

PR #25444 完整报告

sgl-project/sglang

Bundle Scheduler rank/size fields into a frozen ParallelState

合并时间: 2026-05-16 09:23

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/25444>

执行摘要

- 一句话: 将调度器 17 个 rank/size 字段封装为冻结的 ParallelState 值对象
- 推荐动作: 建议仔细阅读 parallel_state_wrapper.py 中 ParallelState 的定义和 scheduler.py 中构建它的逻辑, 理解作者如何通过值对象模式统一管理并行拓扑信息。对于代码评审者, 重点关注是否有任何 self.<rank/size field> 仍留在未修改的文件中 (尤其是条件编译或特定后端路径)。

功能与动机

调度器中分散着大量与并行拓扑相关的整型字段 (tp_rank, pp_size, attn_cp_rank 等), 它们总是成组出现且初始化后不应改变。将它们捆绑进一个冻结的 dataclass, 能明确表示它们属于同一概念模块, 防止误修改, 简化构造函数签名, 并为后续将 ParallelState 作为参数传递给其他组件提供单一事实来源。

实现拆解

1. 定义 ParallelState dataclass: 在新增文件 python/sglang/srt/distributed/parallel_state_wrapper.py 中创建 @dataclass(frozen=True, slots=True, kw_only=True) 类, 包含 17 个字段, 涵盖所有 rank/size 属性。
2. 修改 Scheduler.init: 移除原先逐个 self.tp_rank = tp_rank 的赋值方式, 改为调用 ParallelState(...) 构造并赋值给 self.ps; 从函数参数中删除曾经冗余的 gpu_id (现在由 ps.gpu_id 承载), 同时调整 compute_dp_attention_world_info 调用以使用局部变量。
3. 逐个替换 Mixin 与消费者: 在 scheduler_pp_mixin.py、scheduler_profiler_mixin.py、scheduler_metrics_mixin.py、scheduler_output_processor_mixin.py、scheduler_dp_attn_mixin.py、disaggregation/prefill.py、disaggregation/decode.py、ray/scheduler_actor.py、layers/dp_attention.py 中, 将所有 self.<field> 替换为 self.ps.<field>, 涉及日志、条件判断、构造函数参数以及数据处理逻辑。
4. 修复 Disaggregation 遗漏: 在 scheduler_profiler_mixin.py 的 init_profiler 和 start_profile 中, 将 self.tp_rank 和 self.gpu_id 改为 self.ps.tp_rank 和 self.ps.gpu_id, 同时更新了 RPD profile 路径构建及 TP rank 判断逻辑, 确保与新版 ParallelState 一致。
5. 更新测试: 在 test/registered/unit/observability/test_forward_pass_metrics.py 中添加辅助函数 _make_ps, 用于快速构造具有合理默认值的 ParallelState 实例, 调整现有测试用例使用 scheduler.ps = _make_ps(...) 替代直接设置属性。

关键文件：

- python/sclang/srt/distributed/parallel_state_wrapper.py (模块 分布式状态; 类别 source; 类型 core-logic; 符号 ParallelState) : 新文件, 定义了核心值对象 ParallelState, 是本次重构的重心。
- python/sclang/srt/managers/scheduler.py (模块 调度器; 类别 source; 类型 dependency-wiring) : 调度器主入口, 修改了 __init__ 构造函数: 移除分散的字段赋值, 引入 self.ps 构造。
- test/registered/unit/observability/test_forward_pass_metrics.py (模块 测试; 类别 test ; 类型 test-coverage; 符号 _make_ps) : 测试文件, 添加了 _make_ps 辅助函数, 并更新了测试用例以使用 scheduler.ps。

关键符号: ParallelState, _make_ps

关键源码片段

python/sclang/srt/managers/scheduler.py

调度器主入口, 修改了 __init__ 构造函数: 移除分散的字段赋值, 引入 self.ps 构造。

```
# 在 __init__ 中替换了原先零散的 self.tp_rank = tp_rank 等赋值
```

```
# 改为统一构造 ParallelState 值对象
```

```
self.ps = ParallelState(  
    tp_rank=tp_rank,  
    tp_size=server_args.tp_size,  
    pp_rank=pp_rank,  
    pp_size=server_args.pp_size,  
    dp_rank=dp_rank,  
    dp_size=server_args.dp_size,  
    attn_tp_rank=attn_tp_rank,  
    attn_tp_size=attn_tp_size,  
    attn_cp_rank=attn_cp_rank,  
    attn_cp_size=server_args.attn_cp_size,  
    attn_dp_rank=attn_dp_rank,  
    attn_dp_size=attn_dp_size,  
    moe_ep_rank=moe_ep_rank,  
    moe_ep_size=server_args.ep_size,  
    moe_dp_rank=moe_dp_rank,  
    moe_dp_size=server_args.moe_dp_size,  
    gpu_id=gpu_id,  
)
```

```
# 之后所有消费者通过 self.ps.xxx 访问并行状态
```

评论区精华

本 PR 无常规 review 讨论。作者在 PR body 中说明了使用 `frozen=True, slots=True, kw_only=True` 的设计选择, 强调值对象的不可变性, 并指出通过 squash 修复 disagg 遗漏来保持 commit 与逻辑映射的清晰。

- 暂无高价值评论线程

风险与影响

- 风险：
 - 遗漏引用：重构涉及 14 个文件的机械替换，虽然 diff 覆盖了所有显式出现的字段，但隐式或动态访问（如 `getattr(self, 'tp_rank')` 或通过字符串拼接）可能未被完全替换。PR 作者已专门扫描并修复了 `disagg` 中的两处遗漏，但其他非常用路径（如条件编译 NPU/ROCM 分支）仍需回归测试验证。
 - 初始化顺序：`self.ps` 在 `__init__` 中构造较靠后（在第 400 行附近），若此前有任何方法使用了 `self.ps` 将引发 `AttributeError`。通过检查代码，`__init__` 中在赋值 `self.ps` 之前未调用其他访问并行字段的方法，风险可控。
 - 兼容性：对用户透明，无 API 或配置变更。
- 影响：
 - 开发团队：代码可读性与可维护性提升，并行状态现在具有显式类型和不可变性，减少了误用可能；但所有涉及并行字段的代码需要适应 `self.ps.` 前缀。
 - 功能影响：无功能变化，所有行为应与之前一致。
 - 测试覆盖：核心调度路径由现有多组 CI 覆盖（等待 `run-ci` 标签），测试文件已同步更新。
 - 重构链：本 PR 是“`parallel-state`”重构链的基础步骤，后续 PR（如 #25445）将依赖此 `ParallelState` 类。
 - 风险标记：大规模机械替换，核心路径变更

关联脉络

- PR #25445 Inject ParallelState into ProfilerV2: 本 PR 引入 `ParallelState` 类后，#25445 将其注入到 `ProfilerV2` 组件中，属于同一重构链的后续步骤。
- PR #25446 Fix V2 trace filename collisions when DP/PP/EP enabled: 该 PR 修复了 `ProfilerV2` 中的文件名冲突，与本 PR 的并行状态重构在 `profiler` 路径上存在间接关联。