

# PR #42977 完整报告

vllm-project/vllm

[Parser] Migrate `ResponsesParser` to unified `Parser` interface

合并时间: 2026-06-02 16:50

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/42977>

## 执行摘要

- 一句话: 迁移 ResponsesParser 到统一 Parser 接口
- 推荐动作: 该 PR 值得精读, 展示了将遗留组件迁移到统一接口的典型模式。重点关注 ResponsesParser.process() 的简化和流式处理的缓存设计 (未来改进方向)。建议合并后尽快跟进流式问题的 CI 测试。

## 功能与动机

PR body 指出: 当前 ResponsesParser 直接实例化独立的 ReasoningParser 和 ToolParser, 绕过统一 Parser 类, 而 streaming 路径已正确使用统一 Parser, 造成同一 API 表面两套代码路径使用不同解析器接口。此外, ParsableContext 中 tool\_choice 模式 (forced, required) 未正确处理, 此 PR 一并修复。对应 RFC #32713 的 TODO。

## 实现拆解

1. 合并解析器参数: 在 ResponsesParser.\_\_init\_\_() 中, 将 reasoning\_parser\_cls 和 tool\_parser\_cls 合并为单一的 parser\_cls: type[Parser] | None; 新增 enable\_auto\_tools 和 tool\_call\_id\_type 配置参数。
2. 委托统一解析: ResponsesParser.process() 方法改为调用 self.parser\_instance.extract\_response\_outputs() 一次性获取所有输出项 (reasoning + tool calls), 替代原先分别调用 reasoning\_parser\_instance.extract\_reasoning() 和 tool\_parser\_instance.extract\_tool\_calls() 并手动构造 ResponseReasoningItem 和 tool call 的逻辑。
3. 适配 ParsableContext: 在 responses/context.py 中将 ParsableContext.\_\_init\_\_() 的参数替换为 parser\_cls, 并移除对 reasoning\_parser\_cls 非空的校验; 将 tool\_parser\_cls 的引用替换为 parser\_cls.tool\_parser\_cls 的访问。
4. 调整 serving 层: 在 responses/serving.py 的 \_create\_responses 中创建 ParsableContext 时直接传递 parser\_cls=self.parser, 不再解构; 在 \_render\_next\_turn 中获取工具解析器时改为通过 context.parser\_cls.tool\_parser\_cls。
5. 新增测试: test\_responses\_parser\_unified.py 包含 12 个单元测试, 通过构造 \_NoOpParser、\_ReasoningOnlyParser、\_ToolCallingParser 等桩类, 验证不同场景下统一 Parser 接口的正确性。

关键文件:

- `vllm/entrypoints/openai/parser/responses_parser.py` (模块 解析器; 类别 source; 类型 core-logic; 符号 `ResponsesParser`, `get_responses_parser_for_simple_context`) : 核心改造文件, 将 `ResponsesParser` 从独立使用 `ReasoningParser/ToolParser` 迁移到统一 `Parser` 接口, 是本次 PR 的核心逻辑变更
- `tests/entrypoints/openai/test_responses_parser_unified.py` (模块 测试; 类别 test; 类型 test-coverage; 符号 `_NoOpParser`, `_ReasoningOnlyParser`, `_ToolCallingParser`, `test_responses_parser_noop`) : 新增测试文件, 包含 12 个单元测试和多个桩解析器类, 验证统一 `Parser` 接口的正确性和兼容性
- `vllm/entrypoints/openai/responses/context.py` (模块 上下文; 类别 source; 类型 dependency-wiring; 符号 `ParsableContext`) : 接口适配文件, `ParsableContext` 构造函数签名更改以接收统一 `parser_cls`, 并传递新参数
- `vllm/entrypoints/openai/responses/serving.py` (模块 请求路由; 类别 source; 类型 core-logic; 符号 `_create_responses`, `_render_next_turn`) : 调用方调整, 创建 `ParsableContext` 时直接传递统一 `parser_cls`, 并在下一轮渲染中通过 `parser_cls` 获取工具解析器

关键符号: `ResponsesParser.init`, `ResponsesParser.process`, `get_responses_parser_for_simple_context`, `ParsableContext.init`, `_create_responses`, `_render_next_turn`

## 关键源码片段

### `vllm/entrypoints/openai/parser/responses_parser.py`

核心改造文件, 将 `ResponsesParser` 从独立使用 `ReasoningParser/ToolParser` 迁移到统一 `Parser` 接口, 是本次 PR 的核心逻辑变更

```
# ----- vllm/entrypoints/openai/parser/responses_parser.py ( 关键变更 ) -----
```

```
class ResponsesParser:
    """Incremental parser over completion tokens with reasoning support."""

    def __init__(
        self,
        *,
        tokenizer: TokenizerLike,
        # 原先的参数是 reasoning_parser_cls 和 tool_parser_cls,
        # 现在统一为 parser_cls, 由它接管所有的推理和工具调用解析
        parser_cls: type[Parser] | None,
        response_messages: list[ResponseInputOutputItem],
        request: ResponsesRequest,
        chat_template: str | None,
        chat_template_content_format: ChatTemplateContentFormatOption,
        enable_auto_tools: bool = False, # 新增: 控制自动工具调用行为
        tool_call_id_type: str = "random", # 新增: 工具调用 ID 生成策略
    ):
        self.response_messages = response_messages
```

```

self.num_init_messages = len(response_messages)
self.tokenizer = tokenizer
self.request = request

self.parser_instance: Parser | None = None
if parser_cls is not None:
    chat_template_kwargs = _effective_chat_template_kwargs(
        request,
        chat_template=chat_template,
        chat_template_content_format=chat_template_content_format,
    )
    # 实例化统一的 Parser, 不再区分 reasoning / tool
    self.parser_instance = parser_cls(
        tokenizer,
        tools=request.tools,
        chat_template_kwargs=chat_template_kwargs,
    )

self.enable_auto_tools = enable_auto_tools
self.tool_call_id_type = tool_call_id_type
self.finish_reason: str | None = None

def process(self, output: CompletionOutput) -> "ResponsesParser":
    self.finish_reason = output.finish_reason

    if self.parser_instance is not None:
        # 统一委托给 Parser.extract_response_outputs(),
        # 一次调用即可获得所有 output items (reasoning + tool calls)
        output_items = self.parser_instance.extract_response_outputs(
            model_output=output.text,
            model_output_token_ids=output.token_ids,
            request=self.request,
            enable_auto_tools=self.enable_auto_tools,
            tool_call_id_type=self.tool_call_id_type,
        )
        self.response_messages.extend(output_items)
    else:
        # 没有 parser 时, 直接作为纯文本输出
        if output.text:
            self.response_messages.append(
                ResponseOutputMessage(
                    type="message",
                    id=f"msg_{random_uuid()}",
                    status="completed",
                    role="assistant",
                    content=[
                        ResponseOutputText(
                            annotations=[],
                            type="output_text",

```

```

        text=output.text,
    )
    ],
)
)
return self

```

## tests/entrypoints/openai/test\_responses\_parser\_unified.py

新增测试文件，包含 12 个单元测试和多个桩解析器类，验证统一 Parser 接口的正确性和兼容性

```
# ---- tests/entrypoints/openai/test_responses_parser_unified.py ----
```

```

class _NoOpParser(DelegatingParser):
    """Parser that extracts no reasoning and no tool calls."""
    def is_reasoning_end(self, input_ids: list[int]) -> bool:
        return False
    def extract_content_ids(self, input_ids: list[int]) -> list[int]:
        return input_ids
    def extract_reasoning(self, model_output, request):
        return None, model_output
    def extract_reasoning_streaming(self, *args, **kwargs):
        return None
    def extract_tool_calls(self, model_output, request):
        return ExtractedToolCallInformation(
            tools_called=False, tool_calls=[], content=model_output)
    def extract_tool_calls_streaming(self, *args, **kwargs):
        return None
    def parse_delta(self, *args, **kwargs) -> DeltaMessage | None:
        return None

```

```

class _ToolCallingParser(DelegatingParser):
    """Parser that always returns a hardcoded tool call."""
    def __init__(self, tokenizer, *args, **kwargs):
        super().__init__(tokenizer)
        # 内部包装了一个简单的 stub tool parser
        self._tool_parser = _StubToolParser()
    # ... 其他方法略，关键在 extract_tool_calls 返回固定工具调用

```

```

def test_responses_parser_noop() -> None:
    """Verify that a no-op parser returns only the plain text content."""
    request = MagicMock(spec=ResponsesRequest)
    request.tools = None
    request.tool_choice = "auto"
    request.max_output_tokens = 100

    parser = get_responses_parser_for_simple_context(
        tokenizer=dummy_tokenizer(),

```

```

    parser_cls=_NoOpParser,
    response_messages=[],
    request=request,
    chat_template=None,
    chat_template_content_format="string",
)
output = CompletionOutput(
    index=0, text="Hello world", token_ids=[], finish_reason="stop")
parser.process(output)

# 预期只有一条文本消息，无 reasoning 和 tool calls
assert len(parser.response_messages) == 1
msg = parser.response_messages[0]
assert msg.type == "message"
assert msg.content[0].text == "Hello world"

```

## vllm/entrypoints/openai/responses/context.py

接口适配文件，ParsableContext 构造函数签名更改以接收统一 parser\_cls，并传递新参数

```

# ---- vllm/entrypoints/openai/responses/context.py -----

class ParsableContext(ConversationContext):
    def __init__(
        self,
        *,
        response_messages: list[ResponseInputOutputItem],
        tokenizer: TokenizerLike,
        # 原先使用 reasoning_parser_cls + tool_parser_cls 两个参数,
        # 现在合并为一个 parser_cls, 与 ResponsesParser 保持一致
        parser_cls: type[Parser] | None,
        request: ResponsesRequest,
        available_tools: list[str] | None,
        chat_template: str | None,
        chat_template_content_format: ChatTemplateContentFormatOption,
        enable_auto_tools: bool = False, # 新增
        tool_call_id_type: str = "random", # 新增
    ):
        self.num_prompt_tokens = 0
        self.num_output_tokens = 0
        self.num_cached_tokens = 0
        self.num_reasoning_tokens = 0
        self.all_turn_metrics: list[TurnMetrics] = []

        # 直接传入 parser_cls, 不再做非空校验 (之前要求必须提供 reasoning_parser_cls)
        self.parser = get_responses_parser_for_simple_context(
            tokenizer=tokenizer,
            parser_cls=parser_cls,
            response_messages=response_messages,
            request=request,

```

```
chat_template=chat_template,  
chat_template_content_format=chat_template_content_format,  
enable_auto_tools=enable_auto_tools,  
tool_call_id_type=tool_call_id_type,  
)  
self.parser_cls = parser_cls  
self.request = request  
# ... 其余初始化不变
```

## 评论区精华

`gemini-code-assist[bot]` 指出两个关键问题：

1. 流式处理风险：`extract_response_outputs` 是 stateless 的，期望完整输出，在流式增量调用中可能无法正确识别跨 delta 的标签，建议改用 `parse_delta`。作者回应该问题在重构前已存在，承诺添加 TODO 跟踪。
2. 参数缺失：`enable_auto_tools` 在调用时被硬编码为 `True`，`tool_call_id_type` 完全缺失，可能忽略引擎全局配置。作者在后续 commit (`9f5721e6`) 中通过构造函数参数传递修复了此问题。

最终维护者 `chaunceyjiang` 审核后 LGTM。

- 流式处理风险：`extract_response_outputs` 不适合增量调用 (correctness)：作者承认问题在重构前已存在，承诺添加 TODO 跟踪，后续 PR 再处理。
- 参数缺失：`enable_auto_tools` 和 `tool_call_id_type` 硬编码或缺失 (correctness)：作者在后续 commit (`9f5721e6`) 中通过构造函数参数传递修复了此问题，确保这两个值从 `ResponsesParser` 构造函数传入。

## 风险与影响

• 风险：

1. 流式处理风险：`extract_response_outputs` 在流式场景下可能误判跨 delta 的标签，导致 reasoning 或 tool call 解析错误。此风险在重构前已存在，但统一接口使其更容易暴露。
1. 参数依赖：新增 `enable_auto_tools` 和 `tool_call_id_type` 参数，若调用方未正确传递可能导致行为不一致。当前调用链已串联，但外部使用者若绕过 `ResponsesParser` 直接构造 `ParsableContext` 可能遗漏。
2. 回归风险：`ParsableContext` 实验性功能默认关闭，但签名变更可能影响内部其他实验性调用点。代码搜索未发现其他调用者，风险较低。
3. 测试覆盖不足：新增单元测试未覆盖流式场景，也未测试 `tool_choice` 的 forced/required 模式。  
- 影响：用户影响：无直接用户可见行为变化（默认 `SimpleContext` 不变）。对于启用 `VLLM_USE_EXPERIMENTAL_PARSER_CONTEXT` 的用户，`tool_choice` 的 forced/required 模式现在正确生效。

系统影响：减少解析器接口冗余，统一 Responses API 内部解析路径，后续新增解析器只需实现 `Parser` 接口即可同时支持 chat 和 responses。

团队影响：降低维护成本，代码更容易理解；但引入了流式处理的潜在问题需后续解决。

- 风险标记：流式处理风险，参数依赖，回归风险，实验性功能

## 关联脉络

- PR #44267 [Refactor] Unify reasoning + tool-call parsing behind Parser.parse(): 同一系列的统一解析器重构，将推理和工具调用解析合并到同一个 Parser 接口，本 PR 就是将此统一接口应用到 ResponsesParser。
- PR #44017 [Refactor] Move unstreamed tool-arg flush from serving layer to parser: 也涉及 parser 层重构，将工具参数冲刷逻辑从 serving 层迁移到 parser，与本 PR 的解析器统一目标一致。