

# PR #25249 完整报告

sgl-project/sglang

[NPU]fix:NPUMLATokenToKVPool object has no attribute "kv\_buffer"

合并时间: 2026-05-18 09:09

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/25249>

## 执行摘要

- 一句话: 修复 NPU MLA KV pool 属性缺失及复制方法
- 推荐动作: 建议通读, 特别是 `get_cpu_copy` / `load_cpu_copy` 的覆写模式。对于需要在不同硬件上定制基类行为的场景有参考价值。Review 中关于变量命名统一的讨论值得关注, 反映了 pipeline parallelism 下索引设计的关键细节。

## 功能与动机

NPUMLATokenToKVPool 的 kv buffer 与 MLATokenToKVPool 不同, 因此需要重写 `load_cpu_copy` 和 `get_cpu_copy` (来自 PR body)。

## 实现拆解

1. 在 `memory_pool_npu.py` 中为 `NPUMHATokenToKVPool` 和 `NPUMLATokenToKVPool` 分别新增 `_chunk_copy_npu_to_cpu` 辅助方法, 用于按块将 NPU 上的缓冲区张量异步拷贝到 CPU。
2. 为两个类重写 `get_cpu_copy` 方法: 先同步 NPU 流, 将各层的 k/v 缓冲区展平为 2D 视图 (MHA 用 `head_num/head_dim`, MLA 用 `kv_lora_rank/qk_rope_head_dim`), 然后调用辅助方法执行分块传输。
3. 重写 `load_cpu_copy` 方法: 同样先同步, 然后逐层将 CPU 上的分块数据拷贝回展平后的 NPU 缓冲区。
4. 在 `NPUMLATokenToKVPool` 中额外处理可选的 `index_head_dim` 缓冲区 (`ik_buffer`)。
5. 根据 review 建议统一了返回类型 (全部使用 `list`) 和变量命名 (全部使用 `local_layer_id`), 避免 pipeline parallelism 下的越界问题。

关键文件:

- `python/sglang/srt/hardware_backend/npu/memory_pool_npu.py` (模块 NPU 内存池; 类别 `source`; 类型 `core-logic`; 符号 `_chunk_copy_npu_to_cpu`, `get_cpu_copy`, `load_cpu_copy`): 唯一修改文件, 为 `NPUMHATokenToKVPool` 和 `NPUMLATokenToKVPool` 新增 CPU 卸载复制方法, 修复 `kv_buffer` 属性缺失错误。

关键符号: `_chunk_copy_npu_to_cpu`, `get_cpu_copy`, `load_cpu_copy`

## 关键源码片段

## python/sglang/srt/hardware\_backend/npu/memory\_pool\_npu.py

唯一修改文件，为 NPUMHATokenToKVPool 和 NPUMLATokenToKVPool 新增 CPU 卸载复制方法，修复 kv\_buffer 属性缺失错误。

```
# python/sglang/srt/hardware_backend/npu/memory_pool_npu.py
# NPUMHATokenToKVPool / NPUMLATokenToKVPool 共用辅助函数
def _chunk_copy_npu_to_cpu(self, buf_of_layers, indices):
```

```
    """将 NPU 上的张量按块异步拷贝到 CPU。
```

Args:

buf\_of\_layers: list of list of tensors, 每个元素对应一层的 [k\_buf, v\_buf] 或 [k\_buf, v\_buf, ik\_buf]

indices: 需要拷贝的 page 索引

Returns:

嵌套列表: [layer][chunk] -> list of tensors on CPU

```
    """
```

```
    chunk_size = self.cpu_offloading_chunk_size
```

```
    out = []
```

```
    for tensors_per_layer in buf_of_layers: # [k_buf, v_buf] or [k_buf, v_buf, ik_buf/None]
```

```
        layer_chunks = []
```

```
        for i in range(0, len(indices), chunk_size):
```

```
            ci = indices[i : i + chunk_size]
```

```
            # 过滤 None (MLA 中 ik_buffer 可能为 None)
```

```
            layer_chunks.append(
```

```
                [t[ci].to("cpu", non_blocking=True) for t in tensors_per_layer if t is not None]
```

```
            )
```

```
            out.append(layer_chunks)
```

```
    return out
```

```
# NPUMLA 版本 (部分展示)
```

```
class NPUMLATokenToKVPool(MLATokenToKVPool):
```

```
    # ...
```

```
    def get_cpu_copy(self, indices):
```

```
        torch.npu.synchronize()
```

```
        buf_of_layers = []
```

```
        has_ik = self.index_head_dim is not None
```

```
        for local_layer_id in range(self.layer_num):
```

```
            k_layer = self.k_buffer[local_layer_id].view(-1, 1, self.kv_lora_rank)
```

```
            v_layer = self.v_buffer[local_layer_id].view(-1, 1, self.qk_rope_head_dim)
```

```
            ik_layer = self.index_k_buffer[local_layer_id].view(-1, 1, self.index_head_dim) if has_ik  
            else None
```

```
            buf_of_layers.append([k_layer, v_layer, ik_layer] if has_ik else [k_layer, v_layer])
```

```
        kv_cache_cpu = self._chunk_copy_npu_to_cpu(buf_of_layers, indices)
```

```
        torch.npu.synchronize()
```

```
        return kv_cache_cpu
```

```
    def load_cpu_copy(self, kv_cache_cpu, indices):
```

```
        torch.npu.synchronize()
```

```

chunk_size = self.cpu_offloading_chunk_size
has_ik = self.index_head_dim is not None
for local_layer_id in range(self.layer_num):
    k_layer = self.k_buffer[local_layer_id].view(-1, 1, self.kv_lora_rank)
    v_layer = self.v_buffer[local_layer_id].view(-1, 1, self.qk_rope_head_dim)
    ik_layer = self.index_k_buffer[local_layer_id].view(-1, 1, self.index_head_dim) if has_ik
    else None
    for i in range(0, len(indices), chunk_size):
        ci = indices[i : i + chunk_size]
        tensors_cpu = kv_cache_cpu[local_layer_id][i // chunk_size]
        k_cpu, v_cpu = tensors_cpu[0], tensors_cpu[1]
        k_layer[ci] = k_cpu.to(k_layer.device, non_blocking=True)
        v_layer[ci] = v_cpu.to(v_layer.device, non_blocking=True)
        if has_ik:
            ik_cpu = tensors_cpu[2]
            ik_layer[ci] = ik_cpu.to(ik_layer.device, non_blocking=True)
torch.npu.synchronize()

```

## 评论区精华

Reviewer [whybeyoung](#) 指出两个关键问题：

1) MHA 版本使用 list 而 MLA 版本使用 tuple 存储块拷贝结果，建议统一； 2) MHA 版本使用 `layer_id`（全局索引）而 MLA 使用 `local_layer_id`（PP stage 内偏移），在 pipeline parallelism 下可能越界。作者同意并修改，最终统一为 list 和 `local_layer_id`。

- 统一返回类型（list vs tuple）（style）：作者同意并改为统一使用 list。
- 变量命名统一（`layer_id` vs `local_layer_id`）（correctness）：作者同意修改 MHA 版本统一使用 `local_layer_id`。
- 抽取公共 helper 函数（design）：作者采纳建议，重构代码。

## 风险与影响

- 风险：主要风险： 1) 缺乏单元测试覆盖新复制的逻辑路径，后续重构可能引入回归； 2) 若 NPU 后端启用 pipeline parallelism，`local_layer_id` 的一致性仍需验证； 3) 仅修改单一文件，未联动修改其他调用了 `kv_buffer` 的地方（但 `NPUMLATokenToKVPool` 继承自 `MLATokenToKVPool` 的父类，新增属性不会影响既有 MHA 路径）。
- 影响：影响范围：仅限 NPU 硬件后端使用 MLA 模型（如 DeepSeek 系列）时的 KV 缓存 CPU 卸载功能。修复了 `AttributeError` 崩溃，使 NPU 上的 MLA 推理可正常使用 offload 特性。不影响其他后端或非 MLA 模型。
- 风险标记：缺少测试覆盖，变量命名不一致（已修复）

## 关联脉络

- 暂无明显关联 PR