

PR #24003 完整报告

sgl-project/sclang

[scheduler] Zero gen_throughput and flush KV events on pause

合并时间: 2026-06-02 16:43

原文链接: <http://prhub.com.cn/sgl-project/sclang/pull/24003>

执行摘要

- 一句话: 暂停时清零吞吐并刷新 KV 事件
- 推荐动作: 值得合并, 变更小且聚焦, 解决了明确的监控数据问题。对于关注可观测性的团队, 可以进一步检查暂停期间的其他指标是否也需要类似处理。

功能与动机

当调度器因权重更新等管理操作暂停时, 事件循环在到达 `on_idle` 之前短接, 导致 `gen_throughput` 持续显示上一个非零值, KV 事件也无法刷新, 造成监控面板误导和 KV 事件消费者延迟。

实现拆解

1. 在 `python/sclang/srt/managers/scheduler.py` 的 `pause_generation` 方法末尾, 新增清零 `metrics_reporter.last_gen_throughput` 及强制触发空闲指标日志和 KV 事件刷新的逻辑。
2. 为绕过 `_maybe_log_idle_metrics` 的 30 秒速率限制, 重置 `metrics_collector.last_log_time` 为 0.0。
3. 仅在启用了调度器指标收集时执行指标日志, 否则仅清零吞吐。
4. 在测试文件 `test/registered/unit/managers/test_scheduler_pause_generation.py` 的 `_new_scheduler` 中添加 `metrics_reporter` 和 `kv_events_publisher` 的 Mock 桩, 避免新代码在测试中引发属性错误。

关键文件:

- `python/sclang/srt/managers/scheduler.py` (模块 调度器; 类别 `source`; 类型 `core-logic`) : 核心修改文件, 在 `pause_generation` 方法末尾增加清零吞吐和刷新 KV 事件的代码。
- `test/registered/unit/managers/test_scheduler_pause_generation.py` (模块 测试; 类别 `test`; 类型 `test-coverage`) : 为 `_new_scheduler` 添加 `metrics_reporter` 和 `kv_events_publisher` 的 Mock 桩, 确保新逻辑在测试中不会因属性缺失而失败。

关键符号: `pause_generation`, `_new_scheduler`

关键源码片段

[python/sclang/srt/managers/scheduler.py](#)

核心修改文件，在 `pause_generation` 方法末尾增加清零吞吐和刷新 KV 事件的代码。

```
# python/sqglang/srt/managers/scheduler.py
# 位于 pause_generation 方法末尾，紧跟在原有 batch 清理逻辑之后
self.metrics_reporter.last_gen_throughput = 0.0 # 清零吞吐，防止暂停期间显示旧值
if self.metrics_reporter.current_scheduler_metrics_enabled:
    # 重置速率限制时钟，确保 _maybe_log_idle_metrics 立即执行一次
    self.metrics_reporter.metrics_collector.last_log_time = 0.0
    self.metrics_reporter._maybe_log_idle_metrics() # 强制触发空闲指标记录
self.kv_events_publisher.publish_kv_events() # 立即发布待处理的 KV 事件
```

test/registered/unit/managers/test_scheduler_pause_generation.py

为 `_new_scheduler` 添加 `metrics_reporter` 和 `kv_events_publisher` 的 Mock 桩，确保新逻辑在测试中不会因属性缺失而失败。

```
# test/registered/unit/managers/test_scheduler_pause_generation.py
# 在 _new_scheduler 方法末尾，return 之前添加以下两行：
scheduler.metrics_reporter = MagicMock()
scheduler.metrics_reporter.current_scheduler_metrics_enabled = False
scheduler.kv_events_publisher = MagicMock()
return scheduler
```

评论区精华

Reviewer `alexnaills` 曾建议将刷新逻辑移到 `in_place` 逻辑之前以获得更大收益，但最终方案保留在末尾，且 `alexnaills` 在后续评论中批准了此 PR。此外，第一次审核时 reviewer 指出调度器已有小幅重构，要求拉取最新代码并更新，PR 作者随后 rebased 并调整了属性名为新的 `metrics_reporter` 和 `kv_events_publisher`。

- 刷新逻辑的位置选择 (design): 位置未调整，但 reviewer 最终批准了 PR，说明当前方案可接受。

风险与影响

- 风险：变更仅限于暂停路径，且依赖已有的 `_maybe_log_idle_metrics` 和 `publish_kv_events`，风险极低。唯一的间接影响是重置 `last_log_time` 可能导致下一次正常空闲日志提前触发，但对监控系统无负面影响。
- 影响：直接解决了权重更新等暂停场景下的监控数据误导问题。对用户透明，但对运维和监控系统有正面影响，确保暂停期间吞吐量为零且 KV 事件及时发出。对性能无影响，新增代码仅在暂停时执行一次。
- 风险标记：低风险

关联脉络

- 暂无明显关联 PR