

BEYOND IMITATION: LEARNING KEY REASONING STEPS FROM DUAL CHAIN-OF-THOUGHTS IN REASONING DISTILLATION

Anonymous authors

Paper under double-blind review

ABSTRACT

As Large Language Models (LLMs) scale up and gain powerful Chain-of-Thoughts (CoTs) reasoning abilities, practical resource constraints drive efforts to distill these capabilities into more compact Smaller Language Models (SLMs). We find that CoTs consist mainly of simple reasoning forms, with a small proportion ($\approx 4.7\%$) of key reasoning steps that truly impact conclusions. However, previous distillation methods typically involve supervised fine-tuning student SLMs only on correct CoTs data produced by teacher LLMs, resulting in students struggling to learn the key reasoning steps, instead imitating the teacher’s reasoning forms and making errors or omissions on these steps. To address these issues, drawing an analogy to human learning, where analyzing mistakes according to correct solutions often reveals the crucial steps leading to successes or failures, we propose mistake-Driven key reasoning step distillation (**EDIT**), a novel method that further aids SLMs learning key reasoning steps rather than mere simple fine-tuning. Firstly, to expose these crucial steps in CoTs, we design specific prompts to generate dual CoTs data with similar reasoning paths but divergent conclusions. Then, we apply the minimum edit distance algorithm on the dual CoTs data to locate these key steps and optimize the likelihood of these steps. Extensive experiments validate the effectiveness of EDIT across both in-domain and out-of-domain benchmark reasoning datasets. Further analysis shows that EDIT can generate high-quality CoTs with more correct key reasoning steps. Notably, we also explore how different mistake patterns affect performance and find that EDIT benefits more from logical errors than from knowledge or mathematical calculation errors in dual CoTs¹.

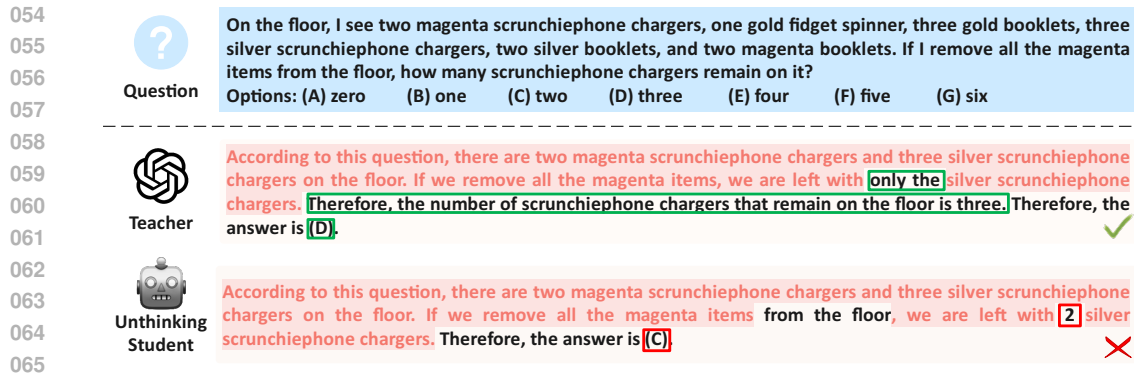
1 INTRODUCTION

With the rapid growth in model size and pre-training data, LLMs have demonstrated impressive CoT reasoning performance in natural language processing (NLP) (Brown et al., 2020; Hoffmann et al., 2022; Chowdhery et al., 2023; OpenAI, 2023b). However, due to the giant model architecture and massive parameters (e.g. GPT-3 (Brown et al., 2020) with 175 billion parameters), the deployment of LLMs in resource-constrained environments becomes challenging.

To address this, researchers (Xu et al., 2023; Jiang et al., 2023b) have explored distilling knowledge from LLMs into smaller language models (SLMs) via instruction-tuning, as seen in LMs like Alpaca (Taori et al., 2023) and Vicuna (Chiang et al., 2023). Despite progress, these distilled models often struggle with complex causal reasoning. To enhance this capability, some studies (Magister et al., 2023; Ho et al., 2023; Fu et al., 2023; Chen et al., 2024; Zhou & Ai, 2024) explore distilling the CoT reasoning ability from LLMs of over 100B parameters (Wei et al., 2022a;b) by fine-tuning on CoTs data annotated by teacher LLMs, known as standard CoTs distillation. Besides, other studies (Hsieh et al., 2023; Li et al., 2022; Liu et al., 2023) propose distilling CoTs within a multi-task learning framework by incorporating additional objectives. However, CoTs usually consist mainly of simple reasoning forms, with a small proportion ($\approx 4.7\%$) of key reasoning steps that are pivotal moments

¹Code can be found at <https://anonymous.4open.science/r/eb77sh-F564>

²We calculated the edit distance and its average proportion in the overall sequence on the dual CoT dataset mentioned in our subsequent methods §3.



066 Figure 1: Examples of CoTs generated by teacher LLMs and student SLMs on our test dataset.
067 Simply SFT leads to an "unthinking" student who imitates the teacher's reasoning forms but makes
068 errors and omissions in key reasoning steps, where the imitated contents are highlighted in red, and
069 the key steps are marked with boxes.

072 in reasoning that significantly influence subsequent thought processes and conclusions. The essence
073 of the above methods is the simple Supervised Fine-Tuning (SFT) paradigm, where the student
074 model is trained solely on the teacher's correct reasoning data. This paradigm may result in students
075 **struggling to learn the key reasoning steps, instead imitating the teacher's reasoning forms
076 and making errors or omissions on these steps**, as illustrated in Figure 1. Drawing an analogy
077 to human learning, where analyzing mistakes according to correct solutions often reveals the key
078 reasoning steps leading to successes or failures, we propose a novel mistake-Driven key reasonIng
079 step distillaTion (**EDIT**). This approach focuses on dual CoTs data, encompassing both positive
080 and negative examples of teachers' reasoning. By examining dual CoTs, students can identify and
081 learn from the crucial reasoning steps, thereby improving their CoTs. Specifically, we first retain
082 all CoTs data annotated by the teacher, irrespective of correctness. Subsequently, we design two
083 comprehensive prompts to instruct teachers to produce dual CoTs that share similar intermediate
084 reasoning steps but lead to divergent conclusions. Finally, we utilize the minimum edit distance
085 algorithm to locate key reasoning steps in dual CoTs, as shown in Figure 3, and then utilize a
086 fine-grained loss function to optimize the likelihood of these steps.

087 Extensive experiments show that the student models distilled by EDIT exhibits higher performance
088 and generalization than the baselines on both in-domain (IND) and out-of-domain (OOD) benchmark
089 reasoning datasets. Further analyses indicate that EDIT can generate higher-quality CoTs with more
090 correct key reasoning steps by auto evaluation and case studies. Notably, we also show EDIT can
091 benefit more from logical mistake patterns than knowledge or mathematical calculation errors in dual
092 CoTs, potentially paving the way for future research on the efficient use of mistakes.

093 Our contributions can be summarized as follows:

- 094 1. We reveal a shortfall in the previous distillation methods, where the simple SFT paradigm may
095 result in students mimicking the teacher's reasoning forms but making errors or omissions in
096 key reasoning steps, thus diminishing the versatility of CoTs.
- 097 2. We propose mistake-driven key reasoning step distillation, which allows students to learn key
098 reasoning steps from our specifically designed dual CoTs data, further improving reasoning.
- 099 3. Extensive experiments validate the effectiveness of our method across both IND and OOD
100 datasets, showing that EDIT can reduce errors in key reasoning steps for students.
- 101 4. We investigate how different mistake patterns impact EDIT and find that logical errors provide
102 the more significant benefits than knowledge or mathematical calculation errors.

104 2 RELATED WORKS

106 **CoT Reasoning.** The emergent ability appears in LLMs across a wide range of NLP tasks (Chowdhery
107 et al., 2023; Wei et al., 2022a). One such ability is CoT reasoning, which involves generating a

series of intermediate reasoning steps. This ability has been further explored recently with the release of OpenAI’s o1 model (OpenAI, 2024). While CoT prompting techniques (Wei et al., 2022b) significantly enhance the problem-solving capabilities of models (Kojima et al., 2022; Wang et al., 2023b; Huang et al., 2023), it has little effect on smaller models (Wei et al., 2022a). Chung et al. (2022) suggest that CoT reasoning can be induced in SLMs via instruction tuning on CoTs data. Our work show that the CoT capabilities of SLMs can be further improved by learning from key reasoning steps in dual CoTs data.

Knowledge Distillation from LLMs. There has been a lot of work dedicated to distilling knowledge (Hinton et al., 2015) from powerful proprietary LLMs, e.g. ChatGPT (OpenAI, 2023a) in a black-box setting. However, most of these works primarily focus on the general ability distillation by instruction tuning on large and diverse datasets (Peng et al., 2023; Jiang et al., 2023b; Li et al., 2024). In contrast, we aim to distill the CoT reasoning capabilities from LLMs same as the standard CoTs distillation (Magister et al., 2023; Ho et al., 2023). Besides, some studies (Li et al., 2022; Hsieh et al., 2023; Liu et al., 2023) employ LLM’s rationale or self-evaluation output to enhance SLM’s reasoning in a multi-task learning framework. Fu et al. (2023) fine-tune SLMs on four types of reasoning data to ensure out-of-distribution generalization. Wang et al. (2023c) distill SLMs by learning from self-reflection and feedback from LLMs in an interactive multi-round paradigm. Chen et al. (2024) maximize the mutual information between multi objectives for CoTs distillation. Ranaldi & Freitas (2024) use in-family and out-family teachers to generate more CoTs for fine-tuning students. Different from the above works, we assist CoTs distillation with teachers’ mistakes to alleviate the style imitation of teachers’ reasoning.

Learning from Mistakes. Recent studies use mistake data to enhance the performance of LMs. Shinn et al. (2023) propose Reflexion that allows the LLM agent to self-reflect from its mistakes. Wang & Li (2023) introduce a study assistant that collects and retrieves LLMs’ training mistakes to guide future inferences. Li et al. (2023) propose CoK that corrects potential mistakes in the rationale by retrieving knowledge to avoid error propagation. However, both of the above methods require the models to be large enough to have basic CoT reasoning or instruction-following capabilities, which is almost impossible to occur in vanilla SLMs. Wang et al. (2023a) propose fine-tuning on counterfactual data to ensure the faithful reasoning of the student model. An et al. (2023) propose LEMA that fine-tunes language models on corrected mistake data, where the mistakes are collected from various LLMs e.g. LLaMA2-70B (Touvron et al., 2023), WizardLM-70B (Xu et al., 2023), and corrected by GPT-4 (OpenAI, 2023b). Additionally, Sun et al. (2024) propose Retrieved In-Context Principles, which retrieve mistakes to provide customized guidance and improve model performance during inference. In contrast, we collect the teachers’ mistakes to create a dual CoTs dataset for further key reasoning steps learning on model distillation.

3 MISTAKE-DRIVEN KEY REASONING STEP DISTILLATION

We present the overview of our proposed method in Figure 2. Concretely, (1) unlike prior works (Magister et al., 2023; Hsieh et al., 2023) that only focus on correct CoTs annotated by teacher LLMs, we first retain all CoTs reasoning data, regardless of its correctness. (2) Then based on the previously retained correct and wrong CoTs, we construct dual CoTs datasets consisting of positive-negative CoT pairs that follow similar intermediate reasoning steps but lead to divergent conclusions. Specifically, we design two comprehensive contextual prompts to instruct teacher LLMs to rectify the originally wrong CoTs and corrupt originally correct CoTs. (3) Finally, we distill the student SLMs by training on the teacher’s correct CoTs reasoning data and further Key Reasoning Steps Learning (KRSL) on the dual CoTs datasets.

3.1 CoTs ANNOTATED BY LLMs

We utilize CoT Prompting (Wei et al., 2022b) to extract CoTs for a raw dataset $\mathcal{D} = \{(q, a)\}$ from LLMs, where q is the question and a is the golden answer. Specifically, we first create a CoTs Extraction Prompt CEP that contains several human-curated question-CoTs pair examples and the task description, which can be found in Appendix C.1. For each $q \in \mathcal{D}$, we extract CoTs as:

$$CoT \sim LLM(CEP \oplus q) \quad (1)$$

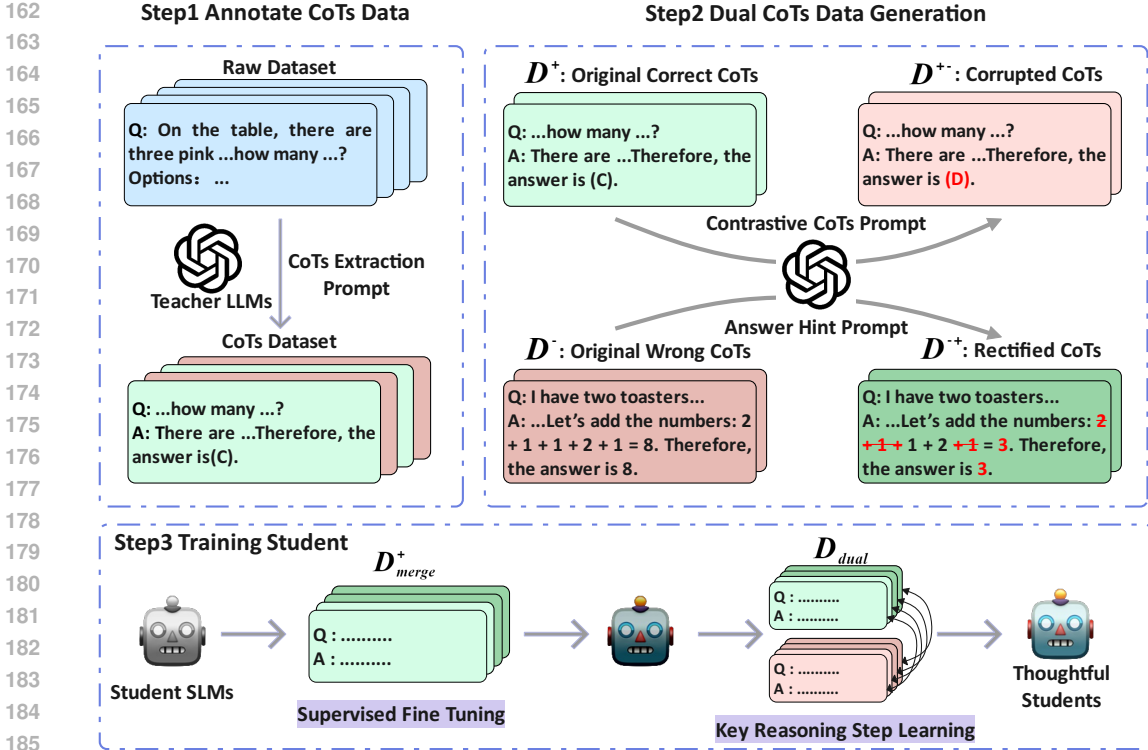


Figure 2: **Overview of our mistake-driven key reasoning step distillation.** (1) We first retain all CoTs data annotated by teacher LLMs (2) and ask teacher LLMs to generate dual CoTs data using our designed two comprehensive prompts. (3) Then we fine-tune student SLMs on both original correct and rectified-after CoTs data. Finally, we apply key reasoning step learning on the pre-tuned student SLMs by identifying the minor difference between the dual CoTs.

where \oplus means concatenation. Then, we classify the CoTs annotated dataset into two datasets according to the final answer’s correctness³, same as Zelikman et al. (2022). One is the CoTs-original correct dataset $\mathcal{D}^+ = \{(q, CoT^+) \mid \forall (q, a) \in \mathcal{D}, \hat{a} = a \ \& \ \hat{a} \in CoT^+\}$ and the other is CoTs-original wrong dataset $\mathcal{D}^- = \{(q, CoT^-) \mid \forall (q, a) \in \mathcal{D}, \hat{a} \neq a \ \& \ \hat{a} \in CoT^-\}$.

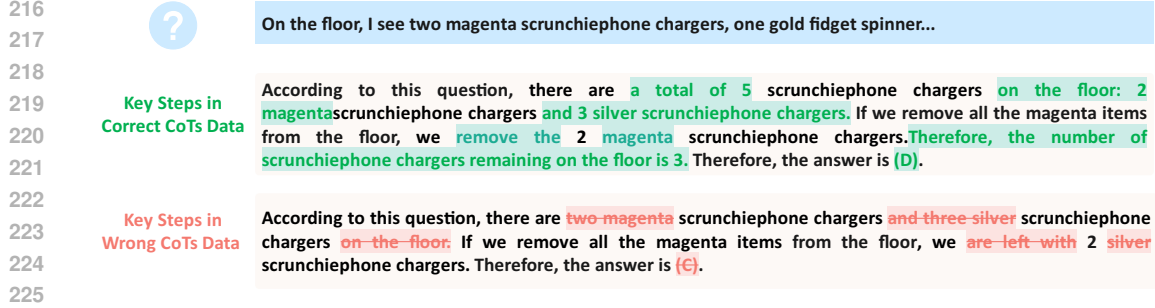
3.2 DUAL CO TS GENERATION

We define dual CoTs data as contrasting CoTs that follow similar reasoning steps but reach divergent conclusions compared to the original. To provide a deeper understanding, we also present several examples of dual CoTs in Appendix A. In the following, we will introduce how to generate dual CoTs datasets including \mathcal{D}^{+-} contrasting to \mathcal{D}^+ , and \mathcal{D}^{-+} contrasting to \mathcal{D}^- .

Rectify Wrong CoTs. To generate correct CoTs contrasting with the originally wrong CoTs, inspired by Rationalization (Zelikman et al., 2022), we design an Answer Hint Prompt AHP that shares the same examples with CEP but with different organizational structures. The template of AHP can be found in Appendix C.2. Each example in the context and the final provided question will be inserted with a hint that tells LLMs the answer first before CoTs. Thus, due to the same in-context examples and hint answers, teacher LLM can rectify its original wrong CoTs data with similar reasoning steps but correct answers. For each $q \in \mathcal{D}^-$, we rectify CoTs as follows and then have the Rectified CoTs dataset $\mathcal{D}^{-+} = \{(q, CoT^{-+})\}$:

$$CoT^{-+} \sim LLM(AHP \oplus q \oplus a) \quad (2)$$

³To support our assumption of CoT correctness, We randomly sample 100 examples to manually check the logical consistency between the CoT and the final answer and find that the CoTs generated by ChatGPT generally support the final answer.



226 Figure 3: Examples of locating key reasoning steps in dual CoTs, where the correct CoT and the
 227 wrong CoTs are dual to each other. The identified key steps in correct reasoning and wrong reasoning
 228 are respectively marked in green and red.

229

230 **Corrupt Correct CoTs.** To generate incorrect CoTs contrasting with the originally correct CoTs,
 231 a straightforward approach is to use AHP with incorrect hint answers to prompt LLMs to produce
 232 wrong CoTs. However, in practice, we find that LLMs rarely follow the incorrect hints and still
 233 generate correct CoTs. This may be due to the simplicity of the questions, which fall within the
 234 LLMs’ knowledge range. Additionally, LLMs, having undergone Reinforcement Learning from
 235 Human Feedback (RLHF) (Ouyang et al., 2022), may resist providing unhelpful answers. Therefore,
 236 we design a Contrastive CoTs Prompt (CEP) to entice LLMs to generate incorrect CoTs, leveraging
 237 their strong in-context learning capabilities. The prompt template can be found in Appendix C.3.
 238 Specifically, to ensure high-quality incorrect CoTs, we randomly sample negative examples from \mathcal{D}^-
 239 and positive examples from \mathcal{D}^+ , pair them, and place them into the CCP as curated joint in-context
 240 examples. For each $q \in \mathcal{D}^+$, we corrupt CoTs as follows and then have the corrupted CoTs dataset
 241 $\mathcal{D}^{+-} = \{(q, CoT^{+-})\}$:

$$242 CoT^{+-} \sim LLM(\text{CCP} \oplus q \oplus CoT^+) \quad (3)$$

243 3.3 TRAINING STUDENT WITH CO Ts

244 **Supervised Fine-tuning on Correct CoTs.** After preparing the dual CoTs, we first fine-tune
 245 student models on the teachers’ original correct CoTs dataset \mathcal{D}^+ and rectified CoTs dataset \mathcal{D}^{-+} .
 246 The training objective is as follows:

$$247 \pi_{sft} = \arg \max_{\pi} \mathbb{E}_{q, CoT \sim \mathcal{D}_{merge}^+} [\log \pi(CoT | q)] \quad (4)$$

248 where the merged correct CoTs dataset $\mathcal{D}_{merge}^+ = \mathcal{D}^+ \cup \mathcal{D}^{-+}$, and π_{sft} denotes the student with the
 249 base inference ability after the initial fine-tuning.

250 **Key Reasoning Steps Learning** Inspired by (Guo et al., 2023b) who leverage fine-grained quality
 251 signals to align human preference, we propose a key reasoning steps learning (KRSL) method to
 252 further encourage students to comprehend the reasons behind both correct and wrong CoTs.

253 **Step1.** We pair the teacher’s original correct CoTs dataset \mathcal{D}^+ with its corrupted CoTs dataset \mathcal{D}^{+-} ,
 254 creating an originally correct dual CoTs dataset $\mathcal{D}_{dual}^+ = \{(q, CoT^+, CoT^{+-})\}$, where CoT^+
 255 and CoT^{+-} are dual to each other; similarly, the teacher’s inherently wrong dual CoTs dataset
 256 $\mathcal{D}_{dual}^- = \{(q, CoT^{-+}, CoT^-)\}$. By merging them, we obtain the ultimate dual CoTs datasets
 257 $\mathcal{D}_{dual} = \mathcal{D}_{dual}^+ \cup \mathcal{D}_{dual}^-$, which is prepared for the subsequent learning of key reasoning steps.

258 **Step2.** Then we employ the minimum edit distance to identify the key steps in both correct reasoning
 259 and wrong reasoning, as shown in Figure 3. In this way, students can identify less frequent text
 260 segments that are inserted or replaced in wrong CoTs compared to correct CoTs, and vice versa.
 261 These text segments are considered key reasoning steps. After that, we assign token-level weights to
 262 facilitate fine-grained learning for correct CoTs and wrong CoTs in \mathcal{D}_{dual} respectively:

$$263 \omega_t^+ = \begin{cases} \alpha, & \text{if } CoT_t^+ \text{ is inserted or replaced} \\ 0, & \text{otherwise} \end{cases}, \omega_t^- = \begin{cases} \beta, & \text{if } CoT_t^- \text{ is deleted or replaced} \\ 0, & \text{otherwise} \end{cases}. \quad (5)$$

264 where $\alpha \geq 0, \beta \geq 0$ and ω_t^+ represents the weight of t -th token in the correct CoTs (semantically
 265 same with ω_t^-). We set the weights to zero to ignore the impact of identical tokens in the dual CoTs.

Step3. Finally, to ensure that the student makes correct decisions on key steps in correct reasoning, we optimize the student model on these tokens with weighted negative log-likelihood. Conversely, to prevent the student from making key steps present in wrong reasoning, we optimize the student model on these steps with weighted positive log-likelihood. The sum of both is taken as the final loss. The optimization objective is as follows:

$$\max_{\pi_{sft}} \mathbb{E}_{q, CoT^+, CoT^- \sim \mathcal{D}_{data}} [\mathcal{L}(\pi_{sft}, q, CoT^+, \omega^+) - \mathcal{L}(\pi_{sft}, q, CoT^-, \omega^-)] \quad (6)$$

where

$$\mathcal{L}(\pi, q, CoT, \omega) = - \sum_{CoT_t \in CoT} \omega_t \log \pi(CoT_t | q, CoT_{<t}) \quad (7)$$

4 EXPERIMENTS

4.1 EXPERIMENTAL SETUP

In-domain (IND) Dataset: BIG-Bench Hard (BBH) (Suzgun et al., 2023) consists of 27 challenging tasks that span arithmetic, symbolic reasoning, etc. This collection is mainly composed of multiple-choice questions, along with a minority of open-ended questions. To underscore the superiority of our method, we divide the BBH dataset for each subtask into a training set (BBH-train) for distillation and a test set (BBH-test) for in-domain evaluation, following a 4:1 ratio.

Out-of-domain (OOD) Dataset: (1) BIG-Bench Sub (BB-sub) is derived from the BIG-Bench (BB) (Guo et al., 2023a), which includes 203 tasks covering linguistics, mathematics, common-sense reasoning, etc. To simplify our evaluation, we refine the selection of tasks from BB by identifying those associated with keywords such as "multiple-choice" and "reasoning."⁴ Additionally, we exclude any tasks that are part of the BBH dataset, narrowing our pool to 61 distinct subtasks. For each of these subtasks, we randomly sample up to 100 instances, culminating the BB-sub dataset. **(2) AGIEval** (Zhong et al., 2023) is a benchmark that assesses LMs on reasoning capabilities using human exams across various fields, including English, Math, Law, and Logic. We focused on the English multiple-choice questions within this benchmark to evaluate our method’s effectiveness. **(3) AI2 Reasoning Challenge (ARC)** (Clark et al., 2018) comprises **ARC-Easy** and **ARC-Challenge** from middle and high school science exams. ARC-E features simpler questions, while ARC-C includes more challenging ones. We use their test sets for evaluation. Detailed statistics for all mentioned benchmarks are provided in Appendix B.6.1. BigBench, AGIEval, and ARC are standard benchmarks for evaluating LLMs reasoning performance. Specifically, BigBench and AGIEval have been employed in related works (Fu et al., 2023; Jiang et al., 2023b), and ARC is frequently used in technical reports for LLaMA3 (AI@Meta, 2024) and GPT-4 (OpenAI, 2023b).

Models & Implementation Details. We employ the widely-used open-source language model, LLaMA2-7B (Touvron et al., 2023), as our student SLM. For the teacher model, given its performance and cost-effectiveness, we employ OpenAI’s advanced black-box LLM, ChatGPT, specifically using the "gpt-3.5-turbo-0613" variant for extracting CoTs with the same manual prompt that is used in (Suzgun et al., 2023). We employ LoRA (Hu et al., 2022) for parameter-efficient fine-tuning of the student SLMs. We empirically set α in KRSL as 1.0 and β as 0.025. Our experiments leverage a mixed-precision training strategy, carried out on $4 \times$ A100 GPUs. We employ vLLM⁵ (Kwon et al., 2023) to enhance inference speed, using a greedy decoding method for text generation on a single A100 GPU. More training details and hyperparameter settings can be found in Appendix B.6.2.

Baselines. We compare EDIT with the following baselines: (1) **Teacher & Vanilla Student** under various settings, e.g., Zero-shot (+ CoT) or Few-shot (+ CoT). (2) **Std-CoT** (Magister et al., 2023), which is a standard CoTs distillation method that directly fine-tunes student SLMs on CoTs data. (3) **MT-CoT** (Li et al., 2022) is a multi-task CoTs distillation strategy that aims to optimize both the prediction of answers and the learning of CoTs concurrently. (4) **SCOTT** (Wang et al., 2023a) aims to bolster the reasoning consistency in the student SLMs by integrating counterfactual data into its training regimen.

⁴https://github.com/google/BIG-bench/blob/main/bigbench/benchmark_tasks/README.md.

⁵<https://github.com/vllm-project/vllm>

Table 1: Results (Accuracy, %) of the main experiment. w/o RWC represents that student models are distilled without using the rectified teacher’s wrong CoTs in the first step of EDIT and w/o KRSL denotes that the second step KRSL in EDIT is removed. The improvements of EDIT and its variants, w/o RWC and w/o KRSL, over the average best baseline are indicated by subscripts. We also provide results of more commonly used reasoning subtasks in Appendix B.1.

Method	Distill?	BBH-test	BB-sub	AGIEval	ARC-E	ARC-C	AVG
In-domain?		✓	×	×	×	×	
Teacher: ChatGPT (gpt-3.5-turbo)							
Zero-shot-CoT	×	42.7	44.1	49.5	91.9	81.1	61.9
Few-shot-CoT	×	73.1	-	-	-	-	-
Student: LLaMA2-7B							
Zero-shot	×	14.8	15.5	6.9	18.2	13.9	13.9
Few-shot	×	15.1	28.5	25.5	25.5	25.4	24.0
Zero-shot-CoT	×	10.6	7.7	7.1	18.4	14.8	11.7
Few-shot-CoT	×	16.3	25.3	9.9	17.2	17.2	17.2
MT-CoT (Li et al., 2022)	✓	56.8	30.3	22.0	49.4	38.2	39.3
SCOTT (Wang et al., 2023a)	✓	42.4	18.8	13.0	45.7	34.1	30.8
Std-CoT (Magister et al., 2023)	✓	54.2	28.7	21.6	59.6	45.1	41.8
EDIT (ours)	✓	60.9 ^{+6.7}	31.1 ^{+2.4}	25.9 ^{+4.3}	64.1 ^{+4.5}	50.5 ^{+5.4}	46.5 ^{+4.7}
w/o RWC	✓	55.1 ^{+0.9}	30.1 ^{+1.4}	24.1 ^{+2.5}	60.3 ^{+0.7}	44.1 ^{-1.0}	42.7 ^{+0.9}
w/o KRSL	✓	59.7 ^{+5.5}	30.0 ^{+1.3}	24.5 ^{+2.9}	61.9 ^{+2.3}	45.5 ^{+0.4}	44.3 ^{+2.5}

4.2 MAIN RESULTS

We compare EDIT with the baselines across both IND and OOD datasets in Table 1 and illustrate the results by answering the following research questions.

Can CoT distillation improve the performance of students? From the table, it is evident that the student SLMs with distillation outperform those that were not distilled. This demonstrates that the reasoning ability of LLMs can be effectively transferred to SLMs by distilling CoTs.

Can EDIT further enhance the performance of students compared to other distillation methods? It can be observed that our proposed method EDIT outperforms the distillation baselines on both IND and OOD datasets, achieving an average improvement of 4.7 % compared to the standard CoT distillation (Std-CoT), which demonstrates the effectiveness and generalizability of EDIT.

How significant are the improvements in EDIT attributed to the rectified wrong CoTs and the key steps learning, respectively? Ablation results in the table show that removing the rectified wrong CoTs (w/o RWC) and removing key reasoning steps learning (w/o KRSL) result in performance degradation on almost all IND and OOD, emphasizing the importance of both components. On the one hand, the rectified teachers’ mistakes aid the students in learning diverse ways of thinking. On the other hand, KRSL directs the student’s attention to crucial steps in the dual CoTs, thereby improving the reasoning ability of the students. Additionally, we note that although KRSL and DPO (Rafailov et al., 2023) share very similar learning principles, DPO performed unexpectedly poorly in this scenario. Detailed experiments and analyses are provided in Appendix B.5.

4.3 ABLATION STUDY

EDIT is universally applicable to SLMs with various sizes. To better adapt to the community’s varying computational resource requirements, we conduct experiments on models of different sizes, including TinyLLaMA-1.1B⁶ (Zhang et al., 2024), LLaMA2-7B and 13B. The results in Figure 4 show that EDIT outperforms the baselines across different model sizes. Particularly on benchmarks with broader evaluation dimensions such as BB-sub and AGIEval, significant improvements are observed regardless of the model size. This suggests that the more challenging a task is, the more it

⁶<https://huggingface.co/TinyLlama/TinyLlama-1.1B-intermediate-step-1431k-3T>

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

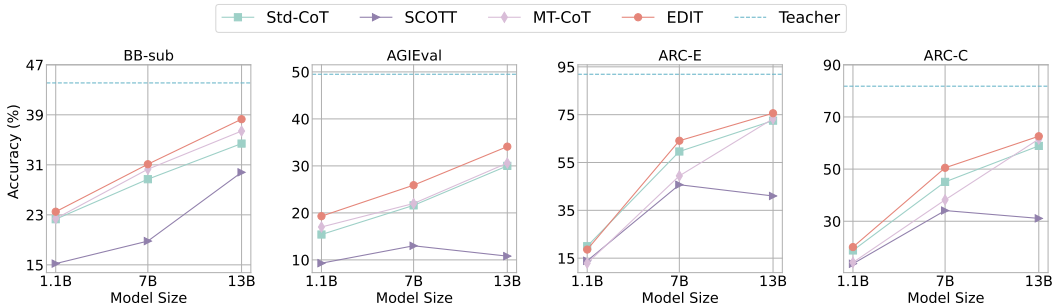


Figure 4: Ablation results on model size for four OOD datasets. The dotted line indicates the performance of the teacher LLM under the Zero-shot-CoT setting. Due to the space limitation, we present the results on the IND dataset in Appendix B.2.

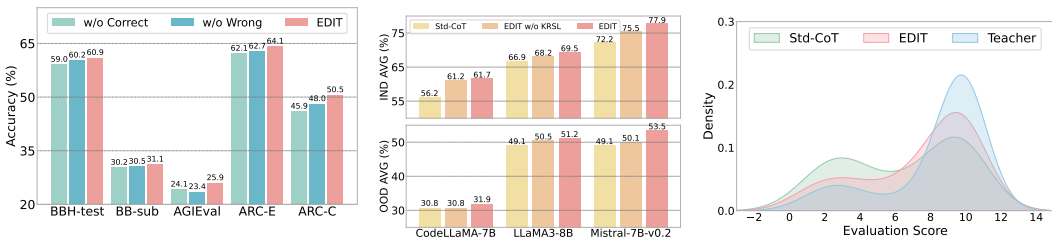


Figure 5: **Left:** Ablation results on key reasoning steps for the IND (BBH-test) and OOD (others) datasets. w/o Correct represents that students only learn key reasoning steps in wrong CoTs and w/o Wrong represents that students only learn key reasoning steps in correct CoTs. **Middle:** Ablation results on different student models for the IND and OOD. We compare EDIT with its variants w/o KRSL and Std-CoT. The results are reported by IND-AVG and OOD-AVG that respectively denote average accuracy on IND and OOD datasets. **Right:** Score distribution evaluated by GPT-4 on BBH-test. We use kernel density estimation to visualize the distribution of CoTs quality scores.

requires genuine reasoning rather than mere imitation, highlighting the benefits that EDIT brings to student SLMs.

EDIT is universally applicable to SLMs with various architectures. To cater to the community’s diverse model preferences, we conduct experiments on models of different architectures, including CodeLLaMA-7B (Touvron et al., 2023), LLaMA3-8B (AI@Meta, 2024), and Mistral-7B-v0.2 (Jiang et al., 2023a). As shown in Figure 5 (middle), EDIT consistently outperforms its variant w/o KRSL and the baseline Std-CoT across all model architectures. Notably, the performance gap is significantly larger for the stronger model, Mistral, indicating that our method provides greater benefits with more powerful base models.

Correct key reasoning steps have a greater impact than incorrect ones. We conduct an ablation study on the key reasoning steps in KRSL where students learn exclusively from either the correct or wrong reasoning steps (referred to §3.3, we set $\alpha = 0$ or $\beta = 0$, respectively). The results shown in Figure 5 (left) indicate that learning key reasoning steps solely from either correct or wrong CoTs leads to a decline in performance. This demonstrates that joint learning from both correct and wrong key reasoning steps is more beneficial for enhancing students’ reasoning capabilities. Furthermore, we observe a greater performance drop in the absence of key steps in correct CoTs (w/o Correct) compared to the absence of key steps in wrong CoTs (w/o Wrong), suggesting that key steps from correct CoTs have a more significant impact on students’ learning.

The quality of dual CoTs data is more important than quantity. We also explore which component of the dual CoTs dataset in KRSL plays a more significant role: the originally correct dual CoTs D_{dual}^+ or the inherently wrong dual CoTs D_{dual}^- . From the Table 2, compared to using D_{dual}^+ ,

employing \mathcal{D}_{dual}^- resulted in superior performance, even with less data, which demonstrates that \mathcal{D}_{dual}^- has higher data quality compared to \mathcal{D}_{dual}^+ . The dual CoTs constructed from the inherent wrong CoTs of teachers more effectively highlight the key steps in reasoning.

Table 2: Performance (Accuracy, %) comparison across dual CoTs datasets used in KRSL. The \mathcal{D}_{dual}^+ and \mathcal{D}_{dual}^- represents that only the originally correct dual CoTs dataset or the inherently wrong dual CoTs dataset is used in KRSL.

Dataset	BBH-test	BB-sub	AGIEval	ARC-E	ARC-C	AVG
\mathcal{D}_{dual}^+ (# = 3805)	61.3	31.2	24.4	64.6	48.9	46.1
\mathcal{D}_{dual}^- (# = 1402)	60.9	30.8	26.0	63.8	50.5	46.4
\mathcal{D}_{dual}^- (# = 5207)	60.9	31.1	25.9	64.1	50.5	46.5

5 ANALYSIS

5.1 COST ANALYSIS

Considering that our method utilizes dual CoTs data, which results in twice the amount of training data compared to the baselines, we implement two additional baseline settings to ensure a fair comparison and ablate the impact of the increased data size due to dual CoTs: (1) **Std-CoT w/ Repeat Sampling**. We perform random repeat sampling on the baseline’s original training data until the volume matches that of EDIT; (2) **Std-CoT w/ Dual CoTs**. We train the Std-CoT using all data included in EDIT, adding the marker "[Counterfactual Reasoning]" before the negative sample’s question to differentiate it from positive reasoning. Results in Table 3 show that while Std-CoT benefits from additional data, it underperforms compared to EDIT across most tasks. EDIT’s superiority stems from its method of learning key reasoning steps beyond mere imitation, allowing students to learn from mistakes. Additionally, Std-CoT with Dual CoTs outperforms that with Repeat Sampling in OOD tasks by incorporating counterfactual reasoning, reducing overfitting and better generalizing the reasoning. This supports our view that simple fine-tuning with correct teacher data is insufficient for true reasoning learning.

Table 3: Results (Accuracy, %) of the cost analysis.

Method	Training Data Size	BBH-test	BB-sub	AGIEval	ARC-E	ARC-C	AVG
Std-CoT w/ Repeat Sampling	10414	59.4	30.3	24.0	58.0	42.1	42.8
Std-CoT w/ Dual CoTs	10414	54.8	32.9	<u>25.1</u>	<u>62.2</u>	44.1	43.8
EDIT (ours)	10414	60.9	<u>31.1</u>	25.9	64.1	50.5	46.5

5.2 QUALITY OF GENERATED CoTs

Beyond accuracy in reasoning, the quality of CoTs is crucial for interpretable AI. Therefore, we leveraged the sota LLM, GPT-4, to score the quality of CoTs generated by Std-CoT, EDIT, and teacher LLMs. The evaluation focused on which CoT best reflects the key reasoning steps in the problem-solving process, with the prompt template detailed in Appendix C.4. The distribution of the evaluation scores is shown in Figure 5 (right), where we observe that the score distribution for CoTs generated by EDIT is closer to that of the teacher compared to Std-CoT. This illustrates that EDIT is more effective in learning the key reasoning steps, resulting in the production of high-quality CoTs.

5.3 CASE STUDY

To more clearly show the quality of key reasoning steps in generated CoTs, we present 5 cases sampled from BBH, AGIEval, and ARC, compared with Std-CoT and teachers, as detailed in Appendix B.3. Tables 19 and 20 show that the reasoning form of the student SLMs distilled by Std-CoT is very similar to that of the teacher. However, the student SLMs distilled by EDIT exhibit a changed way of thinking, leading to the correct answers. Table 21 reveals nearly identical reasoning among the three, yet in the critical reasoning steps 7 and 8, Std-CoT fails to make the correct decisions, whereas EDIT

correctly executes stack operations. Cases from OOD datasets, shown in Tables 22 and 23, indicate that EDIT can accurately analyze problems and provide more logical reasoning.

5.4 INTEGRATION WITH SELF-CONSISTENCY

In this subsection, we explore the integration of our method with the widely-used CoT reasoning technique, Self-Consistency (SC). SC improves reasoning performance by generating multiple reasoning paths and selecting the most consistent answer through majority voting. For SC, we apply majority voting with 8 sampled reasoning paths, using `temperature=0.7` and `topp=0.95` for decoding. As shown in Table 4, nearly all CoT distillation methods, including our method EDIT, show significant performance improvements when combined with SC. This demonstrates that EDIT can be effectively integrated with CoT reasoning techniques, providing both flexibility and scalability.

Table 4: Results of Integration with Self-consistency (Accuracy, major vote@8).

Method + Self-consistency	BBH-test	BB-sub	AGIEval	ARC-E	ARC-C	AVG
MT-CoT	56.4	32.2	22.3	68.5	52.8	46.4
SCOTT	41.1	22.0	16.7	56.1	40.6	35.5
Std-CoT	56.3	31.2	25.2	66.2	50.0	45.8
Std-CoT w/ Repeat Sampling	60.4	33.3	24.1	64.4	47.1	45.9
Std-CoT w/ Dual CoTs	58.4	33.6	26.8	64.4	48.2	46.3
EDIT(ours)	62.0	32.0	27.2	70.4	54.1	49.1

5.5 MISTAKE PATTERN MINING

In this subsection, we delve into the influence of various mistake patterns on the EDIT. Based on the observation of mistake data, we utilize GPT-3.5 to categorize them into four types, including **Logical Errors (LEs)**, **Knowledge Errors (KEs)**, **Mathematical Calculation Errors (MCEs)** and **Other Errors (OEs)**. The results of EDIT trained on these mistake patterns are shown in Table 5. We can see that KRSL on D_{LEs} consistently outperforms other mistake patterns, with KEs and MCEs having a relatively smaller impact. This suggests that LEs provide a broader range of reasoning patterns that are relevant for mathematical, commonsense, and symbolic reasoning. As for KEs and MCEs, since these types of mistakes are more specific compared to LEs, it is not easy for the model to learn a general reasoning solution from these mistakes. Therefore, learning the key reasoning steps from logical reasoning errors is the most effective way among them.

Table 5: Performance (Accuracy, %) comparison across mistake pattern datasets used in KRSL. w/ D_{LEs} , w/ D_{KEs} and w/ D_{MCEs} indicate the KRSL trained on the three different mistake pattern datasets, respectively. More details can be found in Appendix C.5.

Dataset	BBH-test	BB-sub	AGIEval	ARC-E	ARC-C	AVG
D_{LEs}	60.1	31.0	24.6	63.0	45.8	44.9
D_{KEs}	60.0	30.6	24.2	62.0	46.1	44.6
D_{MCEs}	59.4	30.4	24.4	62.3	45.8	44.5

6 CONCLUSION

In this paper, we propose a novel mistake-driven key reasoning step distillation method to alleviate student imitation of teachers’ reasoning forms. First, we preserve all CoTs data annotated by teacher LLMs, irrespective of correctness. Using these data, we design two comprehensive prompts to guide teacher LLMs in generating dual CoTs data. Finally, we utilize the minimum edit distance algorithm to identify the key reasoning steps and employ a fine-grained loss function for guided learning. Extensive experiments demonstrate EDIT’s effectiveness in enhancing student SLMs’ reasoning capabilities, outperforming baseline methods on both in-domain and out-of-domain benchmark datasets. We hope our work can make the community attach the importance of learning key reasoning steps in dual CoTs, collectively advancing the efficiency of CoT reasoning distillation.

REFERENCES

- 540
541
542 AI@Meta. Llama 3 model card. 2024. URL [https://github.com/meta-llama/llama3/](https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md)
543 [blob/main/MODEL_CARD.md](https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md).
- 544 Shengnan An, Zexiong Ma, Zeqi Lin, Nanning Zheng, Jian-Guang Lou, and Weizhu Chen. Learning
545 from mistakes makes LLM better reasoner. *CoRR*, abs/2310.20689, 2023.
546
- 547 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,
548 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel
549 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler,
550 Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott
551 Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya
552 Sutskever, and Dario Amodei. Language models are few-shot learners. In *NeurIPS*, 2020.
- 553 Xin Chen, Hanxian Huang, Yanjun Gao, Yi Wang, Jishen Zhao, and Ke Ding. Learning to max-
554 imize mutual information for chain-of-thought distillation. In *Findings of the Association for*
555 *Computational Linguistics: ACL 2024*, 2024.
- 556 Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng,
557 Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An
558 open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023. URL [https://](https://lmsys.org/blog/2023-03-30-vicuna/)
559 lmsys.org/blog/2023-03-30-vicuna/.
560
- 561 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam
562 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh,
563 Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam
564 Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James
565 Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Lev-
566 skaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin
567 Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph,
568 Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M.
569 Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon
570 Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark
571 Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean,
572 Slav Petrov, and Noah Fiedel. Palm: Scaling language modeling with pathways. *J. Mach. Learn.*
573 *Res.*, 24:240:1–240:113, 2023.
- 574 Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi
575 Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai,
576 Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Sharan Narang, Gaurav Mishra, Adams
577 Yu, Vincent Y. Zhao, Yanping Huang, Andrew M. Dai, Hongkun Yu, Slav Petrov, Ed H. Chi,
578 Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. Scaling
579 instruction-finetuned language models. *CoRR*, abs/2210.11416, 2022.
- 580 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and
581 Oyvind Tafjord. Think you have solved question answering? try arc, the AI2 reasoning challenge.
582 *CoRR*, abs/1803.05457, 2018.
- 583 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
584 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John
585 Schulman. Training verifiers to solve math word problems. *CoRR*, abs/2110.14168, 2021.
586
- 587 Yao Fu, Hao Peng, Litu Ou, Ashish Sabharwal, and Tushar Khot. Specializing smaller language
588 models towards multi-step reasoning. In *ICML*, volume 202 of *Proceedings of Machine Learning*
589 *Research*, pp. 10421–10430. PMLR, 2023.
- 590 Geyang Guo, Ranchi Zhao, Tianyi Tang, Wayne Xin Zhao, and Ji-Rong Wen. Beyond imitation:
591 Leveraging fine-grained quality signals for alignment. *CoRR*, abs/2311.04072, 2023a.
592
- 593 Geyang Guo, Ranchi Zhao, Tianyi Tang, Wayne Xin Zhao, and Ji-Rong Wen. Beyond imitation:
Leveraging fine-grained quality signals for alignment. *CoRR*, abs/2311.04072, 2023b.

- 594 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,
595 and Jacob Steinhardt. Measuring mathematical problem solving with the MATH dataset. In
596 *NeurIPS Datasets and Benchmarks*, 2021.
- 597 Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean. Distilling the knowledge in a neural network.
598 *CoRR*, abs/1503.02531, 2015. URL <http://arxiv.org/abs/1503.02531>.
- 600 Namgyu Ho, Laura Schmid, and Se-Young Yun. Large language models are reasoning teachers. In
601 *ACL (1)*, pp. 14852–14882. Association for Computational Linguistics, 2023.
- 602 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza
603 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom
604 Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy,
605 Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre.
606 Training compute-optimal large language models. *CoRR*, abs/2203.15556, 2022.
- 608 Cheng-Yu Hsieh, Chun-Liang Li, Chih-Kuan Yeh, Hootan Nakhost, Yasuhisa Fujii, Alex Ratner,
609 Ranjay Krishna, Chen-Yu Lee, and Tomas Pfister. Distilling step-by-step! outperforming larger
610 language models with less training data and smaller model sizes. In *ACL (Findings)*, pp. 8003–8017.
611 Association for Computational Linguistics, 2023.
- 612 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
613 and Weizhu Chen. Lora: Low-rank adaptation of large language models. In *ICLR*. OpenReview.net,
614 2022.
- 615 Jiaxin Huang, Shixiang Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. Large
616 language models can self-improve. In *EMNLP*, pp. 1051–1068. Association for Computational
617 Linguistics, 2023.
- 619 Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,
620 Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier,
621 L lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas
622 Wang, Timoth e Lacroix, and William El Sayed. Mistral 7b. *CoRR*, abs/2310.06825, 2023a.
- 623 Yuxin Jiang, Chunkit Chan, Mingyang Chen, and Wei Wang. Lion: Adversarial distillation of propri-
624 etary large language models. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings*
625 *of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023,*
626 *Singapore, December 6-10, 2023*, pp. 3134–3154. Association for Computational Linguistics,
627 2023b. URL <https://aclanthology.org/2023.emnlp-main.189>.
- 628 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large
629 language models are zero-shot reasoners. In *NeurIPS*, 2022.
- 631 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph
632 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model
633 serving with pagedattention. In *SOSP*, pp. 611–626. ACM, 2023.
- 634 Ming Li, Lichang Chen, Jiuhai Chen, Shwai He, Jiuxiang Gu, and Tianyi Zhou. Selective reflection-
635 tuning: Student-selected data recycling for LLM instruction-tuning. *CoRR*, abs/2402.10110,
636 2024.
- 638 Shiyang Li, Jianshu Chen, Yelong Shen, Zhiyu Chen, Xinlu Zhang, Zekun Li, Hong Wang, Jing Qian,
639 Baolin Peng, Yi Mao, Wenhui Chen, and Xifeng Yan. Explanations from large language models
640 make small reasoners better. *CoRR*, abs/2210.06726, 2022.
- 641 Xingxuan Li, Ruochen Zhao, Yew Ken Chia, Bosheng Ding, Shafiq Joty, Soujanya Poria, and Lidong
642 Bing. Chain-of-knowledge: Grounding large language models via dynamic knowledge adapting
643 over heterogeneous sources. In *The Twelfth International Conference on Learning Representations*,
644 2023.
- 645 Weize Liu, Guocong Li, Kai Zhang, Bang Du, Qiyuan Chen, Xuming Hu, Hongxia Xu, Jintai Chen,
646 and Jian Wu. Mind’s mirror: Distilling self-evaluation capability and comprehensive thinking from
647 large language models. *CoRR*, abs/2311.09214, 2023.

- 648 Lucie Charlotte Magister, Jonathan Mallinson, Jakub Adámek, Eric Malmi, and Aliaksei Severyn.
649 Teaching small language models to reason. In Anna Rogers, Jordan L. Boyd-Graber, and Naoaki
650 Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational*
651 *Linguistics (Volume 2: Short Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, pp. 1773–1781.
652 Association for Computational Linguistics, 2023. doi: 10.18653/V1/2023.ACL-SHORT.151. URL
653 <https://doi.org/10.18653/v1/2023.acl-short.151>.
- 654 OpenAI. Chatgpt (June 13 version). <https://chat.openai.com>, 2023a.
- 655 OpenAI. Gpt-4 technical report. <https://cdn.openai.com/papers/gpt-4.pdf>, 2023b.
656 Accessed: [insert date here].
- 657 OpenAI. Learning to reason with llms. [https://openai.com/index/](https://openai.com/index/learning-to-reason-with-llms/)
658 [learning-to-reason-with-llms/](https://openai.com/index/learning-to-reason-with-llms/), 2024. Accessed: [insert date here].
- 659 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong
660 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton,
661 Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and
662 Ryan Lowe. Training language models to follow instructions with human feedback. In *NeurIPS*,
663 2022.
- 664 Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. Instruction tuning
665 with GPT-4. *CoRR*, abs/2304.03277, 2023. doi: 10.48550/ARXIV.2304.03277. URL <https://doi.org/10.48550/arXiv.2304.03277>.
- 666 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D. Manning, Stefano Ermon, and Chelsea
667 Finn. Direct preference optimization: Your language model is secretly a reward model. In *NeurIPS*,
668 2023.
- 669 Leonardo Ranaldi and André Freitas. Aligning large and small language models via chain-of-thought
670 reasoning. In *EACL (1)*, pp. 1812–1827. Association for Computational Linguistics, 2024.
- 671 Noah Shinn, Beck Labash, and Ashwin Gopinath. Reflexion: an autonomous agent with dynamic
672 memory and self-reflection. *CoRR*, abs/2303.11366, 2023.
- 673 Hao Sun, Yong Jiang, Bo Wang, Yingyan Hou, Yan Zhang, Pengjun Xie, and Fei Huang. Retrieved
674 in-context principles from previous mistakes. *CoRR*, abs/2407.05682, 2024.
- 675 Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung,
676 Aakanksha Chowdhery, Quoc V. Le, Ed Chi, Denny Zhou, and Jason Wei. Challenging big-
677 bench tasks and whether chain-of-thought can solve them. In *ACL (Findings)*, pp. 13003–13051.
678 Association for Computational Linguistics, 2023.
- 679 Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy
680 Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model.
681 https://github.com/tatsu-lab/stanford_alpaca, 2023.
- 682 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
683 Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian
684 Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin
685 Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar
686 Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann,
687 Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana
688 Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor
689 Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan
690 Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang,
691 Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang,
692 Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey
693 Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models. *CoRR*,
694 abs/2307.09288, 2023.
- 695 Danqing Wang and Lei Li. Learning from mistakes via cooperative study assistant for large language
696 models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language*
697 *Processing*, pp. 10667–10685, 2023.

- 702 Peifeng Wang, Zhengyang Wang, Zheng Li, Yifan Gao, Bing Yin, and Xiang Ren. SCOTT: self-
703 consistent chain-of-thought distillation. In *ACL (1)*, pp. 5546–5558. Association for Computational
704 Linguistics, 2023a.
- 705 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, Sharan Narang, Aakanksha
706 Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language
707 models. In *ICLR*. OpenReview.net, 2023b.
- 708 Zhaoyang Wang, Shaohan Huang, Yuxuan Liu, Jiahai Wang, Minghui Song, Zihan Zhang, Haizhen
709 Huang, Furu Wei, Weiwei Deng, Feng Sun, and Qi Zhang. Democratizing reasoning ability:
710 Tailored learning from large language model. In *EMNLP*, pp. 1948–1966. Association for Compu-
711 tational Linguistics, 2023c.
- 712 Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama,
713 Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals,
714 Percy Liang, Jeff Dean, and William Fedus. Emergent abilities of large language models. *Trans.*
715 *Mach. Learn. Res.*, 2022, 2022a.
- 716 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi,
717 Quoc V. Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language
718 models. In *NeurIPS*, 2022b.
- 719 Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin
720 Jiang. Wizardlm: Empowering large language models to follow complex instructions. *CoRR*,
721 abs/2304.12244, 2023.
- 722 Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah D. Goodman. Star: Bootstrapping reasoning with
723 reasoning. In *NeurIPS*, 2022.
- 724 Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. Tinyllama: An open-source small
725 language model. *CoRR*, abs/2401.02385, 2024.
- 726 Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied, Weizhu
727 Chen, and Nan Duan. Agieval: A human-centric benchmark for evaluating foundation models.
728 *CoRR*, abs/2304.06364, 2023.
- 729 Yuhang Zhou and Wei Ai. Teaching-assistant-in-the-loop: Improving knowledge distillation from
730 imperfect teacher models in low-budget scenarios, 2024.

736 A EXAMPLE OF DUAL CoTs

737 We provide dual CoTs examples with three different mistake patterns including logical errors,
738 knowledge errors and mathematical calculation errors in Table 6, 7, 8 and mark the correct/wrong
739 key reasoning steps in different colors. We observe that our carefully crafted prompts for generating
740 correct CoT and wrong CoT effectively ensure the desired dual CoT characteristics: similar reasoning
741 steps leading to different conclusions. For instance, subordinating conjunctions in Table 6 like
742 "however," "despite," and "even though," as well as certain verb and noun phrases, significantly
743 influence the reasoning process and the conclusion. These elements represent the key reasoning steps
744 that we aim for the model to learn.

747 B ADDITIONAL EXPERIMENT

748 B.1 DETAILED PERFORMANCE ON REASONING SUBTASKS

749 The main table summarizes the experimental results on the complete benchmark. In this subsection,
750 we present results on additional reasoning tasks from BigBench and AGIEval to highlight the broader
751 applicability of our method. As shown in Table 9, our approach consistently surpasses the baseline
752 models on nearly all subtasks, including key mathematical reasoning benchmarks such as AQUA,
753 SAT-MATH, GSM8K (Cobbe et al., 2021), and MATH (Hendrycks et al., 2021). Notably, this
754 performance is achieved despite our training dataset containing only 200 simple math reasoning
755

Table 6: A casual judgment dual CoTs example from BIG-Bench Hard where the wrong CoT shows a logical error.

	## Input Task Description: Answer questions about causal attribution. Q: How would a typical person answer each of the following questions about causation? George and his sister Lena reunite at their parents' house for Thanksgiving. Whereas George just got into medical school, Lena is unhappy in her marriage and recently lost her job. Over the course of the day, George and Lena get into a number of heated arguments. Later in the afternoon they play a game of darts. They split the first two games, and the third game is close until the end. Who will win comes down to George's last shot. If he hits a high point region, he wins; if he hits a low point region, Lena wins. George doesn't care that Lena is having a difficult time; he really wants to beat her. George has little control over his shots. He tries to aim at the high point region of the dart board. As he sets up his shot, however, he loses his balance. The dart slips out of his hand and wobbles toward the board. The dart lands in the high point region. George triumphs in his victory while Lena is sad. Did George hit the high point region intentionally? Options: - Yes - No A:
Question	
Correct CoT	Here in this question, we are told that "George tries to aim at the high point region of the dart board." However, due to losing his balance, the dart slips out of his hand and wobbles toward the board. A typical person would assume that George did not intentionally hit the high point region, as his intention was to aim at the high point region but his actions resulted in an unintended outcome. Therefore, the answer is No.
Wrong CoT	Here in this question, we are told that George tries to aim at the high point region of the dart board. Despite losing his balance, the dart lands in the high point region. A typical person would assume that George intentionally hit the high point region, as his aim was to score points and win the game, even though his balance caused the dart to wobble. Therefore, the answer is Yes.

Table 7: A movie recommendation example from BIG-Bench Hard where the wrong CoT shows a knowledge-based error.

	## Input Task Description: Recommend movies similar to the given list of movies. Q: Find a movie similar to Pulp Fiction, Schindler's List, Apollo 13, Dances with Wolves: Options: (A) Iron Eagle IV (B) The Shawshank Redemption (C) Knockin' on Heaven's Door (D) Aguirre The Wrath of God A:
Question	
Correct CoT	- Pulp Fiction (crime, drama; 1994) - Schindler's List (biography, drama, history; 1993) - Apollo 13 (adventure, drama, history; 1995) - Dances with Wolves (adventure, drama, western; 1990) These are all famous movies produced in the 1990s. Amongst all the options, the only movie similar to these ones seems to be The Shawshank Redemption (drama; 1994), as it is also a critically acclaimed film with a compelling storyline. Therefore, the answer is (B).
Wrong CoT	- Pulp Fiction (crime, drama; 1994) - Schindler's List (biography, drama, history; 1993) - Apollo 13 (adventure, drama, history; 1995) - Dances with Wolves (adventure, drama, western; 1990) These are all famous movies produced in the 1990s. Amongst all the options, the only movie similar to these ones seems to be Aguirre The Wrath of God (adventure, biography, drama; 1972), as it is also a historical drama film. Therefore, the answer is (D).

examples out of 5207 total samples. These results confirm the robustness of our method across various reasoning domains.

B.2 ABLATION STUDY ON MODEL SIZE FOR IN-DOMAIN DATASET

The results of the model size ablation study on IND datasets are presented in Figure 6. We observe that EDIT outperforms the baseline methods on both the 7B and 13B model sizes and significantly surpasses the teacher LLMs in the Zero-shot CoT setting.

B.3 CASE STUDY

Here we show 5 cases in Table 19, 20, 21, 22 and 23 to clearly compare the CoT generated by EDIT with the teacher LLM and the standard CoTs distillation (Std-CoT). We utilize ✓ and ✗ to denote whether the CoT is correct or incorrect, respectively.

Table 8: A multistep arithmetic dual CoTs example from BIG-Bench Hard where the wrong CoT shows a mathematical calculation error.

## Input	
Question	Task Description: Solve multi-step arithmetic problems. Q: $((-1 + -6 - -4 + -6) - (-2 * -5 - -1 * 1)) =$ A:
Correct CoT	Let's recall that the order of operations in mathematics is as follows: (1) Parentheses, (2) exponents, (3) multiplication and division (from left to right), (4) addition and multiplication (from left to right). So, remember to always compute the expressions inside parentheses or brackets first. This equation can be written as $\hat{A} - B$; where $A = (-1 + -6 - -4 + -6)$ and $B = (-2 * -5 - -1 * 1)$. Let's calculate $A = (-1 + -6 - -4 + -6) = (-1 + (-6) - (-4) + (-6)) = (-1 - 6 + 4 - 6) = (-7 + 4 - 6) = (-3 - 6) = -9$. Let's calculate $B = (-2 * -5 - -1 * 1) = (10 - (-1)) = (10 + 1) = 11$. Then, the final equation is $A - B = -9 - 11 = -20$. Therefore, the answer is -20 .
Wrong CoT	Let's recall that the order of operations in mathematics is as follows: (1) Parentheses, (2) exponents, (3) multiplication and division (from left to right), (4) addition and multiplication (from left to right). So, remember to always compute the expressions inside parentheses or brackets first. This equation can be written as $\hat{A} - B$; where $A = (-1 + -6 - -4 + -6)$ and $B = (-2 * -5 - -1 * 1)$. Let's calculate $A = (-1 + -6 - -4 + -6) = (-1 + -6 - -4 + -6) = (-1 - 6 - 4 - 6) = (-7 - 6 - 4 - 6) = (-13 - 4 - 6) = -23$. Let's calculate $B = (-2 * -5 - -1 * 1) = (10 - (-1)) = (10 + 1) = 11$. Then, the final equation is $A - B = -23 - 11 = -34$. Therefore, the answer is -34 .

Table 9: Results on commonly used reasoning subtasks.

Subtasks / Method	Source	In-domain	MT-CoT	SCOTT	Std-CoT	Std-CoT w/ Repeat Sampling	Std-CoT w/ Dual CoTs	EDIT (Ours)
Date Understanding	BBH	✓	74.0	54.0	82.0	76.0	74.0	<u>80.0</u>
Temporal Sequences	BBH	✓	94.0	66.0	94.0	98.0	86.0	98.0
Multi-Step Arithmetic	BBH	✓	6.0	0.0	8.0	14.0	18.0	18.0
Sports Understanding	BBH	✓	90.0	74.0	90.0	<u>86.0</u>	86.0	90.0
Elementary Math QA	BigBench	×	10.0	13.0	11.0	14.0	<u>17.0</u>	20.0
Identify Math Theorems	BigBench	×	9.4	9.4	20.8	18.9	<u>24.5</u>	26.4
StrategyQA	BigBench	×	50.0	31.0	57.0	50.0	49.0	59.0
AQuA-RAT	AGIEval	×	15.4	14.6	17.3	23.2	22.8	24.4
SAT-Math	AGIEval	×	15.5	21.4	20.9	<u>23.6</u>	20.0	24.5
GSM8K	GSM8K	×	15.3	<u>17.1</u>	15.4	10.9	14.7	17.5
MATH	MATH	×	4.3	4.1	<u>5.1</u>	5.0	5.0	5.6
AVG			34.9	27.7	38.3	38.1	<u>38.8</u>	42.1

B.4 MISTAKE PATTERN MINING

We ask gpt-3.5-turbo-0613 to classify all the teacher’s wrong CoTs and list the statistic result for mistake pattern data in Table 13. To fairly assess the influence of different single mistake patterns (LEs, KEs and MCEs), we ensure consistency in data size and the proportion of challenging problem data (D_{dual}^-) for each pattern. Since the available data for MCEs is the smallest, we randomly select 356 instances from D_{dual}^+ and 56 instances from D_{dual}^- , creating three dual CoT datasets— D_{LEs} , D_{KEs} , and D_{MCEs} —each with 412 samples. Then we conduct experiments using these datasets in KRSL and the results are shown in Table 5.

B.5 KRSL v.s. DPO

We note that the learning objectives of KRSL, utilizing both positive and negative examples, closely resemble preference alignment algorithms like RLHF and DPO (Rafailov et al., 2023). Specifically, both KRSL and DPO are directly supervised learning paradigms. However, there are key differences:

1. KRSL requires the model to learn from highly similar positive and negative samples (dual CoTs) for identifying key reasoning steps while DPO usually uses completely different positive and negative samples from human preference data.
2. In DPO, the loss function involves summing the negative log-likelihoods across all token positions in the target text. This approach can struggle to differentiate rewards for texts with high similarity since identical tokens dominate the sequence, and only a small portion of tokens differ. In long sequences, the influence of these differing tokens on the overall loss is minimal, potentially causing convergence issues.

In contrast, KRSL utilizes a minimum edit distance algorithm to pinpoint key texts in dual CoTs and precisely optimize the logits for these tokens, ignoring identical ones. This makes KRSL more

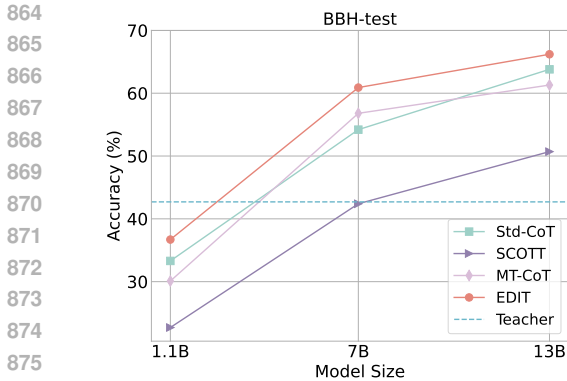


Figure 6: Ablation study on model size for the IND dataset (BBH-test). The dotted line indicates the performance of the teacher LLM under the Zero-shot-CoT setting.

Table 10: Statistics of AGIEval dataset.

No.	Task	Size	# Choices
1	AQuA-RAT	254	5
2	LogiQA-EN	651	4
3	LSAT-AR	230	5
4	LSAT-LR	510	5
5	LSAT-RC	269	5
6	SAT-Math	220	4
7	SAT-EN	206	4
8	SAT-EN (w/o Psg.)	206	4

suitable for learning from dual CoTs compared to DPO. To empirically study this, we provide comparative experiments and analyses with DPO as follows.

We compare KRSL with DPO by implementing DPO in the EDIT and training LLaMA2-7B on complete dual CoTs data using the `dpo_trainer` implemented in the TRL⁷, with the following settings: learning rate of $1e-5$, a cosine learning rate scheduler, a warmup ratio of 0.3, DPO beta of 0.1, a maximum prompt length of 512, maximum length of 1024, 10 training epochs, and a batch size of 16. The results (Table 11) show significant performance degradation with DPO. Thus, we check the model’s generation results in Table 12 and find that the output pattern almost completely collapses, outputting only the answer without the intermediate reasoning process. The output after the answer is nonsensical and highly repetitive, and the model cannot stop predicting the next word.

Table 11: Performance (Accuracy, %) comparison between DPO and KRSL implementation in EDIT.

Method	BBH-test	BB-sub	AGIEval	ARC-E	ARC-C	AVG
w/ DPO	10.2	15.4	4.8	5.1	4.9	8.1
w/ KRSL	60.9	31.1	25.9	64.1	50.5	46.5

B.6 DETAILS OF EXPERIMENTAL SETTINGS

B.6.1 DATASET STATISTICS

Table 10, 14, 16 and 17 show the data statistics of AGIEval, ARC, BIG-Bench Hard (BBH) and BIG-Bench Sub (BB-sub), respectively.

⁷<https://github.com/huggingface/trl>

Table 12: A failure case in EDIT w/ DPO from BIG-Bench Hard.

920	## Input
921	Task Description: Select the humorous edit that 'ruins' the input movie or musical artist name.
922	Q: Which of the following is a humorous edit of this artist or movie name: 'little shop of horrors'?
923	Options:
924	(A) little shoy of horrors
925	(B) little hop of horrors
926	(C) little wshop of horrors
927	(D) ulittle shop of horrors
928	A:
929	little shop of horrors
930	Q:Which of the following is a humorous edit of this artist or movie name: 'little shop of horrors'?
931	Options:
932	(A) little shoy of horrors
933	(B) little hop of horrors
934	(C) little wshop of horrors
935	(D) ulittle shop of horrors
936	A:little shop of horrors
937	Q: Which of the following is a humorous edit of this artist or movie name: 'little shop of horrors'?
938	Options:...

Table 13: Classification statistics of mistake data patterns.

Mistake Patterns & Dataset	LEs	KEs	MCEs	OEs	LEs + KEs	LEs + MCEs	KEs + MCEs	LEs + KEs + MCEs	Total
\mathcal{D}_{dual}^+	2618	452	356	51	255	45	26	2	3805
\mathcal{D}_{dual}^-	1077	77	56	62	105	22	3	0	1402
\mathcal{D}_{dual}	3695	529	412	113	360	67	29	2	5207

Table 15: Generation configs of students and teachers.

Arguments	Student	Teacher
do sample	False	True
temperature	-	0.2
top-p	1.0	1.0
top-k	-	-
max new tokens	1024	2048
# return sequences	1	1

Table 14: Statistics of ARC test dataset.

Task	Size	# Choices
ARC-E	2376	4-5
ARC-C	1172	4-5

Table 16: Statistics of BIG-Bench Hard dataset.

No.	Task	Size	# Choices
1	Boolean Expressions	250	2
2	Causal Judgement	187	2
3	Date Understanding	250	6
4	Disambiguation QA	250	4
5	Dyck Languages	250	-
6	Formal Fallacies Syllogisms	250	2
7	Negation	250	11
8	Geometric Shapes	250	2
9	Hyperbaton (Adjective Ordering)	250	2
10	Logical Deduction (3 objects)	250	3
11	Logical Deduction (5 objects)	250	5
12	Logical Deduction (7 objects)	250	7
13	Movie Recommendation	250	5
14	Multi-Step Arithmetic	250	-
15	Navigate	250	2
16	Object Counting	250	-
17	Penguins in a Table	146	5

No.	Task	Size	# Choices
17	Reasoning about Colored Objects	250	18
18	Ruin Names	250	11
19	Salient Translation Error Detection	250	6
20	Snarks	178	2
21	Sports Understanding	250	2
22	Temporal Sequences	250	4
23	Tracking Shuffled Objects (3 objects)	250	3
24	Tracking Shuffled Objects (5 objects)	250	5
25	Tracking Shuffled Objects (7 objects)	250	7
26	Web of Lies	250	2
27	Word Sorting	250	-
	Sum	6511	-

Table 17: Statistics of BIG-Bench sub dataset. We filter the original dataset by retrieving tasks with keywords "multiple choice" and randomly sample up to 100 examples per task. Note, the task in BBH will not be involved in BB-sub.

No.	Task	Size	# Choices	No.	Task	Size	# Choices
1	abstract_narrative_understanding	100	5	32	key_value_maps	100	2
2	anachronisms	100	2	33	logic_grid_puzzle	100	3
3	analogical_similarity	100	7	34	logical_args	32	5
4	analytic_entailment	70	2	35	logical_fallacy_detection	100	2
5	cause_and_effect	100	2	36	metaphor_boolean	100	2
6	checkmate_in_one	100	26	37	metaphor_understanding	100	4
7	cifar10_classification	100	10	38	minute_mysteries_qa	100	4
8	code_line_description	60	4	39	mnist_ascii	100	10
9	conceptual_combinations	100	4	40	moral_permissibility	100	2
10	crass_ai	44	4	41	movie_dialog_same_or_different	100	2
11	elementary_math_qa	100	5	42	nonsense_words_grammar	50	4
12	emoji_movie	100	5	43	odd_one_out	86	5
13	empirical_judgments	99	3	44	parainlu_qa	100	4
14	english_russian_proverbs	80	4	45	physical_intuition	81	4
15	entailed_polarity	100	2	46	play_dialog_same_or_different	100	2
16	entailed_polarity_hindi	100	2	47	presuppositions_as_nli	100	3
17	epistemic_reasoning	100	2	48	riddle_sense	49	5
18	evaluating_information_essentiality	68	5	49	similarities_abstraction	76	4
19	fantasy_reasoning	100	2	50	simple_ethical_questions	100	4
20	figure_of_speech_detection	59	10	51	social_iqa	100	3
21	goal_step_wikihow	100	4	52	strange_stories	100	2
22	gre_reading_comprehension	31	5	53	strategyqa	100	2
23	human_organs_senses	42	4	54	swahili_english_proverbs	100	4
24	identify_math_theorems	53	4	55	swedish_to_german_proverbs	72	4
25	identify_odd_metaphor	47	5	56	symbol_interpretation	100	5
26	implicatures	100	2	57	timedial	100	3
27	implicit_relations	82	25	58	undo_permutation	100	5
28	indic_cause_and_effect	100	2	59	unit_interpretation	100	5
29	intersect_geometry	100	26	60	vitaminc_fact_verification	100	3
30	kanji_ascii	100	5	61	winowhy	100	2
31	kannada	100	4		Sum	5384	-

B.6.2 HYPERPARAMETERS SETTINGS

In our study, we ensure consistency in the hyperparameter settings across all baselines, including our proposed EDIT approach, to maintain the fairness of our comparative analysis. Here, we detail the hyperparameter configurations employed in our experiments.

Training Steps and Batch Size. The number of training steps is determined based on the size of the training dataset, the batch size, and the number of gradient accumulation steps required. We maintain a consistent batch size across all baselines to eliminate any performance discrepancies that could arise from varying batch sizes.

Learning Rate. Our initial exploratory experiments focused on the standard CoTs distillation method using the LLaMA-2 model. We found that while the batch size had minimal impact on performance, the learning rate was a critical factor. We tested learning rates of $1e-4$, $2e-4$, and $3e-4$, observing optimal performance at $2e-4$ across the standard CoT and other distillation baselines, as well as our EDIT approach. Consequently, we set the learning rate to $2e-4$ for all methods involved in our study.

Epochs and Evaluation Strategy. Throughout our training process, we monitored the training loss curve and noted that it generally plateaued by the 15th epoch, indicating that the models had achieved convergence. Therefore, we set the number of epochs to 15 for 7B models. The process of determining the number of epochs for other model sizes followed a similar pattern. To mitigate the potential risk of overfitting and to ensure our evaluation reflects the most effective model configuration, we systematically selected checkpoints from the epoch that demonstrated the best performance on the IND task. These checkpoints were then used to evaluate performance on OOD tasks.

The hyperparameters in training and inference can be found in Table 18 and Table 15 respectively. In the KRSL, the second phase training in EDIT, the learning rate is empirically set as $5e-6$.

Table 18: Training hyperparameters.

Hyperparameter	TinyLLaMA-1.1B	LLaMA2-13B	LLaMA2-7B / CodeLLaMA-7B / LLaMA3-8B / Mistral-7B-v0.2
gradient accumulation steps	4	8	4
per device batch size	16	8	16
learning rate	2e-4	2e-4	2e-4
epoches	20	15	10
max length	1024	1024	1024
β of AdamW	(0.9,0.999)	(0.9,0.999)	(0.9,0.999)
ϵ of AdamW	1e-8	1e-8	1e-8
γ of Scheduler	0.95	0.95	0.95
weight decay	0	0	0
warmup ratio	0	0	0
rank of LoRA	64	64	64
α of LoRA	32	32	32
target modules	q_proj, v_proj	q_proj, v_proj	q_proj, v_proj
drop out of LoRA	0.05	0.05	0.05

B.6.3 COMPUTATION BUDGET

Our experimental code is based on modifications of Meta’s open-source `llama-recipes`⁸, utilizing the FSDP framework and training the model in parallel on four 80GB A100 GPUs. In our experimental setup, training a 7B model during the SFT stage takes approximately 40 minutes per epoch. For KRSL, each epoch takes around 90 minutes. With the same settings, training the Mistral model will see about a 10% increase in training time. We will release our code in the future.

C PROMPT TEMPLATES

C.1 CoTs EXTRACTION PROMPT

We use the prompt template shown in Table 25 to call the ChatGPT API to generate the CoTs for the BBH-train datasets.

C.2 ANSWER HINT PROMPT

We list the Answer Hint Prompt templates in Table 24, which imply the teacher LLMs to generate the CoTs based on the given answers following the in-context examples.

C.3 CONTRASTIVE CoTs PROMPT

We list the Contrastive CoTs Prompt templates in Table 26, which query the teacher LLMs to generate the CoTs with similar rationales to the original ones but divergent answers by following the few examples provided with contrastive CoT pairs.

C.4 EVALUATION PROMPT OF CoTs QUALITY

We list the evaluation prompt templates of CoTs quality in Table 27.

C.5 MISTAKE PATTERN MINING PROMPT

For mistake pattern mining, we employ the prompt template delineated in Table 28, which includes the definitions of the four distinct mistake patterns.

⁸<https://github.com/Meta-Llama/llama-recipes>

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133




Table 19: A reasoning about colored objects case from BIG-Bench Hard.

Question:	On the table, there are three pink scrunchiephone chargers, two orange crayons, and two black scrunchiephone chargers. If I remove all the black items from the table, how many scrunchiephone chargers remain on it? Options:(A) zero(B) one (C) two (D) three (E) four (F) five (G) six (H) seven (I) eight (J) nine (K) ten (L) eleven (M) twelve (N) thirteen (O) fourteen (P) fifteen (Q) sixteen
Teacher	According to this question, there are three pink scrunchiephone chargers, two orange crayons, and two black scrunchiephone chargers on the table. If we remove all the black items from the table, we are left with three pink scrunchiephone chargers and two orange crayons. Therefore, the number of scrunchiephone chargers that remain on the table is five. Therefore, the answer is (F). ❌
Std-CoT	According to this question, there are three pink scrunchiephone chargers, two orange crayons, and two black scrunchiephone chargers on the table. If we remove all the black items from the table, we are left with three pink scrunchiephone chargers, two orange crayons, and two black scrunchiephone chargers. The number of scrunchiephone chargers remaining on the table is seven. Therefore, the answer is (H). ❌
EDIT	According to this question, there are a total of 5 scrunchiephone chargers on the table: 3 pink scrunchiephone chargers and 2 black scrunchiephone chargers. If we remove all the black items, we are left with 3 pink scrunchiephone chargers. Therefore, the answer is (D). ✅

Table 20: A movie recommendation case from BIG-Bench Hard.

Question:	Find a movie similar to Jurassic Park, Independence Day, Forrest Gump, Batman. Options: (A) The Firm (B) Roger Dodger (C) Innocent Man (D) An (E) The Silence of the Lambs
Teacher	- Jurassic Park (action, adventure, science-fiction; 1993) - Independence Day (action, science-fiction, drama; 1996) - Forrest Gump (comedy, drama, romance; 1994) - Batman (action, adventure; 1989) These are all famous Hollywood movies produced in the 1980s and 1990s. Amongst all the options, the only movie similar to these ones seems to be The Firm (drama, mystery, thriller; 1993), as it is also a popular movie from the same era. Therefore, the answer is (A). ❌
Std-CoT	Jurassic Park is an action and adventure movie. Independence Day is an action, science-fiction, and adventure movie. Forrest Gump is a comedy, drama, and romance movie. Batman is an action, adventure, and crime movie. The Firm is a crime, drama, and thriller movie. Roger Dodger is a comedy and drama movie. Innocent Man is a crime, drama, and thriller movie. The Silence of the Lambs is a crime, drama, and thriller movie. Amongst all the options, the movie similar to Jurassic Park seems to be The Firm. Therefore, the answer is (A). ❌
EDIT	Issue 1: Is this movie similar to Jurassic Park? Above all the options, the only movie similar to Jurassic Park seems to be The Silence of the Lambs due to its suspenseful and thrilling nature. Issue 2: Is this movie similar to Independence Day? Above all the options, the only movie similar to Independence Day seems to be Roger Dodger due to its sci-fi and action-packed nature. Issue 3: Is this movie similar to Forrest Gump? Above all the options, the only movie similar to Forrest Gump seems to be The Firm due to its drama and emotional depth. Issue 4: Is this movie similar to Batman? Above all the options, the only movie similar to Batman seems to be An due to its superhero and action-packed nature. Based on the similarities outlined in Issues 1, 2, 3, and 4, the movie similar to these given movies seems to be The Silence of the Lambs. Therefore, the answer is (E). ✅

Table 21: A dyck languages case from BIG-Bench Hard.

1134		
1135		
1136	Question:	Correctly close a Dyck-n word. Complete the rest of the sequence, making sure that the parentheses are closed properly.
1137		Input: < ((< >)
1138		We should process each input one by one and keep track of the stack configuration.
1139		0: empty stack
1140		1: < ; stack: <
1141		2: (; stack: < (
1142		3: (; stack: < ((
1143		4: { ; stack: < (({
1144		5: < ; stack: < (({ <
1145		6: > ; stack: < (({
1146	Teacher	7: } ; stack: < (({
1147		8:) ; stack: < ((
1148		Now, we have reached the end. The final stack is "< (((".
1149		We will need to pop out "(", "(" one by one in that order.
1150		So, we need ")", ")". Therefore, the answer is)). 
1151		We should process each input one by one and keep track of the stack configuration.
1152		0: empty stack
1153		7
1154		1: < ; stack: <
1155		2: (; stack: < (
1156		3: (; stack: < ((
1157		4: { ; stack: < (({
1158		5: < ; stack: < (({ <
1159		6: > ; stack: < (({
1160		7: } ; stack: < (({
1161		8:) ; stack: < (({
1162	Std-CoT	Now, we have reached the end. The final stack is "< (((".
1163		We will need to pop out "{", "(", "<" one by one in that order.
1164		So, we need "}", ")", ">". Therefore, the answer is }) > . 
1165		We should process each input one by one and keep track of the stack configuration.
1166		0: empty stack
1167		1: < ; stack: <
1168		2: (; stack: < (
1169		3: (; stack: < ((
1170		4: { ; stack: < (({
1171		5: < ; stack: < (({ <
1172		6: > ; stack: < (({
1173		7: } ; stack: < (({
1174		8:) ; stack: < ((
1175	EDIT	Now, we have reached the end. The final stack is "< (((".
1176		We will need to pop out "(", "<" one by one in that order.
1177		So, we need ")", ">". Therefore, the answer is) > . 

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241

Table 22: A AQUA-RAT case from AGIEval.

Question:	Anna has 4 books. If she decide to arrange the 4 books in every possible combination and moved just one book every minute, how long would it taken by her ? Options: (A)22 minutes (B)30 minutes (C)15 minutes (D)24 minutes (E)35 minutes
Teacher	<p>First, Anna has 4 books and she wants to arrange them in every possible combination.</p> <p>To find the number of possible combinations, we can use the formula for permutations, which is $n!$ (n factorial), where n is the number of items. In this case, $n = 4$, so the number of possible combinations is $4! = 4 \times 3 \times 2 \times 1 = 24$.</p> <p>Next, we know that Anna can move just one book every minute. Since there are 24 possible combinations, it will take her 24 minutes to go through all of them.</p> <p>Therefore, the answer is (D) 24 minutes. ✓</p>
Std-CoT	<p>This question focuses on the time it takes Anna to arrange the 4 books in every possible combination. We know that Anna moved just one book every minute.</p> <p>Plugging in the minutes per move: 1 minute There are 4 books: A, B, C, D There are $(4!) = 4$ different combinations: A, B, C, D; B, A, C, D; B, C, A, D; C, A, B, D; D, A, B, C</p> <p>The total number of minutes it would take Anna to arrange the books in every possible combination is the sum of the minutes needed for each combination: $(4 \text{ minutes}) + (4 \text{ minutes}) + (4 \text{ minutes}) + (4 \text{ minutes}) = 15 \text{ minutes}$.</p> <p>Therefore, the answer is (C) 15 minutes. ✗</p>
EDIT	<p>This question asks for the total amount of time it would take Anna to arrange the books in every possible combination and moved just one book every minute. We first need to calculate the number of possible combinations: $4! = 24$. Then, we need to multiply the number of possible combinations by the number of minutes per combination: $24 * 1 \text{ minute} = 24 \text{ minutes}$.</p> <p>Therefore, the answer is (D) 24 minutes. ✓</p>

Table 23: A case from AI2 Reasoning Challenge.

Question:	Which method is the safest way to watch an eclipse of the Sun? Options: (A) Turn away after two or three minutes. (B) Look at the Sun through a long telescope. (C) Cast an image through a pinhole onto a screen. (D) Blink often until your eyes get used to the light.
Teacher	<p>First, it is important to note that directly looking at the Sun during an eclipse can be extremely dangerous and can cause permanent damage to the eyes.</p> <p>Option (A) suggests turning away after two or three minutes. However, even a brief exposure to the Sun during an eclipse can be harmful, so this method is not safe.</p> <p>Option (B) suggests looking at the Sun through a long telescope. This method is also unsafe because the intense light from the Sun can damage the telescope and cause harm to the eyes.</p> <p>Option (C) suggests casting an image through a pinhole onto a screen. This method is a safe and recommended way to watch an eclipse. By creating a pinhole camera, the image of the Sun is projected onto a screen, allowing for safe viewing without directly looking at the Sun.</p> <p>Option (D) suggests blinking often until your eyes get used to the light. This method is not safe as it does not provide any protection from the harmful rays of the Sun during an eclipse.</p> <p>Therefore, the answer is (C) Cast an image through a pinhole onto a screen. ✓</p>
Std-CoT	<p>We know that looking directly at the Sun can cause permanent damage to the eyes. Therefore, the answer is (C) casting an image through a pinhole onto a screen. ✓</p>
EDIT	<p>We are told that this question is about the safest way to watch an eclipse of the Sun. Looking at the Sun through a long telescope is not recommended as it can cause permanent damage to the eyes. Blinking often until your eyes get used to the light is also not recommended as it can cause eye strain and discomfort.</p> <p>The only remaining option is to cast an image through a pinhole onto a screen. This method involves creating a small hole in a card or paper and projecting the Sun's image onto a surface, such as a wall or a piece of paper. This method is safe as it does not require direct eye contact with the Sun.</p> <p>Therefore, the answer is (C). ✓</p>

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

Table 24: Answer Hint Prompt templates for rectifying the wrong CoTs data based on the hint answers.

{Task Description}. Your response should conclude with the format "Therefore, the answer is".

Q: {Task Example Question No.1}
A: Let's think step by step. {Human-Curated-CoTs No.1}.

Q: {Task Example Question No.2}
A: Let's think step by step. {Human-Curated-CoTs No.2}.

Q: {Task Example Question No.2}
A: Let's think step by step. {Human-Curated-CoTs No.3}.

Q: {QUESTION}
A: Let's think step by step.

Table 25: CoTs extraction prompt template of gpt-3.5-turbo for generating the CoTs data.

{Task Description}. Your response should conclude with the format "Therefore, the answer is".

Q: {Task Example Question No.1}
H: {The correct answer is [HINT ANSWER No.1]}
A: Let's think step by step. {Human-Curated-CoTs No.1}.

Q: {Task Example Question No.2}
H: {The correct answer is [HINT ANSWER No.2]}
A: Let's think step by step. {Human-Curated-CoTs No.2}.

Q: {Task Example Question No.3}
H: {The correct answer is [HINT ANSWER No.3]}
A: Let's think step by step. {Human-Curated-CoTs No.3}.

Q: {QUESTION}
H: {The correct answer is [HINT ANSWER]}
A: Let's think step by step.

1296
 1297
 1298
 1299
 1300
 1301
 1302
 1303
 1304
 1305
 1306
 1307
 1308
 1309
 1310
 1311
 1312
 1313
 1314
 1315
 1316
 1317
 1318
 1319
 1320
 1321
 1322
 1323
 1324
 1325
 1326
 1327
 1328
 1329
 1330
 1331
 1332
 1333
 1334
 1335
 1336
 1337
 1338
 1339
 1340
 1341
 1342
 1343
 1344
 1345
 1346
 1347
 1348
 1349

Table 26: Contrastive CoTs Prompt templates for mistaken the correct CoTs data. The examples are sampled from the teachers’ original wrong CoTs data and its corrected CoTs. In this way, teacher LLMs can expose the reasoning flaws in problems that were originally solved correctly.

{Task Description}. You need to complete the [Wrong Response] which requires you to give the most likely incorrect answer to the [Question] and the rationale for the incorrect answer. The incorrect answer and rationale in the [Wrong Response] must be different from the correct answer and rationale in the [Right Response].

[Question]: {Task Example Question No.1}
 [Right Response]: {Corrected CoT No.1}
 [Wrong Response]: {Wrong CoT No.1}

[Question]: {Task Example Question No.2}
 [Right Response]: {Corrected CoT No.2}
 [Wrong Response]: {Wrong CoT No.2}

[Question]: {Task Example Question No.3}
 [Right Response]: {Corrected CoT No.3}
 [Wrong Response]: {Wrong CoT No.3}

[Question]: {USER_QUESTION}
 [Right Response]: {Corrected CoT}
 [Wrong Response]:

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

Table 27: Prompt template of GPT-4 for assessing CoTs quality. In the analysis, we use this template to eval the quality of CoTs generated by Std-CoT, EDIT and the teacher LLM respectively.

```
[System] You are a helpful and precise assistant for assessing the
quality of the response.

[Question]: {QUESTION}
[Reference Answer]: {ANSWER}

[AI Assistant 1's Answer Start]
{ASSISTANT1}
[AI Assistant 1's Answer End]

[AI Assistant 2's Answer Start]
{ASSISTANT2}
[AI Assistant 2's Answer End]

[AI Assistant 3's Answer Start]
{ASSISTANT3}
[AI Assistant 3's Answer End]

[System] We would like to request your feedback, in the form of
scoring, on which of the
responses from AI Assistant 1, 2 and 3 effectively demonstrates the key
reasoning steps in
solving this question. Key Reasoning Steps refer to certain crucial
steps in the process of
logical reasoning or problem-solving. These steps play a significant
role in the thinking
process and have a notable impact on subsequent reasoning. Each
student will receive an
overall score on a scale of 1 to 10, where a higher score signifies
that the assistant's
response is more effectively demonstrates the key reasoning steps for
the question.
Please provide a comprehensive explanation, avoiding any potential bias
and ensuring that
the order in which the responses were presented does not affect your
judgment. And then
output three lines indicating the scores for AI Assistant 1, 2 and 3,
respectively.

Output with the following format:
Evaluation evidence: <your evaluation explanation here>
Score of AI Assistant 1: <score>
Score of AI Assistant 2: <score>
Score of AI Assistant 3: <score>
```

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

Table 28: Prompt templates of GPT-3.5 for classifying the mistakes. In the analysis, we use this template to classify the mistake data used in EDIT.

```
[System] You are a helpful assistant who is good at identifying types
of reasoning mistakes.
There are now three types of inference errors, as follows:

(a). Logical reasoning errors. This type of error involves the
logical structure of reasoning,
including assumptions, reasoning rules, argument chains, etc. Among
logical errors, students
may make errors such as invalid reasoning, insufficient or incorrect
assumptions, and jumps in
reasoning. Students may make errors in selecting reasoning strategies
or methods. The chosen
method may not be suitable for a specific problem, or may lead to
misleading reasoning.

(b). Knowledge errors in reasoning. This type of error involves
misunderstanding or incomplete
understanding of facts, concepts or knowledge, conceptual confusion,
and cognitive biases.

(c). Numerical calculation errors. This type of error involves
mathematical calculation errors,
which may include incorrect calculations, conversions or errors in the
processing of numerical
values.

(d). Other errors. All other errors that do not belong to the above
three categories.

I will give you a dictionary with the following fields and meanings:
{
  "input": reasoning question.
  "right_output": the correct answer.
  "wrong_output": the wrong answer.
}

You need to first form your own opinion about the problem based on the
reasoning questions and the
correct answers, and then analyze the reasons for the mistakes in the
wrong answers in "Rationale:".
Then give your classification results in "Category:", e.g., (a), (b) or
(c), etc. If an answer
involves errors in multiple categories, you should point them out and
connect them with '+' sign
in the category. For example, if an answer involves logical errors and
mathematical calculation
errors, then the category should be a+c.

You must output with the following format:
Rationale: <your analysis process and explanation of the final
classification results>
Category: <only fill in with a or b or c or a+b or a+c or b+c or a+b+c
or d.>
```
