

PR #23568 完整报告

sgl-project/sglang

Parakeet nemotron encoder

合并时间: 2026-04-25 11:00

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/23568>

执行摘要

- 一句话: 为 Nemotron-Nano-VL 模型添加 Parakeet 音频编码器与动态分辨率
- 推荐动作: 值得精读: 动态分辨率预算算法和视频 tubelet 压缩的设计具有良好的通用性, 可复用于其他多模态模型。建议关注:
 1. `compute_dynamic_image_size` 中的 `budget` 调整逻辑 (`factor = sqrt(budget/native_patches)`) 确保了宽高比保持;
 2. `forward_video` 中的 `tubelet` 分组与 `linear embedder` 避免了额外的时序模型, 是一种轻量方案;
 3. `pad_input_ids` 对 `audio` 和 `visual` 的分离处理保证了多模态数据正确填充。变更整体质量较高, 但应尽快补充单元测试覆盖动态分支和音频路径。

功能与动机

Parakeet 是适合的音频编码方案, 官方 vLLM 已有适配参考。NVIDIA Nemotron-Nano-VL 模型需要音频理解能力, 因此集成 Parakeet 编码器。同时, 原有图像处理仅支持固定尺寸缩放 (如 384x384), 导致宽高比失真或冗余 token。动态分辨率通过 `budget` 约束自动选择最佳缩放尺寸, 可减少无效 token 数。视频时间压缩将相邻帧拼合到通道维度, 降低序列长度, 提升长视频处理能力。

实现拆解

1. 新增 Parakeet 音频编码器模块

- `python/sglang/srt/models/parakeet.py`: 定义 `ParakeetProjection` 投影层 (Norm+Linear+ReLU+Linear) 和 `ProjectedParakeet` (封装 HuggingFace 的 `ParakeetEncoder` + 投影层), 提供 `load_weights` 方法加载权重。
- `python/sglang/srt/configs/parakeet.py`: 定义 `ParakeetConfig` (继承 `HFPParakeetEncoderConfig`) 和 `immutable` 的 `ExtractorConfig`, 提供 `from_hf_config` 工厂方法配置转换。

2. 修改主模型 `NemotronH_Nano_VL_V2` (`python/sglang/srt/models/nano_nemotron_vl.py`)

- 在 `__init__` 中条件创建 `self.sound_encoder` (`ProjectedParakeet`), 仅当 `config.sound_config` 存在时。

- 重写 `pad_input_ids`: 分离 `visual` 和 `audio items`, 分别调用 `MultiModalityDataPaddingPatternTokenPairs` 填充, 解决多模态混合填充问题。
- 新增 `extract_feature_dynamic` 支持多尺寸图列表输入, `extract_video_feature_temporal` 支持视频时间压缩, `get_audio_feature` 调用 `sound_encoder` 提取音频特征。

3. 扩展 `MultimodalProcessor` (`python/sglang/srt/multimodal/processors/nano_nemotron_vl.py`)

- 初始化 `ParakeetExtractor` (基于 `ParakeetFeatureExtractor`) 用于音频预处理。
- 新增 `render_audio` 函数将原始音频波形转换为 `input_features`。
- 新增 `render_image_dynamic` 使用动态分辨率缩放, `render_tubelet` 用 `pixel_shuffle` 实现时空压缩。

4. 视频时间压缩 (`python/sglang/srt/models/radio.py`)

- `ViTPatchGenerator` 增加 `video_temporal_patch_size` 和 `separate_video_embedder` 参数。
- 新增 `forward_video` 方法: 按 `temporal_patch_size` 将帧分组为 `tubelet`, 拼合通道维度后通过 `video_embedder` 线性映射回嵌入空间, 实现时间降采样。

5. 动态分辨率工具函数 (`python/sglang/srt/multimodal/internvl_utils.py`)

- `compute_dynamic_image_size`: 根据原始宽高比和 `patch_budget` 计算目标尺寸 (snap 到 `patch_size*ds` 的整数倍), 返回宽高和 token 数。
- `dynamic_resize_image`: 调用 `compute_dynamic_image_size` 进行缩放后生成 `tensor`。
- `compute_budgeted_image_sizes`、`get_video_target_size_and_feature_size` 等配套函数。

6. 其他配套

- 新增 `python/sglang/srt/multimodal/audio_from_video.py` 提供 `extract_audio_from_video_bytes` 从视频中提取音频。
- 配置更新: `nano_nemotron_vl.py config` 增加 `sound_config` 字段; `radio config` 增加 `video_temporal_patch_size`。
- 测试: 未包含对应单元测试, 仅依赖 CI 集成测试。

关键文件:

- `python/sglang/srt/models/parakeet.py` (模块 音频编码器; 类别 `source`; 类型 `data-contract`; 符号 `ParakeetProjection`, `init`, `forward`, `ProjectedParakeet`): 新增 `Parakeet` 音频编码器核心模块, 包含投影层、权重加载和特征提取器, 是整个音频支持的基石
- `python/sglang/srt/models/nano_nemotron_vl.py` (模块 主模型; 类别 `source`; 类型 `data-contract`; 符号 `extract_feature_dynamic`, `extract_video_feature_temporal`, `get_audio_feature`, `is_sound_weights`): 主模型类大量修改: 条件初始化音频编码器、重写 `pad_input_ids` 支持多模态混合填充、新增多个特征提取方法
- `python/sglang/srt/multimodal/internvl_utils.py` (模块 图像处理工具; 类别 `source`; 类型 `core-logic`; 符号 `compute_dynamic_image_size`, `dynamic_resize_image`, `resize_image_to_pixels`, `compute_budgeted_image_sizes`): 新增动态分辨率计算工具函

数，为图像和视频处理提供预算感知缩放

- `python/sglang/srt/models/radio.py` (模块 视觉模型; 类别 `source`; 类型 `data-contract`; 符号 `forward_video`, `_forward_dynamic`, `_forward_video_temporal`) : ViTPatchGenerator 增加视频时间补丁嵌入, 新增 `forward_video` 实现 tubelet 压缩
- `python/sglang/srt/multimodal/processors/nano_nemotron_vl.py` (模块 多模态处理器; 类别 `source`; 类型 `core-logic`; 符号 `render_image_dynamic`, `render_tubelet`, `render_audio`) : 多模态处理器扩展音频支持、动态图像渲染和视频渲染, 是预处理链关键
- `python/sglang/srt/configs/parakeet.py` (模块 配置; 类别 `source`; 类型 `dependency-wiring`; 符号 `ParakeetConfig`, `init`, `from_hf_config`, `ExtractorConfig`) : 新增 `ParakeetConfig` 和 `ExtractorConfig` 配置类, 提供 `from_hf_config` 工厂方法
- `python/sglang/srt/multimodal/audio_from_video.py` (模块 音频提取; 类别 `source`; 类型 `dependency-wiring`; 符号 `extract_audio_from_video_bytes`) : 新增从视频字节流提取音频的工具函数

关键符号: `ProjectedParakeet.init`, `ProjectedParakeet.forward`, `ProjectedParakeet.load_weights`, `ParakeetProjection.init`, `ParakeetProjection.forward`, `ParakeetExtractor.init`, `NemotronH_Nano_VL_V2.init`, `NemotronH_Nano_VL_V2.pad_input_ids`, `NemotronH_Nano_VL_V2.extract_feature_dynamic`, `NemotronH_Nano_VL_V2.extract_video_feature_temporal`, `NemotronH_Nano_VL_V2.get_audio_feature`, `ViTPatchGenerator.init`, `ViTPatchGenerator.forward_video`, `compute_dynamic_image_size`, `dynamic_resize_image`, `resize_image_to_pixels`, `compute_budgeted_image_sizes`, `get_video_target_size_and_feature_size`, `video_to_pixel_values`, `render_image_dynamic`, `render_tubelet`, `render_audio`, `extract_audio_from_video_bytes`

关键源码片段

`python/sglang/srt/models/parakeet.py`

新增 Parakeet 音频编码器核心模块, 包含投影层、权重加载和特征提取器, 是整个音频支持的基石

```
class ProjectedParakeet(nn.Module):
    def __init__(self, config, *, dtype, llm_hidden_size, max_model_len):
        super().__init__()
        # 将 HuggingFace 配置转换为 ParakeetConfig, 注入 LLM 相关参数
        self.config = ParakeetConfig.from_hf_config(
            config, llm_hidden_size=llm_hidden_size, max_model_len=max_model_len
        )
        self.encoder = HFParakeetEncoder(self.config).to(dtype)
        self.projection = ParakeetProjection(self.config).to(dtype)

    def forward(self, input_features, attention_mask=None):
        # 音频特征经 encoder 编码后, 通过两层 MLP 投影到 LLM 隐藏维度
        outputs = self.encoder(
            input_features=input_features, attention_mask=attention_mask
        )
```

```
return self.projection(outputs.last_hidden_state)
```

```
def load_weights(self, weights):
    loaded_params = set()
    params_dict = dict(self.named_parameters())
    buffers_dict = dict(self.named_buffers())
    weights_list = list(weights.items() if isinstance(weights, dict) else weights)
    for name, weight in weights_list:
        # 跳过 feature_extractor (使用 HuggingFace 默认加载), 映射命名前缀
        if name.startswith("sound_encoder.encoder.feature_extractor."):
            continue
        if name.startswith("sound_encoder."):
            target_name = name[len("sound_encoder.")]
        elif name.startswith("sound_projection."):
            target_name = f"projection.{name[len('sound_projection.')]}"
        else:
            continue
        target = params_dict.get(target_name) or buffers_dict.get(target_name)
        if target is not None:
            weight_loader = getattr(target, "weight_loader", default_weight_loader)
            with torch.no_grad():
                weight_loader(target, weight)
            loaded_params.add(target_name)
    return loaded_params
```

python/sclang/srt/models/nano_nemotron_vl.py

主模型类大量修改: 条件初始化音频编码器、重写 pad_input_ids 支持多模态混合填充、新增多个特征提取方法

```
class NemotronH_Nano_VL_V2(EVS):
    def __init__(self, config, quant_config=None, prefix=""):
        super().__init__(config)
        self.downsample_ratio = config.downsample_ratio
        self.language_model = NemotronHForCausalLM(config=config.llm_config, ...)
        self.vision_model = RadioModel(config=config.create_radio_config()).to(...)
        # ...
        self.llm_hidden_size = config.llm_config.hidden_size
        self.model_dtype = self.language_model.config.torch_dtype

        # 条件创建音频编码器
        self.sound_encoder: ProjectedParakeet | None = None
        if getattr(config, "sound_config", None) is not None:
            self.sound_encoder = ProjectedParakeet(
                config.sound_config,
                dtype=self.model_dtype,
                llm_hidden_size=self.llm_hidden_size,
                max_model_len=getattr(config, "max_model_len", 8192),
            )
        self.config = config
```

```

def pad_input_ids(self, input_ids, mm_inputs):
    im_start_id = mm_inputs.im_start_id
    im_end_id = mm_inputs.im_end_id
    # 分离 visual 和 audio items
    visual_items = [item for item in mm_inputs.mm_items if not item.is_audio()]
    audio_items = [item for item in mm_inputs.mm_items if item.is_audio()]

    all_data_offsets = []
    if visual_items:
        mm_inputs.mm_items = visual_items
        helper = MultiModalityDataPaddingPatternTokenPairs([(im_start_id, im_end_id)])
        input_ids = helper.pad_input_tokens(input_ids, mm_inputs)
        all_data_offsets.extend(mm_inputs.data_offsets)
    # 对 audio 使用独立的特殊 token
    audio_start_id = getattr(mm_inputs, "audio_start_id", None)
    audio_end_id = getattr(mm_inputs, "audio_end_id", None)
    if audio_items and audio_start_id is not None and audio_end_id is not None:
        mm_inputs.mm_items = audio_items
        helper = MultiModalityDataPaddingPatternTokenPairs([(audio_start_id, audio_end_id)])
        input_ids = helper.pad_input_tokens(input_ids, mm_inputs)
        all_data_offsets.extend(mm_inputs.data_offsets)
    # 恢复完整列表并更新 offsets
    mm_inputs.mm_items = visual_items + audio_items
    mm_inputs.data_offsets = all_data_offsets
    # ...
    return input_ids

```

评论区精华

gemini-code-assist[bot]在 [nano_nemotron_vl.py:227](#) 指出 `extract_feature` 方法缺少 `mlp1` 投影层，导致输出维度为 `rmsnorm_hidden_size` 而非 `llm_hidden_size`，属于关键缺陷 (critical)。结论：作者在后续提交 [e571e7ff](#) 中修复了该问题，commit message 为 "fix: add missing mlp1 projection in extract_feature"。

gemini-code-assist[bot]在 [parakeet.py:135](#) 指出内联 `import` (`import math`) 违反 PEP8，应移到文件顶部。结论：虽然 commit 中未明确提及，但最终代码中已移除了内联 `import`（根据 head 版本文件可见 `import math` 已在顶部），问题已解决。

- `extract_feature` 缺少 `mlp1` 投影层导致维度错误 (correctness): 作者在后续 commit [e571e7ff](#) 中修复了该问题，添加了 `mlp1` 投影调用
- 内联 `import` 违反 PEP8 (style): 最终代码已修复，`import math` 被移到模块顶部

风险与影响

- 风险:

1. 回归风险: 核心模型 `nano_nemotron_vl.py` 的 `pad_input_ids` 重写、特征提取流程修改 (`extract_feature` 增加投影、新增 `extract_feature_dynamic` 等) 可能影响现有图像 / 视频推理路径，若动态分支未覆盖所有情况可能导致格式异常。

2. 性能风险：动态分辨率缩放计算引入额外 CPU 预处理开销；视频时间压缩可能改变序列长度，影响 Attention 掩码正确性。
3. 兼容性风险：新增 sound_config 配置字段，旧 checkpoint 不包含该字段时不会初始化 sound_encoder，但现有推理路径不受影响；ParakeetConfig 继承自 HF 配置，若 from_hf_config 中 scale_input=False、attention_bias=False 等硬编码可能与未来模型版本不兼容。
4. 测试缺口：本次变更无配套单元测试，仅依靠 CI 端到端测试，回归覆盖不足。 - 影响：
用户侧：用户可以上传音频文件进行推理（需音频编码器配置启用）；动态分辨率使输入的图像根据内容自动缩放，减少 token 量加速推理；视频处理时延降低得益于时间压缩。
系统侧：模型加载稍慢（Parakeet 编码器额外初始化），显存占用略有增加（音频编码器参数约数百 MB）。团队侧：需维护新增的音频编码器模块和动态分辨率逻辑；后续应补全单元测试并监控回归。

- 风险标记：核心模型变更，缺少测试覆盖，新增音频处理依赖，多模态数据流复杂

关联脉络

- 暂无明显关联 PR