

# PR #38127 完整报告

vllm-project/vllm

Various Transformers v5 fixes

合并时间: 2026-03-26 08:10

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/38127>

## 执行摘要

本次 PR 修复了 Transformers v5 版本中的多个配置错误，包括处理 sliding window 设置、清理过时 Olmo3 配置和防止 DeepSeekVL2 参数传递问题。变更影响有限，主要提升代码库兼容性和维护性，适合关注配置处理的工程师快速了解修复模式。

## 功能与动机

变更旨在解决 Transformers v5 升级带来的兼容性问题：

- 修复 Qwen1.5 测试模型中 sliding\_window 设置为 0 的错误，vLLM 使用 None 表示禁用以避免模型长度计算错误。
- 移除 vendored Olmo3 配置，因为该模型已在 Transformers v4.57.0 正式发布，减少代码重复（参考 PR #24534）。
- 防止 DeepSeekVL2Config 传递无效 kv\_lora\_rank 值给 DeepSeekV2Config，避免配置冲突。
- 更新离线模式测试，添加 image\_processing\_utils\_fast 别名以支持新模块。

## 实现拆解

按模块拆解关键改动：

1. 测试模块：在 tests/entrypoints/offline\_mode/test\_offline\_mode.py 中添加正则表达式别名 `r'.+\.\.image_processing_utils_fast$'`，扩展离线测试覆盖。
2. 核心配置：在 vllm/config/model.py 的 `__post_init__` 方法中新增代码块：

```
python if self.get_sliding_window() == 0: self.disable_sliding_window = True self.hf_text_config.sliding_window = None
```

 确保在 `get_and_verify_max_len` 前转换 sliding\_window 值，防止 max\_model\_len 错误计算。
3. 模型配置清理：移除 `vllm/transformers_utils/configs/olmo3.py` 文件，并更新相关导入和注册，例如在 `vllm/model_executor/models/olmo2.py` 中将导入改为 `from transformers import Olmo2Config, Olmo3Config`。
4. DeepSeekVL2 修复：在 `vllm/transformers_utils/configs/deepseek_vl2.py` 的 `__init__` 方法中添加逻辑，检查 language\_config 字典并移除 kv\_lora\_rank 键如果其值为 None。

## 评论区精华

review 讨论中仅有一条高亮评论：

gemini-code-assist[bot] 指出: "Modifying the language\_config dictionary in-place can lead to unexpected side effects for the caller if they reuse the kwargs dictionary. It's safer to work with a copy of the dictionary."

该评论建议使用字典副本来避免潜在副作用，但未被采纳，代码保持原样合并，揭示了配置处理中的设计权衡。

## 风险与影响

风险：

- deepseek\_vl2.py 中直接修改传入字典可能影响调用者，如果 kwargs 被重用。
- sliding\_window 转换逻辑需确保时序正确，否则 max\_model\_len 可能计算错误。
- 移除 Olmo3 文件需验证无残留依赖，以避免导入错误。

影响：

- 对用户透明，无直接影响。
- 改善开发者体验，提升代码库与上游 Transformers 的兼容性，减少维护负担。
- 风险较低，变更集中于配置和测试代码。

## 关联脉络

与历史 PR 的关联主要体现在引用 PR #24534 的讨论，该 PR 可能涉及 Olmo3 配置的早期处理，表明本次变更是清理工作的延续。近期历史 PR 中未见直接相关项，凸显了本次修复的独立性和维护性质。