

# PR #25891 完整报告

sgl-project/sglang

[diffusion] enable inference mode in pipeline executor

合并时间: 2026-05-21 13:20

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/25891>

## 执行摘要

- 一句话: diffusion executor 启用 inference\_mode 加速约 7%
- 推荐动作: 值得精读。设计上采用「全局 inference\_mode + stage 级回退」的 scoped 模式, 优雅处理了不同组件的 version counter 依赖, 为 future 优化提供了可扩展的框架。LoRA 的 `_as_mutable_tensor` 方法也是处理 inference tensor 不可变性的典型模式。

## 功能与动机

Diffusion pipeline 执行时, 大部分 stage (如 DiT、VAE) 不需要梯度且不依赖 tensor version counter, 但在 no\_grad 模式下 torch 仍会维护 version counter 开销。启用 inference\_mode 可以跳过 version counter 维护, 显著降低去噪延迟。PR body 报告 denoise 延迟下降 7.2% (从 152.15ms 到 141.13ms), 且峰值内存不变。

## 实现拆解

1. 导入 torch 和 current\_platform: 在 pipeline\_executor.py 顶部添加 import torch 和 from sglang.multimodal\_gen.runtime.platforms import current\_platform。
2. 包装执行入口: 在 execute\_with\_profiling 和 execute\_group\_with\_profiling 的 profile 上下文内部, 再套一层 with current\_platform.inference\_mode(), 确保整个执行链处于平台指定的 inference 模式。
3. 引入 stage 级别上下文: 新增 \_stage\_execution\_context 和 \_stage\_needs\_version\_counters 静态方法。对于需要 version counter 的阶段 (FSDP 启用、或 CPU-offload 涉及特定组件名如 transformer、text\_encoder、vae 等), 上下文临时回退到 torch.inference\_mode(False), torch.no\_grad(), 使得这些阶段的 hooks 能够正常更新 tensor version。
4. 改造并行执行器: 在 parallel\_executor.py 中将所有 stage 调用替换为 self.run\_stage\_with\_context(...), 该包装方法自动应用 stage 级别的上下文。同步修改 sync\_executor.py 和 layerwise\_offload.py 中需要修改 tensor 的地方, 确保其在 inference\_mode(False) 下执行。
5. 修复 LoRA 合并 / 反合并: 在 linear.py 中添加 \_as\_mutable\_tensor 静态方法, 当 tensor 处于 inference 模式时, 通过 detach().clone() 生成普通 tensor。在 merge\_lora\_weights 和 unmerge\_lora\_weights 的每个修改 point 调用此方法, 避免直接修改 inference tensor。

6. 增加测试覆盖：新建 `test_pipeline_executor.py`（覆盖 `execute/execute_group` 的 `inference_mode` 正确性、`stage` 上下文回退条件、`NoGradPlatform` 兼容）、`test_lora_inference_mode.py`（覆盖 LoRA `merge/unmerge` 对 `inference base weight` 的处理）和扩展现有 `test_layerwise_offload.py`。
7. NPU 平台适配：在 `npu.py` 中添加 `inference_mode` 类方法实现（当前为 `torch.no_grad()`，未来可替换为真正的 NPU `inference` 模式）。

关键文件：

- `python/sglang/multimodal_gen/runtime/pipelines_core/executors/pipeline_executor.py`（模块 执行器；类别 `source`；类型 `core-logic`；符号 `_stage_execution_context`, `_stage_needs_version_counters`, `run_stage_with_context`）：核心变更文件：在 `execute_with_profiling` 和 `execute_group_with_profiling` 中引入 `inference_mode` 包装，新增 `stage` 级上下文管理和版本计数器需求判断，是性能提升和兼容性保障的基石。
- `python/sglang/multimodal_gen/test/unit/test_pipeline_executor.py`（模块 测试；类别 `test`；类型 `test-coverage`；符号 `_RecordingExecutor`, `init`, `execute`, `execute_group`）：新增测试文件，全面覆盖 `execute_with_profiling` 和 `execute_group_with_profiling` 在两种平台（`InferenceTensorPlatform` 和 `NoGradPlatform`）下的 `inference_mode` 行为，以及 `stage` 上下文在多种 `offload/FSDP` 配置下的版本计数器可用性。
- `python/sglang/multimodal_gen/runtime/layers/lora/linear.py`（模块 LoRA；类别 `source`；类型 `core-logic`；符号 `_as_mutable_tensor`）：新增 `_as_mutable_tensor` 工具方法并在 `merge/unmerge` 路径中调用，确保在 `inference_mode` 下也能安全修改权重。这是修复 LoRA 兼容性的关键。
- `python/sglang/multimodal_gen/runtime/pipelines_core/executors/parallel_executor.py`（模块 执行器；类别 `source`；类型 `core-logic`）：将 `stage` 调用改为通过 `run_stage_with_context` 包装，确保每个 `stage` 在正确的上下文（`inference` 或降级）中执行。
- `python/sglang/multimodal_gen/runtime/managers/memory_managers/layerwise_offload.py`（模块 卸载管理；类别 `source`；类型 `core-logic`）：在 `prefetch/release` 操作中显式包装 `torch.inference_mode(False)`, `no_grad()`，确保在 `inference` 模式下也能正常创建空 `tensor` 和修改 `target data`。
- `python/sglang/multimodal_gen/runtime/platforms/npu.py`（模块 平台适配；类别 `source`；类型 `core-logic`；符号 `inference_mode`）：为 NPU 平台添加 `inference_mode` 方法（当前为 `torch.no_grad()`），保持接口一致性，未来可替换为 NPU 原生 `inference` 模式。
- `python/sglang/multimodal_gen/test/unit/test_lora_inference_mode.py`（模块 测试；类别 `test`；类型 `test-coverage`；符号 `test_lora_merge_unmerge_handles_inference_base_weight`）：新增测试，验证 LoRA `merge/unmerge` 在 `base_weight` 为 `inference tensor` 时的正确行为。
- `python/sglang/multimodal_gen/test/unit/test_layerwise_offload.py`（模块 测试；类别 `test`；类型 `test-coverage`；符号 `test_layerwise_offload_uses_normal_tensors_under_inference_mode`）：扩展现有测试，覆盖 `inference_mode` 下 `offload` 操作的正确性。
- `python/sglang/multimodal_gen/runtime/pipelines_core/executors/sync_executor.py`（模块 执行器；类别 `source`；类型 `core-logic`）：将 `stage` 调用替换为

run\_stage\_with\_context, 与 parallel\_executor 保持一致。

- sgl-model-gateway/src/service\_discovery.rs (模块 服务发现; 类别 source; 类型 data-contract) : 与主题无关的微小改动, 可能是 merge 冲突解决或 lint 修复。保留列出的原因是为了完全性。

关键符号: \_stage\_execution\_context, \_stage\_needs\_version\_counters, run\_stage\_with\_context, \_as\_mutable\_tensor

## 关键源码片段

[python/sglang/multimodal\\_gen/runtime/pipelines\\_core/executors/pipeline\\_executor.py](#)

核心变更文件: 在 execute\_with\_profiling 和 execute\_group\_with\_profiling 中引入 inference\_mode 包装, 新增 stage 级上下文管理和版本计数器需求判断, 是性能提升和兼容性保障的基石。

# pipeline\_executor.py - 关键片段: 执行入口包装与 stage 级上下文

```
class PipelineExecutor(ABC):
    def execute_with_profiling(self, stages, batch, server_args) -> OutputBatch:
        with self.profile_execution(batch, dump_rank=0):
            # 在整个执行外部套上 platform 定义的 inference_mode
            with current_platform.inference_mode():
                batch = self.execute(stages, batch, server_args)
            return batch

    def execute_group_with_profiling(self, stages, batches, server_args):
        with self.profile_execution(batches[0], dump_rank=0):
            with current_platform.inference_mode():
                batches = self.execute_group(stages, batches, server_args)
        return batches

    @staticmethod
    @contextlib.contextmanager
    def _stage_execution_context(stage: "PipelineStage", server_args: ServerArgs):
        # 根据 stage 组件与 server_args 判断是否需要 tensor version counter
        if PipelineExecutor._stage_needs_version_counters(stage, server_args):
            # FSDP / CPU-offload 的 hooks 依赖 version counter, 必须退出 inference_mode
            with torch.inference_mode(False), torch.no_grad():
                yield
            return
        # 其余 stage 保持当前上下文中 (通常是 inference_mode)
        yield

    @staticmethod
    def _stage_needs_version_counters(stage, server_args) -> bool:
        # FSDP 总是需要 version counter
        if server_args.use_fsdp_inference:
```

```

    return True
stage_name = stage._active_component_stage_name()
for use in stage.component_uses(server_args, stage_name):
    cname = use.component_name
    # 检查各种 offload 配置涉及的组件名
    if server_args.dit_cpu_offload and cname in ("transformer", "transformer_2", "video_dit",
"audio_dit"):
        return True
    if server_args.text_encoder_cpu_offload and cname.startswith("text_encoder"):
        return True
    if server_args.image_encoder_cpu_offload and cname in ("image_encoder", "condition_
image_encoder"):
        return True
    if server_args.vae_cpu_offload and cname in ("vae", "video_vae", "audio_vae", "vocoder",
"spatial_upsampler", "condition_image_encoder"):
        return True
return False

```

```

def run_stage_with_context(self, stage, payload, server_args, run_stage):
    # 同步执行器 / 并行执行器通过此方法调用 stage, 自动获取正确的上下文
    with self._stage_execution_context(stage, server_args):
        return run_stage(stage, payload)

```

## python/sglang/multimodal\_gen/test/unit/test\_pipeline\_executor.py

新增测试文件, 全面覆盖 `execute_with_profiling` 和 `execute_group_with_profiling` 在两种平台 (InferenceTensorPlatform 和 NoGradPlatform) 下的 `inference_mode` 行为, 以及 `stage` 上下文在多种 `offload/FSDP` 配置下的版本计数器可用性。

# test\_pipeline\_executor.py - 关键测试片段

```

def test_execute_with_profiling_uses_inference_tensor_platform(monkeypatch):
    # 使用真实 torch.inference_mode 平台
    monkeypatch.setattr(pipeline_executor, "current_platform", _InferenceTensorPlatform)
    executor = _RecordingExecutor()
    with torch.inference_mode(False):
        executor.execute_with_profiling([], _batch(), _server_args())
    assert executor.single_inference_mode is True
    assert executor.single_grad_enabled is False

```

```

def test_stage_context_preserves_version_counters_when_needed(server_args, component_
names):
    stage = _ComponentStage(*component_names)
    with torch.inference_mode():
        with PipelineExecutor._stage_execution_context(stage, server_args):
            tensor = torch.ones(1)
            assert torch.is_inference_mode_enabled() is False
            assert torch.is_grad_enabled() is False
            assert tensor._version == 0 # version counter 可用

```

## python/sglang/multimodal\_gen/runtime/layers/lora/linear.py

新增 `_as_mutable_tensor` 工具方法并在 `merge/unmerge` 路径中调用，确保在 `inference_mode` 下也能安全修改权重。这是修复 LoRA 兼容性的关键。

```
# linear.py - LoRA 合并 / 反合并时处理 inference tensor

@staticmethod
def _as_mutable_tensor(tensor: torch.Tensor) -> torch.Tensor:
    # 如果 tensor 处于 inference 模式，无法直接修改；通过 detach().clone() 创建普通 tensor
    if tensor.is_inference():
        # 必须在 non-inference 模式下 clone，否则 clone 结果仍为 inference tensor
        with torch.inference_mode(False):
            return tensor.detach().clone()
    return tensor

# 在 merge_lora_weights 中应用
data = self.base_layer.weight.data.to(get_local_torch_device()).full_tensor()
data = self._as_mutable_tensor(data) # 确保 data 可写
# ... 修改 data ...
unsharded_base_layer.weight = nn.Parameter(
    self._as_mutable_tensor(data.to(current_device, dtype=target_dtype))
)
# 同样在 bias 和 unmerge 路径调用
```

## 评论区精华

Review 中无人工评论，仅 Gemini code-assist bot 自动评论确认无反馈。从 commit 历史看主要有一次设计迭代：最初在全局作用域启用 `inference_mode`，后发现 FSDP/CPU-offload 阶段需要 version counter，于是改为 stage 粒度回退（commit 'scope inference tensor fallback to stages'）。此外还修复了 LoRA 合并时 inference tensor 不可写问题。这些决策体现了细致的兼容性考虑。

- 暂无高价值评论线程

## 风险与影响

- 风险：

1. FSDP/CPU-offload 路径： `_stage_needs_version_counters` 的组件名匹配列表需要随未来新组件同步更新，遗漏可能导致 version counter 缺失引发 FSDP 失败。
2. NPU 平台：当前 `npu.py` 的 `inference_mode` 直接映射为 `torch.no_grad()`，若 NPU 有原生 `inference` 模式需要单独实现，否则无法获得完整加速。
3. LoRA 动态重配：在 `inference_mode` 下调用 `set_lora_weights` 或 `merge_lora_weights` 必须显式回退到 `inference_mode(False)`，否则 `_as_mutable_tensor` 会 clone 造成额外开销（当前测试中明确在 `torch.inference_mode(False)` 下调用的，但外部调用方若在 `inference_mode` 下调用可能引入微小性能退化）。

4. `service_discovery.rs` 的微小改动 (+3/-1) 与主题无关, 可能是 merge 冲突造成, 但无实际影响。 - 影响: 用户 / 系统: 纯 diffusion pipeline 用户受益约 7% 去噪加速, 无峰值内存增加; 使用 FSDP 或 CPU-offload 的用户行为保持不变; NPU 用户性能暂无提升但功能不受影响。团队: 新增 stage 级别上下文抽象, 降低了后续引入其他平台 inference 模式的门槛。测试覆盖完善, 降低回归风险。 - 风险标记: FSDP/CPU-offload 回退条件依赖组件名列表, NPU 平台 `inference_mode` 仅为 `no_grad` 占位, LoRA 动态重配需调用方显式退出 `inference_mode`

## 关联脉络

- 暂无明显关联 PR