

PR #32325 完整报告

vllm-project/vllm

[Model] Add Moondream3 model support(only query and caption skills)

合并时间: 2026-05-01 10:06

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/32325>

执行摘要

- 一句话: 新增 Moondream3 模型, 支持 Query 和 Caption
- 推荐动作: 建议阅读, 尤其是 `reconstruct_from_crops` 函数和 `Moondream3Processor` 的设计, 展示了如何将视觉预处理封装在 `processor` 中, 保持模型核心简洁。此外, Moondream3 的 prefix-LM 实现和 MoE 配置为其他类似模型提供参考。

功能与动机

社区请求支持 Moondream3 模型 (#25215), 该模型具有高吞吐量和优异的多模态性能。PR 实现了基本的 Query 和 Caption 技能, 满足大多数使用场景。

实现拆解

1. 定义模型配置 (`vllm/transformers_utils/configs/moondream3.py`): 创建组合配置类, 包含 Vision 子配置 (`crop_size`、`max_crops` 等) 和 Text 子配置 (MoE、`prefix_attn`、特殊 token ID)。桥接原生命名与 HF 标准属性。
2. 实现核心模型架构 (`vllm/model_executor/models/moondream3.py`):
 - Vision Encoder: 基于重叠 tiling 的 Vision Transformer, 支持多 crop 输入, 使用 `reconstruct_from_crops` 在 patch 级重建特征。内部使用 `MMEncoderAttention` (双向注意力)。
 - Text Decoder: 兼容 vLLM 的 decoder 注意力, 含 MoE、RoPE、tau 缩放。通过 `prefix_attn` 控制前 N 个位置使用双向注意力 (prefix-LM), 其余因果。
 - 多模态嵌入: 将视觉特征与 token 嵌入相加, 填充到 prefix 位置。
 - 抑制答案 token (默认 ID 3), 避免模型输出分隔符。
3. 开发自定义处理器 (`vllm/transformers_utils/processors/moondream3.py`):
 - `Moondream3Processor` 使用独立 tokenizer 仓库 `moondream/starmie-v1`, 包含预处理流水线: `select_tiling` 计算最优 tile 数, 归一化和 BF16 转换。
 - Chat template 根据文本前缀路由生成 Moondream3 特定 prompt (含 `<lendoftextl>`、`<image>`、`<lmd_reserved_0l>` 等特殊 token)。
4. 注册模型到 vLLM 框架: 在 `registry.py` 添加 `Moondream3ForCausalLM` 和 `HfMoondream` 两个架构别名; 在 `configs/__init__.py` 导入配置; 在 `model.py` 的 `is_mm_prefix_lm` 属性中添加 `moondream3`。

5. 编写测试与文档：添加处理器单元测试 (`test_moondream3_processing.py`) 和生成测试 (`test_moondream3_generation.py`, 含 TP 测试)；在 `confstest.py` 和 `model_utils.py` 中添加辅助设施；更新 `supported_models.md` 和 `multimodal_inputs.md`。

关键文件：

- `vllm/model_executor/models/moondream3.py` (模块 模型层；类别 `source`；类型 `data-contract`；符号 `reconstruct_from_crops`, `Moondream3VisionMLP`, `init`, `forward`)：核心模型实现，包含 Vision Encoder、Text Decoder、多模态嵌入和 `forward` 逻辑。
- `vllm/transformers_utils/processors/moondream3.py` (模块 处理器；类别 `source`；类型 `dependency-wiring`；符号 `Moondream3ProcessorKwargs`, `select_tiling`, `Moondream3Processor`, `init`)：自定义处理器，处理图像 `tiling`、归一化、`tokenization` 和 `prompt` 构建。
- `vllm/transformers_utils/configs/moondream3.py` (模块 配置；类别 `source`；类型 `core-logic`；符号 `Moondream3VisionConfig`, `init`, `Moondream3TextConfig`, `Moondream3Config`)：定义模型配置，包括 `vision` 和 `text` 子配置，MoE 参数映射。
- `tests/models/multimodal/processing/test_moondream3.py` (模块 测试；类别 `test`；类型 `test-coverage`；符号 `test_processor_creation`, `test_processor_apply`, `test_processor_pixel_values`, `test_processor_image_token_expansion`)：处理器功能的单元测试，验证 `tiling`、`placeholder` 扩展、像素值等。
- `tests/models/multimodal/generation/test_moondream3.py` (模块 测试；类别 `test`；类型 `test-coverage`；符号 `make_query_prompt`, `make_caption_prompt`, `test_tensor_parallel`, `llm`)：端到端生成测试，验证 `query` 和 `caption` 技能的正确性。
- `vllm/model_executor/models/registry.py` (模块 注册层；类别 `source`；类型 `data-contract`)：注册 `Moondream3ForCausalLM` 和 `HfMoondream` 架构，使 vLLM 能识别模型。
- `tests/confstest.py` (模块 测试配置；类别 `test`；类型 `test-coverage`)：添加 `skip_processor_init` 和 `tokenizer_name` 参数，支持 `moondream3` 测试。
- `tests/models/multimodal/generation/vlm_utils/model_utils.py` (模块 测试工具；类别 `test`；类型 `test-coverage`；符号 `moondream3_processor`, `moondream3_patch_hf_runner`, `processor`, `_normalize_tiling`)：添加 `moondream3_patch_hf_runner`，修补 HF runner 以对齐 vLLM 行为。

关键符号：`reconstruct_from_crops`, `select_tiling`, `Moondream3Processor.call`, `Moondream3Model.forward`, `make_query_prompt`, `make_caption_prompt`, `_encode_vision`, `_normalize_tiling`

关键源码片段

`vllm/model_executor/models/moondream3.py`

核心模型实现，包含 Vision Encoder、Text Decoder、多模态嵌入和 `forward` 逻辑。

```
# 从重叠 crops 中重建特征图
def reconstruct_from_crops(
    crops: torch.Tensor,
```

```
tiling: tuple[int, int],
overlap_margin: int,
patch_size: int = 14,
) -> torch.Tensor:
    """Reconstruct features from overlapping crops.
```

Args:

```
    crops: (N, H, W, D) 的 crop 特征张量。
    tiling: (tiling_h, tiling_w) 的瓦片数。
    overlap_margin: 重叠 margin 的 patch 数。
    patch_size: 每 patch 像素数（默认 14，适用于 SigLIP 风格的 ViT）。
```

Returns:

```
    reconstructed: (output_h, output_w, D) 完整重建特征。
    """
    tiling_h, tiling_w = tiling
    crop_height, crop_width = crops[0].shape[:2]
    margin_pixels = overlap_margin * patch_size

    output_h = (crop_height - 2 * margin_pixels) * tiling_h + 2 * margin_pixels
    output_w = (crop_width - 2 * margin_pixels) * tiling_w + 2 * margin_pixels

    reconstructed = torch.zeros(
        (output_h, output_w, crops[0].shape[2]),
        device=crops[0].device,
        dtype=crops[0].dtype,
    )

    for i, crop in enumerate(crops):
        tile_y = i // tiling_w
        tile_x = i % tiling_w

        # 边界 tile 保留边缘, 非边界裁剪重叠区域
        x_start = 0 if tile_x == 0 else margin_pixels
        x_end = crop_width if tile_x == tiling_w - 1 else crop_width - margin_pixels
        y_start = 0 if tile_y == 0 else margin_pixels
        y_end = crop_height if tile_y == tiling_h - 1 else crop_height - margin_pixels

        out_x = tile_x * (crop_width - 2 * margin_pixels)
        out_y = tile_y * (crop_height - 2 * margin_pixels)

        # 将有效区域放置到重建图对应位置
        reconstructed[
            out_y + y_start : out_y + y_end,
            out_x + x_start : out_x + x_end,
        ] = crop[y_start:y_end, x_start:x_end]

    return reconstructed
```

vllm/transformers_utils/processors/moondream3.py

自定义处理器，处理图像 tiling、归一化、tokenization 和 prompt 构建。

```
import math

def select_tiling(
    height: int, width: int, crop_size: int, max_crops: int
) -> tuple[int, int]:
    """Determine the optimal number of tiles to cover an image.
```

根据图像尺寸和 crop 大小，计算最优的 tiling 网格数 (h, w)。
如果图像小于 crop，则返回 (1, 1)。
否则最小化 tile 数，但不超过 max_crops，并尽量保持正方形。

Args:

height: 原始图像高度。
width: 原始图像宽度。
crop_size: 每个 crop 的尺寸（默认 378）。
max_crops: 最大允许的 crop 数（默认 12）。

Returns:

(tiling_h, tiling_w) 瓦片网格大小。

"""

```
if height <= crop_size or width <= crop_size:
    return (1, 1)
```

```
# 最小需要的 tile 数
```

```
min_h = math.ceil(height / crop_size)
min_w = math.ceil(width / crop_size)
```

```
# 如果最小数已超过限制，按比例缩小
```

```
if min_h * min_w > max_crops:
    ratio = math.sqrt(max_crops / (min_h * min_w))
    return (max(1, math.floor(min_h * ratio)),
            max(1, math.floor(min_w * ratio)))
```

```
# 否则在正方形约束下尽量填满 max_crops
```

```
h_tiles = math.floor(math.sqrt(max_crops * height / width))
w_tiles = math.floor(math.sqrt(max_crops * width / height))
```

```
h_tiles = max(h_tiles, min_h)
```

```
w_tiles = max(w_tiles, min_w)
```

```
# 如果超出，减少一个维度
```

```
if h_tiles * w_tiles > max_crops:
    if w_tiles > h_tiles:
        w_tiles = math.floor(max_crops / h_tiles)
    else:
        h_tiles = math.floor(max_crops / w_tiles)
```

```
return (max(1, h_tiles), max(1, w_tiles))
```

评论区精华

核心 engine 改动取舍: DarkLight1337 反对引入模型钩子和 `model_extra_output`, 认为会增加代码债务。作者 sniper35 最终同意移除 Point/Detect 技能相关核心改动, 只保留 Query/Caption, 避免对 Engine Core 的侵入。

grid_size 计算与权重映射错误: gemini-code-assist 指出 `grid_size` 错误使用了 `enc_n_layers` (应为 `crop_size/patch_size`), 以及 attention 权重 `qkv` 映射错误 (`qkv` 应为 `qkv_proj`)。已修正。

processor 输入标准化: DarkLight1337 要求模型层不应处理多种输入格式, 应由 processor 统一标准化。模型层的输入处理被简化, 仅接受 5d tensor 或 list of 4d tensors。

- 核心 engine 改动 (模型钩子) 是否必要 (design): 移除对 Engine Core 的侵入性修改, 仅保留非侵入模型注册。
- `grid_size` 计算依赖错误 (correctness): 已修正为动态计算。
- 权重名称映射错误 (correctness): 已修正映射规则。
- processor 输入标准化 (design): 简化模型层输入处理, 仅接受 5d tensor 或 list of 4d tensors, 其余由 processor 处理。
- IO Processor plugin 的引入 (design): 最终未引入 IO Processor plugin 改动, 保持原有接口。

风险与影响

- 风险: 新模型稳定性: 首次集成, 虽然通过 HF 对齐测试, 但实际部署可能出现未预期行为, 尤其是 tiling 路径和 prompt 格式化边缘情况。外部 tokenizer 依赖: 使用单独仓库 `moondream/starmie-v1` 作为 tokenizer, 若仓库不可用或变更则模型无法加载。性能开销: 图像预处理涉及重叠 tiling 和特征重建, 高分辨率图像可能增加延迟; MoE 解码器在批处理时有额外计算。Prefix-LM 实现风险: vLLM 的 prefix-LM 机制主要用于少量模型, 固定 prefix 长度 (730) 若配置变更可能导致错位。测试覆盖: 未覆盖所有 tiling 组合和极端情况。
- 影响: 用户: 可直接使用 `vllm.LLM` 加载 Moondream3 模型, 通过 Query 和 Caption prompt 格式进行多模态交互。系统: 新增约 5400 行代码, 包括完整模型、处理器、配置和测试。模型使用独立 tokenizer, 需额外网络获取。能利用 vLLM 的张量并行、pipeline 并行、prefix caching 等特性。团队: 需维护模型对 HF config 的兼容性, 跟踪上游更新。processor 中 `select_tiling` 算法若变化需同步。
- 风险标记: 依赖外部 tokenizer, 新模型稳定性, tiling 计算开销

关联脉络

- PR #32327 Update imports for Moondream3: PR body 中提及可能需要跟随 #32327 更新导入, 实际合并时已适配。