

# PR #23119 完整报告

sgl-project/sglang

[CI] Add per-job uv venv isolation and upgrade CI version to Cuda 13

合并时间: 2026-04-19 20:32

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/23119>

## 执行摘要

- 一句话: 添加每个作业的 uv venv 隔离并升级 CI 到 CUDA 13, 优化依赖管理和环境一致性。
- 推荐动作: 建议技术管理者关注此 PR 的 CI 架构设计, 特别是 uv venv 隔离的实现方式和缓存优化策略, 可作为类似环境管理参考。工程师可精读 `quantization_utils.py` 和 `transformer_load_utils.py` 的变更, 学习 ModelOpt 配置处理和错误恢复模式, 同时注意 `bench_utils.py` 中的性能分析鲁棒性改进。

## 功能与动机

根据 PR 标题和提交历史, 主要动机是改善 CI 环境的稳定性和可维护性。使用 uv venv 隔离可以消除因作业间共享环境导致的依赖冲突, 而升级到 CUDA 13 则为了利用更新的 CUDA 特性和兼容性, 确保测试与最新硬件 (如 H100/H200) 保持一致。

## 实现拆解

1. CI 脚本修改: 更新 `scripts/ci/cuda/ci_install_dependency.sh` 等文件, 添加 `SGLANG_CI_USE_VENV` 开关, 实现每个作业创建独立 venv, 并升级 CUDA 13 依赖安装逻辑, 包括处理 `sgl-kernel` 和 `mooncake` 等包的兼容性。
2. 量化工具增强: 在 `python/sglang/multimodal_gen/runtime/utils/quantization_utils.py` 中添加 `normalize_flat_modelopt_quant_config` 函数, 规范化 ModelOpt 量化配置, 确保与 `diffusers` 兼容; 新增 `_infer_nvfp4_group_size_from_tensors` 用于推断组大小。
3. 加载逻辑优化: 修改 `python/sglang/multimodal_gen/runtime/loader/transformer_load_utils.py`, 引入 `_merge_modelopt_fp4_configs` 函数, 优先使用 `safetensors` 推断的 NVFP4 配置覆盖陈旧的 `config.json` 忽略列表, 提升配置准确性。
4. Diffusers 组件路径处理: 在 `python/sglang/multimodal_gen/runtime/utils/hf_diffusers_utils.py` 中添加 `prepare_diffusers_component_path_for_loading` 函数, 下载组件并修补遗留的扁平 ModelOpt 配置, 确保加载兼容性。
5. 性能分析修复: 更新 `python/sglang/srt/utils/bench_utils.py` 的 `bench_kineto` 函数, 处理空性能分析表的情况, 通过添加 `acc_events=True` 和回退到 `wall-clock` 计时来提高鲁棒性。
6. 测试配套调整: 移除 `test/registered/models/test_transformers_models.py` 中的已弃用 TorchAO 测试, 并将不稳定测试 (如 `HiCache`) 移到 `manual` 文件夹, 确保 CI 通过率; 同时更新模型配置如 `wanvideo.py` 以适配新路径映射。

关键文件:

- `python/sglang/multimodal_gen/runtime/utils/quantization_utils.py` (模块 量化工具; 类别 `source`; 类型 `dependency-wiring`; 符号 `normalize_flat_modelopt_quant_config`, `_infer_nvfp4_group_size_from_tensors`): 关键量化工具文件, 新增 `normalize_flat_modelopt_quant_config` 和 `_infer_nvfp4_group_size_from_tensors` 函数, 用于规范化 `ModelOpt` 配置和推断组大小, 影响扩散模型加载兼容性。
- `python/sglang/multimodal_gen/runtime/loader/transformer_load_utils.py` (模块 加载器; 类别 `source`; 类型 `core-logic`; 符号 `_get_quant_config_name`, `_merge_modelopt_fp4_configs`): 核心加载逻辑文件, 引入 `_merge_modelopt_fp4_configs` 函数, 优先使用 `safetensors` 推断的配置覆盖陈旧的 `config.json`, 提升 `NVFP4` 模型加载准确性。
- `python/sglang/multimodal_gen/runtime/utils/hf_diffusers_utils.py` (模块 `Diffusers` 工具; 类别 `source`; 类型 `dependency-wiring`; 符号 `prepare_diffusers_component_path_for_loading`): `Diffusers` 工具文件, 新增 `prepare_diffusers_component_path_for_loading` 函数, 自动下载组件并修补遗留的扁平 `ModelOpt` 配置, 提升加载兼容性。
- `python/sglang/srt/utils/bench_utils.py` (模块 性能工具; 类别 `source`; 类型 `dependency-wiring`): 性能分析工具文件, 修复 `bench_kineto` 函数以处理空性能分析表, 添加回退到 `wall-clock` 计时, 提高在 `CUDA 13` 等环境下的鲁棒性。
- `scripts/ci/cuda/ci_install_dependency.sh` (模块 `CI` 脚本; 类别 `infra`; 类型 `infrastructure`): 核心 `CI` 依赖安装脚本, 大幅修改以支持 `uv venv` 隔离和 `CUDA 13` 升级, 涉及依赖安装、缓存清理和路径配置。

关键符号: `normalize_flat_modelopt_quant_config`,  
`_infer_nvfp4_group_size_from_tensors`, `_get_quant_config_name`,  
`_merge_modelopt_fp4_configs`, `prepare_diffusers_component_path_for_loading`

## 关键源码片段

### `python/sglang/multimodal_gen/runtime/utils/quantization_utils.py`

关键量化工具文件, 新增 `normalize_flat_modelopt_quant_config` 和 `_infer_nvfp4_group_size_from_tensors` 函数, 用于规范化 `ModelOpt` 配置和推断组大小, 影响扩散模型加载兼容性。

```
def normalize_flat_modelopt_quant_config(
    quant_cfg: dict[str, Any] | None,
) -> dict[str, Any] | None:
    """填充必需的 diffusers 字段用于扁平 ModelOpt 组件配置。"""
    if not isinstance(quant_cfg, dict) or quant_cfg.get("quant_method") != "modelopt":
        return quant_cfg # 若非 ModelOpt 配置, 直接返回
    quant_algo = str(
        quant_cfg.get("quant_algo")
        or quant_cfg.get("quantization", {}).get("quant_algo")
        or ""
    ).upper()
    if not quant_algo:
```

```

    return quant_cfg # 若无量化算法信息, 保持原样
normalized = dict(quant_cfg)
normalized.setdefault("quant_type", quant_algo) # 添加 quant_type 字段以适配 diffusers
return normalized

```

## python/sglang/multimodal\_gen/runtime/loader/transformer\_load\_utils.py

核心加载逻辑文件, 引入 `_merge_modelopt_fp4_configs` 函数, 优先使用 safetensors 推断的配置覆盖陈旧的 config.json, 提升 NVFP4 模型加载准确性。

```

def _merge_modelopt_fp4_configs(
    existing_config: Optional[QuantizationConfig],
    inferred_config: Optional[QuantizationConfig],
) -> Optional[QuantizationConfig]:
    """优先使用 safetensors 推断的 NVFP4 布局覆盖陈旧的 config.json 忽略列表。

    一些 ModelOpt NVFP4 转换器仓库在 config.json 中提供扁平 quantization_config,
    但其 ignore 列表可能滞后于实际检查点内容。safetensors 分片是哪些模块保持 BF16
    回退的真实来源,
    因此当能从分片推断 NVFP4 配置时, 应使用其排除列表, 同时保留显式仓库级开关如 swap_weight_
    nibbles。
    """
    if inferred_config is None:
        return existing_config # 无推断配置, 返回现有配置
    if _get_quant_config_name(inferred_config) != "modelopt_fp4":
        return existing_config or inferred_config # 非 ModelOpt FP4, 合并或返回推断配置
    if existing_config is None:
        return inferred_config # 无现有配置, 直接使用推断配置
    if _get_quant_config_name(existing_config) != "modelopt_fp4":
        return existing_config # 现有配置非 ModelOpt FP4, 保持原样
    # 比较排除列表, 优先使用推断配置的列表
    existing_excludes = getattr(existing_config, "exclude_modules", []) or []
    inferred_excludes = getattr(inferred_config, "exclude_modules", []) or []
    if inferred_excludes != existing_excludes:
        logger.warning(
            "Overriding ModelOpt NVFP4 exclude_modules from config.json with "
            "safetensors-inferred layout (%d -> %d entries).",
            len(existing_excludes),
            len(inferred_excludes),
        )
    # 保留现有配置中的关键属性
    inferred_config.packed_modules_mapping = getattr(
        existing_config, "packed_modules_mapping", {}
    )
    inferred_config.swap_weight_nibbles = getattr(
        existing_config, "swap_weight_nibbles", True
    )
    inferred_config.checkpoint_uses_packed_qkv = getattr(
        inferred_config, "checkpoint_uses_packed_qkv", False
    ) or getattr(existing_config, "checkpoint_uses_packed_qkv", False)

```

```
if getattr(inferred_config, "group_size", None) is None:
    inferred_config.group_size = getattr(existing_config, "group_size", None)
return inferred_config # 返回合并后的配置
```

## 评论区精华

由于 review 评论为空，核心讨论体现在提交历史中的多次迭代和决策。例如，提交消息显示对 FlashInfer JIT 缓存路径稳定性和 deep\_gemm NVRTC 缓存重用的优化，以避免因 per-job venv 路径变化导致的重新编译开销。决策包括回滚无效尝试、添加回退机制和选择性缓存清理。

- 暂无高价值评论线程

## 风险与影响

- 风险：

1. 环境兼容性风险：CUDA 13 升级可能导致现有内核或库行为变化，需全面测试覆盖，如 bench\_utils.py 中性能分析器可能因 CUDA 版本返回空表。
2. CI 性能开销：每个作业创建独立 venv 可能增加启动时间和资源消耗，需监控 CI 作业时长影响。
3. 配置变更影响：量化配置处理逻辑变更（如 normalize\_flat\_modelopt\_quant\_config）可能影响模型加载准确性，需验证向后兼容性，避免生产环境回归。
4. 测试稳定性：移动测试到 manual 文件夹可能掩盖潜在问题，需确保关键路径仍有充分测试覆盖。

- 影响：影响范围：主要针对 CI 流程和开发团队，不影响生产代码核心逻辑。影响程度：中等；CI 环境更稳定和可维护，减少因依赖问题导致的失败，但团队需适应新配置和可能的 CI 时间增加。具体影响：

- 用户：几乎无影响，除非依赖 CI 产出。
- 系统：提升 CI 环境隔离性和一致性，支持 CUDA 13 新特性。
- 团队：简化依赖管理，但需更新本地开发环境以匹配 CI 变更。
- 风险标记：环境兼容性风险，CI 性能开销，配置变更影响

## 关联脉络

- PR #23108 Update CI\_PERMISSIONS: 同属 CI 基础设施变更，涉及权限管理，与本 PR 的 CI 环境升级相辅相成。
- PR #23103 Apply HF transformers patches from sglang init: 涉及依赖管理和补丁应用，与本 PR 的依赖隔离和升级主题相关。