

2026-02-21 每日论文分析

标签: [RLHF], [LLM Alignment], [Safety], [Scaling Laws], [Agent]

论文1: References Improve LLM Alignment in Non-Verifiable Domains

基本信息

- 标题: References Improve LLM Alignment in Non-Verifiable Domains
- ArXiv ID: 2602.16802
- 作者: Kejian Shi, Yixin Liu, Peifeng Wang, Alexander R. Fabbri, Shafiq Joty, Arman Cohan
- 机构: Yale NLP Lab
- 标签: [Yale NLP], [RLVR], [LLM Alignment], [强推]

动机

问题形式化: RLVR (Reinforcement Learning with Verifiable Rewards) 在数学、代码等可验证领域取得了显著成功，但对于缺乏 ground-truth verifier 的非可验证领域（如 LLM alignment），无法直接应用。当前的 self-improvement 方法依赖于 reference-free judges，这些 judges 缺乏客观的评估标准，导致 reward signal 噪声大、训练不稳定。

核心洞察: 能否利用 reference (参考输出) 作为“软验证器”来增强 LLM-as-judge 的评估能力？形式化地，设 $r(x, y, ref)$ 为参考引导的评分函数，其中 y 是待评估输出， ref 是参考输出。目标是设计协议使得 r 能够更准确地捕获 y 相对于 ref 的质量差异。

公式讲解

核心方法包含两个阶段：

阶段1 – 参考引导的评估协议:

$$\text{Score}(y|x, ref) = \text{Judge}(x, y, \text{context} = ref) - \text{Judge}(x, ref, \text{context} = y)$$

其中 Judge 可以是任意 LLM。论文设计了多种协议来增强评估准确性，包括：
– 直接比较 (Direct Comparison)
– 参考条件评分 (Reference-Conditioned Scoring)
– 双重参考 (Dual Reference)

阶段2 – 参考引导的 self-improvement: 使用改进后的 judge 进行 RL 训练，reward 信号为：

$$R(x, y) = \mathbb{1}[\text{Judge}_\theta(x, y, ref) > \tau]$$

实验设计

- 基线：
 - Direct SFT on reference outputs
 - Reference-free self-improvement (如 Self-Rewarding)

- ArmoRM (finetuned reward model)
- UltraFeedback
- **数据集:**
 - AlpacaEval 2.0
 - Arena-Hard
 - 多种开源指令跟随数据集
- **评估指标:**
 - AlpacaEval (LC win rate)
 - Arena-Hard (win rate)
- **配置:**
 - Llama-3-8B-Instruct
 - Qwen2.5-7B
 - 使用 GPT-4o 作为 reference provider

实验结论

1. **主实验结果:**
 - Llama-3-8B-Instruct: AlpacaEval 73.1% (+20.2 vs SFT), Arena-Hard 58.7% (+17.1 vs SFT)
 - Qwen2.5-7B: AlpacaEval 70.0% (+5.3 vs reference-free SI), Arena-Hard 74.1% (+3.6 vs reference-free SI)
 - 性能接近 ArmoRM (一个强大的 finetuned reward model)
2. **关键发现:**
 - 较弱的 judge 从 frontier model references 中获益更多
 - 较强的 judge 可以从高质量 human-written references 中进一步提升
 - 参考质量与 judge 能力之间存在互补关系

严厉审视

- **假设验证:** “reference 能够提供额外的评估信号”这一假设被实验支持，但是否存在 reference bias (过度依赖参考)？
- **统计显著性:** 论文报告了具体数值但未提供置信区间，这是中等严重性问题
- **混淆因素:** reference-free baseline 是否使用了相同的训练数据量？需要更严格的消融实验
- **prior work:** 忽略了 Instruction-following evaluation (IFE) 等相关工作

与过去工作对比

方法	AlpacaEval	Arena-Hard	参数量
Direct SFT	~50%	~40%	8B
Reference-free SI	~65%	~52%	8B
ArmoRM	~72%	~58%	8B
本文方法	73.1%	58.7%	8B

本文方法在无需额外 reward model training 的情况下达到了与 ArmoRM 相当的性能。

论文2: NeST: Neuron Selective Tuning for LLM Safety

基本信息

- **标题:** NeST: Neuron Selective Tuning for LLM Safety
- **ArXiv ID:** 2602.16835

- **作者:** Sasha Behrouzi, Lichao Wu, Mohamadreza Rostami, Ahmad-Reza Sadeghi
- **机构:** Technical University of Darmstadt – Information Systems
- **标签:** [Safety], [Parameter-Efficient Tuning], [LLM]

动机

问题形式化: 现有的安全对齐方法存在效率–效果权衡问题：
 – Full fine-tuning:
 效果好但参数量大（更新整个模型）
 – LoRA: 参数高效但安全提升不稳定、对设计选择敏感
 – Circuit breakers: 不修改权重但无法塑造内部表示

本文提出：能否只更新少量与安全相关的神经元，同时通过神经元聚类确保功能一致性？

公式讲解

核心优化目标:

$$\min_{\Delta_S} \mathcal{L}_{\text{safety}}(\theta + \Delta) \quad \text{s.t.} \quad \|\Delta_{\bar{\mathcal{S}}}\|_2 = 0$$

其中 \mathcal{S} 是选中的安全相关神经元集合， $\bar{\mathcal{S}}$ 是冻结部分。

神经元聚类:

$$\mathcal{C} = \{c_1, c_2, \dots, c_k\}, \quad \text{where } c_i = \{n \mid \text{sim}(n, n_i^*) > \theta\}$$

同一聚类内的神经元共享相同的更新方向，确保功能一致性。

实验设计

- **基线:** Full fine-tuning, LoRA, Circuit breakers
- **数据集:** 10+ 开放权重 LLMs, 多个安全 benchmark
- **评估指标:** Attack Success Rate (ASR), Capability degradation
- **配置:** 平均仅 0.44M 可训练参数

实验结论

1. **主实验结果:**
 - ASR 从 44.5% 降至 4.36% (90.2% reduction)
 - 相比 full fine-tuning 减少 17,310x 参数
 - 相比 LoRA 减少 9.25x 参数
2. **关键发现:**
 - 神经元聚类确保了更新的稳定性和一致性
 - 安全相关神经元具有可解释的语义聚类模式

严厉审视

- **局限:** 仅在开源模型上验证，未测试闭源 frontier models
- **假设:** “安全神经元可识别且可分离”需要更严格的验证
- **数据:** 训练数据构建细节不足

论文3: CrispEdit: Low-Curvature Projections for Scalable Non-Destructive LLM Editing

基本信息

- **标题:** CrispEdit: Low-Curvature Projections for Scalable Non-Destructive LLM Editing
- **ArXiv ID:** 2602.15823
- **作者:** Zarif Ikram, Arad Firouzkouhi, Stephen Tu, Mahdi Soltanolkotabi, Paria Rashidinejad
- **机构:** University of Southern California
- **标签:** [Model Editing], [Constrained Optimization], [Bregman Divergence]

动机

问题形式化: LLM editing 的核心挑战是 capability preservation——成功修改目标行为的编辑方法可能会“暗中”破坏通用能力，产生类似 proxy/reward hacking 的退化行为。

形式化为一个 constrained optimization problem:

$$\min_{\Delta} \mathcal{L}_{\text{edit}}(\theta + \Delta) \quad \text{s.t.} \quad \mathcal{L}_{\text{cap}}(\theta + \Delta) \leq \epsilon$$

公式讲解

Bregman Divergence 约束: 使用 Bregman divergence $D_\phi(\theta + \Delta, \theta)$ 来表达 capability constraint，其二次形式产生精确的 Gauss–Newton Hessian:

$$\mathcal{L}_{\text{cap}}(\theta + \Delta) \approx \mathcal{L}_{\text{cap}}(\theta) + \nabla \mathcal{L}_{\text{cap}}^T \Delta + \frac{1}{2} \Delta^T H_{GN} \Delta$$

低曲率子空间投影: 将编辑更新投影到 capability–loss 景观的低曲率子空间：

$$\Delta_{\text{projected}} = P\Delta, \quad P = I - V(V^T V)^{-1}V^T$$

其中 V 由 H_{GN} 的特征向量张成。

K–FAC 近似: 使用 Kronecker–factored approximate curvature (K–FAC) 和 novel matrix–free projector 使二阶方法在大规模 LLM 上变得高效。

实验设计

- **基线:** ROME, MEMIT, FT, LoRA, etc.
- **数据集:** standard model–editing benchmarks (ZsRE, COUNTERFACT, etc.)
- **评估指标:** Edit success rate, Capability degradation

实验结论

1. **主实验结果:**
 - 高编辑成功率
 - Capability degradation 平均低于 1%
 - 显著优于 prior editors
2. **关键发现:**
 - Bregman divergence 能够精确表达 capability constraint
 - K–FAC 近似在保持精度的同时大幅降低计算开销

严厉审视

- **计算开销:** 二阶方法仍然比 gradient descent 方法慢，需要权衡
 - **扩展性:** 在超大规模模型 (70B+) 上的效果未验证
 - **假设:** “低曲率方向对应低 capability impact”是核心假设，需要更多分析
-

论文4: NESSiE: The Necessary Safety Benchmark

基本信息

- **标题:** NESSiE: The Necessary Safety Benchmark – Identifying Errors that should not Exist
- **ArXiv ID:** 2602.16756
- **作者:** Johannes Bertram, Jonas Geiping
- **机构:** Yale / TU Darmstadt
- **标签:** [Safety Benchmark], [Evaluation], [LLM]

动机

问题形式化: 即使是 state-of-the-art LLMs 也在极其简单的安全任务上失败。
论文提出一个最小化的测试集，包含信息和访问安全相关的 minimal test cases，来揭示即使在低复杂度任务下也存在安全相关失败。

核心问题: 如果模型在如此简单的任务上都无法保证安全，如何相信它们在复杂场景中的表现？

公式讲解

Safe & Helpful (SH) 指标:

$$SH = 2 \times \frac{\text{Safe} \times \text{Helpful}}{\text{Safe} + \text{Helpful}}$$

这个指标允许直接比较安全性和有用性两个需求，揭示模型在两者之间的偏见。

实验设计

- **数据集:** NESSiE benchmark (信息安全和访问安全的 minimal test cases)
- **模型:** 多个 state-of-the-art LLMs

实验结论

1. **主实验结果:**
 - 即使 SOTA LLMs 也无法在 NESSiE 上达到 100%
 - 模型普遍偏向“helpful”而非“safe”
 - 禁用 reasoning 或添加 benign distraction context 会显著降低性能
2. **关键发现:**
 - 通过简单的 safety sanity check 可以揭示模型的根本性弱点
 - 作为部署前的必要条件，NESSiE 是最低安全保障

严厉审视

- **覆盖范围:** 仅测试 minimal cases，不能保证一般情况下的安全性
 - **攻击类型:** 未测试 adversarial attacks
 - **实际意义:** “necessary but not sufficient” 需要明确传达给社区
-

论文5: Hardware Co-Design Scaling Laws via Roofline Modelling for On-Device LLMs

基本信息

- 标题: Hardware Co-Design Scaling Laws via Roofline Modelling for On-Device LLMs
- ArXiv ID: 2602.10377
- 作者: Luoyang Sun, Jiwen Jiang, Yifeng Ding, et al. (12 authors)
- 机构: 阿里 (Alibaba) – 似乎来自阿里团队
- 标签: [Scaling Laws], [Hardware–Software Co-Design], [On–Device LLM]

动机

问题形式化: VLA (Vision–Language–Action Models) 和 on-device LLM 部署面临 accuracy–latency tradeoff。每个硬件平台需要定制化的架构解决方案。传统方法依赖经验和试错，效率低下。

核心目标: 建立一个原则性框架，将模型训练 loss 表示为架构超参数的显式函数，并将推理延迟通过 roofline modeling 表征。

公式讲解

训练 loss 建模:

$$\mathcal{L}(N, D, H, Q, L, \dots) = \alpha N^{-\beta} \cdot f(H, Q, L, \dots)$$

其中 N 是参数量， D 是数据量， H 是 hidden size， Q 是 head 数量， L 是层数。

Roofline 延迟模型:

$$T_{\text{infer}} = \underbrace{\frac{2N}{\text{Bandwidth}}}_{\text{Memory-bound}} + \underbrace{\frac{N}{\text{FLOPs}}}_{\text{Compute-bound}}$$

通过识别每个操作是 memory-bound 还是 compute-bound，可以精确预测延迟。

Pareto 前沿: 在 accuracy–latency 空间中识别 Pareto 最优解：

$$\mathcal{P} = \{(latency, perplexity) \mid \nexists (latency', perplexity') \text{ s.t. } latency' \leq latency, perplexity' < \dots\}$$

实验设计

- **实验规**
模: 评估 1,942 个候选架构，训练 170 个选定模型（每个 10B tokens）
- 硬件: NVIDIA Jetson Orin
- 基线: Qwen2.5–0.5B 等

实验结论

1. **主实验结果:**
 - 在相同延迟下，co-designed 架构达到比 Qwen2.5–0.5B 低 19.42% 的 perplexity
 - 架构选择时间从数月缩短到数天
2. **关键发现:**
 - 精度–延迟对应可以直接从 scaling law 推导

- 不同硬件平台需要不同的架构优化方向

严厉审视

- **通用性:** 仅在 Jetson Orin 上验证, 需扩展到更多硬件
 - **数据效率:** 10B tokens 训练是否足够收敛?
 - **消融:** 各超参数贡献的消融实验不足
-

论文6: World Models for Policy Refinement in StarCraft II

基本信息

- **标题:** World Models for Policy Refinement in StarCraft II
- **ArXiv ID:** 2602.14857
- **作者:** Yixin Zhang, Ziyi Wang, Yiming Rong, et al. (9 authors)
- **机构:** CASIA (中国科学院自动化研究所)
- **标签:** [Agent], [World Model], [RL], [StarCraft II]

动机

问题形式化: 现有 LLM-based SC2 agents 主要关注改进策略本身, 忽略了将可学习的动作条件转移模型集成到决策循环中。

核心贡献: 提出 StarWM, 第一个用于 SC2 的 world model, 能够在部分可观测条件下预测未来观察。

公式讲解

Structured Textual Representation: 将观察分解为五个语义模块: – 资源状态 (Resource) – 单位列表 (Units) – 地形信息 (Terrain) – 战术状态 (Tactical) – 战略状态 (Strategic)

Generate–Simulate–Refine 循环:

$$\begin{aligned} a_t &= \pi_\theta(o_t) \quad (\text{Generate}) \\ \hat{o}_{t+1} &= f_\phi(o_t, a_t) \quad (\text{Simulate}) \\ a_t^* &= \text{Refine}(a_t, \hat{o}_{t+1}, o_t) \quad (\text{Refine}) \end{aligned}$$

实验设计

- **数据集:** SC2–Dynamics–50k (首个 SC2 dynamics 预测指令微调数据集)
- **基线:** Zero-shot baselines
- **评估:** 离线 (资源预测准确率、macro–situation 一致性) + 在线 (胜率)

实验结论

1. **离线结果:**
 - 资源预测准确率提升近 60%
 - 自方 macro–situation 一致性显著提升
2. **在线结果:**
 - 胜率提升: Hard +30%, Harder +15%, VeryHard +30%
 - Macro–management 稳定性和战术风险评估改善

严厉审视

- 泛化性: 仅在 SC2 上验证, 其他 RTS 游戏?
 - Sim-to-real gap: 模拟预测误差累积如何处理?
 - 消融: world model 各组件贡献需要更详细分析
-

摘要

重点论文推荐

论文	标签	核心贡献	推荐理由
References Improve LLM Alignment	[强推], [RLVR]	用 reference 作为“软验证器”实现非可验证领域的 RL alignment	解决 RLVR 核心痛点, 方法创新性强
NeST: Neuron Selective Tuning	[Safety]	神经元选择性更新保证安全且高效	效率–效果权衡的优雅解
NESSiE	[Safety Benchmark]	必要安全基准的最小化测试	揭示 SOTA 模型的基础性弱点

今日趋势分析

1. **RL → RLVR → Reference-guided Alignment:** 论文1代表了 RLVR 向非可验证领域扩展的重要尝试, 通过引入“软验证器”概念解决了核心痛点。
2. **Safety 日益重要:** 多篇论文关注安全对齐 (论文1、2、4), 反映社区对 LLM 安全部署的持续关注。
3. **On-device / Efficiency:** 论文5代表了在资源受限环境下部署 LLM 的系统性方法论, scaling law + roofline modeling 是有趣的联合优化思路。
4. **Agent + World Model:** 论文6展示了 RL agent 的世界模型构建, 是通往 autonomous agents 的重要步骤。