

# Online Reinforcement Learning

## 1. 语言模型的强化学习:



## 2. 在线强化学习工作机制：让模型自主探索更好的响应，流程如下

- a. 准备一批Prompt
- b. 将这批Prompt输入模型
- c. 模型生成相应的Response
- d. 将 (prompt,response) 对输入**奖励函数 (Reward Function)**
- e. 奖励函数为每对 (prompt,response) 打分
- f. 获得(prompt,response,reward) 三元组
- g. 使用这些数据来更新语言模型

更新语言模型的方法很多，本课重点为**PPO (Proximal Policy Optimization 近端策略优化)** 和 **GRPO (Group Relative Policy Optimization 群体相对策略优化)**

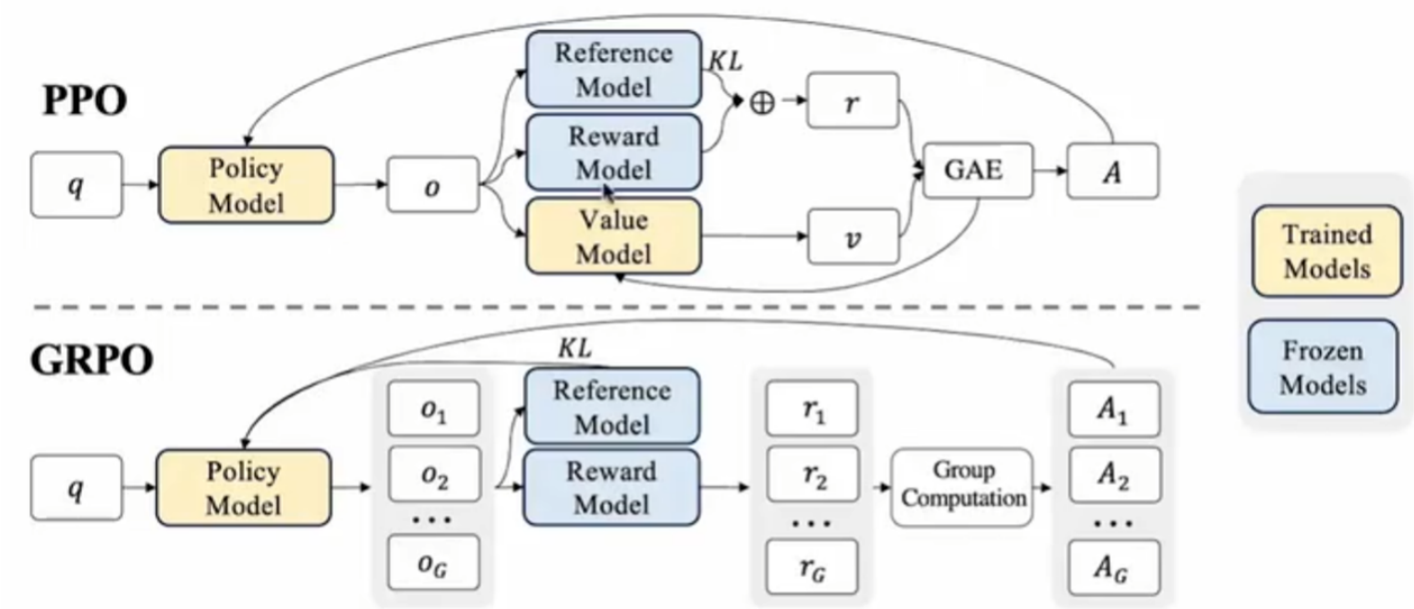
## 3. 奖励函数 (Reward Function)：

训练好的奖励模型  (Reward Model)	<ul style="list-style-type: none"><li>收集多个模型响应，由人类标注选择更优的响应</li><li>使用这些人类偏好数据训练奖励模型</li><li>奖励模型通过优化以下损失函数学习：<div><math display="block">L = \log(\sigma(r_j - r_k))</math></div><ul style="list-style-type: none"><li>若人类认为响应j优于k，则鼓励模型提升<math>r_j</math>，降低<math>r_k</math></li></ul></li></ul>	特点： <ul style="list-style-type: none"><li>通常基于已有的Instruct模型优化</li><li>通过大规模人类或机器生成偏好数据训练</li><li>可应用于开放式任务，如聊天能力、安全性提升等</li><li>在“正确性导向”的任务，如代码、数学、函数调用中可能不够精确</li></ul>
	更适用于“正确性导向”任务	特点：

可验证性奖励 (Verifiable Reward)	<ul style="list-style-type: none"><li>数学任务：验证输出是否与标准答案匹配</li><li>编程任务：通过<b>单元测试（Unit Tests）</b>检验代码执行结果是否正确</li></ul>	<ul style="list-style-type: none"><li>需提前准备<b>真值（Ground Truth）</b>或测试集</li><li>准备成本较高，但奖励信号更精确可靠</li><li>更适合训练推理类模型（Reasoning Models），如代码、数学领域</li></ul>
----------------------------------	---	--

4. PPO VS GRPO

Policy Training in Online RL



方法	工作流程	总结
PPO (Proximal Policy Optimization)	<ol style="list-style-type: none"><li>输入一组<b>查询（queries）</b>（<math>q</math>）</li><li>通过<b>策略模型（Policy Model）</b>（语言模型本身）生成响应</li><li>相应被送入以下3个模块：<ol style="list-style-type: none"><li><b>参考模型（Reference Model）</b>：计算KL散度，限制模型不偏离原始分布</li><li><b>奖励模型（Reward Model）</b>：计算奖励</li><li><b>价值模型（Value Model）</b>或<b>评价者模型（Critic Model）</b>：为每个token分配价值</li></ol></li></ol>	<p>PPO的目标函数：</p> $J_{PPO}(\theta) = E_{q \sim P(q), o \sim \pi_{\theta_{old}}(o   q)} \left[ \frac{1}{ o } \sum_{i=1}^{ o } \min \left[ \frac{\pi_{\theta}(o_i   q, o_{<i})}{\pi_{\theta_{old}}(o_i   q, o_{<i})} A_i, \text{clip} \left( \frac{\pi_{\theta}(o_i   q, o_{<i})}{\pi_{\theta_{old}}(o_i   q, o_{<i})}, 1 - \epsilon, 1 + \epsilon \right) A_i \right] \right]$ <ul style="list-style-type: none"><li>每个token拥有独立的优势值</li><li>反馈粒度更细</li><li>需额外训练价值模型——&gt;占用更多GPU内存</li></ul>

	<div>4. 使用<b>广义优势估计 (Generalized Advantage Estimation, GAE)</b> 来计算每个token的<b>优势函数 (Advantage)</b>，反映该token的贡献</div>	
GRPO (Group Relative Policy Optimization)	<div><div>1. 对每个prompt模型生成多个响应 (<math>O_1, O_2, \dots, O_g</math>)</div><div>2. 对每个响应计算<b>奖励 (Reward)</b> 和<b>参考模型的KL散度</b></div><div>3. 对同一组 (Group) 响应计算<b>相对奖励 (Relative Reward)</b></div><div>4. 将相对响应作为整个响应的优势值</div><div>5. 使用此优势更新策略模型</div></div>	<div><div>• 不再需要<b>价值模型 (Value Model)</b></div><div>• 所有token在同一响应中共享相同优势值</div><div>• 更节省显存，但优势估计较粗糙</div></div>

特征	PPO	GRPO
优势估计	基于价值模型 (Value Model) 的精细估计	基于响应组的相对奖励 (Relative Reward)
计算粒度	每个 Token 拥有独立优势	整个响应共享同一优势
显存需求	较高 (需训练 Critic)	较低 (无 Critic)
样本效率	高 (样本利用率好)	较低 (需更多样本)
奖励适配	适合连续或模型化奖励	适合二元/可验证奖励
应用场景	聊天、对齐、安全优化	数学、代码、推理任务

5. 代码解析

