

PR #42683 完整报告

vllm-project/vllm

[Bugfix][Frontend] streaming tool-call serializer drops first args chunk when name and args share a DeltaMessage

合并时间: 2026-05-28 13:20

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/42683>

执行摘要

- 一句话: 修复 streaming tool-call 序列化器丢弃第一个 args chunk 的 bug
- 推荐动作: 值得阅读, 展示了一个高质量的 bugfix 设计: 通过引入 `split_delta` 函数将复合 delta 原子化, 简化了状态机设计, 同时通过分组按 index 保持 name 和 args 的关联。测试覆盖完善。可以学习这种将复杂问题分解的思路。

功能与动机

根据 PR body, 在 kimi 2.6 通过 `/v1/responses (streaming)` 部署时, 约 60% 的工具调用产生的 `done.arguments` 无法通过 `json.loads` 解析, 因为开头的 `{` 块从 SSE 流中丢失。这是因为 `emit_delta` 只检查 `delta_message.tool_calls[0]`, 而像 `KimiK2ToolParser` 的 `extract_tool_calls_streaming` 会将 name 和 args 分成两个 `DeltaToolCall` 放在同一个 `DeltaMessage` 中, 导致第一个 args chunk 被静默丢弃。

实现拆解

1. 新增 `split_delta` 函数 (`streaming_events.py`): 接受一个 `DeltaMessage`, 根据它包含的字段 (`reasoning`, `content`, `tool_calls`) 分解为若干原子 `DeltaMessage`, 每个只包含一个字段。- 对于 `tool_calls`, 按 index 分组, 同一 index 的多个 `DeltaToolCall` 合并到一个 `DeltaMessage` 中, 确保 name 和 args 保持在同一原子 delta 中。- 如果只有一个字段且 `tool_calls` index 唯一, 则直接返回原 delta。
2. 简化 `emit_delta` 方法 (`streaming_events.py`): 移除原来 `reasoning`→`content` 的特殊处理逻辑, 因为这个场景现在由 `split_delta` 在外部保证不会出现复合字段。
3. 集成到流式循环 (`serving.py`): 在 `_process_simple_streaming_events` 中, 对每个 `delta_message` 先调用 `split_delta`, 然后对每个分解后的原子 delta 分别执行状态机步骤 (`resolve_target_state`, `needs_transition`, `emit_delta`)。
4. 测试配套:
 - 新增 `test_streaming_events.py`: 包含 `TestSplitDelta` 测试 `split_delta` 纯函数, 以及 `TestProcessorCompoundDeltas` 测试通过 `SimpleStreamingEventProcessor` 处理复合 delta 的集成。
 - 在 `test_serving_responses.py` 中新增 `test_compound_content_and_tool_name_args_same_delta` 端到端测试, 验证 `content` 和 `tool name/args` 在同一个 `DeltaMessage` 中

的正确事件序列。

关键文件：

- `vllm/entrypoints/openai/responses/streaming_events.py` (模块 流式事件; 类别 `source`; 类型 `core-logic`; 符号 `split_delta`) : 核心变更文件: 新增 `split_delta` 函数, 修改 `emit_delta` 方法, 移除特殊处理逻辑。
- `vllm/entrypoints/openai/responses/serving.py` (模块 服务编排; 类别 `source`; 类型 `core-logic`) : 集成 `split_delta` 到流式循环中, 修改控制流, 是 `bugfix` 的另一个关键部分。
- `tests/entrypoints/openai/responses/test_streaming_events.py` (模块 流式测试; 类别 `test`; 类型 `test-coverage`; 符号 `_make_tool_call`, `TestSplitDelta`, `test_all_three_fields`, `test_tool_calls_grouped_by_index`) : 新增测试文件, 覆盖 `split_delta` 函数和 `SimpleStreamingEventProcessor` 处理复合 `delta` 的集成。
- `tests/entrypoints/openai/responses/test_serving_responses.py` (模块 服务测试; 类别 `test`; 类型 `test-coverage`; 符号 `test_compound_content_and_tool_name_args_same_delta`) : 新增端到端集成测试, 验证真实流式循环中 `content` 和 `tool name/args` 在同一个 `DeltaMessage` 中的正确处理。

关键符号: `split_delta`, `SimpleStreamingEventProcessor.emit_delta`,
`OpenAIServingResponses._process_simple_streaming_events`

关键源码片段

`vllm/entrypoints/openai/responses/streaming_events.py`

核心变更文件: 新增 `split_delta` 函数, 修改 `emit_delta` 方法, 移除特殊处理逻辑。

```
# 文件: vllm/entrypoints/openai/responses/streaming_events.py

def split_delta(delta: DeltaMessage) -> list[DeltaMessage]:
    # 将复合 `DeltaMessage` 分解为原子 delta, 每个仅含一个字段。
    # Responses API 要求每个 SSE 事件只有一种类型, 因此在进入状态机之前需拆分。
    # 固定顺序: reasoning -> content -> tool_calls (按 index 分组)。
    has_reasoning = delta.reasoning is not None
    has_content = delta.content is not None
    has_tools = bool(delta.tool_calls)
    parts = int(has_reasoning) + int(has_content) + int(has_tools)

    # 如果只有一个字段, 且 tool_calls 的 index 不冲突 ( $\leq 1$  个不同 index)
    if parts <= 1 and (
        not has_tools
        or len({tc.index for tc in delta.tool_calls if tc.index is not None}) <= 1
    ):
        return [delta]

    deltas: list[DeltaMessage] = []
    if has_reasoning:
        deltas.append(DeltaMessage(reasoning=delta.reasoning))
    if has_content:
```

```

    deltas.append(DeltaMessage(content=delta.content))
if has_tools:
    # 按 index 分组, 同一 index 的 tool_calls 合并在一起
    groups: dict[int | None, list[DeltaToolCall]] = {}
    for tc in delta.tool_calls:
        groups.setdefault(tc.index, []).append(tc)
    for tcs in groups.values():
        deltas.append(DeltaMessage(tool_calls=tcs))
return deltas or [delta]

```

```
class SimpleStreamingEventProcessor:
```

```
    # ... 其他代码 ...
```

```

def emit_delta(
    self,
    delta_message: DeltaMessage,
    output: Any | None,
    get_logprobs: Callable[..., list[Logprob] | None] | None = None,
) -> list[StreamingResponsesResponse]:
    # 发射当前状态的增量事件。现在假设 delta_message 已是原子 delta。
    handlers = self._STATE_HANDLERS[self.state.current_state]
    if self.state.current_state == _StateType.TOOL_CALL:
        assert delta_message.tool_calls is not None
        # 遍历所有 tool_calls (可能包含同一 index 的 name 和 args)
        for tc in delta_message.tool_calls:
            if tc.function is not None and tc.function.arguments:
                return handlers.delta_fn(self.state, tc.function.arguments)
        return []
    # 对于 CONTENT / REASONING 状态, 直接取对应字段
    field = (
        delta_message.content
        if self.state.current_state == _StateType.CONTENT
        else delta_message.reasoning
    )
    if field:
        logprobs = get_logprobs(output) if get_logprobs else []
        return handlers.delta_fn(self.state, field, logprobs)
    return []

```

vllm/entrypoints/openai/responses/serving.py

集成 `split_delta` 到流式循环中, 修改控制流, 是 bugfix 的另一个关键部分。

```
# 文件 : vllm/entrypoints/openai/responses/serving.py
```

```
# 在 _process_simple_streaming_events 方法中
```

```
# 原来 : 直接使用 delta_message
```

```
# 现在 : 先拆分再逐个处理
```

```
for dm in split_delta(delta_message):
```

```
target_state, tool_call = processor.resolve_target_state(dm)
if target_state == _StateType.NONE:
    continue

if processor.needs_transition(target_state, tool_call):
    for event in processor.close_current():
        yield _increment_sequence_number_and_return(event)
    for event in processor.open(target_state, tool_call):
        yield _increment_sequence_number_and_return(event)

for event in processor.emit_delta(dm, output, _get_logprobs):
    yield _increment_sequence_number_and_return(event)
```

评论区精华

Review 中 Andreas Karatzas 建议将新测试文件 `test_streaming_events.py` 的内容合并到 `test_serving_responses.py` 的 `TestAutoToolStreaming` 中。sfeng33 回应保留了该文件，理由是对 `split_delta` 作为纯函数进行独立单元测试与集成验证 wiring 有不同目的。最终 Andreas 同意保留两份测试。

- 测试文件结构：是否将新测试合并到现有文件 (design): 双方同意保留两份测试文件，各自覆盖不同层次。

风险与影响

- 风险：主要风险在于 `split_delta` 改变了 `DeltaMessage` 的处理流程，可能对尚未覆盖的 tool parser 或状态机行为产生影响。但测试覆盖了关键路径（复合字段、并行 tool calls、reasoning→content 转换），且 bugfix 本身增加的是通用拆分逻辑，风险可控。没有发现性能或安全风险。
- 影响：直接影响使用 `/v1/responses streaming endpoint` 且 tool parser 会在同一 `DeltaMessage` 中发送 `name` 和 `args` 独立 `DeltaToolCall` 的场景（如 `KimiK2ToolParser`）。修复后，这些场景的工具调用参数不再丢失，SSE 流中 `function_call_arguments.delta` 序列完整。对其他场景无负面影响。
- 风险标记：暂无

关联脉络

- 暂无明显关联 PR