

# PR #42162 完整报告

vllm-project/vllm

Fix Molmo2 image token metadata

合并时间: 2026-05-11 09:19

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/42162>

## 执行摘要

- 一句话: 修复 Molmo2 图像令牌元数据与 HF 处理器不匹配
- 推荐动作: 值得精读的设计决策: 如何将 HF 处理器配置参数传递到 vLLM 的底层令牌生成函数, 保持向后兼容的默认值。讨论中关于测试验证的方法 (TDD, 回归测试确认) 值得借鉴。

## 功能与动机

修复 Issue #38660 中 MolmoWeb 模型在 Triton 注意力内核的 CUDA 断言错误。根本原因是 `_merge_multimodal_embeddings` 中令牌数与嵌入数不匹配, 源于 `build_flat_image_bool_length` 低分辨率令牌布局硬编码。PR body 指出: `prompt image placeholder tokens: 550, multimodal image embeddings: 536`, 缺少 14 个低分辨率行列令牌。

## 实现拆解

步骤:

1. 修改 `build_flat_image_bool_length` 签名, 新增三个参数 (`image_use_col_tokens`, `use_single_crop_col_tokens`, `use_single_crop_start_token`), 默认值保持与之前行为一致 (`True`, `None`, `True`)。
2. 根据这些参数调整 `lengths` 计算 (加入 `low_res_extra` 和 `high_res_extra`)。
3. 重写低分辨率块生成逻辑, 使用循环构造行令牌并重复, 而非之前的平坦补丁填充。
4. 在高分辨率块中条件化插入 `image_col_id`。
5. 在调用点 `patched_call` 中从 HF 处理器传递这些标志。新增测试文件验证 MolmoWeb 风格配置和禁用列令牌场景。

关键文件:

- `vllm/model_executor/models/molmo2.py` (模块 模型层; 类别 `source`; 类型 `data-contract`; 符号 `build_flat_image_bool_length`): 核心修复文件, 修改了 `build_flat_image_bool_length` 函数签名和实现, 新增条件分支以尊重新处理器标志; 同时也更新了调用点。
- `tests/models/multimodal/processing/test_molmo2.py` (模块 测试; 类别 `test`; 类型 `test-coverage`; 符号 `test_build_flat_image_bool_length_matches_molmoweb_processor`)

\_tokens, test\_build\_flat\_image\_bool\_length\_respects\_disabled\_col\_tokens) : 新增测试文件, 验证修复后的函数行为: 测试 MolmoWeb 风格配置的令牌计数 (550 个令牌, 28 个列令牌) 和禁用列令牌场景 (45 个令牌, 0 个列令牌), 确保回归覆盖。

关键符号: build\_flat\_image\_bool\_length, test\_build\_flat\_image\_bool\_length\_matches\_molmo\_web\_processor\_tokens, test\_build\_flat\_image\_bool\_length\_respects\_disabled\_col\_tokens

## 关键源码片段

### tests/models/multimodal/processing/test\_molmo2.py

新增测试文件, 验证修复后的函数行为: 测试 MolmoWeb 风格配置的令牌计数 (550 个令牌, 28 个列令牌) 和禁用列令牌场景 (45 个令牌, 0 个列令牌), 确保回归覆盖。

```
from types import SimpleNamespace
import torch
from vllm.model_executor.models.molmo2 import build_flat_image_bool_length

def test_build_flat_image_bool_length_matches_molmo_web_processor_tokens():
    # 模拟 MolmoWeb 的处理器配置: use_single_crop_start_token=False
    hf_config = SimpleNamespace(
        image_patch_id=151938,
        low_res_image_start_token_id=151940,
        image_start_token_id=151936,
        image_col_id=151939,
        image_end_token_id=151937,
    )
    image_grids = torch.tensor([[14, 14, 14, 23]], dtype=torch.long)

    image_tokens, num_image_tokens = build_flat_image_bool_length(
        image_grids, hf_config,
        image_use_col_tokens=True,
        use_single_crop_col_tokens=None,
        use_single_crop_start_token=False,
    )

    # 预期: 总共 550 个令牌, 其中 28 个列令牌
    assert num_image_tokens.tolist() == [550]
    assert len(image_tokens) == 550
    assert image_tokens[0].item() == hf_config.image_start_token_id
    assert (image_tokens == hf_config.image_col_id).sum().item() == 28

def test_build_flat_image_bool_length_respects_disabled_col_tokens():
    # 测试禁用列令牌且启用低分辨率起始令牌的场景
    hf_config = SimpleNamespace(
        image_patch_id=151938,
        low_res_image_start_token_id=151940,
        image_start_token_id=151936,
        image_col_id=151939,
```

```
        image_end_token_id=151937,
    )
    image_grids = torch.tensor([[2, 3, 5, 7]], dtype=torch.long)

    image_tokens, num_image_tokens = build_flat_image_bool_length(
        image_grids, hf_config,
        image_use_col_tokens=False,
        use_single_crop_col_tokens=False,
        use_single_crop_start_token=True,
    )

    # 预期: 总共 45 个令牌, 0 个列令牌, 起始令牌为 low_res 类型
    assert num_image_tokens.tolist() == [45]
    assert len(image_tokens) == 45
    assert image_tokens[0].item() == hf_config.low_res_image_start_token_id
    assert (image_tokens == hf_config.image_col_id).sum().item() == 0
```

## 评论区精华

主要讨论集中于测试验证。DarkLight1337 询问新测试是否能在修复前失败，作者确认采用 TDD 流程：先添加测试观察到失败，再应用修复。作者还提供了 Im-eval 结果证明对标准 Molmo2-8B 无回归，以及 WebVoyager 性能对比（45.55% vs 46.05%）表明修复接近官方路径。

- 测试验证是否能在修复前失败 (testing): 确认测试能有效捕获回归。
- 对标准 Molmo2-8B 的回归检查 (testing): 无退化，对标准 Molmo2-8B 行为一致。

## 风险与影响

- 风险：主要风险是向后兼容性：新增参数有默认值（True/None/True），未传递时行为与之前保持一致，但需要确认所有调用点已更新（当前仅 patched\_call 一处）。另一个风险是数值变动可能影响其他 Molmo2 变体（如视频处理路径 build\_flat\_video\_bool\_length 未修改，视频未受影响）。测试覆盖了常见和禁用列令牌场景，但缺少端到端模型集成测试。
- 影响：直接影响 Molmo2 系列模型用户，特别是 MolmoWeb 变体（如 WebVoyager 任务）。修复使 vLLM 与 Hugging Face 处理器语义对齐，解锁之前因 CUDA 错误无法使用 MolmoWeb 模型的场景。对标准 Molmo2-8B 用户无影响（行为一致）。范围限于 Molmo2 多模态处理，不涉及其他模块。
- 风险标记：核心路径变更，数据契约调整，缺少端到端集成测试

## 关联脉络

- PR #38660 [Bug]: CUDA assert in triton attention for MolmoWeb models (Molmo2 architecture with different max\_position\_embeddings): 本 PR 修复的原始 Issue，描述了 MolmoWeb 模型在注意力内核中的 CUDA 断言错误，根因是图像令牌元数据不匹配。
- PR #42163 [WIP] MolmoWeb wrapper using Molmo2 architecture: 作者提及的关联 PR，在 Molmo2 基础上添加 MolmoWeb 包装器，依赖本修复。