

PR #24855 完整报告

sgl-project/sglang

[Model] Add MiniCPM-V 4.6 support

合并时间: 2026-05-11 00:24

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/24855>

执行摘要

- 一句话: 添加 MiniCPM-V 4.6 多模态模型支持
- 推荐动作: 该 PR 结构清晰, 目标准确, 适合作为多模态模型集成的参考范例。建议仔细审阅视觉 Transformer 实现中的 MiniCPMV_ViTWindowAttentionMerger 设计 (包含 CPU 端索引构建以避免 PyTorch 2.10+ 设备检查), 以及配置类的 `_build_text_config` 回退逻辑。由于缺少自动化测试, 合并后应尽快补充单元测试和集成测试。

功能与动机

根据 PR 描述, MiniCPM-V 4.6 是 MiniCPM-V 系列的下一代版本, 相比于 4.5, 其架构在视觉编码器中添加了 mid-ViT 压缩 (2x2 窗口注意力 + 2x2 fold)、使用 MLP Merger 替代 Perceiver 重采样器、并将 LLM 骨干更新为 Qwen3.5 混合模型。SGLang 需要新增模型支持以维持功能对齐, 并提供性能与精度基准。

实现拆解

1. 新增视觉 Transformer 模块: 创建 `python/sglang/srt/models/minicpmv_vit.py`, 实现 MiniCPMV_ViTWindowAttentionMerger (mid-ViT 2x2 窗口注意力 + fold)、MiniCPMV_Merger (后编码器 MLP 连接器)、MiniCPMV_VisionEncoder (整体视觉编码器), 直接映射 HuggingFace 参考实现以简化权重加载。
2. 添加配置类: 创建 `python/sglang/srt/configs/minicpmv4_6.py`, 提供 MiniCPMV4_6VisionConfig 和 MiniCPMV4_6Config, 参数化 `insert_layer_id`、`downsample_mode` 等新架构特征。在 `configs/__init__.py` 和 `hf_transformers/common.py` 中注册, 确保 AutoConfig 能正确识别。
3. 集成模型主类: 在 `python/sglang/srt/models/minicpmv.py` 中添加 MiniCPMV4_6 类, 继承 MiniCPMBaseModel。覆盖 `init_vision_module`、`init_resampler` (此处实际省略重采样器)、`get_image_feature`、`forward` 等方法。LLM 部分使用 Qwen3_5ForCausalLM 并处理 `tie_word_embeddings`。
4. 实现图像预处理处理器: 创建 `python/sglang/srt/multimodal/processors/minicpmv4_6.py`, 实现 MiniCPMV4_6ImageProcessor。由于 HuggingFace 尚未提供原生处理器 (预计 5.7+), SGLang 端自行实现图像缩放、分块、NaViT 重排等逻辑, 并通过 `chat-template` 展开完成多模态 token 拼接。

5. 调整服务器参数：修改 `python/sglang/srt/server_args.py`，为 `MiniCPMV4_6ForConditionalGeneration` 架构调用 Mamba 基数缓存处理（`_handle_mamba_radix_cache`），因为 Qwen3.5 环境需要与 Qwen3.5 相同的配置。
6. 数据准备修复：微调 `python/sglang/benchmark/datasets/image.py` 以兼容 4.6 的 image token 生成（小范围控制流调整）。
7. 验证与基准：提供 MMMU val 精度（~37-39%）和 `sglang.bench_serving` 吞吐数据，确认功能正确且性能可接受。

关键文件：

- `python/sglang/srt/models/minicpmv_vit.py`（模块 视觉编码器；类别 source；类型 data-contract；符号 `MiniCPMV_ViTWindowAttentionMerger`, `init`, `get_window_index`, `forward`）：核心视觉 Transformer 实现，包含 mid-ViT 窗口注意力合并器、MLP Merger 和完整 VisionEncoder，是 4.6 架构差异的关键。
- `python/sglang/srt/multimodal/processors/minicpmv4_6.py`（模块 多模态处理器；类别 source；类型 dependency-wiring；符号 `_ensure_divide`, `_to_chw_tensor`, `_resize`, `_divide_to_patches`）：SGLang 自实现的图像预处理处理器，在 HF 原生处理器可用前承担图像缩放、分块、NaViT packing 等任务。
- `python/sglang/srt/models/minicpmv.py`（模块 模型主代码；类别 source；类型 data-contract；符号 `MiniCPMV4_6`, `init`, `init_llm`, `forward`）：模型主文件，添加 `MiniCPMV4_6` 类，集成视觉模块、LLM 骨干、前向逻辑和权重映射。
- `python/sglang/srt/configs/minicpmv4_6.py`（模块 模型配置；类别 source；类型 dependency-wiring；符号 `MiniCPMV4_6VisionConfig`, `init`, `_resolve_text_config_class`, `_build_text_config`）：新增 `MiniCPMV4_6Config` 和 `MiniCPMV4_6VisionConfig` 配置类，注册到 `AutoConfig`，承载 `downsample_mode`、`insert_layer_id` 等关键参数。
- `python/sglang/srt/server_args.py`（模块 服务器参数；类别 source；类型 core-logic）：为 `MiniCPMV4_6ForConditionalGeneration` 架构添加 `mamba radix cache` 配置，保障运行。
- `python/sglang/srt/configs/__init__.py`（模块 配置注册；类别 source；类型 dependency-wiring）：导出 `MiniCPMV4_6Config` 和 `MiniCPMV4_6VisionConfig`，供外部导入。
- `python/sglang/srt/utils/hf_transformers/common.py`（模块 工具函数；类别 source；类型 core-logic）：注册 `minicpmv4_6` 配置到 HuggingFace `AutoConfig` 映射，使 `from_pretrained` 能调用自定义配置。
- `python/sglang/benchmark/datasets/image.py`（模块 基准数据集；类别 source；类型 core-logic）：微调图像数据集生成逻辑以兼容 4.6 的 image token 格式。

关键符号：`MiniCPMV_ViTWindowAttentionMerger.init`,
`MiniCPMV_ViTWindowAttentionMerger.get_window_index`,
`MiniCPMV_ViTWindowAttentionMerger.forward`, `MiniCPMV_Merger.forward`,
`MiniCPMV_VisionEncoder.forward`, `MiniCPMV4_6ImageProcessor.process`,
`MiniCPMV4_6.init_vision_module`, `MiniCPMV4_6.get_image_feature`,
`MiniCPMV4_6.forward`, `MiniCPMV4_6Config.init`, `_resolve_text_config_class`

关键源码片段

[python/sglang/srt/models/minicpmv_vit.py](#)

核心视觉 Transformer 实现，包含 mid-ViT 窗口注意力合并器、MLP Merger 和完整 VisionEncoder，是 4.6 架构差异的关键。

```
class MiniCPMV_ViTWindowAttentionMerger(nn.Module):
    """Mid-ViT 2x2 window attention + 2x2 fold.

    Stage 1: reorder tokens so each 2x2 spatial window becomes 4 contiguous
    tokens; run packed self-attention with one window per cu_seqlens segment;
    un-reorder; add residual. (No length reduction yet.)

    Stage 2: fold each 2x2 window into a single token by concatenating the
    four hidden vectors along channel; pass through ``hidden*4 ->
    intermediate*4 -> hidden`` MLP; add the mean of the four window vectors
    as residual. ``target_sizes`` halves on each axis; ``cu_seqlens`` /
    ``max_seqlens`` are rebuilt for the compressed grid.
    """

    def __init__(
        self,
        config: PretrainedConfig,
        quant_config: Optional[QuantizationConfig] = None,
        prefix: str = "",
    ) -> None:
        super().__init__()
        self.window_kernel_size = (2, 2)
        self.embed_dim = config.hidden_size

        # ``flatten_batch=True``: input is one packed sequence
        # ``(1, sum_windows * window_area, D)`` with cu_seqlens demarcating
        # per-window segments.
        self.self_attn = VisionAttention(
            embed_dim=config.hidden_size,
            num_heads=config.num_attention_heads,
            projection_size=config.hidden_size,
            use_qkv_parallel=True,
            quant_config=quant_config,
            dropout=config.attention_dropout,
            softmax_in_single_precision=True,
            flatten_batch=True,
            prefix=add_prefix("self_attn", prefix),
        )
        self.layer_norm1 = nn.LayerNorm(self.embed_dim, eps=config.layer_norm_eps)

        window_area = self.window_kernel_size[0] * self.window_kernel_size[1]
        hidden_4x = self.embed_dim * window_area
        inter_4x = config.intermediate_size * window_area
```

```

self.pre_norm = nn.LayerNorm(hidden_4x, eps=config.layer_norm_eps)
self.linear_1 = ColumnParallelLinear(
    hidden_4x, inter_4x, bias=True,
    quant_config=quant_config, prefix=add_prefix("linear_1", prefix),
)
self.act = get_act_fn("gelu_pytorch_tanh")
self.linear_2 = RowParallelLinear(
    inter_4x, self.embed_dim, bias=True,
    quant_config=quant_config, prefix=add_prefix("linear_2", prefix),
)

def get_window_index(
    self, target_sizes: torch.Tensor
) -> Tuple[torch.Tensor, torch.Tensor, int]:
    """Return ``(permutation, per-window cu_seqlens, max_seqlens=4)``.

    Kept on CPU because mixing device-bound offsets with CPU arange trips
    strict dtype checks in PyTorch 2.10+.
    """
    window_h, window_w = self.window_kernel_size
    max_seqlens = window_h * window_w # 4

    window_index_list: List[torch.Tensor] = []
    cu_seqlens: List[int] = [0]
    token_offset = 0

    for height, width in target_sizes:
        height, width = int(height), int(width)
        # ... build permutation ...
        # (rest of implementation omitted for brevity; full source in PR)

```

评论区精华

该 PR 未引发实质性讨论。仅有维护者 JustinTong0323 的批准评论 'Great Job!!'，表明代码质量被认可，无争议。

- 模型架构对齐确认 (design): 架构变更被接受，无需修改。

风险与影响

- 风险:
 - 模型集成风险: 手动实现的视觉编码器与 MLP Merger 需与 HF 权重完全对齐，任何细微差异可能导致精度下降。作者提供了 MMMU 基准验证，但未包含单元测试覆盖。
 - 依赖上游风险: 当前图像处理器为 SGLang 自实现，未来 HF 发布原生处理器后需要切换，可能引入接口不兼容。
 - 运行时配置风险: 需要用户手动指定 `--dtype bfloat16`，否则可能因 `config.json` 缺省 `dtype` 导致 Triton kernel 加载失败（已确认 reproducible）。

- 缺少测试：未包含自动化测试文件，回归风险依赖 CI 手动触发。
- 影响：
 - 用户影响：用户现可使用 SGLang 运行 MiniCPM-V 4.6 模型，支持图像 / 视频 / 文本多模态输入，并可通过 `chat_template_kwargs` 控制推理模式。
 - 系统影响：新增约 1500 行代码，涉及视觉编码、配置、预处理等多个模块，对现有架构无侵入性（主要增量添加）。
 - 团队影响：增加长期维护负担，需跟踪 HF 上游更新并及时对齐。
 - 风险标记：缺少自动化测试，依赖上游 transformers 版本，需要显式 dtype 配置

关联脉络

- PR #24876 MiniCPM-V 4.6 cookbook: 该 PR 是直接配套的使用文档 PR，提供更详细的用法示例。