

PR #20700 完整报告

sgl-project/sclang

fix(serving_chat): catch TypeError from tojson on Jinja2 Undefined variables

合并时间: 2026-05-23 00:43

原文链接: <http://prhub.com.cn/sgl-project/sclang/pull/20700>

执行摘要

- 一句话: 修复 Jinja2 模板 tojson 未定义变量错误处理
- 推荐动作: 该 PR 值不值得精读取决于是否负责聊天模板或错误处理模块。对于其他开发者而言, 这是一个教科书级的防御性错误处理改进, 值得了解但不必须深入。

功能与动机

当聊天模板对 Jinja2 Undefined 变量使用 tojson 过滤器时 (例如 SGLang 的 fallback 机制移除工具函数的包装器后, 模板仍期望 tool.function), json.dumps 会抛出 TypeError 而非 jinja2.TemplateError。现有的 except jinja2.TemplateError 无法捕获 TypeError, 导致异常传播到 serving_base.py 的通用 except Exception 处理器, 返回 500 内部服务器错误, 而不是预期的 400 请求错误。

实现拆解

1. 修改异常捕获类型: 在 python/sclang/srt/entrypoints/openai/serving_chat.py 的 _apply_jinja_template 方法中, 将第 758 行的 except jinja2.TemplateError as template_error 改为 except (jinja2.TemplateError, TypeError) as template_error。
2. 更新注释: 同步修改注释, 说明新增的 TypeError 场景 (例如 tojson 过滤器作用于 Jinja2 Undefined 变量)。
3. 保持后续逻辑不变: 捕获到异常后仍然抛出 ValueError, 由上层统一处理为 400 Bad Request。

关键文件:

- python/sclang/srt/entrypoints/openai/serving_chat.py (模块 请求路由; 类别 source; 类型 core-logic) : 核心文件: 修改了 _apply_jinja_template 方法的异常捕获类型, 增加 TypeError 以覆盖 tojson 过滤器在 Undefined 变量上的错误。

关键符号: _apply_jinja_template

关键源码片段

[python/sclang/srt/entrypoints/openai/serving_chat.py](#)

核心文件: 修改了 _apply_jinja_template 方法的异常捕获类型, 增加 TypeError 以覆盖 tojson 过滤器在 Undefined 变量上的错误。

```
# python/sglang/srt/entrypoints/openai/serving_chat.py
# 在 _apply_jinja_template 方法的 fallback 分支中
    except (jinja2.TemplateError, TypeError) as template_error:
        # 模板错误 (如 Jinja 模板中的 raise_exception)
        # 以及 TypeError (如 tojson 过滤器作用于 Jinja2 Undefined 变量)
        # 都应视为客户端错误, 返回 400 BadRequest
        raise ValueError(str(template_error)) from template_error
```

评论区精华

无技术争辩。Gemini Code Assist 的自动审查确认该修复正确，JustinTong0323 直接批准合并。

- 暂无高价值评论线程

风险与影响

- 风险：风险极低。仅修改异常捕获路径，不涉及正常请求处理、模型推理或性能关键路径。唯一潜在风险是误抓了预期之外的其他 `TypeError`，但该 `except` 块仅在第一次 `apply_chat_template` 失败后的 fallback 尝试中，且 `TypeError` 在此上下文中极大概率来源于模板错误，整体可控。
- 影响：影响范围：仅影响聊天模板处理中的异常路径。对于使用 `tojson` 过滤器的模板（如 `MiniMax-M2.5`）且因 fallback 导致变量未定义时，返回状态码从 500 变为 400，用户体验更准确。
- 风险标记：暂无

关联脉络

- 暂无明显关联 PR