

PR #22858 完整报告

sgl-project/sglang

[VLM] Enable per-image ViT cache and avoid TP CUDA context creation for Kimi-K2.5

合并时间: 2026-04-16 01:14

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/22858>

执行摘要

- 一句话: 修复 Kimi-K2.5 多模态模型在 TP 并行时每个 rank 在 device 0 上重复创建 CUDA 上下文的内存浪费问题。
- 推荐动作: 该 PR 值得精读, 重点关注其如何通过简单的数据移动 (CPU 卸载) 和键名标准化解决跨进程 CUDA 上下文重复初始化的深层问题。设计决策包括: 1) 优先内存优化而非微秒级数据传输开销; 2) 清理未使用代码以简化维护; 3) 遵循 SGL 标准键名以启用未来功能。建议结合多模态数据处理流程和 TP 通信机制理解变更。

功能与动机

根据 PR body 描述, 问题根因是 `KimiGPUProcessorWrapper._gpu_call` 返回的 `grid_thws` 张量构造在 `cuda:0` 上, 且不在 `FEATURE_NAMES` 列表中, 因此在跨进程传输 (`tokenizer_manager` → scheduler TP0 → 其他 TP rank) 时, Pickle 序列化会触发 PyTorch 的 CUDA IPC 机制, 强制每个接收进程在 device 0 上打开 CUDA IPC 句柄, 从而在每个 TP rank 初始化完整的 CUDA 上下文, 造成显著内存浪费 (H100/A100 上每个约 500 MiB)。

实现拆解

1. 数据键名标准化与 CPU 卸载: 在 `python/sglang/srt/multimodal/processors/kimi_k25.py` 的 `_gpu_call` 方法中, 将 GPU 预处理生成的 `grid_thws` 张量通过 `.cpu()` 移至 CPU, 并将返回字典的键从 `"grid_thws"` 改为 `"image_grid_thw"`, 以匹配 SGL 标准键名, 使下游 `get_new_expanded_mm_items()` 能按图像拆分以实现缓存粒度。
2. 模型端数据键名适配: 在 `python/sglang/srt/models/kimi_k25.py` 的 `KimiK25ForConditionalGeneration.get_image_feature` 方法中, 将访问 `items` 的字段从 `item.grid_thws` 改为 `item.image_grid_thw`, 与处理器输出保持一致。
3. 清理未使用代码: 删除同一文件中的 `vision_tower_forward_auto` 函数及其依赖的常量 `KIMIV_VT_INFER_MAX_PATCH_NUM`, 因为这些代码路径在 Kimi-K2.5 的当前实现中未被使用, 简化了代码库并消除了潜在的维护负担。
4. 测试与配置配套: 本次改动未包含直接的测试文件变更, 但涉及核心数据处理逻辑, 需确保现有多模态测试 (如涉及 Kimi-K2.5 的测试) 覆盖修改后的数据流。

关键文件:

- `python/sglang/srt/multimodal/processors/kimi_k25.py` (模块 多模态处理器; 类别 source; 类型 core-logic; 符号 `_gpu_call`): 这是数据预处理入口, 修改了 GPU 预处理路

径的输出键名并将 grid_thws 张量移至 CPU，直接解决了 CUDA 上下文重复创建的问题。

- python/sglang/srt/models/kimi_k25.py (模块 模型层; 类别 source; 类型 data-contract ; 符号 KimiK25ForConditionalGeneration.get_image_feature, vision_tower_forward_auto) : 模型端适配处理器的数据键名变更, 并删除了未使用的 vision_tower_forward_auto 函数, 简化了代码结构。

关键符号: _gpu_call, get_image_feature, vision_tower_forward_auto

关键源码片段

python/sglang/srt/multimodal/processors/kimi_k25.py

这是数据预处理入口, 修改了 GPU 预处理路径的输出键名并将 grid_thws 张量移至 CPU, 直接解决了 CUDA 上下文重复创建的问题。

```
def _gpu_call(self, text, images):
    # ... 前面的图像预处理和分词逻辑 ...

    # 4. GPU image preprocessing
    image_mean, image_std_inv = self._get_gpu_norm_tensors()
    pixel_values, grid_thws = _gpu_preprocess_images(
        images, resize_configs, image_mean, image_std_inv, self._patch_size
    )

    grid_thws = grid_thws.cpu() # 关键修复: 将 grid_thws 张量从 GPU 移至
    CPU, 避免后续跨进程传输时触发 CUDA IPC

    return {
        "input_ids": text_inputs["input_ids"],
        "pixel_values": pixel_values, # pixel_values 仍保留在 GPU 上, 用于后续视觉塔推理
        # 使用 SGL 标准键名, 以便 get_new_expanded_mm_items() 能按图像拆分, 实现每图像 ViT
        # 缓存粒度
        "image_grid_thw": grid_thws,
    }
```

python/sglang/srt/models/kimi_k25.py

模型端适配处理器的数据键名变更, 并删除了未使用的 vision_tower_forward_auto 函数, 简化了代码结构。

```
def get_image_feature(self, items: List[MultimodalDataItem]) -> torch.Tensor:
    device = self.vision_tower.device
    target_dtype = self.vision_tower.patch_embed.proj.weight.dtype
    pixel_values = torch.cat([item.feature for item in items], dim=0).to(
        device=device, dtype=target_dtype
    )
    # 适配处理器输出的键名变更: 从 item.grid_thws 改为 item.image_grid_thw
    grid_thws = torch.concat([item.image_grid_thw for item in items], dim=0).to(
        device # 注意: grid_thws 已从 CPU 移回 GPU, 但这是在单个进程内, 不会触发跨进程 CUDA
        IPC
    )
```

```
if self.use_data_parallel:
    image_embeds = run_dp_sharded_mrope_vision_model(
        self.vision_tower,
        pixel_values,
        grid_thws.tolist(),
        rope_type="rope_2d",
    )
    image_features = self.mm_projector(image_embeds)
    return image_features

# 非数据并行路径: 直接调用 vision_tower 和投影器
image_embeds = self.vision_tower(pixel_values, grid_thws)
proj_out = mm_projection_auto(self.mm_projector, image_embeds)
return torch.cat(proj_out, dim=0)
```

评论区精华

在 Issue 评论中, [wcsjtu](#) 提出了一个关键疑问: 为什么在 `get_image_feature` 方法中直接调用 `self.vision_tower` 而不是使用 `vision_tower_forward_auto`? 这可能在大批量图像 (`len(items)` 很大) 时导致 CUDA OOM。作者未在 PR 讨论中直接回应, 但通过删除 `vision_tower_forward_auto` 函数, 暗示该函数在当前上下文中已不再需要或存在设计权衡。Review 中 [mickqian](#) 批准了 PR, 未提出进一步讨论。

- 直接调用 `vision_tower` 与自动批处理的权衡 (design): PR 通过删除 `vision_tower_forward_auto` 函数隐含了该函数在当前上下文中不再需要, 但未直接回应 OOM 风险。

风险与影响

- 风险: 技术风险:
 1. 回归风险: 修改数据键名 (`grid_thws` → `image_grid_thw`) 可能破坏依赖旧键名的下游代码, 但 PR 同步更新了模型端的引用, 且注释说明是为了启用每图像 ViT 缓存, 风险可控。
 2. 性能影响: 将 `grid_thws` 移至 CPU 可能引入额外的 CPU-GPU 数据传输开销, 但相比避免每个 TP rank 重复创建 CUDA 上下文的内存节省 (每个约 500 MiB), 收益显著。
 3. 兼容性: 删除 `vision_tower_forward_auto` 函数可能影响其他模型或未来扩展, 但该函数未被当前 Kimi-K2.5 使用, 且 PR 未提及其他依赖, 风险较低。
 4. 测试覆盖: 缺少针对修改的直接单元测试, 需依赖现有集成测试验证功能正确性。
- 影响: 影响范围:
 1. 用户影响: 对终端用户透明, 但会减少多模态推理时的内存占用, 提升系统稳定性, 尤其在大规模 TP 并行场景。
 2. 系统影响: 优化了 Kimi-K2.5 模型在多 GPU 环境下的内存使用, 避免不必要的 CUDA 上下文重复初始化, 降低 OOM 风险。
 3. 团队影响: 统一了数据键名规范 (`image_grid_thw`), 促进了代码一致性, 并为每图像 ViT 缓存功能铺平道路。

- 风险标记: 核心路径变更, 缺少测试覆盖, 数据契约改动

关联脉络

- PR #22490 [EPD][VLM] Support Kimi VL EPD: 同属 Kimi 多模态模型功能扩展, 涉及相似的多模态处理器和模型文件 (如 kimi_k25.py), 可参考其数据处理模式。
- PR #22690 [diffusion] model: Properly validate device for Mistral 3 attention: 类似设备验证和跨平台 (AMD/NVIDIA) 支持问题, 涉及 CUDA 上下文管理。