

PR #26583 完整报告

sgl-project/sglang

[Utils] Support configure log level at runtime

合并时间: 2026-05-30 05:49

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/26583>

执行摘要

- 一句话: 运行时动态配置日志级别
- 推荐动作: 值得阅读, 展示了一种轻量的跨进程运行时配置模式。对于生产部署, 建议在 HTTP 入口层增加服务端输入校验, 避免仅依赖 CLI 约束。

功能与动机

动机: 在运维现场排查问题时, 能够随时查看 debug 日志而不需要重启服务。PR 作者指出, "It would be useful to check debug log when troubleshooting live server."

实现拆解

1. 在 `io_struct.py` 的 `ConfigureLoggingReq` 数据类中新增 `log_level: Optional[str]` 字段。
2. 在 `tokenizer_manager.py` 的 `configure_logging` 方法中增加对 `obj.log_level` 的判断: 如果不为 `None`, 则调用 `logging.getLogger().setLevel(obj.log_level.upper())` 设置根日志级别; 若级别非法则抛出异常返回给调用者 (HTTP API 会返回错误), 仅当设置成功后才将请求对象通过 `self.send_to_scheduler.send_pyobj(obj)` 转发给 `scheduler` 进程。
3. 在 `scheduler.py` 中新增 `ConfigureLoggingReq` 的分发路由, 并实现 `configure_logging` 方法: 收到后通过 `logging.getLogger().setLevel()` 设置级别, 然后通过 `self.ipc_channels.send_to_detokenizer.send_output` 继续转发给 `detokenizer` 进程。
4. 在 `detokenizer_manager.py` 中新增 `ConfigureLoggingReq` 的处理入口 `handle_configure_logging_req`, 同样设置根日志级别。
5. 在 `configure_logging.py` 的命令行参数中增加 `--log-level`, 可选值为 `debug/info/warning/error/critical`, 将用户的输入填入 `payload["log_level"]`。
6. 更新了 `observability.mdx` 文档, 增加运行时修改日志级别示例。

关键文件:

- `python/sglang/srt/managers/tokenizer_manager.py` (模块 入口管理; 类别 `source`; 类型 `core-logic`; 符号 `configure_logging`): 核心入口: 配置日志级别并决定是否转发给其他进程。
- `python/sglang/srt/managers/scheduler.py` (模块 调度器; 类别 `source`; 类型 `core-logic`; 符号 `configure_logging`): 调度进程处理日志级别设置并转发给 `detokenizer`。

- python/sclang/srt/managers/detokenizer_manager.py (模块 解码器; 类别 source; 类型 core-logic; 符号 handle_configure_logging_req) : 解码器进程设置日志级别。
- python/sclang/srt/managers/configure_logging.py (模块 CLI 工具; 类别 source; 类型 configuration; 符号 main) : CLI 工具增加 --log-level 参数。
- python/sclang/srt/managers/io_struct.py (模块 数据结构; 类别 source; 类型 data-contract; 符号 ConfigureLoggingReq) : 数据类新增 log_level 字段。
- docs_new/docs/advanced_features/observability.mdx (模块 文档; 类别 other; 类型 documentation) : 文档更新, 增加运行时修改日志级别示例。

关键符号: TokenizerManager.configure_logging, Scheduler.configure_logging, DetokenizerManager.handle_configure_logging_req

关键源码片段

python/sclang/srt/managers/tokenizer_manager.py

核心入口: 配置日志级别并决定是否转发给其他进程。

```
def configure_logging(self, obj: ConfigureLoggingReq):
    self.request_logger.configure(
        log_requests=obj.log_requests,
        log_requests_level=obj.log_requests_level,
        log_requests_format=obj.log_requests_format,
    )
    if obj.dump_requests_folder is not None:
        self.dump_requests_folder = obj.dump_requests_folder
    if obj.dump_requests_threshold is not None:
        self.dump_requests_threshold = obj.dump_requests_threshold
    if obj.dump_requests_exclude_meta_keys is not None:
        self.dump_requests_exclude_meta_keys = list(
            obj.dump_requests_exclude_meta_keys
        )
    if obj.crash_dump_folder is not None:
        self.crash_dump_folder = obj.crash_dump_folder
    if obj.log_level is not None:
        # setLevel() 如果 level 是非法字符串会抛出 ValueError。
        # 故意不捕获, 让异常传播到 HTTP API 层,
        # 从而只将合法请求继续转发给 scheduler 和 detokenizer。
        logging.getLogger().setLevel(obj.log_level.upper())
        self.send_to_scheduler.send_pyobj(obj) # 仅当设置成功才转发
    logging.info(f"Config logging: {obj=}")
```

评论区精华

Review 中 gemini-code-assist 指出直接 `setLevel` 可能因非法字符串抛出 `ValueError`, 建议在 `tokenizer`、`scheduler`、`detokenizer` 三处均加 `try-except` 避免进程崩溃。作者 stepinto 回复: 在 `tokenizer` 中故意不捕获异常, 让异常传播回 HTTP 调用者, 从而只将合法请求发送给 `scheduler/detokenizer`; CLI 工具已用 `choices` 限制有效值。最终 reviewer hnyls2002 批准

合并。

- 日志级别异常处理安全性讨论 (security): 保持当前设计: tokenizer 让异常向上传播, 不转发非法请求; scheduler/detokenizer 信任来自 tokenizer 的请求。
- CLI 参数 choices 限制 (design): stepinto 修改代码, 添加 choices 参数, 限制为 debug/info/warning/error/critical。

风险与影响

- 风险: 主要风险在于: 尽管 CLI 限制了参数, 但 HTTP API 接口 /configure_logging 仍可接收任意字符串作为 log_level。如果 tokenizer 进程收到非法级别, 将抛出 ValueError, 未捕获的情况下会返回 500 错误给客户端, 进程不会崩溃 (因为异常在 asyncio 处理器中被捕获)。但若 tokenizer 因其他原因绕过 (如直接通过 IPC 发送请求), scheduler/detokenizer 收到无效级别可能崩溃。当前实现假设所有来自 tokenizer 的请求均已通过合法性校验, 这种 IPC 信任假设在安全性敏感场景下需额外加固。
- 影响: 影响小范围: 用户可在不重启服务时调整日志级别, 便于排查问题。系统无性能影响。团队需注意未来如果 HTTP API 直接暴露, 需增加服务器端校验。
- 风险标记: IPC 信任假设, 缺少服务端输入校验

关联脉络

- 暂无明显关联 PR