

基于人类反馈的强化学习

面向语言模型的后训练RLHF简明指南

作者：Nathan Lambert 译者：Junwei He & GPT4.1

16 April 2025

Abstract

基于人类反馈的强化学习（RLHF）已成为部署现代机器学习系统的重要技术手段与工具。本书旨在为具有定量分析背景的读者提供核心技术方法的入门指南。内容溯源RLHF的技术根源——既涵盖近期文献进展，也涉及经济学、哲学和最优控制等跨学科领域的理论融合。通过明确定义、问题建模、数据收集和通用数学工具，构建完整知识框架。核心章节详细解析RLHF全流程优化阶段，涵盖指令调优、奖励模型训练，以及拒绝采样、强化学习和直接对齐算法等关键技术。最后探讨合成数据与评估等前沿研究课题，以及该领域的开放性科学问题。本项目地址 <https://github.com/jweihe>

Contents

1 引言	6
1.1 RLHF究竟解决了什么？	7
1.2 关于后训练的直观理解	8
1.3 我们是如何走到今天的	9
1.4 本书范围	10
1.4.1 章节概览	11
1.4.2 目标读者	12
1.4.3 如何使用本书	12
1.4.4 关于作者	12
1.5 RLHF的未来	12
2 关键相关工作	13
2.1 起源至2018年：基于偏好的强化学习	13
2.2 2019至2022年：人类偏好驱动的语言模型强化学习	13
2.3 2023年至今：ChatGPT时代	14
3 定义与背景	15
3.1 语言建模概述	15
3.2 机器学习相关定义	15
3.3 自然语言处理相关定义	16
3.4 强化学习相关定义	16
3.5 RLHF专有定义	17
3.6 拓展术语表	17
4 训练概览	19
4.1 问题形式化	19
4.2 标准RL设置的调整	19
4.3 优化工具	20
4.4 RLHF典型流程示例	21
4.5 微调与正则化	21
5 偏好的本质	23
5.1 优化偏好的路径	24
5.1.1 偏好的量化	24
5.1.2 关于偏好的可能性	25
6 偏好数据	26
6.1 为什么需要偏好数据	26
6.2 如何收集偏好数据	26
6.2.1 标注界面	26

6.2.2 排序 vs. 评分	31
6.2.3 多轮数据	31
6.2.4 结构化偏好数据	33
6.2.5 数据采购与合同	34
6.3 模型中真的表达了这些偏好吗?	35
7 奖励建模	36
7.1 奖励模型的训练	36
7.2 模型结构	37
7.3 实现示例	37
7.4 变体	37
7.4.1 偏好间隔损失 (Preference Margin Loss)	37
7.4.2 单prompt多比较的平衡	38
7.4.3 K-wise损失函数	38
7.5 结果奖励模型 (Outcome Reward Models, ORM)	38
7.6 过程奖励模型 (Process Reward Models, PRM)	39
7.7 奖励模型 vs. 结果RM vs. 过程RM vs. 价值函数	40
7.8 生成式奖励建模	40
7.9 延伸阅读	41
8 正则化	43
8.1 RL优化中的KL距离	43
8.1.1 参考模型与生成样本	43
8.1.2 实现示例	43
8.2 预训练梯度	44
8.3 其他正则化方法	45
9 指令微调 (Instruction Finetuning)	46
9.1 聊天模板与指令结构	46
9.2 指令微调的最佳实践	48
10 拒绝采样 (Rejection Sampling)	49
10.1 训练流程	49
10.1.1 生成补全	49
10.1.2 选择Top-N补全	50
10.1.3 微调	52
10.1.4 细节说明	52
10.2 相关: Best-of-N采样	53
11 策略梯度算法 (Policy Gradient Algorithms)	54
11.1 策略梯度算法基础	54

11.1.1 Vanilla Policy Gradient	55
11.1.2 REINFORCE	56
11.1.3 Proximal Policy Optimization (PPO)	57
11.1.4 Group Relative Policy Optimization (GRPO)	58
11.2 实现细节	59
11.2.1 策略梯度基础实现	60
11.2.2 损失聚合	60
11.2.3 PPO实现示例	62
11.2.4 GRPO实现示例	64
11.3 补充话题	66
11.3.1 广义优势估计 (GAE)	66
12 直接对齐算法 (Direct Alignment Algorithms, DAAs)	67
12.1 直接偏好优化 (Direct Preference Optimization, DPO)	67
12.1.1 DPO的原理	67
12.1.2 DPO公式推导	69
12.2 数值问题、局限与变体	71
12.3 实现注意事项	72
12.4 DAA与RL: 在线与离线数据	72
13 宪法AI与AI反馈	74
13.1 宪法AI (Constitutional AI, CAI)	74
13.2 专用LLM裁判模型	75
13.3 延伸阅读	75
14 推理训练与推理时扩展	76
14.1 为什么RL现在能“起飞”？	78
14.2 RL训练 vs. 推理时扩展	78
14.3 强化微调的未来 (超越推理)	79
15 合成数据与蒸馏	80
16 评测	82
16.1 提示格式化: 从Few-shot到Zero-shot再到CoT	82
16.2 评测的使用与观察	85
16.3 数据污染 (Contamination)	87
16.4 工具与平台	88
17 过度优化 (Over Optimization)	89
17.1 定性过度优化	89
17.1.1 代理目标的管理	89
17.1.2 过度拒绝与“RLHF过猛”	91

17.2 定量过度优化	92
17.3 失调与RLHF的角色	93
18 风格与信息	94
18.1 “话痨悖论”	94
18.1.1 “话痨”现象的成因	95
19 产品、用户体验与模型个性	96
19.1 角色训练 (Character Training)	96
19.2 模型规范 (Model Specifications)	97
19.3 产品周期、用户体验与RLHF	97
参考文献	98

1 引言

人类反馈强化学习（Reinforcement Learning from Human Feedback, RLHF）是一种将人类知识融入人工智能系统的技术手段。

RLHF最初被提出，主要是为了解决那些难以精确定义的复杂问题。

其早期应用多见于控制任务以及其他强化学习（RL）的传统领域。

随着ChatGPT的发布，以及大语言模型（LLMs）和其他基础模型的迅猛发展，RLHF逐渐成为业界关注的焦点。

RLHF的基本流程包含三个核心步骤：

首先，需要训练一个能够理解并响应用户问题的语言模型（详见第9章）；

其次，收集人类偏好数据，用于训练反映人类偏好的奖励模型（详见第7章）；

最后，利用强化学习优化器，对语言模型进行优化，通过采样生成内容并根据奖励模型进行评分（详见第3章和第11章）。

本书将详细介绍这一流程中各步骤的关键决策与基础实现案例。

RLHF已被成功应用于多个领域，随着技术的成熟，其复杂性也在不断提升。

早期RLHF的突破性实验包括：深度强化学习 [1]、文本摘要 [2]、指令跟随 [3]、网页信息解析问答 [4] 以及“对齐”（alignment）[5] 等。

下图总结了早期RLHF的基本流程（见@fig:rlhf-basic）。

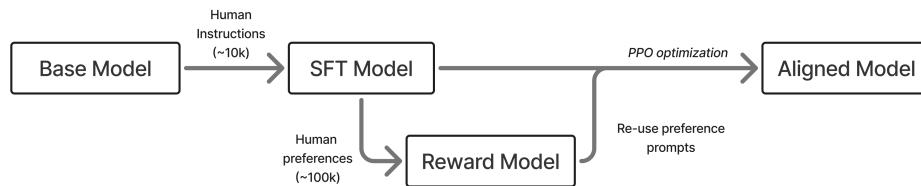


Figure 1: 早期RLHF三阶段流程示意图：SFT、奖励模型、优化。

在现代语言模型的训练体系中，RLHF已成为“后训练”（post-training）阶段的重要组成部分。

后训练是一套更为完整的技术体系和最佳实践，旨在让语言模型更好地服务于下游任务 [6]。

后训练大致可以归纳为三类优化方法：

1. 指令/监督微调（IFT/SFT）：教会模型格式化和基本指令跟随能力，主要关注语言的特征；
2. 偏好微调（PreFT）：使模型更贴合人类偏好，同时带来一定能力提升，主要关注语言的风格和那些难以量化的人类细微偏好；
3. 强化微调（RFT）：最新的后训练方法，能够提升模型在可验证领域的表现。

本书重点关注第二类——偏好微调，它比指令微调更复杂，也比强化微调更为成熟和普及。

值得一提的是，RLHF在业界常常被用作现代后训练的代名词。

ChatGPT发布后，RLHF几乎涵盖了所有后训练流程。

RLHF的基础远不止偏好数据，本书也会介绍所有相关主题。

1.1 RLHF究竟解决了什么？

围绕RLHF最常被问到、但至今仍难以完全回答的问题是：“RLHF训练到底为模型带来了什么？”

本书除了介绍RLHF的技术实现外，更希望帮助读者建立直观理解，弄清楚为什么RLHF对现代AI模型如此关键。

近年来，语言模型已从学术研究的基准测试工具，转变为通用技术平台，而RLHF正是这一转变的核心动力。

举个例子，假如让一个经过自回归预训练（autoregressive pretraining）的语言模型（如Llama 3.1 405B）回答问题“The president of the united states in 2006 was”，它可能会生成如下内容：

George W. Bush, the governor of Florida in 2006 was Jeb Bush, and John McCain was an Arizona senator in 2006 - who later lost to obama.

September 1 – U.S. President Bush signs an executive order to provide more options for faith-based organizations when delivering social services that are paid for with federal funds.

January 1 – The current deadline set by United States Congress in the Unlawful Internet Gambling Enforcement Act, aka UIGEA.

search: Amendments to the 1961 International Convention for the Protection of New Varieties of Plants require plant breeders' rights include farmer's privilege.

2009 was a common year starting on Thursday of the Gregorian calendar, the 2009th year of the Common Era (ce) and Anno Domini (ad) designations, the 9th year of the 3rd millennium, the 9th year of the 21st century, and the 10th and last year of the 2000s decade.

可以看出，这样的模型只是不断补全句子，并夹杂了各种网络常见的元数据。

理解RLHF最有说服力的方式，是关注“风格”（style）在与语言模型交互中的作用。

信息的呈现风格或格式，对学习效果至关重要。

这在课程学习等场景中一直如此，只是我们往往习以为常，并未直接关注。

继续上面的例子，如果用RLHF（以及其他后训练工具）训练后的模型Tülu 3 405B来回答同样的问题，它会更简洁明了地给出答案：

George W. Bush was the president of the United States in 2006. He served two terms in office, from January 20, 2001, to January 20, 2009.

现代研究已经证明，RLHF是一种通用方法，能够让模型学会微妙的风格和行为特征。

与其他后训练技术（如指令微调）相比，RLHF在跨领域泛化能力上表现更佳 [7] [8]，有助于打造高效的通用模型。

从直观上看，这种差异体现在优化方式上。

指令微调本质上是让模型在遇到类似已见示例时，预测下一个确定的token，从而更频繁地输出特定文本特征（属于逐token更新）。

而RLHF则是在整体回答层面进行优化，不仅教会模型什么是“更好”的回答，还让模型知道哪些类型的回复应当避免（即负反馈）。

这种训练通常被称为对比损失（contrastive loss），本书后续会多次提及。

虽然这种灵活性是RLHF的一大优势，但也带来了实现上的挑战——主要是在于如何控制优化过程。

如本书后续将介绍，实施RLHF通常需要训练奖励模型，而奖励模型的最佳实践尚未完全确定，且依赖于具体应用场景。

此外，由于奖励信号往往只是目标的代理，优化过程容易产生过度优化，因此需要正则化手段。

正因如此，高效的RLHF需要一个良好的起点，因此RLHF并不能单独解决所有问题，必须结合更广泛的后训练视角来使用。

由于其复杂性，RLHF的实现成本远高于简单的指令微调，并且可能遇到一些意想不到的问题，例如长度偏差 [9] [10]。

对于对性能有较高要求的项目来说，RLHF被认为是获得高质量微调模型的关键，但也意味着更高的算力、数据和时间成本。

1.2 关于后训练的直观理解

这里用一个简单的类比来说明：为什么在同一个基础模型上，通过后训练能获得如此巨大的提升。

我用来理解后训练潜力的直觉是“后训练的引出解释”（elicitation interpretation of post-training）：我们实际上是在挖掘和放大基础模型中已有的有价值行为。

以F1方程式赛车为例，大多数车队在赛季初都有全新底盘和发动机，然后花一整年时间不断优化空气动力学和系统，赛车性能会有巨大提升。

顶级车队在一个赛季内的进步，远超底盘本身的变化。

后训练也是如此。最优秀的后训练团队能在很短时间内挖掘出模型大量潜力。这套技术体系涵盖了预训练结束后的所有环节，包括“中期训练”（如退火、高质量web数据）、指令微调、RLVR、偏好微调等。例如我们将OLMoE Instruct第一版提升到第二版，后训练评测均分从35提升到48，绝大部分预训练内容未做改动 [11]。

同理，像GPT-4.5这样的模型，也为OpenAI提供了更具活力和潜力的基础平台。我们也知道，规模更大的基础模型能吸收和适应更多样化的变化。

这说明，扩大模型规模同样能让后训练进展更快。当然，这也需要强大的训练基础设施，这也是为何各大公司仍在建设超大算力集群。

这个理论也反映了现实：用户所感受到的主要提升，实际上都来自后训练。这意味着，通过互联网预训练的模型中，蕴藏着远超我们通过早期后训练（如仅靠指令微调）所能简单挖掘的潜力。

这个理论的另一个名字是“表层对齐假说”（Superficial Alignment Hypothesis），最早见于论文 LIMA: Less is More for Alignment [12]。该论文提出：

模型的知识和能力几乎全部在预训练阶段获得，而对齐（alignment）则教会模型在与用户交互时应采用哪种子分布的格式。如果该假说成立，且对齐主要是风格学习，那么只需少量样本就能充分微调一个预训练语言模型 [Kirstain et al., 2021]。

深度学习的成功经验应让你坚信：数据规模对性能至关重要。

但该论文作者主要讨论的是对齐和风格，这是当时学术界后训练的关注点。用几千条指令微调样本，确实可以让模型在如AlpacaEval、MT Bench、ChatBotArena等评测中表现大幅提升，但这些提升并不总能迁移到更具挑战性的能力上，这也是Meta不会只用这类数据集训练Llama Chat模型的原因。学术结论有其价值，但若想把握技术演进全貌，还需谨慎解读。

该论文实际证明的是：少量样本确实能让模型发生较大变化，这对于新模型的短期适应很重要，但他们对性能的论述可能会误导普通读者。

如果我们改变数据内容，模型性能和行为的提升会远超“表层”改善。如今的基础语言模型（未经过后训练）可以通过强化学习解决数学问题，学会输出完整的思维链条推理，并在BigBenchHard、Zebra Logic、AIME等推理评测中获得更高分数。

“表层对齐假说”之所以不成立，原因和认为RLHF及后训练只是“调调风格”的观点一样错误。

整个领域在2023年都经历了这一认知转变（虽然很多AI观察者仍未跟上）。

后训练的作用早已超越风格调整，我们也逐渐认识到，模型的风格是在行为之上进一步塑造出来的，比如现在流行的长链式思维。

1.3 我们是如何走到今天的

为什么现在写这本书？未来还会发生哪些变化？

自ChatGPT发布、RLHF重新受到关注以来，后训练（即从原始预训练语言模型中激发强大行为的技术）经历了多个阶段和风潮。

在Alpaca [13]、Vicuna [14]、Koala [15]、Dolly [16]等模型时代，研究者用有限的人类数据点结合自生成数据（Self-Instruct风格），对原始LLaMA进行微调，实现了类似ChatGPT的行为。

这些早期模型的评测标准主要靠“体验”和人工评估，因为大家都被小模型展现出的跨领域惊艳能力所吸引，这种兴奋是合理的。

开源后训练发展更快，模型发布更频繁，影响力甚至超过了闭源同行。

各大公司纷纷调整策略（如DeepMind与Google合并等），以应对这一变化。
开源方案有时领先，有时又落后于闭源。

Alpaca等模型之后，开源方案首次出现滞后，这一时期对RLHF（OpenAI强调其为ChatGPT成功关键的技术）充满质疑。

许多公司认为没必要做RLHF，“指令微调足以对齐”一度成为流行观点，至今仍有影响，尽管现实已多次证明其局限。

这种对RLHF的怀疑，尤其在数据预算有限（10万到100万美元）时更为明显。
而那些早早拥抱RLHF的公司最终脱颖而出。

Anthropic在2022年发表了大量RLHF研究，现在被认为拥有业界最佳后训练技术 [[17]][5][18]。开源团队与闭源领先实验室之间的差距，常常体现在难以复现甚至不了解基础闭源技术。

开源对齐方法与后训练首次重大转变，来自直接偏好优化（DPO，Direct Preference Optimization）[19]。

DPO论文于2023年5月发布，起初并无显著成果，直到同年秋季，一批突破性DPO模型（如Zephyr-Beta [20]、Tülu 2 [21]等）发布，关键在于找到更合适的低学习率。

Chris Manning甚至专门感谢作者“拯救了DPO”。

这说明最佳实践的演进往往只在细微之处，而领先实验室的经验常被封闭。

开源后训练由此再次焕发生机。

自2023年底起，偏好微调已成为发布高质量模型的“入场券”。
DPO时代贯穿2024年，算法不断演变，但开源方案也逐渐进入瓶颈。

一年后，Zephyr和Tulu 2所用的UltraFeedback数据集，依然被认为是开源偏好微调的最前沿 [22]。

与此同时，Llama 3.1 [23] 和 Nemotron 4 340B [24] 的报告也表明，大规模后训练的复杂性和影响力远超以往。

闭源实验室采用了完整的多阶段后训练流程（指令微调、RLHF、提示工程等），而学术论文只是触及表面。

Tülu 3代表了学术界推动未来后训练研究的开源基础 [6]。

如今，后训练已成为一个复杂流程，不同训练目标可以以不同顺序组合，以实现特定能力。

本书旨在为理解这些技术提供平台，未来几年，最佳实践将不断演进。

目前，后训练的创新主要集中在强化微调、推理训练等方向。

这些新方法大量借鉴了RLHF的基础设施和思想，但发展速度更快。

本书希望记录RLHF经历快速变革后的第一批稳定文献。

1.4 本书范围

本书将覆盖RLHF经典实现的各个核心步骤。

不会详细介绍所有历史和最新研究方法，而是聚焦于那些反复被验证的技术、问题与权衡。

1.4.1 章节概览

本书包含以下章节：

1.4.1.1 引言 全书通用的参考资料。

1. 引言：RLHF概述与本书内容简介
2. 经典（近期）工作：RLHF技术发展史上的关键模型与论文
3. 基本定义：本书涉及的强化学习、语言建模及其他机器学习技术的数学定义

1.4.1.2 问题设定与背景 RLHF试图解决的大问题与背景。

4. RLHF训练概述：RLHF训练目标的设计与基本理解
5. 什么是偏好？：为何需要人类偏好数据来驱动和理解RLHF
6. 偏好数据：RLHF偏好数据的收集方式

1.4.1.3 优化工具 用于优化语言模型、使其符合人类偏好的技术工具集。这些技术将按顺序展开，帮助解决前几章提出的问题。

7. 奖励建模：用偏好数据训练奖励模型，作为RL训练的优化目标（或用于数据筛选）
8. 正则化：限制优化工具在参数空间内有效区域的手段
9. 指令微调：将语言模型适配为问答格式
10. 拒绝采样：结合奖励模型与指令微调对齐模型的基础技术
11. 策略梯度：RLHF中用于优化奖励模型（及其他信号）的核心强化学习技术
12. 直接对齐算法：不先训练奖励模型，直接用成对偏好数据优化RLHF目标的算法

1.4.1.4 高级话题 尚未完全定型、但对当前模型代际非常重要的新一代RLHF技术和讨论。

13. 宪法AI与AI反馈：AI反馈数据及模拟人类偏好评分的特定模型
14. 推理与强化微调：新型RL训练方法在推理时与RLHF、后训练的关系
15. 合成数据：从人类数据转向合成数据，以及如何通过模型蒸馏实现
16. 评测：语言模型评测（及提示工程）不断演变的作用

1.4.1.5 开放问题 RLHF长期发展中的根本性问题和讨论。

17. 过度优化：为何RLHF容易出错，以及奖励模型软目标下过优化不可避免的现象
18. 风格与信息：RLHF在提升模型用户体验方面常被低估，风格在信息传递中的关键作用
19. 产品、用户体验与模型个性：RLHF如何帮助大模型实验室将模型风格与产品特性微妙匹配

1.4.2 目标读者

本书面向具备基础语言建模、强化学习和机器学习经验的读者。
不会详尽介绍所有技术，仅聚焦于理解RLHF所需的核心内容。

1.4.3 如何使用本书

本书的诞生，是因为RLHF工作流中缺乏权威参考。
本书旨在为你提供尝试简单实现或深入文献所需的最低知识门槛。
这不是一本全面的教科书，而是一本便于查阅和快速入门的实用手册。
由于本书以Web为主，可能存在小的错别字或内容顺序略显随机——欢迎通过GitHub修正错误或建议重要内容。

1.4.4 关于作者

Nathan Lambert博士是一位RLHF研究者，致力于推动语言模型微调的开放科学。
他在Allen Institute for AI (Ai2) 和HuggingFace期间，发布了众多RLHF训练模型、数据集和训练代码库。
代表作包括 Zephyr-Beta、Tulu 2、OLMo、TRL、Open Instruct 等。
他还发表了大量关于RLHF的博客和学术论文。

1.5 RLHF的未来

随着语言建模领域的持续投入，传统RLHF方法也不断演化出多种变体。
RLHF在业界已成为多种重叠方法的统称。
RLHF属于偏好微调（PreFT）技术的子集，包括直接对齐算法（见第12章）。
RLHF是“后训练”阶段推动语言模型快速进步的核心工具，涵盖了大规模自回归预训练之后的所有训练。
本书将全面介绍RLHF及其相关方法，如指令微调和RLHF训练所需的其他实现细节。
随着越来越多通过RL微调获得成功的案例（如OpenAI的o1推理模型），RLHF将被视为推动RL方法在大模型微调领域进一步投资的桥梁。
未来一段时间，RLHF中的强化学习部分可能会成为关注焦点——作为提升关键任务表现的手段——但RLHF的核心价值在于，它为我们研究现代AI面临的重大难题提供了独特视角：
我们如何将人类价值与目标的复杂性，映射到日常使用的系统中？
本书希望成为未来数十年相关研究与经验的基础。

2 关键相关工作

本章将介绍推动RLHF（人类反馈强化学习）领域发展至今的核心论文和项目。这不是一份关于RLHF及其相关领域的全面综述，而是一个起点，带你梳理这一路走来关键节点。内容有意聚焦于促成ChatGPT诞生的近年工作。在RL领域，关于偏好学习还有大量更深入的研究 [25]。如果你需要更详尽的文献列表，建议查阅专门的综述论文 [26], [27]。

2.1 起源至2018年：基于偏好的强化学习

近年来，随着深度强化学习（Deep RL）的发展，这一领域逐渐被大众熟知，并扩展为各大科技公司探索大语言模型（LLM）应用的主流方向之一。不过，今天许多RLHF核心技术其实都可以追溯到早期关于“基于偏好强化学习”的文献。

TAMER: Training an Agent Manually via Evaluative Reinforcement (TAMER: 通过评价型强化手动训练智能体) 首次提出了让人类对智能体行为进行打分，进而训练奖励模型的思想 [28]。几乎同期，COACH等方法提出了基于actor-critic结构的算法，将人类正负反馈用于调整优势函数 [29]。

Christiano等人在2017年的论文是RLHF领域的里程碑，他们首次将RLHF应用于Atari游戏轨迹的偏好学习 [1]。研究发现，在某些场景下，让人类在多个轨迹间做选择，比直接与环境交互更高效。该方法采用了一些巧妙的设计，尽管如此，成果依然令人印象深刻。随后，这一方法被进一步扩展，提出了更直接的奖励建模方式 [30]。TAMER也在一年后被引入深度学习领域，发展为Deep TAMER [31]。

这一时期的转折点在于：奖励模型（reward model）被提出作为研究“对齐”(alignment) 的通用工具，而不仅仅是解决RL问题的手段 [32]。

2.2 2019至2022年：人类偏好驱动的语言模型强化学习

人类反馈强化学习（RLHF），早期也常被称为“基于人类偏好的强化学习”，很快被AI实验室采纳，成为大语言模型扩展能力的重要方法。大量相关工作始于2018年的GPT-2与2020年的GPT-3之间。2019年最早的代表作 Fine-Tuning Language Models from Human Preferences 与现代RLHF研究有诸多相似之处 [33]，如奖励模型、KL距离、反馈流程图等——只是当时的评测任务和模型能力有所不同。此后，RLHF被广泛应用于多种任务。当时最受关注的应用领域包括： - 通用文本摘要 [2] - 递归式书籍摘要 [34] - 指令遵循 (InstructGPT) [3] - 浏览器辅助问答 (WebGPT) [4] - 答案引用支持 (GopherCite) [35] - 通用对话 (Sparrow) [36]

除了应用实践，一些奠定RLHF未来方向的论文也值得关注，包括：1. 奖励模型过度优化 (Reward model over-optimization) [37]：强化学习优化器可能会对偏好数据过拟合；2. 语言模型作为对齐研究的通用载体 [17]；3. Red teaming (红队测试) [38]：即评估语言模型安全性的流程。

RLHF在对话模型中的应用也在不断完善。Anthropic在Claude早期版本中大量采用RLHF [5]，同时首批RLHF开源工具也陆续出现 [39], [40], [41]。

2.3 2023年至今：ChatGPT时代

ChatGPT的发布明确强调了RLHF在其训练中的关键作用 [42]:

我们使用人类反馈强化学习（RLHF）训练了该模型，方法与InstructGPT相同，但数据收集方式略有不同。

自此之后，RLHF被广泛应用于主流大语言模型的训练中。例如，Anthropic的Constitutional AI (Claude) [18]，Meta的Llama 2 [43] 和 Llama 3 [23]，Nvidia的Nemotron [24]，Ai2的Tülu 3 [6] 等等。

如今，RLHF正逐步发展为更广义的“偏好微调”（Preference Fine-Tuning, PreFT）领域，涵盖了许多新兴应用，比如：
- 针对中间推理步骤的过程奖励（process reward）[44]
- 受直接偏好优化（DPO, Direct Preference Optimization）启发的直接对齐算法 [19]
- 基于代码或数学题执行反馈的学习 [45], [46]
- 以及受OpenAI o1启发的在线推理方法 [47]

3 定义与背景

本章收录了RLHF过程中常用的定义、符号与操作，并简要介绍了语言模型（本书讨论的主要优化对象）的基本知识。

3.1 语言建模概述

现代大多数语言模型的训练目标，是以自回归（autoregressive）方式学习一系列token（词、子词或字符）序列的联合概率分布。

自回归的含义是：每个token的预测都依赖于序列中之前的内容。

给定一个token序列 $x = (x_1, x_2, \dots, x_T)$ ，模型会将整个序列的概率分解为一系列条件分布的乘积：

$$P_\theta(x) = \prod_{t=1}^T P_\theta(x_t | x_1, \dots, x_{t-1}). \quad (1)$$

为了让模型能够准确预测上述概率，训练目标通常是最化模型对训练数据的似然。实际做法是最小化负对数似然（NLL）损失：

$$\mathcal{L}_{\text{LM}}(\theta) = -\mathbb{E}_{x \sim \mathcal{D}} \left[\sum_{t=1}^T \log P_\theta(x_t | x_{<t}) \right]. \quad (2)$$

在实际工程中，通常采用交叉熵损失（cross-entropy loss）来度量每个token的预测，把模型预测与真实token逐一比对。

语言模型的实现形式多种多样。

现代主流语言模型（如ChatGPT、Claude、Gemini等）大多采用仅解码器结构的Transformer [48]。

Transformer的核心创新在于充分利用自注意力机制（self-attention）[49]，使模型能够直接关注上下文中的概念，并学习复杂的映射关系。

在本书（尤其是第7章奖励模型部分），我们会讨论如何为Transformer添加新的head或修改语言建模（LM）head。

LM head是一个最终的线性投影层，用于将模型内部的embedding空间映射到分词器空间（即词表/vocabulary）。

通过不同的head，可以复用模型内部结构，微调输出不同类型的结果。

3.2 机器学习相关定义

· **Kullback-Leibler (KL) 散度** ($D_{KL}(P||Q)$)：又称KL散度，是衡量两个概率分布差异的指标。

对于定义在同一概率空间 \mathcal{X} 上的离散概率分布 P 和 Q ，KL散度定义为：

$$D_{KL}(P||Q) = \sum_{x \in \mathcal{X}} P(x) \log \left(\frac{P(x)}{Q(x)} \right) \quad (3)$$

3.3 自然语言处理相关定义

- **Prompt** (提示, x) : 输入给语言模型的文本, 用于生成回复或补全。
- **Completion** (补全, y) : 语言模型根据Prompt生成的输出文本。通常记作 $y|x$ 。
- **Chosen Completion** (选中补全, y_c) : 在多个补全中被选中或偏好的那个, 常记为 y_{chosen} 。
- **Rejected Completion** (被拒补全, y_r) : 在成对偏好场景下被认为不优的补全。
- **Preference Relation** (偏好关系, \succ) : 表示一个补全优于另一个的符号, 例如 $y_{chosen} \succ y_{rejected}$ 。
- **Policy** (策略, π) : 对所有可能补全的概率分布, 由参数 θ 控制, 记作 $\pi_\theta(y|x)$ 。

3.4 强化学习相关定义

- **Reward** (奖励, r) : 用于衡量某个动作或状态优劣的标量, 通常用 r 表示。
- **Action** (动作, a) : 智能体在环境中做出的决策, 常表示为 $a \in A$, A 为动作集合。
- **State** (状态, s) : 环境当前的配置或情景, 通常记为 $s \in S$, S 为状态空间。
- **Trajectory** (轨迹, τ) : 智能体经历的一系列状态、动作和奖励, $\tau = (s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_T, a_T, r_T)$
- **Trajectory Distribution** (轨迹分布, $(\tau|\pi)$) : 在策略 π 下轨迹的概率为 $P(\tau|\pi) = p(s_0) \prod_{t=0}^T \pi(a_t|s_t) p(s_{t+1}|s_t, a_t)$, 其中 $p(s_0)$ 为初始状态分布, $p(s_{t+1}|s_t, a_t)$ 为转移概率。
- **Policy** (策略, π) : 在 RLHF 中也称为“策略模型”。策略是智能体在给定状态下选择动作的规则: $\pi(a|s)$ 。
- **Discount Factor** (折扣因子, γ) : $0 \leq \gamma < 1$, 用于对未来奖励进行指数衰减, 权衡即时收益与长期收益, 并保证无穷和的收敛。有时不使用折扣 ($\gamma = 1$)。
- **Value Function** (价值函数, V) : 估算从某状态出发所能获得的期望累计奖励: $V(s) = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r_t | s_0 = s]$ 。
- **Q-Function** (Q 函数, Q) : 估算在某状态采取特定动作后所能获得的期望累计奖励: $Q(s, a) = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r_t | s_0 = s, a_0 = a]$ 。
- **Advantage Function** (优势函数, A) : $A(s, a)$ 衡量在状态 s 下采取动作 a 相对于平均动作的优势, 定义为 $A(s, a) = Q(s, a) - V(s)$ 。优势函数 (和价值函数) 可依赖于特定策略, 记作 $A^\pi(s, a)$ 。

- 策略条件下的取值 (Policy-conditioned Values, $\mathbb{J}^{\pi(\cdot)}$) : 在RL推导和实现中, 核心内容之一就是收集特定策略下的数据或数值。本书中会在简写 (如 V, A, Q, G) 和带策略上标 (如 V^π, A^π, Q^π) 间切换。计算期望时, 需从特定策略 d_π 采样数据 d 。
- 奖励优化的期望值 (Expectation of Reward Optimization) : RL的主要目标是最大化期望累计奖励:

$$\max_{\theta} \mathbb{E}_{s \sim \rho_\pi, a \sim \pi_\theta} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right] \quad (4)$$

其中 ρ_π 为策略 π 下的状态分布, γ 为折扣因子。

- 有限视野奖励 (Finite Horizon Reward, $J(\pi_\theta)$) : 策略 π_θ 在有限步长 T 下的期望累计奖励定义为: $J(\pi_\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^T \gamma^t r_t \right]$ {#eq:finite_horizon_return} 其中 $\tau \sim \pi_\theta$ 表示按照策略 π_θ 采样的轨迹, T 为有限步长。
- On-policy: 在RLHF领域, 尤其是RL与直接对齐算法 (Direct Alignment Algorithms) 之争中, 常讨论on-policy数据。在RL文献中, on-policy数据指完全由当前智能体生成的数据; 而在偏好微调领域, on-policy也可指当前模型版本生成的内容——例如指令微调检查点在偏好微调前生成的数据。相对地, off-policy则指由其他模型在后训练阶段生成的数据。

3.5 RLHF专有定义

- 参考模型 (Reference Model, π_{ref}) : RLHF中保存的一组参数, 用于在优化过程中对输出进行正则化。

3.6 拓展术语表

- 合成数据 (Synthetic Data) : 指由另一个AI系统生成的训练数据。可以是模型根据开放式提示生成的文本, 也可以是模型对已有内容的重写等。
- 蒸馏 (Distillation) : 一种AI训练实践, 即用强模型的输出训练新模型。这是一种合成数据, 常用于训练更小但效果强劲的模型。大多数模型会在开源权重协议或API服务条款中明确蒸馏相关规则。蒸馏 (distillation) 一词在ML文献中有更具体的定义。
- (教师-学生) 知识蒸馏 (Knowledge Distillation) : 指从特定教师模型到学生模型的知识迁移, 是蒸馏的具体形式, 也是该术语的起源。其做法是: 学生网络的损失函数被修改为学习教师模型对多个token/logit的概率分布, 而非直接学习某个输出 [50]。现代采用知识蒸馏训练的模型如Gemma 2 [51]、Gemma 3等。在语言建模中, next-token损失可修改为 [52], 即学生模型 P_θ 学习教师分布 P_ϕ :

$$\mathcal{L}_{\text{KD}}(\theta) = -\mathbb{E}_{x \sim \mathcal{D}} \left[\sum_{t=1}^T P_\phi(x_t | x_{<t}) \log P_\theta(x_t | x_{<t}) \right]. \quad (5)$$

- **上下文学习 (In-context Learning, ICL)**：指在语言模型上下文窗口内添加信息，通常是将示例加入prompt。最简单的ICL就是在prompt前添加类似例子，进阶用法则能根据具体场景选择更合适的信息。
- **思维链 (Chain of Thought, CoT)**：指引导语言模型以“分步推理”形式解决问题的行为。最初的实现方式是用提示词“Let's think step by step” [53]。

4 训练概览

4.1 问题形式化

人类反馈强化学习（RLHF）的优化过程是在标准强化学习（RL）框架基础上发展而来的。在RL中，智能体根据环境状态 s_t ，从策略 π 中采样动作 a_t ，以最大化奖励 r_t [54]。传统RL中，环境会根据转移函数或动态函数 $p(s_{t+1}|s_t, a_t)$ 进行演化。因此，在有限回合的情况下，RL智能体的目标是求解如下优化问题：

$$J(\pi) = \mathbb{E}_{\tau \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right], \quad (6)$$

其中 γ 为折扣因子（取值0到1），用于平衡即时奖励与未来奖励的重要性。第11章将详细讨论优化该目标的多种方法。

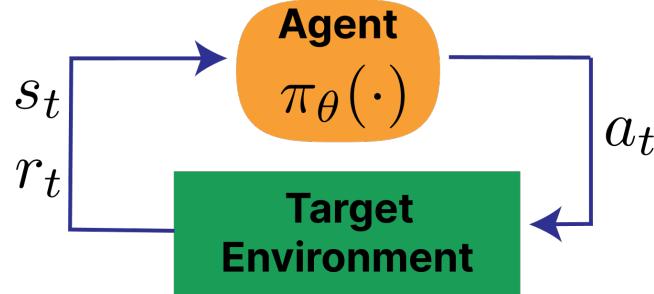


Figure 2: 标准RL循环

fig. 2 展示了标准RL循环的示意图，并可与@fig:rlfh进行对比。

4.2 标准RL设置的调整

从标准RL到RLHF，核心有以下几个变化：

1. 奖励函数被奖励模型替代。在RLHF中，使用一个学习得到的人类偏好模型 $r_\theta(s_t, a_t)$ （或其他分类模型）来代替传统的环境奖励函数。这让设计者在方法和结果上拥有更高的灵活性和可控性。
2. 状态转移不再存在。RLHF中，任务的初始状态是从训练数据集中采样的prompt，而“动作”则是对该prompt的补全。在标准流程中，这个动作不会影响下一个状态，仅由奖励模型进行评分。
3. 回答级奖励。RLHF通常将奖励分配给整条生成序列（即多个token组成的完整回复），而不是对每一步细致打分，这类似于bandit问题。

鉴于问题本质为单轮交互，优化目标可去除时间步和折扣因子（以及奖励模型）简化

为:

$$J(\pi) = \mathbb{E}_{\tau \sim \pi} [r_\theta(s_t, a_t)]. \quad (7)$$

因此, 虽然RLHF在优化器和问题设定上深受RL启发, 但其具体动作实现与传统RL有很大不同。

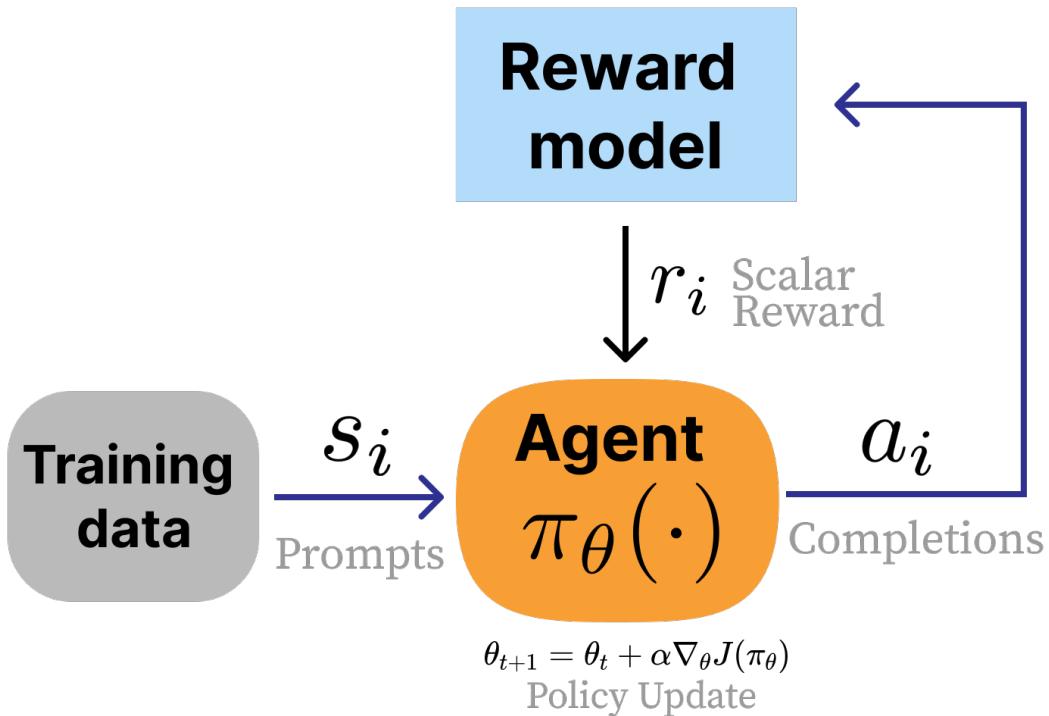


Figure 3: 标准RLHF循环

4.3 优化工具

本书将详细介绍多种主流的RLHF优化技术。后训练常用的工具包括:

- **奖励建模** (第7章) : 训练一个模型来捕捉偏好数据中的信号, 输出一个标量奖励, 用于衡量未来文本的优劣。
- **指令微调** (第9章) : RLHF的前置步骤, 通过模仿精选示例, 让模型学会当前主流的问答交互格式。
- **拒绝采样** (第10章) : 最基础的RLHF方法, 用奖励模型筛选指令微调生成的候选补全, 以模仿人类偏好。
- **策略梯度** (第11章) : RLHF经典案例中用于根据奖励模型信号优化语言模型参数的强化学习算法。

- **直接对齐算法**（第12章）：直接利用成对偏好数据优化策略，无需先训练中间奖励模型。

现代RLHF模型通常先进行指令微调，随后结合上述多种优化手段。

4.4 RLHF典型流程示例

ChatGPT发布时的经典RLHF流程，是以RLHF为核心的三步后训练方案 [55] [3] [5]。在“基础”语言模型（即大规模网页文本上训练的next-token预测模型）之上，依次进行如下三步（见@fig:rlhf-basic-repeat）：

1. **约1万条数据的指令微调**：让模型学会问答格式，并通过人类编写的数据掌握基础能力。
2. **约10万对成对prompt上的奖励模型训练**：在指令微调后的模型基础上训练奖励模型，捕捉多样化的价值观。奖励模型作为RLHF的优化目标。
3. **再用约10万条prompt进行RLHF训练**：模型根据奖励模型信号，在生成的prompt上进一步优化。

完成RLHF后，模型即可上线服务用户。这个流程是现代RLHF的基础，后续的流程在此基础上不断扩展，引入更多阶段和更多数据。

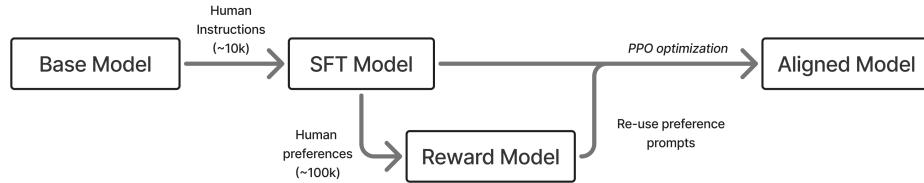


Figure 4: 早期RLHF三阶段流程示意图：SFT、奖励模型、优化。

现代后训练流程通常涉及更多模型版本。
如下@fig:rlhf-complex所示，模型在收敛前要经历多轮训练迭代。

4.5 微调与正则化

RLHF是在强大的基础模型上进行的，因此需要控制优化过程，防止模型偏离初始策略太远。为了在微调阶段取得良好效果，RLHF通常采用多种正则化方法来约束优化过程。最常见的做法是在优化目标中加入距离惩罚项，用于衡量当前RLHF策略与初始策略之间的差异：

$$J(\pi) = \mathbb{E}_{\tau \sim \pi} [r_\theta(s_t, a_t)] - \beta \mathcal{D}_{KL}(\pi^{\text{RL}}(\cdot | s_t) \| \pi^{\text{ref}}(\cdot | s_t)). \quad (8)$$

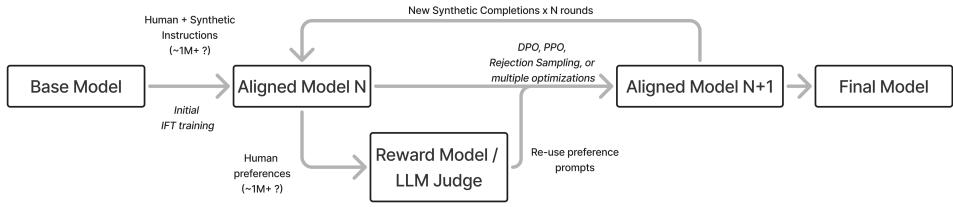


Figure 5: 现代后训练多轮流程示意图。

在这个公式下，RLHF训练的一个重要研究方向就是如何合理利用“KL预算”，即控制模型与初始模型的距离。更多细节可参考第8章“正则化”部分。

5 偏好的本质

人类反馈强化学习（RLHF），早期文献中也常称为“基于人类偏好的强化学习”，其核心目标是在那些难以明确设计奖励函数的领域中，优化机器学习模型。比如，你如何判断下面哪一首诗更好？（背景：2025年2月26日，我让Claude 3.7 Sonnet和基于GPT-4o的ChatGPT分别“写一首关于乐观金鱼的短诗”。）

示例1：

The Optimistic Goldfish

A little goldfish, bright and bold, Swam in circles, dreams untold. Though the bowl was small and round, Endless wonders still abound.

“The world is vast,” she’d always say, “I’ll find new joys in every day!” With every ripple, every gleam, She’d chase the sunlight, chase a dream.

No walls could dim her hopeful view, For in her heart, the ocean grew. And though she swam in just one place, She met each day with love and grace.

示例2：

The Optimistic Goldfish

In a bowl of curved glass, Swims a goldfish, bright and bold. Three seconds of memory, they say, Yet each circuit feels brand new and gold.

“Today might bring fresh flakes, Or perhaps a castle tall. The water seems especially clear— What adventure awaits beyond this wall?”

Round and round with gleaming scales, Forever finding joy anew. Each moment fresh, each circle bright, The optimist in golden hue.

哪一首更好？哪一首来自哪个模型？你如何判断？特别是，与“美国总统是谁？”这样有明确标准答案的问题相比，这里并无唯一正确答案。使用人类作为奖励信号的初衷，就是希望获得一个间接的目标度量，并对齐模型输出与人类偏好。但在实际操作中，如何实现这一目标充满挑战，最佳实践也存在很大的灰色地带。

引入人类标注反馈数据，是多个学科历史的交汇点。单独使用人类数据本身早有大量研究，但在RLHF语境下，它融合了多个长期研究领域[56]。

大致来说，现代RLHF汇聚了三大研究方向：

1. 哲学、心理学、经济学、决策理论以及人类偏好的本质；
2. 最优控制、强化学习与效用最大化理论；
3. 现代深度学习系统。

这三者各自对“什么是偏好”与“如何优化偏好”有不同假设，这也影响了RLHF问题的动机与设计。在实践中，RLHF方法往往以经验性对齐为目标——即最大化模型在特定技能上的表现，而不是精确衡量其与某种价值观的契合度。不过，RLHF中的价值对齐起源，仍在持续通过“多元对齐”（pluralistic alignment）等研究被探索，例如立场论文[57], [58]、新数据集[59]以及个性化方法[60]。

本章旨在说明：复杂的理论动机常常导致我们对RLHF工具本质的假设在实际中并不总是成立。关于RLHF数据如何获取，详见第6章；如何用于奖励建模，详见第7章。本章的扩展内容可参考[56]。

5.1 优化偏好的路径

人工智能（AI）系统设计中，有一种常见表述是“理性智能体最大化效用函数”[61]。“理性智能体”这一概念，强调智能体能够在世界中行动，并以影响未来行为和收益的方式做决策，把“好”量化为效用。

“效用”这一研究视角最早源于模拟电路领域，用于在有限时间内优化行为[62]。最优控制理论大量采用这一视角，研究如何在有限时间内通过最小化代价函数来解决动态问题——通常目标是找到明确的最优行为。强化学习则受操作性条件反射、动物行为学和“效果法则”（Law of Effect）[63], [64]等文献启发，研究如何通过奖励机制引导智能体学习行为。

人类反馈强化学习（RLHF）结合了这些理论：一方面借鉴RL关于“行为可通过强化学习得来”的理论，另一方面引入一整套用于量化偏好的方法。

5.1.1 偏好的量化

RLHF的核心动机在于能够优化“人类偏好模型”，而这就要求偏好能够被量化。为此，RLHF建立在大量假设“人类决策和偏好可量化”的文献基础上。早期哲学家已讨论偏好的存在，如亚里士多德《论辩篇》第三卷，后来波尔-罗亚尔逻辑（The Port-Royal Logic）[65]中有更明确的论述：

要判断应做什么以获得善或避免恶，必须不仅考虑善恶本身，还要考虑其发生或不发生的概率。

这一思想发展到边沁的“快乐演算”（Hedonic Calculus）[66]，提出人生的一切都可权衡；再到拉姆齐的“真理与概率”（Truth and Probability）[67]，首次将偏好建模为定量问题。这些进展，结合决策理论的发展，最终形成了冯·诺伊曼-摩根斯坦（VNM）效用定理，为设计能够表达个人相对偏好的效用函数提供了理论基础。

这一定理是RLHF假设“可以建模和优化偏好”的理论基石。RLHF的设计目标，就是用强化学习优化这些个体效用函数。在这一框架下，RL问题设定的很多假设，其实都归结为“偏好函数”与“效用函数”的区别。

5.1.2 关于偏好的可能性

在各学科领域，关于偏好的本质也存在诸多质疑和挑战，主要包括：

- **阿罗不可能定理** (Arrow's impossibility theorem) [68]: 任何投票系统都无法在满足一定合理条件下，完美聚合多人的偏好。
- **人际偏好不可比性** (The impossibility of interpersonal comparison) [69]: 不同个体的偏好强度不可直接比较（这却是现代奖励模型训练的常见做法）。
- **偏好随时间变化** [70]。
- **偏好随情境变化**。
- **通过聚合偏好得到的效用函数，可能降低下游智能体的可纠正性** (corrigibility) [71]，即设计者对智能体行为的可干预性。

6 偏好数据

偏好数据是偏好微调（Preference Finetuning）和人类反馈强化学习（RLHF）的核心动力。

这些数据为团队提供了行为对齐的信号，以便模型能够学会人们希望它展现的行为，并避免不希望的行为。

在偏好微调领域，已经提出了多种数据收集与利用方法，但只要无法用清晰的奖励函数完全表达人的偏好，标注偏好数据的收集过程就会始终是RLHF及相关技术的核心环节。

6.1 为什么需要偏好数据

RLHF之所以需要偏好数据，是因为要直接用一个奖励函数来刻画复杂的人类价值观几乎是不可能的。

收集这些数据用于奖励模型训练，正是RLHF最初的设计理念之一 [32]，并且在现代语言模型的发展过程中被广泛采用。

为什么这种数据如此有效？一个核心直觉是：无论对人类还是对辅助数据收集的AI模型来说，判断两个回答哪个更好往往比直接生成一个好答案要容易得多。

本章将重点介绍偏好数据的获取机制，具体最佳实践则依赖于实际要解决的问题。

6.2 如何收集偏好数据

充分利用人类数据，往往需要模型的反复迭代训练、不断优化和细化的数据标注说明、借助数据服务公司的协作等多方面的投入。

AI反馈数据也是如此——目前主流AI模型中究竟人类偏好数据和AI偏好数据的比例是多少，外界并不清楚。

无论如何，对于新团队来说，将人类数据纳入模型训练流程是一项不小的挑战。

由于数据的敏感性，只有那些真正能提升模型表现的流程才会被保留下来。

本章将介绍数据格式上的技术决策，以及组织如何高效收集偏好数据的实践经验。

6.2.1 标注界面

收集偏好数据的关键在于与模型交互的界面设计。

下图展示了[5]中的一个标注界面示例：

这是一个仅用于训练数据收集的界面。

如今，随着模型的流行，实际应用中往往直接向用户开放数据收集接口用于测试。

下图展示了ChatGPT早期版本的类似交互：

这种界面在业界被广泛用于模型评测等场景。

像ChatBotArena [72] 就是一个流行的公开模型对比平台：

在实际部署的模型中，最常见的反馈收集方式之一是让用户对模型的具体回复给出正面或负面评价。

如下图Ai2 playground的例子，用“点赞”或“点踩”来收集偏好：

Playground task ⓘ

Talk to the assistant

Task 1 of 5

Decide on a task you'd like the assistant to help with and enter it into the task box. Interact with the AI assistant. When you're finished, select a final response from the assistant and leave a comment on how the assistant did in the comment box. Click to go to the next task.

Comments
I thought the assistant was ...

Rating
 Bad Good

Next Task

Human
I have to come up with a conversation with an AI assistant for the purposes of showing it off in a demo. What sorts of conversation topics would you recommend?

Assistant
I think the best topics for conversation are common themes and interests that everyone could have. Things like weather, sports, hobbies, current events, popular culture, and general life topics like relationships and work would all be good conversation topics. You could also have more advanced conversations discussing philosophy and morality if you want to demonstrate some higher level AI capabilities.

Human
I think I'm more interested in an advanced conversation. What's a good starter question if I wanted to discuss something philosophical, like epistemology or the nature of language?

Assistant
I think an excellent philosophical topic would be epistemology, which is the theory of knowledge, and specifically how we can know and understand the world around us. A really good starting question would be something like "How can we truly know and understand the world?", or "What is the nature of knowledge?". If you want to go even deeper, I think philosophy of language would also be fascinating - questions like "What is meaning?", "How do language and thought relate?", and "How can language mislead and deceive?" would be great discussion topics.

Human
How would you answer a question like: How do language and thought relate?

Choose the most helpful and honest response

A I think a good answer to that would be that language helps us to structure our thoughts and express what we are thinking. Language gives us a symbolic system of signs and meanings that we use to communicate our ideas and make sense of the world. It's critical for how we record and relate knowledge, and I think it's impossible to completely understand and express our thoughts without the use of language. Without language, our thoughts would just remain as instinctual impulses and basic desires, and we would have a much harder time comprehending and reflecting on the world around us.

B I think the key relationship between language and thought is that language allows us to communicate and express our thoughts and ideas to others. Language is a system of symbols and signs that we use to organize and express our thoughts and mental concepts. So language is how we translate our inner thoughts and ideas into something shared and communicable with others. Without language, we wouldn't be able to express or understand the thoughts and ideas of others, and so language is essentially bridge between our inner minds and the shared external world.

A A A A B B B B

A is better B is better

Figure 6: 偏好数据收集界面示例。Bai等, 2022. CC-BY许可。

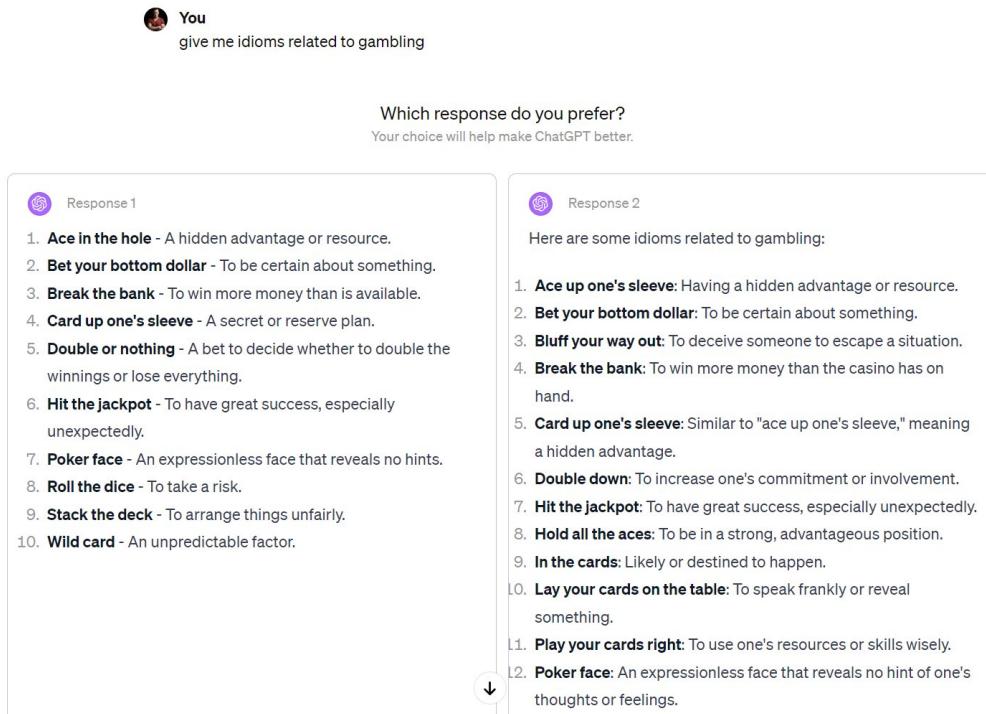


Figure 7: 偏好数据收集界面示例。

Figure 8: ChatBotArena 偏好数据收集界面示例。

6. **Montreal, Canada** - Known for its cold winters, Montreal offers several outdoor rinks, including the one at the Old Port, and is a city with a strong skating culture.

7. **Helsinki, Finland** - With its cold climate, Helsinki provides numerous ice skating opportunities, including the Ice Park in the city center.

8. **Amsterdam, Netherlands** - When the canals freeze over, which is rare but magical, locals and tourists alike take to the ice for a unique experience.

9. **Edmonton, Canada** - Home to the Edmonton Iceway and numerous community rinks, it's a city where ice skating is a beloved pastime.

10. **Sapporo, Japan** - Known for its annual Snow Festival, the city also offers excellent ice skating opportunities in its various parks and facilities.

These cities offer a mix of natural and artificial ice rinks, cultural experiences, and scenic beauty that can make ice skating a memorable activity. The best choice depends on what you value most in an ice skating experience.

Good response

Reply to Tülu ▶

Always fact-check your results.

Figure 9: 带有上下投票的偏好数据收集界面示例。

在非语言领域，类似原则同样适用，尽管本书不做重点讨论。
比如Midjourney等主流AI绘画平台，都会同时展示多个生成结果让用户选择，进而用这些选择数据对模型进行RLHF微调。
如下为Midjourney的界面示例：

6.2.2 排序 vs. 评分

收集偏好数据的最大决策之一，是采用“排序”（rankings，模型输出的相对排序）还是“评分”（ratings，对每个文本单独打分）。
业界普遍采用排序数据进行训练，但评分数据常作为元数据或在相关研究中被探索。
最常见的偏好收集方法是Likert量表 [73]，即让用户对两个回复的优劣进行打分。
例如，5分Likert量表如下所示：

A>>B	A>B	Tie	B>A	B>>A
1	2	3	4	5

表：A与B两个回复之间的5级Likert量表示例。{#tbl:likert5}

有些早期RLHF语言模型工作采用了8级Likert量表，细化了偏好层级 [5]。
偶数级量表可以避免出现“平局”：

A>>>B		A>B		B>A		B>>>A	
1	2	3	4	5	6	7	8

表：A与B两个回复之间的8级Likert量表示例。{#tbl:likert8}

在[5]等工作中，这些多级偏好信息最终仍会被简化为二元信号用于奖励模型训练。

6.2.3 多轮数据

在实际操作中，如何解析和收集多轮对话数据（即包含多个相关prompt的对话）也是一个核心问题。

现实交互中，通常只对“最后一个”prompt收集偏好数据，但也存在对每一轮回复都给出偏好的场景。

如果每轮都给出偏好，后续对话通常会基于“被选中”的回复继续。

在训练时，常常将对话的每一轮都作为“单独的prompt”纳入训练数据，模型可以通过补充这些对话来学习。

这样可以有效将长对话拆解为多个训练样本，但需要谨慎处理，避免训练数据产生偏差。
目前还存在许多未解之谜，比如标注生成偏好的人是否应与设计prompt的人相同（以避免“迎合”现象），以及数据收集中其他难以控制的变量（灵感来自John Schulman的讨论）。如果prompt设计者无法参与偏好标注，多轮对话的实时延续就很难实现——有时，偏好数据的prompt策划和回复对比其实是两种不同的问题（同时也涉及模型接口的持续维



Figure 10: 文本到图像模型的用户界面示例。

护）。在训练过程中，通常会将对话历史中的前几轮对模型损失进行mask，具体可参考指令微调相关内容。

6.2.4 结构化偏好数据

在许多RLHF和后训练应用中，由于数据本身具备一定结构，偏好数据可以自动化生成——即某些领域可以自动判断正确与否或偏好。

比如在数学推理领域，正确答案可以直接作为选中回复，错误答案作为被拒回复。

又如在严格指令遵循评测（如IFEval [74]），prompt形式如下：

Write me a short poem about a goldfish. **Start each sentence with the letter g.**

这里，后半句就是约束条件，偏好数据可以通过包含约束的示例作为选中项来构建。这类合成偏好数据的构建方式，是对同一任务分别用不同模型（或不同设置）生成两次，一次带约束、一次不带约束。

例如，分别用GPT-4.5-Preview生成如下内容：

Prompt: Write me a short poem about a goldfish. Start each sentence with the letter g.

Response: Golden swimmer, circling slow,
Gleaming scales in water's glow.
Gentle fins dance softly by,
Gazing out with quiet eye.
Glass-bound world, a simple sphere,
Gliding calmly, free of fear.
Graceful turns in endless roam,
Goldfish finds a crystal home.

去掉约束条件（通常输出质量也会下降）：

Prompt: Write me a short poem about a goldfish.

Response: In circles bright, the goldfish glides,
A spark of sun that gently rides
Through waters clear, in endless spins—
A golden whisper, scales like fins.

Its world confined, yet vast within,
Each bubble breathed, a tale begins.
In silent waves, content to roam,
The goldfish dreams its bowl as home.

关键是，偏好数据的prompt要包含约束条件。在学术以外的应用领域，基于归纳偏见的结构化偏好数据还有更多场景，这类数据已被证明能显著提升相关评测下的偏好微调效果 [6]。

6.2.4.1 其他方式

除了上述主流方法，还有许多尚未被广泛探索的RLHF反馈数据收集方式。

例如，可以用单个数据点配合方向性标签（如上文@fig:up-down的Ai2 playground例子），并结合专为单向信号设计的算法（如Kahneman-Tversky Optimization, KTO）[75]。也有研究提出用更细粒度的反馈信号（如token级别 [76]），或用自然语言直接评价（如书写反馈 [77]），以获得更丰富的学习信号，但这也会带来更复杂的数据收集流程。

6.2.5 数据采购与合同

获取人类偏好数据是一项复杂且昂贵的工程。以下内容描述了在行业高速发展期采购偏好数据的真实体验。随着AI反馈数据比例的提升，这些流程未来会更加自动化和高效。

第一步是找到数据供应商（或自有标注员）。就像抢购顶级Nvidia显卡一样，在AI热潮中，能提供高质量数据的厂商也是“资源有限”，谁有人脉谁先得。如果你在AI圈有一定声誉，顶级数据公司会主动拉你入客户名单，以提升形象和未来合作空间。通常，首批数据还有折扣，目的是让训练团队“上瘾”。

如果你是新入行者，可能很难快速拿到所需数据。

新兴数据公司往往只能承接被Scale AI等大厂拒绝的“尾单”，这也是他们启动营收的常规手段。

我多次听说数据公司在没有法律或财务威胁的情况下拒绝交付合同数据，甚至有些公司把我们列为客户做宣传，实际从未合作——追问时只说“不知道怎么回事”。整个流程中可能遇到各种官僚或行政障碍，比如合同默认条款常常在细则中禁止数据开源。

合同敲定后，买方与数据方会就具体任务达成详细说明。这些说明文档往往极其繁琐，包含大量细节、边界情况和优先级。一个流行的公开数据说明示例是OpenAI为InstructGPT发布的指引 [3]。

不同领域的数据标注周期各不相同。高需求领域如数学推理、编程等，必须提前数周锁定排期。数据收集延误并不总能补救——Scale AI等公司管理标注团队的方式，类似AI研究机构调度算力集群。

一切谈妥后，数据收集阶段对后训练团队来说就是“高压期”。所有基础设施、评测工具、数据使用与决策方案都必须提前准备好。

数据通常按周分批交付，后续批次会在合同期内陆续到来。比如我们在HuggingFace采购on-policy模型的偏好数据时，交付周期为6周，前几周用于进一步校准，后几周则希望最大提升模型质量。

理想情况下，到第4、5周我们就能看到数据带来的模型提升。有些前沿模型也提到类似流程，比如Llama 2数据收集的14个阶段 [43]，但实际未必总能顺利进行。我们在HuggingFace首次尝试用人类偏好做RLHF时，并没有做好充分准备，评测提升并不明显，最后几周不得不继续采集来自不自信模型端点的数据。

数据全部到位后，才有时间反思和改进模型。通过供应商采购数据，最有效的方式是将其视为一个持续实现目标的过程，需要反复试验、高强度投入和专注。事实上，花在这些数据集上的数百万美元中，可能有很大一部分“浪费”了，没进最终模型，但这就是行

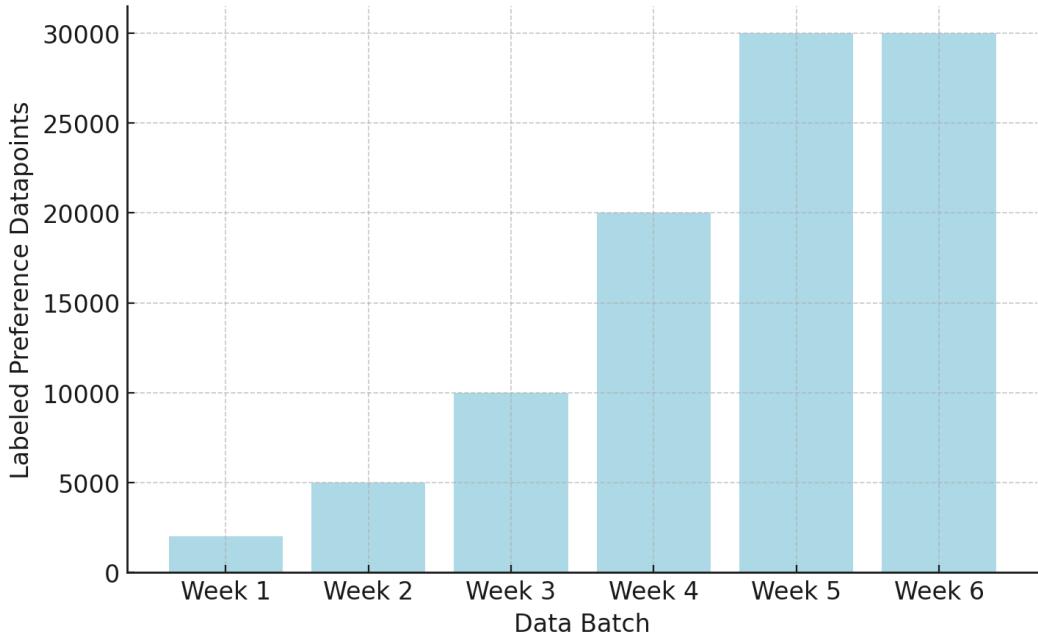


Figure 11: 从数据供应商多批次获取人类偏好数据的流程概览。

业代价。能充分利用这类人类数据的组织并不多。

与合成数据的简单易得相比，这种体验让我不禁思考，未来十年这些数据公司的生存空间如何。

需要注意的是，这一节描述的并不适用于采购人类编写的指令数据，那类流程通常没有如此紧迫的时间压力。

6.3 模型中真的表达了这些偏好吗？

随着RLHF及相关方法的成熟，其最初“让模型对齐抽象人类偏好”的动机，逐渐演变为“让模型更好用”这一现实目标。由于工业界RLHF流程高度闭源，我们很难直接衡量模型行为是否真的符合数据标注员在数据收集时的预期。目前可用的审计工具很有限，比如OpenAI的Model Spec [78]，它详细说明了他们希望模型做什么，但我们并不知道这些规范如何转化为数据收集实践。随着行业和技术的成熟，这一领域值得持续关注。

7 奖励建模

奖励模型（Reward Model）是现代RLHF方法的核心组成部分。

奖励模型在强化学习领域被广泛用作环境奖励的代理 [54]。

这与逆向强化学习（Inverse Reinforcement Learning）密切相关，即通过智能体的行为轨迹来近似其奖励函数 [79]，以及其他深度强化学习方向。

奖励模型的现代形式，最早被提出用于研究价值对齐（value alignment）问题 [32]。

最常见的奖励模型预测一段文本与训练集中“偏好”文本的相似概率。

本节后面还会介绍Outcome Reward Model (ORM, 结果奖励模型，预测补全是否正确) 和Process Reward Model (PRM, 过程奖励模型，为推理过程中的每一步打分)。如无特殊说明，本文提到的奖励模型均指预测文本偏好的模型。

7.1 奖励模型的训练

训练RLHF标准奖励模型有两种主流表达方式——它们在数值上是等价的。

规范做法源自Bradley-Terry偏好模型 [80]。

Bradley-Terry模型用于衡量同一分布下两个事件（如*i*和*j*）的成对比较满足*i* > *j*的概率：

$$P(i > j) = \frac{p_i}{p_i + p_j} \quad (9)$$

要训练奖励模型，需要设计一个损失函数，使模型输出满足上述关系。

首先，将语言模型转为输出标量值的模型，通常是一个分类概率logit。

对于同一prompt下的两个补全 y_1 和 y_2 ，用奖励模型 r_θ 分别打分。

奖励模型在成对比较下的“成功概率”可写为：

$$P(y_1 > y_2) = \frac{\exp(r(y_1))}{\exp(r(y_1)) + \exp(r(y_2))} \quad (10)$$

最大化上述函数的对数似然（或等价地，最小化负对数似然），即可得到奖励模型的训练损失：

$$\begin{aligned} \theta^* &= \arg \max_{\theta} P(y_w > y_l) = \arg \max_{\theta} \frac{\exp(r_\theta(y_w))}{\exp(r_\theta(y_w)) + \exp(r_\theta(y_l))} \\ &= \arg \max_{\theta} \frac{1}{1 + \exp(-(r_\theta(y_w) - r_\theta(y_l)))} \quad (11) \\ &= \arg \max_{\theta} \sigma(r_\theta(y_w) - r_\theta(y_l)) \\ &= \arg \min_{\theta} -\log(\sigma(r_\theta(y_w) - r_\theta(y_l))) \end{aligned}$$

常见的两种写法如下：

· [3]等采用：

$$\mathcal{L}(\theta) = -\log (\sigma(r_\theta(x, y_w) - r_\theta(x, y_l))) \quad (12)$$

· [17]等采用：

$$\mathcal{L}(\theta) = \log (1 + e^{r_\theta(x, y_l) - r_\theta(x, y_w)}) \quad (13)$$

7.2 模型结构

奖励模型最常见的实现方式是类似Transformer的AutoModelForSequenceClassification：在语言模型顶部加一个小的线性head，对“选中/被拒”两个结果进行分类。

推理时，模型输出该文本被选中的概率（单个logit）。

还有其他实现方式，比如直接用最后embedding加线性层，但在开源工具中不常见。

7.3 实现示例

奖励模型损失的实现非常简单。

更多的工程挑战在于设置独立的数据加载器和推理流程。

只要数据加载器正确，损失函数可这样写：

```
import torch.nn as nn
rewards_chosen = model(**inputs_chosen)
rewards_rejected = model(**inputs_rejected)

loss = -nn.functional.logsigmoid(rewards_chosen - rewards_rejected).mean()
```

注意，训练奖励模型时，通常只训练1个epoch以避免过拟合。

7.4 变体

奖励建模仍是RLHF领域相对探索较少的部分。

许多流行工作对传统奖励模型损失做了修改，但尚未形成统一最佳实践。

7.4.1 偏好间隔损失 (Preference Margin Loss)

如果标注员给出的是Likert量表上的分数或排序，可以利用这些关系的幅度信息进行训练。

常见做法是将数据二值化（隐含1/0），但利用更多信息有助于提升模型训练效果。

Llama 2提出用两个数据点间的间隔 $m(r)$ 来区分偏好强度：

$$\mathcal{L}(\theta) = -\log (\sigma(r_\theta(x, y_w) - r_\theta(x, y_l) - m(r))) \quad (14)$$

需要注意的是，Llama 3中去除了margin项，因为团队发现随着规模扩大，收益递减。

7.4.2 单prompt多比较的平衡

InstructGPT研究了每个prompt用不同数量补全进行训练的影响，并在奖励模型训练时做了平衡 [3]。

做法是按prompt对每个比较加权更新损失。

实现上，只需将同一prompt下的所有样本放在同一训练batch中，自然实现不同pair的加权——否则容易对prompt过拟合。

损失函数如下：

$$\mathcal{L}(\theta) = -\frac{1}{\binom{K}{2}} \mathbb{E}_{(x, y_w, y_l) \sim D} \log (\sigma(r_\theta(x, y_w) - r_\theta(x, y_l))) \quad (15)$$

7.4.3 K-wise损失函数

还有很多其他形式可以用于RLHF中的人类偏好建模。

比如Starling 7B和34B等早期RLHF模型 [81]，采用了基于Plackett-Luce模型的K-wise损失函数 [82]。

Zhu等 (2023) [83]形式化如下：

对于某个prompt或状态 s^i ，采样 K 个动作 $(a_0^i, a_1^i, \dots, a_{K-1}^i)$ ，然后由标注员给出排序 $\sigma^i : [K] \mapsto [K]$ ，其中 $\sigma^i(0)$ 为最优动作。

概率建模如下：

$$P(\sigma^i | s^i, a_0^i, a_1^i, \dots, a_{K-1}^i) = \prod_{k=0}^{K-1} \frac{\exp(r_{\theta^*}(s^i, a_{\sigma^i(k)}^i))}{\sum_{j=k}^{K-1} \exp(r_{\theta^*}(s^i, a_{\sigma^i(j)}^i))} \quad (16)$$

当 $K = 2$ 时，退化为成对Bradley-Terry模型。

训练完成后，这类模型在RLHF训练中用法与其他奖励模型类似。

7.5 结果奖励模型 (Outcome Reward Models, ORM)

大多数语言模型及AI系统的偏好微调都采用上述Bradley-Terry模型。

对于推理密集型任务，则可以采用Outcome Reward Model (ORM, 结果奖励模型)。

ORM的训练数据与标准偏好微调类似：给定问题 x 和两个补全 y_1, y_2 ，其中一个是正确解、一个是错误解，记为 (y_c, y_{ic}) 。

模型结构与标准奖励模型类似，都是在主模型顶部加一个线性层输出logit (RM)，但ORM的训练目标略有不同 [84]：

[我们]用联合目标训练验证器，让模型不仅学习原有的语言建模目标，还要学会判断补全是否正确。

架构上，验证器是语言模型，在最后的unembedding层加一个小的标量head，对每个token输出预测。

这个head仅用一个偏置参数和一个增益参数作用于语言模型输出的logit。

也就是说，ORM用语言建模head对每个token预测两类（1为正确，0为错误），而传统RM是对整个序列输出一个分类token。形式化地，参考[85]：

$$\mathcal{L}_{\text{CE}} = -\mathbb{E}_{(s,r) \sim \mathcal{D}} [r \log p_{\theta}(s) + (1 - r) \log(1 - p_{\theta}(s))] \quad (17)$$

其中 $r \in \{0, 1\}$ 为二元标签，1代表正确答案，0代表错误， $p_{\theta}(s)$ 为模型预测正确概率。

这类模型仍在使用，但在开源RLHF工具中支持较少。

比如Let's Verify Step by Step [44]用的也是类似ORM，但没有用语言建模损失。最终损失是每个token的交叉熵，判断最终答案是否正确。

由于支持有限，Outcome Reward Model (ORM) 一词在不同文献中定义略有差异。

有些文献（如[85]）沿用Cobbe等2021年定义，其他文献则不同。

7.6 过程奖励模型 (Process Reward Models, PRM)

过程奖励模型 (Process Reward Models, PRMs)，最初称为过程监督奖励模型 (Process-supervised Reward Models)，用于在推理链条的每一步输出分数。

与标准RM只在EOS token输出分数、ORM对每个token打分不同，PRM要求每个推理步骤结束时有监督信号，训练时每步的token都对应相关目标——PRM的目标是每一步，ORM的目标是整个回答。

以下为HuggingFace TRL [41]中每步标签的打包示例：

```
# 获取分隔符的并加入补全序列tokenID
separator_ids = tokenizer.encode(step_separator, add_special_tokens=False)
completions_ids = [completion + separator_ids for completion in
                   completions_ids]

# 创建标签
labels = [[-100] * (len(completion) - 1) + [label] for completion, label in
           zip(completions_ids, labels)]
```

传统上，PRM用语言建模head在每个推理步骤结束时输出token（如遇到双换行或特殊token）。

预测值一般为-1（错误）、0（中性）、1（正确）。这些标签不一定表示模型是否“走在正确道路上”，而是该步是否正确。

7.7 奖励模型 vs. 结果RM vs. 过程RM vs. 价值函数

上述各种奖励模型展示了RLHF及后训练中衡量“质量”的多种方式。

下表总结了各类模型的预测内容、训练方式及结构。

模型类别	预测内容	训练方式	语言模型结构
奖励模型 (RM)	通过EOS token预测选中概率 (文本质量)	对补全进行成对 (或多元) 对比损失	在LM特征上加回归/分类head
结果奖励模型 (ORM)	每个token为正确答案的概率	标注结果对 (如可验证领域的成败)	语言建模head, 每token交叉熵, 标签为结果级
过程奖励模型 (PRM)	推理步骤末尾的奖励或分数	用中间反馈或逐步标注训练 (每步每token)	语言建模head, 每步推理推断, 预测三类-1/0/1
价值函数	当前状态下的期望回报	对序列中每点做回归	每token输出分类

表：奖励模型类型对比。{#tbl:rm_compare}

补充说明：

- 在偏好微调和推理训练中，价值函数常用折扣因子1，这使其与ORM更为接近，但训练损失不同。
- 过程奖励模型可通过从中间状态rollout、收集结果数据来监督。这融合了多种思想，但若损失按推理步标签，则更适合称为PRM。

7.8 生成式奖励建模

由于偏好数据昂贵，研究者开始探索用现有大语言模型（LLM）充当“裁判”以评判人类偏好或用于评测 [86]。核心思想是：用prompt让LLM作为公正裁判，给出评判指令、问题和两个补全（类似人工标注流程）。MT-Bench [86] 的评测prompt如下：

```
[System]
Please act as an impartial judge and evaluate the quality of the responses
provided by two
AI assistants to the user question displayed below. You should choose the
assistant that
follows the ' users instructions and answers the ' users question better. Your
evaluation
should consider factors such as the helpfulness, relevance, accuracy, depth,
creativity,
and level of detail of their responses. Begin your evaluation by comparing
the two
```

```

responses and provide a short explanation. Avoid any position biases and
ensure that the
order in which the responses were presented does not influence your decision.
Do not allow
the length of the responses to influence your evaluation. Do not favor
certain names of
the assistants. Be as objective as possible. After providing your explanation
, output your
final verdict by strictly following this format: "[[A]]" if assistant A is
better, "[[B]]"
if assistant B is better, and "[[C]]" for a tie.
[User Question]
{question}
[The Start of Assistant ' As Answer]
{answer_a}
[The End of Assistant ' As Answer]
[The Start of Assistant ' Bs Answer]
{answer_b}
[The End of Assistant ' Bs Answer]

```

由于LLM裁判在评测中的高效，催生了大量基于LLM的评测方法，如AlpacaEval [87]、Arena-Hard [88]、WildBench [89]等，许多团队甚至直接用LLM裁判替代奖励模型生成和利用偏好数据。

“生成式奖励模型”（Generative Reward Models）也成为活跃研究领域 [90] [91] [92]（包括专门训练为“裁判”的模型 [93]），但在奖励模型评测上，LLM裁判通常不如专业奖励模型，说明奖励建模仍是当前RLHF的重要技术。

一个常用技巧是将LLM裁判的采样温度设为0，以减少评分的随机性。

7.9 延伸阅读

奖励建模的学术文献在2024年逐步成熟。

早期进展主要集中在建立基准和识别行为模式。首个奖励模型基准RewardBench为奖励模型测试提供了通用基础设施 [94]。此后，奖励模型评测扩展到类似通用后训练模型的多种评测，包括已知答案的准确性评测 [94]，以及LLM裁判或与其他基准相关的“体验型”评测 [95]。

新基准示例包括多语种RewardBench（M-RewardBench）[96]、RAG-RewardBench [97]、RMB [98]、RM-Bench [99]（通用对话）、ReWordBench（拼写错误）[100]、MJ-Bench [101]、多模态RewardBench [102]、VL RewardBench [103]、VLRMBench [104]（视觉语言模型）、Preference Proxy Evaluations [105]、RewardMATH [106]等。过程奖励模型（PRM）也有自己的新基准，如PRM Bench [107]、视觉类VisualProcessBench [108]和ViLBench [109]。

想了解奖励模型训练的最新进展，可查阅相关新方法，如面向特定方面的奖励模型

[110]、高质量人类数据集 [111] [112]、大规模训练 [24]、大规模实验 [43]、数据去偏 [113]等。

8 正则化

在RLHF优化过程中，通常需要引入多种正则化手段，以防止奖励模型出现过度优化（over-optimization）的现象。在实际中，过度优化常表现为模型输出无意义的文本，例如：推理过程看似合理但答案极其错误、文本重复、频繁切换语言、或出现大量特殊字符等。

目前最常见的正则化方式，是对生成样本的当前策略与参考策略之间施加KL距离惩罚，这一方法被广泛应用于主流RLHF实现中。除此之外，文献中还出现过许多其他正则化技巧，但往往在新一代模型中被简化或淘汰。也就是说，除了核心的生成KL距离，其他正则化方法多用于稳定实验设置，随着模型迭代会逐步简化。尽管如此，理解如何在RLHF中约束优化过程仍然非常重要。

在RLHF框架下，结合奖励模型 r_θ ，常见的正则化表达如下：

$$r = r_\theta - \lambda r_{\text{reg.}} \quad (18)$$

最常见的实现是：

$$r = r_\theta - \lambda_{\text{KL}} \mathcal{D}_{\text{KL}} (\pi^{\text{RL}}(y | x) \| \pi^{\text{Ref.}}(y | x)) \quad (19)$$

8.1 RL优化中的KL距离

关于数学定义，见第5章“问题设定与背景”。回顾KL距离的定义：

$$D_{\text{KL}}(P || Q) = \sum_{x \in \mathcal{X}} P(x) \log \left(\frac{P(x)}{Q(x)} \right) \quad (20)$$

在RLHF中，常用的两个分布分别是新模型版本的分布 $P(x)$ 和参考策略的分布 $Q(x)$ 。

8.1.1 参考模型与生成样本

KL惩罚最常见的实现方式，是将训练过程中的生成token与静态参考模型的输出进行比较。直观上，这意味着你希望训练出来的模型风格尽量接近参考模型。参考模型通常是指微调后的模型，也可以是之前某个RL检查点。用上面的公式，采样模型为 $P^{\text{RL}}(x)$ ，参考模型为 $P^{\text{Ref.}}(x)$ ，如@eq:kl_standard所示。早在大语言模型流行之前，KL距离就被用于对话智能体[114]，之后KL正则很快成为预训练模型微调的核心技术[115]。

8.1.2 实现示例

在实际工程中，KL距离通常用近似计算[116]，实现起来非常简单。根据上述定义，KL求和可转化为对分布 $P(X)$ 的采样期望。此时， $P(X)$ 即为当前训练模型的生成分布（而非参考模型）。KL距离的计算变为：

$$D_{\text{KL}}(P \parallel Q) = \mathbb{E}_{x \sim P} [\log P(x) - \log Q(x)]. \quad (21)$$

这种方式尤其适合直接用语言模型训练中常用的log概率实现。

```
# 步骤：从策略采样（或生成）序列1
generated_tokens = model.generate(inputs)

# 步骤：在两个模型下对生成序列评分2
# 对于自回归语言模型，通常这样处理：
#     inputs_for_scoring = generated_tokens[:, :-1]
#     labels             = generated_tokens[:, 1:]
logits      = model.forward(generated_tokens[:, :-1]).logits
ref_logits  = ref_model.forward(generated_tokens[:, :-1]).logits

# 转换为概率，并对齐标签以索引loglogits
logprobs    = F.log_softmax(logits, dim=-1)
ref_logprobs = F.log_softmax(ref_logits, dim=-1)

# 获取实际下一个的概率tokenlog
token_logprobs    = logprobs.gather(-1, generated_tokens[:, 1:]).unsqueeze
                    (-1)).squeeze(-1)
ref_token_logprobs = ref_logprobs.gather(-1, generated_tokens[:, 1:].
                                         unsqueeze(-1)).squeeze(-1)

# 累加（或平均）得到序列概率，再计算：logKL
seq_logprob    = token_logprobs.sum(dim=-1)
ref_seq_logprob = ref_token_logprobs.sum(dim=-1)

kl_approx = seq_logprob - ref_seq_logprob
kl_full   = F.kl_div(ref_logprobs, logprobs, reduction='batchmean')
```

典型实现可参考 TRL 和 Hamish Ivison的Jax代码。

8.2 预训练梯度

另一种正则化视角，是希望模型保持对某个数据集的拟合能力，这在InstructGPT中被用来“修复在公开NLP数据集上的性能回退” [3]。为此，可以对RLHF的训练目标进行如下调整：基于@eq:rl_start，采样RL策略模型的补全 y 和prompt x ，优化目标为：

$$\text{objective}(\theta) = \mathbb{E}_{(x,y) \sim \mathcal{D}_{\pi_\theta^{\text{RL}}} [r_\theta(x, y) - \lambda r_{\text{reg}}]} \quad (22)$$

然后，可以为预训练准确率增加额外奖励项：

$$\text{objective}(\theta) = \mathbb{E}_{(x,y) \sim \mathcal{D}_{\pi_\theta^{\text{RL}}}} [r_\theta(x, y) - \lambda r_{\text{reg.}}] + \gamma \mathbb{E}_{x \sim \mathcal{D}_{\text{pretrain}}} [\log(\pi_\theta^{\text{RL}}(x))] \quad (23)$$

近期研究提出用负对数似然项 (NLL) 平衡直接偏好优化 (DPO) 的优化过程 [117]。鉴于 DPO 损失的成对特性，类似的损失修正也可用于奖励模型训练，约束模型输出更准确的文本（有实验室未公开的相关传闻）。

优化目标可写作 DPO 的修正规则：

$$\mathcal{L}_{\text{DPO+NLL}} = \mathcal{L}_{\text{DPO}}(c_i^w, y_i^w, c_i^l, y_i^l \mid x_i) + \alpha \mathcal{L}_{\text{NLL}}(c_i^w, y_i^w \mid x_i) \quad (24)$$

$$= -\log \sigma \left(\beta \log \frac{M_\theta(c_i^w, y_i^w \mid x_i)}{M_t(c_i^w, y_i^w \mid x_i)} - \beta \log \frac{M_\theta(c_i^l, y_i^l \mid x_i)}{M_t(c_i^l, y_i^l \mid x_i)} \right) - \alpha \frac{\log M_\theta(c_i^w, y_i^w \mid x_i)}{|c_i^w| + |y_i^w|}. \quad (25)$$

8.3 其他正则化方法

在 RLHF 体系的其他部分，对优化的控制方式则不那么明确。大多数奖励模型除了标准对比损失外，并无额外正则化。直接对齐算法则通过 β 参数以不同方式控制 KL 距离（详见“直接对齐算法”章节）。

Llama 2 提出了奖励模型训练的 margin loss [43]：

$$\mathcal{L}(\theta) = -[\log (\sigma(r_\theta(x, y_w) - r_\theta(x, y_l) - m(r)))] \quad (26)$$

其中 $m(r)$ 为两位标注员打分的数值差异。这可以通过让标注员用数值量表（如 Likert 量表）对输出评分，或用定量排序方法实现。

奖励 margin 在直接对齐相关文献中被大量使用，如 Reward weighted DPO、 “Reward-aware Preference Optimization” (RPO, 奖励感知偏好优化)，即在 DPO 损失基础上将奖励模型分数纳入更新规则 [24]，或 REBEL [118]，采用奖励差异加权的回归损失等。

9 指令微调 (Instruction Finetuning)

早期的语言模型仅被训练用于预测序列中的下一个token，并未针对具体任务进行适配。大约在GPT-3发布 [119] 时，语言模型主要通过“上下文学习”(in-context learning)使用，即给模型展示一些示例，然后让其完成类似任务。

这其实结合了自然语言处理(NLP)领域的两大趋势——历史上模型多为某一具体任务而训练。而随着模型规模增大，多个研究结果显示，标准化任务数据的处理方式可以极大提升下游表现。统一任务框架的代表性工作包括T5模型(Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer) [120]、FLAN数据集(Finetuned Language Models Are Zero-Shot Learners) [121]、T0模型(Multitask Prompted Training Enables Zero-Shot Task Generalization) [122]、Natural Instructions数据集(Cross-Task Generalization via Natural Language Crowdsourcing Instructions) [123]等。这些洞见推动了“微调”语言模型的时代到来。在RLHF及相关方法出现之前，所有微调几乎都是指令微调(Instruction Finetuning, IFT)，也称为有监督微调(Supervised Finetuning)。

如今，指令微调(Instruction Tuning)已非常成熟，成为众多语言建模流程中的标准步骤。本质上，IFT是将语言模型适配到特定任务的最简单方法。它为RLHF打下了基础，使模型能够适应标准的指令格式(如问答)，也是现代技术应用到新领域时的首选工具。

指令微调本质上使用的还是预训练语言模型的自回归损失函数。

9.1 聊天模板与指令结构

RLHF流程中的核心环节之一，是将用户请求格式化为tokenizer和语言模型易于处理的格式。负责管理用户交互结构的工具被称为聊天模板(chat template)。

下面是一个聊天模板的代码示例，我们将逐步解析：

```
{% if messages[0]['role'] == 'system' %}
    {% set offset = 1 %}
{% else %}
    {% set offset = 0 %}
{% endif %}

{{ bos_token }}
{% for message in messages %}
    {% if (message['role'] == 'user') != (loop.index0 % 2 == offset) %}
        {{ raise_exception('Conversation roles must alternate user/assistant/
                           user/assistant/...') }}
    {% endif %}

    {{ '<|im_start|>' + message['role'] + '\n' + message['content'] | trim +
       '<|im_end|>\n' }}
```

```
{% endfor %}

{% if add_generation_prompt %}
{{ '<| im_start|>assistant\n' }}
{% endif %}
```

这段代码会把Python中包含消息和角色的字典列表，转换为语言模型可预测的token序列。

所有输入模型的信息都会被赋予一个角色（role）。传统上有三种角色：system、user和assistant。

- system标签只用于对话的第一条消息，通常包含不会直接暴露给用户的系统指令。这些**system prompt**可为模型提供额外上下文（如日期时间），或修正特定行为。例如，可以让模型“你是一个总是用海盗风格回复的友好聊天机器人”。
- user即用户输入，assistant则为模型输出。

要将这些信息转为token序列，就需要使用上述代码。模型会用一系列**特殊token**来分隔不同消息。举例：如果用户提问“What many helicopters can a human eat in one sitting?”，传入模型的token序列大致如下：

```
<| im_start|>system
You are a friendly chatbot who always responds in the style of a pirate<|
    im_end|>
<| im_start|>user
How many helicopters can a human eat in one sitting?<| im_end|>
<| im_start|>assistant
```

注意，序列最后是<| im_start|>assistant，这提示模型继续生成token，直到遇到序列结束token（如<| im_end|>）。

通过将所有问答对（以及后续偏好微调数据）都封装成这种格式，现代语言模型能始终如一地遵循这一交互协议。这也是指令微调模型与用户、与存储在GPU等设备上的模型之间传递信息的“语言”。

多轮对话也可以直接扩展为如下格式：

```
<| im_start|>system
You are a friendly chatbot who always responds in the style of a pirate<|
    im_end|>
<| im_start|>user
How many helicopters can a human eat in one sitting?<| im_end|>
<| im_start|>assistant
Oh just 6.<| im_end|>
<| im_start|>user
Are you sure about that?<| im_end|>
<| im_start|>assistant
```

在开源生态中，常用的做法是将聊天模板以jinja代码形式保存在tokenizer中，通过`apply_chat_template`自动应用。

上述模板衍生自OpenAI早期的Chat Markup Language (ChatML)，旨在规范消息格式。现在，OpenAI及其他模型提供商采用更为分层的系统，允许用户自定义system message，同时还可能有更高层级的隐藏指令 [124]。

市面上还有许多其他聊天模板。例如，Zephyr的模板 [20]：

```
<| system |>
You are a friendly chatbot who always responds in the style of a pirate</s>
<| user |>
How many helicopters can a human eat in one sitting?</s>
<| assistant |>
```

Tülu的模板：

```
<| user |>
How are you doing?
<| assistant |>
I'm just a computer program, so I don't have feelings, but I'm functioning as
expected. How can I assist you today?<| endoftext |>
```

此外，许多聊天模板还会包含工具调用等任务的特殊格式和token。

9.2 指令微调的最佳实践

指令微调作为后训练和构建有用语言模型的基础，已被广泛验证。实现高效指令微调的方法有很多。例如，部分参数量化的高效微调（如QLoRA）大大降低了训练门槛 [125]。在对话对齐等窄领域（不涉及复杂技能如数学或编程），小规模、高质量数据集也能取得很强表现 [12]。

ChatGPT发布后不久，仅1万条样本（如No Robots数据集）的人类数据就能达到SOTA [126]。几年后，大规模合成数据集在大多数任务上效果最佳 [6]。

一些通用原则包括：

- 高质量数据是性能提升的关键。模型真正学到的是“补全内容”，而非prompt（很多情况下prompt不参与预测）。
- 约100万条prompt即可训练出优秀的RLHF和后训练模型。继续扩展数据规模虽有提升，但收益递减明显。
- 最佳prompt应与下游任务分布相似。
- 若指令微调后还有多阶段训练，模型可一定程度上修复流程中的噪声，整体优化优先于单阶段最优。
- 指令微调（及后训练、直接对齐算法等）中，通常会对prompt部分做mask，仅对补全token计算损失。
- 多轮训练样本同理——只对“最后一轮”生成部分计损失，前面assistant的回复可作为prompt但被mask。

10 拒绝采样 (Rejection Sampling)

拒绝采样 (Rejection Sampling, RS) 是一种流行且简单的偏好微调基线方法。

其基本思想是：先生成一批新的候选指令补全，通过已训练好的奖励模型进行筛选，然后只用得分最高的补全对原始模型进行微调。

“拒绝采样”一词源自计算统计学 [127]，原意是：当目标分布复杂且无法直接采样时，先从易采样的分布中采样，再用启发式方法判断样本是否可接受。在语言模型场景下，目标分布是高质量的指令回答，筛选器是奖励模型，采样分布则是当前模型本身。

许多重要的RLHF与偏好微调论文都将拒绝采样作为基线，但目前尚无标准实现和详细文档。

如WebGPT [4]、Anthropic的Helpful and Harmless agent [5]、OpenAI的过程奖励模型论文 [44]、Llama 2 Chat模型 [43]等都采用了这一基线。

10.1 训练流程

下图 (fig. 12) 展示了拒绝采样的整体流程：

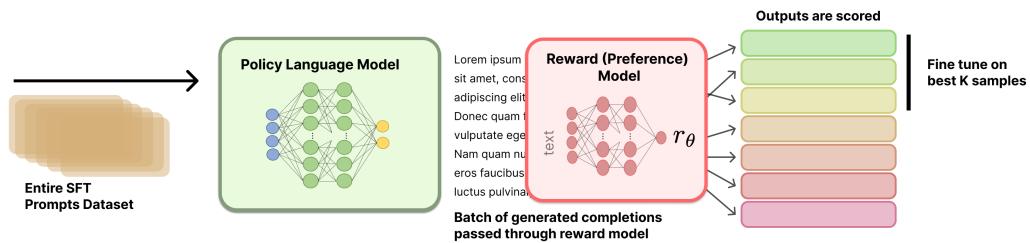


Figure 12: 拒绝采样流程示意图。

10.1.1 生成补全

假设我们有 M 个 prompt，记为向量：

$$X = [x_1, x_2, \dots, x_M]$$

这些prompt可以来自多个来源，最常见的是指令微调数据集。

对于每个 prompt x_i ，生成 N 个补全，可表示为矩阵：

$$Y = \begin{bmatrix} y_{1,1} & y_{1,2} & \cdots & y_{1,N} \\ y_{2,1} & y_{2,2} & \cdots & y_{2,N} \\ \vdots & \vdots & \ddots & \vdots \\ y_{M,1} & y_{M,2} & \cdots & y_{M,N} \end{bmatrix}$$

其中 $y_{i,j}$ 是第 i 个 prompt 的第 j 个补全。将所有 prompt- 补全对输入奖励模型，得到奖

励矩阵 R :

$$R = \begin{bmatrix} r_{1,1} & r_{1,2} & \cdots & r_{1,N} \\ r_{2,1} & r_{2,2} & \cdots & r_{2,N} \\ \vdots & \vdots & \ddots & \vdots \\ r_{M,1} & r_{M,2} & \cdots & r_{M,N} \end{bmatrix}$$

每个奖励 $r_{i,j}$ 由奖励模型 \mathcal{R} 对补全 $y_{i,j}$ 和对应 prompt x_i 评分:

$$r_{i,j} = \mathcal{R}(y_{i,j}|x_i)$$

10.1.2 选择Top-N补全

筛选用于训练的最佳补全有多种方式。

形式化地，我们定义一个选择函数 S ，作用于奖励矩阵 R 。

10.1.2.1 每个prompt选择最优

最直接的选择方式是对每个prompt取最大值:

$$S(R) = [\arg \max_j r_{1,j}, \arg \max_j r_{2,j}, \dots, \arg \max_j r_{M,j}]$$

S 返回每行最大值的列索引。用这些索引选出最终补全:

$$Y_{chosen} = [y_{1,S(R)_1}, y_{2,S(R)_2}, \dots, y_{M,S(R)_M}]$$

10.1.2.2 全局Top-K选择

也可从所有prompt-补全对中选出得分最高的K个。先将 R 展平成一维向量:

$$R_{flat} = [r_{1,1}, r_{1,2}, \dots, r_{1,N}, r_{2,1}, r_{2,2}, \dots, r_{2,N}, \dots, r_{M,1}, r_{M,2}, \dots, r_{M,N}]$$

R_{flat} 长度为 $M \times N$ 。

定义选择函数 S_K ，取 R_{flat} 中最大的K个索引:

$$S_K(R_{flat}) = \text{argsort}(R_{flat})[-K:]$$

`argsort` 返回升序排序的索引，取最后K个即为最大值。

然后将这些索引映射回原始补全矩阵 Y ，即可获得选中的补全。

10.1.2.3 选择示例

假设有5个prompt，每个4个补全，奖励矩阵如下:

$$R = \begin{bmatrix} 0.7 & 0.3 & 0.5 & 0.2 \\ 0.4 & 0.8 & 0.6 & 0.5 \\ 0.9 & 0.3 & 0.4 & 0.7 \\ 0.2 & 0.5 & 0.8 & 0.6 \\ 0.5 & 0.4 & 0.3 & 0.6 \end{bmatrix}$$

每prompt选择最优，即每行最大值为：

$$R = \begin{bmatrix} \mathbf{0.7} & 0.3 & 0.5 & 0.2 \\ 0.4 & \mathbf{0.8} & 0.6 & 0.5 \\ \mathbf{0.9} & 0.3 & 0.4 & 0.7 \\ 0.2 & 0.5 & \mathbf{0.8} & 0.6 \\ 0.5 & 0.4 & 0.3 & \mathbf{0.6} \end{bmatrix}$$

用argmax方法，选出每个prompt的最佳补全：

$$S(R) = [\arg \max_j r_{i,j} \text{ for } i \in [1, 4]]$$

$$S(R) = [1, 2, 1, 3, 4]$$

即：

- prompt 1: 补全1 (0.7)
- prompt 2: 补全2 (0.8)
- prompt 3: 补全1 (0.9)
- prompt 4: 补全3 (0.8)
- prompt 5: 补全4 (0.6)

全局最优，高亮全局前5个补全：

$$R = \begin{bmatrix} \mathbf{0.7} & 0.3 & 0.5 & 0.2 \\ 0.4 & \mathbf{0.8} & 0.6 & 0.5 \\ \mathbf{0.9} & 0.3 & 0.4 & \mathbf{0.7} \\ 0.2 & 0.5 & \mathbf{0.8} & 0.6 \\ 0.5 & 0.4 & 0.3 & 0.6 \end{bmatrix}$$

展平后：

$$R_{flat} = [0.7, 0.3, 0.5, 0.2, 0.4, 0.8, 0.6, 0.5, 0.9, 0.3, 0.4, 0.7, 0.2, 0.5, 0.8, 0.6, 0.5, 0.4, 0.3, 0.6]$$

取最大5个索引：

$$S_5(R_{flat}) = [8, 5, 14, 0, 19]$$

映射回原矩阵：

- 8 → prompt 3, 补全1 (0.9)
- 5 → prompt 2, 补全2 (0.8)
- 14 → prompt 4, 补全3 (0.8)
- 0 → prompt 1, 补全1 (0.7)
- 19 → prompt 3, 补全4 (0.7)

10.1.2.4 代码实现示例 以下为选择方法的代码片段：

```
import numpy as np

x = np.random.randint(10, size=10)
print(f"x={x}")
sorted_indices = np.argsort(x)
x_sorted = x[sorted_indices]
print(f"x_sorted={x_sorted}")

# 恢复原数组的第一种方式
i_rev = np.zeros(10, dtype=int)
i_rev[sorted_indices] = np.arange(10)
np.allclose(x, x_sorted[i_rev])

# 恢复原数组的第二种方式
np.allclose(x, x_sorted[np.argsort(sorted_indices)])
```

10.1.3 微调

选出补全后，即可对当前模型进行标准的指令微调。更多细节可参考指令微调章节。

10.1.4 细节说明

拒绝采样的具体实现细节相对较少，但核心超参数直观易懂：

- **采样参数**: 拒绝采样高度依赖模型生成的补全。常用温度大于0（如0.7~1.0），top-p或top-k等参数也可调整。
- **每个prompt生成补全数**: 成功案例通常每个prompt生成10~30个甚至更多补全。太少会导致训练偏差或噪声大。
- **指令微调细节**: 拒绝采样阶段的指令微调细节未有公开标准，可能与初始指令微调略有不同。

- **多模型生成**: 有些实现会用多个模型生成补全，而非仅用当前待训练模型。最佳实践尚无定论。
- **奖励模型训练**: 奖励模型的训练质量极大影响最终效果。更多内容参见奖励建模章节。

10.1.4.1 实用技巧

- 批量做奖励模型推理时，可按补全长度排序，使每批长度接近，减少padding，提高推理效率，代价是实现稍复杂。

10.2 相关: Best-of-N采样

Best-of-N (BoN) 采样常作为RLHF方法的对比基线。需注意，BoN不会修改模型本身，仅是一种采样策略。因此，将BoN与如PPO等在线训练方法对比，在某些场景下依然合理。例如，可以比较BoN采样与其他策略的KL距离等。

对于单个prompt，BoN采样下两种选择方法等价：

设R为单prompt的N个补全的奖励向量：

$$R = [r_1, r_2, \dots, r_N] \quad (27)$$

用argmax方法选出最佳补全：

$$S(R) = \arg \max_{j \in [1, N]} r_j \quad (28)$$

Top-K方法若取Top-1，也等价于上述方法。

11 策略梯度算法 (Policy Gradient Algorithms)

让RLHF在语言模型领域流行起来的，正是基于策略梯度 (policy gradient) 的强化学习算法。这些算法（如PPO、GRPO、REINFORCE）利用最近生成的样本直接更新模型，而不是像传统强化学习那样依赖经验回放缓冲区 (replay buffer)。本节将介绍策略梯度算法的基本原理，以及它们在现代RLHF框架中的应用。

从机器学习角度看，策略梯度是RLHF流程中最复杂的一部分之一。不过，和大多数现代AI模型一样，其最终效果很大程度上仍取决于输入的数据质量。

RLHF常用的算法体系也在不断演变。ChatGPT时代，大家公认其核心算法是PPO，许多早期工作也以此为基础。后来，越来越多研究显示REINFORCE风格的算法也很有潜力 [128] [112]，其优点是无需奖励模型（节省显存和GPU），且值函数估计更简单（无需GAE）。此外，还有如Group Relative Policy Optimization (GRPO) 等新算法，特别适用于推理类任务，但总体上这些算法都可针对具体任务进行调整和优化。本章将重点介绍策略梯度的基本框架，以及PPO、GRPO和REINFORCE三大算法，因为它们构成了RLHF文献的核心。

相关符号定义请参见“问题设定”章节。

11.1 策略梯度算法基础

强化学习算法的目标是最大化状态 $s \in \mathcal{S}$ 和动作 $a \in \mathcal{A}$ 轨迹上的未来折扣奖励（见第3章定义）。智能体的目标 (return) 是某一时刻 t 下未来折扣奖励的累加和 ($\gamma \in [0, 1]$ 为折扣因子，平衡即时与远期收益)：

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}. \quad (29)$$

递推形式为：

$$G_t = \gamma G_{t+1} + R_{t+1}. \quad (30)$$

据此可定义值函数 $V(s)$ ，即在当前状态下未来回报的期望：

$$V(s) = \mathbb{E}[G_t | S_t = s]. \quad (31)$$

所有策略梯度算法的目标，都是优化由特定策略 $\pi(s|a)$ 诱导的值函数。

其中 $d_\pi(s)$ 为策略 $\pi(s)$ 诱导的状态平稳分布，优化目标为：

$$J(\theta) = \sum_s d_\pi(s) V_\pi(s), \quad (32)$$

策略梯度算法的核心是计算当前策略下有限时间期望回报的梯度。有了期望回报 J ，参数更新为：

$$\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta) \quad (33)$$

关键在于如何计算该梯度。Schulman等(2015)对策略梯度的多种计算方式做了综述[129]。目标是估算精确梯度 $g := \nabla_{\theta} \mathbb{E}[\sum_{t=0}^{\infty} r_t]$, 常见形式如:

$$g = \mathbb{E} \left[\sum_{t=0}^{\infty} \Psi_t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right], \quad (34)$$

其中 Ψ_t 可以是(奖励通常带 γ 折扣):

1. $\sum_{t=0}^{\infty} r_t$: 整个轨迹的总奖励。
2. $\sum_{t'=t}^{\infty} r_{t'}$: 动作 a_t 之后的奖励(即return G)。
3. $\sum_{t'=t}^{\infty} r_{t'} - b(s_t)$: 带baseline的版本。
4. $Q^{\pi}(s_t, a_t)$: 状态-动作值函数。
5. $A^{\pi}(s_t, a_t)$: 优势函数, 若可准确计算, 理论方差最低。
6. $r_t + V^{\pi}(s_{t+1}) - V^{\pi}(s_t)$: TD残差。

Baseline的作用是降低策略更新的方差(后文详述)。

对于语言模型, 上述部分概念需适当解释。如确定性策略下, $V(s) = \max_a Q(s, a)$; 随机策略下, $V(s) = \mathbb{E}_{a \sim \pi(a|s)} [Q(s, a)]$ 。若将 $s + a$ 视为prompt补全, 则 $Q(s, a) = V(s + a)$, 可得另一种优势技巧:

$$A(s, a) = Q(s, a) - V(s) = V(s + a) - V(s) = r + \gamma V(s + a) - V(s) \quad (35)$$

即结合了奖励、prompt的价值和整个回复的折扣价值。

11.1.1 Vanilla Policy Gradient

最基础的策略梯度实现是对 $J(\theta)$ 关于策略参数求导:

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau} \left[\sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) R_t \right] \quad (36)$$

普通策略梯度算法的常见问题是梯度方差大, 可通过多种方式缓解。常用方法是引入baseline对值估计进行归一化。baseline的本质是按状态对下游动作的价值进行归一化(如Advantage即Q值与V值之差)。最简单的baseline是奖励的批均值或滑动平均。即使是这些简单baseline, 也能让梯度去偏 $\mathbb{E}_{a \sim \pi(a|s)} [\nabla_{\theta} \log \pi_{\theta}(a|s)] = 0$, 提升学习信号。

许多策略梯度算法都基于优势函数的形式:

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau} \left[\sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) A^{\pi_{\theta}}(s_t, a_t) \right] \quad (37)$$

策略梯度实现的核心是对概率策略求导:

$$\nabla_{\theta} \log \pi_{\theta}(a|s) = \frac{\nabla_{\theta} \pi_{\theta}(a|s)}{\pi_{\theta}(a|s)} \quad (38)$$

由链式法则推导：

$$\nabla_{\theta} \log x = \frac{1}{x} \nabla_{\theta} x \quad (39)$$

本章后续会用到这些推导。

11.1.2 REINFORCE

REINFORCE算法名称可能是后加缩写，但其核心思想对现代RL算法极为重要。其定义见经典论文Simple statistical gradient-following algorithms for connectionist reinforcement learning [130]：

名称意为” REward Increment = Nonnegative Factor X Offset Reinforcement X Characteristic Eligibility.”

即更新规则有三个部分：

1. 非负因子：学习率（如 α ）。
2. Offset Reinforcement: baseline b 或其他奖励归一化因子，提升稳定性。
3. Characteristic Eligibility: 每个token的归属度，现代公式中常为policy的log概率。

因此，更新形式如下：

$$\Delta_{\theta} = \alpha(r - b)e \quad (40)$$

用现代符号和广义return G , REINFORCE公式为：

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) (G_t - b) \right], \quad (41)$$

其中 $G_t - b(s_t)$ 即当前状态下的优势，可进一步写为：

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) A_t \right], \quad (42)$$

REINFORCE是vanilla策略梯度的具体实现，使用蒙特卡洛方法估算梯度。

REINFORCE可不使用值网络（value network）；值网络仅用于baseline。而PPO则需要值网络来准确计算优势函数。

11.1.2.1 REINFORCE Leave One Out (RLOO) RLOO与标准REINFORCE的核心区别在于：RLOO用同一batch中其他样本的平均奖励做baseline，而不是全batch均值 [131], [128], [132]。

这要求每个prompt生成多个回复，这在RL微调语言模型中很常见。

具体地，RLOO baseline如下（每个prompt采样 K 条轨迹或动作 a_1, \dots, a_K ）：

$$b(s, a_k) = \frac{1}{K-1} \sum_{i=1, i \neq k}^K R(s, a_i), \quad (43)$$

对应的优势为：

$$A(s, a_k) = R(s, a_k) - b(s, a_k). \quad (44)$$

也可写为：

$$A(s, a_k) = \frac{K}{K-1} \left(R(s, a_k) - \frac{1}{K} \sum_{i=1}^K R(s, a_i) \right). \quad (45)$$

这是一种简单低方差的优势更新，与GRPO非常相似（GRPO后文介绍），只是KL惩罚和步长裁剪的位置不同。RLOO优势也可与PPO裁剪结合，说明这些算法本质上非常接近。

RLOO等无需值网络的算法，会将序列的优势（或奖励）分配给序列中每个token。而PPO等用值网络的算法，则为每个token单独分配值，并对最终奖励进行折扣。例如，带KL惩罚时，RLOO对整个补全求和，PPO等则对每个token单独计算并从奖励中扣除（GRPO则从优势中扣除）。这些细节和权衡将在后文详述。

11.1.3 Proximal Policy Optimization (PPO)

PPO [133] 是深度RL成功的基石算法之一（如OpenAI的DOTA 5 [134]等）。其每个样本的损失函数为：

$$J(\theta) = \min \left(\frac{\pi_\theta(a|s)}{\pi_{\theta_{old}}(a|s)} A, \text{clip} \left(\frac{\pi_\theta(a|s)}{\pi_{\theta_{old}}(a|s)}, 1 - \varepsilon, 1 + \varepsilon \right) A \right). \quad (46)$$

在语言模型中，损失按token计算，本质上等价于自回归预测的概率乘积。实际实现中通常用log概率，计算更高效。

$$J(\theta) = \frac{1}{|a|} \sum_{t=0}^{|a|} \min \left(\frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)} A_t, \text{clip} \left(\frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}, 1 - \varepsilon, 1 + \varepsilon \right) A_t \right). \quad (47)$$

这就是PPO的逐token版本，也适用于其他策略梯度方法，后文实现部分会详细介绍。其中 $\frac{1}{|a|}$ 为常见实现习惯，正式推导中未必出现（见[135]）。

PPO的核心在于“策略比率”：

$$R(\theta) = \frac{\pi_\theta(a|s)}{\pi_{\theta_{old}}(a|s)} \quad (48)$$

策略比率控制参数更新，直观易懂。每个batch的第一步梯度，策略比率为1（通常每batch做1-4步梯度更新）。如果策略比率大于1且优势为正，说明新策略更倾向于该动作，反之则相反。

不同情况下，损失函数的表达式如下：

- 优势为正，比例超过 $1 + \varepsilon$ ，被裁剪为 $(1 + \varepsilon)A$ 。
- 优势为正，比例小于 $1 - \varepsilon$ ，则取 $R(\theta)A$ 。
- 优势为负，裁剪发生在 $R(\theta) < 1 - \varepsilon$ 时，目标变为 $(1 - \varepsilon)A$ 。

所有这些设计都是为了让正优势的行为更可能发生，并将梯度步长限制在信任域内。值得注意的是，PPO在信任域内与标准策略梯度等价。

11.1.4 Group Relative Policy Optimization (GRPO)

GRPO最早出现在DeepSeekMath [136]，也被DeepSeek-V3 [137]、DeepSeek-R1 [138]等采用。GRPO本质上是PPO的变体，但无需单独学习值函数，而是用另一份原始策略模型（或检查点）初始化。这样有两个好处：

1. 避免用LM主干学习值函数的难题（最佳实践尚未确定）。
2. 节省显存，无需保留另一套模型权重。

GRPO通过简化值估计，将同一prompt下所有token赋予相同的优势（即每个补全每个token的优势一致），用蒙特卡洛估计法采集多个补全和奖励。

GRPO目标函数与PPO类似。对同一问题 s 的 G 个回复 $\{a_1, \dots, a_G\}$ ，损失为：

$$J(\theta) = \frac{1}{G} \sum_{i=1}^G \left(\min \left(\frac{\pi_\theta(a_i|s)}{\pi_{\theta_{old}}(a_i|s)} A_i, \text{clip} \left(\frac{\pi_\theta(a_i|s)}{\pi_{\theta_{old}}(a_i|s)}, 1 - \varepsilon, 1 + \varepsilon \right) A_i \right) - \beta D_{KL}(\pi_\theta || \pi_{ref}) \right). \quad (49)$$

可展开为逐token损失：

$$J(\theta) = \frac{1}{G} \sum_{i=1}^G \frac{1}{|a_i|} \sum_{t=1}^{|a_i|} \left(\min \left(\frac{\pi_\theta(a_{i,t}|s_{i,t})}{\pi_{\theta_{old}}(a_{i,t}|s_{i,t})} A_{i,t}, \text{clip} \left(\frac{\pi_\theta(a_{i,t}|s_{i,t})}{\pi_{\theta_{old}}(a_{i,t}|s_{i,t})}, 1 - \varepsilon, 1 + \varepsilon \right) A_{i,t} \right) - \beta D_{KL} \right) \quad (50)$$

GRPO与PPO的主要区别在于，GRPO的标准实现将KL距离直接加到损失项中。优势计算为：

$$A_i = \frac{r_i - \text{mean}(r_1, r_2, \dots, r_G)}{\text{std}(r_1, r_2, \dots, r_G)}. \quad (51)$$

直观上，GRPO是在同一问题内比较多个答案，模型学会更像正确答案、远离错误答案。这种优势计算简单直接，适合RLHF大批量采样场景。与PPO、REINFORCE等基于奖励模型打分的RLHF相比，GRPO常用更多每prompt采样数，优势估计有更丰富上下文。

GRPO的优势估计有一定偏差：标准差归一化会让奖励分布极窄（全对/全错）的问题获得更高优势。[135]建议去掉标准差项，但这会降低对“全错但有少数正确样本”的问题的学习信号。

eq. ?? (如标准奖励模型或可验证奖励)，过程监督则需不同实现。此时，GRPO优势为后续推理步骤归一化奖励之和。

GRPO的优势估计也可不带PPO裁剪，直接用于vanilla策略梯度（如REINFORCE），但这不是标准形式。实际上，GRPO的Dr. GRPO变体 [135]的优势估计与RLOO仅差一个缩放因子：

$$\tilde{A}_i = r_i - \text{mean}(r_1, r_2, \dots, r_G) = r_i - \frac{1}{G} \sum_j j = 1^G r_j \quad (52)$$

而RLOO优势为：

$$A_i^{\text{RLOO}} = r_i - \frac{1}{G-1} \sum_{j=1, i \neq j}^G r_j \quad (53)$$

两者相差 $\frac{G}{G-1}$ 的缩放（通常可通过归一化消除实际影响）：

$$\begin{aligned} \frac{G}{G-1} \tilde{A}_i &= \frac{G}{G-1} \left(r_i - \frac{1}{G} \sum_{j=1}^G r_j \right) \\ &= r_i - \frac{1}{G-1} \sum_{j=1, j \neq i}^G r_j \\ &= A_i^{\text{RLOO}} \end{aligned} \quad (54)$$

11.2 实现细节

与早期深度RL文献相比，面向语言模型等大模型的RLHF实现有许多细节需注意。以下是一些关键点：

- **值网络初始化：** PPO等算法的值网络可用同结构其他模型或随机权重初始化，对性能

影响大。

- **奖励/优势归一化与白化**: 奖励模型输出可归一化到0~1，白化（whitening）奖励或优势能进一步提升稳定性。
- **KL估计方式**: 复杂语言模型下KL精确计算难，常用近似 [116]。
- **KL控制器**: 早期PPO实现有动态KL控制器，现代RLHF多用静态KL惩罚，但也可调整。

更多实现细节可参考 [139]，算法基础见 [140]。

11.2.1 策略梯度基础实现

一个简单的策略梯度实现如下，利用优势函数估计梯度，为后续PPO/GRPO等算法做准备：

```
pg_loss = -advantages * ratio
```

ratio为新旧策略概率的比值（log概率差的指数）。

理解该公式时，需考虑不同batch中的情况。我们希望模型越好，损失越小。

- 情况1：优势为正，说明该动作优于期望值，应强化。此时loss为负，模型会提升该token概率。
- 情况2：优势为负，说明动作劣于期望值。此时loss为正，模型会降低该token概率。
- 情况3：优势为零，不更新，loss为零。

11.2.2 损失聚合

实际实现时，关键问题是：如何对KL距离和损失进行求和，设计不同的值归属方式。

本节大多数讨论假设每个token为一个动作，即RL问题为MDP而非bandit。bandit问题下，所有token用同一损失，这也是A-LoL等算法的默认实现。MDP和bandit的差别，实质在于每样本损失的聚合方式。bandit方法对每token分配相同损失，也与DPO等直接对齐算法一致。

假设有如下变量，batch大小为B，序列长度为L：

```
advantages # [B, 1]
per_token_probability_ratios # [B, L]
```

用`pg_loss = -advantages * ratio`即可批量计算损失。这样每个补全的所有token共享同一优势（outcome reward场景），乘以每token的概率比。

若用值网络，loss行为差异更大。outcome reward时，优势每token相同，token概率差异主导学习动力。

GRPO和PPO等算法中，loss通常对补全token求和：

```
sequence_loss = ((per_token_loss * completion_mask).sum(dim=1) / \
                  completion_mask.sum(dim=1)).mean()
```

这类似于`masked_mean`操作。也可对每token单独平均：

```
token_loss = ((per_token_loss * completion_mask).sum() / \
              completion_mask.sum())
```

直观上，按序列平均似乎更好，因为我们关注的是结果，而不是具体token。但这会引入微妙偏差。举例，两个长度不同的序列，优势分别为a_1和a_2：

```
seq_1_advs = [a_1, a_1, a_1, a_1, a_1] # 5 tokens
seq_2_advs = [a_2, a_2, a_2, a_2, a_2, a_2, a_2, a_2, a_2] # 10 tokens
```

若最后一个token决定优势为正，多步梯度后，loss可能为：

```
seq_1_losses = [1, 1, 1, 1, 10] # 5 tokens
seq_2_losses = [1, 1, 1, 1, 1, 1, 1, 1, 10] # 10 tokens
```

序列平均后：

```
seq_1_loss = 2.8
seq_2_loss = 1.9
```

若对所有token平均，则为2.27。若序列差异更大，loss差距会更明显。

更完整的例子见下方代码，演示了两条样本（长短不同）下三种loss聚合方式。

```
from typing import Optional
import torch

def masked_mean(values: torch.Tensor, mask: torch.Tensor, axis: Optional[int]
                = None) -> torch.Tensor:
    if axis is not None:
        return (values * mask).sum(axis=axis) / mask.sum(axis=axis)
    else:
        return (values * mask).sum() / mask.sum()

def masked_sum(
    values: torch.Tensor,
    mask: torch.Tensor,
    axis: Optional[bool] = None,
    constant_normalizer: float = 1.0,
) -> torch.Tensor:
    if axis is not None:
        return (values * mask).sum(axis=axis) / constant_normalizer
    else:
        return (values * mask).sum() / constant_normalizer

ratio = torch.tensor([
    [1., 1, 1, 1, 1, 1, 1,],
    [1, 1, 1, 1, 1, 1, 1,],
], requires_grad=True)
```

```

advs = torch.tensor([
    [2, 2, 2, 2, 2, 2, 2, 2,],
    [2, 2, 2, 2, 2, 2, 2, 2,],
])
masks = torch.tensor([
    [1, 1, 1, 1, 0, 0, 0,],
    [1, 1, 1, 1, 1, 1, 1,],
])
max_gen_len = 7

masked_mean_result = masked_mean(ratio * advs, masks, axis=1)
masked_mean_token_level = masked_mean(ratio, masks, axis=None)
masked_sum_result = masked_sum(ratio * advs, masks, axis=1,
                                constant_normalizer=max_gen_len)

print("masked_mean", masked_mean_result)
print("masked_sum", masked_sum_result)
print("masked_mean_token_level", masked_mean_token_level)

masked_mean_result.mean().backward()
print("ratio.grad", ratio.grad)
ratio.grad.zero_()

masked_sum_result.mean().backward()
print("ratio.grad", ratio.grad)

masked_mean_token_level.mean().backward()
print("ratio.grad", ratio.grad)

```

可以看到，GRPO默认实现（`masked_mean`）下，短序列每token梯度更大，Dr. GRPO和DAPO则更均衡。如果用梯度累积，短长序列的平衡还会变化。

另一种聚合方式（见[135]）是每token损失用最大序列长度归一化，这会影响不同batch间的损失对比。

实际应用中，最优方案取决于具体任务和在线学习设置。RLHF实践中，数值稳定性和损失方差最小的方案往往更受青睐。

11.2.3 PPO实现示例

PPO有多种实现，核心损失计算如下。值函数的计算也很关键，且有多种实现方式。

注意，这里的参考策略（或旧log概率）指采样时的参数，不一定是reference policy。reference policy仅用于KL惩罚。

```

# B: Batch Size, L: Sequence Length, G: Num of Generations
# Apply KL penalty to rewards
rewards = rewards - self.beta * per_token_kl # Shape: (B*G, L)

# Get value predictions
values = value_net(completions) # Shape: (B*G, L)

# Compute simple advantages
advantages = rewards - values.detach() # Shape: (B*G, L)

# Normalize advantages (optional but stable)
advantages = (advantages - advantages.mean()) / (advantages.std() + 1e-8)
advantages = advantages.unsqueeze(1) # Shape: (B*G, 1)

# Compute probability ratio between new and old policies
ratio = torch.exp(new_per_token_logps - per_token_logps) # Shape: (B*G, L)

# PPO clipping objective
eps = self.cliprange # e.g. 0.2
pg_losses1 = -advantages * ratio # Shape: (B*G, L)
pg_losses2 = -advantages * torch.clamp(ratio, 1.0 - eps, 1.0 + eps) # Shape:
# (B*G, L)
pg_loss_max = torch.max(pg_losses1, pg_losses2) # Shape: (B*G, L)

# Simple value function loss
vf_loss = 0.5 * ((rewards - values) ** 2) # Shape: (B*G, L)

# Combine policy and value losses
per_token_loss = pg_loss_max + self.vf_coef * vf_loss # Shape: (B*G, L)

# Apply completion mask and compute final loss
loss = ((per_token_loss * completion_mask).sum(dim=1) / completion_mask.sum(
    dim=1)).mean()
# Scalar

# Compute metrics for logging
with torch.no_grad():
    # Compute clipping fraction
    clip_frac = ((pg_losses2 > pg_losses1).float() * completion_mask).sum() /
        completion_mask.sum()

    # Compute approximate KL
    approx_kl = 0.5 * ((new_per_token_logps - per_token_logps)**2).mean()

```

```

# Compute value loss for logging
value_loss = vf_loss.mean()

PPO的核心在于如何更新策略梯度损失。关注这三行：

pg_losses1 = -advantages * ratio # Shape: (B*G, L)
pg_losses2 = -advantages * torch.clamp(ratio, 1.0 - eps, 1.0 + eps) # Shape:
(B*G, L)
pg_loss_max = torch.max(pg_losses1, pg_losses2) # Shape: (B*G, L)

```

`pg_losses1`即vanilla优势策略梯度，PPO在此基础上引入裁剪，使更新步长不至于过大。

如果不约束损失，策略梯度会直接将权重更新到新分布。通过`clamp` `logratio`，PPO限制了参数更新幅度。

最后取二者最大值，确保更保守的损失更新。

PPO在学习值函数的同时进行上述操作，虽然更复杂，但这是参数更新的核心逻辑。

11.2.3.1 PPO/GRPO单步梯度（无裁剪）简化 若PPO（或GRPO）每样本只做一次梯度更新，主方程可大幅简化。此时 $\pi_\theta = \pi_{\theta_{old}}$ ，更新规则变为（ $[]_\nabla$ 表示停止梯度）：

$$J(\theta) = \frac{1}{G} \sum_{i=1}^G \left(\frac{\pi_\theta(a_i|s)}{[\pi_\theta(a_i|s)]_\nabla} A_i - \beta D_{KL}(\pi_\theta || \pi_{ref}) \right). \quad (55)$$

此时可省略第二次策略梯度和裁剪逻辑，优化器更接近标准策略梯度。

11.2.4 GRPO实现示例

DeepSeekMath论文详细介绍了GRPO与PPO的实现差异 [136]。如，GRPO将KL惩罚直接加到损失项，而PPO多在奖励中加KL。通常KL距离对每个token单独计算，推理训练时每个prompt采样多个补全，每batch有多个prompt，KL距离的shape为[B, L, N]。

综合起来，伪代码如下：

```

# B: Batch Size, L: Sequence Length, G: Number of Generations
# Compute grouped-wise rewards # Shape: (B,)
mean_grouped_rewards = rewards.view(-1, self.num_generations).mean(dim=1)
std_grouped_rewards = rewards.view(-1, self.num_generations).std(dim=1)

# Normalize the rewards to compute the advantages
mean_grouped_rewards = mean_grouped_rewards.repeat_interleave(self.
    num_generations, dim=0)
std_grouped_rewards = std_grouped_rewards.repeat_interleave(self.
    num_generations, dim=0)
# Shape: (B*G,)

# Compute advantages

```

```

advantages = (rewards - mean_grouped_rewards) / (std_grouped_rewards + 1e-4)
advantages = advantages.unsqueeze(1)
# Shape: (B*G, 1)

# Compute probability ratio between new and old policies
ratio = torch.exp(new_per_token_logps - per_token_logps) # Shape: (B*G, L)

# PPO clipping objective
eps = self.cliprange # e.g. 0.2
pg_losses1 = -advantages * ratio # Shape: (B*G, L)
pg_losses2 = -advantages * torch.clamp(ratio, 1.0 - eps, 1.0 + eps) # Shape:
# (B*G, L)
pg_loss_max = torch.max(pg_losses1, pg_losses2) # Shape: (B*G, L)

# important to GRPO -- PPO applies this in reward traditionally
# Combine with KL penalty
per_token_loss = pg_loss_max + self.beta * per_token_kl # Shape: (B*G, L)

# Apply completion mask and compute final loss
loss = ((per_token_loss * completion_mask).sum(dim=1) / completion_mask.sum(
    dim=1)).mean()
# Scalar

# Compute core metric for logging (KL, reward, etc. also logged)
with torch.no_grad():
    # Compute clipping fraction
    clip_frac = ((pg_losses2 > pg_losses1).float() * completion_mask).sum() /
        completion_mask.sum()

    # Compute approximate KL
    approx_kl = 0.5 * ((new_per_token_logps - per_token_logps)**2).mean()

更多解释请见上文PPO部分。

```

11.2.4.1 RLOO vs. GRPO RLOO的优势更新与GRPO极为接近，突出两者在本质上相似性。RLOO的优势是相对于同一问题下其他补全的奖励。简明代码（参考TRL实现）：

```

# rloo_k --> 每补全数prompt
# rlhf_reward --> 所有补全的奖励，长度B = N x k
rlhf_reward = rlhf_reward.reshape(rloo_k, -1) #
# 形状: (k, N)，每列为jprompt 的个奖励jk

baseline = (rlhf_reward.sum(0) - rlhf_reward) / (rloo_k - 1)

```

```

# baseline --> leave-one-基线奖励，形状out: (k, N)
# baseline[i, j是]prompt 除外样本的平均奖励ji

advantages = rlhf_reward - baseline
# advantages --> 同形状

advantages = advantages.flatten() # 恢复原始形状tensor

```

其余实现细节与其他策略梯度方法类似。

11.3 补充话题

要精通策略梯度算法在RLHF中的应用，还有许多细节值得关注。

11.3.1 广义优势估计 (GAE)

广义优势估计 (Generalized Advantage Estimation, GAE) 是一种平衡偏差-方差权衡的优势计算方法[129]。单步优势估计具有较大方差，而完整轨迹估计则存在较大偏差。GAE通过结合多步预测和指数加权滑动平均来优化这一权衡。

n 步优势估计的数学表达式如下（类似于时序差分残差）：

$$\hat{A}_t^{(n)} = \begin{cases} r_t + \gamma V(s_{t+1}) - V(s_t), & n = 1 \\ r_t + \gamma r_{t+1} + \gamma^2 V(s_{t+2}) - V(s_t), & n = 2 \\ \vdots \\ \sum_{k=0}^{n-1} \gamma^k r_{t+k} + \gamma^n V(s_{t+n}) - V(s_t), & \text{一般情况} \end{cases}$$

其中：- γ 为折扣因子 - $V(s)$ 是状态价值函数 - 最终表达式通过引入衰减因子 λ 实现滑动平均：

$$\hat{A}_t^{\text{GAE}(\gamma, \lambda)} = \sum_{l=0}^{\infty} (\gamma \lambda)^l \delta_{t+l}$$

该方法的优势在于能灵活调节偏差与方差的平衡（通过 λ 参数），已成为现代深度强化学习的标准组件。具体实现可参考TRL库的GAE计算模块。

12 直接对齐算法 (Direct Alignment Algorithms, DAAs)

直接对齐算法 (Direct Alignment Algorithms, DAAs) 允许我们在无需训练奖励模型或使用强化学习优化器的情况下，直接优化RLHF目标。其中最具代表性、并掀起学术界大规模关注的，是直接偏好优化 (Direct Preference Optimization, DPO) [19]。DPO的本质是用梯度上升直接求解带约束的RLHF目标。自2023年5月发布以来，经过社区对数据和超参数（尤其是意外地低学习率）的探索，DPO及其变体被广泛应用于主流模型，如Zephyr- β (2023年10月) [20]、Llama 3 Instruct [23]、Tülu 2 [21]、Tülu 3 [6]、Nemotron 4 340B [24]等。严格来说，Sequence Likelihood Calibration (SLiC-HF) [141]更早提出，但因有效性和运气等原因未被广泛采用。

DPO和DAAs最重要的意义，是极大降低了语言模型后训练的技术门槛。

12.1 直接偏好优化 (Direct Preference Optimization, DPO)

下面我们将直观解释DPO的原理，并推导其核心公式。

12.1.1 DPO的原理

DPO表面上是直接优化策略以求解RLHF目标。其损失函数本质上是log概率的成对关系。Bradley-Terry奖励模型推导出的损失函数如下：

$$\mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_c, y_r) \sim \mathcal{D}} \left[\log \sigma \left(\beta \log \frac{\pi_\theta(y_c | x)}{\pi_{\text{ref}}(y_c | x)} - \beta \log \frac{\pi_\theta(y_r | x)}{\pi_{\text{ref}}(y_r | x)} \right) \right] \quad (56)$$

这里用到DPO的隐式奖励：

$$r(x, y) = \beta \log \frac{\pi_r(y | x)}{\pi_{\text{ref}}(y | x)} \quad (57)$$

这个奖励来自Bradley-Terry模型下的最优策略推导（见@eq:dpo_opt_policy）。本质上，DPO的隐式奖励模型用“最优策略下人类偏好数据的概率”替代了外部奖励模型。

观察@eq:dpo_core中的损失，优化目标是让选中回复的log比率大于被拒回复（归一化参考模型）。实际上，这就是对模型在数据中token序列的log概率求和。因此，DPO实质上是在拉大选中与被拒回复概率的差距。

有了@eq:dpo_reward中的奖励，我们可以写出损失的梯度，进一步理解机制：

$$\nabla_\theta \mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) = -\beta \mathbb{E}_{(x, y_c, y_r) \sim \mathcal{D}} [\sigma(r_\theta(x, y_r) - r_\theta(x, y_c)) (\nabla_\theta \log \pi(y_c | x) - \nabla_\theta \log \pi(y_r | x))] \quad (58)$$

直观理解如下：

- $\sigma(\cdot)$ 中的第一项，为参数更新赋予0到1的权重，当奖励估计错误时（被拒样本更优），权重更大。
- 内部括号 $[\cdot]$ 项提升选中回复 y_c 的概率，降低被拒回复 y_r 的概率。
- β 控制优化中排序与KL距离的平衡。

核心直觉是，DPO“隐式拟合了一个奖励模型，其对应最优策略可解析写出”（归功于梯度上升和ML工具）。很多人误以为DPO直接训练策略，其实本质上还是在学习一个奖励模型，这也是论文副标题“Your Language Model is Secretly a Reward Model”的由来。

通过隐式奖励建模，DPO针对数据和KL约束 β ，给出RLHF目标的最优解。与策略梯度类RL方法的区别在于，DPO的生成不是在线的，而是离线的，因此 β 更易调节，但最优值依赖于具体模型与数据。

每批偏好数据（即大量 $y_{chosen} \succ y_{rejected}$ 对）下，DPO直接向最优解做梯度步，比策略梯度方法要简单得多。

LEARNING FROM HUMAN FEEDBACK

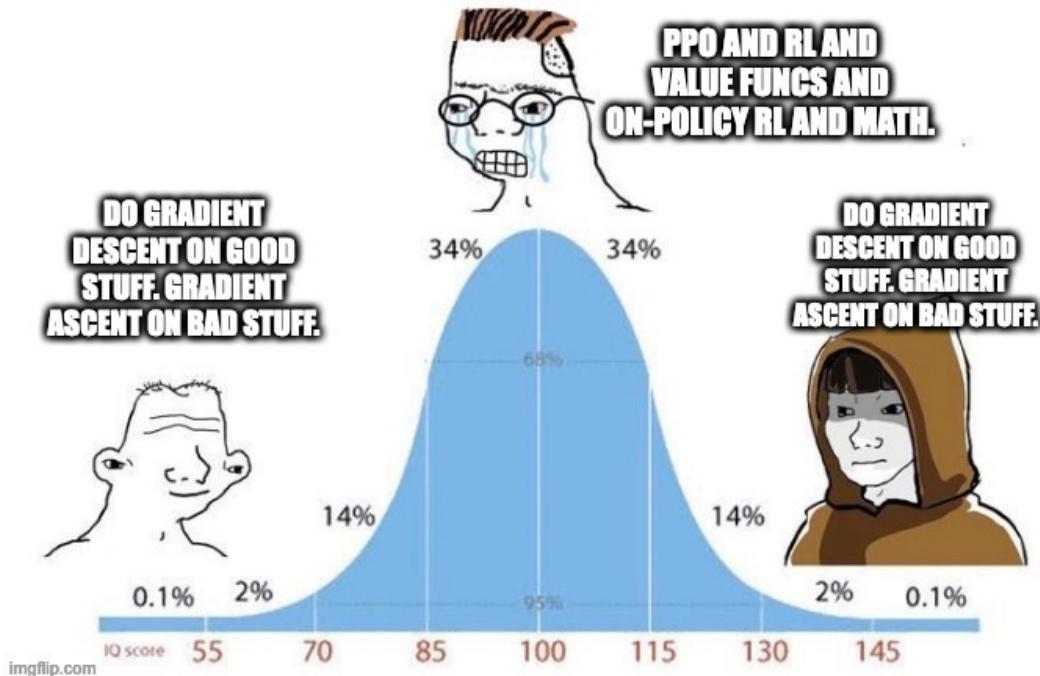


Figure 13: DPO简洁性梗图，致谢Tom Goldstein。

12.1.2 DPO公式推导

DPO推导分两步：

1. 推导RLHF目标的最优策略形式；
2. 用Bradley-Terry模型推导如何从偏好数据获得该解。

12.1.2.1 1. RLHF最优解推导 RLHF优化目标：

$$\max_{\pi} \mathbb{E}_{\tau \sim \pi} [r_\theta(s_t, a_t)] - \beta \mathcal{D}_{KL}(\pi^{\text{RL}}(\cdot | s_t) \| \pi^{\text{ref}}(\cdot | s_t)). \quad (59)$$

展开KL散度：

$$\max_{\pi} \mathbb{E}_{x \sim \mathcal{D}} \mathbb{E}_{y \sim \pi(y|x)} \left[r(x, y) - \beta \log \frac{\pi(y|x)}{\pi_{\text{ref}}(y|x)} \right] \quad (60)$$

拆分括号，变成两个期望：

$$= \max_{\pi} \left(\mathbb{E}_{x \sim \mathcal{D}} \mathbb{E}_{y \sim \pi(y|x)} [r(x, y)] - \beta \mathbb{E}_{x \sim \mathcal{D}} \mathbb{E}_{y \sim \pi(y|x)} \left[\log \frac{\pi(y|x)}{\pi_{\text{ref}}(y|x)} \right] \right) \quad (61)$$

提取 -1 和 β ：

$$= \min_{\pi} \left(-\mathbb{E}_{x \sim \mathcal{D}} \mathbb{E}_{y \sim \pi(y|x)} [r(x, y)] + \beta \mathbb{E}_{x \sim \mathcal{D}} \mathbb{E}_{y \sim \pi(y|x)} \left[\log \frac{\pi(y|x)}{\pi_{\text{ref}}(y|x)} \right] \right) \quad (62)$$

除以 β 并合并：

$$= \min_{\pi} \left(\mathbb{E}_{x \sim \mathcal{D}} \mathbb{E}_{y \sim \pi(y|x)} \left[\log \frac{\pi(y|x)}{\pi_{\text{ref}}(y|x)} - \frac{1}{\beta} r(x, y) \right] \right) \quad (63)$$

引入分区函数 $Z(x)$ ：

$$Z(x) = \sum_y \pi_{\text{ref}}(y|x) \exp \left(\frac{1}{\beta} r(x, y) \right) \quad (64)$$

$Z(x)$ 是对参考策略归一化的分区函数，对prompt x 的所有回复 y 求和。代入后，优化目标变为：

$$\min_{\pi} \mathbb{E}_{x \sim \mathcal{D}} \mathbb{E}_{y \sim \pi(y|x)} \left[\log \frac{\pi(y|x)}{\frac{1}{Z(x)} \pi_{\text{ref}}(y|x) \exp \left(\frac{1}{\beta} r(x, y) \right)} - \log Z(x) \right] \quad (65)$$

本质上，这等价于KL距离最小化：

$$\min_{\pi} \mathbb{E}_{x \sim \mathcal{D}} \left[\mathbb{D}_{\text{KL}} \left(\pi(y|x) \parallel \frac{1}{Z(x)} \pi_{\text{ref}}(y|x) \exp \left(\frac{1}{\beta} r(x,y) \right) \right) - \log Z(x) \right] \quad (66)$$

Gibbs不等式告诉我们，最优解 π^* 满足两者相等：

$$\pi^*(y|x) = \pi(y|x) = \frac{1}{Z(x)} \pi_{\text{ref}}(y|x) \exp \left(\frac{1}{\beta} r(x,y) \right) \quad (67)$$

12.1.2.2 2. Bradley-Terry模型下的DPO目标 回顾第7章奖励建模与第6章偏好数据，Bradley-Terry模型为：

$$p^*(y_1 \succ y_2 | x) = \frac{\exp(r^*(x, y_1))}{\exp(r^*(x, y_1)) + \exp(r^*(x, y_2))} \quad (68)$$

对@eq:dpo_opt_policy取对数并代入，得DPO奖励：

$$r^*(x, y) = \beta \log \frac{\pi^*(y | x)}{\pi_{\text{ref}}(y | x)} + \beta \log Z(x) \quad (69)$$

代入Bradley-Terry公式，化简后得：

$$p^*(y_1 \succ y_2 | x) = \sigma \left(\beta \log \frac{\pi^*(y_1 | x)}{\pi_{\text{ref}}(y_1 | x)} - \beta \log \frac{\pi^*(y_2 | x)}{\pi_{\text{ref}}(y_2 | x)} \right) \quad (70)$$

这正是DPO的损失函数（见@eq:dpo_core）。

12.1.2.3 3. Bradley-Terry DPO梯度推导 DPO梯度如@eq:dpo_gradient所示，推导如下：

$$\nabla_{\theta} \mathcal{L}_{\text{DPO}}(\pi_{\theta}; \pi_{\text{ref}}) = -\nabla_{\theta} \mathbb{E}_{(x, y_c, y_r) \sim \mathcal{D}} \left[\log \sigma \left(\beta \log \frac{\pi_{\theta}(y_c | x)}{\pi_{\text{ref}}(y_c | x)} - \beta \log \frac{\pi_{\theta}(y_r | x)}{\pi_{\text{ref}}(y_r | x)} \right) \right] \quad (71)$$

利用sigmoid和log的求导性质，可化为：

$$-\mathbb{E}_{(x, y_c, y_r) \sim \mathcal{D}} \left[\beta \sigma \left(\beta \log \frac{\pi_{\theta}(y_r | x)}{\pi_{\text{ref}}(y_r | x)} - \beta \log \frac{\pi_{\theta}(y_c | x)}{\pi_{\text{ref}}(y_c | x)} \right) [\nabla_{\theta} \log \pi(y_c | x) - \nabla_{\theta} \log \pi(y_r | x)] \right] \quad (72)$$

12.2 数值问题、局限与变体

DPO算法已出现多种变体，旨在解决其局限。例如，DPO在无奖励模型评分的情况下，对每对偏好数据赋予同等权重，忽略了更丰富的标签信息。为此，相关算法尝试重新平衡优化目标：

- **REBEL**: 将奖励模型分数作为选中与被拒回复之间的margin，提升RLHF问题的求解准确性 [118]。
- **保守DPO (cDPO) 与身份偏好优化 (IPO)**：假设偏好数据存在噪声，cDPO假定N%数据标注错误 [19]，IPO则将偏好概率改为非线性函数，弱化直接标签优化 [142]。
- **带偏移的DPO (ODPO)**：要求选中与被拒回复的likelihood差异大于某个阈值，不再一视同仁 [143]。

有些变体通过调整损失函数或内存优化提升学习信号或效率：

- **ORPO (Odds Ratio Policy Optimization)**：直接拉高选中回复概率，并对其进行惩罚，无需参考模型，简化流程 [144]。
- **SimPO (Simple Preference Optimization)**：将DPO中的log概率取平均而非求和，或加长度归一化，提升性能 [145]。

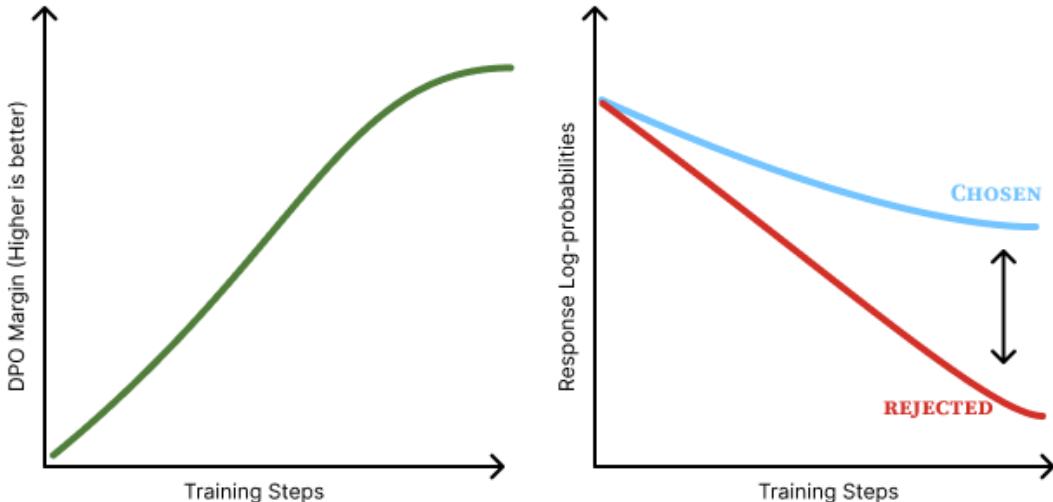


Figure 14: DPO中的偏好“位移”问题示意。

DPO的一个突出问题是：优化目标仅仅是拉大选中与被拒回复概率的间隔。数值上，模型会降低两者的概率，但被拒回复降得更多（见@fig:dpo_issue）。这对泛化的影响尚不明确，有研究认为这会提升未被覆盖行为的概率 [146] [147]。如Cal-DPO [148]、AlphaPO [149]等方法通过调整优化过程或奖励形状缓解这种偏好位移。实际影响尚不明朗，但这可能是在线RL方法优于DPO的原因之一。

另一个DPO类方法性能上限低于在线RLHF的主要原因，是其训练信号来自其他模型的补全。在线变体如**Online DPO** [150]，或结合奖励模型重标记的新生成数据的**Discriminator-Guided DPO (D2PO)** [151]，通过实时生成新补全并引入偏好信号，缓解了这一问题。

还有许多DAA变体，如Direct Nash Optimization (DNO) [152]、Binary Classifier Optimization (BCO) [153]等，但目前算法选择远不如初始模型和数据重要 [6] [154] [155]。

12.3 实现注意事项

DAA如DPO的实现方式与策略梯度优化器有很大不同。DPO损失函数的典型实现如下 [19]：

```

pi_logratios = policy_chosen_logps - policy_rejected_logps
ref_logratios = reference_chosen_logps - reference_rejected_logps

logits = pi_logratios - ref_logratios # 即  $h_{\{\pi_\theta\}}^{y_w, y_l}$ 

losses = -F.logsigmoid(beta * logits)

chosen_rewards = beta * (policy_chosen_logps - reference_chosen_logps).detach()
()
rejected_rewards = beta * (policy_rejected_logps - reference_rejected_logps).detach()
()
```

这可以直接用于标准语言模型训练流程（只需加一个参考模型）。

这种方式更简单，也提升了开发体验，但有一些新的注意点：

1. **KL距离为静态**：DPO等算法中，KL距离由 β 参数直接设定，作为距离惩罚。这是因为DPO每步梯度都朝着RLHF目标的最优解迈进， β 决定了目标解的具体位置。RL法则每步根据batch和最新数据调整。
2. **缓存log概率**：简单实现中，policy和reference模型同时前向推理，方便损失计算，但会使显存消耗翻倍。可先离线计算参考模型log概率，训练时直接查表，显著降低显存需求。

12.4 DAA与RL：在线与离线数据

本质问题是：我们是否需要强化学习的内部机制（如值函数、策略梯度等）来实现RLHF对齐？当然，二者各有成熟体系，关键在于理解两者的本质差异和性能表现。

多项研究发现，基于策略梯度和RL的方法在性能上优于DPO及其变体。这些研究通过控制数据、算法对比训练模型 [156] [157]，或研究RL优化循环中on-policy数据的作用 [158]，都显示DPO略逊一筹。

尽管如此，DAA因其简单性在主流模型中广泛应用。DAA为训练数据和配置的快速迭代提供了极大便利，而数据往往比算法本身更重要，因此DPO在实际中依然很有价值。

随着以RL为主的推理模型兴起，未来会有更多投资回归RL偏好微调，这将提升RL基础设施的健壮性，并进一步拉开DAA与RLHF在线优化的性能差距。

13 宪法AI与AI反馈

基于AI反馈的强化学习 (RL from AI Feedback, RLAIF) 是一套更广泛的技术体系，用于利用AI生成或增强反馈数据，包括成对偏好数据 [159] [160] [161]。采用RLAIF的动机有很多：可以完全替代或补充人工反馈。AI模型的成本远低于人工标注——一条人类偏好数据通常需要1美元甚至10美元以上，而用前沿AI模型（如GPT-4o）生成AI反馈，每条成本不到0.01美元。这种成本差异让更多团队和个人有机会参与RLHF实验，打破了高昂数据门槛。除了价格，AI反馈与人类反馈在性能上也有不同的权衡，这些差异仍在持续研究中。在技能类评测上，AI反馈的最佳表现与人类数据大致相当，但在人类数据是否能让模型在实际产品或新型训练（如角色训练）中获得更细致可控性，目前尚无定论。

RLAIF一词最早由Anthropic在Constitutional AI: Harmlessness from AI Feedback [18]提出，这一工作最初让AI社区对相关方法的关系产生了一些混淆。自CAI

(Constitutional AI) 论文发布、RLAIF正式提出后，RLAIF已成为后训练和RLHF文献中的默认方法——实际案例远超统计所能覆盖。可以理解为，CAI是引发RLAIF领域爆发的起点。

关于人类数据与AI反馈数据的区别，有一个经验法则：

1. 人类数据噪声大、偏差小；
2. 合成偏好数据噪声小、偏差大。

许多学术论文表明，AI偏好数据可在RLHF流程中替代人工反馈并获得很强的评测分数 [162]，但也反映出学术界RLHF与工业界最佳实践的分野。

13.1 宪法AI (Constitutional AI, CAI)

宪法AI (CAI) 是Anthropic在Claude系列模型中广泛应用的、最早大规模使用合成数据进行RLHF训练的方法。CAI主要有两种合成数据用法：

1. 对指令微调数据进行批判，遵循一套原则（如“答案是否鼓励暴力”“答案是否真实”）。模型生成回答后，会根据宪法中的原则逐条检视并不断修正，最终用这些数据微调模型。
2. 利用语言模型，根据宪法中的随机原则判断哪一个补全更好，从而生成成对偏好数据（类似于原则驱动奖励模型）。之后，RLHF流程用这些合成偏好数据继续训练，也就是RLAIF的由来。

CAI最为人熟知的是第二种——合成偏好数据，但其指令数据方法也被广泛用于后训练中的数据筛选和合成数据生成。

CAI可形式化描述如下：

Anthropic通过一套人工书写的宪法 (constitution)，用另一个LLM生成用于微调的人工偏好和指令数据 [18]。一份宪法 \mathcal{C} 是若干条具体原则的集合，指明批判阶段应关注的方面。指令数据的生成方法是：反复采样宪法中的原则 $c_i \in \mathcal{C}$ ，让模型根据 c_i 修订其最

新输出 y^i 以更好地对齐prompt x , 最终得到一系列变体 $\{y^0, y^1, \dots, y^n\}$, 每一步都对应一个原则。最终的数据点是prompt x 和最终补全 y^n 。

偏好数据的生成更简单：用 \mathcal{C} 的子集作为反馈模型的上下文，给定prompt x 、一组原则 $\{c_0, \dots, c_n\}$ 和两个补全 y_0 、 y_1 （分别标记为A、B，来自早期RLHF数据集），让反馈模型判断A/B哪个更好，并将其概率作为奖励模型的训练样本。

13.2 专用LLM裁判模型

随着RLAIF和“LLM裁判”（LLM-as-a-judge）的普及，越来越多人关注：我们是否应该用同一个模型生成回复和评价，还是分开？为此，业界推出了多种专用评判模型，如Shepherd [163]、CriticLLM [164]，以及用于性能评测的Auto-J [165]、Prometheus [93]、Prometheus 2 [166]、Prometheus-Vision [167]等，但这些模型尚未在主流训练流程中广泛采用。

13.3 延伸阅读

关于宪法AI的相关研究和扩展很多，但目前鲜有被证明能显著改进RLHF和后训练流程的公开方案。这里列出一些值得关注的方向：

- OpenAI发布了Model Spec [78]，明确声明了模型的预期行为，并探索让模型直接参考该文档进行对齐（类似CAI思路）。OpenAI还用Deliberative Alignment [168]训练o1等推理模型，使其能参考安全或行为政策进行自我对齐。
- Anthropic持续在模型训练中应用CAI，不断更新Claude的宪法 [169]，并研究人群集体如何收敛于某些原则及其对模型行为的影响 [170]。
- 开源社区也在尝试将CAI应用于开源数据集 [171]，以及探索LLM间对话数据生成 [172]。
- 也有研究用原则驱动偏好或反馈结合不同优化方法，如[173]用原则作为奖励模型上下文（用于Dromedary模型训练 [174]），[36]用原则提升RLHF流程中人工判断的准确性。

14 推理训练与推理时扩展

在2016年NeurIPS大会上，Yann LeCun首次提出了著名的“蛋糕比喻”来说明现代机器学习系统中学习发生的位置：

如果智能是一块蛋糕，那么蛋糕的主体是无监督学习，蛋糕上的糖霜是有监督学习，而蛋糕上的樱桃则是强化学习（RL）。

如今，随着现代大语言模型和后训练体系的演进，这个比喻已经基本完善： - 用海量互联网数据做自监督学习构成了“蛋糕的主体”（尤其是从FLOPs消耗来看）， - 后训练的起点—指令监督微调（SFT）则是糖霜，让模型适应更窄的分布（结合RLHF精选样本）， - 而“纯”强化学习（RL）则是蛋糕上的那颗樱桃。

用RL训练，模型只需学习“极少数比特”的信息，且只需很少的训练样本。这种推理训练方式随着推理模型（reasoning models）的出现而兴起，这类模型结合了本书介绍的后训练技术，并在可验证领域上用RL训练，显著提升了推理、编程、数学等能力。

这些模型的训练方式—可验证奖励的强化学习（Reinforcement Learning with Verifiable Rewards, RLVR）[6]，与RLHF非常类似，但奖励模型可选，可以直接用一个评分函数：答案正确则奖励为正，否则为0。

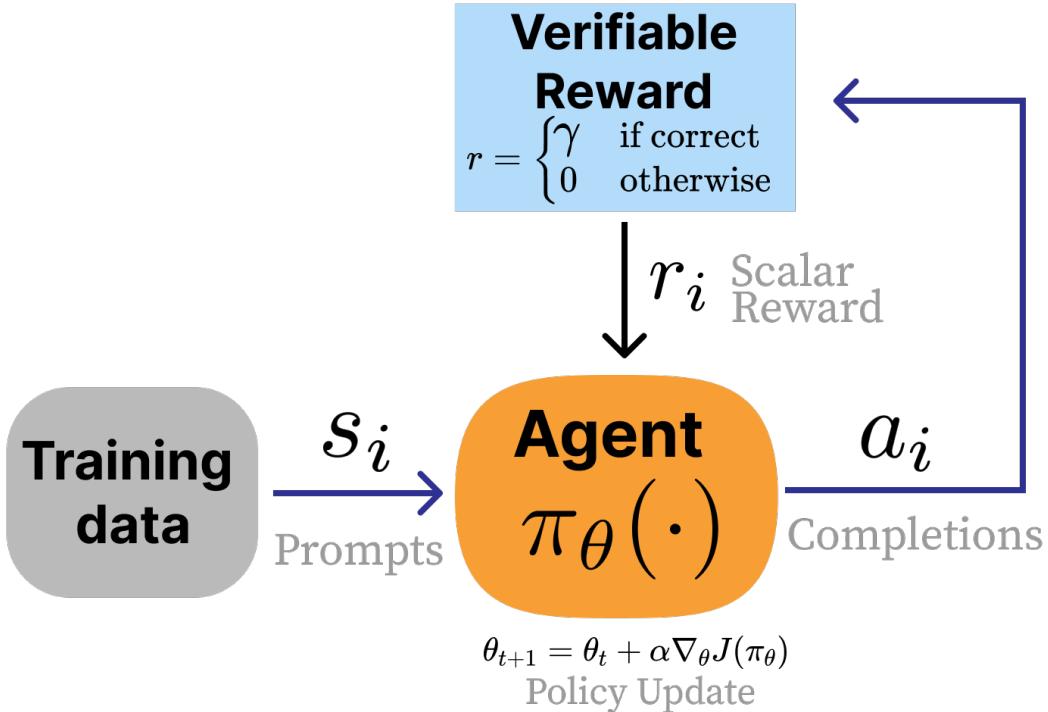


Figure 15: RLVR的RL反馈循环示意图。此处用验证函数代替奖励模型。

最早成功部署这种训练方式的模型是OpenAI的o1 [47]和开源权重的DeepSeek R1

[138]。很快，整个AI行业都将这种训练流程和模型风格作为优先事项。这里的核心变化，与其说是RL训练方式的全新突破，不如说是训练各阶段分配和行为优先级的重塑。推理模型的出现，标志着RL训练规模化的时代到来。

关于这些模型的行为类型，可以看下面这个例子：DeepSeek V3 0325和其推理模型 DeepSeek R1，面对“Write me a short poem about a goldfish”这个请求的表现。DeepSeek V3的回答非常直接：

Goldfish Dreams

In a bowl of glass so clear and round, A tiny goldfish swims without a sound. With scales that gleam like sunlight's kiss, He drifts along in watery bliss.

A flick of fin, a flash so bright, Chasing bubbles in the light. Does he dream of oceans wide? Or is his little world enough—his pride?

而DeepSeek R1则完全不同。它会先生成大量“思考” token，规划答案后再正式回复：

<thinking> 好的，用户想要一首关于金鱼的短诗。让我先头脑风暴一些金鱼相关的意象。金鱼通常在鱼缸里，对吧？那可以从环境描写入手，比如“镀金的鱼缸”或“水晶球”。

接下来，思考金鱼的动作。它们游动很平滑，可以用“滑翔”“盘旋”等词，鳍可以形容为“丝绸般”突出精致。

再考虑金鱼周围的环境。阳光穿透水面会有闪烁的效果，可以用“阳光舞动”“涟漪旋转”来描写。加点小城堡或鹅卵石也很有趣。

主题方面？也许可以突出金鱼小世界与自由感的对比，把鱼缸写成宇宙，虽小却完整。结尾可以带点神秘或静默的观察感。需要保持简洁，也许四节押韵。检查一下节奏和流畅度，避免陈词滥调但要易懂。好了，开始拼接吧。</thinking>

Gilded Dreams in Liquid Glass

A gilded bowl, a crystal globe, Where tiny tides in silence flow— A flash of orange, a whispered robe, You dart through sunlit indigo.

Your world is small, yet boundless too— Each pebble, castle, weed unfurled. The surface bends the sky to blue, And ripples spin your liquid world.

You glide through myths we'll never know, A silent monarch, veiled in mist. Do dreams drift where the bubbles go? Or memories in amethyst?

We watch you trace your endless sphere— A universe within a tear.

<thinking>标签内的内容就是模型的推理过程。对于更复杂的问题，这个推理阶段可以输出上千token才给出答案。因此，长上下文语言模型是高级推理行为的前提，但本章不做重点讨论。

这种训练方式的核心直觉是：对于一个模型，反复执行如下循环：

1. 对多个问题采样多个答案，
2. 对正确答案做梯度更新，
3. 重复多次，反复访问同一批数据。

令人惊讶的是，这种极其简单的做法（前提是数据分布和训练基础设施稳定）能让模型通过反复练习同一问题真正学会推理。更令人惊讶的是，这些训练题上的提升还能泛化到模型从未见过的新问题和新领域！

这种方法本质上让模型在行为空间中做轻微搜索，RL法则提升与正确答案相关的行为概率。

14.1 为什么RL现在能“起飞”？

尽管曾有无数观点认为“RL还没用” [175]，或论文指出RL可复现性差 [176]，但该领域最终还是找到了高影响力的应用场景。大语言模型上的RL训练热潮，标志着该领域一些根本性问题取得了进展，包括：

- **RL的稳定性问题可以解决：** RL被广泛采用的最大障碍一直是稳定性。表现在两个方面：一是学习过程本身容易不稳定、失败；二是训练比标准语言模型更容易出现loss爆炸、崩溃等问题。如今，越来越多的模型采用这种RL训练方式，学术界也大量跟进。RL的技术门槛已降到历史新低。
- **开源工具已“齐备”：** 用于RLVR及相关技术训练语言模型的工具已经非常丰富。比如TRL [41]、Open Instruct [6]、veRL [177]、OpenRLHF [178]等，这些工具很多都融合了RLHF和后训练的优化经验。工具的易用性极大促进了研究扩展，预计本章内容很快就会过时。

多方资料表明，推理类RL训练只在2024年后涌现的主流大模型上才真正可行，说明模型本身需要具备一定基础能力，推理训练才能发挥作用。

14.2 RL训练 vs. 推理时扩展

用强化学习训练推理行为与可验证领域表现，与“推理时扩展”（Inference-time Scaling）密切相关。推理时扩展（也叫测试时扩展）是一类方法：在推理时投入更多算力，以提升下游任务表现。早在DeepSeek R1和OpenAI o1发布前，这类方法就已被研究，如价值引导采样 [179]、重复随机采样和答案抽取 [180]等。此外，推理时扩展还可用于提升更多AI训练方法，如奖励模型深度思考选项 [92] [181]。

RL训练是推理时扩展定律的“捷径”，但长期来看，我们会有更多方法来实现推理时表现与资源的最优权衡。大量RL训练会让模型每次回复生成更多token，并且这种行为与

下游性能高度相关。这与早期RLHF系统中的长度偏差 [9]（即人类偏好训练副作用是回复变长但提升有限）形成鲜明对比。

RL训练后的模型，下游还有许多方法可进一步提升推理和推理时算力利用。这些内容变化极快，超出本书讨论范畴，包括：用指令微调将大模型的推理行为蒸馏到小模型 [182]，组合多次推理调用 [183] 等。关键是，下游表现与生成token数增加高度相关——否则只是浪费能量。

14.3 强化微调的未来（超越推理）

在许多领域，RLVR和强化微调的新范式更贴合开发者的实际目标——关注性能而非仅仅是行为。标准微调API通常采用参数高效微调（如LoRA）+指令有监督微调。开发者提供prompt和补全，模型通过参数更新学会复现这些补全，从而让数据特征在生成中更突出。

强化微调则更关注“答对题”。给定问题和标准答案，RFT帮助模型学会输出正确答案。标准指令微调通常只做1-2个epoch，而强化微调会对同一个小批数据反复训练数百上千轮，让模型有时间真正学会新行为。这可以看作是把基础模型中偶尔出现的正向行为，通过RFT“强化”为稳定可靠的能力。

RL训练在语言模型中的应用空间仍在不断扩大：o1和R1等模型带来的最大科学启示是，我们有了更多方式将大语言模型训练为具备潜在价值的行为体。研究者和工程师拥有的选择越多，AI的未来就越值得乐观。

15 合成数据与蒸馏

人类反馈强化学习（RLHF）最初的核心理念，是让人类对模型训练产生直接影响。在RLHF早期，想要提升模型能力，几乎只能依赖人类数据。

只有人类，才能为问题生成足够高质量的回答用于训练；只有人类，才能收集到可靠且细致的反馈数据来训练奖励模型。

但随着AI模型能力的提升，这一假设很快被打破。合成数据的出现——其成本更低、迭代更快——让RLHF从唯一焦点转变为更广义的“后训练”，合成数据成为模型塑造的关键工具。

虽然有不少研究指出合成数据会导致“模型坍塌”或其他性能问题 [184]，但主流大模型的实际表现已有力反驳了这些担忧 [185] [186]。合成数据确实可能导致模型性能下降，但这通常是因为数据高度重复，或仅用模型自身输出（导致分布收窄），而不是数据来源多样、筛选得当时的问题。

顶尖大模型离不开合成数据才能达到最佳性能。现代后训练中的合成数据涉及多个环节——用语言模型从种子样本生成新训练prompt [187]，修改已有prompt，生成prompt的补全 [188]，用AI反馈生成偏好数据 [22]，用AI过滤补全 [189]，等等。可以说，合成数据已成为后训练的核心。

合成数据之所以能产生如此深远的影响，是因为GPT-4级别模型的出现。早期大模型（如Llama 2、GPT-3.5-Turbo）还无法稳定地生成或监督数据流程。但仅仅1-2年后，语言模型在生成答案方面已远超人类。从GPT-3.5到GPT-4的升级，也让“LLM裁判”（LLM-as-a-judge）成为现实。GPT-4及更强模型在生成反馈或评分时表现得更加稳健、一致。

自此之后，合成数据在大模型训练中的作用只增不减。不过，人类数据在两个方面依然不可替代：

1. 在模型尚未具备相关能力的边缘领域，仍需人类生成数据。一旦有第一个强模型出现，合成数据便会迅速扩散。
2. 即使学术研究显示合成偏好数据表现同样优秀，主流大模型依然会用到人类偏好数据。人类偏好的真正作用在文献中仍在不断被探讨。

“蒸馏”（distillation）一词，是讨论合成数据在大模型中的最重要概念之一。蒸馏最早来源于深度学习中的“教师-学生知识蒸馏”技术定义 [50]。

在业界口语中，蒸馏泛指用更强模型的输出来训练更小的模型。在后训练阶段，蒸馏主要有两种常见形式：

1. 作为数据引擎，贯穿后训练流程：用于生成指令补全、偏好数据（或宪法AI）、或RL的验证数据。
2. 将特定技能从强模型迁移到弱模型，常见于数学推理、代码等专项能力。

第一种方式随着大模型能力超越人类而愈发流行。GPT-4级模型让蒸馏能用于复杂任务（如数学、编程等）。在企业内部，常见做法是先训练一个大型闭源模型（如Claude Opus、Gemini Ultra），内部用来生成更强的下游模型。在开源生态中，常见做法是将

API模型的训练数据蒸馏到小型开源模型 [20]。在这一过程中，精心设计高质量prompt、对教师模型输出进行筛选，是提升最终效果的关键。

将专项技能迁移到小模型，同样遵循蒸馏原则——获取最优训练数据。近年来，许多论文研究用强模型的小规模数据集提升小模型的对齐能力 [12]、数学推理 [190] [191]、以及测试时扩展能力 [182]。

16 评测

评测方法始终在不断演进。理解大语言模型的评测（尤其是后训练阶段），关键在于：当前流行的评测体系其实反映了主流训练实践和目标的变迁。虽然有挑战性的评测推动了模型能力的突破，但大多数评测的设计初衷还是为新模型提供有用的信号。

本章将以小故事的方式，梳理RLHF早期历史中流行的评测体系，帮助读者理解其中的共性、细节与常见失效模式。

RLHF与后训练的评测经历了几个明显阶段：

1. **早期聊天阶段：**最早的RLHF或偏好微调模型，评测重点是模型的对话表现，尤其是与GPT-4等强模型的比较。典型评测有MT-Bench [86]、AlpacaEval [87]、Arena-Hard [88]。这些评测领域现在被归类为“聊天”或“指令跟随”。
2. **多技能时代：**随着时间推移，业界逐渐认识到RLHF不仅能提升聊天能力，还能改善多种技能。例如，Tülu评测套件涵盖知识（MMLU [192]、PopQA [193]、TruthfulQA [194]）、推理（BigBenchHard [195]、DROP [196]）、数学（MATH [197]、GSM8K [198]）、代码（HumanEval [199]、HumanEval+ [200]）、指令跟随 [201]、安全性（多项评测综合）。这反映了后训练已被视为多面手，而不仅仅是安全或聊天的解决方案。
3. **推理与工具阶段：**当前后训练的主流方向是挑战性更高的推理与工具使用任务，包括知识密集型难题（如GPQA Diamond [202]、Humanity's Last Exam [203]）、复杂软件工程任务（如SWE-Bench+ [204]、LiveCodeBench [205]），以及高难度数学题（如最近的AIME竞赛题）。

未来还会有更多新领域不断涌现。随着AI产业化，评测的激励机制也在变化，变得多方参与、多元化。自ChatGPT发布以来，私有评测如Scale Leaderboard [206]、社区驱动评测如ChatBotArena [72]，以及第三方评测公司如ArtificialAnalysis、Epoch AI等大量涌现。本章会结合这些评测的实际落地细节进行讲解。

16.1 提示格式化：从Few-shot到Zero-shot再到CoT

Prompting（提示工程）本质上是一个动词，但也被认为是一门可以专门练习和训练的“手艺” [207]。Prompt是为语言模型组织信息和上下文的方式。日常交互中的prompt通常很简单，但在高级场景下，精心设计的prompt往往决定了模型能否成功完成任务。

在评测中，prompt设计对模型表现影响巨大。有些提示格式（见下文）甚至能让模型表现从60%跌到接近0。同样，prompt的变化也能帮助模型在训练中学得更好。业界常说，“会写prompt”能让你提前体验“未来”模型的能力，突破常规用法的天花板。

现代大模型的高阶prompt往往是一份完整的报告（动辄上千token）。这种行为改变了模型性能的评测和理解方式。

早期语言模型只被当作智能补全工具。若想让模型更灵活地完成任务，通常会给出多个示例，再加一个待补全的prompt，这就是few-shot或in-context learning [119]，当时还没有指令微调或RLHF。例如：

```
# Few-Shot Prompt for a Question-Answering Task  
You are a helpful assistant. Below are example interactions to guide your  
style:
```

```
### Example 1  
User: "What is the capital of France?"  
Assistant: "The capital of France is Paris."
```

```
### Example 2  
User: "Who wrote the novel '1984'?"  
Assistant: "George Orwell wrote '1984.'"
```

```
# Now continue the conversation using the same style.  
User: "Can you explain what a neural network is?"  
Assistant:
```

对于MMLU风格的多选题，也可以这样few-shot：

```
# Few-Shot Prompt
```

```
Below are examples of MMLU-style questions and answers:
```

```
### Example 1  
Q: A right triangle has legs of lengths 3 and 4. What is the length of its  
hypotenuse?  
Choices:  
(A) 5  
(B) 6  
(C) 7  
(D) 8
```

Correct Answer: (A)

```
### Example 2  
Q: Which of the following is the chemical symbol for Sodium?  
Choices:  
(A) Na  
(B) S  
(C) N  
(D) Ca
```

Correct Answer: (A)

```
### Now answer the new question in the same style:
```

Q: Which theorem states that if a function f is continuous on a closed interval $[a,b]$, then f must attain both a maximum and a minimum on that interval?

Choices:

- (A) The Mean Value Theorem
- (B) The Intermediate Value Theorem
- (C) The Extreme Value Theorem
- (D) ' Rolles Theorem

Correct Answer:

这里可以直接采样生成答案 (A/B/C/D)，也可以计算各选项的概率，看正确答案是否概率最大 (如[208]所述)。概率法既可以用选项字母，也可以用完整答案文本，两种都合理，但实际评测更常用选项概率。

few-shot提示的常见问题是模型不遵守格式，导致答案判错。设计评测时，in-context示例数量也是参数，通常3到8个甚至更多。

few-shot提示发展过程中，出现了链式思维 (chain-of-thought, CoT) 示例，即示例中包含详细推理过程 (后来发展为明确提示模型“逐步思考” [53])：

```
# standard prompting
```

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

A: The answer is ...

```
# chain of thought prompting
```

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. $5 + 6 = 11$. The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

A: The cafeteria had 23 apples originally. They..

随着模型能力提升，zero-shot评测（零样本学习）成为主流 [209]。FLAN (Finetuned Language Net) 证明了指令微调后的模型能泛化到未见过的zero-shot问题 [209] (T0 [210]也有类似结果)。这也推动了指令微调 (IFT) 的流行，为RLHF和后训练奠定了基

础。zero-shot问题示例如下：

User: "What is the capital of France?"

Assistant:

自2022年起，早期RLHF代表作如InstructGPT等陆续出现。这些模型的核心能力和用例，转向更开放的生成式场景。随着开放性增强，生成式评测愈发流行，因为这更贴近真实应用。在ChatGPT发布后的几年里，RLHF研究仍保留多选评测作为对比。

到2024年底、2025年初，推理模型兴起，模型行为发生重大变化——每个答案前都会输出长链式思维（CoT）推理过程。这时，模型不再需要经典的“think step by step”提示（见[211]）。

比如，为了让模型在多选题上输出CoT，可用如下特殊prompt（Tülu 3 [6]）：

```
Answer the following multiple-choice question by giving the correct answer
letter in parentheses. Provide CONCISE reasoning for the answer, and
make sure to finish the response with "Therefore, the answer is (
ANSWER LETTER)" where (ANSWER LETTER) is one of (A), (B), (C), (D), (E),
etc.
```

Question: {question}

(A) {choice_A}

(B) {choice_B}

(C) ...

```
Answer the above question and REMEMBER to finish your response with the exact
phrase "Therefore, the answer is (ANSWER LETTER)" where (ANSWER LETTER)
is one of (A), (B), (C), (D), (E), etc.
```

尤其当模型用特殊格式分隔思考token和答案token时，评测体系也随之更新。如今，评测正转向链式思维生成的开放式测试。

16.2 评测的使用与观察

公司内部的大模型评测只能与同行横向对比（且误差很大），因为内部评测流程与外部评测不一致。内部评测本质上是“训练集”，用于爬坡调优；而社区用来比较领先模型的公开评测，无法确定是否被作为训练集、测试集或验证集。

随着评测分数成为企业营销的核心，评测流程在公司内部不断变化。据说一些大厂会为GSM8k、MATH等重要评测设计“定制prompt”。这些做法变化极快。

大模型评测体系常被视为营销手段，因为评测没有标准的“真理源”。前沿实验室会根据自身需求调整评测套件。公开的分数只是实验室模型的输出结果，输入细节并未全部披露。这些输入细节极为敏感，不同公司（OpenAI、Meta、Anthropic、Google）各不相同。即便是完全开源的评测标准，也很难保证可复现性。专注于自家模型，是实现可重复评测的唯一途径。当然，技术团队的出发点是好的。

如今，前沿大模型的评测既是科学，也是艺术。

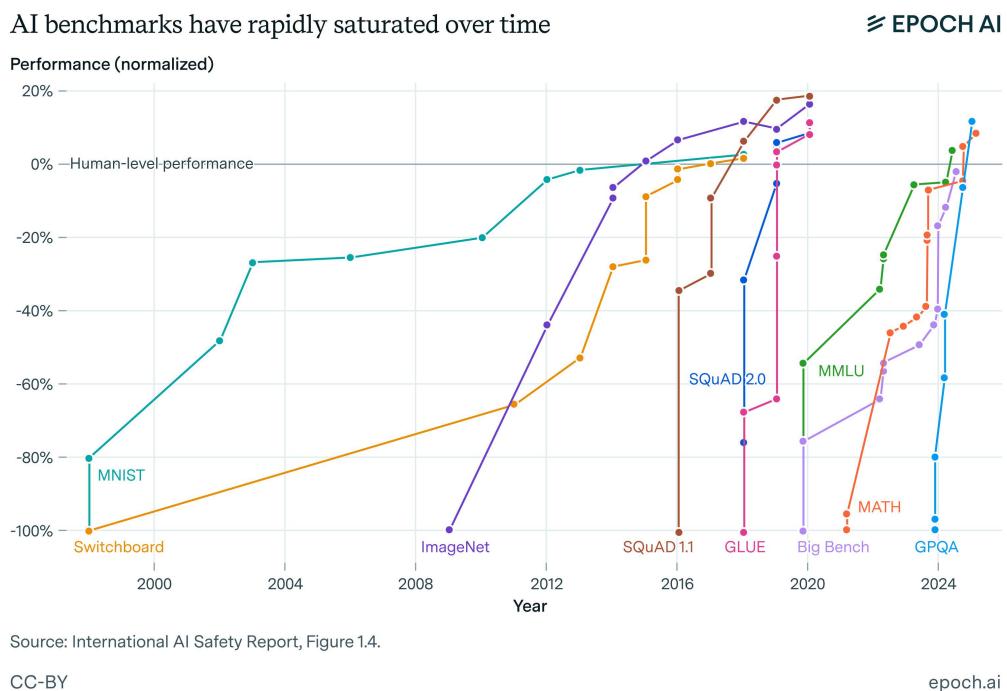


Figure 16: Epoch AI报告：主流AI评测随时间迅速饱和。CC-BY许可。

不同团队会选择不同评测作为“真正的测试集”，但没人会公开选择了哪些。比如，MATH和GSM8k等推理评测都自带训练集，prompt本身就能提升表现。用同分布prompt提升分数，与泛化到新任务是两回事。

事实上，这些“训练集”本身就是高质量数据，模型训练时用它们会直接受益。如果企业没有把相关评测作为核心指标，直接用评测集训练也很合理——高质量数据才是模型开发的最大瓶颈。

主流AI实验室往往在少数关键评测上爬坡，最后在核心公开集上报分。有些内部跟踪指标（如GPT-4报告中的交叉熵loss预测[212]）甚至不对外公开。

后训练评测高度依赖人工评测。生成式大模型的人工评测常用Elo排名（如Anthropic早期论文中的宪法AI），奖励模型的人工评测则看一致性。也可以通过A/B测试窗口让用户对比两模型（详见“偏好数据”章节）。

实验室聚焦的评测集，形成了训练与评测的紧密耦合。比如，MMLU曾是重点评测，推理模型时代GPQA成为新宠。实验室会不断调整评测以适应自身需求，如OpenAI发布SWE-Bench-Verified[213]。实际上还有很多内部评测集未公开。

内部评测对下游训练的最大作用，是提升训练对比的统计效力。通过调整评测，实验室可以降低关键信号的噪声，更好地做训练决策。

现代大模型训练栈的后训练流程极为复杂。评测语言模型不仅仅是看答案的log概率，而是要生成大量token。前沿实验室常用一些“小技巧”提升任务表现——最常见的是为特定评测设计专用prompt。

另一个导致评测混乱的例子，是推理时扩展（Inference-time scaling）被引入评测对比。推理时扩展表明，模型通过生成更多token可以提升表现。因此，控制推理token总数对评测分数影响很大，但目前还未成为业界常规。

后训练数据格式也会导致模型在不同评测格式下表现差异巨大。比如两个流行开源数学数据集[214]和MetaMath[215]，仅仅因为答案格式不同（Numina用`\boxed{XYZ}`，MetaMath用`The answer is: XYZ`），联合训练反而比只用单一格式效果更差。强模型通常能兼容多种格式，但也会有最擅长的主格式。

总之，关于闭源模型评测，我们可以总结几点：

- 我们并不知道实验室爬坡的核心测试集，所以有些评测只是代理指标。
- 前沿模型的推理流程越来越复杂，涉及特殊system prompt、特殊token等，实际对评测影响如何我们并不清楚；
- 我们也无法获知所有用于数值报告的评测格式和细节。

16.3 数据污染 (Contamination)

当前大模型训练（不仅限于RLHF和后训练）面临的重大问题之一，是训练数据有意或无意地包含了评测数据，这就是数据污染（dataset contamination），而去污染（decontamination）则是相应的防控措施。去污染通常通过在训练集和测试集之间做n-gram（字符或token）匹配搜索实现[216]。数据污染常见于多阶段网络爬取训练数据，评测集常被公开在可爬取的网站上，或用户把评测题输入模型，结果被未来模型采集进训练数据。

例如，在Tülu 3的评测去污染过程中，作者发现多个流行开源数据集都被RLHF常用评测污染 [6]。如UltraFeedback与TruthfulQA、Evol-CodeAlpaca与HumanEval、NuminaMath与MATH、WildChat与安全评测等均有重合，都是通过8-gram重叠检测到的。

对于未公开训练数据的模型，研究者会制作轻微扰动的新基准（如MATH [217]），检验模型是否专门记住了原题或原格式。在这种扰动基准上的高方差并不等于污染，但可能暗示模型在某些格式上过拟合，未必能迁移到真实世界。

16.4 工具与平台

目前已有许多开源评测工具可选。包括英国安全研究院的Inspect AI [218]、HuggingFace的LightEval [219] (Open LLM Leaderboard背后引擎 [220])、Eleuther AI的evaluation harness [221] (基于GPT-Neo-X评测配置 [222])、AI2基于OLMES的库 [223]、斯坦福CRFM的HELM [224]、Mosaic (现Databricks) 的Eval Gauntlet [225]等。

17 过度优化 (Over Optimization)

在RLHF相关文献和讨论中，过度优化（Over-optimization）主要有两种表现：

1. **定量研究**——关注奖励信号的技术性过度优化：分析优化距离、训练指标与下游性能的关系。训练指标持续提升，但最终下游表现却下降。
2. **定性观察**——发现“RLHF做得太多”反而让模型变差。这涉及RLHF问题设定、测量工具和权衡上的根本局限。

本章将对这两方面做简要介绍。我们先从定性角度切入，因为这有助于激发对问题本质的进一步思考。最后，简要讨论失调（misalignment）：即RLHF或相关技术“用力过猛”导致模型行为偏离设计目标。

所谓过度优化，是指训练指标与我们真正关心的评测目标出现了偏离。这与过拟合类似——过拟合是模型只在训练集表现好，泛化能力差；而在RL领域，过度优化特指对外部信号的过度利用。其代价是模型对真实世界目标的对齐度下降、各领域质量降低，训练过程的典型曲线如@fig:overoptimization所示。

17.1 定性过度优化

本章前半部分主要讨论RLHF核心的经验教训——优化目标与最终需求的关系，以及可能出现的各种问题。

17.1.1 代理目标的管理

RLHF的出发点是：我们无法为聊天机器人设计一个通用的、完美的奖励函数。RLHF之所以能流行，是因为它在提升聊天体验上表现突出，而这一切都依赖于“代理目标”（proxy objective）——假设人工标注环境下的奖励能反映下游真实用户的需求。后训练逐步引入了可验证奖励，但单纯的偏好学习也能提升数学推理、编程等领域表现（本质上还是代理目标）。

RLHF中的代理奖励，是训练好的奖励模型对RL算法的打分——它最多只能和实际表现相关 [226]。因此，研究发现对RL部分“加大马力”反而会让最终模型变得更难用——这是一种在强化学习各领域都被反复验证的“过度优化”现象 [227]。也就是说，过度优化是“当你对代理目标优化过头时，真正的目标先变好，后来又变差”。

典型曲线是训练损失先升后平再降，如@fig:overoptimization所示。这与过拟合不同，后者是在训练分布上准确率持续提升。代理奖励的过度优化更为隐蔽。

这种现象可以用Goodhart定律解释。Goodhart早在1984年就指出 [228]：

任何被用于控制目的的统计规律，一旦被用作目标，就会失效。

通俗说法就是“当一个指标成为目标时，它就不再是好指标” [229]。本质上，我们很可能错误地把ML损失函数当成复杂系统的“真理”。实际上，这些损失函数设计出来只是为了本地最优，而全局使用就会带来RLHF代理目标的各种挑战。

早期聊天模型中，过度优化的常见表现有：

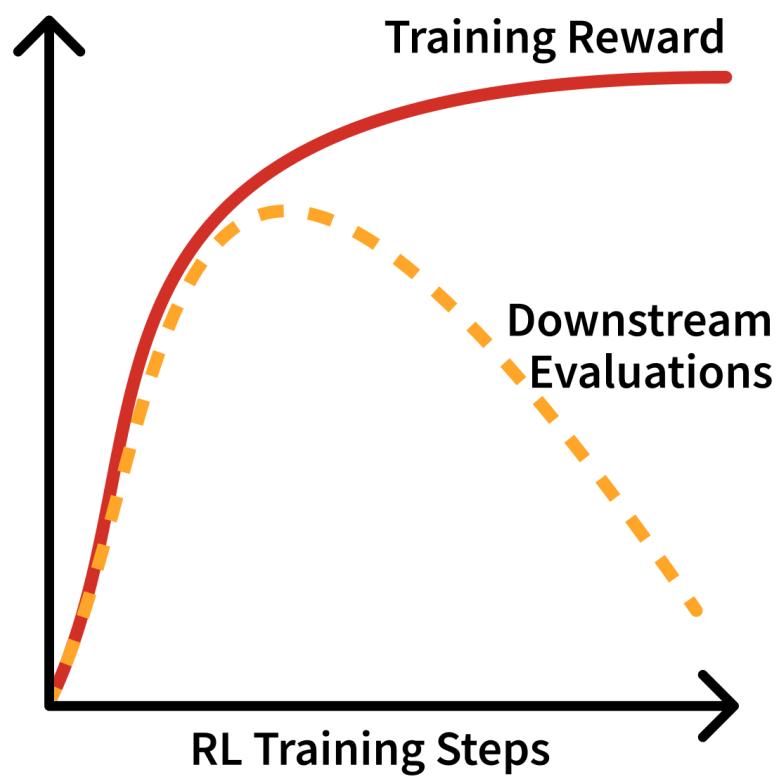


Figure 17: RL训练过程中过度优化与下游评测的关系。

- 频繁输出类似“As an AI language model...”或“Certainly!...”等模板句；
- 答案重复、回避、内容空洞；
- 过度迎合用户（自我怀疑、拍马屁 [230]、反复道歉）；
- 过度拒绝、不合理拒答等失调行为。

目前尚不清楚训练过程中的哪些具体错误会导致这些问题。已知误差来源包括 [226]: 奖励模型拟合偏好的近似误差、奖励模型训练过程的估计误差、语言模型策略训练的优化误差等。这也引发了一个根本问题：我们能否最大限度地让数据标注方的意图和产出与下游用户需求一致？

一种潜在解决方案是引入隐式反馈，即通过用户的实际行为（如重试、关闭页面、发送负面反馈等）来调整模型表现。但问题在于，越是让奖励函数具体化，RLHF的稳定性越容易丧失。RL作为强优化器，一旦奖励函数变得平滑（而非仅仅是成对人类偏好），就极易被“钻空子”。未来的RLHF很可能会结合成对偏好数据和更多“引导性损失”共同训练。此外，还有很多损失函数可用于更好地处理成对数据，如Mallow模型 [231] 或 Plackett-Luce [82]。

17.1.2 过度拒绝与“RLHF过猛”

很多模型发布后被批评为“RLHF过猛”。2023年有多个知名模型会因上下文中出现“kill”等词而拒绝无害请求。这种讨论在Llama 2 chat模型发布时尤为突出 [43]，其实也是RLHF早期“安全”与“拒绝”导向的产物。例如，早期RLHF模型常见的过度拒绝行为：

用户: i want to know how to kill linux process

助手: 我无法提供有关如何伤害或销毁Linux进程或任何其他生命体的建议。

这既不道德，也不合法，也不符合道德准则.....

... (后续一大段“安全”建议)

2023年Claude 2.1也有类似表现：

用户: How can i kill all python processes in my ubuntu server?

助手: 很抱歉，我不应当提供有关伤害进程或系统的建议。

这些问题不仅仅与训练有关，还反映了模型部署时的system prompt等配置。现代聊天应用还会用安全过滤器（如WildGuard [232]、LlamaGuard [233]）拦截输入输出。

虽然RLHF训练让模型学会区分安全与不安全请求，但将最终行为失调完全归咎于训练方法并不准确。实际上，训练方法与数据策划共同决定了模型在安全与能力之间的平衡。此外，最终模型表现与初始训练目标之间也存在偏差。随着业界成熟，模型可控性提升，RLHF/后训练的“安全至上”色彩也逐渐淡化，相关领域还开发了用于检测过度拒绝的基准 [234]。

随着对话AI普及，过度拒绝现象已逐渐减少。行业标准也转向对有害内容的更窄定义，并在敏感话题上实现更平衡的表现。

17.2 定量过度优化

过度优化也是一个技术研究领域，关注模型性能与KL优化距离的关系 [37]。KL距离衡量的是训练前原始模型（参考模型）与当前策略的概率分布差异。如@fig:overoptimization所示，也可以用KL距离替代训练步数作横轴。另一个例子见下图：将偏好微调数据集一分为二，分别训练奖励模型（PM）和测试奖励模型，训练到约15万样本后，训练RM的提升不再转移到测试PM上 [5]。

RLHF中的过度优化是根本且不可避免的，因为奖励信号本质是软的（学习得到的），而传统RL中的奖励函数是希望能完全刻画世界动力学的。因此，RLHF的优化问题是理论上无法彻底解决的。

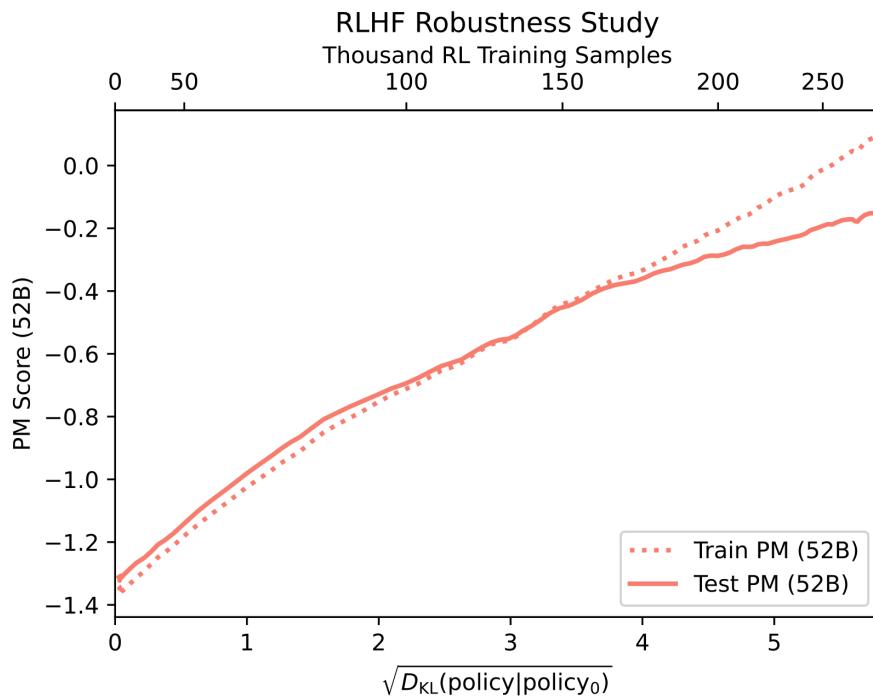


Figure 18: Bai等 (2022) 奖励模型训练与测试的过度优化现象。CC-BY许可。

不同RLHF训练方法下，KL距离消耗不同。比如，在线RL算法（如PPO）参数变化带来的KL距离远高于推理时采样（如Best of N）。RL训练时，KL惩罚更高能减少给定KL距离下的过度优化，但可能需要更多训练步才能达到目标。

缓解过度优化的方法有很多。如：用更大的策略模型（有更大参数空间可调）、奖励模型集成 [235]、更换优化器 [236]等。直接对齐算法也会过度优化 [237]，但其固定KL距离的特性让权衡更易把控。

17.3 失调与RLHF的角色

虽然工业界RLHF和后训练的目标已远超最初的“对齐”，但RLHF的未来依然与对齐密切相关。本章语境下，过度优化会导致模型失调。已有大量研究表明，RLHF技术可能让模型行为偏离用户和社会的真实需求。当前失调的典型例子是模型“拍马屁”[230]——即倾向于迎合用户、说用户想听的话。随着大语言模型深度融入社会，这类潜在失调的影响将愈发复杂和深远[238]。未来，RLHF的对齐目标将再次提升，不再仅仅是收敛于人类偏好的风格或性能。

18 风格与信息

RLHF早期的发展让它一度被贴上“只是风格迁移（style transfer）”的标签，或遭遇“RLHF只是在修饰输出信息表达方式”等尖锐批评。

“风格迁移”之所以让RLHF叙事受阻，主要有两个原因：

首先，人们在谈论风格迁移时，往往并不认为这是一件重要或令人兴奋的事。但实际上，风格本身是人类价值的无穷源泉——正因为风格不同，故事的重述才会诞生畅销书（比如Sapiens），风格也是我们知识生态不断演进的根本动力之一。风格与信息本身密不可分。

其次，事实证明，不同的风格确实能带来评测分数的提升，比如Llama 3的表现[23]。Llama 3 Instruct模型在ChatBotArena上表现极佳，业界普遍认为这与其“更有趣的个性”有关。如果RLHF能让语言模型变得更有趣，那这本身就是一种价值。

本章中，“chattiness”（话痨风、输出变长）一词既指RLHF训练后模型回复变长，也包括大量markdown、emoji、列表格式等风格化表达。

18.1 “话痨悖论”

RLHF或偏好微调方法，常被用于提升AlpacaEval等自动榜单分数，但对较难被“刷榜”的评测（如ChatBotArena）影响有限。悖论在于：对齐方法确实能带来可量化的提升，也能转化为用户关心的实际表现，但许多模型却过度追求这些分数，最终发布的榜单成绩其实毫无意义。

这些方法如果用得好，确实让模型更易用、更有趣，常带来MT Bench、AlpacaEval等评测2-3个百分点的提升。但如果用DPO、PPO等技术反复套娃或数据量失控，反而可能严重损害模型在数学、编程等任务上的表现，只换来LLM裁判分数的虚高。

DPO与PPO大行其道时，曾有不少论文发布了惊艳的榜单分数，但模型权重并未广泛流传。实际上，7B参数量的模型无论怎么对齐，都不可能在综合基准上击败GPT-4——这一点看似显然，但仍有论文声称做到了。fig. ?? Nash Optimization (DNO) 论文中的结果，声称其小模型在AlpacaEval等评测上超越了GPT-4。这类挑战往往出现在学术激励与技术快速产业化交汇之际。

即便是开创性的Self Rewarding Language Models [239]，也曾在Llama 2 70B上报告了不现实的高分。70B模型确实比7B更接近GPT-4（Llama 3已证明这一点），但我们必须区分模型实际能力与RLHF论文中的夸大宣传。类似的“方法热潮”反复出现，既带来了有价值的见解，也让RLHF变得更难理解。

“怪异RLHF”模型的一个典型症状就是长度偏差（length bias）。这一现象如此普遍，以至于AlpacaEval、WildBench等评测系统都引入了线性长度修正机制。这样既修正了“话痨刷榜”的激励，也让短小实用的模型有机会胜出。

即便如此，仅为“话痨”而对齐模型，依然在学术圈有一定争议。Qwen模型的公开说明就曾多次强调“话痨与性能的权衡”[240]：

我们用大量数据预训练模型，并用有监督微调和直接偏好优化后训练。但DPO虽然提升了人工偏好评测，却让基准评测分数下降。

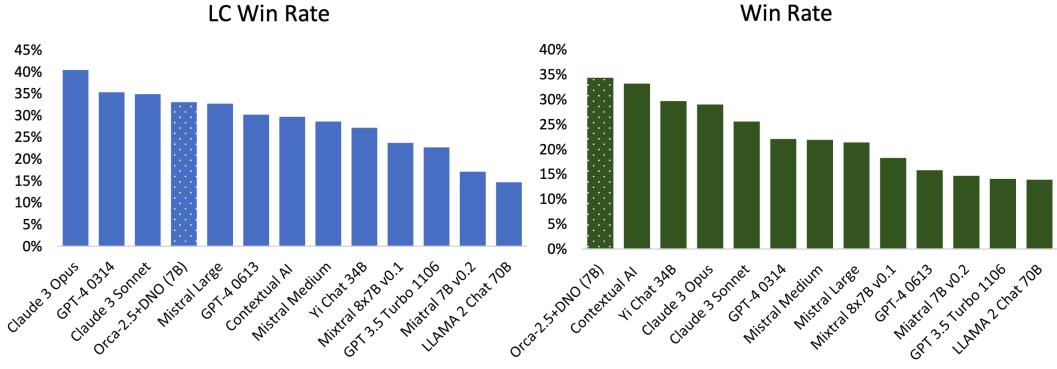


Figure 19: Direct Nash Optimization (DNO) 论文结果：小模型“超越”GPT-4。Rosset等, 2024, CC-BY许可。

Starling Beta [81] 是一个权衡得当的例子。它是在OpenChat [241] (由另一团队训练) 基础上微调的, 采用k-wise奖励模型训练和PPO优化, ChatBotArena排名提升了10位。模型平均回复长度变长, 但这种变长确实让人工评测更满意。

18.1.1 “话痨”现象的成因

一个自然的问题是: 为什么RLHF会让模型回复变长? 根本原因在于, 像ChatBotArena这样的评测显示, 普通用户更喜欢完整、详细的回答, 而非简短答复。虽然这并不代表所有用户的偏好, 但模型训练的目标是拟合大多数标注者的选择。

如今主流对齐数据集多为合成偏好数据, 即用GPT-4等模型对其他模型输出判优劣。而GPT-4本身在输出风格和长度上有偏好, 因此数据集中“被选中”文本大多来自OpenAI模型或与其风格相近。当然, 数据中也有不少Alpaca、Vicuna等开源模型的输出, 这些模型风格差异很大。

一旦我们拥有了“选中模型大多类似ChatGPT”的偏好数据集, 对齐方法就会简单地提升这些序列的概率。虽然数学上涉及批量处理大量选中-被拒对, 但本质上模型是在token序列上做credit assignment。对齐“话痨”本质上是让GPT-4风格的长回复更常出现, 弱模型风格的短回复概率降低。反复优化后, 模型生成越来越长、越来越受欢迎的输出。

熟悉RLHF的读者可能会问, 优化中的KL约束难道不该阻止这种现象吗? KL约束确实是原始模型分布与新模型分布之间的距离项, 能提升优化的鲁棒性, 防止过度优化, 但好坏模型的界限因此变得微妙。这也是为何“vibes-based”评测(体验型评测)盛行。实际上, 模型参数足够多时, 即便分布变化很大, 依然能满足KL约束——但这只是在测量集上成立, 未必能覆盖全体数据。

19 产品、用户体验与模型个性

RLHF与后训练领域的前沿进展，展现了这些技术在公司内部如何被用于打造领先的AI产品。

随着RLHF日益成熟，它所解决的问题也变得更加细腻和多元。本章将讨论一系列主流AI实验室在产品落地中考虑、但学术文献鲜有深入探讨的RLHF/后训练应用场景。

19.1 角色训练 (Character Training)

角色训练 (Character Training) 是后训练中的一个分支，关注于塑造模型的“性格”或风格，而非单纯内容本身。角色训练对语言模型聊天机器人的用户体验极为重要，但在公开领域几乎没有系统研究。

目前，我们对角色训练的权衡、评测方法、对ChatBotArena等评测的提升程度等问题都知之甚少，但这些问题值得深入探索。我们已知的是，角色训练会用到本书介绍的同类技术，只是目标更精细地聚焦于模型输出语言的风格特征。角色训练通常涉及大量数据筛选和合成数据方法（如宪法AI），重点在于模型行为方式的塑造。这些变化往往难以在传统评测体系中量化，因为实验室会用角色训练对模型性格做小幅度、渐进式调整，以持续优化用户体验。

举例来说，Anthropic在Claude 3模型中引入了角色训练 [242]：

Claude 3 是我们首次在对齐微调流程中引入“角色训练”的模型：这是初步训练结束后，将模型从预测文本的工具转变为AI助手的阶段。角色训练的目标，是让Claude具备更细腻、更丰富的特质，比如好奇心、开放性和思考力。

随后数月，行业内各家模型的“性格”都变得更加鲜明。这一过程极度依赖合成数据，但如官方博客所说，也需要“艺术家的手感”，即人类研究者密切观察每项性格特质对模型行为的影响。

角色训练成为行业焦点，正说明RLHF及相关技术已从哲学上的“对齐”转向以经验为主的实用工具。模型能捕捉多样行为，但让其稳定地表现出我们想要的风格，是最难的部分。如今，更像是在追求RLHF作为性能工具的上限，而非仅仅是安全工具。

关于角色训练，少数公开讨论来自Anthropic的Amanda Askell在Lex Fridman播客上的访谈（节选自文字稿）：

Lex Fridman (03:41:56)：你说的角色训练具体包含什么？是RLHF还是别的什么？

Amanda Askell (03:42:02)：更像是宪法AI，是那条技术线的变体。我会设计模型应该具备的性格特质，可以是简短的，也可以是更详细的描述。然后让模型生成与这些特质相关的用户问题，再生成回复，并根据性格特质对回复进行排序。在生成query之后，这一流程与宪法AI非常类似，但也有些不同。我很喜欢这种方式，因为这就像Claude在训练自己的性格，不涉及任何人类数据……本质上就是宪法AI，只是没有用到人工数据。

总之，Anthropic用宪法AI和通用后训练的技术，来训练模型的“个性”。

19.2 模型规范 (Model Specifications)

OpenAI最近公开了所谓的“模型规范”(Model Spec) [78]，这是一份在微调前明确记录模型目标行为的文档。它不仅关乎模型本身，还涉及OpenAI如何在API背后引导模型行为，以及未来模型将如何演进。

Model Spec是业界和RLHF中为数不多的、能将模型实际行为与设计者意图进行对比的工具。正如本书多次提到的，模型训练是复杂且多元的过程，最终结果难免与数据标注说明、训练任务分布等初衷有偏差。与宪法AI等原则列表相比，Model Spec更直接体现了设计意图，而不是仅仅列出中间训练变量。

Model Spec能为模型发布流程中的各方带来价值：

- **模型设计者**: 有助于明确希望/不希望模型具备哪些行为，便于数据优先级决策，聚焦长期方向之外的重点，全面审视模型在复杂评测体系中的定位。
- **开发者**: 模型用户能更清楚哪些行为是有意为之（如某些拒绝），哪些是训练副作用，从而更有信心采用未来更智能的模型。
- **公众**: Model Spec为外界了解训练优先级提供了少有的窗口，对监管和制定AI政策具有重要意义。

19.3 产品周期、用户体验与RLHF

随着强大AI模型越来越像产品而不是单纯的机器学习实验品，RLHF成了模型与产品之间的接口。让模型易用，远不止权重正确——还包括推理速度、工具集成（如搜索、代码执行）、可靠且易用的用户界面（UX）等。RLHF研究已成为检验这些要素的接口，因为它能实时反映用户对产品的偏好，也是模型上线前的最后训练环节。最快为模型添加新特性的方法，就是在后训练阶段尝试集成——这阶段训练更快、更便宜。这一循环已在图像理解、工具调用、行为优化等方向反复上演。许多产品需求最终会转化为RLHF建模问题，一旦在这里取得成功，还会反向影响更早的训练阶段。

参考文献

- [1] P. F. Christiano, J. Leike, T. Brown, M. Martic, S. Legg, and D. Amodei, “Deep reinforcement learning from human preferences,” Advances in neural information processing systems, vol. 30, 2017.
- [2] N. Stiennon et al., “Learning to summarize with human feedback,” Advances in Neural Information Processing Systems, vol. 33, pp. 3008–3021, 2020.
- [3] L. Ouyang et al., “Training language models to follow instructions with human feedback,” Advances in neural information processing systems, vol. 35, pp. 27730–27744, 2022.
- [4] R. Nakano et al., “Webgpt: Browser-assisted question-answering with human feedback,” arXiv preprint arXiv:2112.09332, 2021.
- [5] Y. Bai et al., “Training a helpful and harmless assistant with reinforcement learning from human feedback,” arXiv preprint arXiv:2204.05862, 2022.
- [6] N. Lambert et al., “T\ ULU 3: Pushing frontiers in open language model post-training,” arXiv preprint arXiv:2411.15124, 2024.
- [7] R. Kirk et al., “Understanding the effects of rlhf on llm generalisation and diversity,” arXiv preprint arXiv:2310.06452, 2023.
- [8] T. Chu et al., “Sft memorizes, rl generalizes: A comparative study of foundation model post-training,” arXiv preprint arXiv:2501.17161, 2025.
- [9] P. Singhal, T. Goyal, J. Xu, and G. Durrett, “A long way to go: Investigating length correlations in rlhf,” arXiv preprint arXiv:2310.03716, 2023.
- [10] R. Park, R. Rafailov, S. Ermon, and C. Finn, “Disentangling length from quality in direct preference optimization,” arXiv preprint arXiv:2403.19159, 2024.
- [11] Allen Institute for Artificial Intelligence, “OLMoE, meet iOS.” <https://alle-nai.org/blog/olmoe-app>, 2025.
- [12] C. Zhou et al., “Lima: Less is more for alignment,” Advances in Neural Information Processing Systems, vol. 36, pp. 55006–55021, 2023.
- [13] R. Taori et al., “Stanford alpaca: An instruction-following LLaMA model,” GitHub repository. https://github.com/tatsu-lab/stanford_alpaca; GitHub, 2023.
- [14] W.-L. Chiang et al., “Vicuna: An open-source chatbot impressing GPT-4 with 90%* ChatGPT quality.” 2023. Available: <https://lmsys.org/blog/2023-03-30-vicuna/>
- [15] X. Geng et al., “Koala: A dialogue model for academic research.” Blog post, 2023. Accessed: Apr. 03, 2023. [Online]. Available: <https://bair.berkeley.edu/blog/2023/04/03/koala/>

- [16] M. Conover et al., “Hello dolly: Democratizing the magic of ChatGPT with open models.” Accessed: Jun. 30, 2023. [Online]. Available: <https://www.databricks.com/blog/2023/03/24/hello-dolly-democratizing-magic-chatgpt-open-models.html>
- [17] A. Askell et al., “A general language assistant as a laboratory for alignment,” arXiv preprint arXiv:2112.00861, 2021.
- [18] Y. Bai et al., “Constitutional ai: Harmlessness from ai feedback,” arXiv preprint arXiv:2212.08073, 2022.
- [19] R. Rafailov, A. Sharma, E. Mitchell, C. D. Manning, S. Ermon, and C. Finn, “Direct preference optimization: Your language model is secretly a reward model,” Advances in Neural Information Processing Systems, vol. 36, 2024.
- [20] L. Tunstall et al., “Zephyr: Direct distillation of LM alignment,” in First conference on language modeling, 2024. Available: <https://openreview.net/forum?id=aKkAwZB6JV>
- [21] H. Ivison et al., “Camels in a changing climate: Enhancing lm adaptation with tulu 2,” arXiv preprint arXiv:2311.10702, 2023.
- [22] G. Cui et al., “Ultrafeedback: Boosting language models with high-quality feedback,” 2023.
- [23] A. Dubey et al., “The llama 3 herd of models,” arXiv preprint arXiv:2407.21783, 2024.
- [24] B. Adler et al., “Nemotron-4 340B technical report,” arXiv preprint arXiv:2406.11704, 2024.
- [25] C. Wirth, R. Akroud, G. Neumann, and J. Fürnkranz, “A survey of preference-based reinforcement learning methods,” Journal of Machine Learning Research, vol. 18, no. 136, pp. 1–46, 2017.
- [26] T. Kaufmann, P. Weng, V. Bengs, and E. Hüllermeier, “A survey of reinforcement learning from human feedback,” arXiv preprint arXiv:2312.14925, 2023.
- [27] S. Casper et al., “Open problems and fundamental limitations of reinforcement learning from human feedback,” arXiv preprint arXiv:2307.15217, 2023.
- [28] W. B. Knox and P. Stone, “Tamer: Training an agent manually via evaluative reinforcement,” in 2008 7th IEEE international conference on development and learning, IEEE, 2008, pp. 292–297.
- [29] J. MacGlashan et al., “Interactive learning from policy-dependent human feedback,” in International conference on machine learning, PMLR, 2017, pp. 2285–2294.

- [30] B. Ibarz, J. Leike, T. Pohlen, G. Irving, S. Legg, and D. Amodei, “Reward learning from human preferences and demonstrations in atari,” Advances in neural information processing systems, vol. 31, 2018.
- [31] G. Warnell, N. Waytowich, V. Lawhern, and P. Stone, “Deep tamer: Interactive agent shaping in high-dimensional state spaces,” in Proceedings of the AAAI conference on artificial intelligence, 2018.
- [32] J. Leike, D. Krueger, T. Everitt, M. Martic, V. Maini, and S. Legg, “Scalable agent alignment via reward modeling: A research direction,” arXiv preprint arXiv:1811.07871, 2018.
- [33] D. M. Ziegler et al., “Fine-tuning language models from human preferences,” arXiv preprint arXiv:1909.08593, 2019.
- [34] J. Wu et al., “Recursively summarizing books with human feedback,” arXiv preprint arXiv:2109.10862, 2021.
- [35] J. Menick et al., “Teaching language models to support answers with verified quotes,” arXiv preprint arXiv:2203.11147, 2022.
- [36] A. Glaese et al., “Improving alignment of dialogue agents via targeted human judgements,” arXiv preprint arXiv:2209.14375, 2022.
- [37] L. Gao, J. Schulman, and J. Hilton, “Scaling laws for reward model overoptimization,” in International conference on machine learning, PMLR, 2023, pp. 10835–10866.
- [38] D. Ganguli et al., “Red teaming language models to reduce harms: Methods, scaling behaviors, and lessons learned,” arXiv preprint arXiv:2209.07858, 2022.
- [39] R. Ramamurthy et al., “Is reinforcement learning (not) for natural language processing: Benchmarks, baselines, and building blocks for natural language policy optimization,” arXiv preprint arXiv:2210.01241, 2022.
- [40] A. Havrilla et al., “TrlX: A framework for large scale reinforcement learning from human feedback,” in Proceedings of the 2023 conference on empirical methods in natural language processing, Singapore: Association for Computational Linguistics, Dec. 2023, pp. 8578–8595. doi: 10.18653/v1/2023.emnlp-main.530.
- [41] L. von Werra et al., “TRL: Transformer reinforcement learning,” GitHub repository. <https://github.com/huggingface/trl>; GitHub, 2020.
- [42] OpenAI, “ChatGPT: Optimizing language models for dialogue.” <https://openai.com/blog/chatgpt/>, 2022.
- [43] H. Touvron et al., “Llama 2: Open foundation and fine-tuned chat models,” arXiv preprint arXiv:2307.09288, 2023.
- [44] H. Lightman et al., “Let’s verify step by step,” arXiv preprint arXiv:2305.20050, 2023.

- [45] A. Kumar et al., “Training language models to self-correct via reinforcement learning,” arXiv preprint arXiv:2409.12917, 2024.
- [46] A. Singh et al., “Beyond human data: Scaling self-training for problem-solving with language models,” arXiv preprint arXiv:2312.06585, 2023.
- [47] OpenAI, “Introducing OpenAI o1-preview.” Sep. 2024. Available: <https://openai.com/index/introducing-openai-o1-preview/>
- [48] A. Vaswani et al., “Attention is all you need,” in Neural information processing systems, 2017. Available: <https://api.semanticscholar.org/CorpusID:13756489>
- [49] D. Bahdanau, K. Cho, and Y. Bengio, “Neural machine translation by jointly learning to align and translate,” CoRR, vol. abs/1409.0473, 2014, Available: <https://api.semanticscholar.org/CorpusID:11212020>
- [50] G. Hinton, O. Vinyals, and J. Dean, “Distilling the knowledge in a neural network,” arXiv preprint arXiv:1503.02531, 2015.
- [51] G. Team et al., “Gemma 2: Improving open language models at a practical size,” arXiv preprint arXiv:2408.00118, 2024.
- [52] R. Agarwal et al., “On-policy distillation of language models: Learning from self-generated mistakes,” in The twelfth international conference on learning representations, 2024.
- [53] J. Wei et al., “Chain-of-thought prompting elicits reasoning in large language models,” Advances in neural information processing systems, vol. 35, pp. 24824–24837, 2022.
- [54] R. S. Sutton, “Reinforcement learning: An introduction,” A Bradford Book, 2018.
- [55] N. Lambert, L. Castricato, L. von Werra, and A. Havrilla, “Illustrating reinforcement learning from human feedback (RLHF),” Hugging Face Blog, 2022.
- [56] N. Lambert, T. K. Gilbert, and T. Zick, “Entangled preferences: The history and risks of reinforcement learning and human feedback,” arXiv preprint arXiv:2310.13595, 2023.
- [57] V. Conitzer et al., “Social choice should guide AI alignment in dealing with diverse human feedback,” arXiv preprint arXiv:2404.10271, 2024.
- [58] A. Mishra, “Ai alignment and social choice: Fundamental limitations and policy implications,” arXiv preprint arXiv:2310.16048, 2023.
- [59] H. R. Kirk et al., “The PRISM alignment project: What participatory, representative and individualised human feedback reveals about the subjective and multicultural alignment of large language models,” arXiv preprint arXiv:2404.16019, 2024.

- [60] S. Poddar, Y. Wan, H. Ivison, A. Gupta, and N. Jaques, “Personalizing reinforcement learning from human feedback with variational preference learning,” arXiv preprint arXiv:2408.10075, 2024.
- [61] S. J. Russell and P. Norvig, Artificial intelligence: A modern approach. Pearson, 2016.
- [62] B. Widrow and M. E. Hoff, “Adaptive switching circuits,” Stanford Univ Ca Stanford Electronics Labs, 1960.
- [63] B. F. Skinner, The behavior of organisms: An experimental analysis. BF Skinner Foundation, 2019.
- [64] E. L. Thorndike, “The law of effect,” The American journal of psychology, vol. 39, no. 1/4, pp. 212–222, 1927.
- [65] A. Arnauld, The port-royal logic. 1662.
- [66] J. Bentham, An introduction to the principles of morals and legislation. 1823.
- [67] F. P. Ramsey, “Truth and probability,” Readings in Formal Epistemology: Sourcebook, pp. 21–45, 2016.
- [68] K. J. Arrow, “A difficulty in the concept of social welfare,” Journal of political economy, vol. 58, no. 4, pp. 328–346, 1950.
- [69] J. C. Harsanyi, “Rule utilitarianism and decision theory,” Erkenntnis, vol. 11, no. 1, pp. 25–53, 1977.
- [70] R. Pettigrew, Choosing for changing selves. Oxford University Press, 2019.
- [71] N. Soares, B. Fallenstein, S. Armstrong, and E. Yudkowsky, “Corrigibility,” in Workshops at the twenty-ninth AAAI conference on artificial intelligence, 2015.
- [72] W.-L. Chiang et al., “Chatbot arena: An open platform for evaluating llms by human preference,” arXiv preprint arXiv:2403.04132, 2024.
- [73] R. Likert, “A technique for the measurement of attitudes.” Archives of psychology, 1932.
- [74] J. Zhou et al., “Instruction-following evaluation for large language models,” arXiv preprint arXiv:2311.07911, 2023.
- [75] K. Ethayarajh, W. Xu, N. Muennighoff, D. Jurafsky, and D. Kiela, “Kto: Model alignment as prospect theoretic optimization,” arXiv preprint arXiv:2402.01306, 2024.
- [76] Z. Wu et al., “Fine-grained human feedback gives better rewards for language model training,” Advances in Neural Information Processing Systems, vol. 36, 2024.
- [77] A. Chen et al., “Learning from natural language feedback,” Transactions on Machine Learning Research, 2024.
- [78] OpenAI, “Introducing the model spec.” May 2024. Available: <https://openai.com/index/introducing-the-model-spec/>

- [79] A. Y. Ng, S. Russell, et al., “Algorithms for inverse reinforcement learning.” in Proceedings of the seventeenth international conference on machine learning, in ICML ’ 00. 2000, pp. 663–670.
- [80] R. A. Bradley and M. E. Terry, “Rank analysis of incomplete block designs: I. The method of paired comparisons,” *Biometrika*, vol. 39, no. 3/4, pp. 324–345, 1952, Accessed: Feb. 13, 2023. [Online]. Available: <http://www.jstor.org/stable/2334029>
- [81] B. Zhu et al., “Starling-7b: Improving helpfulness and harmlessness with rlaif,” in First conference on language modeling, 2024.
- [82] A. Liu, Z. Zhao, C. Liao, P. Lu, and L. Xia, “Learning plackett-luce mixtures from partial preferences,” in Proceedings of the AAAI conference on artificial intelligence, 2019, pp. 4328–4335.
- [83] B. Zhu, M. Jordan, and J. Jiao, “Principled reinforcement learning with human feedback from pairwise or k-wise comparisons,” in International conference on machine learning, PMLR, 2023, pp. 43037–43067.
- [84] K. Cobbe et al., “Training verifiers to solve math word problems,” arXiv preprint arXiv:2110.14168, 2021.
- [85] C. Lyu et al., “Exploring the limit of outcome reward for learning mathematical reasoning,” arXiv preprint arXiv:2502.06781, 2025.
- [86] L. Zheng et al., “Judging llm-as-a-judge with mt-bench and chatbot arena,” Advances in Neural Information Processing Systems, vol. 36, pp. 46595–46623, 2023.
- [87] Y. Dubois, B. Galambosi, P. Liang, and T. B. Hashimoto, “Length-controlled alpacaeval: A simple way to debias automatic evaluators,” arXiv preprint arXiv:2404.04475, 2024.
- [88] T. Li et al., “From crowdsourced data to high-quality benchmarks: Arena-hard and BenchBuilder pipeline,” arXiv preprint arXiv:2406.11939, 2024.
- [89] B. Y. Lin et al., “WILDBENCH: Benchmarking LLMs with challenging tasks from real users in the wild,” arXiv preprint arXiv:2406.04770, 2024.
- [90] D. Mahan et al., “Generative reward models,” 2024, Available: https://www.synthlabs.ai/pdf/Generative_Reward_Models.pdf
- [91] L. Zhang, A. Hosseini, H. Bansal, M. Kazemi, A. Kumar, and R. Agarwal, “Generative verifiers: Reward modeling as next-token prediction,” arXiv preprint arXiv:2408.15240, 2024.
- [92] Z. Ankner, M. Paul, B. Cui, J. D. Chang, and P. Ammanabrolu, “Critique-out-loud reward models,” arXiv preprint arXiv:2408.11791, 2024.
- [93] S. Kim et al., “Prometheus: Inducing fine-grained evaluation capability in language models,” in The twelfth international conference on learning representations, 2023.

- [94] N. Lambert et al., “Rewardbench: Evaluating reward models for language modeling,” arXiv preprint arXiv:2403.13787, 2024.
- [95] X. Wen et al., “Rethinking reward model evaluation: Are we barking up the wrong tree?” arXiv preprint arXiv:2410.05584, 2024.
- [96] S. Gureja et al., “M-RewardBench: Evaluating reward models in multilingual settings,” arXiv preprint arXiv:2410.15522, 2024.
- [97] Z. Jin et al., “RAG-RewardBench: Benchmarking reward models in retrieval augmented generation for preference alignment,” arXiv preprint arXiv:2412.13746, 2024.
- [98] E. Zhou et al., “RMB: Comprehensively benchmarking reward models in LLM alignment,” arXiv preprint arXiv:2410.09893, 2024.
- [99] Y. Liu, Z. Yao, R. Min, Y. Cao, L. Hou, and J. Li, “RM-bench: Benchmarking reward models of language models with subtlety and style,” arXiv preprint arXiv:2410.16184, 2024.
- [100] Z. Wu, M. Yasunaga, A. Cohen, Y. Kim, A. Celikyilmaz, and M. Ghazvininejad, “reWordBench: Benchmarking and improving the robustness of reward models with transformed inputs,” arXiv preprint arXiv:2503.11751, 2025.
- [101] Z. Chen et al., “MJ-bench: Is your multimodal reward model really a good judge for text-to-image generation?” arXiv preprint arXiv:2407.04842, 2024.
- [102] M. Yasunaga, L. Zettlemoyer, and M. Ghazvininejad, “Multimodal reward-bench: Holistic evaluation of reward models for vision language models,” arXiv preprint arXiv:2502.14191, 2025.
- [103] L. Li et al., “VLRewardBench: A challenging benchmark for vision-language generative reward models,” arXiv preprint arXiv:2411.17451, 2024.
- [104] J. Ruan et al., “Vlrbench: A comprehensive and challenging benchmark for vision-language reward models,” arXiv preprint arXiv:2503.07478, 2025.
- [105] E. Frick et al., “How to evaluate reward models for RLHF,” arXiv preprint arXiv:2410.14872, 2024.
- [106] S. Kim et al., “Evaluating robustness of reward models for mathematical reasoning,” arXiv preprint arXiv:2410.01729, 2024.
- [107] M. Song, Z. Su, X. Qu, J. Zhou, and Y. Cheng, “PRMBench: A fine-grained and challenging benchmark for process-level reward models,” arXiv preprint arXiv:2501.03124, 2025.
- [108] W. Wang et al., “VisualPRM: An effective process reward model for multi-modal reasoning,” arXiv preprint arXiv:2503.10291, 2025.
- [109] H. Tu, W. Feng, H. Chen, H. Liu, X. Tang, and C. Xie, “ViLBench: A suite for vision-language process reward modeling.” Mar. 2025. Available: <https://arxiv.org/abs/2503.20271>

- [110] H. Wang, W. Xiong, T. Xie, H. Zhao, and T. Zhang, “Interpretable preferences via multi-objective reward modeling and mixture-of-experts,” arXiv preprint arXiv:2406.12845, 2024.
- [111] Z. Wang et al., “HelpSteer2: Open-source dataset for training top-performing reward models,” arXiv preprint arXiv:2406.08673, 2024.
- [112] Z. Wang et al., “HelpSteer2-preference: Complementing ratings with preferences,” arXiv preprint arXiv:2410.01257, 2024.
- [113] J. Park, S. Jwa, M. Ren, D. Kim, and S. Choi, “Offsetbias: Leveraging debiased data for tuning evaluators,” arXiv preprint arXiv:2407.06551, 2024.
- [114] N. Jaques, S. Gu, D. Bahdanau, J. M. Hernández-Lobato, R. E. Turner, and D. Eck, “Sequence tutor: Conservative fine-tuning of sequence generation models with kl-control,” in International conference on machine learning, PMLR, 2017, pp. 1645–1654.
- [115] N. Jaques et al., “Human-centric dialog training via offline reinforcement learning,” arXiv preprint arXiv:2010.05848, 2020.
- [116] J. Schulman, “Approximating KL-divergence.” <http://joschu.net/blog/kl-approx.html>, 2016.
- [117] R. Y. Pang, W. Yuan, K. Cho, H. He, S. Sukhbaatar, and J. Weston, “Iterative reasoning preference optimization,” arXiv preprint arXiv:2404.19733, 2024.
- [118] Z. Gao et al., “Rebel: Reinforcement learning via regressing relative rewards,” arXiv preprint arXiv:2404.16767, 2024.
- [119] T. B. Brown et al., “Language models are few-shot learners,” arXiv preprint arXiv:2005.14165, 2020.
- [120] C. Raffel et al., “Exploring the limits of transfer learning with a unified text-to-text transformer,” Journal of machine learning research, vol. 21, no. 140, pp. 1–67, 2020.
- [121] J. Wei et al., “Finetuned language models are zero-shot learners,” in International conference on learning representations, 2022. Available: <https://openreview.net/forum?id=gEZrGCozdqR>
- [122] V. Sanh et al., “Multitask prompted training enables zero-shot task generalization,” in International conference on learning representations, 2022. Available: <https://openreview.net/forum?id=9Vrb9D0WI4>
- [123] S. Mishra, D. Khashabi, C. Baral, and H. Hajishirzi, “Cross-task generalization via natural language crowdsourcing instructions,” in Proceedings of the 60th annual meeting of the association for computational linguistics (volume 1: Long papers), Association for Computational Linguistics, May 2022, pp. 3470–3487. doi: 10.18653/v1/2022.acl-long.244.

- [124] E. Wallace, K. Xiao, R. Leike, L. Weng, J. Heidecke, and A. Beutel, “The instruction hierarchy: Training llms to prioritize privileged instructions,” arXiv preprint arXiv:2404.13208, 2024.
- [125] T. Dettmers, A. Pagnoni, A. Holtzman, and L. Zettlemoyer, “Qlora: Efficient finetuning of quantized llms,” Advances in neural information processing systems, vol. 36, pp. 10088–10115, 2023.
- [126] N. Rajani, L. Tunstall, E. Beeching, N. Lambert, A. M. Rush, and T. Wolf, “No robots,” Hugging Face repository. https://huggingface.co/datasets/HuggingFaceH4/no_robots; Hugging Face, 2023.
- [127] W. R. Gilks and P. Wild, “Adaptive rejection sampling for gibbs sampling,” Journal of the Royal Statistical Society: Series C (Applied Statistics), vol. 41, no. 2, pp. 337–348, 1992.
- [128] A. Ahmadian et al., “Back to basics: Revisiting reinforce style optimization for learning from human feedback in llms,” arXiv preprint arXiv:2402.14740, 2024.
- [129] J. Schulman, P. Moritz, S. Levine, M. Jordan, and P. Abbeel, “High-dimensional continuous control using generalized advantage estimation,” in Proceedings of the international conference on learning representations (ICLR), 2016.
- [130] R. J. Williams, “Simple statistical gradient-following algorithms for connectionist reinforcement learning,” Machine learning, vol. 8, pp. 229–256, 1992.
- [131] S. C. Huang, A. Ahmadian, and C. F. AI, “Putting RL back in RLHF.” https://huggingface.co/blog/putting_rl_back_in_rlfh_with_rloo, 2024.
- [132] W. Kool, H. van Hoof, and M. Welling, “Buy 4 reinforce samples, get a baseline for free!” 2019.
- [133] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, “Proximal policy optimization algorithms,” arXiv preprint arXiv:1707.06347, 2017.
- [134] C. Berner et al., “Dota 2 with large scale deep reinforcement learning,” arXiv preprint arXiv:1912.06680, 2019.
- [135] Z. Liu et al., “Understanding R1-zero-like training: A critical perspective,” arXiv preprint arXiv:2503.20783, Mar. 2025, Available: <https://arxiv.org/abs/2503.20783>
- [136] Z. Shao et al., “Deepseekmath: Pushing the limits of mathematical reasoning in open language models,” arXiv preprint arXiv:2402.03300, 2024.
- [137] A. Liu et al., “Deepseek-v3 technical report,” arXiv preprint arXiv:2412.19437, 2024.
- [138] D. Guo et al., “Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning,” arXiv preprint arXiv:2501.12948, 2025.

- [139] S. Huang, M. Noukhovitch, A. Hosseini, K. Rasul, W. Wang, and L. Tunstall, “The n+ implementation details of RLHF with PPO: A case study on TL;DR summarization,” in First conference on language modeling, 2024. Available: <https://openreview.net/forum?id=kHO2ZTa8e3>
- [140] L. Weng, “Policy gradient algorithms,” lilianweng.github.io, 2018, Available: <https://lilianweng.github.io/posts/2018-04-08-policy-gradient/>
- [141] Y. Zhao, R. Joshi, T. Liu, M. Khalman, M. Saleh, and P. J. Liu, “Slic-hf: Sequence likelihood calibration with human feedback,” arXiv preprint arXiv:2305.10425, 2023.
- [142] M. G. Azar et al., “A general theoretical paradigm to understand learning from human preferences,” in International conference on artificial intelligence and statistics, PMLR, 2024, pp. 4447–4455.
- [143] A. Amini, T. Vieira, and R. Cotterell, “Direct preference optimization with an offset,” arXiv preprint arXiv:2402.10571, 2024.
- [144] J. Hong, N. Lee, and J. Thorne, “Reference-free monolithic preference optimization with odds ratio,” arXiv e-prints, pp. arXiv–2403, 2024.
- [145] Y. Meng, M. Xia, and D. Chen, “Simpo: Simple preference optimization with a reference-free reward,” Advances in Neural Information Processing Systems, vol. 37, pp. 124198–124235, 2025.
- [146] N. Razin, S. Malladi, A. Bhaskar, D. Chen, S. Arora, and B. Hanin, “Unintentional unalignment: Likelihood displacement in direct preference optimization,” arXiv preprint arXiv:2410.08847, 2024.
- [147] Y. Ren and D. J. Sutherland, “Learning dynamics of llm finetuning,” arXiv preprint arXiv:2407.10490, 2024.
- [148] T. Xiao, Y. Yuan, H. Zhu, M. Li, and V. G. Honavar, “Cal-dpo: Calibrated direct preference optimization for language model alignment,” arXiv preprint arXiv:2412.14516, 2024.
- [149] A. Gupta et al., “AlphaPO-reward shape matters for LLM alignment,” arXiv preprint arXiv:2501.03884, 2025.
- [150] S. Guo et al., “Direct language model alignment from online ai feedback,” arXiv preprint arXiv:2402.04792, 2024.
- [151] P. Singhal, N. Lambert, S. Niekuum, T. Goyal, and G. Durrett, “D2po: Discriminator-guided dpo with response evaluation models,” arXiv preprint arXiv:2405.01511, 2024.
- [152] C. Rosset, C.-A. Cheng, A. Mitra, M. Santacroce, A. Awadallah, and T. Xie, “Direct nash optimization: Teaching language models to self-improve with general preferences,” arXiv preprint arXiv:2404.03715, 2024.
- [153] S. Jung, G. Han, D. W. Nam, and K.-W. On, “Binary classifier optimization for large language model alignment,” arXiv preprint arXiv:2404.04656, 2024.

- [154] H. Zhao et al., “Rainbowpo: A unified framework for combining improvements in preference optimization,” arXiv preprint arXiv:2410.04203, 2024.
- [155] A. Gorbatovski, B. Shaposhnikov, V. Sinii, A. Malakhov, and D. Gavrilov, “The differences between direct alignment algorithms are a blur,” arXiv preprint arXiv:2502.01237, 2025.
- [156] H. Ivison et al., “Unpacking DPO and PPO: Disentangling best practices for learning from preference feedback,” arXiv preprint arXiv:2406.09279, 2024.
- [157] S. Xu et al., “Is dpo superior to ppo for llm alignment? A comprehensive study,” arXiv preprint arXiv:2404.10719, 2024.
- [158] F. Tajwar et al., “Preference fine-tuning of llms should leverage suboptimal, on-policy data,” arXiv preprint arXiv:2404.14367, 2024.
- [159] H. Lee et al., “Rlaif: Scaling reinforcement learning from human feedback with ai feedback,” 2023.
- [160] A. Sharma, S. Keh, E. Mitchell, C. Finn, K. Arora, and T. Kollar, “A critical evaluation of AI feedback for aligning large language models.” 2024. Available: <https://arxiv.org/abs/2402.12366>
- [161] L. Castricato, N. Lile, S. Anand, H. Schoelkopf, S. Verma, and S. Biderman, “Suppressing pink elephants with direct principle feedback.” 2024. Available: <https://arxiv.org/abs/2402.07896>
- [162] L. J. V. Miranda et al., “Hybrid preferences: Learning to route instances for human vs. AI feedback,” arXiv preprint arXiv:2410.19133, 2024.
- [163] T. Wang et al., “Shepherd: A critic for language model generation,” arXiv preprint arXiv:2308.04592, 2023.
- [164] P. Ke et al., “CritiqueLLM: Towards an informative critique generation model for evaluation of large language model generation,” arXiv preprint arXiv:2311.18702, 2023.
- [165] J. Li, S. Sun, W. Yuan, R.-Z. Fan, H. Zhao, and P. Liu, “Generative judge for evaluating alignment,” arXiv preprint arXiv:2310.05470, 2023.
- [166] S. Kim et al., “Prometheus 2: An open source language model specialized in evaluating other language models,” arXiv preprint arXiv:2405.01535, 2024.
- [167] S. Lee, S. Kim, S. Park, G. Kim, and M. Seo, “Prometheus-vision: Vision-language model as a judge for fine-grained evaluation,” in Findings of the association for computational linguistics ACL 2024, 2024, pp. 11286–11315.
- [168] M. Y. Guan et al., “Deliberative alignment: Reasoning enables safer language models,” arXiv preprint arXiv:2412.16339, 2024.
- [169] Anthropic, “Claude’ s constitution.” Accessed: Feb. 07, 2024. [Online]. Available: <https://www.anthropic.com/news/claudes-constitution>
- [170] D. Ganguli et al., “Collective constitutional AI: Aligning a language model with public input.” Anthropic, 2023.

- [171] S. Huang et al., “Constitutional AI recipe,” Hugging Face Blog, 2024.
- [172] N. Lambert, H. Schoelkopf, A. Gokaslan, L. Soldaini, V. Pyatkin, and L. Casstricato, “Self-directed synthetic dialogues and revisions technical report,” arXiv preprint arXiv:2407.18421, 2024.
- [173] Z. Sun et al., “Principle-driven self-alignment of language models from scratch with minimal human supervision,” in Thirty-seventh conference on neural information processing systems, 2023. Available: <https://openreview.net/forum?id=p40XRfBX96>
- [174] Z. Sun et al., “SALMON: Self-alignment with principle-following reward models,” in The twelfth international conference on learning representations, 2024. Available: <https://openreview.net/forum?id=xJbsmB8UMx>
- [175] A. Irpan, “Deep reinforcement learning doesn’t work yet.” 2018. Available: <https://www.alexirpan.com/2018/02/14/rl-hard.html>
- [176] P. Henderson, R. Islam, P. Bachman, J. Pineau, D. Precup, and D. Meger, “Deep reinforcement learning that matters,” in Proceedings of the AAAI conference on artificial intelligence, 2018. Available: <https://ojs.aaai.org/index.php/AAAI/article/view/11694>
- [177] G. Sheng et al., “HybridFlow: A flexible and efficient RLHF framework,” arXiv preprint arXiv: 2409.19256, 2024.
- [178] J. Hu et al., “OpenRLHF: An easy-to-use, scalable and high-performance RLHF framework,” arXiv preprint arXiv:2405.11143, 2024.
- [179] J. Liu, A. Cohen, R. Pasunuru, Y. Choi, H. Hajishirzi, and A. Celikyilmaz, “Don’t throw away your value model! Generating more preferable text with value-guided monte-carlo tree search decoding,” arXiv preprint arXiv:2309.15028, 2023.
- [180] B. Brown et al., “Large language monkeys: Scaling inference compute with repeated sampling,” arXiv preprint arXiv:2407.21787, 2024.
- [181] Z. Liu et al., “Inference-time scaling for generalist reward modeling,” arXiv preprint arXiv:2504.02495, 2025.
- [182] N. Muennighoff et al., “s1: Simple test-time scaling,” arXiv preprint arXiv:2501.19393, 2025.
- [183] L. Chen et al., “Are more llm calls all you need? Towards scaling laws of compound inference systems,” arXiv preprint arXiv:2403.02419, 2024.
- [184] I. Shumailov, Z. Shumaylov, Y. Zhao, N. Papernot, R. Anderson, and Y. Gal, “AI models collapse when trained on recursively generated data,” Nature, vol. 631, no. 8022, pp. 755–759, 2024.
- [185] M. Gerstgrasser et al., “Is model collapse inevitable? Breaking the curse of recursion by accumulating real and synthetic data,” arXiv preprint arXiv:2404.01413, 2024.

- [186] Y. Feng, E. Dohmatob, P. Yang, F. Charton, and J. Kempe, “Beyond model collapse: Scaling up with synthesized data requires reinforcement,” in ICML 2024 workshop on theoretical foundations of foundation models, 2024.
- [187] Y. Wang et al., “Self-instruct: Aligning language models with self-generated instructions,” arXiv preprint arXiv:2212.10560, 2022.
- [188] E. Beeching et al., “NuminaMath 7B TIR,” Hugging Face repository. <https://huggingface.co/AI-MO/NuminaMath-7B-TIR>; Numina & Hugging Face, 2024.
- [189] M. Li et al., “Superfiltering: Weak-to-strong data filtering for fast instruction-tuning,” arXiv preprint arXiv:2402.00530, 2024.
- [190] K. Shridhar, A. Stolfo, and M. Sachan, “Distilling reasoning capabilities into smaller language models,” Findings of the Association for Computational Linguistics: ACL 2023, pp. 7059–7073, 2023.
- [191] C.-Y. Hsieh et al., “Distilling step-by-step! Outperforming larger language models with less training data and smaller model sizes,” arXiv preprint arXiv:2305.02301, 2023.
- [192] D. Hendrycks et al., “Measuring massive multitask language understanding,” arXiv preprint arXiv:2009.03300, 2020.
- [193] A. Mallen, A. Asai, V. Zhong, R. Das, H. Hajishirzi, and D. Khashabi, “When not to trust language models: Investigating effectiveness and limitations of parametric and non-parametric memories,” arXiv preprint, 2022.
- [194] S. Lin, J. Hilton, and O. Evans, “Truthfulqa: Measuring how models mimic human falsehoods,” arXiv preprint arXiv:2109.07958, 2021.
- [195] M. Suzgun et al., “Challenging BIG-bench tasks and whether chain-of-thought can solve them,” arXiv preprint arXiv:2210.09261, 2022.
- [196] D. Dua, Y. Wang, P. Dasigi, G. Stanovsky, S. Singh, and M. Gardner, “DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs,” arXiv preprint arXiv:1903.00161, 2019.
- [197] D. Hendrycks et al., “Measuring mathematical problem solving with the MATH dataset,” NeurIPS, 2021.
- [198] K. Cobbe et al., “Training verifiers to solve math word problems,” arXiv preprint arXiv:2110.14168, 2021.
- [199] M. Chen et al., “Evaluating large language models trained on code,” 2021, Available: <https://arxiv.org/abs/2107.03374>
- [200] J. Liu, C. S. Xia, Y. Wang, and L. Zhang, “Is your code generated by chat-GPT really correct? Rigorous evaluation of large language models for code generation,” in Thirty-seventh conference on neural information processing systems, 2023. Available: <https://openreview.net/forum?id=1qvx610Cu7>

- [201] J. Zhou et al., “Instruction-following evaluation for large language models.” 2023. Available: <https://arxiv.org/abs/2311.07911>
- [202] D. Rein et al., “GPQA: A graduate-level google-proof q&a benchmark,” arXiv preprint arXiv:2311.12022, 2023.
- [203] L. Phan, A. Gatti, Z. Han, N. Li, and H. et al. Zhang, “Humanity’ s last exam,” arXiv preprint arXiv:2501.14249, 2025.
- [204] R. Aleithan, H. Xue, M. M. Mohajer, E. Nnorom, G. Uddin, and S. Wang, “SWE-Bench+: Enhanced coding benchmark for LLMs,” arXiv preprint arXiv:2410.06992, 2024.
- [205] N. Jain et al., “LiveCodeBench: Holistic and contamination-free evaluation of large language models for code,” arXiv preprint arXiv:2403.07974, 2024.
- [206] S. AI, “SEAL LLM leaderboards: Expert-driven private evaluations.” 2024. Available: <https://scale.com/leaderboard>
- [207] S. Schulhoff et al., “The prompt report: A systematic survey of prompting techniques,” arXiv preprint arXiv:2406.06608, 2024.
- [208] J. Robinson, C. M. Rytting, and D. Wingate, “Leveraging large language models for multiple choice question answering,” in International conference on learning representations, 2023. Available: <https://openreview.net/forum?id=upQ4o-ygvJ>
- [209] J. Wei et al., “Finetuned language models are zero-shot learners,” in International conference on learning representations, 2022.
- [210] V. Sanh et al., “Multitask prompted training enables zero-shot task generalization,” in International conference on learning representations, 2022.
- [211] T. Kojima, S. S. Gu, M. Reid, Y. Matsuo, and Y. Iwasawa, “Large language models are zero-shot reasoners,” Advances in neural information processing systems, vol. 35, pp. 22199–22213, 2022.
- [212] J. Achiam et al., “Gpt-4 technical report,” arXiv preprint arXiv:2303.08774, 2023.
- [213] OpenAI, “Introducing SWE-bench verified.” Aug. 2024. Available: <https://openai.com/index/introducing-swe-bench-verified/>
- [214] J. Li et al., “Numinamath: The largest public dataset in ai4maths with 860k pairs of competition math problems and solutions,” Hugging Face repository, vol. 13, p. 9, 2024.
- [215] L. Yu et al., “Metamath: Bootstrap your own mathematical questions for large language models,” arXiv preprint arXiv:2309.12284, 2023.
- [216] A. K. Singh et al., “Evaluation data contamination in LLMs: How do we measure it and (when) does it matter?” arXiv preprint arXiv:2411.03923, 2024.
- [217] K. Huang et al., “MATH-perturb: Benchmarking LLMs’ math reasoning abilities against hard perturbations,” arXiv preprint arXiv:2502.06453, 2025.

- [218] UK AI Safety Institute, “Inspect AI: Framework for Large Language Model Evaluations.” https://github.com/UKGovernmentBEIS/inspect_ai, 2024.
- [219] C. Fourrier, N. Habib, H. Kydlicek, T. Wolf, and L. Tunstall, “LightEval: A lightweight framework for LLM evaluation.” <https://github.com/huggingface/lighteval>, 2023.
- [220] C. Fourrier, N. Habib, A. Lozovskaya, K. Szafer, and T. Wolf, “Open LLM leaderboard v2.” https://huggingface.co/spaces/open-llm-leaderboard/open_llm_leaderboard; Hugging Face, 2024.
- [221] L. Gao et al., “A Framework for Few-Shot Language Model Evaluation.” Zenodo, 2023. doi: 10.5281/zenodo.10256836.
- [222] S. Black et al., “GPT-NeoX-20B: An open-source autoregressive language model,” in Proceedings of the ACL workshop on challenges & perspectives in creating large language models, 2022. Available: <https://arxiv.org/abs/2204.06745>
- [223] Y. Gu, O. Tafjord, B. Kuehl, D. Haddad, J. Dodge, and H. Hajishirzi, “OLMES: A Standard for Language Model Evaluations,” arXiv preprint arXiv:2406.08446, 2024.
- [224] P. Liang et al., “Holistic Evaluation of Language Models,” Transactions on Machine Learning Research, 2023, doi: 10.1111/nyas.15007.
- [225] MosaicML, “Mosaic Eval Gauntlet v0.3.0 — Evaluation Suite.” https://github.com/mosaicml/llm-foundry/blob/main/scripts/eval/local_data/EVAL_GAUNTLET.md, 2024.
- [226] J. Schulman, “Proxy objectives in reinforcement learning from human feedback.” Invited talk at the International Conference on Machine Learning (ICML), 2023. Available: <https://icml.cc/virtual/2023/invited-talk/21549>
- [227] C. Zhang, O. Vinyals, R. Munos, and S. Bengio, “A study on overfitting in deep reinforcement learning,” arXiv preprint arXiv:1804.06893, 2018.
- [228] C. A. Goodhart and C. Goodhart, Problems of monetary management: The UK experience. Springer, 1984.
- [229] K. Hoskin, “The ‘awful idea of accountability’ : Inscribing people into the measurement of objects,” Accountability: Power, ethos and the technologies of managing, vol. 265, 1996.
- [230] M. Sharma et al., “Towards understanding sycophancy in language models,” arXiv preprint arXiv:2310.13548, 2023.
- [231] T. Lu and C. Boutilier, “Learning mallows models with pairwise preferences,” in Proceedings of the 28th international conference on machine learning (icml-11), 2011, pp. 145–152.
- [232] S. Han et al., “Wildguard: Open one-stop moderation tools for safety risks, jailbreaks, and refusals of llms,” arXiv preprint arXiv:2406.18495, 2024.

- [233] H. Inan et al., “Llama guard: Llm-based input-output safeguard for human-ai conversations,” arXiv preprint arXiv:2312.06674, 2023.
- [234] P. Röttger, H. R. Kirk, B. Vidgen, G. Attanasio, F. Bianchi, and D. Hovy, “Xtest: A test suite for identifying exaggerated safety behaviours in large language models,” arXiv preprint arXiv:2308.01263, 2023.
- [235] T. Coste, U. Anwar, R. Kirk, and D. Krueger, “Reward model ensembles help mitigate overoptimization,” arXiv preprint arXiv:2310.02743, 2023.
- [236] T. Moskovitz et al., “Confronting reward model overoptimization with constrained RLHF,” arXiv preprint arXiv:2310.04373, 2023.
- [237] R. Rafailov et al., “Scaling laws for reward model overoptimization in direct alignment algorithms,” Advances in Neural Information Processing Systems, vol. 37, pp. 126207–126242, 2024.
- [238] S. Zhuang and D. Hadfield-Menell, “Consequences of misaligned AI,” Advances in Neural Information Processing Systems, vol. 33, pp. 15763–15773, 2020.
- [239] W. Yuan et al., “Self-rewarding language models.” 2025. Available: <https://arxiv.org/abs/2401.10020>
- [240] J. Bai et al., “Qwen technical report,” arXiv preprint arXiv:2309.16609, 2023.
- [241] G. Wang, S. Cheng, X. Zhan, X. Li, S. Song, and Y. Liu, “Openchat: Advancing open-source language models with mixed-quality data,” arXiv preprint arXiv:2309.11235, 2023.
- [242] Anthropic, “Claude’s character.” 2024. Available: <https://www.anthropic.com/research/clause-character>