

# PR #22778 完整报告

sgl-project/sglang

[BUGFIX]Fix Ascend backend pre-allocated range in NPU Graph Mode.

合并时间: 2026-04-24 01:23

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/22778>

## 执行摘要

- 一句话: 修复 NPU Graph Mode 中 `block_tables` 预分配不足问题
- 推荐动作: 该 PR 值得相关开发者精读, 尤其是维护 NPU 后端和 Graph Mode 特性的工程师。它揭示了一个常见的边界问题: 预分配内存时未考虑推测解码产生的额外令牌。设计决策上体现了“在初始化阶段一次性预留足够内存”的容错思路, 避免了运行时动态增长的复杂性。

## 功能与动机

修复 NPU Graph Mode 下启用 `--speculative-num-draft-tokens` 后, 输出 token 长度达到系统最大上下文长度时, 因预分配内存不足导致的间发性异常。PR body 中描述了该错误为“`tensor size to be copied` 与 `graph` 初始化时预分配的大小不匹配”。

## 实现拆解

### 1. 修改 `init_cuda_graph_state` 方法中的预分配计算

- 涉及文件: `python/sglang/srt/hardware_backend/npu/attention/ascend_backend.py`
- 关键符号: `init_cuda_graph_state`, `speculative_num_draft_tokens`, `block_tables`, `block_tables_swa`
- 具体变更: 引入变量 `total_context_len`, 在原上下文长度的基础上加上 `self.speculative_num_draft_tokens` (如果该值不为 `None`), 然后用 `total_context_len // page_size` 计算块数量。
- 原因: 原代码使用 `self.max_context_len + self.page_size - 1` 计算最大上下文长度, 未考虑推测解码产生的额外令牌, 导致预分配的 `block_tables` 容量不足。当实际序列长度达到 `max_context_len + speculative_num_draft_tokens` 时, `copy_` 操作会访问超出预分配范围的张量, 引发异常。

### 2. 消除代码重复

- 变更后的代码在 `block_tables` 和 `block_tables_swa` 两个预分配张量中都使用了 `total_context_len // page_size`, 避免了两个地方分别写不同计算逻辑可能引入的不一致。

### 3. 测试与验证

- 没有新增单元测试，仅通过手动构造超过预分配大小的场景验证：使用 QwQ-32B 模型，`--context-length 32768`，`--speculative-num-draft-tokens 16`，在极限流式输出的压力测试中确认请求能正常返回且无崩溃。

关键文件：

- `python/sglang/srt/hardware_backend/npu/attention/ascend_backend.py`（模块 NPU 后端；类别 `source`；类型 `core-logic`；符号 `init_cuda_graph_state`）：这是唯一的变更文件。其中 `init_cuda_graph_state` 方法在初始化 NPU Graph Mode 的预分配张量时，需要包含推测解码产生的额外令牌数，否则在长序列下会发生内存越界。

关键符号：`init_cuda_graph_state`

## 关键源码片段

### `python/sglang/srt/hardware_backend/npu/attention/ascend_backend.py`

这是唯一的变更文件。其中 `init_cuda_graph_state` 方法在初始化 NPU Graph Mode 的预分配张量时，需要包含推测解码产生的额外令牌数，否则在长序列下会发生内存越界。

```
# 文件：python/sglang/srt/hardware_backend/npu/attention/ascend_backend.py
# 方法：init_cuda_graph_state
# 关键修改：在计算 block_tables 的预分配大小时加入推测令牌数
```

```
def init_cuda_graph_state(self, max_bs: int, max_num_tokens: int):
    # 计算总上下文长度：原 max_context_len + page_size - 1（用于向上取整）
    total_context_len = self.max_context_len + self.page_size - 1
    # 如果启用了推测解码，额外加上推测令牌数量，避免预分配不足
    if self.speculative_num_draft_tokens is not None:
        total_context_len += self.speculative_num_draft_tokens

    # 预分配 block_tables，形状为 (max_bs, 块数)
    self.graph_metadata = {
        "block_tables": torch.empty(
            (max_bs, total_context_len // self.page_size), # 原式不包含推测令牌
            dtype=torch.int32,
            device=self.device,
        ),
    }
    # 对于混合 SWA 模型，同样预分配 block_tables_swa
    if self.is_hybrid_swa:
        self.graph_metadata["block_tables_swa"] = torch.empty(
            (max_bs, total_context_len // self.page_size),
            dtype=torch.int32,
            device=self.device,
        )
```

## 评论区精华

评审建议：提取公共变量

- 讨论者: gemini-code-assist[bot]
- 原文: 建议将块数计算提取为 `max_num_blocks = (total_context_len + self.page_size - 1) // self.page_size` 变量, 用于 `block_tables` 和 `block_tables_swa`, 以减少重复并提升可读性。
- 采纳情况: 该建议在最终合并版本中未完全采用——变量 `total_context_len` 被保留, 但计算式直接使用了 `total_context_len // self.page_size` (实际上 `total_context_len` 已包含 `self.page_size - 1` 的预处理)。评审意见虽合理, 但作者选择了更直接的方式, 同样达到了消除重复的效果。
- 结论: 已解决, 评审人 iforgetmyname 批准了 PR。
- 提取公共变量以减少重复 (design): 未完全采纳建议, 但已通过 `total_context_len` 变量间接实现了代码复用。评审人批准了 PR。

## 风险与影响

- 风险: ### 回归风险
- 低: 改动仅涉及 `block_tables` 和 `block_tables_swa` 预分配尺寸的计算, 将 `max_context_len` 扩展为 `max_context_len + speculative_num_draft_tokens` (若启用推测解码)。原始行为 (无推测解码) 完全不变, 因此不会对非 NPU 或非推测场景造成回归。

## 性能风险

- 极低: 预分配的张量尺寸略有增加 (最多增加推测令牌数对应的块数), 对内存和运行时影响可忽略不计。

## 兼容性风险

- 低: 仅改动 Ascend 后端 NPU Graph Mode 特有的初始化路径, 不影响其他硬件后端或普通推理模式。
- 影响: ### 用户影响
- 正面: NPU 用户在启用推测解码 (如 NGRAM 算法) 和 Graph Mode 时, 不会再遇到偶发的内存越界崩溃, 长上下文推理稳定性显著提升。

## 系统影响

- 低: 仅影响 Ascend NPU 后端, 其他硬件 (如 CUDA、AMD) 不受影响。改动范围小, 仅一个文件 7 行变更。

## 团队影响

- 低: 无需额外配置或迁移。代码逻辑直观, 维护简单。
- 风险标记: 核心路径变更, 缺少测试覆盖

## 关联脉络

- PR #23572 [Diffusion][NPU][Bugfix] Ascend\_fa crashes when sequence parallelism is used.: 同为 NPU 后端的 bug 修复, 涉及 Ascend 注意力后端, 与本 PR 有类似的问题域 (NPU Graph Mode 下的内存分配)。