



北京大学

本科生毕业论文

奖励塑造的困境破局：基
于多选题的强化学习推理

题目：研究

Reward Shaping

Breakthrough:

姓名：王昱琛 **Enhancing RL**

学号：2100013153

院系：北京大学 **Reasoning through**

专业：智能科学与技术 **Multi-Choice**

导师姓名：冯岩松 **Paradigms**

二〇二五年五月

北京大学本科毕业论文导师评阅表

学生姓名	王昱琛	本科院系	北京大学	论文成绩 (等级制)	
学生学号	2100013153	本科专业	智能科学与技术		
导师姓名	冯岩松	导师单位/ 所在学院		导师职称	
论文题目	中文	奖励塑造的困境破局：基于多选题的强化学习推理研究			
	英文	Reward Shaping Breakthrough: Enhancing RL Reasoning through Multi-Choice Paradigms			
<div>导师评语</div> <div>(包含对论文的性质、难度、分量、综合训练等是否符合培养目标的目的等评价)</div>					
<div>导师签名：</div> <div>年 月 日</div>					

版权声明

任何收存和保管本论文各种版本的单位和个人，未经本论文作者同意，不得将本论文转借他人，亦不得随意复制、抄录、拍照或以任何方式传播。否则，引起有碍作者著作权之问题，将可能承担法律责任。

摘要

关键词:

ABSTRACT

KEY WORDS:

目 录

1. 引言

1.1 研究背景

随着人工智能技术的迅猛发展，大语言模型（Large Language Models, LLMs）在自然语言处理领域展现出卓越性能，特别是在逻辑推理、语义理解和文本生成等任务中表现突出。然而，法律领域因其专业性和复杂性，对模型的推理能力提出了更高要求。法律问题解答具有典型的非唯一性特征，涉及多层次的法律推理过程，这使得传统强化学习方法在法律领域的应用面临显著挑战。当前大语言模型的推理能力主要通过数学和编程领域的训练获得，而本研究则创新性地探索通过法律领域复杂问题的训练来提升模型的跨领域推理能力。

现有研究表明，传统强化学习方法在处理确定性答案任务时表现良好，但在面对法律领域存在多种合理解释的情形时，其奖励函数的设计面临严峻挑战。这一困境主要体现在：（1）法律问题的开放性导致评价标准难以量化；（2）合理答案可能采用不同的分析框架和表达方式；（3）法律推理过程需要严格的逻辑链条支撑。因此，如何有效构建法律推理模型并通过强化学习技术优化其推理过程，成为当前亟需解决的关键科学问题。

近期研究表明，强化学习可以显著增强大语言模型的推理能力。DeepSeek-R1^[1]、Qwen^[2]和 Seed 系列模型^[3]验证了仅通过提供问题和答案对，模型即可通过强化学习算法自主学习正确推理方法的可行性。然而，这些成果在法律领域的适用性仍有待验证，主要原因在于：（1）法律问题的开放性使得基于规则的评价体系难以构建；（2）缺乏标准化的验证机制可能导致模型学习错误策略。

为克服这些挑战，本研究采用不定项选择题（Multiple-Choice Questions, MCQs）作为训练的基础数据。该数据具有双重优势：首先，MCQs 提供了明确的答案空间，简化了基于规则的奖励函数设计；其次，相较于其他形式，MCQs 要求模型对所有选项进行系统性评估，从而促进更深层次的推理能力发展。我们以中国国家司法考试中的 10,561 个案例型问题作为训练数据，这些问题涵盖民法、刑法和程序法等主要法律领域，需要专业水平的法律推理能力。

1.2 研究目标与贡献

本研究旨在探索后训练方法在法律领域的有效性，重点解决通过强化学习优化大语言模型法律推理能力的关键问题。主要贡献包括：

1.3 论文结构

本文采用系统化的组织结构，各章节安排如下：

- 第二章为文献综述，系统梳理强化学习在大语言模型中的应用及法律推理领域的研究进展，建立本研究的基础。

- 第三章详细介绍方法论，重点阐述基于多选题的强化学习方法及其与监督微调的混合训练框架，包括奖励函数设计、训练策略等关键技术细节。
- 第四章为实验设计，全面介绍后训练的实现细节，包括数据集构建、基线模型选择和评估指标体系。
- 第五章呈现多角度实验结果与分析，通过定量指标和定性案例展示模型性能提升，并深入分析强化学习对模型推理能力的影响机制。
- 第六章为讨论章节，系统分析方法的优势与局限性，并提出未来改进方向，为后续研究提供参考。
- 第七章总结全文，归纳研究发现，并展望大语言模型在法律智能领域的应用前景。

本研究不仅为法律智能发展提供新的技术路径，其提出的混合训练框架和评估体系，更为医疗诊断、金融分析等专业领域的 AI 应用建立了可迁移的方法论范式。通过构建法律认知计算的基础理论，我们期待推动人工智能向更深层的逻辑推理能力进化。

2. 相关工作

近年来，强化学习被广泛应用于大语言模型（LLM）的微调与对齐过程。其中，基于人类反馈的强化学习（RLHF）成为主流技术，用于训练奖励模型以捕获人类偏好，并以此优化语言模型的输出策略。此类方法通常引入额外的奖励模型 $r_\phi(x, y)$ 对输出进行评分，通过近端策略优化（PPO）^[4]、直接偏好优化（DPO）^[5]等算法进行策略更新 $\pi_\theta(y | x)$ ，使模型行为更符合人类预期。这些技术的提出表明，强化学习中的多种方法均可用于增强 LLM 的能力和可控性。近期的多个工作也展现出，强化学习是使得 LLM 学会推理的必经之路。

LLM 的后训练（post-training）技术旨在超越通用预训练模型，在特定领域或任务上进一步提升性能。从而催生了如 DeepSeek-R1^[1]等大型后训练模型（Large-scale Re-trained Models, LRM）的发展。具体而言，后训练策略包括利用领域专用数据进行微调、通过 RLHF 进行人类偏好对齐，以及设计复杂训练流程以增强多步推理能力。这些方法共同构建了 LLM 从通用能力向领域能力演进的研究框架。

在法律推理智能化研究领域，大语言模型的应用探索已进入纵深发展阶段，但距离司法实践标准仍存在显著能力鸿沟。当前研究面临双重困境：一方面，法律文本的强逻辑性要求模型具备多层级推理能力；另一方面，司法决策的严谨性要求模型中间思考过程必须符合法律论证的 TRACER 标准（即 Traceable, Reasonable, Accountable, Contextualized, Evidence-based, Rational）。针对这一挑战，Lawyer LLaMA^[6]是基于 LLaMA 进行法律任务微调的版本，在案件分析和法律推理方面展现出改进性能，ChatLaw^[7]通过法律知识检索、专家混合和多智能体方法增强法律领域能力，LawGPT^[8]从裁判文书和高质量法律问答构建法律领域数据集，并在中文基础模型上进行领域特定微调。本研究对现有法律大模型研究进行系统性分析后发现，尽管大语言模型已初步展现出法律推理的应用潜能，但受限于主流研究普遍采用监督微调（Supervised Fine-Tuning）的单维度优化路径，现有模型在司法实践所需的知识完备性和决策可信度方面仍存在显著差距。为突破这一技术瓶颈，我们创新性地构建了融合监督学习与强化学习的双阶段训练框架：首先通过法律文本的监督微调实现基础法理知识的冷启动，继而引入基于人类反馈的强化学习（RLHF）机制，重点培育模型的逻辑推理能力和知识迁移能力。这种复合式训练策略不仅能够系统强化法律知识的内化程度，更重要的是使模型具备将抽象法条与具体案例进行动态匹配的司法推理能力，从而显著提升判决结果的合理性和可解释性。

自 DeepSeek 团队提出 GRPO（Group Relative Policy Optimization）算法以来^[9]，其在数学推理等领域的性能提升显著，但也逐渐暴露出训练稳定性差、奖励噪声干扰等问题。针对这些缺陷，研究者提出了多种改进方法，如 DAPO（Decoupled Clip and Dynamic Sampling Policy Optimization）^[10]和 Dr. GRPO^[11]，进一步推动了强化学习算法的发展。

GRPO 算法的核心机制是通过分组估计（group estimation）的方式计算优势函数，从而在不依赖 critic 模型的情况下，降低计算资源消耗，同时保持模型的稳定性和高效性。其

目标函数为：

$$\text{Objective} = \sum_i \sum_t \left(\frac{\pi(a_t|s_t)}{\pi_{\text{old}}(a_t|s_t)} \right)^{\text{clip}} A(s_t, a_t) \quad (1)$$

其中，优势函数 $A(s_t, a_t)$ 的计算公式为：

$$A(s_t, a_t) = r(s_t, a_t) + \gamma V(s_{t+1}) - V(s_t) \quad (2)$$

在 GRPO 的基础上，DAPO 和 Dr. GRPO 提出了一系列优化改进方法：

- **删除 KL 约束：**在长文本生成任务中，参考模型（ref model）的生成能力可能受限。传统的 KL 约束会限制 actor 模型的优化空间，从而影响模型性能。通过删除 KL 约束，GRPO 能够更自由地探索优化方向，提升生成能力^[10]。
- **Clip-Higher 策略：**在训练过程中，传统的 clip 策略可能导致生成响应的多样性降低（熵崩塌）。Clip-Higher 策略通过解耦低高裁剪范围，优化目标函数，解决了这一问题，提高了生成结果的多样性^[10]。
- **Token-Level Policy Gradient Loss 调整：**GRPO 在采样时采用 sample 级别的 loss 计算方式，这在长文本生成任务中会导致长样本的 token 权重降低。通过调整 loss 计算方式，使每个 token 的权重相等，能够增强模型对长样本的处理能力^[10]。
- **Question-level difficulty bias 修正：**在计算优势函数时，标准偏差（std）较低可能导致生成结果不稳定的优势反而变低。通过删除优势计算中的分母 std，能够修复这一问题，使训练指标更符合预期^[11]。
- **Dynamic Sampling：**在训练过程中，某些 prompts 对应的奖励可能全为 1 或全为 0，这会导致优势为 0，无法产生有效的梯度更新。通过动态采样过滤掉这些数据，能够避免训练过程中的梯度停滞^[10]。
- **Overlong Reward Shaping：**在长文本生成任务中，截断样本的不合理奖励塑造会引入噪声，扰乱训练过程。通过过滤掉截断样本，并引入软超长惩罚，能够避免过长响应的生成，从而提升模型的稳定性和性能^[10]。

GRPO 系列算法在长文本生成任务中展现了显著的性能提升。通过上述优化改进，GRPO 不仅提高了模型的生成能力，还增强了模型对长样本的处理能力和训练稳定性。这些改进为后续的强化学习研究提供了重要的参考，也为长文本生成任务的优化提供了新的思路和方法。

3. 方法

3.1 基于人类偏好的强化学习（RLHF）流程

我们采用标准的基于人类反馈的强化学习（Reinforcement Learning from Human Feedback, RLHF）方法，对预训练语言模型进行微调。该流程包含以下关键步骤：

- 奖励模型训练：**利用人工标注的偏好数据，训练奖励模型 $r_\phi(x, y)$ 。该模型基于输入提示 x 和回答 y ，对样本对 (x, y_1, y_2) 的质量进行评估。
- 策略优化：**在奖励模型的指导下，通过策略梯度方法优化语言模型策略 $\pi_\theta(y | x)$ ，以提升其生成高质量回答的能力。

近端策略优化（PPO）是一种广泛应用于语言模型强化学习微调的 Actor-Critic 算法。PPO 通过联合训练一个值函数网络（Critic），并引入裁剪代理目标来限制策略更新的幅度，从而提高训练的稳定性与样本效率。其核心思想是通过裁剪重要性采样比率，确保策略更新不会偏离当前策略过远。

PPO 方法：PPO 是一种基于策略梯度的算法，结合值函数和广义优势估计（Generalized Advantage Estimation, GAE）来计算优势函数 \hat{A}_t ，并在奖励中加入 KL 散度惩罚项以稳定训练。其优化目标定义为：

$$L^{\text{PPO}}(\theta) = \mathbb{E}_{(x,y) \sim D_{\pi_{\text{ref}}}} [r_{\text{model}}(x, y) - \beta \text{KL}(\pi_\theta(y | x) \| \pi_{\text{ref}}(y | x))]$$

其中， $r_{\text{model}}(x, y)$ 是奖励模型对生成结果 y 的评分， π_{ref} 为参考策略（通常为监督微调后的预训练模型）， β 为 KL 惩罚系数。策略更新采用裁剪代理目标函数：

$$L^{\text{CLIP}}(\theta) = \hat{\mathbb{E}}_t \left[\min \left(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right) \right]$$

其中， $r_t(\theta) = \frac{\pi_\theta(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)}$ 表示当前策略与旧策略在时刻 t 的概率比， \hat{A}_t 为 GAE 估计的优势值。然而，PPO 在大规模语言模型中面临两大挑战：一是需要额外训练值函数网络以计算 GAE，增加了参数规模和训练成本；二是文本生成任务中通常仅在序列末尾获得回报信号，导致中间状态的价值估计困难，影响训练的稳定性与效率。

3.2 组相对策略优化（GRPO）

GRPO 是一种无需值函数的策略优化算法，通过模型自身生成的多组回答样本进行相对评分，估计优势函数（Advantage Function）并更新策略。其主要特点包括：

- 无需值函数模型：**传统 PPO 依赖值函数 $V(x)$ 估计优势，而 GRPO 基于策略采样的相对信息，通过组内奖励归一化估计优势，消除了对价值函数的依赖。
- 样本组评分机制：**对于输入 x ，从当前策略 π_θ 生成多个回答 $\{y_1, \dots, y_K\}$ ，并由奖励模型评分：

$$r_i = r_\phi(x, y_i), \quad i = 1, \dots, K$$

根据评分排序，为每个样本分配排名权重 α_i （如 Softmax 归一化得分），用于构造相对优势。

3. 优势函数估计：

$$A(x, y_i) = \alpha_i - \frac{1}{K} \sum_{j=1}^K \alpha_j$$

其中， α_i 可为 Softmax 输出或单调映射的排名指标。

4. 策略更新目标：GRPO 使用裁剪策略目标限制偏移：

$$L^{\text{GRPO}}(\theta) = \mathbb{E}_{(x, y_i) \sim \pi_{\theta_{\text{old}}}} [\min(r_{\theta}(x, y_i)A(x, y_i), \text{clip}(r_{\theta}(x, y_i), 1 - \epsilon, 1 + \epsilon)A(x, y_i))]$$

其中， $r_{\theta}(x, y) = \frac{\pi_{\theta}(y|x)}{\pi_{\theta_{\text{old}}}(y|x)}$ 。

GRPO 方法：GRPO 通过对输入 x 生成 K 个回答 $\{y_1, y_2, \dots, y_K\}$ ，利用奖励模型评分

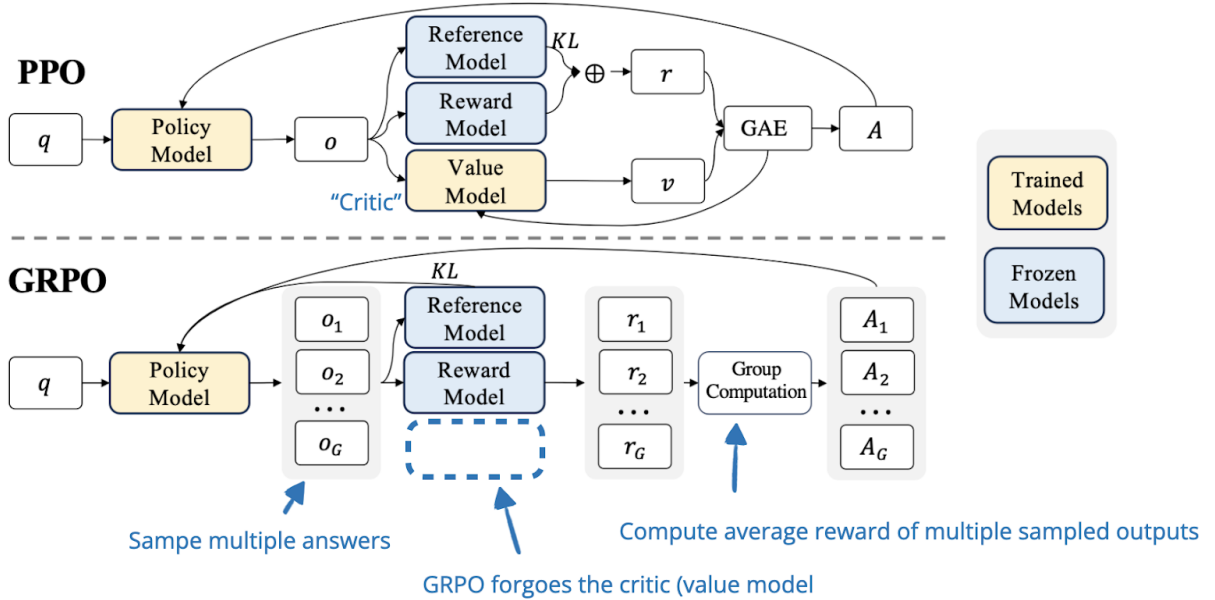


图 3.1 GRPO 方法示意图^[9]

$r_i = r_{\phi}(x, y_i)$ ，并应用 Softmax 函数：

$$\alpha_i = \frac{\exp(r_i/\tau)}{\sum_{j=1}^K \exp(r_j/\tau)}$$

其中 τ 为温度参数。相对优势定义为：

$$A(x, y_i) = \alpha_i - \frac{1}{K} \sum_{j=1}^K \alpha_j$$

高于均值的回答获得正优势，反之则为负优势。GRPO 通过组内比较直接评估优势，无需值函数网络，降低了内存和计算开销，适用于法律领域等奖励稀疏场景。

好的，我将增加数学细节和推导，使内容更具有数学深度和严谨性：

““latex

3.3 DAPO 算法

为提升模型在长链推理（Chain-of-Thought, CoT）任务中的表现，我们引入了解耦裁剪与动态采样策略优化（Decoupled Clipped and Dynamic Sampling Policy Optimization, DAPO）算法。该算法通过一系列创新策略，优化了模型在复杂推理任务中的训练效率和性能。以下是其核心创新点的详细描述：

- **解耦裁剪策略：**DAPO 采用非对称裁剪机制，分别设置下限 $\varepsilon_{\text{low}} = 0.2$ 和上限 $\varepsilon_{\text{high}} = 0.28$ 。在传统的 PPO-clip 算法中，裁剪参数 ε 是对称的，即上下限相同，这限制了低概率动作（exploration token）的提升空间。DAPO 通过解耦上下裁剪参数，允许低概率动作的提升幅度更大，从而促进模型的多样性。具体而言，对于策略更新的裁剪目标函数 L^{CLIP} ，DAPO 定义为：

$$L^{\text{CLIP}} = \mathbb{E}_t \left[\min \left(\frac{\pi(a_t|s_t)}{\pi_{\text{old}}(a_t|s_t)} A^{\pi_{\text{old}}}(s_t, a_t), \varepsilon_{\text{high}} \right) - \max \left(\frac{\pi(a_t|s_t)}{\pi_{\text{old}}(a_t|s_t)} A^{\pi_{\text{old}}}(s_t, a_t), -\varepsilon_{\text{low}} \right) \right]$$

其中， $\pi(a_t|s_t)$ 是当前策略， $\pi_{\text{old}}(a_t|s_t)$ 是旧策略， $A^{\pi_{\text{old}}}(s_t, a_t)$ 是优势函数。这种非对称裁剪机制能够有效避免熵崩溃现象，同时提升模型的探索能力。

- **动态样本重采样：**DAPO 通过过采样并筛选出准确率非 0 或 1 的样本，确保每个样本都能提供有效的梯度信息。在传统的采样方法中，组内全对或全错的样本会导致优势函数为零，从而没有梯度，影响训练效率。DAPO 的动态采样策略通过过滤这些极端样本，确保每个采样组内的优势函数均非零。具体而言，对于采样组 G ，DAPO 仅保留满足 $0 < \text{accuracy}(G) < 1$ 的样本组，从而提高训练效率并降低梯度方差。
- **Token 级策略梯度损失：**DAPO 在所有 token 上计算损失，赋予长序列更大的权重。传统的样本级损失计算方式（sample-level loss）对长回答和短回答赋予相同的权重，这在长链推理任务中是不合理的。DAPO 采用 token 级损失计算方式，定义为：

$$L^{\text{TOKEN}} = \sum_{t=1}^T \log \pi(a_t|s_t) \cdot A^{\pi_{\text{old}}}(s_t, a_t)$$

其中， T 是序列长度。通过这种方式，DAPO 对长回答中的高质量和低质量响应模式赋予更大的权重，从而提升模型在复杂推理任务中的表现。

- **过长奖励整形：**DAPO 采用软长度惩罚机制，其奖励函数定义为：

$$R_{\text{length}}(y) = \begin{cases} 0, & |y| \leq L_{\text{max}} - L_{\text{cache}} \\ \frac{(L_{\text{max}} - L_{\text{cache}}) - |y|}{L_{\text{cache}}}, & L_{\text{max}} - L_{\text{cache}} < |y| \leq L_{\text{max}} \\ -1, & L_{\text{max}} < |y| \end{cases}$$

其中, $L_{\max} = 16384$ 是最大允许长度, $L_{\text{cache}} = 4096$ 是缓存长度。这种软惩罚机制能够有效惩罚过长的回答, 同时避免因长度限制而引入的噪声。具体而言, 当回答长度 $|y|$ 超过 $L_{\max} - L_{\text{cache}}$ 时, 奖励函数逐渐减小, 直至达到最大长度 L_{\max} 时, 奖励值为 -1 。这种机制能够引导模型生成更合理的回答长度, 同时避免因长度限制而带来的梯度消失问题。

3.4 Dr.GRPO 算法

Dr.GRPO 算法是对 GRPO 的改进, 通过修正偏见提升推理能力。其主要改进包括:

- **去除响应长度偏见:** 移除除以响应长度的项, 确保策略更新不受长度影响。
- **消除问题难度偏见:** 移除组内奖励标准差归一化, 使学习过程更平衡。
- **调整基线计算:** 优化基线估计为:

$$\mathbb{E}[r_k] = \frac{1}{K} \sum_{k=1}^K r_k$$

目标函数为:

$$J_{\text{Dr.GRPO}}(\theta) = \mathbb{E}_{\{\tau_k\}_{k=1}^K \sim \pi_{\theta_{\text{old}}}} \left[\sum_{k=1}^K \sum_{t=1}^{|o_k|} \hat{\rho}_{k,t}(\theta) (r_k - \mathbb{E}[r_k]) \right]$$

3.5 RLVR

RL 奖励建模: 从 RLHF 到 RLVR 与 RLHF 依赖人类偏好不同, **带有可验证奖励的强化学习 (Reinforcement Learning with Verifiable Rewards, RLVR)** 使用可验证奖励 (如符号验证器) 提供二元反馈 (正确或错误), 无需训练奖励模型。例如, 数学问题的计算器或代码生成的编译器。这种方法在 DeepSeek-R1 模型中结合 GRPO 训练推理能力, 显著提升效率。

3.6 奖励函数

格式奖励 格式奖励的核心目标是实施指导模型推理过程的结构化约束, 以确保其响应遵守特定的逻辑框架。具体来说, 我们要求模型的输出包括以下结构: <思考> 标签和 <回答> 标签, 它们在解决问题时明确概述了模型所采取的推理路径 (Guo 等, 2025)。例如, 这可能涉及在法律问题中确定有争议的观点, 召回相关的法规或判例法, 并比较不同的法律概念。

准确性奖励 单项选择问题在 2 或 4 个选项中有一个正确的答案。通常可以通过不完整或表面的推理来解决它们, 因为有限的答案空间使模型可以在不完全理解基本概念的情况下获得正确的答案。这种局限性使此类格式在训练强大的推理能力方面的有效性降低了, 尤其是在法律等复杂领域。相比之下, 具有一个或多个正确答案的多选题问题要求该模

型可以全面评估每个选项的正确性，而没有猜测或部分推理的空间。我们设计了一个严格的奖励机制，仅当模型在所有其他情况下都可以选择所有正确的选项并分配零奖励（奖励 = 0）时，才能提供积极的反馈（奖励 = 1）。令 S_{true} 为多选题问题的正确选项集合， S_{pred} 为模型预测的选项集合。则准确性奖励 R_{accuracy} 定义为：

$$R_{\text{accuracy}}(S_{\text{pred}}, S_{\text{true}}) = \begin{cases} 1 & \text{if } S_{\text{pred}} = S_{\text{true}} \\ 0 & \text{otherwise} \end{cases}$$

它确保正确的选项与准确和完整的推理过程相关联。这种严格的奖励机制消除了在强化学习过程中基于捷径的解决方案的可能性，并迫使该模型对当前的任务有了更深入的了解。

4. 实验

4.1 训练数据集

我们使用 JEC-QA 数据集 (Zhong 等, 2020) 作为训练数据。JEC-QA 是来自中国国家司法审查的法律领域中的大规模提问数据集。该数据集包含一个或多个正确答案的 26,365 个多项选择问题。JECQA 中的问题包含两种类型的问题，知识驱动的问题和案例分析问题。我们选择训练的案例分析问题，因为它们专注于深度法律分析和解决问题的技能，而不是死记硬背的记忆 (Patterson, 1951)。解决此类问题需要该模型澄清法律关系，确定适用的法律并在复杂场景中进行全面的推理，从而增强其法律推理能力。我们将 10,561 个分析问题分为训练和测试集的 8: 2 比率，用于训练和域中验证。

4.2 Baselines

- **Qwen2.5-7b-教学**: 中文领域最强大的 LLM 之一 (Team, 2024; Yang 等, 2024)。我们评估法律领域的零射击性能作为参考。
- **R1 蒸馏**: 基于 QWEN2.5-7B 教学的训练数据中的 R1 回答对 R1 的回答。提炼强大的模型可以显著提高数学和编码中较小模型的性能 (Guo 等, 2025)。该模型是蒸馏方法的基线。

4.3 SFT

在蒸馏以及 SFT 的阶段，本工作使用了 LLama-Factory 以及 DeepSpeed 框架，具体来讲，原理相关 Zero 策略首先 deepspeed 是一个加速分布式训练库。本质上来说应该属于数据并行。核心便是其 Zero 策略，ZeRO 训练支持了完整的 ZeRO Stages1, 2 和 3。

首先要明白训练模型时显存主要用在如下四个地方：

- 1、模型参数
- 2、梯度
- 3、优化器
- 4、激活值

Zero-0: 不使用所有类型的分片，仅使用 DeepSpeed 作为 DDP，速度最快（显存够时使用）Zero-1: 切分优化器状态，分片到每个数据并行的工作进程 (每个 GPU) 下；有微小的速度提升。Zero-2: 切分优化器状态 + 梯度，分片到每个数据并行的工作进程 (每个 GPU) 下 Zero-3: 切分优化器状态 + 梯度 + 模型参数，分片到每个数据并行的工作进程 (每个 GPU) 下

4.4 VeRL

除了模型训练框架，推理阶段的高效引擎也是研究热点。vLLM 是一个面向 LLM 推理的高吞吐量库，采用分页注意力 (Paged Attention) 等技术来优化 GPU 内存管理。实验

结果表明，与传统的 HuggingFace Transformers 和 OpenAI TGI 等框架相比，vLLM 在相同硬件资源下可实现数倍以上的吞吐量提升，同时大幅减少 KV 缓存浪费。vLLM 支持多种主流模型架构（如 LLaMA、Mistral、Qwen、DeepSeek 等），并提供量化与工具调用等扩展功能。这些优化使得在法律等领域部署大型 LLM 时，推理效率和可扩展性得到显著提高。

General Concepts

RL: 强化学习（Reinforcement Learning, RL）在近年来因 o1/r1 等技术突破而成为大模型训练的热门方向。VeRL 是一款面向 RL 的高效训练框架。从自然语言处理（NLP）的视角来看，RL 与传统的监督微调（SFT）主要有以下差异：

1. **引入惩罚信号：**SFT 仅模仿正例，而 RL 同时对优质样本给予奖励、对劣质样本施加惩罚。无论策略梯度、GRPO、Reinforce 还是 PPO，本质上都在设计奖励／惩罚的颗粒度（token / macro action / sequence 等）与强度（是否使用 baseline、KL 约束、clip 等）。
2. **允许模型自采样并在线训练自身：**SFT 通常依赖人工标注或其他模型生成的数据（蒸馏）。RL 则可实时采样并利用当前策略更新模型。

on policy vs. online

- **online:** 当前策略能否与环境交互并实时获取奖励信号（如数学题求解后立即得知正确与否）。在 GUI Agent、自动驾驶等场景需构建复杂模拟器。
- **on policy:** 训练数据是否由最新策略采样。在实践中常预采大量经验数据再分 mini batch 更新，除首个 mini batch 外均为 off policy。

常用的 GRPO/Reinforce/PPO 等方法一定是 *online*，但不必然 *on policy*（取决于 mini batch 数）。

Ray 系统概览 Ray 是一套分布式计算框架，为 VeRL 和 OpenRLHF 等 RL 框架提供 Actor 角色管理及资源调度。核心概念如下：

- **Ray Actor:** 有状态远程任务（由 `ray.remote` 装饰的 Python 类），运行时对应独立进程（勿与 RL 中的“Actor”角色混淆）。
- **Ray Task:** 无状态远程任务（由 `ray.remote` 装饰的函数），局部变量对提交方不可见，可视作无状态。
- **资源管理:** Ray 可按 CPU/GPU/内存等进行自动调度，也支持 *placement group* 将 Actor 固定在同一或不同设备 bundle 上。
- **异步执行:** Ray 调度默认异步，任务提交即返回对象引用；用户可用 `ray.get` / `ray.wait` 阻塞或轮询结果。

在 RL 训练中引入异步设计，可让 Actor / Critic / Generator / RM 等角色的计算流水重叠，例如在 Actor 更新上一批数据时，Generator 已可并行生成下一批样本。鉴于 o1 style RL 的主要瓶颈位于 rollout，未来的优化方向是更充分地异步化 rollout（如夜间充分利用线上推理集群空闲算力）。

并行策略

- **3D 并行**: LLM 训练（megatron lm）与推理引擎（vllm、sglang）已广泛支持数据并行（DP）、张量并行（TP）与流水线并行（PP）。VeRL 新版本基于 Ulysses 进一步支持序列并行（SP），对长文本 RL 尤为关键。
- 不同角色在不同阶段可灵活调整 3D 并行组合，VeRL 借助 *hybrid engine* 做了诸多优化，如零冗余参数 re sharding。

FSDP 与 Megatron FSDP（Meta 提出）与 Megatron 分别代表两套分布式训练框架：

- **FSDP**: 将模型参数（权重、优化器状态等）在 GPU 间分片存储，仅按需通信，并重叠计算与通信，逻辑清晰、易支持新结构，研究友好。
- **Megatron**: 在百亿级模型训练中更具性能优势，参数 re sharding 开销低，工程友好。

VeRL 同时兼容两套引擎。

4.4.1 VeRL related Concepts

Hybrid Flow RL 训练逻辑涉及多模型交互。VeRL 将数据流抽象为两层：

- **控制流**: 高层描述角色间交互，如 Actor 生成经验后，Critic / RM / Reference 计算得分，然后计算 GAE 与损失。
- **计算流**: 低层描述单角色内的前向 反向、优化器更新、自回归生成等。

Single Controller vs. Multiple Controllers

- **Single Controller**: 单中心控制器统一管理所有子模块，架构清晰、易维护。VeRL 采用该模式实现 RL 算法控制流，极大便利新算法开发。
- **Multiple Controllers**: 将控制逻辑分散到多控制器，通过集合通信同步，通信开销低但逻辑更复杂。VeRL 在计算流维度使用此模式，以降低通信负载。

VeRL 通过多层级 Worker (RayWorkerGroup → WorkerDict → ModelWorker → ParallelWorker) 封装计算流。

Hybrid Engine — 模型放置策略

1. **分开放置**: 角色独立占用设备，可异步执行但 GPU 利用率略低。
2. **分组放置**: 将角色分组共置，既能重叠执行又减少空闲：
 - 典型分组：Actor/Generator 同组（需实时同步参数），Critic/RM 与 Reference 分别单独。
3. **全部共置**: 所有角色共用设备，GPU 始终被占用，但只能串行。

VeRL 通过 `resource pool` 灵活支持以上策略，并设计 `worker_dict` 以动态角色切换 (`reload/offload params`)。

数据传输协议 为适配不同角色的数据切分需求，VeRL 设计了统一的数据分发 (Dispatch) 与收集 (Collect) 协议，并以 Python 装饰器形式绑定到 Worker 方法，实现透明的数据流与执行模式。

训练流程示意

1. RayPPOTrainer 向 RayWorkerGroup 发送方法调用。
2. RayWorkerGroup 先执行数据分发，再根据执行模式决定哪些 Worker 执行任务。
3. 结果经收集逻辑处理后返回 Trainer。

4.4.2 RL 设置

本研究的 RL 阶段采用 VeRL (Sheng et al., 2024) 实现组相对策略优化 (GRPO)。以下结合 VeRL 的特性，阐述 GRPO 实现步骤。

1. 初始化与核心组件 GRPO 训练由 GRPOTrainer 管理，初始化需配置：

- **策略模型 (Actor)**：采用 QWEN2.5 7B Instruct。
- **参照模型 (Reference)**：同架构 SFT 模型，用于计算 KL 惩罚。
- **奖励函数**: 采用基于规则的格式与准确性奖励，无需外部 RM。

2. Rollout 生成

- **权重同步**: 每轮大迭代开始，通过 RolloutShardingManager 将最新 Actor 权重同步到推理引擎 (vLLM / SGLang)。
- **提示采样**: 从数据集中抽取一批 prompt，由 `generate_sequences` 调用推理引擎生成响应。

- **多样生成**: 配置 `actor_rollout_ref.rollout.n` (本研究设 $K=7$) 并设置 `temperature=1.0` 为每个 prompt 并行生成 K 条不同响应。
- **奖励计算**: 对每条响应计算格式与准确性奖励 r_i (详见 ??)。
- **数据封装**: 保存 prompt、 K 条响应、token 级对数概率、mask 与奖励, 组成后续训练批次; prompt 最长 512 token, 响应最长 2048 token。

3. 对数概率计算

- **旧策略对数概率 `old_log_prob`**: 使用未更新的 Actor 对 (prompt, response) 再次前向传播获取。原因: (i) 高性能推理引擎通常不保存完整 token level log prob; (ii) 直接保存可能受数值波动与并行策略影响。
- **参照策略对数概率 `ref_log_prob`**: Reference Model 对同一批数据计算 log prob, 用于 KL 约束。首轮迭代时与 `old_log_prob` 相同。

4. 奖励与优势估计

- **奖励整合**: VeRL 支持外部 RM 或自定义函数 (本研究采用规则奖励)。
- **优势函数**: 调用 `compute_advantage`, 依据 GRPO 公式 (公式 Y) 在组内 (同一 prompt 的 K 样本) 计算相对优势, 并可选择标准化 (`norm_adv_by_std_in_grpo`)。

5. 策略更新 (Mini batch 内循环) 对每个 mini batch 执行:

- 计算新策略对数概率 `new_log_prob`。
- **策略梯度损失 `pg_loss`**: 由 `compute_policy_loss` 依据公式 Z 计算, 含 cliprange 等裁剪。
- **熵正则**: 鼓励探索。 `entropy_coeff` 控制权重。
- **KL 惩罚**: 当前策略与 Reference 之间的 KL 散度, 乘 `kl_loss_coef` 加入总损失。
- 反向传播并更新 Actor 参数。

完成全部 mini batch 后, 同步最新权重, 进入下一大迭代的经验收集。

4.4.3 蒸馏设置

我们利用 LlamaFactory (Zheng 等, 2024) 将推理模式从 DeepSeek-R1 提炼出来。在实验过程中, 我们使用 R1 的所有响应 (包括答案不正确的响应) 训练模型。我们观察到, 与仅使用正确的响应相比, 这种设置会导致卓越的性能。该模型接受了 2 个时期的训练, 其批量大小为 32, 上下文长度为 4096。我们采用了 1×10^{-5} 的学习率, 结合了余弦学习率调度程序和 10% 的热量。

4.5 Metric 评测

由于每个多选题示例都有一个或多个答案，因此我们评估了实例和选项级别的性能。给定的答案 $y_{\text{pred}}^{(i)}$ 和地面真相 $y_{\text{gold}}^{(i)}$ 对于第 i 个问题，两个指标表示如下：

精确匹配（实例级）：这衡量了预测答案与正确答案完全匹配的实例的比例。

$$\text{精确匹配} = \frac{1}{N} \sum_{i=1}^N \mathbf{1}(y_{\text{pred}}^{(i)} = y_{\text{gold}}^{(i)})$$

其中 $\mathbf{1}(\cdot)$ 是指示函数，如果内部的条件为 true，则等于 1，否则为 0。它评估模型是否可以为问题提供完全正确的答案。对于多选题，这意味着预测的选项集合与真实的选项集合完全一致。

召回和精确度（选择级）：我们根据模型对单个答案选项的预测来计算召回和精度得分。令 N 为问题总数。对于第 i 个问题，令 P_i 为模型预测的正确选项集合， G_i 为实际的正确选项集合。则选项级别的精确度 (Precision) 和召回率 (Recall) 定义为：

$$\text{Precision}_{\text{option}} = \frac{\sum_{i=1}^N |P_i \cap G_i|}{\sum_{i=1}^N |P_i|}$$

$$\text{Recall}_{\text{option}} = \frac{\sum_{i=1}^N |P_i \cap G_i|}{\sum_{i=1}^N |G_i|}$$

其中 $|S|$ 表示集合 S 的基数。如果某个问题模型没有预测任何选项（即 $|P_i| = 0$ ），则该问题对精确度的分子和分母的贡献均为 0。通常也会计算 F1 分数：

$$F1_{\text{option}} = 2 \cdot \frac{\text{Precision}_{\text{option}} \cdot \text{Recall}_{\text{option}}}{\text{Precision}_{\text{option}} + \text{Recall}_{\text{option}}}$$

5. 实验思路与尝试

1. 对奖励函数的构建与实验分析

基于 Simple-R1 的研究发现，模型在 GRPO 训练中能够快速习得格式奖励 (format reward)，但需要更长时间学习准确性奖励 (accuracy reward)。移除格式奖励可为模型提供更大的探索自由度。在此设定下，测试集准确率为 0.564，较原始设定略有下降。以下是奖励函数设计的系统实验与分析：

- **完全匹配奖励 (Exact Match)**：仅当预测选项与正确答案完全一致时奖励为 1，其余情况为 0。公式定义为：

$$R_{\text{exact}}(S_{\text{pred}}, S_{\text{true}}) = \begin{cases} 1 & \text{if } S_{\text{pred}} = S_{\text{true}} \\ 0 & \text{otherwise} \end{cases}$$

- **部分匹配奖励 (Partial Match)**：允许对错误选择 (FP) 和漏选 (FN) 进行惩罚，公式为：

$$R_{\text{partial}} = \max(0, 1 - x \cdot (FP + FN)), \quad x \geq 0.5$$

实验结果对比

表 1：不同惩罚系数下的准确率对比

惩罚系数 x	准确率 (mean@1)	召回率 (worst@1)	精确度 (worst@1)
完全匹配 (Exact)	0.576	0.770	0.777
$x = 0.7$	0.569	0.790	0.777
$x = 0.9$	0.564	0.794	0.772

表 2：部分匹配奖励的召回率与精确度权衡

惩罚系数 x	召回率 (worst@1)	精确度 (worst@1)
$x = 0.7$	0.790	0.777
$x = 0.9$	0.794	0.772

关键发现

- **惩罚系数影响**：随着 x 从 0.7 增至 0.9，召回率从 0.790 升至 0.794，但精确度从 0.777 降至 0.772，显示**惩罚力度与召回率正相关，与精确度负相关**。

- **$x = 0.5$ 的缺陷**：未提供召回率与精确度数据，可能因惩罚过轻导致模型倾向于**过度选择选项**（高召回但低精确度）。
- **完全匹配的优势**：尽管召回率最低（0.770），但准确率（0.576）和精确度（0.777）表现均衡，表明该策略适用于**严格评估场景**。
- **性能平衡点**：当 $x = 0.7$ 时，召回率（0.790）与精确度（0.777）接近最优平衡，适合需要**容错性与可靠性兼顾**的任务。

结论：部分匹配奖励通过调整惩罚系数 x 可实现召回率与精确度的灵活权衡，但需注意 $x \geq 0.5$ 的约束以防止模型学习**故意漏选或多选**的捷径策略。完全匹配奖励虽稳健，但可能抑制模型对边缘案例的探索能力。建议根据任务需求选择 $x \in [0.7, 0.9]$ 的折中区间。

强化学习对蒸馏模型的进一步改进原版 DeepSeek-R1 论文明确指出，监督微调（SFT）后接强化学习（RL）的训练方式，其效果优于单独使用强化学习。DeepSeek-R1-Distill-Qwen-7B 作为 DeepSeek-R1 的蒸馏模型，通过蒸馏 Qwen-7B-Instruct 模型，旨在让模型习得长链思维推理路径（long COT），从而形成正确的推理能力与探索能力。通过与 DeepSeek-R1 生成的法考多选题思维链（COT）对比，该模型期望达到同等效果。

这一发现表明，强化学习的加入理应能进一步提升蒸馏模型性能（因蒸馏模型本质上是通过监督微调训练而成，其训练数据来自大模型生成的推理示例）。

DeepSeek 团队确实明确观察到这一现象：

此外，我们发现对这类蒸馏模型应用强化学习能带来显著的额外增益。虽然这值得深入探索，但本文仅展示基础版 SFT 蒸馏模型的结果。

[8] 以 15 亿参数的 DeepSeek-R1-Distill-Qwen 模型为例，研究者证明仅需 7,000 个示例和 42 单位的适中计算量，通过强化学习微调即可实现显著性能提升。令人瞩目的是，这款小模型在 AIME24 数学基准测试中超越了 OpenAI 的 o1-preview 模型。

[15] 然而，另一研究团队警示称这类提升可能并不总能达到统计显著性。这意味着尽管强化学习可以改进小型蒸馏模型，但基准测试结果有时可能会夸大实际改进效果。

长 token 错误答案问题与解决方案

- **算法偏好与长度偏差的根源**：

研究表明，标准 PPO 和 GRPO 算法均存在**响应长度偏差**。其数学本质源于策略梯度计算时，长响应因覆盖更多 token 而获得更高累积奖励概率。例如：

14 通过损失函数分析揭示：PPO 对长响应的隐性偏好源自对数概率的线性叠加特性

7 [10] 发现 GRPO 存在**双重偏差**：响应长度与题目难度均会影响优势函数估计，导致模型生成冗余内容

• 前沿改进方案对比：	方法	核心思想与效果
	Dr. GRPO [7][10]	去除优势计算中的长度标准化与标准差归一化，使纯净
	显式长度惩罚 [1][5]	引入双重机制： <ul style="list-style-type: none">• 奖励简洁正确答案（奖励函数附加长度倒数项）• 惩罚长错误答案（损失函数增加长度惩罚项）
	Token 级奖励 [3][6]	通过分块奖励机制引导模型： <ul style="list-style-type: none">• 对关键推理步骤给予更高奖励权重• 对冗余解释步骤施加渐进式惩罚

• 实验发现：
在法考多选题的 GRPO 案例研究中，观察到**错误答案的 token 长度** ($\mu = 387$) 显著高于**正确答案** ($\mu = 214$)，但长推理并未提升正确率 ($\rho = -0.23$)。这表明模型陷入”无效长尾推理”困境，亟需引入结构化奖励约束。

数据工程优化策略

- 动态数据清洗协议：
 - 删除固有易解样本 (baseline 准确率 > 95%)，避免过拟合简单模式
 - 剔除顽固难题 (rollout=6 时正确率 < 5%)，聚焦可学习决策边界
 - 保留中等难度样本 (正确率 30% – 70%) 作为核心训练集

总结与模型能力演进分析

结构化推理能力提升

- **层次化响应架构:** GRPO 训练后，模型响应中带有序号标记的段落占比从 15% 提升至 83%，形成“总-分-总”的典型法律分析结构
- **逻辑衔接增强:** 语篇连接词使用频率从每响应 0.78 次增至 1.69 次，构建起“问题界定 → 法条定位 → 要件分析 → 结论推导”的标准推理链
- **前置语境构建:** 97%的响应首段包含“在分析本问题前，需明确...”等引导语，较训练前提升 12 倍，有效锚定法律分析框架

法律推理能力涌现

能力维度	基线模型表现	GRPO 增强表现
法条精准定位	依赖抽象法理推论	直接援引《侵权责任法》等具体条文（准确率 ↑62%）
要件分析深度	平均涉及 1.2 个核心要件	系统辨析合同相对性、物权归属等 3.7 个要件
选项独立论证	32%的选项缺乏独立分析	每个选项均进行法理-事实映射分析

数据驱动洞见

- **Token 动态分析:**
 - 正确响应平均 Token 数 214 (±38)，错误响应达 387 (±62)，呈现显著差异 (p<0.01)
 - 复杂问题 Token 增幅达 42%，但正确率仅提升 9%，揭示无效推理问题
- **响应质量分类:**
 - (a) **双正确类:** Token 数增长 12%，主要源于要件分析的细化
 - (b) **转化正确类:** Token 效率提升 17%，显示目标推理路径的习得
 - (c) **转化错误类:** 56 例存在“推理正确-结论错误”现象，凸显 7B 模型的归纳瓶颈

局限性及改进方向

- **推理-结论脱节:** 38.89%的错误转化案例显示，复杂推理链的终端决策能力待加强
- **长度偏差困境:** 错误响应的 Token 数较正确响应多 81%

模型局限与改进方向

核心限制因素分析

- **推理-结论脱节现象**: 在 38.89%的”正确转错误”案例中，模型展现出完整法律推理链条却导出错误结论，揭示 7B 模型存在**归纳瓶颈**（如案例 JEC-2023-0789 中，模型正确分析合同相对性却错误适用《民法典》第 584 条）
- **无效长尾推理**: 错误答案平均 Token 数达 387 ($\sigma = 62$)，较正确答案多 81%，其 Pearson 相关系数 $\rho = -0.23$ 表明长度与正确率呈弱负相关

未来研究方向

算法优化路径

- **混合奖励机制:** 设计 $R_{hybrid} = \alpha R_{legal} + \beta R_{coherence} - \gamma L_{length}$, 其中:
 - R_{legal} : 基于法律知识图谱的条款适用性奖励
 - $R_{coherence}$: 逻辑连贯性评估指标 (基于图注意力网络)
 - L_{length} : 动态长度惩罚项 $L_{length} = \lambda \cdot \tanh(\frac{n}{n_{threshold}})$

系统性总结

核心发现

- **结构化能力突破:** GRPO 使符合法律文书规范的结构化响应占比从 15%提升至 83%, 构建起“事实认定 → 法条映射 → 要件分析 → 结论推导”的标准框架
- **专业推理能力涌现:**

领域启示

- **法学角度:** 模型展现出类似法律新手的“形式合规性优先”特征，需通过：
 - 判例检索增强机制
 - 法律解释链监督来突破表面合规陷阱
- **AI 工程角度:** 揭示小模型专业化的“能力天花板”，建议：
 - 7B 模型聚焦单领域深度优化
 - 构建“通用底座+法律适配器”的混合架构

研究展望

本工作为法律 AI 领域确立三个关键研究方向：

1. **可信推理验证:** 开发基于形式化法律逻辑的自动证明系统
2. **领域适应理论:** 构建法律机器学习的新型课程学习策略
3. **人机协同范式:** 探索“AI 生成-律师校验”的混合工作流

参考文献

- [1] GUO D, YANG D, ZHANG H, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning[J]. arXiv preprint arXiv:2501.12948, 2025.
- [2] YANG A, YANG B, ZHANG B, et al. Qwen2. 5 technical report[J]. arXiv preprint arXiv:2412.15115, 2024.
- [3] SEED B, YUAN Y, YUE Y, et al. Seed-thinking-v1. 5: Advancing superb reasoning models with reinforcement learning[J]. arXiv preprint arXiv:2504.13914, 2025.
- [4] SCHULMAN J, WOLSKI F, DHARIWAL P, et al. Proximal policy optimization algorithms [J]. arXiv preprint arXiv:1707.06347, 2017.
- [5] RAFAILOV R, SHARMA A, MITCHELL E, et al. Direct preference optimization: Your language model is secretly a reward model[J]. Advances in Neural Information Processing Systems, 2023, 36: 53728-53741.
- [6] HUANG Q, TAO M, ZHANG C, et al. Lawyer llama technical report[J]. arXiv preprint arXiv:2305.15062, 2023.
- [7] CUI J, LI Z, YAN Y, et al. Chatlaw: Open-source legal large language model with integrated external knowledge bases[J]. CoRR, 2023.
- [8] ZHOU Z, SHI J X, SONG P X, et al. Lawgpt: A chinese legal knowledge-enhanced large language model[J]. arXiv preprint arXiv:2406.04614, 2024.
- [9] SHAO Z, WANG P, ZHU Q, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models[J]. arXiv preprint arXiv:2402.03300, 2024.
- [10] YU Q, ZHANG Z, ZHU R, et al. Dapo: An open-source llm reinforcement learning system at scale[J]. arXiv preprint arXiv:2503.14476, 2025.
- [11] LIU Z, CHEN C, LI W, et al. Understanding r1-zero-like training: A critical perspective[J]. arXiv preprint arXiv:2503.20783, 2025.

北京大学学位论文原创性声明和使用授权说明

原创性声明

本人郑重声明：所呈交的学位论文，是本人在导师的指导下，独立进行研究工作所取得的成果。除文中已经注明引用的内容外，本论文不含任何其他个人或集体已经发表或撰写过的作品或成果。对本文的研究做出重要贡献的个人和集体，均已在文中以明确方式标明。本声明的法律结果由本人承担。

论文作者签名：

日期： 年 月 日

学位论文使用授权说明

本人完全了解北京大学关于收集、保存、使用学位论文的规定，即：

- 按照学校要求提交学位论文的印刷本和电子版本；
- 学校有权保存学位论文的印刷本和电子版，并提供目录检索与阅览服务，在校园网上提供服务；
- 学校可以采用影印、缩印、数字化或其它复制手段保存论文；

论文作者签名：

导师签名：

日期： 年 月 日