

PR #5713 完整报告

verl-project/verl

[3/n][reward] feat: flowgrp0 - support image-based rewards (rule-based & genrm)

合并时间: 2026-03-26 15:54

原文链接: <http://prhub.com.cn/verl-project/verl/pull/5713>

执行摘要

本 PR 为 verl 仓库添加图像奖励支持，是 FlowGRPO 训练的关键扩展。通过新增 VisualRewardManager 类和修改奖励循环，实现了生成奖励模型（如 OCR）和规则奖励（如 JPEG 压缩性），使图像生成模型训练成为可能。变更影响奖励系统架构，提升了多模态能力，并提供测试验证。

功能与动机

动机源于之前奖励模型不支持图像 / 视频输入，限制了 FlowGRPO 训练的应用。PR body 指出这是 #4639 和 #5297 的后续工作，旨在支持 QwenImage 训练。Issue 评论中 yyDing1 确认了此需求，表示之前实现缺少对视觉输入的支持。

实现拆解

- 新增 VisualRewardManager 类：位于 verl/experimental/reward_loop/reward_manager/visual.py，继承自 RewardManagerBase，专门处理视觉响应奖励计算，支持异步和同步评分。
- 修改奖励循环：在 verl/experimental/reward_loop/reward_loop.py 中，_preprocess_reward_inputs 方法现在根据响应维度（3D 为图像，其他为文本）区分输入类型，实现多模态预处理。
- 奖励函数实现：
 - compute_score_ocr：位于 examples/flowgrp0_trainer/reward_fn.py，使用生成奖励模型进行 OCR 评分，基于 Levenshtein 距离计算匹配度。
 - jpeg_compressibility.compute_score：位于 verl/utils/reward_score/jpeg_compressibility.py，提供基于 JPEG 压缩性的规则奖励。
- 工具函数：新增 pil_image_to_base64 和 prepare_query_for_multi_modal（verl/utils/experimental/reward_utils.py），支持图像到 base64 转换和多模态查询构建。
- 测试与 CI：新增单元测试文件 tests/experimental/reward_loop/test_visual_reward_manager.py，并更新 CI workflow（.github/workflows/reward_model_vllm.yml）以集成测试。

评论区精华

- ZeroDivisionError 风险：gemini-code-assist[bot] 指出 compute_score_ocr 函数在 ground_truth 为空字符串时可能引发除以零错误，建议添加边界检查。作者在后续提交中

可能已修复，但 review 中未详细说明。

- 设计决策：SamitHuang 提议将 ImageRewardManager 重命名为 VisualRewardManager 以覆盖图像和视频，作者采纳并修改，体现了设计的前瞻性。
- 代码结构优化：SamitHuang 建议通过 get_default_compute_score 函数改进 load_reward_manager 逻辑，作者实现此函数，简化了配置依赖。
- 奖励函数放置：讨论规则奖励应放在 jpeg_compressibility.py，而 OCR 作为生成奖励模型示例，强调了模块化设计原则。

风险与影响

风险分析：

- 正确性风险：compute_score_ocr 函数可能存在边界条件未处理，如空字符串导致的崩溃。
- 配置风险：奖励管理器类型依赖字符串配置，错误配置可能引发运行时问题。
- 性能风险：异步图像转换和网络请求可能增加延迟或内存消耗。
- 测试风险：新增测试仅覆盖基本场景，缺乏全面边缘案例验证。

影响分析：

- 用户影响：研究人员现在可以使用图像奖励进行 FlowGRPO 训练，扩展了应用范围，尤其是图像生成领域。
- 系统影响：奖励系统变得更加复杂，需要维护多模态处理逻辑，可能影响系统稳定性和性能。
- 团队影响：开发人员需学习新模块，可能增加代码审查和维护成本。

关联脉络

本 PR 是系列工作的一部分，关联 PR #4639、#5297、#5616 和 #5716，共同推进 FlowGRPO 训练能力，特别是针对 QwenImage。从近期历史 PR 看，仓库正专注于扩展训练功能（如 #5745 的教师模型支持），本 PR 在多模态奖励方面填补了空白，显示了架构向更通用方向演进。