

每日论文分析 – 2026–02–17

标签: [RLVR], [数据筛选], [Agent], [Benchmark]

1. Experiential Reinforcement Learning

标题: Experiential Reinforcement Learning **作者:** Taiwei Shi, Sihao Chen, Bowen Jiang, Linxin Song, Longqi Yang, Jieyu Zhao **机构:** Microsoft
ArXiv ID: 2602.13949 **标签:** [微软], [强推]

动机

问题形式化: 给定一个任务 T , 传统 RL 需要模型在稀疏、延迟的奖励信号下学习。设 r_t 为时刻 t 的奖励, 传统的 RL 目标函数为最大化期望累积奖励 $\mathbf{E}_\tau \left[\sum_{t=0}^T \gamma^t r_t \right]$ 。然而在真实场景中, 奖励 r_t 往往是稀疏的——只有任务完成或失败时才给出反馈。这导致信用分配 (credit assignment) 问题异常困难。

为什么是现在: RLHF → DPO → GRPO → RLVR 的演进路线已经证明了显式反馈的重要性。然而, self-reflection 在 RL 训练中的角色尚未被充分探索。这类似于物理学中的”准粒子”概念——我们需要显式地建模”反馈粒子”如何改变系统状态。

核心假设

显式假设: – 显式的 reflection 步骤可以帮助模型更有效地将稀疏反馈转化为行为改变 – 两阶段尝试 (initial attempt → refined attempt) 比单次尝试更有效

隐式假设: – 环境反馈是可获得的 (稀疏但存在) – 模型具备生成 reflection 的能力

评估: 假设有合理性。人类学习明确依赖内省, RL 中的 Hindsight Experience Replay (HER) 也有类似思想。实验结果支持: +81% 在复杂多步环境, +11% 在工具使用推理任务。

技术贡献

1. **Experience–Reflection–Consolidation 循环**: 将原始的 single-step RL 扩展为三阶段循环
2. **Structured Behavioral Revision**: 不是让模型”隐式推断”失败原因，而是让模型显式生成 reflection 来指导下一步行动
3. **部署时无额外推理开销**: consolidation 后，policy 直接改进，无需 reflection 模块

实验设计

- **基线**: PPO, GRPO, DPO 等强 RL 基线
- **环境**: Sparse-reward control environments, agentic reasoning benchmarks
- **指标**: 任务成功率, 学习效率

消融实验

关键组件贡献未详细披露，但从实验结果看，reflection 步骤是核心。

潜在影响

物理学类比: 这项工作类似于在 RL 系统中引入”显式相位”。传统 RL 是”相干态”，依赖隐式干涉；ERL 引入”激发态” (reflection)，允许系统显式地改变基态。

若成立: Self-reflection 将成为 RL 训练的标准组件，类似 CoT 成为推理的标准。

严厉审视

- **假设验证**: 实验仅在控制环境和 agentic benchmark 上验证，未在真实 LLM 训练中验证
 - **可扩展性**: 两阶段生成会增加训练时间和计算成本
 - **对比缺失**: 未与 Implicit REINFORCE 等方法对比
-

2. Embed–RL: Reinforcement Learning for Reasoning–Driven Multimodal Embeddings

标题: Embed–RL: Reinforcement Learning for Reasoning–Driven Multimodal Embeddings **作者:** Haonan Jiang, Yuji Wang, Yongjie Zhu, Xin Lu, Wenyu Qin, Meng Wang, Pengfei Wan, Yansong Tang **机构:** Tsinghua University **ArXiv ID:** 2602.13823 **标签:** [清华], [有趣但有缺陷]

动机

问题形式化: 给定一个 query 和一个 target, 构建一个 multimodal embedding e 使得相似度 $sim(e_q, e_t)$ 与任务相关。当前方法使用 CoT 生成文本分析, 但这些 reasoning 与 retrieval 目标无关。

数学描述: 优化目标不是简单的对比学习 $\mathcal{L} = -\log \frac{\exp(sim(q, t^+)/\tau)}{\sum \exp(sim(q, t)/\tau)}$, 而是要最大化 reasoning quality 与 retrieval quality 的联合分布一致性。

核心假设

1. **Traceability CoT (T–CoT):** 生成的 reasoning trace 应该包含可追溯的 multimodal cues
2. **Embedder–Guided RL:** Embedder 可以提供显式的 reward signal 指导 Reasoner

评估: 假设新颖但有一定道理。T–CoT 的思想类似于在 retrieval 中引入”注意力证据”。

技术贡献

1. **EG–RL 框架:** Embedder 提供 reward, 训练 Reasoner
2. **T–CoT:** 提取 multimodal cues 作为 retrieval 的证据
3. **在有限资源下超越 baseline**

潜在问题

- **Reward 设计:** Embedder 提供的 reward 是否真的与 retrieval 质量相关? 未详细说明
- **多模态对齐:** multimodal cues 的提取是否可靠?
- **与现有方法的对比:** 缺乏与 E5, BGE 等成熟方法的对比

评价

[清华] 机构背景强，方向有趣 (RL + Multimodal)，但实验规模和对比不够充分。

3. A Critical Look at Targeted Instruction Selection

标题: A Critical Look at Targeted Instruction Selection: Disentangling What Matters (and What Doesn't) **作者:** Nihal V. Nayak, Paula Rodriguez-Diaz, Neha Hulkund, Sara Beery, David Alvarez-Melis **机构:** Harvard DCML **ArXiv ID:** 2602.14696 **标签:** [哈佛], [强推], [数据构建]

动机

问题形式化: 给定一个大型 instruction 数据池 \mathcal{D} 和一个小的目标任务 query 集 \mathcal{Q} ，选择 \mathcal{D} 的子集 \mathcal{S} 来微调 LLM，使得在 \mathcal{Q} 上表现最好。

数学形式化: 目标是最小化选中的 subset \mathcal{S} 与 query 集 \mathcal{Q} 之间的”距离”：

$$\min_{\mathcal{S} \subseteq \mathcal{D}, |\mathcal{S}|=k} \text{dist}(\mathcal{S}, \mathcal{Q})$$

核心假设

关键发现: 只有 gradient-based 的数据表示能够一致地预测性能。现有方法的差异主要来自： 1. 数据表示 (data representation): semantic embedding vs gradient-based 2. 选择算法 (selection algorithm): greedy, importance sampling, etc.

技术贡献

1. **统一框架:** 将多种选择算法统一为 approximate distance minimization
2. **理论保证:** 提供了 generalization bound
3. **实践指导:** 低预算时用 gradient-based + greedy round-robin，高预算时差异消失

实验设计

- **基线:** 多种 instruction selection 方法
- **数据集:** 多个 LLM 和任务
- **指标:** 性能预测一致性

评价

[哈佛] [强推] — 这是数据构建方向的关键工作。对于用户的”数据筛选方法”兴趣高度相关。

与过去工作对比

- 之前的工作往往 entanglement 各组件
 - 缺乏系统性的对比框架
 - 未考虑 zero-shot baseline
-

4. WebWorld: A Large-Scale World Model for Web Agent Training

标题: WebWorld: A Large-Scale World Model for Web Agent Training
作者: Zikai Xiao, Jianhong Tu, Chuhang Zou, Yuxin Zuo, Zhi Li, Peng Wang, Bowen Yu, Fei Huang, Junyang Lin, Zuozhu Liu
机构: Qwen (Alibaba)
ArXiv ID: 2602.14721
标签: [阿里], [Agent], [Benchmark]

动机

问题: Web agent 训练需要大量 trajectories, 但真实网络训练受限于延迟、速率限制和安全风险。

形式化: 构建一个世界模型 W 来模拟网络环境, 使得 agent 可以在 W 中高效训练, 然后迁移到真实网络。

核心假设

1. 大规模 (1M+ interactions) 训练可以产生有效的世界模型
2. 世界模型可以泛化到不同域 (code, GUI, game)

技术贡献

1. 首个大规模开放网络模拟器: 1M+ open-web interactions
2. WebWorld-Bench: 9 维度 dual metrics
3. 跨域泛化: 代码、GUI、游戏环境

实验结果

- 内在评估: WebWorld 性能 comparable to Gemini-3-Pro

- 外在评估: Qwen3–14B 在 WebArena 上 +9.2%
- 推理时搜索: 超越 GPT–5 作为世界模型

评价

[阿里] — 大规模数据构建 + Agent 训练的代表作。1M trajectories 是显著规模。

与过去工作对比

工作	规模	环境
WebArena	~10K	封闭
WebShop	~2K	模拟
WebWorld	1M+	开放

总结

论文	方向	评分	标签
Experiential RL	RLVR	★★★★★	[微软], [强推]
Targeted Instruction Selection	数据构建	★★★★★	[哈佛], [强推]
Embed–RL	RL + Multimodal	★★★	[清华], [有趣但有缺陷]
WebWorld	Agent/Benchmark	★★★★★	[阿里]

关键发现

1. **RLVR 演进:** ERL 展示了 explicit self–reflection 在 RL 中的价值，延续了从 RLHF → DPO → GRPO → RLVR 的演进路线
2. **数据筛选:** Harvard 工作揭示了 gradient–based representation 在 instruction selection 中的核心作用
3. **Agent 训练:** WebWorld 证明了大规模世界模型训练的可行性

明日预告

- 关注 DeepSeek/OpenAI 的新 RL 工作

- 关注 MoRL (Peking) 的运动推理