

PR #15236 完整报告

sgl-project/sglang

[CI] Add consistency test in CI

合并时间: 2026-04-07 09:50

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/15236>

执行摘要

本 PR 为 SGLang 仓库的 CI 管道新增了 diffusion 模型输出一致性测试，通过生成 ground truth 并基于 CLIP 余弦相似度进行验证，旨在防止代码变更导致的精度回归。关键改动涉及测试框架扩展、阈值配置定义和 LoRA 权重加载优化，影响范围主要覆盖 multimodal diffusion 测试模块和 CI 基础设施，是一个有意义的测试基础设施改进。

功能与动机

为什么做: PR body 明确指出，现有 CI 管道缺乏严格的精度检查点 (guardrails)，这可能导致意外合并引起精度下降或模型性能回归的 PR。例如，diffusion 模型生成输出时，微小的代码变更可能引发图像或视频质量波动，而现有测试未能充分捕获。因此，引入一致性测试作为安全网，确保模型输出在代码演进中保持稳定。

关键表述引用: 从 PR body: "CI pipeline lacks strict accuracy checkpoints (guardrails). This risks accidentally merging PRs that cause precision degradation or model performance regression."

实现拆解

实现按模块拆解如下:

- CI 工作流层: 在 `.github/workflows/pr-test-amd.yml` 和 `.github/workflows/pr-test-amd-rocm720.yml` 中添加 `SGLANG_SKIP_CONSISTENCY` 环境变量，允许在特定平台 (如 AMD) 跳过检查，减少 CI 开销。

```
yaml -e SGLANG_SKIP_CONSISTENCY=1
```
- 测试框架层: 在 `python/sglang/multimodal_gen/test/server/test_server_common.py` 中新增 `_validate_consistency` 方法，核心逻辑包括:
 - 检查环境变量跳过条件。
 - 加载 ground truth (本地或从 `sglang-ci-data` 仓库远程获取)。
 - 使用 CLIP 模型计算输出与 ground truth 的余弦相似度，并对比 SSIM、PSNR 和平均绝对差。
 - 根据阈值判断通过与否，失败时输出详细错误信息。
- 配置管理层: 新增 `python/sglang/multimodal_gen/test/server/consistency_threshold.json` 文件，定义各测试案例的阈值，例如:

```
json "flux_image_t2i": { "clip_threshold": 0.92, "ssim_threshold": 0.95, "psnr_threshold": 24.0, "mean_abs_diff_threshold": 8.0 }
```

4. LoRA 扩展层：修改 `lora_pipeline.py`、`server_args.py` 和 `hf_diffusers_utils.py`，添加 `lora_weight_name` 参数，支持从多文件 LoRA 仓库中确定性加载特定权重文件，避免随机选择导致的不一致。
5. 测试工具层：在 `test_utils.py` 中增强函数如 `compare_with_gt`，并新增 `test_consistency_metrics.py` 单元测试，验证像素级指标的正确性。

评论区精华

由于 review 评论为空，主要讨论体现在 Issue 评论中：

- 协作流程调整：维护者 `mickqian` 多次要求 `rebase` 和重跑 CI（例如 `"/rerun-failed-ci"`），贡献者 `Prozac614` 及时响应，显示高效协作。

`mickqian`: "please rebse" `Prozac614`: "Done"

- 外部依赖处理：`mickqian` 提到需先修复 PR #22059 中的 `flux` 问题，凸显跨 PR 依赖的挑战。

`mickqian`: "we need to fix all the flux issues in #22059 before we proceed"

- 社区协助：`shljessie` 主动提供帮助解决合并冲突和生成 `ground truth` 文件，体现团队协作精神。

`shljessie`: "I'd be happy to help finish this PR! I can resolve the merge conflicts and generate the missing ground truth files."

风险与影响

具体风险：

1. 外部依赖风险：CLIP 模型和 `ground truth` 文件托管在外部仓库（如 Hugging Face 和 `snglang-ci-data`），网络波动或仓库变更可能导致测试失败。
2. 阈值配置敏感：`consistency_threshold.json` 中的值需根据模型和硬件动态调整；例如，"`wan2_1_t2v_1_3b_lora_1gpu`" 案例的 `clip_threshold` 仅为 0.54，表明某些场景容忍度低，易引发误报。
3. 性能开销：一致性检查涉及图像 / 视频帧提取、CLIP 嵌入计算和远程文件加载，可能显著增加 CI 运行时间，尤其在资源受限环境中。
4. 兼容性影响：LoRA 权重名称的添加可能破坏现有使用多文件 LoRA 的脚本，需通过默认值 (`weight_name=None`) 保持向后兼容。

影响评估：

- 用户影响：间接提升 `diffusion` 模型输出的可靠性，防止回归影响应用质量。
- 系统影响：CI 更严格，但增加了维护负担（如管理 `ground truth` 文件）；LoRA 扩展提高了部署灵活性。
- 团队影响：开发者需学习新测试流程，贡献时可能需操作外部仓库，增加了入门门槛。

关联脉络

从近期历史 PR 分析，本 PR 是测试和一致性保障演进的一部分：

- 相关 PR: PR #21849 (VLM 一致性修复) 和 PR #22194 (CI 阈值调整) 都涉及 consistency 标签, 显示团队持续加强测试可靠性。
- 功能线演进: 本 PR 专注于 diffusion 模型, 与仓库中其他 speculative-decoding、multimodal 测试 PR (如 #22199、#21425) 形成互补, 共同构建全面的质量保障体系。
- Issue 关联: 讨论中提及 PR #22059 的 flux 问题, 表明本 PR 的推进依赖更大范围的问题解决, 反映了仓库中跨模块测试的复杂性。