# Language Model Self-improvement by Reinforcement Learning Contemplation without External Supervision

**Jing-Cheng Pang**[1,2]**, Kaiyuan Li**[1,2]**, Pengyuan Wang**[1,2]**, Xiong-Hui Chen**[1,2]**,**
**Jiacheng Xu**[1]**, Zongzhang Zhang**[1] **& Yang Yu**[1,2*]
[1] **National Key Laboratory for Novel Software Technology, Nanjing University,**
**China & School of Artificial Intelligence, Nanjing University, China;** [2]**Polixir.ai**

## Abstract

Language model self-improvement (LMSI) techniques have recently gained significant attention as they improve language models without requiring external supervision. A notable approach is reinforcement learning from AI feedback (RLAIF), which trains a reward model based on AI preference data and employs reinforcement learning (RL) algorithm to train the language model. However, RLAIF relies on a heuristic assumption that the AI model is able to provide effective feedback, which requires the language model to possess solid capability. In this paper, we present a novel LMSI method, Reinforcement Learning Contemplation (RLC). We disclose that it is simpler for language models to evaluate the text than to generate it, even for small models under 1B parameters. Leveraging the gap between the evaluation and generation, RLC evaluates the generated answers and updates language model using RL to maximize self-evaluation scores. We demonstrate the effectiveness of RLC on a wide range of challenging tasks, including reasoning, summarization, conditioned generation and emotion recognition, resulting in a increase in the answering accuracy ($31.23\% \rightarrow 37.09$) for BigBench-hard reasoning tasks, and a rise in BERTScore for CNN/Daily Mail summarization tasks. In addition, RLC can be applied to models of different sizes (80M to 3B), showcasing its broad applicability. We further verify that when training on larger scale of dataset, RLC improves language model's evaluation and generation ability on unseen tasks simultaneously, enabling a general capability improvement without external supervision.

## 1. Introduction

Scaling large language models (LLM) has led to remarkable performance in various natural language processing (NLP) tasks, such as language understanding, machine translation, and question answering (Zhao et al., 2023; Liu et al., 2023; Goldberg, 2016). However, training LLMs at scale often require extensive high-quality supervised datasets (Huang et al., 2023). To rectify this drawback, unsupervised training of LLMs has gained significant interest recently, offering the potential for training LLMs without relying on external supervision. This paper calls these methods *language model self-improvement* (LMSI) methods.

Two notable LMSI works are RLCAI (Bai et al., 2022) and RLAIF (Lee et al., 2024). RLCAI employs a model to evaluate which of the two answers better adheres to a set of rules or principles, subsequently creating a dataset of AI preferences. The preference data is then utilized to train a preference model, which provides a reward signal for training language models with reinforcement learning (RL). RLCAI aims to enhance specific properties of generated text, such as helpfulness and harmlessness. On the other hand, RLAIF follows a

---

∗. Corresponding: yuy@nju.edu.cn.

similar training process (i.e., reward model + RL) but empirically validates that AI feedback can perform comparably to human feedback in summarization tasks. However, these works are built based on the heuristic assumption that AI models can provide effective feedback or correct errors, which requires a strong capability of the language models and such capability is not sufficiently justified in their works. Besides, the requirements to collect preference data and train an extra reward model introduce added costs.

In contrast to these LMSI methods, humans possess the inherent ability for self-improvement through reflection, as indicated by contemplative science in human psychology (Van Gordon et al., 2022). For example, people can read and evaluate their articles and improve their writing skills. This human reflection process does not necessitate the explicit learning of an additional reward model to provide feedback. Thus we expect to mine the analogous self-improvement capabilities within LLM. In this paper, we propose a novel LMSI method, which is grounded on the observation that it is simpler for a language model to evaluate a sentence than to generate it. An intuitive example is that while writing an attractive story can be challenging, identifying the generated text is relatively easy. Fig. 1 illustrates such disparity between text generation and evaluation for language models. We comprehensively investigate the performance gap between text evaluation and text generation of language models across various NLP tasks. Leveraging this inherent performance gap within language models, we present **R**einforcement **L**earning **C**ontemplation (RLC). Given an unlabelled question set, the language model generates answers to unlabeled samples and directly evaluates the quality of these answers. The model is then updated using RL to optimize for maximum evaluation scores. Our proposed method employs self-evaluation results as rewards and utilizes reinforcement learning to optimize the language model. We refer to this process as *reinforcement learning contemplation*.
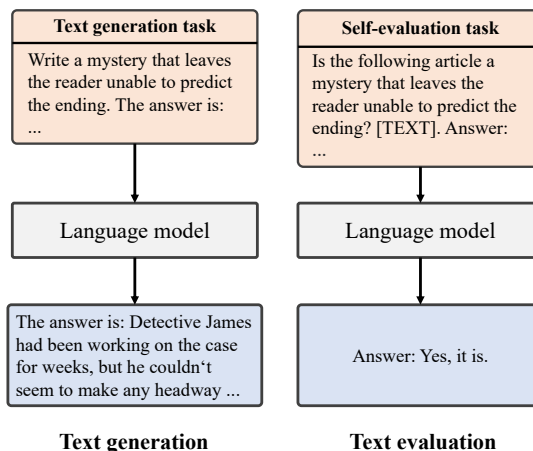


Figure 1: A comparison between text generation and evaluation. Text evaluation involves evaluating and analyzing certain properties of existing text, while generation requires the creation of entirely new text.

The contribution of this work can be summarized as follows: This work explains why the LMSI methods like RLAIF can work effectively. We suggest that the source of self-improvement comes from that text evaluation is simpler than generation for LLMs. This performance gap can be leveraged to boost the text generation of LLMs with various sizes. To our knowledge, this is the first work that tries to explore the underlying reasons for language model self-improvement. Besides, we propose a novel LMSI method for improving LLMs without external supervision, eliminating the need for training reward models. Lastly, we conduct comprehensive experiments to demonstrate that LLM, trained with RLC method, can self-improve its ability to solve a wide range of natural language processing problems. We also present that RLC can be applied to LLMs with a parameter range of 80M to 3B, and the trained LLM generalizes well to new and unseen datasets when scaling up the training dataset, demonstrating the extensive applicability of the proposed method.

## 2. Related Work

**Language model self-improvement.** LMSI refers to methods of training language models without any external supervision, which recently emerged as the the wind vane for enhancing the capability of LLM. These methods can be broadly classified into two categories, namely *noninvasive* and *invasive* methods. Noninvasive methods utilize AI feedback to guide the LLM to revise its output during the generation proc This approach involves a repeated process consisting of refinement and correction of the model's output, based on the LLM self-evaluation (Huang et al., 2023; Patel et al., 2024; Weng et al., 2023; Raman et al., 2022). For example, the re-prompting (Raman et al., 2022) method identifies errors in the current plan generated by the LLM and revises the current step if an error is detected. ReAct (Yao et al., 2022) and Self-Refine (Madaan et al., 2023) propose a specific design to correct the self-generated errors. On the other hand, invasive methods update the LLM parameters (Hu et al., 2024) using AI feedback as reward signals. Recent advancements in this field include RLAIF (Lee et al., 2024), which trains a reward model based on AI preference data and trains the LLM with RL algorithm. In contrast, prior studies relied mostly on comprehensive experiment results, this work explains why such a self-evaluation mechanism works well: we provide evidence that evaluation tasks are simpler than generation tasks for LLMs, and this difference can be leveraged to boost the performance of LLMs of various sizes.

**Train language model with unlabeled data.** Learning from unlabelled data is a promising approach for enhancing the model's efficiency. Self-training is a popular technique in this field, which entails assigning pseudo labels from a learned classifier to unlabeled data. These pseudo-labeled examples are then utilized to enhance the initial model training (RoyChowdhury et al., 2019; Chen et al., 2021; He et al., 2020). In recent years (Huang et al., 2023), self-training has been employed to fine-tune large-scale models such as PaLM (Chowdhery et al., 2023) with 540B parameters, and it has shown immense potential. However, it fails to present efficiency on LLMs with smaller sizes. In contrast to methods that utilize LLMs to generate pseudo-labels, our approach employs RL to train LLMs, which has proven more effective than supervised learning (Stiennon et al., 2020).

**Train language models with RL.** RL has demonstrated significant success in training language models (Stiennon et al., 2020; Jang et al., 2022; Erdem et al., 2022; Pang et al., 2023). For instance, some early studies train LLMs by utilizing algorithmically defined reward functions for specific NLP tasks, such as BLEU score for translation (Wu et al., 2016; Nguyen et al., 2017) and ROUGE score for summarization (Wu & Hu, 2018). In a departure from the heuristic definition of reward functions, another approach (Cho et al., 2018) involves learning to evaluate text coherence and utilizing the learned model to provide rewards for summarization and long-form generation tasks. Incorporating human preferences into language models using RL has recently gained popularity. A series of works (Ziegler et al., 2019; Ouyang et al., 2022; OpenAI, 2023) develop reward models reflecting human preferences and train language models using RL algorithms. While this work also uses RL to update the LLM, the difference is that our proposed method does not require a manually defined or trained reward model. In contrast, our method eliminates the need for external labels or reward models. Instead, reward signals come from the LLM self-evaluation results.

## 3. Preliminary and Problem Formulation

This section introduces the preliminary information and notations required for presenting our method and formally formulates the LMSI problem studied in this paper.

### 3.1 Preliminary

We begin with a vocabulary $\Sigma$ and an LLM $\mathcal{M}$ which takes a token sequence of the question $q = \{q_0, ..., q_n\}$ as input and predicts the next token using autoregressive modeling: $o_{t+1} = \mathcal{M}(q, \{o_0, \ldots, o_t\})$, where $q \in \Sigma^n$ and $o_t \in \Sigma$. To fine-tune LLM with RL, we can view this problem by Markov Decision Process (MDP) (Puterman, 1994; Sutton & Barto, 1998), which is described as a tuple $(\mathcal{S}, \mathcal{A}, P, \mathcal{R}, \gamma, d_0)$:

- State space $\mathcal{S}$: the space of input token sequences $q \cup \{o_0, o_1, \ldots, o_t\}$. $s_0$ is defined as the token sequence of question $q$.

- Action space $\mathcal{A}$: the space of tokens $o_t$.

- Reward function $\mathcal{R}(q, \{o_0, \ldots, o_t\})$: a score that reflects the quality of the generated answer to the question, which can be obtained from human feedback or a predefined criterion. The reward is typically given when the complete answer has been generated.

- Transition $P$: $s_{t+1} = s_t \cup o_{t+1}$.

- Initial distribution $d_0$: the distribution of question $q$.

Here, the LLM $\mathcal{M}$ acts as a policy mapping from state space to the probability space over action space. The objective of RL is to train the policy to maximize the expected return:

$$\mathbb{E}\left[ \sum_{t=0}^{\infty} \gamma^t \mathcal{R}(q, o) \big| q \sim d_0, o_t \sim \mathcal{M}(\cdot|s_t) \right]. \tag{1}$$

In RLC, we formulate the LLM text generation problem as an MDP, analogous to the one described above, with the reward function derived from evaluation results. To simplify the notation, we use $o \sim \mathcal{M}(q)$ to represent the autoregressive sampling of a complete answer $o$ from the language model $\mathcal{M}$, based on the input question $q$.

### 3.2 Problem Formulation

Let us consider a pre-trained language model $\mathcal{M}$ and an unlabeled question dataset $\mathcal{D}^{\text{train}} = \{q_i\}_{i=1}^{|\mathcal{D}|}$, where $q_i$ is denoted as a sequence of tokens, $q_i = \{o_i^0, o_i^1, \cdots, o_i^t\}$, indicating a textual question to answer. The objective is to update the pre-trained model $\mathcal{M}$ to maximize the answer quality to these questions, without the requirement of any external supervision. This objective can be formally expressed as:

$$\max_{\mathcal{M}} \mathbb{E}_{q \sim \mathcal{D}^{\text{train}}}[Q(q, o)|o \sim \mathcal{M}(q)], \tag{2}$$

where $Q$ stands for the quality metric for the answer $o$ to the question $q$, such as answer accuracy for the judgment question. The objective presented in Eq. 2 defines the task of improving the language model's performance on the provided dataset. More generally, we expect the answer quality on the questions out of the training question set can be improved:

$Q(q, \mathcal{M}(q)) \leq Q(q, \tilde{\mathcal{M}}(q))$ for $q \notin \mathcal{D}^{\mathrm{train}}$, where $\tilde{\mathcal{M}}$ denotes the updated language model. The underlying motivation behind this objective is the expectation that as the scale of training increases, the language model will exhibit general improvements across a diverse array of question domains. Such generalization can indicate that the improvements are not confined to the training data but extend to new, unseen questions, thereby increasing the model's utility in more broad applications.

## 4. Evaluation is Simpler than Generation

We first establish the fundamental idea that supports our method: text evaluation is simpler than text generation for LLMs. We dive into three key topics in subsequent subsections: (1) a comparison of LLMs' ability in text generation and evaluation; (2) the correlation between evaluation results and the established evaluation metrics; and (3) the potential for improving LLMs by self-evaluation.

### 4.1  Comparison of Evaluation and Generation

We conduct experiments to compare the text generation and evaluation abilities of LLMs using the CommonGen (Lin et al., 2020) task, which involves generating a sentence that describes an everyday scenario based on a given set of common concepts such as *{dog, frisbee, catch, and throw}*. Specifically, we use FLAN-T5 (Chung et al., 2022) as the LLM to generate text based on common conceptual questions provided by CommonGen and evaluate the accuracy of the generated text through human evaluation. In addition, we assess the LLM's text evaluation ability by using the same LLM to evaluate whether the generated text meets the given concepts. Appendix B shows more experiment details (e.g., the prompts we use). As the experiment results presented in Fig. 2, we observe that the evaluation accuracy exceeds the generation accuracy in all scales of models. Especially when the parameter size of the model is small (see FLAN-T5-Large/XL), the evaluation accuracy significantly outperforms the text generation by 15%. These results indicate that it is simpler to self-evaluate the generated text than to generate high-quality text that meets contextual requirements. However, it is essential to note that the evaluation accuracy is affected by the quality of the target of evaluation, and this experiment is only a preliminary demonstration of the LLM's ability to self-evaluate. We conduct more experiments to verify the text evaluation ability of LLM further, as presented in the following subsections.
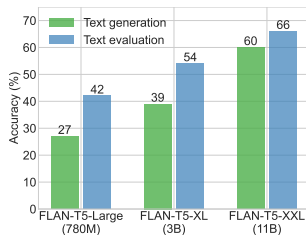


| Task \ Metric | BLEU | ROUGE | BERTScore |
|---|---|---|---|
| CNN/Daily Mail | 0.26 | 0.16 | 0.23 |
| IWSLT 2017 | 0.21 | 0.28 | 0.29 |

Figure 2: Comparison of the text generation and evaluation.

Table 1: The correlation coefficient between evaluation and different metrics. The evaluation result correlates positively with all three metrics. The experiments are conducted with FLAN-T5-XL.

## 4.2 Correlation Between Evaluation and Established Metrics

This section provides an analysis of the correlation between text evaluation and established metrics using two datasets: the CNN/Daily Mail dataset (Nallapati et al., 2016) for text summarization and the IWSLT 2017 dataset (Cettolo et al., 2017) for text translation. To find the relationship between evaluation and established metrics, LLM generates two answers for each question in these tasks and then self-evaluates to determine which answer is better. Additionally, we utilize three metrics, namely BLEU (Papineni et al., 2002), ROUGE (Lin, 2004), and BERTScore (Zhang et al., 2020), to compare the two answers, calculate their correlation coefficient with evaluation, and present the findings in Tab. 1. As demonstrated by the results, the evaluation exhibits a positive correlation with all the considered metrics, indicating a consistent measurement of the quality of the generated text, and it is a reliable way to measure the quality of the generated text.

## 4.3 Potential for Self-improvement

The previous section shows a positive correlation between evaluation and established metrics. However, it remains to be seen whether and how text evaluation can be utilized to improve text generation. To investigate this problem, we design a text generation strategy that utilizes the text evaluation ability of LLM: the LLM re-samples the answer based on the self-evaluation result, which will be explained in the following paragraphs. We conduct experiments on two NLP benchmarks.

|          | Reasoning about Colored Objects | Logical Deduction (7) | Tracking Shuffled Objects (5) | Object Counting | Tracking Shuffled Objects (3) | Geometric Shapes |
|----------|---------------------------------|-----------------------|-------------------------------|-----------------|-------------------------------|------------------|
| w/o SE   | 43.2%                           | 30.3%                 | **13.7%**                     | **22.6%**       | 25.7%                         | 9.6%             |
| w/ SE    | **44.7%**                       | **34.5%**             | 10.6%                         | 22.5%           | **28.5%**                     | **10.5%**        |

|          | Web of Lies | Sports Understanding | Logical Deduction (3) | Logical Deduction (5) | Penguins in a Table | Navigate |
|----------|-------------|----------------------|-----------------------|-----------------------|---------------------|----------|
| w/o SE   | 50.5%       | **55.1%**            | 50.7%                 | 29.7%                 | 30.8%               | **51.1%**|
| w/ SE    | **51.0%**   | **55.1%**            | **57.2%**             | **34.8%**             | **34.0%**           | 49.9%    |

Table 2: Comparison of the answer accuracy between answer generation with/without self-evaluation. Full results on all 27 BigBench tasks are presented in Appendix C.3.

One such benchmark is Bigbench-hard (Srivastava et al., 2022), which includes multiple reasoning tasks consisting of multiple-choice and direct-generation tasks. We compare two answer generation strategies: (1) **w/o SE**: the answer is a directly deterministic output of the LLM, and (2) **w/ SE**: LLM generates an answer and evaluates its correctness. If the evaluation indicates the answer is incorrect, LLM re-generates an answer as the final output. We tested these two strategies on multiple Bigbench-hard tasks, and the results are presented in Tab. 2. The experiment results demonstrate that with evaluation, the answer accuracy outperforms or is comparable to that of direct answer generation on 10 of 12 evaluation tasks. This result justifies using evaluation to help LLM improve answer accuracy.

Furthermore, we also conduct experiments on two text summarization tasks, CNN/Daily Mail and BBC (Gupta et al., 2022). As it is not meaningful to evaluate the correctness of
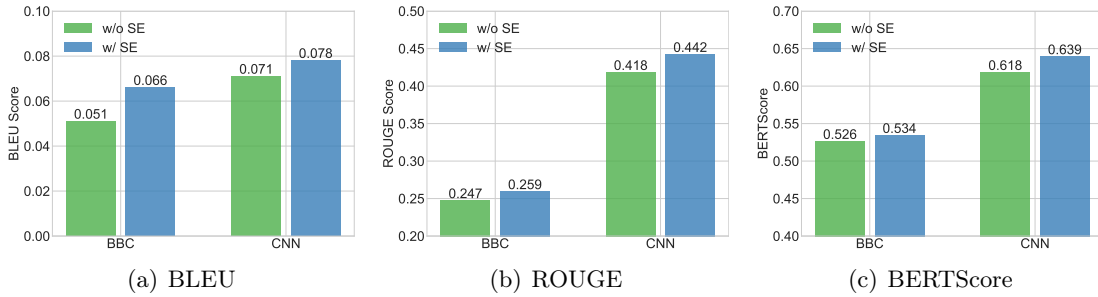
Figure 3: Comparison of generation with/without self-evaluation on summarization tasks.

generated summarization, we use a different approach to utilize self-evaluation ability of LLM in this experiment: (1) **w/ SE**: LLM samples three different answers and evaluates which answer is the best one as the final answer, and (2) **w/o SE**: LLM samples three different answers, and we present the average score of the three answers. As shown in Fig. 3, the generated answers have higher scores under all three metrics when self-evaluation is used. This result suggests that evaluation can potentially improve the quality of the generated text, which serves as a stepping stone for building RLC method.

## 5. Self-improvement by Reinforcement Learning Contemplation

In the previous section, we observe a gap between text evaluation and generation for LLMs and present an opportunity to enhance the generation abilities of LLMs. For example, LLMs can effectively evaluate the sentences they generate. These evaluation results serve as learning signals to guide the LLMs toward improved generation ability.

### 5.1 Method Overview

In this section, we will elaborate on our proposed LMSI method, RLC, that improves LLM ability without external labels based on self-evaluation. We are given a pre-trained language model $\mathcal{M}$ and an unlabeled training dataset $\mathcal{D}^{\mathrm{train}} = \{q_i\}_{i=1}^{|\mathcal{D}|}$. The overview of our method is illustrated in Fig. 4, which iterates through the following steps:

- Gather question-answer pair $(q, o)$ via $q \sim \mathcal{D}^{\mathrm{train}}$, $o \sim \mathcal{M}$.
- Evaluate the question-answer and obtains the reward $r$ via Eq. 3.
- Self-improvement with reinforcement learning.

The overall workflow of RLC method is summarized in Algorithm 1. Next, we elaborate on each step of RLC method.

### 5.2 Details of Reinforcement Learning Contemplation

**Answer generation to unlabeled questions.** RLC randomly samples a batch of questions from $\mathcal{D}^{\mathrm{train}}$ and asks the LLM to generate answers to these questions. For different types of questions, we use different prompts to generate answers, as described in Appendix B.4. For reasoning problems, we use the Chain-of-Thought (CoT) prompt (Wei et al., 2022), which
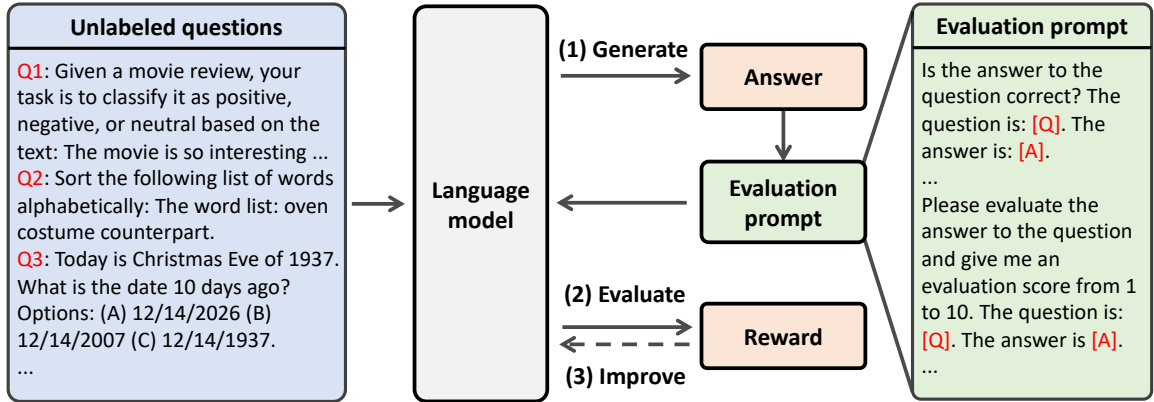
Figure 4: Overall training procedure of RLC, which iterates through three steps: (1) Answer generation to the unlabeled questions. (2) Self-evaluation by asking LM using *evaluation prompt*, with the evaluation results as the reward. (3) Update the language model to maximize the reward using reinforcement learning algorithms. The solid lines represent the data flow, while the dashed line represents the update of LLM parameters.

has been shown to effectively improve LLM performance in previous studies (Wang et al., 2023; Huang et al., 2023).

**Self-evaluation as the reward.** After gathering the question-answer pair $(q, o)$, RLC self-evaluate the generated answer to produce the reward for RL training:

$$R(q, o) = \phi(\mathcal{M}(p_{\text{EP}}, q, o)), \tag{3}$$

where $\phi$ is a processing function that processes the LLM output to the numerical reward and $p_{\text{EP}}$ is the prompt used for self-evaluation. Compared to previous works that refine/revise the generated text based on the evaluation results, RLC uses an *invasive learning style*, i.e., updating the language model to improve its generation ability. The advantage of such invasive training is the quick response speed of the trained model: it does not need to repeatedly generate and revise the answer.

RLC uses two types of prompts for evaluating the self-generated text: (1) Correctness Evaluation Prompt (**CEP**): *"Is the answer to the question correct? The question is: [Q]. The answer is: [A]"*, and (2) Quality Evaluation Prompt (**QEP**): *"Please evaluate the answer to the question and give me an evaluation score from 1 to 10. The question is: [Q]. The answer is [A]"*. Depending on the type of question, either CEP or QEP is utilized to evaluate the generated text.

RLC applies CEP to assess the factual accuracy of the generated text, which is essential for tasks such as question-answering and reasoning. The CEP prompts LLMs to verify the answer's correctness and identify any factual errors. Given a question $q$ and the answer $o$ generated by the LLM, the reward $R(q, o)$ is a 0-1 value based on the evaluation result. On the other hand, RLC employs QEP to evaluate the overall effectiveness and coherence of the generated text in terms of its writing, structure, and style. Quality is often evaluated on a scale, with higher scores indicating the text presents better effectiveness and fits some given requirements more. The QEP prompts LLMs to rate the text's quality on a scale of 1 to 10

---

**Algorithm 1** Reinforcement Learning Contemplation (RLC)

---

**Input**: a pre-trained LLM $\mathcal{M}$, a pre-trained LLM for self-evaluation $\mathcal{M}^* = \mathcal{M}$, an unlabelled dataset $\mathcal{D}^{train} = \{q_i\}_{i=1}^{|\mathcal{D}|}$.
**Output**: the trained LLM $\mathcal{M}$.

1: **while** training not complete **do**
2:     Sample questions $\{q_i\}$ from dataset $D^{train}$.
3:     Use $\mathcal{M}$ to sample answers $\{o_i\}$ to the questions, using stochastic output of LLM.
4:     Calculate reward $R(q, o)$ according to self-evaluation results (Eq. 3).
5:     Update $\mathcal{M}$ with RL algorithm.
6: **end while**
7: **return** the trained LLM $\mathcal{M}$.

---

based on how well it answers the question. This prompt type helps solve generation tasks like translation and summarization. During training, the reward distribution may change as the LLM is updated. We use the initial pre-trained LLM $\mathcal{M}^*$ for self-evaluation while keeping it fixed to ensure stable training.

**Self-improvement through reinforcement learning.** With the evaluation reward, the LLM $\mathcal{M}$ can be updated using any RL algorithm to maximize the reward. In our case, we employ the Proximal Policy Optimization (PPO) algorithm (Schulman et al., 2017), which has demonstrated promising results in applications (Huang et al., 2022). To ensure better exploration, we apply entropy regularization, which prevents the sampling phase from converging too early. Additionally, we use the Kullback-Leibler (KL) divergence to prevent $\mathcal{M}$ from deviating too far from the initial pre-trained LLM.

## 6. Experiments

We conduct a series of experiments to demonstrate the effectiveness of our proposed RLC method. Our experiments encompass the following topics: (1) comparison of the proposed method with baseline methods on various NLP tasks (Section 6.2); (2) the applicability of the proposed method to different model sizes (Section 6.3); and (3) general improvement on unseen tasks (Section 6.4). We first introduce our experimental setup.

### 6.1 Experiment Setup

**Benchmarks for evaluation.** We use a wide range of established benchmarks to evaluate the performance of the proposed method. **(1) Reasoning**: Reasoning is one of the most essential ability of LLMs (Huang & Chang, 2023). BigBench-Hard dataset (Srivastava et al., 2022) consists of 27 challenging reasoning tasks, which are pretty diverse, including reasoning the final results of a sequence of behaviors, understanding dates, and completing tasks that require simple arithmetic calculations. In our experiments, we follow HIR (Zhang et al., 2023) to use 12 challenging tasks from the BigBench-Hard datasets[1], which covers judgments, multiple choices, and text generation tasks. See Tab. 3 for examples of the reasoning tasks used in our experiments. **(2) Article generation**: CommonGen (Lin et al., 2020)

---

1. Detailed descriptions about the tasks are in `https://github.com/google/BIG-bench/blob/main`.

benchmark requires constructing a coherent sentence that describes everyday scenarios using a predefined set of concepts. CommonGen presents significant challenges, as it demands two primary skills: (1) utilizing background knowledge for relational reasoning and (2) effectively managing combination generalization of concealed concept combinations. **(3) Emotion recognition**: Emotion (Saravia et al., 2018) is a dataset consisting of messages with six basic emotions: anger, fear, joy, love, sadness, and surprise. The objective is to recognize the emotion based given the message. **(4) Summarization**: CNN/Daily Mail (Nallapati et al., 2016) is a popular dataset used for text summarization tasks. This dataset covers a wider range of topics than the BBC dataset, including politics, sports, and entertainment. In contrast, the BBC dataset focuses more on news and current events.

| Task type | Example input | Example output |
| --- | --- | --- |
| Judgement | "Jamey lies. Vernell says Jamey tells the truth. Millicent says Vernell lies. Rashida says Millicent ... Does Vina tell the truth?" | "Vernell is lying because Vernell is telling ... The answer: Yes." |
| Text generation | "I have a chicken, a snail, three dogs, a snake, and a goat. How many animals do I have?" | "...So I have a chicken + snail + four dogs + a snake + a goat = 7 animals. So the answer is 7." |
| MultiChoice | "Here is a table where the first line is a header and each subsequent line is a penguin: name, age, ... How many penguins are there in the tables? Options: (A) 1 (B) 2 (C) 3 (D) 4 (E) 5." | "There are 4 penguins in the first table and 4 giraffes in the second table. The answer: (D). " |

Table 3: Examples of inputs and outputs for the tasks in our experiments.

**Baselines for comparison.** We compare RLC with representative baselines that improve LLM without supervised data. These methods include: (1) **RLAIF** trains a reward model based on AI preference data and employs a RL algorithm to train LLM. We follow implementation of RLAIF (Lee et al., 2024), which utilizes a smaller LLM as reward model. In the experiment, we use GPT2 (Radford et al., 2019); (2) Reinforcement Learning Fine-Tuning (**RLFT**) fine-tunes the LLM using RL and employs true dataset label as the reward. We use answer correctness as the reward for multiple-choice, judgement, fill-in-the-blank tasks, and BERTScore for summarization task; (3) **Self-train** (Huang et al., 2023) utilizes the SC method to generate "high-confidence" answers for unlabelled questions as the pseudo-labels. It then fine-tunes the LLM on the generated data with supervised learning; (4) Direct Generation (**DG**) directly generates the answer with the deterministic output of the LLM; (5) Self-consistency (**SC**) (Wang et al., 2023) samples a diverse set of reasoning paths instead of solely relying on the greedy LLM output. It then selects the most consistent answer by marginalizing out the sampled reasoning paths. In our experiments, we use the number of sample paths as 40; (6) **Self-refine** (Madaan et al., 2023) employs a LLM to generate an initial output. This output is evaluated from multiple perspectives and revised by the

| | Reasoning about Colored Objects | Logical Deduction (7) | Tracking Shuffled Objects (5) | Object Counting | Tracking Shuffled Objects (3) | Geometric Shapes |
|---|---|---|---|---|---|---|
| RLFT | 32.1% | 45.7% | 12.4% | 42.6% | 33.6% | 18.9% |
| DG | 32.0% | 35.2% | 12.4% | 31.9% | 31.2% | 5.2% |
| SC | **39.6**% | 27.6% | 12.4% | 24.0% | 33.6% | 15.6% |
| Self-train | 19.5% | 13.1% | **15.5%** | 11.7% | 33.1% | 12.4% |
| Self-refine | 25.2% | 13.2% | 8.0% | 18.0% | 25.2% | 10.0% |
| Best-of-N | 26.8% | 12.8% | 12.1% | 14.0% | 30.0% | 8.4% |
| RLAIF | 30.4% | 36.9% | 11.4% | 32.5% | 32.8% | 14.0% |
| RLC | 35.0% | **39.2%** | 12.2% | **35.4%** | **33.6%** | **17.8%** |

| | Web of Lies | Sports Understanding | Logical Deduction (3) | Logical Deduction (5) | Penguins in a Table | Navigate |
|---|---|---|---|---|---|---|
| RLFT | 72.2% | 68.8% | 58.6% | 41.9% | 44.2% | 55.6% |
| DG | 43.6% | 53.2% | 39.6% | 28.4% | 15.7% | 46.4% |
| SC | 48.4% | 53.6% | 42.8% | 30.8% | **35.2**% | **62.8**% |
| Self-train | 51.1% | 51.1% | 34.0% | 18.4% | 19.7% | 48.7% |
| Self-refine | 47.2% | 50.0% | 28.4% | 17.2% | 17.8% | 46.0% |
| Best-of-N | 50.0% | **59.2**% | 42.0% | 22.0% | 17.8% | 45.2% |
| RLAIF | 52.1% | 56.1% | 22.0% | 33.7% | 19.8% | 48.8% |
| RLC | **52.9%** | 53.5% | **44.0%** | **34.6%** | 29.8% | 57.1% |

Table 4: The answer accuracy of RLC and baseline methods on the Bigbench-hard tasks. Each value represents the average answer accuracy of the last three training iterations. The highest performing value among methods without external labels is highlighted in **bold**.

same LLM. (7) **Best-of-N** (OpenAI, 2022) generates N answers and uses the same model to evaluate and choose the best answer. Some may be concerned about the fairness of the comparison, as Best-of-N and SC generate multiple answers and select one best answer from the answer candidates, while the other method only generates one answer as the final answer. It is important to highlight that our experimental focus lies in assessing the performance gains of various methods in enhancing language models. This, the evaluation primarily measures the performance disparity between DG.

**Implementation details.** We utilize PPO to train LLM for 6,000 gradient steps on each unlabeled dataset, with a training batch size of 12. The PPO implementation is from the open-sourced trlx repository (CarperAI, 2020). We implement RLC using QEP for summarization tasks and CEP for other types of tasks. Unless otherwise specified, we use FLAN-T5-Large (780M parameters) as the LLM for evaluation. All reported results are averaged over three random trials, except for RLAIF with one seed, and the experiments can be conducted using two GTX 3090 graphics cards, each with 24GB memory. We provide specific hyper-parameters and more detailed implementation descriptions in Tab. 8 in Appendix B.3. All RL-based methods use the same setting for the RL algorithm. In our experiments, we consistently employ the same prompts for all baseline methods and the RLC. To prompt the model to output high-quality answers, we utilize the CoT prompts, specifically, "Let's think step by step." A comprehensive list of the prompts used in our experiments can be found in Table 9 of the Appendix.

## 6.2 Comparison Analysis on Various Tasks

**Results on reasoning tasks**. Reasoning is one of the most fundamental capabilities of LLMs. We conduct experiments on extensive reasoning tasks in the BigBench-Hard dataset. Tab. 4 presents the answer accuracy of the LLMs trained with RLC and other baseline methods. We see RLC outperforms all the baselines that do not use external supervision. In particular, the RLC improve answer accuracy over DG from 31.23% to 37.09%. RLFT serves as the performance upper bound as it fine-tunes LLM with ground-truth supervision. On some tasks (e.g., Reasoning about Colored Objects), RLC even catches up with the score of RLFT. This result potentially indicates that RLC method effectively identifies the proper direction for self-improvement, contributing to the close performance of the RLFT method. SC presents good performance on these tasks and obtains the highest score on 3 of 12 tasks. Note that SC is also compatible with RLC, which can be used to improve the algorithm performance further. Besides, RLC can achieve comparable or better performance than Best-of-N and self-refine, while RLC takes less computation time during the generation process. We also compare RLC with RLAIF, and results show that RLC achieve a better performance. This can be attributed that the reward model with a small parameter size can not provide effective preference results (Lee et al., 2024). We depict the portion of RL training curves in Fig. 10, where the answer accuracy improves with the RL training.



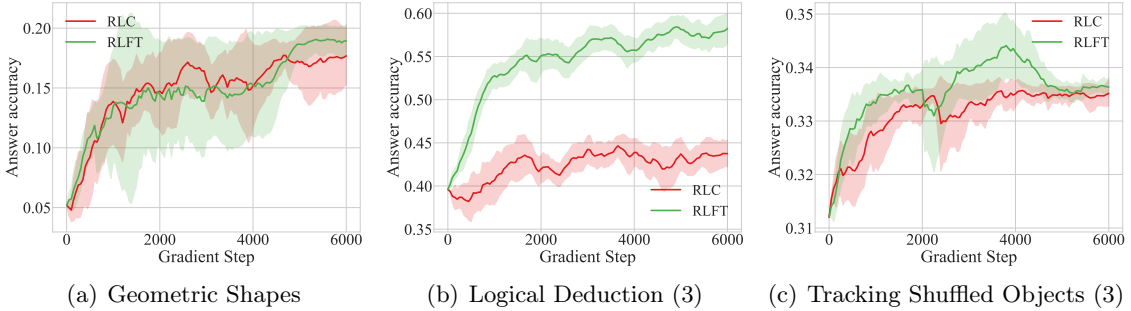| (a) Geometric Shapes | (b) Logical Deduction (3) | (c) Tracking Shuffled Objects (3) |

Figure 5: Training curves of reinforcement learning contemplation on BigBench-hard tasks. The shaded area represents the standard deviation over three seeds. Full results are presented in Fig. 10 in Appendix.

**Results on summarization task.** We evaluate the performance of RLC on a challenging text generation benchmark: CNN/Daily Mail. In this experiment, the Self-train method is implemented as follows: We utilize the initial LLM to generate the answers as the supervision label to the questions in the dataset, after which the LLM is fine-tuned on this newly generated data. As illustrated in Fig. 6, RLC enhances the BERTScore from 0.886 to 0.899 in the summarization task, while self-train hardly improves the model performance on this dataset. Unlike the BigBench tasks, where the generated text is relatively short, the LLM
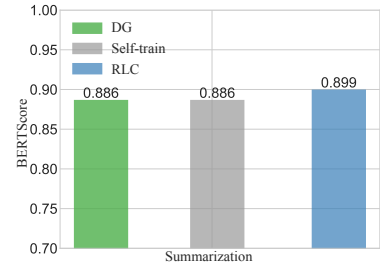


Figure 6: Performance on summarization task.

is required to produce longer text for these two tasks. The experimental results demonstrate that RLC effectively improves the text generation capabilities of the LLM.

**Results on article generation and emotion recognition tasks.** Tab. 5 presents the experimental results on CommonGen and Emotion datasets with FLAN-T5 (780M and 3B). These tasks are particularly important in the realm of natural language processing as they evaluate the foundational capabilities of LLMs in comprehending and generating coherent text. The results suggest that RLC effectively improves LLMs with 780M and 3B parameters. for the Emotion classification task, RLC improves the performance of the 780M parameter model from 57.9% to 61.9%, which approaches the accuracy achieved by the 3B parameter model. We present the training curves in Fig. 7 in Appendix.

|     | CommonGen (780M) | CommonGen (3B) | Emotion (780M) | Emotion (3B) |
|-----|------------------|----------------|----------------|--------------|
| DG  | 45.5%            | 31.5%          | 57.9%          | 63.7%        |
| RLC | 47.5% (↑2)       | 33.5% (↑2)     | 61.9% (↑4)     | 64.7% (↑1)   |

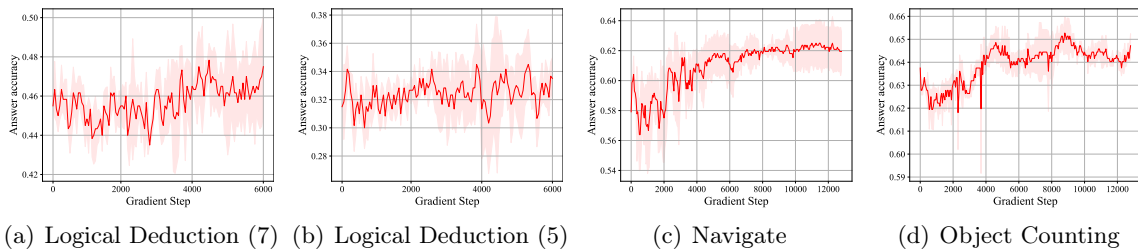Table 5: Answer accuracy on CommonGen and Emotion benchmark.



(a) Logical Deduction (7)  (b) Logical Deduction (5)  (c) Navigate  (d) Object Counting

Figure 7: Training curves of RLC method on CommonGen and Emotion.

## 6.3 Evaluation on Model Size Variations

We conduct experiments to assess the performance of RLC across various language model sizes. We select three distinct models: FLAN-T5-Small, FLAN-T5-Base, and FLAN-T5-Large, containing 80M, 250M, and 780M parameters, respectively. We train these models using RLC on three challenging BigBench tasks, with the final scores presented in Fig. 8. In general, RLC effectively enhances performance across different scales of language models. Notably, when the parameter size is small (80M) and the base score is low, the language model exhibits a significant improvement.

## 6.4 General Improvement by Scale Training

Previous experiments present RLC can effectively improve LLM answer quality on a specific dataset. We further explore an important question: can RLC improve LLM's general ability? To investigate this, we scale up the training by training LLM on 8 datasets (Geometric Shape, Logical Deduction (3), Logical Deduction (5), Tracking Shuffled Objects 3, Penguins in a Table, Navigate, Web of Lies, Temporal Sequences) simultaneously and evaluating the trained LLM on 5 other unseen tasks (Logical Deduction (7), Reasoning about Colored Objects, Formal Fallacies, Ruin Names and Emotion).

(a) Tracking Shuffled Objects (3)    (b) Penguins in a Table    (c) Object Counting
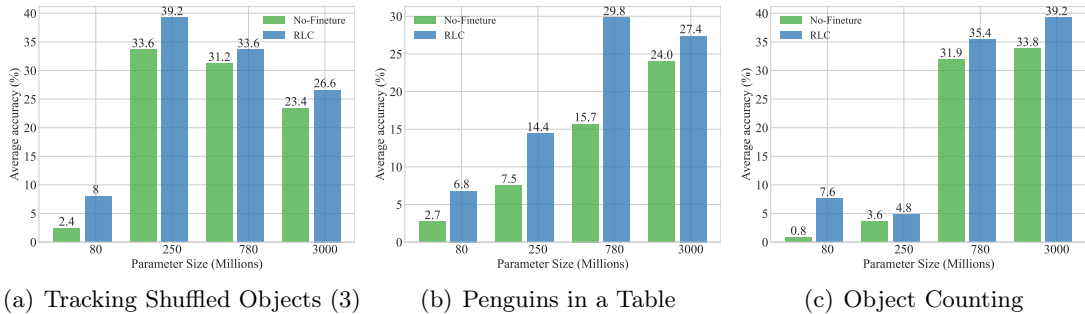
Figure 8: Performance of RLC on different sizes of language models. RLC can improve the answer accuracy over the initial pre-trained model without supervision.

**Performance on unseen tasks.** As the evaluation results depicted in Tab. 6, the LLM generally presents an improvement on the unseen tasks. There are several key findings: First, LLM shows improvement across tasks of various subjects. Notably, there is a marked improvement in performance on Logical Deduction (7), which involves deducing the sequence of objects, and on Ruin Names, which requires the humorous alteration of movie or artist names. Besides, we observe that improvements achieved by LLMs can be transferred between tasks that share similarities. For instance, training with RLC on Logical Deduction (5) leads to improvement on the more complex Logical Deduction (7), despite the latter involving a greater number of objects to reason about. Last, we note cross-category improvement in LLM capabilities. This is supported by a 4.2-point increase in performance on the Emotion task, which is not part of the BigBench-Hard dataset. These findings underscore the potential of RLC for application across a wider array of datasets, thereby improving the general ability of LLMs.

| | Logical Deduction (7) | Reasoning about Colored Objects | Formal Fallacies | Ruin Names | Emotion | Average |
|---|---|---|---|---|---|---|
| DG | 34.8 | 31.6 | 52.4 | 20.0 | 57.9 | 39.3 |
| RLC | 36.4 (↑ 1.6) | 30.8 (↓ 0.8) | 54.0 (↑ 1.6) | 21.2 (↑ 1.2) | 62.1 (↑ 4.2) | 40.9 (↑ 1.6) |

Table 6: Performance on 5 unseen tasks. The LLM (FLAN-T5, 780M) is trained with RLC on 8 BigBench-Hard tasks and is evaluated on 5 unseen tasks.

**Discussions about continual improvement.** Another interesting question we would like to discuss is *Can RLC enable language model's continual self-improvement?* We have observed that RLC improves LLM's performance on some unseen tasks by scale training. Hopefully, we expect that with an increased training scale, LLM can improve its ability in various domains, including its evaluation ability. In turn, the evaluation ability contributes to boosting LLM's generation ability (as previous results show), thus facilitating continual improvement. Fig. 9 shows the change of



Figure 9: Generation and evaluation accuracy change during the training.

14

evaluation/generation accuracy on unseen tasks during the scale training. The experiment results suggest that RLC not only effectively improves LLM's text generation, but also its evaluation accuracy. Besides, the evaluation accuracy presents an overall upward trend during the training, which indicates a potential for further improving language model's generation, achieving a continual improvement.

## 7. Conclusion

This paper introduces a novel approach to enable LLM to self-improve: reinforcement learning contemplation (RLC). We demonstrate that it is simpler for LLM to evaluate the text than to generate it, and thus the evaluation can serve as an powerful tool to boost generation. By utilizing evaluation ability, LLMs can obtain valuable insights into their performance and identify areas for improvement. Our experiments demonstrate that RLC effectively improves LLM performance across various text generation tasks. Moreover, we illustrate the potential for the RLC method to scale with varying model sizes, training datasets, achieving generation ability improvement on unseen tasks. However, some things could be improved in our approach. First, like previous unsupervised methods for training language models, RLC requires an unlabeled dataset to generate answers and facilitate self-improvement. Such dataset can be viewed as a specific improvement direction (e.g., reasoning ability). It would be worthwhile to explore training LLMs without relying on the datasets. For instance, the LLM identifies a improvement direction and self-improves towards such direction without rely on the unlabeled dataset. Besides, we primarily evaluate RLC on LLMs eaual to or less than 3B parameters due to computational limitation. Future experiments should include larger language models to better demonstrate the method's application range. We hope this work could provide new insights into LMSI. Future research could address the above questions to develop more efficient and effective training methods for language models.

## Appendix A. Discussion

### A.1 The Usage of Unlabelled Questions Dataset

In Section 5, we introduce RLC, assuming the availability of a training dataset $\mathcal{D}^{\text{train}}$ consisting of unlabeled questions. This raises two questions: (1) what data type does the dataset contain? and (2) does RLC work without an unlabeled dataset?

Concerning the first question, the training dataset in our experiments involves questions in $\mathcal{D}^{\text{train}}$ that adhere to an *instruction + problem* format. For example, a question might be "Summarize the following article (instruction). [TEXT] (problem)". In this way, RLC can train an LLM using any open-source datasets in the NLP community by appending instructional text to the problems within these datasets. Tab. 7 presents the instructions we used in our experiments, where '[TEXT]' denotes the questions/text in the original dataset. This way of using a dataset is similar to instruction learning (Lou et al., 2023), which utilizes task instructions to guide the learning system in rapidly adapting to various new tasks. We posit that as the data in $\mathcal{D}^{\text{train}}$ becomes more diverse and abundant, RLC can be employed to train an LLM with generalization ability. Our experiments in Section 6.3 offer preliminary validation of this outcome.

| Task | Instruction |
|------|-------------|
| BigBench-hard | [TEXT] Let's think step by step. |
| Text summarization | Please give a summary of the following text. Text: [TEXT] Answer: |

Table 7: Instructions for different tasks in our experiments.

As for the second question, although our experiments utilize an unlabeled dataset to train the LLM, we are interested in exploring whether RLC can be applied without a dataset. To achieve this, we need an objective, which serves as an evaluation criterion for assessing the quality of the generated text, to train the language model. For example, we might expect an LLM to consistently generate positive or polite text. In such a case, we can prompt the LLM to generate text randomly from the start token and evaluate its positivity or politeness using either the CEP or the QEP. This approach allows the LLM to optimize toward the desired attribute. However, it is essential to note that some training objectives or text attributes might be challenging to evaluate, making them unsuitable for self-improvement using the reinforcement learning contemplation method.

### A.2 Sampling Strategy

In our experiments, RLC employs a random sampling approach to select questions from the dataset. Although this method is effective, alternative sampling techniques could enhance the process. For instance, prioritizing questions with more uncertainty could be beneficial. The field of active learning offers many practical sampling methods (Xu, Sun, & Zhang, 2013) that could potentially augment learning efficiency. We posit that these active sampling strategies

align well with the concept of *self-improvement*. However, in this study, our primary focus lies in the design of the LMSI method, for which we maintain the use of random sampling.

## Appendix B. Experiment Details

In our experiments, we utilize the open-sourced RL repository, trlx, to implement the reinforcement learning contemplation. All experiments are run three times with different random seeds. This section will present more experiment details omitted in the main body due to the space limitation, including the evaluation task we use, the baselines, the prompts in different experiments, the hyper-parameters for reproducibility, etc.

### B.1 Tasks for Evaluation

In our experiments, we use five challenging benchmarks in the NLP domain to conduct various experiments to support our method. This section gives a detailed introduction to these benchmarks, which are omitted in the main body.

**CommonGen** (Lin et al., 2020) is a task that focuses on constrained text generation and includes a benchmark dataset. Its primary objective is to assess a machine's ability to generate common sense reasoning. The task requires constructing a coherent sentence that describes everyday scenarios using a predefined set of concepts. CommonGen presents significant challenges, as it demands two primary skills: (1) utilizing background knowledge for relational reasoning and (2) effectively managing combination generalization of concealed concept combinations. We employ CommonGen to evaluate the text generation and self-evaluation ability of LLM, wherein the LLM is tasked with generating a sentence based on four concepts simultaneously.

**BigBench-hard** (Srivastava et al., 2022) consists of 27 challenging tasks designed to evaluate the reasoning abilities of language models. These tasks present increased difficulty due to their complexity. In our experiments, we employ 12 challenging tasks that focus on evaluating different aspects of reasoning ability:

- Reasoning about Colored Objects (multiple choices): Answer simple questions about the colors of objects on a surface.
- Logical Deduction (multiple choices): Deduce the order of a sequence of objects.
- Tracking Shuffled Objects (multiple choices): Determine the final positions of a set of objects given their initial positions and a description of a sequence of swaps.
- Object Counting (text generation): Questions that involve enumerating objects of different types and asking the model to count them.
- Geometric Shapes (text generation): Name geometric shapes from their SVG paths.
- Web of Lies (judgement): Evaluate a random boolean function expressed as a word problem.
- Sports Understanding (judgement): Determine whether an artificially constructed sentence relating to sports is plausible or implausible.
- Penguins in a Table (text generation): Answer questions about a table of penguins and their attributes.

- Navigate (judgement): Given a series of navigation instructions, determine whether one would end up back at the starting point.

Note that specific tasks (e.g., Logical Deduction) encompass the same topic but are presented at varying difficulty levels, and we introduce these tasks with varying levels together.

The **CNN/Daily Mail** (Nallapati et al., 2016) dataset serves as a widely recognized benchmark for text summarization. Comprising over 300,000 news articles from CNN and the Daily Mail and corresponding human-written summaries, this dataset has been instrumental in training and evaluating various text summarization models, including extractive and abstractive methods. Each article features several highlights summarizing its main points, making the dataset ideal for training and testing automatic text summarization models that aim to generate concise versions of the original text while retaining crucial information.

The **BBC** (Gupta et al., 2022) dataset, created by the BBC News website, is another widely used resource for text summarization. The dataset has been pre-processed and annotated with human-generated summaries, encompassing approximately 2,225 news articles on diverse topics such as politics, entertainment, technology, and sports. These concise summaries, typically 3 to 4 long sentences, make the dataset invaluable for training and evaluating text summarization models.

### B.2    Baselines

In our experiments, we compare RLC with representative baseline methods. This section presents these baselines and their implementation details.

(1) SC: Self-Consistency. SC enhances the LLM's answers without updating its parameters. Its motivation is similar to ensemble methods, which sample multiple answers simultaneously and vote among them to select the most consistent one. SC has demonstrated significant capabilities in reasoning problems. In our experiments, the sampling path (i.e., the number of answers that LLM samples) is set to 40.

(2) Self-train: Language Model Self-train. LMSI combines the SC method with supervised learning fine-tuning. It first uses SC to generate answers for all questions in the dataset and then fine-tunes the LLM using supervised learning. In contrast to LMSI, which employs supervision signals to enhance language models, RLC utilizes reward signals. Reward signals provide a more abstract and general performance measure, enabling reinforcement learning algorithms to learn from fewer examples and generalize better to new situations. Furthermore, reward signals allow encoding of more complex and diverse settings, potentially leading to more creative and adaptive answers.

(3) RLFT: Reinforcement Learning Fine-Tuning. RLFT trains the LLMs using Oracle metrics as rewards, such as answer correctness or BERTScore. We include this baseline to verify the effectiveness and performance upper bound of RLC. In our experiments, we implement RLFT and RLC using the same code framework. The primary difference is that RLC employs self-evaluation results as rewards.

(4) Self-Refine: Self-Refine is a method designed to enhance the initial outputs generated by Large Language Models (LLMs) through an iterative feedback and refinement process. The core concept involves generating an initial output using an LLM, after which the same LLM evaluates and refines its output iteratively. Similarly, Self-Refine operates without

the need for supervised training data. However, due to Flan-T5-Large being too small to generate effective feedback, we have pruned the feedback section and directly inputted the results generated by refinement into LLM to obtain the following result.

(5) Best-of-N: Best-of-N sample n times and take the one that scores the highest according to the proxy objective. In our experiments, the number of answers LLM generates (i.e., N) is set to 40. And then, we use the same LLM to evaluate the 40 generated answers as the proxy objective. Then, we choose the answer with the highest score as the final answer.

(6) RLAIF: The RLAIF framework is similar to RLHF, but the difference is that an off-the-shelf LLM, instead of humans, labels the preferences. In our experiments, we use GPT2-Large as the off-the-shelf LLM to train the reward model, which is the same as (Lee et al., 2024).

## B.3   Hyper-parameters

Tab. 8 presents the hyper-parameters used in our experiments.

Table 8: Hyper-parameters in our experiments.

| Hyper-parameters | Value |
|---|---|
| PPO epoch | 4 |
| PPO clip ratios | 0.2 |
| PPO $\lambda$ | 0.95 |
| batch size | 12 |
| value function coefficient | 1.0 |
| learning rate | 1e-4 |
| discount factor $\gamma$ | 0.99 |
| temperature for LLM exploration | 1 |
| top_k of LLM | 50 |
| top_p of LLM | 0.95 |
| Sampling path of SC/LMSI | 3 |
| LMSI learning rate | 5e-4 |
| LMSI batch size | 2 |
| SC sampling path | 40 |

## B.4   Prompts in Our Experiments

In various experiments, different prompts serve specific experimental objectives. A comprehensive overview of these prompts can be found in Tab. 9. We employ unique instructional text to generate answers for questions across diverse datasets, as illustrated in Tab. 7.

## B.5   Metrics used in Our Experiments

In our experiments, we present various experimental results under different metrics. For BigBench, we utilize **accuracy** to judge the correctness of the generated answer compared to the reference answer. In translation and summarization tasks, we consider the following metrics: **BLEU** is a reference-based metric that evaluates the similarity between a machine-generated output and one or more reference outputs. **BERTScore** is a reference-less metric

| Experiment | Prompt |
|---|---|
| Comparison of the text generation and self-evaluation (Section 4.1) | Consider a task that needs to generate a coherent sentence describing an everyday scenario using all following concepts. You will be given a few concepts and a sentence, please tell me whether the task is done. If you think the task is done, reply yes. If you think the task is not done, reply no. Concepts: [CONCEPT]. Sentence: [SENTENCE]. |
| Correlation with self-evaluation and established metrics (Translation) (Section 4.2) | Suppose you are a reviewer of the text translation. You will be given two translations of a text, please tell me which one is better according to the conciseness and integrality of the translation. If you think Translation (1) is better, reply (1). If you think Translation (2) is better, reply (2). Text:[TASK] Translation (1): [Translation_1] Translation (2): [Translation_2]. Which one is better? |
| Correlation with self-evaluation and established metrics (Summarization) (Section 4.2) | Suppose you are a reviewer of the text summary. You will be given two summaries of a text. Please tell me which one is better according to the conciseness, integrality of the summary. If you think Summary (1) is better, reply (1). If you think Summary (2) is better, reply (2). Text: [TASK] Summary (1): [Summary_1] Summary (2): [Summary_2]. Which one is better. |
| CEP for BigBench-Hard (Section 6.2) | Is the answer to the question correct? The question is: [Q]. The answer is: [A] |
| QEP for Summarization (Section 6.2) | Please help me evaluate the summary results of the following text. Only give a score from 1 to 10, without explanation. Text: [Q] Summary: [A] |
| Generation prompt for CommonGen (Section 6.2) | Please help me generate a coherent sentence describing an everyday scenario using all following concepts.\nWords: words |
| CEP for CommonGen (Section 6.2) | Answer the following no/yes question.\nIf the sentence contains all the given words?\nWords: words\nSentence: answer |
| Generation prompt for Emotion (Section 6.2) | Please help me judge what is the emotion of the given text.\nChoose from emotions\nText: text |
| CEP for Emotion (Section 6.2) | Answer the following no/yes question.\nIf 'emotion' describes the following text's emotion?\nText: {text} |

Table 9: Prompts used in different experiments.

that assesses the similarity between the embeddings of a machine-generated output and a reference answer. **ROUGE** is another reference-based metric that measures the overlap between the generated output and one or more reference outputs. BLEU emphasizes precision, while ROUGE focuses on recall.

## Appendix C. Additional Experimental Results

### C.1 Training Curves of RLC Method

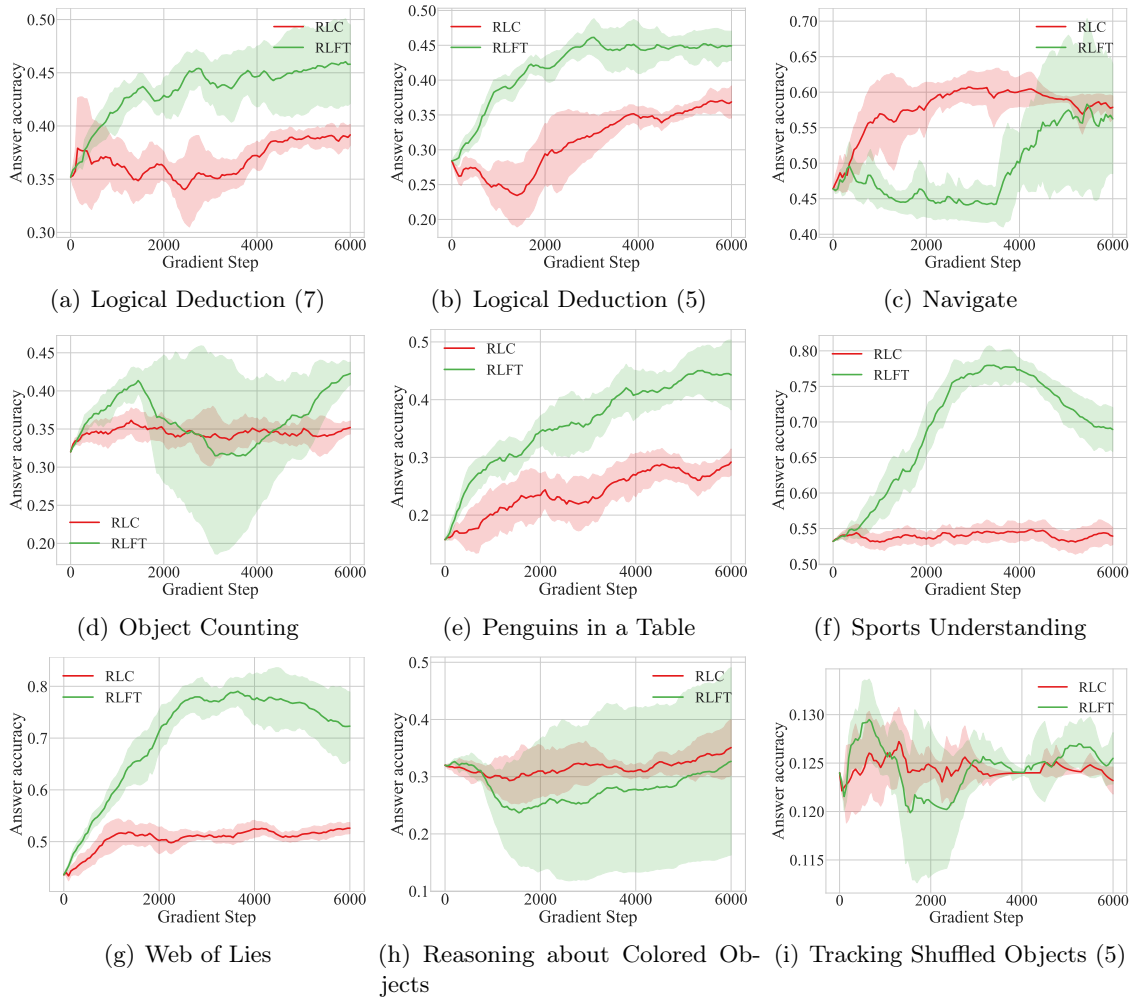Fig. 10 presents additional training curves of RLC method omitted in the main body.



Figure 10: Training curves of reinforcement learning contemplation on BigBench-hard tasks. The shaded area represents the standard deviation over three seeds.

## C.2  Correlation Between Answer Accuracy and Training Return

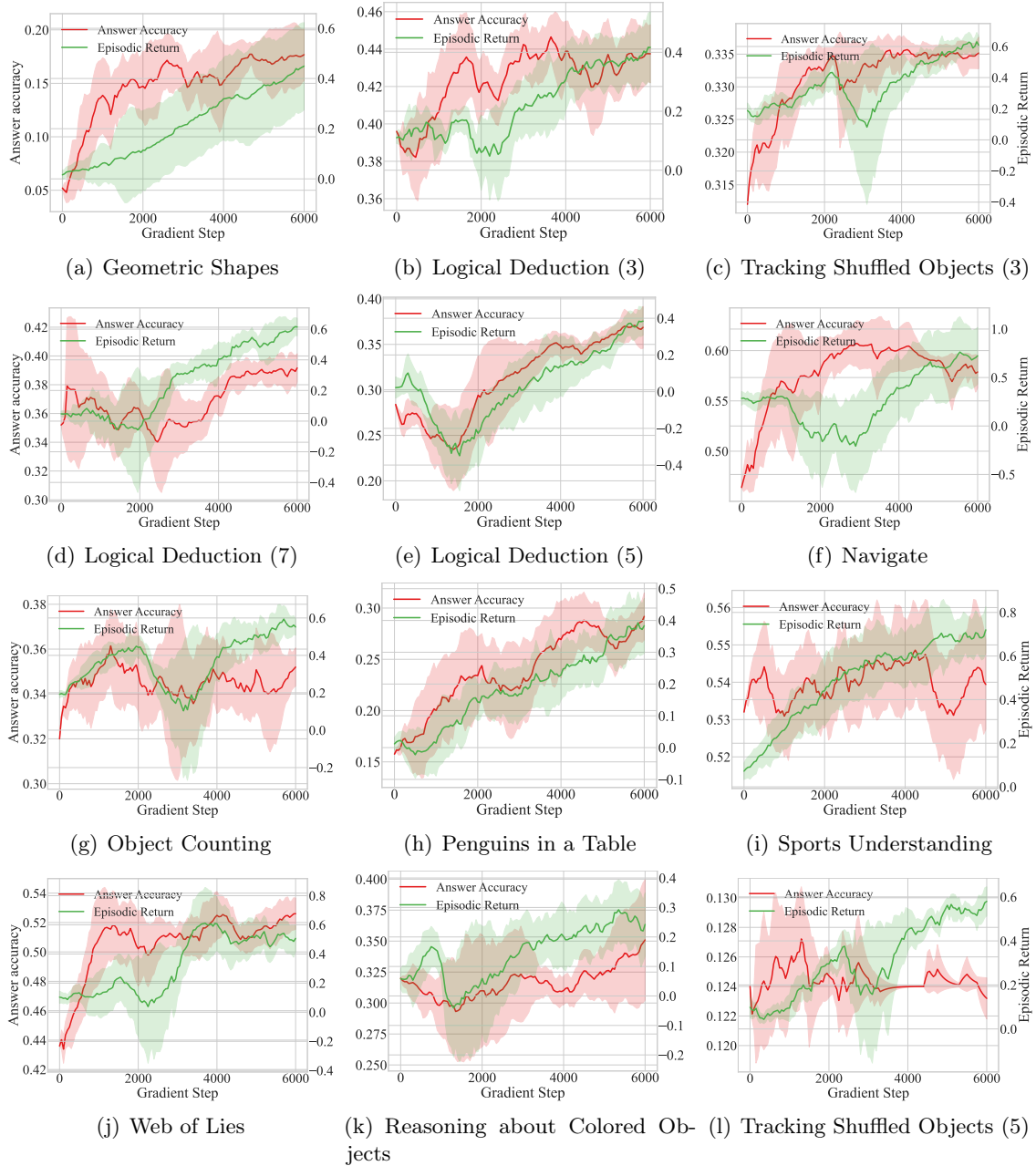Fig. 11 presents all the changes of answer accuracy and episodic return during RLC training process.



Figure 11: Training curves of the answer accuracy and episodic return on BigBench-hard tasks.
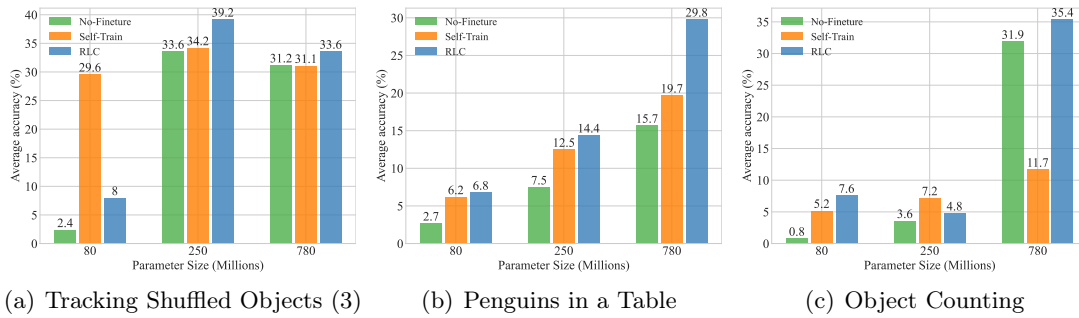
## C.3 More Experiment Results about Self-evaluation Ability Verification

In Section 4.3, we examine the self-evaluation ability of LLMs and their potential for self-improvement. Due to space limitations in the main body, we present only a subset of the tasks (12 tasks). This section provides the complete results for all 27 tasks on the BigBench-hard benchmark, as shown in Tab. 10.

| | Reasoning about Colored Objects | Logical Deduction (7) | Tracking Shuffled Objects (5) | Object Counting |
|---|---|---|---|---|
| w/o SE | 43.2% | 30.3% | **13.7%** | **22.6%** |
| w/ SE | **44.7%** | **34.4%** | 10.7% | **22.5%** |
| | Web of Lies | Sports Understanding | Logical Deduction (3) | Logical Deduction (5) |
| w/o SE | 50.5% | **55.1%** | 50.7% | 29.7% |
| w/ SE | **51.0%** | **55.1%** | **57.2%** | **34.8%** |
| | Hyperbaton | Formal Fallacies | Date Understanding | Causal Judgement |
| w/o SE | 55.3% | **62.8%** | 29.9% | 56.0% |
| w/ SE | **63.1%** | 60.9% | **32.5%** | **57.4%** |
| | Boolean Expressions | Ruin Names | Tracking Shuffled Objects (7) | Temporal Sequences |
| w/o SE | 50.8% | **27.1%** | **9.9%** | **22.8%** |
| w/ SE | **57.5%** | 26.8% | 9.2% | **22.8%** |
| | Tracking Shuffled Objects (3) | Geometric Shapes | Snarks | Navigate |
| w/o SE | 25.7% | 9.6% | **65.5%** | **51.1%** |
| w/ SE | **28.5%** | **10.5%** | 59.4% | 49.9% |
| | Penguins in a Table | Disambiguation QA | Multistep Arithmetic Two | Word Sorting |
| w/o SE | 30.8% | 55.3% | **3.5%** | **1.1%** |
| w/ SE | **34.0%** | **58.5%** | **3.5%** | **1.1%** |
| | Dyck Languages | Salient Translation Error Detection | Movie Recommendation | Average |
| w/o SE | 27.6% | 26.5% | 35.5% | 34.9% |
| w/ SE | **28.4%** | **34.7%** | **41.2%** | **36.7%** |

Table 10: Accuracy of answer generation with and without self-evaluation on 27 BigBench tasks.

## C.4    Self Train on Different Model Sizes



(a) Tracking Shuffled Objects (3)    (b) Penguins in a Table    (c) Object Counting

Figure 12: Performance of RLC  and Self-train on different sizes of language models.

## C.5    Examples of the LLM Output

Tab. 11-14 presents examples of LLM output before/after training with RLC.

| Question | Before training | After training |
|---|---|---|
| Please give a summary of the following text. Text: Kyle Naughton has been ruled out for the rest of the season after David Meyler's red card challenge on Saturday. The Hull midfielder was shown a straight red for the hefty challenge that left Naughton in hospital at the weekend, though neither Garry Monk nor Steve Bruce condemned Meyler at the time. A scan on Monday revealed Naughton suffered ligament damage and the right-back will face a six-week lay-off. David Meyler (left) lunges in high on Kyle Naughton and was sent off the poor challenge on Saturday. The full back writhes in pain following the incident and faces a six-week lay off with ankle ligament damage. A club statement said: 'An initial scan on Saturday ruled out a fracture, but a follow up scan last night confirmed ligament damage that will keep the January signing from Spurs out for around six weeks. 'With just seven league games left, it's a blow for Naughton and the Swans with the right-back continuing to impress with every appearance.' Angel Rangel will return to the side to face Everton on Saturday. Rangel said: 'It's not good because Kyle has an injury and I don't wish bad on anyone, but from my point of view I always try and be ready so if the manager needs me I'll try my best.' Naughton will miss the rest of the season and had to helped off the pitch by the Swansea team physios. Meyler looks shocked as referee Andre Marriner shows him a straight red, a decision that was justified. Answer: | [BELU score: 0.887] David Meyler lunged in high on Kyle Naughton and was sent off. The Hull midfielder was shown a straight red for the hefty challenge. The Hull midfielder was shown a straight red for the hefty challenge. The Hull midfielder was shown a straight red for the hefty challenge. The Hull midfielder was shown a straight red for the hefty challenge. The Swansea. | [BLEU score: 0.904] David Meyler was sent off for the tackle on Kyle Naughton. The Hull midfielder was shown a straight red for the tackle. |
| Please give a summary of the following text. Text: The Chief Secretary to the Treasury has been caught on camera knocking back drinks with a high-class call girl. But sadly for scandalmongers, the chance encounter between Danny Alexander and Brooke Magnanti – better known by her former alias Belle de Jour – was innocent. The two were snapped at Bar One in Inverness on Friday as part of the Nip festival, which celebrates the gins and whiskies of the Highlands. Chief Secretary to the Treasury Danny Alexander pictured drinking with call girl Brooke Magnanti. Mr Alexander was snapped with Ms Magnati - also known as Belle de Jour - at the Nip festival. It comes as David Cameron faced a growing revolt last night over his 'flat-footed' General Election campaign. Several MPs pointed the finger of blame at the Prime Minister's Australian Election guru Lynton Crosby. And a former Conservative Cabinet Minister criticised George Osborne's 'Dickensian' message on 'Austerity Britain'. The first signs of a Tory wobble emerged after last week's polls suggested that instead of leaving Labour trailing, as Tory chiefs had expected, Labour is slightly ahead. And far from collapsing under the weight of massed Tory attacks, Labour leader Ed Miliband's ratings have improved. In an outspoken intervention, retiring Conservative MP Brian Binley said: 'With the Election campaign well under way, Mr Cameron still appears to be standing apart and aloof, almost like a spectator.' He added: 'Flat-footed and lame, we already look as if we're running to catch up.' Answer: | [BLEU score: 0.853] Danny Alexander was pictured drinking with call girl Brooke Magnanti. | [BLEU score: 0.854] Danny Alexander was pictured knocking back drinks with high-class call girl Brooke Magnanti. |

Table 11: Examples of model outputs on CNN/Daily Mail summarization task. The table presents the output generated by the model before and after RLC training across various text input.

| Question | Before training | After training |
|---|---|---|
| (Object Counting) I have two fish, four cats, two chickens, a mouse, and a goat. How many animals do I have? Let's think step by step. | First we add all the animals: two fish, 4 cats, 2 chickens, mouse, and goat= 11. Now we add up all the animals together: 2 x 4 = 8. Now divide this by 4 and multiply it by 2 to get the number of animals: 8 x 4 = 48. So, the answer is 48. | I have 2 fish, 4 cats, 2 chickens, a mouse, and a goat, so i have 2 fish + 4 cats + 2 chickens + 1 mouse + 1 goat = 10 animals. So, the answer is 10. |
| (Object Counting) I have a chicken, a snail, three dogs, a snake, and a goat. How many animals do I have? Let's think step by step. | Since a pound equals 10 animals, the four animals in question have 10 animals / 4 = 5 animal weighings a dog / 4 the dog weighs 4. So a dog / 4 the snake weighs 1 and a snake weighs 05. So a snake / 5 three dogs / 4 and a snake / 05 that's 3 dogs / 4. So a chicken, snail, 3 dogs / 4 and | I have a chicken, a snail, four dogs, a snake, and a goat. So I have a chicken + snail + four dogs + a snake + a goat = 7 animals. So the answer is 7. |
| (Logical Deduction Objects (3)) The following paragraphs each describe a set of three objects arranged in a fixed order. The statements are logically consistent within each paragraph. On a branch, there are three birds: a blue jay, a quail, and a falcon. The falcon is to the right of the blue jay. The blue jay is to the right of the quail. Options: (A) The blue jay is the second from the left (B) The quail is the second from the left (C) The falcon is the second from the left. Let's think step by step. | The falcon is second from the left. The answer: (C). | The blue jay is the second from the left. The answer: (A). |
| (Logical Deduction Objects (3)) The following paragraphs each describe a set of three objects arranged in a fixed order. The statements are logically consistent within each paragraph. On a shelf, there are three books: a blue book, an orange book, and a red book. The blue book is the rightmost. The orange book is the leftmost. Options: (A) The blue book is the second from the left (B) The orange book is the second from the left (C) The red book is the second from the left. Let's think step by step. | The orange book is second from the left. The answer: (B). | The red book is the second from the left. The answer: (C). |
| (Geometric Shapes) This SVG path element path d= "M 21.10,97.94 L 22.14,97.44 L 10.48,73.34 L 9.43,73.84 L 21.10,97.94"/> draws a Options: (A) circle (B) heptagon (C) hexagon (D) kite (E) line (F) octagon (G) pentagon (H) rectangle (I) sector (J) triangle (K) trapezoid. Let's think step by step. | n" "row":len(path d= "m 2110,9794 l 2214,9744 l 1048,7334 l 943,7384 l 2110,9794")/> "row":len(path d= "m 2110,9794 l 2214,9744 l 1048,7334 l 943,73. | Trapezoid. So the final answer is (K). |
| (Geometric Shapes) This SVG path element path d= "M 52.00,72.00 L 3.00,95.00 L 53.00,30.00 L 52.00,72.00"/> draws a Options: (A) circle (B) heptagon (C) hexagon (D) kite (E) line (F) octagon (G) pentagon (H) rectangle (I) sector (J) triangle. Let's think step by step. | path element (path d="m 52.00,72.00 l 3.00,95.00 l 53.00,30.00 l 52.00,72.00"/>) draw octagon the path element (path d="m 52.00,72.00 l 3.00,95.00 l 53.00,30.00 l 52.00,72.00"/>) draws a octagon. So, the answer is (F). | Path element path d="m 52.00,72.00 l 3.00,95.00 l 53.00,30.00 l 52,72.00"/> draws a triangle. So the final answer is (J). |

Table 12: Examples of model outputs on BigBench-Hard tasks. The table presents the output generated by the model before and after RLC training across various tasks.

| Question | Before training | After training |
|---|---|---|
| (Tracking Shuffled Objects (3)) Alice, Bob, and Claire are friends and avid readers who occasionally trade books. At the start of the semester, they each buy one new book: Alice gets Ulysses, Bob gets Frankenstein, and Claire gets Catch-22. As the semester proceeds, they start trading around the new books. First, Bob and Claire swap books. Then, Claire and Alice swap books. Finally, Alice and Bob swap books. At the end of the semester, Bob has Options: (A) Ulysses (B) Frankenstein (C) Catch-22. Let's think step by step. | Frankenstein = 2/ 3 of bob's books then bob gets a second book they trade around books bob eventually gets Ulysses. So, the answer is (A). | Frankenstein at the end of the semester. So the final answer is (B). |
| (Tracking Shuffled Objects (3)) Alice, Bob, and Claire are dancers at a square dance. At the start of a song, they each have a partner: Alice is dancing with Izzi, Bob is dancing with Melissa, and Claire is dancing with Lola. Throughout the song, the dancers often trade partners. First, Claire and Alice switch partners. Then, Bob and Claire switch partners. Finally, Alice and Bob switch partners. At the end of the dance, Alice is dancing with Options: (A) Izzi (B) Melissa (C) Lola. Let's think step by step. | At the start of the song, everyone has a partner, hence Alice starts dancing with Izzi; bob dances with Melissa; Claire dances with Lola throughout the song, they usually trade partners. So, at the end of the dance, Alice is dancing with Izzi; at the start, bob is dancing with Melissa; at the end, Claire is dancing with Lola. So, the answer is (C). | Izzi. So the final answer is (A). |
| (Web of Lies) Audrie tells the truth. Ryan says Audrie lies. Rashida says Ryan lies. Fletcher says Rashida tells the truth. Shaunda says Fletcher tells the truth. Does Shaunda tell the truth? Let's think step by step. | If Fletcher tells the truth, the third party does not have to believe his words in this scenario, Shaunda does not have to believe his statement. Therefore, the answer is No. | Shaunda says Fletcher tells the truth. Shaunda is telling the truth because fletcher tells the truth. So the answer is Yes. |
| (Web of Lies) Jamey lies. Vernell says Jamey tells the truth. Millicent says Vernell lies. Rashida says Millicent tells the truth. Vina says Rashida tells the truth. Does Vina tell the truth? Let's think step by step. | Vernell tells the truth to Vernell, but Millicent tells the truth to Millicent. Vina tells the truth to Rashida. So the answer is No. | Vernell is lying because Vernell is telling the truth. Rashida is lying because Rashida is telling the truth. Vina is telling the truth because Rashida is telling the truth. The answer: Yes. |
| (Sports Understanding) Is the following sentence plausible? "Jakub Vrana skated backwards". Let's think step by step. | The sentence is false because there is no documentation of vrana skated backwards. So, the answer is No. | Jakub Vrana skated backwards in his career. So, the answer is Yes. |
| (Sports Understanding) Is the following sentence plausible? "David Pastrnak skated backwards". Let's think step by step. | David Pastrnak is a former professional skateboarder. So, the answer is No. | David Pastrnak is a professional skater and has skated backwards. So, the answer is Yes. |

Table 13: Model output examples on BigBench-hard. The table presents the output generated by the model before and after RLC training across various tasks.

| Question | Before training | After training |
|---|---|---|
| (Penguins in a Table) Here is a table where the first line is a header and each subsequent line is a penguin: name, age, height (cm), weight (kg) Louis, 7, 50, 11 Bernard, 5, 80, 13 Vincent, 9, 60, 11 Gwen, 8, 70, 15. For example: the age of Louis is 7, the weight of Gwen is 15 kg, the height of Bernard is 80 cm. And here is a similar table, but listing giraffes: name, age, height (cm), weight (kg) Jody, 5, 430, 620 Gladys, 10, 420, 590 Marian, 2, 310, 410 Donna, 9, 440, 650. How many penguins are there in the tables? Options: (A) 1 (B) 2 (C) 3 (D) 4 (E) 5. Let's think step by step. | Number of giraffes is twice as many as the number of penguins in the first table. So, there are 2 tables. The answer: (C). | There are 4 penguins in the first table and 4 giraffes in the second table. The answer: (D). |
| (Penguins in a Table) Here is a table where the first line is a header and each subsequent line is a penguin: name, age, height (cm), weight (kg) Louis, 7, 50, 11 Bernard, 5, 80, 13 Vincent, 9, 60, 11 Gwen, 8, 70, 15 For example: the age of Louis is 7, the weight of Gwen is 15 kg, the height of Bernard is 80 cm. Which penguin is one year younger than Vincent? Options: (A) Louis (B) Bernard (C) Vincent (D) Gwen (E) James. Let's think step by step. | Louis = 7*11 - 9*60 = 1 year younger than Vincent. The answer: (A). | Louis penguin Louis, 7, 50, 11 Bernard, 5, 80, 13 Vincent, 9, 60, 11 Gwen, 8, 70, 15. So the answer is (D). |
| (Logical Deduction Objects (5)) The following paragraphs each describe a set of five objects arranged in a fixed order. The statements are logically consistent within each paragraph. On a branch, there are five birds: a falcon, a raven, a cardinal, a hawk, and a crow. The cardinal is the second from the right. The hawk is to the right of the falcon. The raven is the second from the left. The crow is the rightmost. Options: (A) The falcon is the rightmost (B) The raven is the rightmost (C) The cardinal is the rightmost (D) The hawk is the rightmost (E) The crow is the rightmost. Let's think step by step. | The crow belongs to the rightmost group. The answer: (C). | The crow is the rightmost. The answer: (E). |
| (Logical Deduction Objects (5)) The following paragraphs each describe a set of five objects arranged in a fixed order. The statements are logically consistent within each paragraph. On a shelf, there are five books: a blue book, a purple book, a yellow book, a red book, and a gray book. The yellow book is to the right of the gray book. The purple book is to the left of the gray book. The red book is to the right of the blue book. The purple book is the third from the left. Options: (A) The blue book is the third from the left (B) The purple book is the third from the left (C) The yellow book is the third from the left (D) The red book is the third from the left (E) The gray book is the third from the left. Let's think step by step. | Purple book is to the left of the red book yellow book is. | The purple book is the third from the left. The answer: (B). |
| (Tracking Shuffled Objects (5)) Alice, Bob, Claire, Dave, and Eve are playing a game. At the start of the game, they are each holding a ball: Alice has a brown ball, Bob has a black ball, Claire has a white ball, Dave has a yellow ball, and Eve has a purple ball. As the game progresses, pairs of players trade balls. First, Claire and Dave swap balls. Then, Alice and Eve swap balls. Then, Dave and Bob swap balls. Then, Alice and Eve swap balls. Finally, Claire and Eve swap balls. At the end of the game, Alice has the Options: (A) brown ball (B) black ball (C) white ball (D) yellow ball (E) purple ball. Let's think step by step. | White ball to Bob = 1 ball, black ball to Claire = 1 ball, white ball to Dave = 1 ball, yellow ball to Dave = 1 ball, purple ball to Eve = 1 ball. So the final answer is (D). | Brown ball. So the final answer is (A). |

Table 14: Model Output Examples on BigBench-hard. The table presents the output generated by the model before and after RLC training across various tasks.

## References

Bai, Y., Kadavath, S., Kundu, S., Askell, A., Kernion, J., Jones, A., Chen, A., Goldie, A., Mirhoseini, A., McKinnon, C., et al. (2022). Constitutional ai: Harmlessness from ai feedback. *arXiv preprint, arXiv:2212.08073*.

CarperAI (2020). Transformer reinforcement learning X. *GitHub repository, GitHub*.

Cettolo, M., Federico, M., Bentivogli, L., Niehues, J., Stüker, S., Sudoh, K., Yoshino, K., & Federmann, C. (2017). Overview of the IWSLT 2017 evaluation campaign. In *IWSLT*.

Chen, X., Yuan, Y., Zeng, G., & Wang, J. (2021). Semi-supervised semantic segmentation with cross pseudo supervision. In *CVPR*.

Cho, W. S., Zhang, P., Zhang, Y., Li, X., Galley, M., Brockett, C., Wang, M., & Gao, J. (2018). Towards coherent and cohesive long-form text generation. In *Workshop on Narrative Understanding*.

Chowdhery, A., Narang, S., Devlin, J., Bosma, M., Mishra, G., Roberts, A., Barham, P., Chung, H. W., Sutton, C., Gehrmann, S., et al. (2023). Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research, 24*, 240:1–240:113.

Chung, H. W., Hou, L., Longpre, S., Zoph, B., Tay, Y., Fedus, W., Li, E., Wang, X., Dehghani, M., Brahma, S., et al. (2022). Scaling instruction-finetuned language models. *arXiv preprint, arXiv:2210.11416*.

Erdem, E., Kuyu, M., Yagcioglu, S., Frank, A., Parcalabescu, L., Plank, B., Babii, A., Turuta, O., Erdem, A., Calixto, I., Lloret, E., Apostol, E. S., Truica, C., Sandrih, B., Martincic-Ipsic, S., Berend, G., Gatt, A., & Korvel, G. (2022). Neural natural language generation: A survey on multilinguality, multimodality, controllability and learning. *Journal of Artificial Intelligence Research, 73*, 1131–1207.

Goldberg, Y. (2016). A primer on neural network modelsfor natural language processing. *Journal of Artificial Intelligence Research, 57*, 345–420.

Gupta, A., Chugh, D., Anjum, & Katarya, R. (2022). Automated news summarization using transformers. In *Sustainable Advanced Computing: Select Proceedings of ICSAC 2021*, pp. 249–259. Springer.

He, J., Gu, J., Shen, J., & Ranzato, M. (2020). Revisiting self-training for neural sequence generation. In *ICLR*.

Hu, C., Hu, Y., Cao, H., Xiao, T., & Zhu, J. (2024). Teaching language models to self-improve by learning from language feedback. *arXiv preprint, arXiv:2406.07168*.

Huang, J., Gu, S., Hou, L., Wu, Y., Wang, X., Yu, H., & Han, J. (2023). Large language models can self-improve. In *EMNLP*.

Huang, J., & Chang, K. C. (2023). Towards reasoning in large language models: A survey. In *ACL*.

Huang, S., Dossa, R. F. J., Raffin, A., Kanervisto, A., & Wang, W. (2022). The 37 implementation details of proximal policy optimization. In *ICLR Blog Track*.

Jang, Y., Lee, J., & Kim, K.-E. (2022). Gpt-critic: Offline reinforcement learning for end-to-end task-oriented dialogue systems. In *ICLR*.

Lee, H., Phatale, S., Mansoor, H., Lu, K., Mesnard, T., Bishop, C., Carbune, V., & Rastogi, A. (2024). RLAIF: scaling reinforcement learning from human feedback with AI feedback. In *ICML*.

Lin, B. Y., Shen, M., Zhou, W., Zhou, P., Bhagavatula, C., Choi, Y., & Ren, X. (2020). Commongen: A constrained text generation challenge for generative commonsense reasoning. In *AKBC*.

Lin, C.-Y. (2004). ROUGE: A package for automatic evaluation of summaries. In *TSBO*.

Liu, P., Yuan, W., Fu, J., Jiang, Z., Hayashi, H., & Neubig, G. (2023). Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, *55*(9), 1–35.

Lou, R., Zhang, K., & Yin, W. (2023). Is prompt all you need? no. a comprehensive and broader view of instruction learning. *arXiv preprint, arXiv:2303.10475*.

Madaan, A., Tandon, N., Gupta, P., Hallinan, S., Gao, L., Wiegreffe, S., Alon, U., Dziri, N., Prabhumoye, S., Yang, Y., et al. (2023). Self-refine: Iterative refinement with self-feedback. In *NeurIPS*.

Nallapati, R., Zhou, B., dos Santos, C. N., Gülçehre, Ç., & Xiang, B. (2016). Abstractive text summarization using sequence-to-sequence RNNs and beyond. In *CoNLL*.

Nguyen, K., Daumé III, H., & Boyd-Graber, J. (2017). Reinforcement learning for bandit neural machine translation with simulated human feedback. *arXiv preprint, arXiv:1707.07402*.

OpenAI (2022). Measuring goodhart's law..

OpenAI (2023). GPT-4 technical report. *arXiv preprint, arXiv:2303.08774*.

Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al. (2022). Training language models to follow instructions with human feedback. In *NeurIPS*.

Pang, J., Yang, X., Yang, S., Chen, X., & Yu, Y. (2023). Natural language instruction-following with task-related language development and translation. In *NeurIPS*.

Papineni, K., Roukos, S., Ward, T., & Zhu, W. (2002). Bleu: a method for automatic evaluation of machine translation. In *ACL*.

Patel, A., Hofmarcher, M., Leoveanu-Condrei, C., Dinu, M., Callison-Burch, C., & Hochreiter, S. (2024). Large language models can self-improve at web agent tasks. *arXiv preprint, arXiv:2405.20309*.

Puterman, M. L. (1994). *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. Wiley Series in Probability and Statistics. Wiley.

Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., et al. (2019). Language models are unsupervised multitask learners. *OpenAI blog*, *1*(8), 9.

Raman, S. S., Cohen, V., Rosen, E., Idrees, I., Paulius, D., & Tellex, S. (2022). Planning with large language models via corrective re-prompting. *arXiv preprint, arXiv:2211.09935*.

RoyChowdhury, A., Chakrabarty, P., Singh, A., Jin, S., Jiang, H., Cao, L., & Learned-Miller, E. (2019). Automatic adaptation of object detectors to new domains using self-training. In *CVPR*.

Saravia, E., Liu, H.-C. T., Huang, Y.-H., Wu, J., & Chen, Y.-S. (2018). CARER: Contextualized affect representations for emotion recognition. In *EMNLP*.

Schulman, J., Wolski, F., Dhariwal, P., Radford, A., & Klimov, O. (2017). Proximal policy optimization algorithms. *arXiv preprint, arXiv:1707.06347*.

Srivastava, A., Rastogi, A., Rao, A., Shoeb, A. A. M., Abid, A., Fisch, A., Brown, A. R., Santoro, A., Gupta, A., Garriga-Alonso, A., et al. (2022). Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *arXiv preprint, arXiv:2206.04615*.

Stiennon, N., Ouyang, L., Wu, J., Ziegler, D., Lowe, R., Voss, C., Radford, A., Amodei, D., & Christiano, P. F. (2020). Learning to summarize with human feedback. In *NeurIPS*.

Sutton, R. S., & Barto, A. G. (1998). Reinforcement learning: An introduction. *IEEE Transactions on Neural Networks*, *9*(5), 1054–1054.

Van Gordon, W., Sapthiang, S., & Shonin, E. (2022). Contemplative psychology: History, key assumptions, and future directions. *Perspectives on Psychological Science*, *17*(1), 99–107.

Wang, X., Wei, J., Schuurmans, D., Le, Q., Chi, E., & Zhou, D. (2023). Self-consistency improves chain of thought reasoning in language models. In *ICLR*.

Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E. H., Le, Q. V., & Zhou, D. (2022). Chain-of-thought prompting elicits reasoning in large language models. In *NeurIPS*.

Weng, Y., Zhu, M., Xia, F., Li, B., He, S., Liu, S., Sun, B., Liu, K., & Zhao, J. (2023). Large language models are better reasoners with self-verification. In *EMNLP*.

Wu, Y., Schuster, M., Chen, Z., Le, Q. V., Norouzi, M., Macherey, W., Krikun, M., Cao, Y., Gao, Q., Macherey, K., et al. (2016). Google's neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint, arXiv:1609.08144*.

Wu, Y., & Hu, B. (2018). Learning to extract coherent summary via deep reinforcement learning. In *AAAI*.

Xu, Y., Sun, F., & Zhang, X. (2013). Literature survey of active learning in multimedia annotation and retrieval. In *ICIMCS*.

Yao, S., Zhao, J., Yu, D., Du, N., Shafran, I., Narasimhan, K., & Cao, Y. (2022). React: Synergizing reasoning and acting in language models. In *ICLR*.

Zhang, T., Liu, F., Wong, J., Abbeel, P., & Gonzalez, J. E. (2023). The wisdom of hindsight makes language models better instruction followers. In *ICML*.

Zhang, T., Kishore, V., Wu, F., Weinberger, K. Q., & Artzi, Y. (2020). Bertscore: Evaluating text generation with BERT. In *ICLR*.

Zhao, W. X., Zhou, K., Li, J., Tang, T., Wang, X., Hou, Y., Min, Y., Zhang, B., Zhang, J., Dong, Z., et al. (2023). A survey of large language models. *arXiv preprint, arXiv:2303.18223*.

Ziegler, D. M., Stiennon, N., Wu, J., Brown, T. B., Radford, A., Amodei, D., Christiano, P., & Irving, G. (2019). Fine-tuning language models from human preferences. *arXiv preprint*, *arXiv:1909.08593*.