

PR #27167 完整报告

sgl-project/sglang

[Model] Support encoder-free unified Text/Vision/Audio model

合并时间: 2026-06-03 23:58

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/27167>

执行摘要

- 一句话: 支持 Gemma4 Unified 编码器自由多模态模型
- 推荐动作: 值得精读, 特别是研究如何在不重构核心架构的前提下, 通过继承和轻量投影支持新多模态模型的设计模式。MTP 投机解码的集成方式也值得参考。建议在合并后关注 CI 结果, 并补充性能基准测试。

功能与动机

为 SGLang 添加对 gemma4_unified 模型族的支持 (如 google/gemma-4-12B-it), 这是一种编码器自由的统一文本 + 视觉 + 音频架构。视觉和音频使用轻量嵌入器直接投影原始像素块 / 波形帧到 token 流, 无需单独的视觉塔或音频 conformer。

实现拆解

实现包含以下主要步骤:

1. 核心模型定义: 新增 python/sglang/srt/models/gemma4_unified.py, 定义 Gemma4UnifiedForConditionalGeneration, 继承自 Gemma4ForConditionalGeneration, 重写模态特定特征提取和权重加载方式。
2. 视觉 / 音频嵌入器: 实现 Gemma4UnifiedVisionEmbedder (LN -> Dense -> LN + 因子化位置编码) 和 Gemma4UnifiedMultimodalEmbedder (RMSNorm -> Linear), 将原始像素和音频帧直接映射到 LM 空间。
3. 多模态处理器: 新增 python/sglang/srt/multimodal/processors/gemma4_unified.py, 继承 Gemma4SGLangProcessor, 调整音频填充倍数为 audio_samples_per_token (640)。
4. 配置注册: 在 config.py 中扩展 model_type 检查范围, 包含 gemma4_unified 和 gemma4_unified_assistant, 并添加 SWA 属性反转逻辑及 eoa_token_id 别名。同时在 model_config.py 和 base_processor.py 中注册模型类型和处理器映射。
5. 投机解码 (MTP) 集成: 在 gemma4_mtp.py 中新增 Gemma4UnifiedAssistantForCausalLM, 在 speculative_hook.py 中扩展算法识别逻辑, 使 NEXTN/EAGLE 自动提升为 FROZEN_KV_MTP。
6. 聊天入口点: 在 serving_chat.py 中扩展 is_gemma4 门控, 使推理 token 在流式输出时不被跳过。
7. 文档与部署选择器: 更新 cookbook 和部署选择器, 添加 12B 变体的硬件需求、启动命令和 MTP 配置说明。

关键文件：

- `python/sglang/srt/models/gemma4_unified.py` (模块 模型定义; 类别 `source`; 类型 `data-contract`; 符号 `Gemma4UnifiedVisionEmbedder`, `init`, `forward`, `Gemma4UnifiedMultimodalEmbedder`) : 核心模型实现, 包含 `Gemma4UnifiedForConditionalGeneration`、视觉和音频嵌入器, 是 PR 的主体变更。
- `python/sglang/srt/multimodal/processors/gemma4_unified.py` (模块 处理器; 类别 `source`; 类型 `dependency-wiring`; 符号 `Gemma4UnifiedSGLangProcessor`, `_get_audio_pad_multiple`) : 多模态处理器, 继承自 `Gemma4` 处理器并调整音频填充倍数, 是模型与数据预处理之间的桥梁。
- `python/sglang/srt/utils/hf_transformers/config.py` (模块 配置解析; 类别 `source`; 类型 `core-logic`) : 配置解析器扩展, 为 `gemma4_unified` 类型添加 `SWA` 属性反转和 `eoq_token_id` 别名, 影响模型加载阶段。
- `python/sglang/srt/models/gemma4_mtp.py` (模块 MTP 投机解码; 类别 `source`; 类型 `data-contract`; 符号 `Gemma4UnifiedAssistantForCausalLM`) : 添加 `Gemma4UnifiedAssistantForCausalLM` MTP 辅助类, 使统一模型支持投机解码。
- `python/sglang/srt/arg_groups/speculative_hook.py` (模块 投机钩子; 类别 `source`; 类型 `core-logic`) : 扩展投机算法别名解析, 识别 `Gemma4UnifiedAssistantForCausalLM` 并提升为 `FROZEN_KV_MTP`。
- `docs_new/cookbook/autoregressive/Google/Gemma4.mdx` (模块 部署文档; 类别 `docs`; 类型 `documentation`) : 更新 `cookbook` 文档, 添加 12B 变体的兼容性表格、硬件要求和启动命令。
- `python/sglang/srt/configs/model_config.py` (模块 模型配置; 类别 `source`; 类型 `data-contract`) : 注册 `gemma4_unified` 模型类型到配置映射表。
- `docs_new/src/snippets/autoregressive/gemma4-deployment.jsx` (模块 部署选择器; 类别 `docs`; 类型 `dependency-wiring`) : 部署选择器组件, 新增 12B 变体的选项和命令生成逻辑。
- `python/sglang/srt/entrypoints/openai/serving_chat.py` (模块 聊天入口; 类别 `source`; 类型 `core-logic`) : 扩展 `is_gemma4` 门控以包含 `gemma4_unified`, 确保推理 token 在流式输出中保留。
- `python/sglang/srt/multimodal/processors/base_processor.py` (模块 基础处理器; 类别 `source`; 类型 `dependency-wiring`) : 注册 `Gemma4UnifiedSGLangProcessor` 到处理器映射表。
- `python/sglang/srt/server_args.py` (模块 服务器参数; 类别 `source`; 类型 `core-logic`) : 注册 `gemma4_unified` 为有效模型类型, 通过 `server_args` 校验。

关键符号: `Gemma4UnifiedVisionEmbedder.forward`,
`Gemma4UnifiedMultimodalEmbedder.forward`, `Gemma4UnifiedForConditionalGeneration`,
`get_image_feature`, `Gemma4UnifiedForConditionalGeneration.get_audio_feature`,
`Gemma4UnifiedSGLangProcessor._get_audio_pad_multiple`,
`_resolve_speculative_algorithm_alias`

关键源码片段

python/sglang/srt/models/gemma4_unified.py

核心模型实现，包含 `Gemma4UnifiedForConditionalGeneration`、视觉和音频嵌入器，是 PR 的主体变更。

```
class Gemma4UnifiedVisionEmbedder(nn.Module):
    """Encoder-free vision embedder.

    Projects raw merged pixel patches ``(..., model_patch_size**2 * 3)`` into
    ``mm_embed_dim`` via ``LN1 -> Dense -> LN2``, adds factorized 2D positional
    embeddings, and applies a final ``LN``. Mirrors HF
    ``Gemma4UnifiedVisionEmbedder``; runs on the first PP rank only, so it uses
    plain (un-sharded) ``nn`` modules.
    """

    def __init__(self, config):
        super().__init__()
        patch_dim = config.model_patch_size**2 * 3 # 48*48*3 = 6912
        mm_embed_dim = config.mm_embed_dim

        self.patch_ln1 = nn.LayerNorm(patch_dim)
        self.patch_dense = nn.Linear(patch_dim, mm_embed_dim)
        self.patch_ln2 = nn.LayerNorm(mm_embed_dim)

        # Factorized 2D positional embedding table: (mm_posemb_size, 2, mm_embed_dim)
        self.pos_embedding = nn.Parameter(
            torch.zeros(config.mm_posemb_size, 2, mm_embed_dim)
        )
        self.pos_norm = nn.LayerNorm(mm_embed_dim)

    def forward(
        self, pixel_values: torch.Tensor, image_position_ids: torch.Tensor
    ) -> torch.Tensor:
        # pixel_values: (B, num_patches, patch_dim); image_position_ids: (B, num_patches, 2)
        hidden_states = self.patch_ln1(pixel_values.to(self.patch_dense.weight.dtype))
        hidden_states = self.patch_dense(hidden_states)
        hidden_states = self.patch_ln2(hidden_states)

        clamped = image_position_ids.clamp(min=0).long()
        valid = (image_position_ids != -1).to(self.pos_embedding.dtype).unsqueeze(-1)
        axes = torch.arange(2, device=image_position_ids.device)
        pos_embs = (self.pos_embedding[clamped, axes] * valid).sum(-2)
        hidden_states = hidden_states + pos_embs
        hidden_states = self.pos_norm(hidden_states)
        return hidden_states

class Gemma4UnifiedMultimodalEmbedder(nn.Module):
    """Shared vision/audio projection: ``RMSNorm(no scale) -> Linear`` to LM space.
```

Both the vision and audio configs expose ``output_proj_dims`` (the projection input dim) and ``rms_norm_eps``. ``embedding_pre_projection_norm`` has no learnable scale, so the only checkpoint tensor is ``embedding_projection.weight``.

```
def __init__(self, config, projection_config):
    super().__init__()
    # projection_config is either the vision or audio sub-config
    output_proj_dims = projection_config.output_proj_dims
    rms_norm_eps = getattr(projection_config, "rms_norm_eps", 1e-6)
    self.pre_proj_norm = Gemma4RMSNorm(output_proj_dims, eps=rms_norm_eps, no_scale=True)
    self.proj = nn.Linear(output_proj_dims, config.mm_embed_dim, bias=False)

def forward(self, hidden_states: torch.Tensor) -> torch.Tensor:
    hidden_states = self.pre_proj_norm(hidden_states)
    hidden_states = self.proj(hidden_states)
    return hidden_states
```

评论区精华

PR 未产生实质性 review 讨论；仅包含机器人自动回复（每日配额警告）和作者触发 CI 的指令，以及合并者提醒关注 CI 状态。

- 无实质讨论 (other): 无需处理。

风险与影响

- 风险：

1. 测试覆盖不足：未新增专门测试文件，仅依赖 CI 中的已有测试和作者声称的精度验证，存在回归风险。
2. 视觉注意力后端依赖：B200 (sm100) 默认使用 trtllm_mha 注意力后端，仅支持因果注意力；双向图像注意力需要显式指定 --attention-backend triton，否则图像质量会下降。
3. MTP 集成风险：修改了投机算法别名解析逻辑，可能影响现有 EAGLE/EAGLE3 模型的自动检测，需确保 Gemma4UnifiedAssistantForCausalLM 不干扰其他架构。
4. 配置兼容性：在 config.py 中为 gemma4_unified 添加了 eoa_token_id 别名逻辑，可能影响其他模型类型（如 gemma4 本身已有 eoa_token_index 但未处理），但实际影响局限于新模型。
5. 性能风险：视觉 / 音频投影路径未经性能基准测试，可能在高吞吐场景下成为瓶颈。 - 影响：对用户：可部署 google/gemma-4-12B-it 及未来类似编码器自由的多模态模型，支持文本、图像、视频和音频输入。对系统：新增约 500 行模型相关代码，增加模型注册和处理器映射，不影响现有模型。对团队：需维护该新模型路径及后续变体，但代码复用度高（继承自现有 gemma4 实现），维护成本较低。对部署：更新了文档和部署选择器，降低用户上手难度。 - 风险标记：缺少测试覆盖，视觉注意力后端依赖 triton，MTP 集成风险，配置兼容性

关联脉络

- PR #27171 [Docs] Update unified Text/Vision/Audio model cookbook: install + sgl-eval accuracy: 修改了同一文档文件 (Gemma4.mdx) , 为本 PR 的文档变更提供前驱更新。