

PR #22549 完整报告

sgl-project/sglang

Fix broken streaming response with --incremental-streaming-output

合并时间: 2026-04-13 06:05

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/22549>

执行摘要

- 一句话: 修复 --incremental-streaming-output 流式响应文本乱码问题
- 推荐动作: 该 PR 值得精读, 尤其是理解增量流式与累积流式的设计差异。关注 `_generate_chat_stream` 中的条件判断逻辑, 这是修复的核心。同时, 回归测试展示了如何模拟增量流式场景, 对测试编写有参考价值。

功能与动机

根据 PR body 和关联 Issue #22510, 用户在使用 --incremental-streaming-output 参数时, 所有模型的流式响应都会产生乱码文本 (如 "Helloeagesg")。Issue 中报告 Gemma 模型出现 "rge"、"age"、"eable" 等错误分块, 导致最终组装的消息不可用。根本原因是 #21037 (c37200f5e) 变更后, `tokenizer_manager` 在启用 `incremental_streaming_output` 时返回增量文本 (`state.text[last_text_offset:]`), 但 `serving_chat.py` 未相应更新, 仍按累积缓冲区长度切片, 产生垃圾片段。

实现拆解

实现分为两个关键部分: 1. 在 `python/sglang/srt/entrypoints/openai/serving_chat.py` 的 `_generate_chat_stream` 函数中, 添加条件判断: 当 `incremental_streaming_output` 启用时, 直接使用 `content["text"]` 作为增量文本, 否则按原逻辑切片。2. 在 `test/registered/openai_server/basic/test_serving_chat.py` 中添加 `test_incremental_streaming_output_delta` 回归测试, 模拟增量流式场景并验证文本正确组装。

关键文件:

- `python/sglang/srt/entrypoints/openai/serving_chat.py` (模块 `openai_server`): 修复流式响应文本组装的核心逻辑, 根据 `incremental_streaming_output` 标志正确处理增量文本
- `test/registered/openai_server/basic/test_serving_chat.py` (模块 `openai_server_test`): 添加回归测试 `test_incremental_streaming_output_delta`, 确保增量流式场景下文本正确组装, 防止问题复现

关键符号: `_generate_chat_stream`, `test_incremental_streaming_output_delta`

评论区精华

Review 中没有实质性技术讨论，两位 reviewer (alexnaills 和 JustinTong0323) 均直接批准。PR body 中详细分析了根本原因和修复方案，Issue #22510 提供了具体的复现步骤和错误现象。

- 增量流式文本组装逻辑修复 (correctness): 在 `serving_chat` 中添加条件判断，当 `incremental_streaming_output` 启用时直接使用 `content["text"]` 作为增量文本

风险与影响

- 风险：风险较低但需注意：1. 核心路径变更：修改了流式响应的核心逻辑 `_generate_chat_stream`，虽然变更简单，但涉及所有模型的流式输出，需确保非增量流式模式不受影响。2. 条件判断依赖 `server_args.incremental_streaming_output` 标志，需确保该标志在 `tokenizer_manager` 和 `serving_chat` 中一致。3. 回归测试覆盖了增量流式场景，但未测试混合模式（部分请求启用、部分禁用）的并发场景。
- 影响：影响范围：1. 用户影响：修复后，所有使用 `--incremental-streaming-output` 参数的用户将获得正确的流式响应文本，消除乱码问题，提升用户体验。2. 系统影响：仅影响流式聊天响应的文本组装逻辑，不改变模型推理、调度等核心组件。3. 团队影响：作为关键 bugfix，需尽快合并到主分支，避免影响生产环境。
- 风险标记：核心路径变更，条件判断依赖外部标志

关联脉络

- PR #21037 未知（根据 PR body 引用）：PR body 指出该 PR (c37200f5e) 变更了 `tokenizer_manager` 的文本输出逻辑，导致当前 bug 的根本原因
- PR #22567 [tokenizer] eliminate $O(n^2)$ copy in non-incremental streaming: 同样修改了 `tokenizer_manager.py`，涉及流式输出性能优化，可能与当前 PR 的增量流式逻辑相关