常，默认关闭无影响。
  - 测试覆盖: 单元测试全面，但缺少真实二级后端（如存储）的端到端验证。
- 影响:
  - 用户: 默认行为不变，仅当配置 `secondary_tiers` 时激活多层卸载; `ExampleSecondaryTier` 可用于快速体验。
  - 系统: 新增约 1800 行代码，修改 6 个现有文件但均为轻量兼容分支。
  - 团队: 为开发真实二级后端（存储、网络）奠定架构基础，需维护抽象接口一致性。
  - 风险标记: 零拷贝内存安全风险，实验性功能，缺少真实二级后端验证

## 关联脉络

- PR #38260 [RFC]: Multi-tier KV offloading via the vLLM offloading connector: 本 PR 直接实现该 RFC 的设计方案。
- PR #33689 [Issue] File structure reorganization for kv\_offload: 文件组织结构参考此 issue 建议。
- PR #41200 [PR] Use Sequence for keys in LookupStoreSpec: 依赖此 PR 的 keys 类型改进, 避免不必要的列表转换。