

# PR #26022 完整报告

sgl-project/sglang

Group ScheduleBatch and ForwardBatch fields by data-flow role

合并时间: 2026-05-28 15:11

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/26022>

## 执行摘要

- 一句话: 按数据流角色重组批量数据结构字段声明
- 推荐动作: 建议所有涉及推理调度和模型执行的开发者快速浏览此 PR, 以了解 ScheduleBatch 和 ForwardBatch 的新分组约定。该约定有助于在后续开发中保持字段组织一致性, 避免随意添加字段。

## 功能与动机

PR 作者指出需要对 `ScheduleBatch` 和 `ForwardBatch` 的字段声明进行重组, 按照数据流角色划分章节 (如核心请求列表、全局配置、批处理变体状态、传递给 `ForwardBatch` 的 GPU 张量等), 并移除死字段 `is_hybrid_swa`, 以降低理解难度, 防止未来开发者在错误位置添加新字段。

## 实现拆解

1. 移除死字段: 在 `schedule_batch.py` 中删除 `ScheduleBatch.is_hybrid_swa` 字段, 该字段仅被赋值而无任何读取者, 可从分配器类型重新计算。
2. 重组 `ScheduleBatch` 字段: 按数据流角色将约 30 个字段划分为若干组: 核心请求列表 (`reqs`)、全局配置与共享资源 (内存池、模型配置、设备)、批变体调度器状态 (`batch_is_full`、分块预填充、DP 注意力、分割预填充等)、传递给 `ForwardBatch` 的 GPU 张量 (`input_ids`、`req_pool_indices`、`seq_lens` 等)、一次性 forward 覆盖 (`forward_mode`、`sampling_info` 等)。并为每组添加醒目的章节注释。
3. 重组 `ForwardBatch` 字段: 按字段来源分组: 必需核心输入 (模式、批量大小、张量)、从 `ScheduleBatch` 借用的 GPU 张量 (标记跨流克隆目标)、从 `ScheduleBatch` 借用的配置 / 标志、从 `reqs` 派生的主机元数据、从一次性覆盖解析的字段、前向派生的全属字段 (`FB-owned`)、运行时填充字段。同样添加章节注释。
4. 配套调整: 在 `schedule_batch.py` 中移除了对 `SWATokenToKVPoolAllocator` 的导入, 因 `is_hybrid_swa` 被移除后不再需要此引用。`forward_batch_info.py` 中 `ForwardBatch` 的 `dataclass` 字段顺序也相应调整, 但未更改字段默认值或类型。

关键文件:

- `python/sglang/srt/managers/schedule_batch.py` (模块 调度器; 类别 `source`; 类型 `refactor`; 符号 `ScheduleBatch`): 核心调度批处理数据类, 本次对其所有字段进行按数据流角色分组, 并移除了已废弃的 `is_hybrid_swa` 字段。

- `python/sglang/srt/model_executor/forward_batch_info.py` (模块 前向批处理; 类别 `source`; 类型 `refactor`; 符号 `ForwardBatch`): 前向批处理数据类, 同样按来源分组, 明确标识了从 `ScheduleBatch` 借用、派生、运行时填充等类别的字段。

关键符号: 未识别

## 关键源码片段

### `python/sglang/srt/managers/schedule_batch.py`

核心调度批处理数据类, 本次对其所有字段进行按数据流角色分组, 并移除了已废弃的 `is_hybrid_swa` 字段。

```
@dataclasses.dataclass
class ScheduleBatch(ScheduleBatchDisaggregationDecodeMixin):
    # (docstring omitted)

    # === Core: request list ===
    reqs: List[Req]

    # === Global config and shared resources ===
    req_to_token_pool: ReqToTokenPool = None
    token_to_kv_pool_allocator: BaseTokenToKVPoolAllocator = None
    tree_cache: BasePrefixCache = None
    model_config: ModelConfig = None
    enable_overlap: bool = False
    device: str = "cuda"
    hisparse_coordinator: Optional[HiSparseCoordinator] = None

    # === Batch-variant scheduler state ===
    batch_is_full: bool = False
    chunked_req: Optional[Req] = None
    contains_last_prefill_chunk: bool = True
    inner_idle_batch: Optional[ScheduleBatch] = None
    decoding_reqs: List[Req] = None
    split_index: int = 0
    split_prefill_finished: bool = False
    split_forward_count: int = 1
    split_forward_batch: ForwardBatch = None
    # ... (remaining fields omitted for brevity)
```

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

前向批处理数据类, 同样按来源分组, 明确标识了从 `ScheduleBatch` 借用、派生、运行时填充等类别的字段。

```
@dataclass
class ForwardBatch(ForwardBatchDeepSeekMHAMixin):
    # (docstring omitted)

    # === Required core inputs ===
```

```
forward_mode: ForwardMode
batch_size: int
input_ids: torch.Tensor
req_pool_indices: torch.Tensor
seq_lens: torch.Tensor
out_cache_loc: torch.Tensor
seq_lens_sum: int

# === Borrowed from ScheduleBatch: GPU tensors ===
orig_seq_lens: Optional[torch.Tensor] = None
mamba_track_indices: Optional[torch.Tensor] = None
mamba_track_mask: Optional[torch.Tensor] = None
mamba_track_seq_lens: Optional[torch.Tensor] = None
mamba_cow_src_indices: Optional[torch.Tensor] = None
mamba_cow_dst_indices: Optional[torch.Tensor] = None
mamba_clear_indices: Optional[torch.Tensor] = None
input_embeds: Optional[torch.Tensor] = None
replace_embeds: Optional[torch.Tensor] = None
replace_positions: Optional[torch.Tensor] = None
token_type_ids: Optional[torch.Tensor] = None
encoder_lens: Optional[torch.Tensor] = None
encoder_out_cache_loc: Optional[torch.Tensor] = None
# ... (remaining groups omitted)
```

## 评论区精华

本 PR 未触发实质性的 Review 讨论，仅有一条来自 bot 的配额警告，无人工审核意见。PR 作者自行合并。

- 暂无高价值评论线程

## 风险与影响

- 风险：本次变更为纯注释和字段顺序调整，未修改任何运行时逻辑或字段值。回归风险极低。但若未来有人依赖字段在 dataclass 定义中的位置进行反射或序列化，可能受影响。不过 Python dataclass 的 field 顺序默认按定义顺序，但通常不应依赖。整体风险可忽略。
- 影响：对用户无影响。对团队而言，提高了两个核心数据类的可读性，降低了新成员理解门槛。但若后续有合入冲突，字段位置变化可能导致合并复杂度略微上升。
- 风险标记：暂无

## 关联脉络

- 暂无明显关联 PR