

# PR #24648 完整报告

sgl-project/sglang

[NIXL][XPU] Fix uint64 overflow for mismatched P/D TP sizes (e.g. prefill\_tp=1, decode\_tp=2)

合并时间: 2026-05-12 11:02

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/24648>

## 执行摘要

- 一句话: 修复 XPU 上 uint64 溢出导致 KV 缓存通信失败
- 推荐动作: 值得合并, 修复目标准确, 改动极小且安全。建议在后续版本中补充针对不匹配 TP 大小的集成测试, 以覆盖回归。

## 功能与动机

修复 Intel XPU 上 P/D 分离部署中不匹配 TP 大小 (prefill\_tp=1, decode\_tp>1) 时 KV 缓存地址因高位为 1 导致 int64 溢出, 使 send\_kvcache\_slice 无法正常工作的 bug。PR body 明确描述了触发条件。

## 实现拆解

1. 在 python/sglang/srt/disaggregation/nixl/conn.py 的 send\_kvcache\_slice 方法中, 将 self.kv\_args.kv\_data\_ptrs 和 dst\_kv\_ptrs 的构建从默认 dtype 改为 np.uint64, 确保指针值不被溢出。
2. 同时将 prefill\_kv\_indices、dst\_kv\_indices 和 token\_offsets 的 dtype 从 np.int64 改为 np.uint64, 因为这些变量也参与地址计算, 需要保持无符号语义。
3. 这些改动仅影响 XPU 平台, 对 CUDA 平台无副作用 (因为 CUDA 地址通常高位为 0, int64 也能正确表示)。

关键文件:

- python/sglang/srt/disaggregation/nixl/conn.py (模块 KV 通信; 类别 source; 类型 core-logic): 核心修复文件, 修改 send\_kvcache\_slice 方法中的数据类型, 解决 XPU 上 KV 缓存地址溢出问题。

关键符号: send\_kvcache\_slice

## 关键源码片段

[python/sglang/srt/disaggregation/nixl/conn.py](#)

核心修复文件, 修改 send\_kvcache\_slice 方法中的数据类型, 解决 XPU 上 KV 缓存地址溢出问题。

# 在 send\_kvcache\_slice 方法中, 关键变更如下:

# torch.int 在 Intel XPU 上对高位为 1 的地址 (如 0xffff81ab54e01000) 会溢出,

```
# 使用 np.uint64 确保地址不被截断或解释为负数。

# 变更前:
# self.get_mha_kv_ptrs_with_pp(self.kv_args.kv_data_ptrs, dst_kv_ptrs)
# 变更后: 显式指定 uint64
kv_data_ptrs = np.array(self.kv_args.kv_data_ptrs, dtype=np.uint64)
dst_kv_ptrs = np.array(dst_kv_ptrs, dtype=np.uint64)
src_k_ptrs, src_v_ptrs, dst_k_ptrs, dst_v_ptrs, layers_current_pp_stage = (
    self.get_mha_kv_ptrs_with_pp(kv_data_ptrs, dst_kv_ptrs)
)

# 索引与偏移数组同样需要 uint64 以保持一致性
prefill_indices = np.asarray(prefill_kv_indices, dtype=np.uint64)
dst_indices = np.asarray(dst_kv_indices, dtype=np.uint64)
token_offsets = np.arange(page_size, dtype=np.uint64)
```

## 评论区精华

无实质讨论，仅 PR 创建者 mingfeima 评论指出此 PR 与 #24188 类似。Reviewer ShangmingCai 直接 approve，无其他争议。

- 暂无高价值评论线程

## 风险与影响

- 风险：
  1. 回归风险低：改动仅将涉及指针运算的数组 dtype 从 int64 改为 uint64，语义更精确，对正常地址（高位为 0）无任何影响。
  2. 缺少新增测试：本 PR 未附带测试用例来覆盖 prefill\_tp=1, decode\_tp=2 的场景，长期来看存在回归风险。
  3. 兼容性：uint64 在 Python 中可能因平台差异导致大整数精度问题，但此处均为 numpy 数组内部计算，风险可控。- 影响：影响范围：仅涉及 Intel XPU 上 P/D 分离部署中 TP 大小不匹配的场景。修复后该配置可以正常使用，否则会因溢出报错。对 CUDA 和其他平台无影响。- 风险标记：缺少测试覆盖

## 关联脉络

- PR #24188 类似修复（被提及）：mingfeima 指出本 PR 与 #24188 类似，可能涉及同一函数或同类型溢出问题。