

PR #26797 完整报告

sgl-project/sglang

[core] Compute token_type_ids in ForwardBatch.init_new

合并时间: 2026-05-31 15:54

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/26797>

执行摘要

- 一句话: 将 token_type_ids 计算挪入 ForwardBatch
- 推荐动作: 值得精读。该 PR 展示了如何通过重构保持 ScheduleBatch 的职责纯洁性 (只做调度编排), 将前向相关的设备张量构建下沉到 ForwardBatch, 是流管理和职责分离的良好实践。但需注意 review 中提出的性能建议尚未解决, 可在后续跟进。

功能与动机

PR body 明确说明: "Move the cross-encoder token_type_ids device-tensor build out of ScheduleBatch.prepare_for_extend into ForwardBatch.init_new, so the H2D runs on the forward stream and SB no longer carries a forward-only field." 这是一个明确的职责分离和流管理优化。

实现拆解

1. 在 forward_batch_info.py 的 init_new 中移除了 token_type_ids=batch.token_type_ids 的传递, 改为从 batch.reqs 直接收集并构建设备张量。
2. 将原有的 _maybe_init_prefill_only 方法重命名为 _maybe_init_non_generation_fields, 扩大了其职责: 除了原有的 dimensions、return_pooled_hidden_states、multi_item_delimiter_indices 外, 新增了对 token_type_ids 的处理。
3. 在 schedule_batch.py 的 ScheduleBatch 类定义中移除了 token_type_ids 字段声明, 并在 prepare_for_extend 中删除了对应的收集、构造张量以及赋值操作。清理后的代码中, prepare_for_extend 不再持有或构建该张量。
4. 为支持 token_type_ids 的 pinned memory 分配, 在 forward_batch_info.py 的 import 中新增了 is_pin_memory_available 的导入。

关键文件:

- python/sglang/srt/model_executor/forward_batch_info.py (模块 前向批处理; 类别 source; 类型 core-logic; 符号 _maybe_init_prefill_only, _maybe_init_non_generation_fields): 核心变更文件, 重命名并扩展 _maybe_init_prefill_only 为 _maybe_init_non_generation_fields, 在其中新增 token_type_ids 设备张量构建逻辑。
- python/sglang/srt/managers/schedule_batch.py (模块 调度批处理; 类别 source; 类型 core-logic): 移除了 token_type_ids 字段声明、收集及构建逻辑, 对应职责移除。

关键符号: `_maybe_init_prefill_only`, `_maybe_init_non_generation_fields`, `prepare_for_extend`, `init_new`

关键源码片段

`python/sglang/srt/model_executor/forward_batch_info.py`

核心变更文件, 重命名并扩展 `_maybe_init_prefill_only` 为 `_maybe_init_non_generation_fields`, 在其中新增 `token_type_ids` 设备张量构建逻辑。

```
# python/sglang/srt/model_executor/forward_batch_info.py

def _maybe_init_non_generation_fields(self, batch: ScheduleBatch):
    """Derive non-generation (max_new_tokens==0) forward fields from reqs.

    token_type_ids gates on presence, not is_prefill_only: a missing
    tensor makes bert/roberta silently fall back to zeros.
    """
    if self.is_prefill_only:
        # 原有逻辑: dimensions (Matryoshka), return_pooled_hidden_states, multi_item_delimiter_
        # indices
        if batch.model_config.is_matryoshka and any(
            r.dimensions is not None for r in batch.reqs
        ):
            self.dimensions = [
                r.dimensions if r.dimensions else batch.model_config.hidden_size
                for r in batch.reqs
            ]
            self.return_pooled_hidden_states = any(
                r.return_pooled_hidden_states for r in batch.reqs
            )
            if get_global_server_args().enable_mis and any(
                r.multi_item_delimiter_indices is not None for r in batch.reqs
            ):
                assert all(
                    r.multi_item_delimiter_indices is not None for r in batch.reqs
                ), "MIS batch must have delimiter indices on every request"
                self.multi_item_delimiter_indices = [
                    torch.tensor(r.multi_item_delimiter_indices, dtype=torch.int64)
                    for r in batch.reqs
                ]

    # 新增: 从每个 req 收集 token_type_ids, 若存在则构建设备张量
    # 注意: 当前实现使用 sum(list_of_lists, []) 展平, 复杂度 O(N^2)
    token_type_ids = [
        r.token_type_ids for r in batch.reqs if r.token_type_ids is not None
    ]
    if token_type_ids:
        self.token_type_ids = torch.tensor(
            sum(token_type_ids, []),
```

```
dtype=torch.int64,  
pin_memory=is_pin_memory_available(self.device),  
)to(self.device, non_blocking=True)
```

评论区精华

gemini-code-assist[bot] 在 review 中提出性能优化建议：原来使用 `sum(list_of_lists, [])` 展平列表是 $O(N^2)$ 的反模式，建议改为嵌套列表推导式 ($O(N)$ 复杂度)。该建议未被作者采纳或回复，PR 就已合并。

- 使用 $O(N^2)$ `sum` 展平列表的性能问题 (performance): 建议未被采用或回复，PR 直接合并。

风险与影响

- 风险：
 1. 性能风险：当前实现仍使用 `sum(list_of_lists, [])` 展平 `token_type_ids`，对于大批量请求可能存在 $O(N^2)$ 的性能问题，但通常 `token_type_ids` 长度较小，实际影响有限。
 2. 回归风险：`token_type_ids` 的构建逻辑从 `prepare_for_extend` 移到了 `init_new`，且只有 `is_prefill_only` 时才会执行，需要确认所有使用 `token_type_ids` 的场景（如上文提到的 `bert/roberta`）均满足此条件。
 3. 兼容性风险：`ScheduleBatch` 不再暴露 `token_type_ids` 字段，任何直接访问该字段的外部调用都会出错。 - 影响：影响范围较小，仅涉及两个核心文件。主要影响使用 `cross-encoder`（如 `bert/roberta`）且启用了 `is_prefill_only` 模式的请求路径。对于普通 `decode` 路径无影响。 - 风险标记：潜在性能反模式，缺失测试覆盖

关联脉络

- 暂无明显关联 PR