

# PR #22606 完整报告

sgl-project/sglang

[serving] replace  $O(n^2)$  stream\_buffer string concat with integer offset

合并时间: 2026-04-15 05:48

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/22606>

## 执行摘要

- 一句话: 优化流式响应性能, 将字符串拼接替换为整数偏移量以消除  $O(n^2)$  开销。
- 推荐动作: 该 PR 值得精读, 展示了如何将  $O(n^2)$  字符串操作优化为  $O(1)$  整数操作的设计决策。重点关注状态表示从字符串到整数的转变, 以及首次块检测逻辑的显式化处理。同时, review 中关于 echo 场景的偏移量修正揭示了流式处理中常见的时间状态跟踪陷阱, 具有借鉴意义。

## 功能与动机

PR body 明确指出原实现在 `serving_chat.py` 和 `serving_completions.py` 中, 流式 delta 计算通过字符串拼接 (`stream_buffer + delta`) 在每个块上累积全文, 导致每次迭代  $O(n)$  开销, 整个流式响应生命周期形成  $O(n^2)$  总复杂度。目标是通过整数偏移量替换字符串缓冲区, 将关键操作优化为  $O(1)$ 。

## 实现拆解

1. 状态跟踪重构: 在 `serving_completions.py` 和 `serving_chat.py` 中, 将 `stream_buffers` 字典 (存储累积字符串) 替换为 `stream_offsets` 字典 (存储整数偏移量)。
2. 首次块检测逻辑更新: 原依赖 `not stream_buffer` (空字符串为假) 检测首次块, 但整数 0 也是假值。改为显式检查 `index not in stream_offsets`, 更清晰且兼容新表示。
3. Delta 计算优化: 原使用 `text[len(stream_buffer):]` 计算增量, 现改为 `text[offset:]`, 直接从当前偏移量切片。偏移更新从 `stream_buffers[index] = stream_buffer + delta` 改为 `stream_offsets[index] = len(content["text"])`, 避免字符串拼接。
4. Echo 处理修正: 根据 review 反馈, 修复了 `serving_completions.py` 中启用 echo 时偏移量计算错误。原提交错误地使用 `len(text)` (包含 echo 文本), 但后续迭代的 `content["text"]` 不包含 echo 文本。最终采纳建议使用 `len(content["text"])` 确保偏移量正确跟踪模型原始输出。
5. 测试配套: PR body 提到“现有流式测试应通过”, 但未包含直接测试文件变更, 依赖现有测试验证行为不变。

关键文件:

- `python/sglang/srt/entrypoints/openai/serving_completions.py` (模块 OpenAI 入口; 类别 source; 类型 core-logic; 符号 `_generate_completion_stream`): OpenAI 补全流式响应的核心入口, 重构了状态跟踪和 delta 计算逻辑。

- python/sglang/srt/entrypoints/openai/serving\_chat.py (模块 OpenAI 入口; 类别 source; 类型 core-logic; 符号 \_generate\_chat\_stream) : OpenAI 聊天流式响应的核心入口, 同步优化了 delta 计算逻辑。

关键符号: \_generate\_completion\_stream, \_generate\_chat\_stream

## 关键源码片段

### python/sglang/srt/entrypoints/openai/serving\_completions.py

OpenAI 补全流式响应的核心入口, 重构了状态跟踪和 delta 计算逻辑。

```
async def _generate_completion_stream(
    self,
    adapted_request: GenerateReqInput,
    request: CompletionRequest,
    raw_request: Request,
) -> AsyncGenerator[str, None]:
    """Generate streaming completion response"""
    created = int(time.time())

    # State tracking for streaming
    stream_offsets = {} # 替换原 stream_buffers = {}, 存储整数偏移量而非累积字符串
    n_prev_tokens = {}

    # ... 其他状态初始化 ...

    async for content in self.tokenizer_manager.generate_request(
        adapted_request, raw_request
    ):
        index = content.get("index", 0)
        text = content["text"]
        # ... 其他元数据提取 ...

        is_first_chunk = index not in stream_offsets # 显式检测首次块, 替代原 not stream_buffer
        offset = stream_offsets.get(index, 0) # 获取当前偏移量, 默认为0

        # Handle echo for first chunk
        if is_first_chunk: # 原条件为 if not stream_buffer:
            if request.echo:
                echo_text = self._get_echo_text(request, index)
                text = echo_text + text # echo文本仅在第一块本地预置

        # ... logprobs处理 ...

        # Generate delta
        delta = text[offset:] # 直接切片, 替代原 text[len(stream_buffer):]
        stream_offsets[index] = len(content["text"]) # 存储模型原始文本长度, 确保偏移量正确跟踪
        # 注意: 使用 content["text"] 而非 text, 避免echo文本干扰
        finish_reason = content["meta_info"].get("finish_reason", None)
```

```
# ... 后续流式响应构建 ...
```

## python/slang/srt/entrypoints/openai/serving\_chat.py

OpenAI 聊天流式响应的核心入口，同步优化了 delta 计算逻辑。

```
async def _generate_chat_stream(
    self,
    adapted_request: GenerateReqInput,
    request: ChatCompletionRequest,
    raw_request: Request,
) -> AsyncGenerator[str, None]:
    """Generate streaming chat completion response"""
    # ... 解析器初始化 ...

    # State tracking for streaming
    is_firsts = {}
    stream_offsets = {} # 替换原 stream_buffers = {}
    n_prev_tokens = {}
    has_tool_calls = {}
    finish_reasons = {}

    # ... 其他状态初始化 ...

    async for content in self.tokenizer_manager.generate_request(
        adapted_request, raw_request
    ):
        index = content.get("index", 0)
        # ... 元数据提取和首次块角色处理 ...

        offset = stream_offsets.get(index, 0) # 获取偏移量，替代原 stream_buffer = stream_
        buffers.get(index, "")
        if self.tokenizer_manager.server_args.incremental_streaming_output:
            delta = content["text"] # 增量流模式下，text直接作为delta
        else:
            delta = content["text"][offset:] # 非增量流模式，切片提取增量
        stream_offsets[index] = len(content["text"]) # 更新偏移量为当前文本全长

    # ... 后续推理内容处理和流式响应构建 ...
```

## 评论区精华

review 中主要讨论集中在正确性风险：

- gemini-code-assist[bot] 指出高严重性 bug：当启用 request.echo 时，偏移量计算错误地包含了预置的 echo 文本长度。因为 echo 文本仅在第一块本地预置，后续迭代的模型输出 content["text"] 中不包含它，导致偏移量过大，后续块会跳过部分生成内容。
- 解决方案：建议使用 len(content["text"]) 而非 len(text) 存储偏移量，确保正确跟踪模型生成输出进度。该建议被采纳，最终提交修正了此问题。

- hnyls2002 批准合并，未提出额外疑虑。
- Echo 场景下偏移量计算错误 (correctness): 采纳建议，改为 `stream_offsets[index] = len(content["text"])`，确保偏移量正确跟踪模型原始输出。

## 风险与影响

- 风险：1. 回归风险：核心风险是流式响应内容截断或重复。review 中已识别并修复了 echo 场景下的偏移量计算错误，但其他边缘场景（如增量流输出模式 `incremental_streaming_output`、多索引并发）仍需现有测试覆盖。2. 性能风险：变更旨在提升性能，但整数操作替代字符串拼接，在极短文本场景下微优化收益可能不明显，不过无负面性能影响。3. 兼容性风险：行为应保持不变，但首次块检测逻辑从隐式空字符串检查改为显式字典成员检查，逻辑等价但更健壮，无兼容性破坏。4. 测试覆盖不足：PR 未添加新测试，依赖现有流式测试。若现有测试不充分，可能遗漏边界条件。
- 影响：1. 用户影响：流式响应客户端应感知不到变化，但服务端性能提升，尤其长文本生成场景可减少内存开销和 CPU 时间。2. 系统影响：优化核心流式路径，降低  $O(n^2)$  复杂度为  $O(1)$  操作，提升系统可扩展性，减少长会话资源消耗。3. 团队影响：代码更清晰，首次块检测逻辑显式化，便于后续维护。为类似流式处理场景提供性能优化范式。
- 风险标记：核心路径变更，正确性风险已修复，依赖现有测试覆盖

## 关联脉络

- PR #22755 Rename `_alive_streaming_session_count`; use `_is_streaming` helper: 同属流式会话管理和状态跟踪优化，涉及 `scheduling` 和 `consistency` 标签。
- PR #22753 Fix streaming session busy-check double-counting via `active_pool_idx`: 同属流式会话相关 bugfix，涉及 `scheduling` 和 `kv-cache` 模块。