

# PR #24470 完整报告

sgl-project/sglang

Cache empty MatchResult in RadixCache

合并时间: 2026-05-08 08:13

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/24470>

## 执行摘要

- 一句话: 缓存空 MatchResult 避免重复分配
- 推荐动作: 值得合并, 变更简洁且安全。建议关注 HiRadixCache 场景下 `_empty_match_result.device` 与运行时 `device` 的一致性, 必要时在 HiRadixCache 中 override `reset()` 确保 tensor 创建在正确 device 上。

## 功能与动机

现有代码中每次 `match_prefix()` 返回空结果时都会构造一个新的 MatchResult 和 `torch.empty` tensor。PR body 指出通过缓存一个实例可避免重复分配, 且 MatchResult 作为不可变 NamedTuple 共享一个实例是安全的。

## 实现拆解

1. `RadixCache.reset()` 中创建缓存实例: 在 `reset()` 末尾添加 `self._empty_match_result = MatchResult(...)` 构造, 包含一个长度 0 的 int64 tensor (创建在对应 device 上) 以及指向 `root_node` 的 `last_device_node/last_host_node`。此缓存实例随 `root_node` 重置而重建, 保证引用正确。
2. `match_prefix()` 中使用缓存: 在 `RadixCache.match_prefix()` 中, 原本的局部函数 `empty_match_result()` 和 `torch.empty(...)` 调用被直接替换为 `self._empty_match_result` 或 `self._empty_match_result.device_indices`。HiRadixCache.match\_prefix() 同理。
3. `init_load_back()` 中使用缓存: HiRadixCache.init\_load\_back() 中 fallback 返回的 `torch.empty(...)` 替换为 `self._empty_match_result.device_indices`。
4. 设备类型容错修复: `RadixCache.__init__` 中原本直接赋值 `self.device = self.token_to_kv_pool_allocator.device`, 但 `allocator` 的 `device` 属性可能是 mock 对象或 str (如 "cuda:0") 而非 `torch.device`。PR 增加类型检查: 若为 str 或 `torch.device` 则包装为 `torch.device`, 否则 fallback 到 CPU。这一修复是独立于缓存的改进。

关键文件:

- `python/sglang/srt/mem_cache/radix_cache.py` (模块 缓存层; 类别 source; 类型 core-logic; 符号 `empty_match_result`): 核心修改: 在 `reset()` 中创建 `_empty_match_result` 并替换 `match_prefix()` 中的空结果构造逻辑; 同时修复 `device` 类型容错。

- python/sglang/srt/mem\_cache/hiradix\_cache.py (模块 缓存层; 类别 source; 类型 core-logic; 符号 empty\_match\_result) : 配套修改: match\_prefix() 和 init\_load\_back() 中使用父类的 \_empty\_match\_result 替换局部 empty tensor 构造。

关键符号: RadixCache.reset, RadixCache.match\_prefix, HiRadixCache.match\_prefix, HiRadixCache.init\_load\_back

## 关键源码片段

### python/sglang/srt/mem\_cache/radix\_cache.py

核心修改: 在 reset() 中创建 \_empty\_match\_result 并替换 match\_prefix() 中的空结果构造逻辑; 同时修复 device 类型容错。

```
# python/sglang/srt/mem_cache/radix_cache.py (head)
class RadixCache(BasePrefixCache):
    def __init__(self, params: CacheInitParams):
        # ... 原有 init 逻辑 ...
        if self.token_to_kv_pool_allocator:
            dev = self.token_to_kv_pool_allocator.device
            # 兼容 mock allocator (非 torch.device) 和 str 类型
            if isinstance(dev, (str, torch.device)):
                self.device = torch.device(dev)
            else:
                self.device = torch.device("cpu")
        else:
            self.device = torch.device("cpu")
        # ...
        self.reset()

    def reset(self):
        # ... 初始化 root_node ...
        self.evictable_size_ = 0
        self.protected_size_ = 0
        self.evictable_leaves.clear()
        # 预创建空 MatchResult 实例, 避免每次 match_prefix 返回空时重复构造
        self._empty_match_result = MatchResult(
            device_indices=torch.empty(
                (0,),
                dtype=torch.int64,
                device=self.device,
            ),
            last_device_node=self.root_node,
            last_host_node=self.root_node,
        )
        self._record_all_cleared_event()

    def match_prefix(self, params: MatchPrefixParams) -> MatchResult:
        key = params.key
        key, _ = key.maybe_to_bigram_view(self.is_eagle)
```

```

if self.disable or len(key) == 0:
    # 直接返回缓存的空结果, 不再构造
    return self._empty_match_result
key = key.page_aligned(self.page_size)
if len(key) == 0:
    return self._empty_match_result
value, last_node = self._match_prefix_helper(self.root_node, key)
if value:
    value = torch.cat(value)
else:
    # 复用缓存 tensor 而非重新分配
    value = self._empty_match_result.device_indices
return MatchResult(
    device_indices=value,
    last_device_node=last_node,
    last_host_node=self.root_node,
)

```

### python/sglang/srt/mem\_cache/hiradix\_cache.py

配套修改: `match_prefix()` 和 `init_load_back()` 中使用父类的 `_empty_match_result` 替换局部 `empty tensor` 构造。

```

# python/sglang/srt/mem_cache/hiradix_cache.py (head)
class HiRadixCache(RadixCache):
    def init_load_back(self, params: InitLoadBackParams):
        last_node = params.last_host_node
        mem_quota = params.mem_quota
        if last_node.evicted:
            # ... 加载逻辑 ...
        # fallback: 复用父类缓存的空设备索引 tensor
        return (
            self._empty_match_result.device_indices,
            last_node,
        )

    def match_prefix(self, params: MatchPrefixParams):
        if self.disable:
            # 直接返回父类缓存的空结果
            return self._empty_match_result
        key = params.key
        key, _ = key.maybe_to_bigram_view(self.is_eagle)
        key = key.page_aligned(self.page_size)
        if len(key) == 0:
            return self._empty_match_result
        value, last_node = self._match_prefix_helper(self.root_node, key)
        if value:
            value = torch.cat(value)
        else:
            value = self._empty_match_result.device_indices

```

```
# ... 后续处理 ...
return MatchResult(
    device_indices=value,
    last_device_node=last_node,
    last_host_node=last_host_node,
    host_hit_length=host_hit_length,
)
```

## 评论区精华

该 PR 无 reviewer 评论。提交历史中 merrymercy 在第二次和第三次提交中修复了 device 类型处理的 bug，说明作者在实现过程中发现了模拟场景下的缺陷并做了针对性修复。

- 暂无高价值评论线程

## 风险与影响

- 风险：

1. 共享 `_empty_match_result` 实例是安全的，因为 `MatchResult` 是 `NamedTuple`（不可变）。但有调用方修改返回的 `device_indices` tensor 或 `last_device_node/node` 状态的风险——虽然 `NamedTuple` 字段不可变，但 tensor 本身可变（如 in-place 操作）。当前代码中所有返回空结果的路径均只读取 `device_indices` 和节点引用，不会修改，风险可控。
2. `HiRadixCache` 的 `_empty_match_result` 从父类 `RadixCache.reset()` 继承而来，其 `device` 是否正确（默认为 CPU）需确认；当前 `HiRadixCache.match_prefix()` 中使用 `self._empty_match_result` 时假设 `device` 兼容，若 `HiRadixCache` 实际 `device` 为 GPU 而父类中 `cached tensor` 在 CPU 上，可能引入设备不一致问题。- 影响：影响范围：仅影响 `RadixCache` 和 `HiRadixCache` 的空匹配结果路径。对正常前缀匹配路径无性能影响（只在未命中时减少一次 tensor 分配）。降低每请求的 Python 对象分配和 GC 压力，在大量短请求场景下有一定正面效果。兼容性：无用户可见 API 变更。- 风险标记：共享可变 tensor 引用，子类 `device` 一致性需确认

## 关联脉络

- 暂无明显关联 PR