

PR #23649 完整报告

sgl-project/sglang

[diffusion] support LoRA for LTX2.3

合并时间: 2026-04-25 01:52

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/23649>

执行摘要

- 一句话: 为 LTX2.3 添加 LoRA 支持与多条件图像
- 推荐动作: 值得精读, 特别关注 `linear.py` 中的 LoRA 权重管理重构和多条件图像的 SP 支持设计; 但需注意 `review` 指出的两个正确性风险, 若未修复应尽快跟进。文档片段 `ltx-deployment.jsx` 作为交互式配置示例, 可用于其他部署场景。

功能与动机

PR 描述了使用 LoRA 权重进行生成的示例命令, 但未显式说明动机。从实现分析, 主要动机是支持 LTX-2.3-Transition-LORA 这类微调权重, 使用户可以在不改变基模型的情况下通过 LoRA 实现视频风格迁移或过渡效果。

实现拆解

1. LoRA 权重管理重构: `linear.py` 引入 `LoRAWeightEntry` 元组, 每个条目包含 `rank`、`alpha` 等缩放参数; 添加 `_should_merge_in_fp32` 环境变量控制, 允许在合并 LoRA 时提升精度。
2. 多条件图像支持: `image_encoding.py` 和 `ltx_2_denoising.py` 将单图像输入泛化为列表, 支持首帧和尾帧条件; 添加 `_normalize_ltx2_image_paths`、`_normalize_ltx2_image_latents`、`_get_ltx2_condition_spans` 等函数处理条件跨度的计算与 SP 语义。
3. 低显存优化: `ltx_2_pipeline.py` 新增 `prepare_upsample_after_stage1` 和 `_release_stage1_for_low_vram`, 在 `snapshot` 低显存模式下提前释放 `stage-1` 以腾出空间用于 `stage-2` 和 `upsample`。
4. LoRA 应用日志改进: `lora_pipeline.py` 按适配器跟踪缺失层, 提供更清晰的覆盖报告, 帮助诊断 LoRA 权重未加载问题。
5. 文档与部署示例: 新增 `ltx-deployment.jsx` 交互式组件和 `LTX.mdx` 手册, 展示包含 LoRA 选型的部署命令生成器。
6. 测试配置清理: 删除旧的性能基线、精度配置和 GPU 案例文件, 简化测试套件。

关键文件:

- `docs_new/src/snippets/diffusion/ltx-deployment.jsx` (模块 部署交互; 类别 `source`; 类型 `core-logic`; 符号 `LTXDeployment`, `getInitialState`, `checkDarkMode`, `availableLoras`): 新增交互式部署命令生成器组件, 包含 LoRA 选型逻辑和模型配置, 是用户可见的核心变更

入口之一。

- `python/sglang/multimodal_gen/runtime/pipelines_core/stages/ltx_2_denoising.py` (模块 去噪阶段; 类别 `source`; 类型 `core-logic`; 符号 `_normalize_ltx2_condition_latents`, `_get_ltx2_condition_spans`, `_should_apply_ltx2_ti2v`) : 核心去噪阶段, 新增多条件图像跨度计算函数, 处理 SP 切分下的条件帧映射。
- `python/sglang/multimodal_gen/runtime/pipelines_core/stages/image_encoding.py` (模块 图像编码; 类别 `source`; 类型 `core-logic`; 符号 `_normalize_ltx2_image_paths`, `_normalize_ltx2_image_latents`) : 图像编码阶段, 将单图像处理扩展为多图像列表, 支持首帧和尾帧条件。
- `python/sglang/multimodal_gen/runtime/layers/lora/linear.py` (模块 LoRA 层; 类别 `source`; 类型 `core-logic`; 符号 `_should_merge_in_fp32`) : LoRA 层核心重构: 引入 `LoRAWeightEntry` 元组存储 `per-adapter` 缩放参数, 新增 FP32 合并控制路径。
- `python/sglang/multimodal_gen/runtime/pipelines/ltx_2_pipeline.py` (模块 管道管理; 类别 `source`; 类型 `core-logic`; 符号 `prepare_upsample_after_stage1`, `_release_stage1_for_low_vram`, `prepare_ltx2_upsample_after_stage1`) : 管道设备管理增强, 新增低显存释放方法并同步到外层接口, 提升 `upsample` 阶段 VRAM 效率。
- `python/sglang/multimodal_gen/runtime/pipelines_core/lora_pipeline.py` (模块 LoRA 流水线; 类别 `source`; 类型 `core-logic`) : LoRA 应用逻辑改进, 按适配器跟踪缺失层, 提供详细覆盖日志, 便于用户调试。
- `python/sglang/multimodal_gen/runtime/pipelines_core/stages/upsampling.py` (模块 上采样; 类别 `source`; 类型 `core-logic`) : 配合低显存优化, 调整 `upsample` 阶段的设备管理逻辑。

关键符号: `LTXDeployment`, `getInitialState`, `checkDarkMode`, `availableLoras`, `handleRadioChange`, `nextSupported`, `handleLoraToggle`, `getDeviceMode`, `_normalize_ltx2_condition_latents`, `_get_ltx2_condition_spans`, `_should_apply_ltx2_ti2v`, `_normalize_ltx2_image_paths`, `_normalize_ltx2_image_latents`, `prepare_upsample_after_stage1`, `_release_stage1_for_low_vram`, `prepare_ltx2_upsample_after_stage1`, `_should_merge_in_fp32`

关键源码片段

docs_new/src/snippets/diffusion/ltx-deployment.jsx

新增交互式部署命令生成器组件, 包含 LoRA 选型逻辑和模型配置, 是用户可见的核心变更入口之一。

```
export const LTXDeployment = () => {
  // 硬件、模型、管道选项配置
  const options = { /* ... */ };

  // 模型配置, 包括支持的 LoRA 列表
  const modelConfigs = {
    ltx2: {
      repoId: 'Lightricks/LTX-2',
      pipelines: {
```

```

    'one-stage': 'LTX2Pipeline',
    'two-stage': 'LTX2TwoStagePipeline',
  },
  supportedLoras: [], // 无 LoRA 支持
},
ltx23: {
  repoId: 'Lightricks/LTX-2.3',
  pipelines: {
    'one-stage': 'LTX2Pipeline',
    'two-stage': 'LTX2TwoStagePipeline',
    'two-stage-hq': 'LTX2TwoStageHQPipeline',
  },
  supportedLoras: [ // 定义 LoRA 元数据
    {
      id: 'transition',
      path: 'valiantcat/LTX-2.3-Transition-LORA',
      weightName: 'ltx2.3-transition.safetensors',
      validPipelines: ['two-stage', 'two-stage-hq'], // 仅双阶段管道可用
    },
  ],
},
};

```

// 根据当前选择的模型和管道过滤可用 LoRA

```

const availableLoras = (() => {
  const config = modelConfigs[values.model];
  return (config?.supportedLoras || []).filter((lora) =>
    lora.validPipelines.includes(values.pipeline)
  );
})();

```

// 单选变更时同步维护 LoRA 选择合法性

```

const handleRadioChange = (optionName, itemId) => {
  setValues((prev) => {
    const next = { ...prev, [optionName]: itemId };
    // 若管道与模型不匹配, 回退到 two-stage
    if (!options.pipeline.items.some(
      (item) => item.id === next.pipeline && item.validModels.includes(next.model)
    )) {
      next.pipeline = 'two-stage';
    }
    // 若当前 LoRA 在新配置下无效, 重置为 none
    const config = modelConfigs[next.model];
    const nextSupported = (config?.supportedLoras || [])
      .filter((lora) => lora.validPipelines.includes(next.pipeline));
    const isValid = nextSupported.some((lora) => lora.path === prev.selectedLoraPath);
    if (!isValid) {
      next.selectedLoraPath = 'none';
    }
  })
};

```

```

        return next;
    });
};
};

```

python/sclang/multimodal_gen/runtime/pipelines_core/stages/ltx_2_denoising.py

核心去噪阶段，新增多条件图像跨度计算函数，处理 SP 切分下的条件帧映射。

```

@classmethod
def _get_ltx2_condition_spans(
    cls,
    batch: Req,
    latents: torch.Tensor,
    image_latent: torch.Tensor | list[torch.Tensor] | None,
    num_img_tokens: int,
) -> list[tuple[int, torch.Tensor]]:
    if num_img_tokens <= 0:
        return []
    if not (isinstance(latents, torch.Tensor) and latents.ndim == 3):
        raise ValueError("LTX-2 TI2V expects packed token latents [B, S, D].")

    # 1. 规范化：将单个或列表统一为列表
    condition_latents = cls._normalize_ltx2_condition_latents(image_latent)
    if not condition_latents:
        return []
    if len(condition_latents) > 2:
        raise ValueError("LTX-2 TI2V currently supports at most two conditioning images.")

    # 2. 验证每个条件 latent 的维度
    for cond in condition_latents:
        if not (isinstance(cond, torch.Tensor) and cond.ndim == 3):
            raise ValueError("Expected LTX-2 conditioning latents to be packed tensors [B, S, D].")
        if int(cond.shape[1]) < int(num_img_tokens):
            raise ValueError("LTX-2 conditioning latent is shorter than one frame token span.")

    # 3. 根据是否经过 SP 切分计算 span 位置
    did_sp_shard = bool(getattr(batch, "did_sp_shard_latents", False))
    if not did_sp_shard:
        # 非 SP：条件帧在序列首（和尾）
        if len(condition_latents) == 1:
            return [(0, condition_latents[0])]
        return [
            (0, condition_latents[0]),
            (int(latents.shape[1]) - int(num_img_tokens), condition_latents[1]),
        ]

    # SP 场景：根据 local_start_frame 判断当前分片是否包含首帧 / 尾帧
    tokens_per_frame = int(getattr(batch, "sp_video_tokens_per_frame", 0))

```

```

if tokens_per_frame <= 0:
    raise ValueError("SP-sharded LTX-2 TI2V requires batch.sp_video_tokens_per_frame.")
if int(num_img_tokens) != int(tokens_per_frame):
    raise ValueError("LTX-2 conditioning token count must match one latent frame when using SP.")

raw_shape = getattr(batch, "raw_latent_shape", None)
if raw_shape is None:
    raise ValueError("SP-sharded LTX-2 TI2V requires batch.raw_latent_shape.")
global_seq_len = int(raw_shape[1])
if global_seq_len % tokens_per_frame != 0:
    raise ValueError("SP-sharded LTX-2 TI2V expected raw seq_len divisible by tokens_per_frame.")

global_num_frames = global_seq_len // tokens_per_frame
local_start_frame = int(getattr(batch, "sp_video_start_frame", 0))
local_num_frames = int(getattr(batch, "sp_video_latent_num_frames", 0))
local_end_frame = local_start_frame + local_num_frames

spans: list[tuple[int, torch.Tensor]] = []
if local_start_frame == 0:
    spans.append((0, condition_latents[0]))

if len(condition_latents) == 2:
    last_global_frame = global_num_frames - 1
    if local_start_frame <= last_global_frame < local_end_frame:
        local_last_frame = last_global_frame - local_start_frame
        spans.append(
            (local_last_frame * tokens_per_frame, condition_latents[1])
        )

return spans

```

[python/sglang/multimodal_gen/runtime/layers/lora/linear.py](#)

LoRA 层核心重构：引入 LoRAWeightEntry 元组存储 per-adapter 缩放参数，新增 FP32 合并控制路径。

```

# LoRA 权重条目：包含 A/B 参数、路径、强度、rank 和 alpha
LoRAWeightEntry = tuple[
    torch.nn.Parameter,
    torch.nn.Parameter,
    str | None,
    float,
    int | None,
    int | None,
]

def _should_merge_in_fp32(self, lora_list: list[LoRAWeightEntry]) -> bool:
    # 环境变量开关，默认关闭

```

```

if os.getenv("SGLANG_DIFFUSION_LORA_MERGE_FP32", "0") != "1":
    return False
# 跳过蒸馏 LoRA (已在高精度训练时融合)
for _, _, lora_path, _, _, _ in lora_list:
    if lora_path and "distilled-lora" in lora_path.lower():
        return False
return True

def _merge_lora_into_data(self, data: torch.Tensor, lora_list: list[LoRAWeightEntry]) -> None:
    # 遍历所有 LoRA 适配器, 按顺序合并
    for lora_A, lora_B, _, lora_strength, lora_rank, lora_alpha in lora_list:
        lora_A_sliced = self.slice_lora_a_weights(lora_A.to(data))
        lora_B_sliced = self.slice_lora_b_weights(lora_B.to(data))
        scale = lora_strength
        # 使用 per-adapter 的 rank/alpha 计算缩放比例
        if lora_alpha is not None and lora_rank is not None and lora_alpha != lora_rank:
            scale *= lora_alpha / lora_rank
        # ... 合并逻辑 (注意: 此处若 lora_B_sliced 是非 tensor 的 tuple 会出错)

```

评论区精华

Review 中 gemini-code-assist[bot] 指出两个 high 优先级问题:

1. Non-tensor LoRA 权重合并: 当 `slice_lora_b_weights` 返回 tuple (例如 `QKVParallelLinearWithLoRA`) 时, 代码直接对 tuple 执行矩阵乘法, 将引发 `TypeError`。
 2. DTensor 合并路径中的 shard 错误: 在 FSDP DTensor 路径中, `_merge_lora_into_data` 在全量 base weight (`full_tensor()`) 上合并, 但仍使用 `slice_lora_*` 对 LoRA 权重进行切分, 导致权重更新不完整或形状不匹配。两个问题尚未看到明确回复, 需确认是否在合并前已修复。
- Non-tensor LoRA 权重合并错误 (correctness): 未在评论中看到回复或修复, 需后续确认。
 - DTensor 合并路径中的 LoRA shard 错误 (correctness): 未在评论中看到回复或修复, 需后续确认。

风险与影响

- 风险:
 1. LoRA 合并正确性风险: review 指出的两个问题如果未修复, 可能导致 LoRA 权重合并后模型精度异常或运行时崩溃。文件: `linear.py` 第 206、302 行。
 2. 多条件图像兼容性风险: `batch.image_path` 此前为标量路径, 现支持列表; 若上游仍传递标量, `_normalize_ltx2_image_paths` 可兼容, 但其他消费方可能未适配列表格式。
 3. 测试覆盖缺失: 删除了性能基线等测试配置, 但没有补充对应的 LoRA 或多条件图像的测试用例, 返回回归风险升高。- 影响: 对用户: LTX2.3 用户现可通过 `--lora-path` 和 `--lora-weight-name` 参数加载自定义 LoRA 权重, 支持过渡效果等风格化视频生成; 多条件图像用户可指定首帧和尾帧以控制视频起始与结束内容。对系统: 新增的低显存释放策略可降低 VRAM 峰值占用, 但 `_should_merge_in_fp32` 可能增加计算开销。对团队: 代码结构上 LoRA 层已泛化, 为后续支持更多扩散模型的 LoRA 提供基础。- 风险标记

: LoRA 合并正确性风险，多图像 API 兼容性，测试覆盖缩减

关联脉络

- PR #23624 [diffusion] fix: unify LTX-2.3 HQ codepath gates for all LTX-2.3 variants: 同属 LTX-2.3 扩散模型改进，修改了同为 `denoising_av.py` 和 `ltx_2_pipeline.py` 等文件，补充了高质量管道路径。
- PR #23366 [diffusion] model: support LTX2.3 high quality pipeline: 为 LTX-2.3 添加了高质量两阶段流水线，是本 PR LoRA 支持的基础前提。
- PR #23151 [Diffusion] add per-step rollout options for SDE and trajectory capture: 涉及扩散模型调度器组件，本 PR 中 `ltx_2_denoising.py` 的修改与该文件的 SDE 步骤逻辑有一定关联。