# Logic-RL: Unleashing LLM Reasoning with Rule-Based Reinforcement Learning

**Lingkai Zu, Xiyue Peng, Liyu Yang,** ShanghaiTech University

## BACKGROUND

Improving the reasoning capabilities of language models to better assist in mathematical reasoning, code generation, and other real-world applications has become a prominent research focus. Although large language models have demonstrated remarkable performance in these tasks, the more successful cases are all large-scale models supported by substantial computational resources. Aligning smaller-scale models offers a promising approach to investigating language model alignment when computational resources are limited. This raises the question:
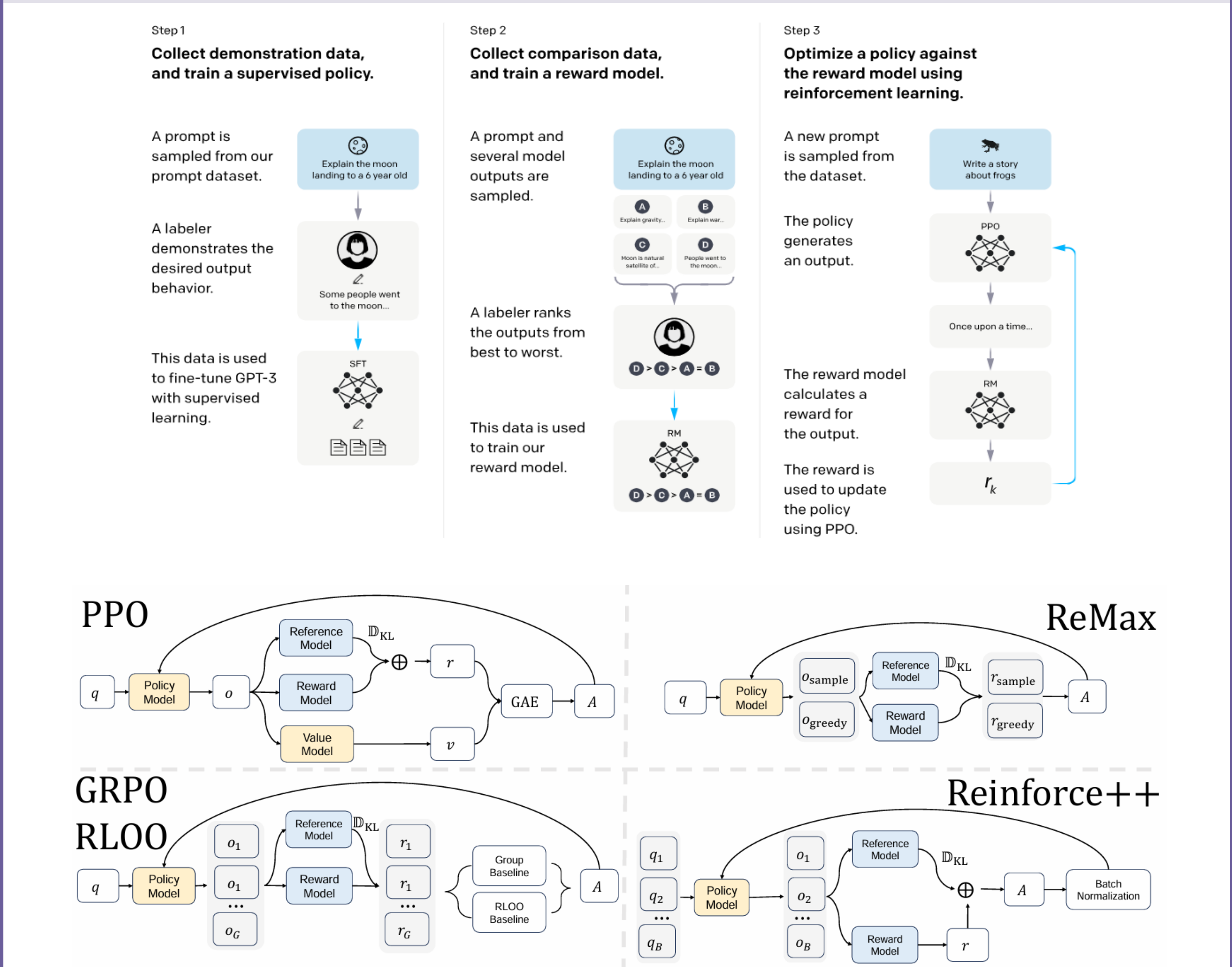
**Can similar reasoning abilities emerge in smaller-scale models?**

## SYNTHETIC LOGIC PUZZLES FRAMEWORK

| An example of a Knights and Knaves (K&K) puzzle | |
|---|---|
| Problem | Knights always tell the truth, and knaves always lie. You meet 2 inhabitants: Zoey, and Oliver. Zoey: "*Oliver is not a knight*". Oliver: "*Oliver is a knight if and only if Zoey is a knave*". So who is a knight and who is a knave? |
| Solution | (1)Zoey is a knave; (2)Oliver is a knight. |

| Reward Type | Reward Value | Condition |
|---|---|---|
| $s_{answer}$ | 2 | The answer fully matches the ground truth |
| | -1.5 | The answer partially matches the ground truth |
| | -2 | The answer cannot be parsed or missing |
| $s_{format}$ | 1 | The format is correct |
| | -1 | The format is incorrect |

## POST TRAINING METHOD OVERVIEW



## EXPERIMENTS RESULTS AND DISCUSSION

### Efficient Consumption Resource



Timing v.s. Step

GPU Memory Allocated (%)

### Poor Reasoning Capacity



Average Reward

All Correct Ratio

### Is Format Reward a Blessing or a Curse?



Average Length of Response

| Quiz | A very special island is inhabited only by knights and knaves. Knights always tell the truth, and knaves always lie. You meet 3 inhabitants: Scarlett, David, and Liam. According to Scarlett, "Liam is a knave if and only if Liam is a knight." David stated, "If Scarlett is a knave then Liam is a knave." Liam stated, "Scarlett is not a knave." So who is a knight and who is a knave? |
|---|---|
| **Algorithm** | **Response** |
| GRPO | \<think\>Scarlett is a knight, David is a knight, Liam is a knight \</think\>\<answer\> Scarlett is a knight, David is a knight, Liam is a knight \</answer\> |
| ReMax | \<think\>Reasoning process: \</think\>\<answer\> (1) Scarlett : knave, (2) David : knight, (3) Liam : knave \</answer\> |
| REINFORCE++ | \<think\>Identify doubts \</think\>\<answer\> penelope is a knave, aiden is a knight, amelia is a knave \</answer\> |
| RLOO | \<think\>\</think\>\<answer\> (1) Penelope is a knight, Aiden is a knight, Amelia is a knight. \</answer\> |

- Whether small-scale models possess reasoning capabilities?
- What kind of rewards should we design to elicit reasoning from small-scale models?

identifying and rewarding certain thinking tokens