

PR #26047 完整报告

sgl-project/sglang

Add --disable-attn-tp-gather opt-out for model-managed SP

合并时间: 2026-05-25 06:53

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/26047>

执行摘要

- 一句话: 新增 opt-out 开关跳过 SP gather 路径
- 推荐动作: 该 PR 值得精读, 尤其是需要优化解码性能或集成新模型的工程师:
 - 学习如何通过短小精悍的配置项绕过对性能有害的通用路径。
 - 关注 `require_attn_tp_gather` 的短路模式, 可作为类似优化的模板。

功能与动机

部分模型在 attention 内部自行执行 `all_gather/reduce_scatter`, 不消费上游 `gathered_buffer`, 但调度器侧的 `attn_tp_gather` 机制仍会将 `num_tokens` 填充至 `attn_tp_size` 并预分配缓冲, 导致 CUDA graph 捕获时 kernel autotuner (如 FlashInfer) 选错 kernel 变体, 在小 batch 时显著退化 decode TPOT。此前唯一绕过方式是废弃 `--moe-a2a-backend`, 但也禁用了模型原本的 MoE A2A 路径。

实现拆解

1. 新增配置项 (`python/sglang/srt/server_args.py`) - 在 `ServerArgs` 类中新增字段 `disable_attn_tp_gather: bool = False`。 - 在 CLI 参数注册中添加 `--disable-attn-tp-gather`, action 为 `store_true`, 默认关闭。 - 帮助文本通用, 不包含模型特定表述。
2. 短路判断逻辑 (`python/sglang/srt/utils/common.py`) - 在 `require_attn_tp_gather` 函数开头增加短路: 若 `server_args.disable_attn_tp_gather` 为 True, 直接返回 False。 - 该短路效果级联至 `require_gathered_buffer` 和 `require_mlp_sync`, 从而跳过:
 - `cuda_graph_runner` 中的 `gathered_buffer` 分配;
 - `adjust_num_token_non_padded_for_attn_tp` 的 forward 执行;
 - 调度器 `prepare_mlp_sync_batch` 中的 `gloo:all_gather`。
3. 无行为改变: 默认开关关闭, 不影响现有逻辑。

关键文件:

- `python/sglang/srt/server_args.py` (模块 配置; 类别 source; 类型 configuration) : 新增 CLI 参数和 `ServerArgs` 字段, 作为开关入口
- `python/sglang/srt/utils/common.py` (模块 调度器; 类别 source; 类型 core-logic) : 在 `require_attn_tp_gather` 中实现短路逻辑, 是实际生效点

关键符号: require_attn_tp_gather

关键源码片段

python/sglang/srt/server_args.py

新增 CLI 参数和 ServerArgs 字段, 作为开关入口

```
# python/sglang/srt/server_args.py (partial)
```

```
class ServerArgs:
```

```
...
```

```
# 新增: 用于模型自管理 SP 时绕过调度器侧的 attn_tp_gather 路径
```

```
disable_attn_tp_gather: bool = False # 默认不启用
```

```
...
```

```
@classmethod
```

```
def add_cli_args(cls, parser: argparse.ArgumentParser):
```

```
...
```

```
parser.add_argument(
```

```
    "--disable-attn-tp-gather",
```

```
    action="store_true",
```

```
    help=(
```

```
        "Disable scheduler-side attn_tp_gather (the upstream SP path "
```

```
        "that pads num_tokens to attn_tp_size and pre-allocates a gathered "
```

```
        "buffer). Use for models that manage SP scatter/gather at the "
```

```
        "model level (e.g., perform their own all_gather/reduce_scatter "
```

```
        "inside attention) and do not consume the upstream gathered_buffer. "
```

```
        "Without this, the cuda graph runner pads num_tokens to attn_tp_size, "
```

```
        "which can cause kernel autotuners to select wrong-sized variants "
```

```
        "at small batches."
```

```
    ),
```

```
)
```

```
...
```

python/sglang/srt/utils/common.py

在 require_attn_tp_gather 中实现短路逻辑, 是实际生效点

```
# python/sglang/srt/utils/common.py (partial)
```

```
def require_attn_tp_gather(server_args: ServerArgs):
```

```
    """
```

```
    Check if the input of attention is scattered.
```

```
    """
```

```
# 新增短路: 若用户显式禁用了 attn_tp_gather, 则直接返回 False
```

```
# 适用于模型在 attention 内部自行管理 SP 的场景
```

```
if server_args.disable_attn_tp_gather:
```

```
    return False
```

```
from sglang.srt.layers.moe.utils import get_moe_a2a_backend
```

```
assert server_args.moe_dense_tp_size in [1, None]
if not get_moe_a2a_backend().is_none() or server_args.moe_dense_tp_size == 1:
    if server_args.enable_dp_attention:
        return server_args.dp_size < server_args.tp_size
    else:
        return True
else:
    return False
```

评论区精华

该 PR 无 review 评论讨论。仅有 hanming-lu 的批准，无进一步技术讨论。

- 暂无高价值评论线程

风险与影响

- 风险：
 - 回归风险低：新参数默认 False，不影响现有逻辑。
 - 误用风险：若用户误对需要 SP gather 的模型启用该开关，可能导致正确性问题（如 attention 输入未正确同步）。但该风险已有 help text 说明。
 - 覆盖不足：无自动化测试验证开关与行为的一致性。
- 影响：
 - 影响范围：仅影响显式启用 `--disable-attn-tp-gather` 的使用者，主要面向模型开发者或高性能调优场景。
 - 性能影响：对自管理 SP 的模型（如 DeepSeek-V4），BS=1 下 TPOT 消除 29% 回归，大 batch 下无影响。
 - 团队影响：为后续模型开发者提供标准化开关，避免需要 hack 内部逻辑。
 - 风险标记：无自动化测试覆盖，误用可能导致正确性问题

关联脉络

- PR #25898 [AMD] Dsv4/pr1 fix run time issue: 与 DeepSeek-V4 模型优化相关，该模型可能是自管理 SP 的典型案例
- PR #25948 [dsv4] support eplb: 同为 DeepSeek-V4 相关优化，可能受益于该开关
- PR #26218 suppress cutlass-dsl noisy warning: 同样涉及 common.py 的修改，但无直接关联