

# PR #41991 完整报告

vllm-project/vllm

[Bugfix][Gemma4] Fix infinite loop and array boundary issues in tool parser

合并时间: 2026-05-09 05:24

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/41991>

## 执行摘要

- 一句话: 修复 Gemma4 工具解析器死循环和数组边界问题
- 推荐动作: 建议精读此 PR, 尤其是零进度保护的防御性编码风格, 适用于类似自定义解析器的健壮性提升。

## 功能与动机

通过对抗性输入模糊测试发现 Gemma4 工具解析器在 `_parse_gemma4_args` 和 `_parse_gemma4_array` 中存在三个健壮性问题, 包括无限循环和解析错误, 影响工具调用的可靠性。

## 实现拆解

1. 修复嵌套数组字符串跳转: 在 `_parse_gemma4_array` 的嵌套数组分支 (`elif arr_str[i] == "["`) 中, 添加了字符串定界符 (`STRING_DELIM`) 检测和跳过逻辑, 与已有的嵌套对象分支对称, 避免括号出现在字符串内时导致计数器错误。
2. 添加零进度保护: 在 `_parse_gemma4_args` 和 `_parse_gemma4_array` 的纯值分支中, 添加 `if i == val_start` 保护, 当游标未前进时记录警告并跳出循环, 防止恶意输入导致死循环。
3. 新增回归测试: 在 `tests/tool_parsers/test_gemma4_tool_parser.py` 中添加三个带超时标记的测试用例, 分别覆盖嵌套数组括号在字符串内、纯值分支死循环等场景。

关键文件:

- `vllm/tool_parsers/gemma4_tool_parser.py` (模块 工具解析器; 类别 `source`; 类型 `core-logic`; 符号 `_parse_gemma4_args`, `_parse_gemma4_array`): 核心解析逻辑, 包含所有修复: 嵌套数组字符串跳转和两个纯值分支零进度保护。
- `tests/tool_parsers/test_gemma4_tool_parser.py` (模块 工具解析器测试; 类别 `test`; 类型 `test-coverage`; 符号 `test_malformed_partial_array`, `test_string_element_with_closing_bracket`, `test_stray_closing_bracket`): 新增三个回归测试用例, 覆盖修复的三种场景, 带超时标记确保死循环被检测。

关键符号: `_parse_gemma4_args`, `_parse_gemma4_array`

## 关键源码片段

[vllm/tool\\_parsers/gemma4\\_tool\\_parser.py](#)

核心解析逻辑，包含所有修复：嵌套数组字符串跳转和两个纯值分支零进度保护。

```
# _parse_gemma4_array 中嵌套数组分支（新增字符串跳过）
# 在 while 循环中，先检查字符串定界符，跳过后再处理括号
if arr_str[i:].startswith(String_Delimiter):
    i += len(String_Delimiter)
    nd = arr_str.find(String_Delimiter, i)
    i = nd + len(String_Delimiter) if nd != -1 else n
    continue

# 纯值分支零进度保护（两个函数各一处）
val_start = i
while i < n and arr_str[i] not in ("(", ")", "{", "}"):
    i += 1
# ... 省略原有 partial 分支 ...
if i == val_start:
    logger.warning("Gemma4 parser made no progress at position %d; aborting on malformed
input.", i)
    break
```

## 评论区精华

审查者 bbrowning 确认测试在变更前失败、变更后通过，并指出纯值分支的保护在某些路径可能不可达，但作为防御性编码保留是合理的。最终批准合并。

- 暂无高价值评论线程

## 风险与影响

- 风险：变更集中于解析器边界条件处理，影响范围小。风险在于防御性保护可能在某些合法但罕见的流式输入中提前终止解析，但超时保护可避免死循环。
- 影响：影响 Gemma4 模型工具解析功能，修复潜在无限循环和解析错误。对正常输入无影响，仅提高对异常输入 / 对抗性样本的鲁棒性。
- 风险标记：防御性编码，边界路径变化

## 关联脉络

- PR #39311 未知（PR body 提及）：PR body 指出该 PR 部分解决了嵌套数组字符串定界符问题，但此 PR 合并了更多修复和测试。