

PR #23978 完整报告

sgl-project/sclang

Add engine_type label to tokenizer manager metrics

合并时间: 2026-04-29 10:52

原文链接: <http://prhub.com.cn/sgl-project/sclang/pull/23978>

执行摘要

- 一句话: 为 Tokenizer 指标添加 engine_type 标签
- 推荐动作: 值得精读, 展示了如何通过提取公共方法消除重复代码并保持配置一致性, 适合作为代码复用和可观测性增强的参考。

功能与动机

Tokenizer manager metrics 缺少 engine_type 标签, 导致无法在监控面板中按 prefill/decode/unified 角色过滤。PR body 明确指出此问题并通过添加标签解决。

实现拆解

1. 在 DisaggregationMode 枚举中新增 to_engine_type() 静态方法, 将 if/else 逻辑封装为单一调用点。
2. 在 scheduler_metrics_mixin.py 中将原有的内联 if/else 替换为对 DisaggregationMode.to_engine_type() 的调用, 减少代码重复。
3. 在 tokenizer_manager.py 的 init_metric_collector_watchdog 中计算 engine_type 并加入 labels 字典, 使 Tokenizer 指标带上该标签。

关键文件:

- python/sclang/srt/disaggregation/utils.py (模块 调度器; 类别 source; 类型 core-logic ; 符号 to_engine_type) : 新增了 DisaggregationMode.to_engine_type() 静态方法, 将分散的 engine_type 计算逻辑统一集中管理。
- python/sclang/srt/observability/scheduler_metrics_mixin.py (模块 调度器; 类别 source ; 类型 core-logic) : 将原有的 if/else engine_type 计算替换为调用统一方法, 消除重复。
- python/sclang/srt/managers/tokenizer_manager.py (模块 管理器; 类别 source; 类型 core-logic) : 在 Tokenizer 指标初始化时添加 engine_type 标签, 使 Tokenizer 指标可被过滤。

关键符号: DisaggregationMode.to_engine_type

关键源码片段

[python/sclang/srt/disaggregation/utils.py](#)

新增了 `DisaggregationMode.to_engine_type()` 静态方法，将分散的 `engine_type` 计算逻辑统一集中管理。

```
class DisaggregationMode(Enum):
    NULL = "null"
    PREFILL = "prefill"
    DECODE = "decode"

    @staticmethod
    def to_engine_type(mode: str) -> str:
        # 将 disaggregation_mode 字符串转换为 engine_type 标签值
        if mode == DisaggregationMode.PREFILL.value:
            return "prefill"
        elif mode == DisaggregationMode.DECODE.value:
            return "decode"
        # 非 PD 场景默认 unified
        return "unified"
```

评论区精华

无 review 讨论。

- 暂无高价值评论线程

风险与影响

- 风险：风险极低。仅修改标签生成逻辑，不改变指标数值或上报路径。但需注意 `tokenizer_metrics_allowed_custom_labels` 和 `extra_metric_labels` 的交互，避免标签覆盖问题。另外，`fake_bootstrap_room_counter` 被误删（可能与本 PR 无关），已通过 diff 确认，但未在 PR 中说明。
- 影响：影响范围小，仅涉及可观测性模块。对最终用户无影响；对运维人员而言，Tokenizer 指标现在可以通过 `engine_type` 过滤，提升监控灵活性。
- 风险标记：代码删除未说明

关联脉络

- PR #23530 [Spec] Fix spec_accept_rate and unify accept/draft naming: 同样修改了 `scheduler_metrics_mixin.py` 和 `tokenizer_manager.py`，与本 PR 涉及可观测性指标的统一和清理。