

AbstRaL: Augmenting LLMs’ Reasoning by Reinforcing Abstract Thinking

Silin Gao^{1,2}, Antoine Bosselut², Samy Bengio¹, Emmanuel Abbe^{1,2}
¹Apple ²EPFL

Abstract

Recent studies have shown that large language models (LLMs), especially smaller ones, often lack robustness in their reasoning. I.e., they tend to experience performance drops when faced with distribution shifts, such as changes to numerical or nominal variables, or insertions of distracting clauses. A possible strategy to address this involves generating synthetic data to further “instantiate” reasoning problems on potential variations. In contrast, our approach focuses on “abstracting” reasoning problems. This not only helps counteract distribution shifts but also facilitates the connection to symbolic tools for deriving solutions. We find that this abstraction process is better acquired through reinforcement learning (RL) than just supervised fine-tuning, which often fails to produce faithful abstractions. Our method, AbstRaL—which promotes abstract reasoning in LLMs using RL on granular abstraction data—significantly mitigates performance degradation on recent GSM perturbation benchmarks.

1 Introduction

The ability of reasoning, which involves the integration of knowledge to derive dynamic conclusions rather than direct recourse to memorized information (e.g., [46, 1]), is an essential quality for artificial general intelligence [48]. Toward this end, recently developed large language models (LLMs) have been equipped with impressive reasoning capabilities, within the general scope [13, 41, 12] or in specialized domains such as mathematics [31, 42].

However, most LLMs, especially smaller ones¹, still face the challenge of robustness when reasoning, which still have considerable room for improvement in out-of-distribution (OOD) generalization. Recent works [22, 19, 24] have shown that even in simple grade school math tasks, LLMs suffer performance degradation when facing perturbations and distribution shifts. In particular, LLMs can be prone to reasoning errors on instantiation shifts such as when numerical or nominal variables are altered, even though the LLMs can respond to the original question correctly. On more challenging interferential shifts, where a distracting (topic-related but useless) condition is added to the question, LLMs suffer even more drastic performance drops.

To improve the robustness of reasoning, a possible learning strategy [3] is to synthesize more instances of a reasoning problem that are varied in surface-form contexts but following the same abstraction. In this paper, instead of scaling up the training instances (which can also be computationally expensive), we teach LLMs to **directly learn an abstraction underlying the reasoning problems and thereby learn to reason in a manner that is invariant to contextual distribution shifts.**

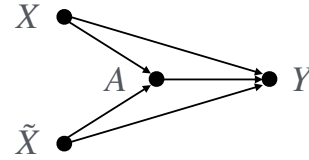


Figure 1: Two paraphrased queries X and \tilde{X} , having same solution Y , can be both handled by a common abstraction A .

¹as verified by our analysis in §5

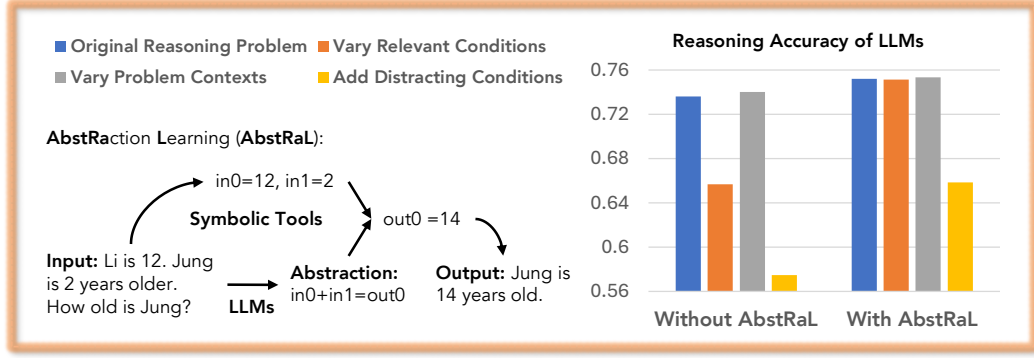


Figure 2: Our **Abstraction Learning (AbstRaL)** method effectively improves the reasoning robustness of LLMs, especially facing the variations of relevant input conditions and the interference of distracting conditions. We present average accuracy of all our tested LLMs on different GSM-Plus [19] testing sets, including the original GSM8K testing set (Original Reasoning Problem), the testing sets with numerical variations (Vary Input Conditions), averaged across three portions (digit expansion, integer-decimal-fraction conversion and numerical substitution), the testing set with problem rephrasing (Vary Problem Contexts) and with distractor insertion (Add Distracting Conditions).

We propose the reinforced **Abstraction Learning** framework, **AbstRaL**, as shown in Figure 2, which first teaches LLMs to generate abstraction of input problem, and then connects the abstraction with symbolic tools to stably derive the output solution, regardless of the specific input contexts. The learning of abstraction relies on our **Granularly-decomposed Abstract Reasoning (GranularAR)** data distilled from an oracle LLM, which integrates symbolic reasoning within socratic problem decomposition [32] and chain-of-thought (CoT) explanations [36]. On top of supervised fine-tuning (SFT), AbstRaL uses reinforcement learning (RL) with a new set of model-free rewards to further improve the faithfulness of generated abstraction.

We tested AbstRaL on two benchmarks that evaluate the robustness of mathematical reasoning, GSM-Symbolic [22] and GSM-Plus [19]. Experimental results on various seed LLMs consistently demonstrate that AbstRaL effectively improves the reasoning robustness of LLMs. As shown in Figure 2, AbstRaL almost reverts the performance drop of LLMs caused by variations of relevant input conditions, and also significantly mitigates the interference of distracting conditions added to the perturbed testing samples.

2 Preliminary: Learning Strategies to Improve Reasoning Robustness

We assume that every reasoning data sample, consisting of input question (or query) \mathcal{X} and output answer (or response) \mathcal{Y} , is an instantiation of an underlying symbolic abstraction \mathcal{A} that represents high-level knowledge or reasoning schema. For example, in Figure 3, the problem of calculating Jung’s age, according to how much older he is than Li, is based on the abstract arithmetic rule of adding two numbers. A robust reasoner is supposed to master the abstraction \mathcal{A} , and therefore stably give a faithful answer \mathcal{Y} to any question \mathcal{X} implicitly derived from \mathcal{A} , rather than overfitting to only a subset of instances of \mathcal{A} and vulnerable to distribution shifts that go beyond the subset.

A common strategy to improve the reasoning robustness of LLMs is to augment the learning data by synthesizing more instances $\{(\mathcal{X}', \mathcal{Y}'), (\mathcal{X}'', \mathcal{Y}''), \dots\}$ of abstraction \mathcal{A} , with paraphrasing [9, 49] or templates [3]. For example, as shown in Figure 3 (a), the numbers and names appearing in an instance can be replaced with other values to create a new instance, which has a different problem context, but follows the same abstract arithmetic rule. As verified by previous study [3], this learning strategy requires a large amount of synthetic data augmentation, to effectively boost LLMs’ grasp of the high-level abstraction, thus being able to resist the interference of surface-form variations.

In this work, instead, we focus on the strategy [15, 10] of teaching LLMs to directly learn the abstraction \mathcal{A} underlying each instance $(\mathcal{X}, \mathcal{Y})$, and connect \mathcal{A} with symbolic tools such as an equation solver, to steadily derive the answer \mathcal{Y} to the question \mathcal{X} , as illustrated in Figure 3 (b).

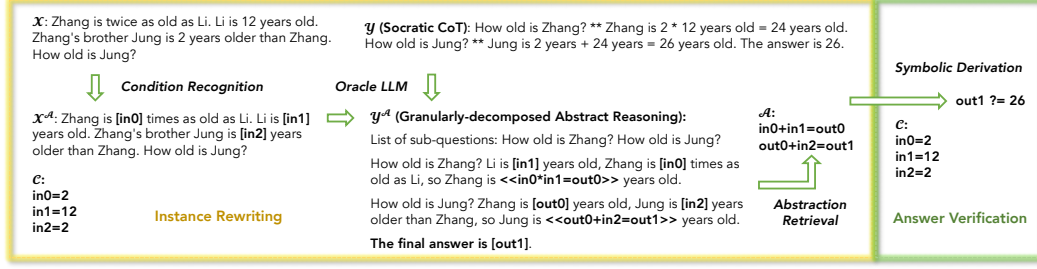


Figure 4: Overview of GranularAR training data construction, which consists of an instance rewriting procedure to rewrite existing socratic CoT data (\mathcal{X}, \mathcal{Y}) into fine-grained abstract reasoning data ($\mathcal{X}^A, \mathcal{C}, \mathcal{Y}^A, \mathcal{A}$), followed by an answer LLM verification procedure to check the correctness of rewriting.

3.1 Granularly-decomposed Abstract Reasoning (GranularAR) Data

Motivations LLMs have learned various fine-grained reasoning strategies at either the pre-training [44, 5] or the post-training [17] phase, such as chain-of-thought (CoT) [36] and socratic problem decomposition [32] as representatives. In our GranularAR training data, we integrate abstract reasoning with these pre-learned strategies, which enables LLMs to gradually construct the abstraction in a fine-grained reasoning chain, as shown in Figure 4 \mathcal{Y}^A . Such data format is close to the pre-training or the post-training data distribution, and therefore eases the difficulty of adapting LLMs to our new abstract reasoning manner.

GranularAR Format The answer \mathcal{Y}^A of GranularAR first decomposes the question \mathcal{X}^A into a list of sub-questions, which enables a holistic planning of step-by-step reasoning. Based on that, \mathcal{Y}^A answers each sub-question with CoT and abstract symbols, where it first quotes relevant input conditions (or answers to previous sub-questions), and then derives the answer with quoted symbols. Finally, \mathcal{Y}^A draws a conclusion to clarify which output abstract symbol represents the final answer.

Data Construction Figure 4 illustrates how we construct the GranularAR training data. We first conduct the condition recognition described in §3 to formulate the conditions \mathcal{C} from the question \mathcal{X} and creates the abstract question \mathcal{X}^A . Based on that, we prompt an oracle LLM to rewrite the gold socratic CoT answer \mathcal{Y} into our desired abstract answer \mathcal{Y}^A . The abstract question \mathcal{X}^A is also fed into the oracle LLM, to complement the problem contexts and clarify the abstract symbols of input variables. Given the distilled abstract answer \mathcal{Y}^A , we then conduct the abstraction retrieval to get the de-contextualized abstraction \mathcal{A} . Finally, the symbolic derivation step is performed to verify whether \mathcal{A} along with \mathcal{C} can derive the correct final answer stated in \mathcal{Y} . We only keep the rewritten instances that pass the answer verification.

3.2 Learning of Abstract Reasoning

Motivations Previous study [10] has shown that LLMs are poor at following in-context demonstrations to reason in abstract manner, indicating that the learning of abstract reasoning requires training with proper supervision, rather than relying on only in-context instructions and examples. A straightforward way is to train LLMs with supervised fine-tuning (SFT). However, although SFT on abstract data can teach LLMs decent abstract reasoning formats, its auto-regressive training objective also forces LLMs to learn the specific contexts of each training sample. This hinders LLMs from learning more general abstract thinking strategy, which leads to frequent test-time failure of generating an abstraction that is aligned with the problem, skewed by the new contexts in the testing data³, as shown by our results in §5. Therefore, we propose to conduct reinforcement learning (RL) on top of SFT, to augment LLMs' capability of constructing faithful abstractions.

Supervised Fine-Tuning We fine-tune LLMs to auto-regressively generate our constructed GranularAR answer \mathcal{Y}^A based on the input question \mathcal{X}^A , simply with the causal language modeling loss of predicting each token in \mathcal{Y}^A based on former tokens.

³On the other hand, learning to directly generate de-contextualized abstraction is rather hard for LLMs that are pre-trained on mostly contextualized natural language corpus, as verified by our analysis in §5, which motivates our use of a fine-grained framework to still incorporate contexts in the learning of abstract reasoning.

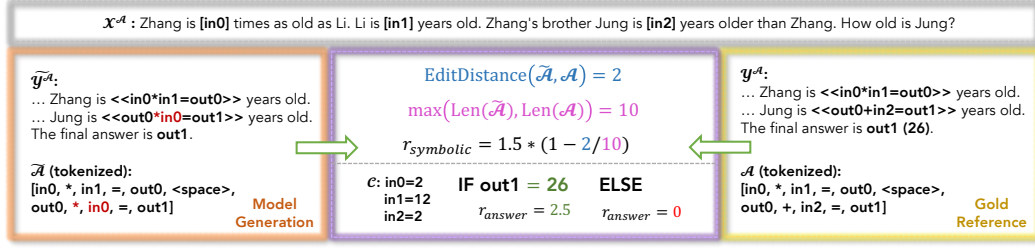


Figure 5: Illustration of the abstraction rewards in our reinforcement learning approach, including the symbolic distance reward $r_{symbolic}$ and the answer correctness reward r_{answer} .

Reinforcement Learning with Abstraction Rewards On top of SFT, we develop a RL approach to further improve the faithfulness of abstraction generated by LLMs. Our RL approach proposes a new set of rewards to closely rate the generated abstraction on two aspects. First, given the model created abstraction \tilde{A} retrieved from its generated answer \tilde{y}^A , we perform the symbolic derivation to check whether \tilde{A} can derive the correct final answer (denoted as Ans) with the conditions C given in the gold reference. If yes, a positive reward $r_{correct}$ (hyperparameter) is granted to the model, otherwise zero reward is given. We denote this **answer correctness** reward as $r_{answer}(\tilde{A}, C, \text{Ans})$. Second, we more granularly measure how \tilde{A} is aligned with the expected abstraction A retrieved from the gold answer y^A . Specifically, we split \tilde{A} and A into lists of symbolic tokens, where each token is either an abstract symbol (created in conditions C or abstract answer \tilde{y}^A or y^A) that represents an input or output variable (such as in0), or a pre-defined operator that connects variables (such as =) or separates derivations (such as <space>), as shown in Figure 5. Based on that, we calculate a **symbolic distance** reward:

$$r_{symbolic}(\tilde{A}, A) = r_{max} \cdot \left(1 - \text{EditDistance}(\tilde{A}, A) / \max_{a \in \{\tilde{A}, A\}} \text{Len}(a)\right) \quad (1)$$

where r_{max} denotes the maximum reward hyperparameter, $\text{EditDistance}(\tilde{A}, A)$ denotes the list-wise edit distance [18] between tokenized abstractions \tilde{A} and A , and $\text{Len}(\cdot)$ calculates the length of a list, used for normalizing the edit distance into the range of 0 to 1.⁴ A higher symbolic distance reward is granted to the model if \tilde{A} is closer (or more similar) to A , which gives the model more fine-grained learning signal of how far away it is from creating the correct abstraction. Figure 5 shows an example case of our abstraction rewards. Note that our proposed rewards do not require a pre-trained reward model, but just comparison to the gold reference.⁵ We plug our abstraction rewards with GRPO [31], an advanced RL algorithm for incentivizing reasoning capability in LLMs [13]. The formulation of GRPO with our proposed rewards is described in Appendix A.

4 Experimental Settings

Implementations of AbstRaL Framework Our experimental study applies AbstRaL to the task of mathematical reasoning, a representative domain for reasoning robustness research [22, 19, 24].

In particular, we create few-shot examples to prompt a Llama-3.3-70B-Instruct [12] model to accomplish the first **condition recognition** step. As shown in Figure 4, the LLM is tasked to label the numerical values in the input question \mathcal{X} with square brackets “[]”, and then sequentially replace each labeled value with an abstract symbol indexing as “in0”, “in1”, etc., to create the abstract question \mathcal{X}^A . Note that for implicit numerical values such as “one hundred” and “twice”, we also prompt the LLM to convert them into explicit format “100” and “2 times”, so that the numbers can be labeled and replaced. Meanwhile, the LLM is asked to use equations to record the replacements, e.g., “in0 = 2”, in order to create the conditions C .

We then tune various LLMs on our constructed GranularAR data (§3.1), with our SFT and RL scheme (§3.2), to perform the core **abstract reasoning** step. Our tested LLMs include Llama3 model series

⁴The maximum possible edit distance is the length of the longest list.

⁵This reduces computational cost and avoids potential preference biases [31] introduced by the reward model.

[12] (Llama-3.2-1B-Instruct, Llama-3.2-3B-Instruct and Llama-3.1-8B-Instruct), Qwen2.5 model series [41] (Qwen2.5-0.5B-Instruct, Qwen2.5-1.5B-Instruct, Qwen2.5-3B-Instruct, Qwen2.5-7B-Instruct and its math specialized version Qwen2.5-Math-7B-Instruct) and Mathstral-7B⁶. All LLMs are supposed to follow our Granular training data to index the derived output variables in their generations as “out0”, “out1”, etc., and highlight each abstract math derivation in double angle brackets such as “<<in0 * in1 = out0>>”. LLMs are also supposed to generate a fixed statement “The final answer is [outN].”, to clarify the output variable that represents the final answer. At inference phase, LLMs use greedy decoding to generate their abstract answers.

In the **abstraction retrieval** step, we simply use a regex-matching script to extract all math derivations that are enclosed in double angle brackets in the model generated answer $\widetilde{\mathcal{Y}}^{\mathcal{A}}$ (or $\mathcal{Y}^{\mathcal{A}}$ in our Granular data construction). All extracted math derivations form the problem abstraction $\widetilde{\mathcal{A}}$ (or \mathcal{A}). During the tokenization of the abstraction (used for calculating the symbolic distance reward in our RL approach), the tokenized lists of derivations are concatenated with a special token “<space>”, as shown in Figure 5. To perform the final **symbolic derivation**, we treat the output math derivations in $\widetilde{\mathcal{A}}$ (or \mathcal{A}) and the input conditions in \mathcal{C} jointly as a system of equations, which is fed into a SymPy⁷ equation solver to derive the final answer.

For Granular training **data construction**, we use Llama-3.3-70B-Instruct as the oracle LLM, and prompt it with few-shot examples to rewrite existing socratic CoT data. We use the socratic version of GSM8K [8] training set from OpenAI⁸ as the seed data for rewriting, which contains 7473 grade school math word problems. After rewriting and filtering, 6386 problems are finally kept for training.

Evaluation Datasets We evaluate our method on two datasets derived from the GSM8K testing samples, which are GSM-Symbolic [22] and GSM-Plus [19].

GSM-Symbolic manually constructs problem templates from 100 GSM8K testing samples, and uses the templates to create new problems where the numbers or names or both in the original problem are varied to different values, denoted as **Vary Num.**, **Vary Name** and **Vary Both**, respectively. We follow GSM-Symbolic to conduct 50 rounds of evaluation, each round creating 1 problem from each template (so models are tested on 100 new problems per round), and measure the average performance across these 50 rounds, *i.e.*, the mean (with standard deviation) of accuracy, and check whether it matches the performance on the original 100 (**Origin 100**) GSM8K problems.

GSM-Plus creates different variations of the full GSM8K testing set instead, where each varied testing set contains all 1319 GSM8K testing problems. For each type of variation, 1 varied sample is created for each original problem, which leads to a single round of evaluation per variation. We test our method on a subset of GSM-Plus variations⁹, including digit expansion (**Digit Ex.**) that adds more digits to a number (*e.g.*, from 16 to 1600), integer-decimal-fraction conversion (**Int-Dec-Fra**) that changes the type of a number (*e.g.*, from 2 to 2.5), numerical substitution (**Num. Sub.**) that replaces a number with another same-digit one (*e.g.*, from 16 to 20), rephrasing the question to check problem understanding (**Rephrase**), and distractor insertion (**Distract**) that adds topic-related but useless conditions, compared with model performance on the **Original** GSM8K testing set.

Baseline Methods We compare AbstRaL to several baseline reasoning methods.

CoT-8S prompts LLMs with the demonstration template suggested by GSM-Symbolic [22] and the common 8-shot examples¹⁰ used for GSM8K evaluation, to generate CoT answer to the input question \mathcal{X} . The last number (typically after “The final answer is”) is extracted as the final answer.

CoT-RL tunes LLMs on the non-rewritten GSM8K training data (\mathcal{X}, \mathcal{Y}) (socratic CoT version¹¹ with 7473 samples), using the same SFT and RL algorithm (GRPO [31]) as our approach. Since CoT-RL does not generate abstractions in the reasoning chains, only our answer correctness reward r_{answer}

⁶<https://huggingface.co/mistralai/Mathstral-7B-v0.1>

⁷<https://github.com/sympy/sympy>

⁸<https://huggingface.co/datasets/openai/gsm8k/tree/main/socratic>

⁹Some of GSM-Plus and GSM-Symbolic variations that modify the underlying abstraction, such as adding a useful condition to complicate a problem, are excluded and beyond the scope of our study.

¹⁰https://github.com/EleutherAI/lm-evaluation-harness/blob/main/lm_eval/tasks/gsm8k/gsm8k-cot.yaml

¹¹We also did a pilot study of training LLMs with non-socratic CoT data, which achieves similar results.

Table 1: Evaluation results on **GSM-Symbolic**. Δ denotes the relative percentage of drop comparing performance on **Vary Both** to performance on **Origin 100**. Best results on each model are **bold**, where lower is better for Δ . Standard deviation (std) of multi-round evaluation results are in brackets, where lowest std on each model are underlined.

Model	Method	Vary Num.	Vary Name	Vary Both	Origin 100	Δ (%)
Llama-3.2-1B-Instruct	CoT-8S	0.3864 (0.030)	0.4250 (0.026)	0.3890 (0.038)	0.4800	18.96
	CoT-RL	0.4788 (0.031)	0.5342 (0.024)	0.4488 (0.038)	0.5700	21.26
	CoA	0.4534 (0.020)	0.4372 (0.026)	0.4242 (0.030)	0.4600	7.78
	AbstRaL	0.5912 (<u>0.016</u>)	0.5838 (<u>0.023</u>)	0.5804 (<u>0.027</u>)	0.6000	3.27
Llama-3.1-8B-Instruct	CoT-8S	0.8440 (0.026)	0.8746 (0.021)	0.8236 (0.032)	0.8700	5.33
	CoT-RL	0.7784 (0.029)	0.8710 (0.014)	0.7540 (0.026)	0.8300	9.16
	CoA	0.7240 (0.019)	0.7086 (0.019)	0.6944 (0.025)	0.7200	3.56
	AbstRaL	0.8686 (<u>0.013</u>)	0.8672 (0.018)	0.8620 (<u>0.023</u>)	0.8700	0.92
Qwen2.5-0.5B-Instruct	CoT-8S	0.3394 (0.039)	0.3724 (<u>0.024</u>)	0.3398 (0.033)	0.3800	10.58
	CoT-RL	0.3192 (0.025)	0.3948 (0.032)	0.3228 (0.032)	0.3500	7.77
	CoA	0.2866 (0.025)	0.3060 (0.026)	0.2872 (0.026)	0.2900	0.97
	AbstRaL	0.4396 (<u>0.015</u>)	0.4416 (0.026)	0.4456 (<u>0.025</u>)	0.4400	-1.27
Qwen2.5-Math-7B-Instruct	CoT-8S	0.8956 (0.021)	0.9108 (0.018)	0.8766 (0.023)	0.9500	7.73
	CoT-RL	0.8942 (0.022)	0.9154 (<u>0.012</u>)	0.8812 (0.021)	0.9600	8.21
	CoA	0.7122 (0.028)	0.6976 (0.031)	0.6970 (0.033)	0.7100	1.83
	AbstRaL	0.9066 (<u>0.015</u>)	0.9014 (0.013)	0.9022 (<u>0.016</u>)	0.9100	0.86

(§3.2) is used as the learning signal in RL, which checks whether the final answer (last number) extracted from the generated CoT matches the gold answer.

CoA [10] is another abstract reasoning method that also fine-tunes LLMs to generate reasoning chains with abstract symbols, and plugs the generation with SymPy equation solver to derive the final answer. However, different from AbstRaL, only the output numbers in CoA reasoning chains are represented by abstract symbols, without abstracting the input conditions (\mathcal{X}^A and \mathcal{C}) and using our granularly decomposed reasoning chains (GranularAR), and the learning of CoA is based on only SFT, without integrating proper RL. Similar to our GranularAR training data construction, CoA prompts an oracle LLM to rewrite existing CoT data into abstract format, whose output numbers are replaced with abstract symbols, and then also uses the equation solver to verify the correctness of rewritten samples. For fair comparison, we use the same oracle LLM (Llama-3.3-70B-Instruct) and seed data (socratic CoT version of GSM8K training set) to construct the CoA training data.

5 Experimental Results

Table 1 shows some of our representative evaluation results on GSM-Symbolic dataset. We report the performances of the smallest and strongest models in our tested Llama3 and Qwen2.5 series, while leave the results of other models in Appendix B. We find that AbstRaL effectively improves the reasoning robustness of all models ranging from 0.5B to 8B sizes, with respect to the variations of both numbers and names (Vary Both), demonstrating better generalization to distribution shifts. Specifically, compared to baseline methods, models with AbstRaL achieve consistently better (mean) accuracy on Vary Both samples, with overall lower standard deviation across different testing rounds. Besides, models with AbstRaL suffer less performance drop (Δ) when transferring from Origin 100 to Vary Both, especially compared to CoT-based methods that do not use any abstract reasoning. Interestingly, on large-size Llama (8B) and Qwen (7B) models, although AbstRaL scores lower than 8-shot prompting (CoT-8S) on Origin 100 and Vary Name, it outperforms the prompting method on the Vary Num. and Vary Both perturbed samples. This implies that learning with AbstRaL may mitigate LLMs’ overfitting to the existing input conditions (or numbers), caused by potential data contamination [21, 38] at the pre-training stage.¹² By contrast, CoT-RL (without abstract reasoning) and CoA (without abstracting input conditions and RL on our GranularAR data) often fail to improve the tested models, which indicates that our fine-grained abstract reasoning framework with RL is more stable at adapting LLMs to robust reasoning manner.

¹²Note that with AbstRaL, although the input numbers are pre-extracted and abstracted into symbols, there is still minor gap between performances on Vary Num. and Origin 100. This is because minor errors may occur when our condition recognition tool (based on a LLM) extracts input numbers in the varied questions.

Table 2: Evaluation results on **GSM-Plus**. Best results on each model are **bold**.

Model	Method	Digit Ex.	Int-Dec-Fra	Num. Sub.	Rephrase	Distract	Original
Qwen2.5-0.5B-Instruct	CoT-8S	0.3601	0.2866	0.3958	0.4359	0.2267	0.4238
	CoT-RL	0.3336	0.2373	0.3571	0.4079	0.1524	0.3798
	CoA	0.2745	0.2267	0.2782	0.3161	0.1266	0.3033
	AbstRaL	0.4670	0.4663	0.4670	0.4625	0.3654	0.4670
Qwen2.5-Math-7B-Instruct	CoT-8S	0.8552	0.8218	0.8453	0.9052	0.7627	0.9181
	CoT-RL	0.8757	0.8522	0.8506	0.9037	0.8150	0.9340
	CoA	0.7460	0.6907	0.7180	0.7642	0.5709	0.7809
	AbstRaL	0.8916	0.8886	0.8916	0.8992	0.8226	0.8916

Table 2 presents our representative evaluation results on GSM-Plus dataset (full results are in Appendix B), which tests more diverse perturbations on the full GSM8K testing set. Similar to the results on GSM-Symbolic (Vary Num. and Vary Name), AbstRaL almost reverts the performance degradation caused by variations of input numbers (Digit Ex., Int-Dec-Fra and Num. Sub.), and also maintains robustness to contextual variations (Rephrase) comparable to baseline methods. Perhaps more interestingly, AbstRaL significantly mitigates the interference of distracting conditions added to the problems (Distract), while by comparison, models with baseline reasoning methods all score more drastically lower when transferring from Original to Distract. AbstRaL’s improvement on Distract is largely due to learning on its granularly-decomposed abstract reasoning (GranulAR) data (verified in our ablation study below), which enables overall planning of the reasoning steps at the start of the generation, and rethinking of useful input conditions at each reasoning step.

Ablation Study One natural concern of our method is whether the improvements of AbstRaL are due to the learning of abstract reasoning, or just because of integrating powerful symbolic tools, *i.e.*, Llama-3.3-70B-Instruct prompted for condition recognition, regex-matching script for abstraction retrieval and SymPy equation solver for symbolic derivation. To clarify this, we conduct an ablation study (**w/o Tools**), which replaces all above symbolic tools in our framework with the LLM itself, to accomplish each step of our pipeline with step-specific prompts. Results in Table 3 show that LLMs without additional tools suffer only minor performance drop, indicating that the learning of abstract reasoning in our framework is the major contributor to the improvements.

We also investigate the effect of using AbstRaL’s pipeline to construct the abstraction within contexts, *i.e.*, $\mathcal{X} \rightarrow \mathcal{X}^{\mathcal{A}} \rightarrow \mathcal{Y}^{\mathcal{A}} \rightarrow \mathcal{A}$, by testing an ablated framework (**w/o Contexts**), where LLMs are trained (with SFT and the same RL approach) to directly generate the abstraction \mathcal{A} based on the abstract input question $\mathcal{X}^{\mathcal{A}}$, *i.e.*, $\mathcal{X} \rightarrow \mathcal{X}^{\mathcal{A}} \rightarrow \mathcal{A}$ without explicitly generating $\mathcal{Y}^{\mathcal{A}}$ and post-processing it by abstraction retrieval.¹³ As shown in Table 3, LLMs that learn inferences without intermediate contexts fall far behind the unablated ones. This indicates that narrowing down the gap between symbolic abstraction and natural language modeling is important to faithful abstract thinking, which is in line with the globality degree [1] view of reasoning.

Our third ablation studies the contribution of RL in AbstRaL, by totally ablating the RL (**w/o RL**) or removing only the symbolic distance reward (**w/o $r_{symbolic}$**). On both the original and perturbed testing sets, we find that LLMs without RL score far lower than reinforcement-learned ones. Without $r_{symbolic}$ that hints how close the generated abstraction is to the correct one, LLMs also suffer significant performance drop, and thus under-perform the prompting baseline (CoT-8S). These findings demonstrate that acquiring faithful abstract reasoning requires careful learning, *e.g.*, via a proper RL approach and a milestone-style reward that closely monitors the progress of learning.

Lastly, we ablate the granularly-decomposed abstract reasoning (**w/o GranulAR**) used as AbstRaL’s training data format. We alter the gold abstract answer $\mathcal{Y}^{\mathcal{A}}$ in our constructed training data back to the standard socratic CoT format (as shown in Figure 4), where we remove listing sub-questions at the start of reasoning chain and quoting of input conditions at each reasoning step, while just keep the abstract derivations that forms the abstraction \mathcal{A} . LLMs trained on this ablated data suffer drastic performance drop on the Distract testing set, indicating that they are vulnerable to useless distractors. This verifies that our adopted fine-grained reasoning format plays an essential role in identifying interference conditions, duo to planning the reasoning steps and useful input conditions in each step.

¹³We also did a pilot ablation study of training LLMs to generate the abstraction \mathcal{A} directly from the original input question \mathcal{X} , *i.e.*, $\mathcal{X} \rightarrow \mathcal{A}$ without constructing intermediate $\mathcal{X}^{\mathcal{A}}$ and $\mathcal{Y}^{\mathcal{A}}$, which achieves similar results.

Table 3: Results of ablation study on GSM-Symbolic and GSM-Plus. **Num. Pert.** denotes average on the three GSM-Plus testing sets that perturb input numbers (Digit Ex., Int-Dec-Fra and Num. Sub.).

Model	Method	GSM-Symbolic			GSM-Plus			
		Vary Both	Origin 100	Δ (%)	Num. Pert.	Rephrase	Distract	Original
Qwen2.5-0.5B-Instruct	AbstRaL	0.4456	0.4400	-1.27	0.4668	0.4625	0.3654	0.4670
	- w/o Tools	0.4362	0.4400	0.86	0.4599	0.4564	0.3556	0.4610
	- w/o Contexts	0.2306	0.3000	23.13	0.2702	0.2881	0.1653	0.2866
	- w/o RL	0.2956	0.3600	17.89	0.3803	0.3230	0.2828	0.3571
	- w/o $r_{symbolic}$	0.3760	0.3900	3.59	0.4167	0.3995	0.3131	0.3882
	- w/o GranuLAR	0.4262	0.4200	-1.48	0.4236	0.4230	0.2358	0.4238
Qwen2.5-Math-7B-Instruct	AbstRaL	0.9022	0.9100	0.86	0.8906	0.8992	0.8226	0.8916
	- w/o Tools	0.8996	0.9100	1.14	0.8860	0.9037	0.8180	0.8893
	- w/o Contexts	0.6926	0.7100	2.45	0.6715	0.6936	0.5513	0.7065
	- w/o RL	0.7854	0.8200	4.22	0.8211	0.7998	0.7445	0.8226
	- w/o $r_{symbolic}$	0.8322	0.8700	4.34	0.8577	0.8234	0.7809	0.8590
	- w/o GranuLAR	0.8814	0.8900	0.97	0.8835	0.8976	0.6679	0.8855

6 Related Work

Reasoning Robustness Recent advances in LLM reasoning have also spotted considerable robustness challenges [47], reflected by benchmarks in symbolic reasoning, such as logical [2] and mathematical [34, 19, 33, 22] reasoning, and in factual or commonsense reasoning [43, 24, 25]. All above benchmarks consistently reveal that LLMs are easily perturbed by test-time data distribution shifts. Inspired by improving robustness via data augmentation [28], previous works [9, 27, 49, 3] use various data synthesis techniques to expand the coverage of training samples, and thus anticipate potential distribution shifts, which naturally increases the computational cost of developing a LLM. In this work, we aim to improve reasoning robustness of LLMs by incentivizing their abstract thinking, instead of replying on larger amount of reasoning instantiations.

Abstract Thinking and Planning Abstract thinking is an essential component of general fluid intelligence [6], and is also the key to human cognition and reasoning [29]. It requires making inferences based on abstract fundamental rules or concepts [4], rather than just memorizing a probabilistic pattern matching [11, 39]. Recently, various reasoning (or data) formats have been proposed for LLMs to learn abstract thinking, such as abstraction-of-thought (AoT) [15], chain-of-abstraction (CoA) [10], etc. However, all above methods lack a proper learning scheme beyond plain SFT, to stably adapt LLMs to abstract thinking. On the other hand, planning is also a basic human reasoning skill [37], which benefits LLM reasoning. Typical planning methods, such as chain-of-thought (CoT) [36], tree-of-thought (ToT) [45] and socratic problem decomposition [32], are widely adopted in LLMs to improve reasoning. Our work develops a better learning scheme of abstract thinking (based on RL), and also integrates it within the power of planning.

Reinforcement Learning RL is a popular technique used in recent development of LLMs [31, 13], to boost reasoning capabilities. A representative RL approach is PPO [30], typically used for learning from human feedback (RLHF) [7, 23] via a reward model pre-trained on human preference annotations. As a step forward, DPO [26] simplifies the PPO implementation by optimizing the direct preference of policy model, which gets rid of pre-training an additional reward model. Our method adopts GRPO [31] with model-free rewards (not relying on preference feedback), which further cuts off the training of value model (used for advantage estimation) by using a group relative advantage.

7 Conclusion

This paper proposes a method, AbstRaL, to promote abstract thinking in large language models (LLMs). AbstRaL is designed to improve the reasoning robustness of LLMs, based on the natural principle that abstract thinking leads to reasoning steps that are more invariant to surface-form variations. Our abstraction mechanism is implemented through a proper reinforcement learning (RL) framework, where model-free rewards are derived from newly designed rationales GranuLAR that blend socratic chain-of-thought with augmented granularity. This enables both the de-contextualization of problems and the integration of symbolic tools. We evaluate AbstRaL on recent GSM perturbation benchmarks and show that it effectively mitigates the performance degradation caused by instantiation and inferential shifts.

8 Theoretical open problems

It would be interesting to find a theoretical setting to study the impact of abstraction on the sample/time complexity of learning, the OOD robustness and the model size requirements for learning. For instance, formalizing data distributions of triplets (A, X, Y) , where $X \rightarrow Y$ represents a target of a reasoning problem X with solution Y , and A an abstraction of X that removes context/informalization (with the same solution Y). One has to define properly the notion of abstraction, distinguishing potentially de-contextualization from reasoning resolution (e.g., A should not be Y). Is there an appropriate framework (including a learning model class) under which statements of the following kind could be made rigorous: (i) the sample complexity of learning $A \rightarrow Y$ is lower than $X \rightarrow Y$ (a consequence of defining abstraction and instantiation processes formally and appropriately)? (ii) learning $X \rightarrow Y$ by relying on a proper abstraction A leads to improved robustness to properly modeled instantiation shifts? (iii) learning via abstraction can be achieved with models of smaller sizes?

Limitations

We acknowledge a few limitations in our work. First, the datasets used for testing our method cannot have exhaustive coverage of all real-world reasoning scenarios. We instead consider the most representative domain, *i.e.*, grade school mathematics, which is a common and typical domain for studying reasoning robustness [19, 24, 22], and use English as a primary language in our testing. For future work, the testbed of our method can be extended to more domains such as high-school competition mathematics [14, 33] and commonsense (or factual) reasoning [43, 24, 25], and to more languages. Second, our test of reasoning robustness is scoped to instantiation and inferential shifts, based on the assumption that the tested perturbations do not modify the underlying abstraction, *i.e.*, the abstract math derivations. Future work can extend our study to perturbations on the abstraction, such as altering “A is M years old, B is N years **older** than A, how old is B?” (M+N) to “A is M years old, B is N years **younger** than A, how old is B?” (M-N), to test the robustness of generalization to similar reasoning strategies. Furthermore, our method is tested on the setting of tuning the full LLMs, which requires considerable computing resources. More efficient model training schemes, such as LoRA [16], can be applied in future work. Lastly, all LLMs in our experiments use greedy decoding to generate inferences, which leaves room for future work to test our method on more advanced decoding strategies, such as self-consistency [35] decoding.

Acknowledgements

Antoine Bosselut gratefully acknowledges the support of the Swiss National Science Foundation (No. 215390), Innosuisse (PFFS-21-29), the EPFL Center for Imaging, Sony Group Corporation, and a Meta LLM Evaluation Research Grant.

A GRPO with Abstraction Rewards

Our proposed AbstRaL framework adopts GRPO [31] as the RL algorithm to train LLMs (on top of supervised fine-tuning) on the task of abstract reasoning, *i.e.*, generating the abstract answer \mathcal{Y}^A based on the input abstract question \mathcal{X}^A .

For each input question \mathcal{X}^A , GRPO samples a group of output answers $\{\widetilde{\mathcal{Y}}_1^A, \widetilde{\mathcal{Y}}_2^A, \dots, \widetilde{\mathcal{Y}}_G^A\}$ from the current (old policy) model $\pi_{\theta_{\text{old}}}$, and optimizes the (policy) model π_{θ} by maximizing the objective:

$$\frac{1}{G} \sum_{i=1}^G \left(\min \left(\frac{\pi_{\theta}(\widetilde{\mathcal{Y}}_i^A | \mathcal{X}^A)}{\pi_{\theta_{\text{old}}}(\widetilde{\mathcal{Y}}_i^A | \mathcal{X}^A)} R_i, \text{clip} \left(\frac{\pi_{\theta}(\widetilde{\mathcal{Y}}_i^A | \mathcal{X}^A)}{\pi_{\theta_{\text{old}}}(\widetilde{\mathcal{Y}}_i^A | \mathcal{X}^A)}, 1 - \varepsilon, 1 + \varepsilon \right) R_i \right) - \beta \mathbb{D}_{KL}(\pi_{\theta} || \pi_{\text{ref}}) \right) \quad (2)$$

where ε and β are hyperparameters. We set the reference policy π_{ref} as the model trained with only supervised fine-tuning (SFT), which is used for calculating the KL divergence regularization:

$$\mathbb{D}_{KL}(\pi_{\theta} || \pi_{\text{ref}}) = \frac{\pi_{\text{ref}}(\widetilde{\mathcal{Y}}_i^A | \mathcal{X}^A)}{\pi_{\theta}(\widetilde{\mathcal{Y}}_i^A | \mathcal{X}^A)} - \log \frac{\pi_{\text{ref}}(\widetilde{\mathcal{Y}}_i^A | \mathcal{X}^A)}{\pi_{\theta}(\widetilde{\mathcal{Y}}_i^A | \mathcal{X}^A)} - 1 \quad (3)$$

Table 4: Full evaluation results on GSM-Symbolic. Δ denotes the relative percentage of drop comparing performance on **Vary Both** to performance on **Origin 100**. Best results on each model are **bold**, where lower is better for Δ . Standard deviation (std) of multi-round evaluation results are in brackets, where lowest std values on each model are underlined. The best (bold) results with * are significantly better than their corresponding second-best results, with significant test p-value < 0.05 .

Model	Method	Vary Num.	Vary Name	Vary Both	Origin 100	Δ (%)
Llama-3.2-1B-Instruct	CoT-8S	0.3864 (0.030)	0.4250 (0.026)	0.3890 (0.038)	0.4800	18.96
	CoT-RL	0.4788 (0.031)	0.5342 (0.024)	0.4488 (0.038)	0.5700	21.26
	CoA	0.4534 (0.020)	0.4372 (0.026)	0.4242 (0.030)	0.4600	7.78
	AbstRaL	0.5912* (0.016)	0.5838* (0.023)	0.5804* (0.027)	0.6000*	3.27
Llama-3.2-3B-Instruct	CoT-8S	0.7218 (0.027)	0.7638 (0.027)	0.7114 (0.031)	0.8400*	15.31
	CoT-RL	0.7056 (0.027)	0.7516 (0.023)	0.6898 (0.026)	0.8000	13.78
	CoA	0.6450 (0.019)	0.6802 (0.021)	0.6760 (0.026)	0.6800	0.59
	AbstRaL	0.7960* (0.014)	0.7982* (0.020)	0.7946* (0.023)	0.8000	0.68
Llama-3.1-8B-Instruct	CoT-8S	0.8440 (0.026)	0.8746 (0.021)	0.8236 (0.032)	0.8700	5.33
	CoT-RL	0.7784 (0.029)	0.8710 (0.014)	0.7540 (0.026)	0.8300	9.16
	CoA	0.7240 (0.019)	0.7086 (0.019)	0.6944 (0.025)	0.7200	3.56
	AbstRaL	0.8686* (0.013)	0.8672 (0.018)	0.8620* (0.023)	0.8700	0.92
Qwen2.5-0.5B-Instruct	CoT-8S	0.3394 (0.039)	0.3724 (0.024)	0.3398 (0.033)	0.3800	10.58
	CoT-RL	0.3192 (0.025)	0.3948 (0.032)	0.3228 (0.032)	0.3500	7.77
	CoA	0.2866 (0.025)	0.3060 (0.026)	0.2872 (0.026)	0.2900	0.97
	AbstRaL	0.4396* (0.015)	0.4416* (0.026)	0.4456* (0.025)	0.4400*	-1.27
Qwen2.5-1.5B-Instruct	CoT-8S	0.5728 (0.032)	0.6416 (0.030)	0.5752 (0.033)	0.6600	12.85
	CoT-RL	0.5296 (0.037)	0.5830 (0.034)	0.5126 (0.034)	0.5600	8.46
	CoA	0.4680 (0.027)	0.4942 (0.029)	0.4656 (0.027)	0.5100	8.71
	AbstRaL	0.6444* (0.018)	0.6414 (0.028)	0.6416* (0.025)	0.6500	1.29
Qwen2.5-3B-Instruct	CoT-8S	0.7222 (0.037)	0.7820 (0.025)	0.7256 (0.027)	0.8200*	11.51
	CoT-RL	0.7150 (0.036)	0.7706 (0.025)	0.6888 (0.038)	0.7900	12.81
	CoA	0.6424 (0.030)	0.6134 (0.030)	0.6234 (0.034)	0.6500	4.09
	AbstRaL	0.7842* (0.014)	0.7852 (0.024)	0.7834* (0.024)	0.7900	0.84
Qwen2.5-7B-Instruct	CoT-8S	0.8726 (0.026)	0.9230 (0.016)	0.8740 (0.023)	0.9200*	5.00
	CoT-RL	0.7770 (0.033)	0.8170 (0.024)	0.7928 (0.034)	0.8500	6.73
	CoA	0.7310 (0.023)	0.7408 (0.027)	0.7414 (0.029)	0.7600	2.45
	AbstRaL	0.9022* (0.014)	0.9248 (0.017)	0.8834* (0.019)	0.8900	0.74
Qwen2.5-Math-7B-Instruct	CoT-8S	0.8956 (0.021)	0.9108 (0.018)	0.8766 (0.023)	0.9500	7.73
	CoT-RL	0.8942 (0.022)	0.9154 (0.012)	0.8812 (0.021)	0.9600	8.21
	CoA	0.7122 (0.028)	0.6976 (0.031)	0.6970 (0.033)	0.7100	1.83
	AbstRaL	0.9066* (0.015)	0.9014 (0.013)	0.9022* (0.016)	0.9100	0.86
Mathstral-7B-v0.1	CoT-8S	0.7876 (0.024)	0.8084 (0.018)	0.7604 (0.031)	0.8000	4.95
	CoT-RL	0.8082 (0.018)	0.7986 (0.021)	0.7688 (0.025)	0.7800	1.44
	CoA	0.7506 (0.031)	0.7740 (0.028)	0.7402 (0.027)	0.7500	1.31
	AbstRaL	0.8100* (0.012)	0.8214* (0.017)	0.8228* (0.019)	0.8100	-1.58

R_i is the group relative advantage granted to the abstraction $\tilde{\mathcal{A}}_i$ retrieved from each sampled output answer $\tilde{\mathcal{Y}}_i^{\mathcal{A}}$, which applies a group normalization on our proposed abstraction rewards r_{answer} and $r_{symbolic}$ defined in §3.2, with reference to the input conditions \mathcal{C} and gold answer Ans, and to the gold abstraction \mathcal{A} retrieved from the gold response $\mathcal{Y}^{\mathcal{A}}$, respectively:

$$R_i = \frac{r_i - \text{mean}(\{r_1, r_2, \dots, r_G\})}{\text{std}(\{r_1, r_2, \dots, r_G\})}, \quad r_i = r_{answer}(\tilde{\mathcal{A}}_i, \mathcal{C}, \text{Ans}) + r_{symbolic}(\tilde{\mathcal{A}}_i, \mathcal{A}) \quad (4)$$

B Full Experimental Results

Table 4 and 5 present the evaluation results of all our tested LLMs on GSM-Symbolic and GSM-Plus datasets, respectively. Results on all LLMs consistently demonstrate that AbstRaL effectively improves reasoning robustness when generalizing to both instantiation and interferential shifts. On each tested LLM, we also conduct the bootstrap statistical significant test [40] between the best and second-best results, and highlight the best results (with *) if they are significantly better than their corresponding second-best results with significant test p-value < 0.05 .

Table 5: Full evaluation results on GSM-Plus. Best results on each model are **bold**. The best (bold) results with * are significantly better than their corresponding second-best results, with significant test p-value < 0.05 .

Model	Method	Digit Ex.	Int-Dec-Fra	Num. Sub.	Rephrase	Distract	Original
Llama-3.2-1B-Instruct	CoT-8S	0.3889	0.2934	0.4321	0.5042	0.2790	0.4519
	CoT-RL	0.4064	0.3063	0.4572	0.5299	0.2646	0.5125
	CoA	0.3798	0.3237	0.3889	0.4139	0.2032	0.4261
	AbstRaL	0.5641*	0.5626*	0.5641*	0.5633*	0.4359*	0.5641*
Llama-3.2-3B-Instruct	CoT-8S	0.7142	0.6277	0.7172	0.7877	0.6073	0.7953
	CoT-RL	0.6808	0.5406	0.7043	0.7718	0.5398	0.7763
	CoA	0.6224	0.5444	0.6126	0.6793	0.4177	0.6626
	AbstRaL	0.7968*	0.7945*	0.7968*	0.7923	0.6755*	0.7968
Llama-3.1-8B-Instruct	CoT-8S	0.7938	0.7445	0.7854	0.8461	0.7400	0.8567
	CoT-RL	0.7240	0.5914	0.7187	0.8309	0.5466	0.8234
	CoA	0.7202	0.6398	0.7035	0.7657	0.5344	0.7497
	AbstRaL	0.8506*	0.8476*	0.8506*	0.8514	0.7854*	0.8506
Qwen2.5-0.5B-Instruct	CoT-8S	0.3601	0.2866	0.3958	0.4359	0.2267	0.4238
	CoT-RL	0.3336	0.2373	0.3571	0.4079	0.1524	0.3798
	CoA	0.2745	0.2267	0.2782	0.3161	0.1266	0.3033
	AbstRaL	0.4670*	0.4663*	0.4670*	0.4625*	0.3654*	0.4670*
Qwen2.5-1.5B-Instruct	CoT-8S	0.6096	0.5421	0.6194	0.6793	0.4488	0.6702
	CoT-RL	0.5527	0.4632	0.5754	0.6520	0.3950	0.6202
	CoA	0.4602	0.3700	0.4511	0.5186	0.2631	0.5019
	AbstRaL	0.6778*	0.6763*	0.6778*	0.6755	0.5777*	0.6778
Qwen2.5-3B-Instruct	CoT-8S	0.7726	0.7074	0.7430	0.8218	0.6346	0.8120
	CoT-RL	0.7149	0.6353	0.7331	0.7817	0.5277	0.7726
	CoA	0.6232	0.5497	0.5989	0.6679	0.4200	0.6702
	AbstRaL	0.8158*	0.8128*	0.8158*	0.8249	0.7036*	0.8158
Qwen2.5-7B-Instruct	CoT-8S	0.8400	0.8067	0.8324	0.8779	0.7938	0.8901
	CoT-RL	0.7582	0.6854	0.7369	0.8097	0.6331	0.8036
	CoA	0.7195	0.6528	0.6755	0.7657	0.5588	0.7597
	AbstRaL	0.8825*	0.8795*	0.8825*	0.8870*	0.7953	0.8825
Qwen2.5-Math-7B-Instruct	CoT-8S	0.8552	0.8218	0.8453	0.9052	0.7627	0.9181
	CoT-RL	0.8757	0.8522	0.8506	0.9037	0.8150	0.9340*
	CoA	0.7460	0.6907	0.7180	0.7642	0.5709	0.7809
	AbstRaL	0.8916*	0.8886*	0.8916*	0.8992	0.8226	0.8916
Mathstral-7B-v0.1	CoT-8S	0.7544	0.6892	0.7604	0.8029	0.6808	0.8074
	CoT-RL	0.7665	0.7111	0.7331	0.8089	0.5542	0.7953
	CoA	0.7149	0.6422	0.6831	0.7483	0.5246	0.7619
	AbstRaL	0.8241*	0.8211*	0.8241*	0.8247*	0.7657*	0.8241*

C Qualitative Analysis

Table 6 and 7 present two mathematical reasoning examples (respectively on GSM-Symbolic and GSM-Plus datasets) of our tested strongest LLM Qwen2.5-Math-7B-Instruct, using either baseline CoT-8S method or our AbstRaL. The LLM with AbstRaL performs more stable math derivations when facing the variations of relevant conditions (Table 6), and achieves more robust reasoning chains when dealing with the inserted distracting condition (Table 7). We include more detailed analysis in the corresponding table captions.

D AbstRaL Implementation Details

This section includes more details of how we implement our AbstRaL framework.

For the first condition recognition step of AbstRaL, we prompt a Llama-3.3-70B-Instruct [12] model to replace the numerical values in the input question \mathcal{X} with abstract symbols enclosed in square brackets, to create the abstract question \mathcal{X}^A . A set of conditions \mathcal{C} is also created to record the replacements as a list of equations. Table 8 presents the instruction and few-shot examples used as the task demonstration in our prompting. 4 NVIDIA A100-SXM4 (80GB) GPUs were used for running the condition recognition step based on Llama-3.3-70B-Instruct, which took about 36 hours to process all training and testing data samples.

Based on the abstract question \mathcal{X}^A and the gold socratic CoT response \mathcal{Y} to the question, we also prompt Llama-3.3-70B-Instruct to rewrite \mathcal{Y} into our granularly-decomposed abstract reasoning (GranularAR) format \mathcal{Y}^A , which is used as the training data for LLMs to learn the abstract reasoning step in AbstRaL. To facilitate the rewriting, we employ a two-step pipeline. First, the Llama model is prompted to replace the numerical values in \mathcal{Y} with abstract symbols, by either quoting the abstract symbols in \mathcal{X}^A for input values or creating new abstract symbols for derived output values. All derivations in the response are supposed to be enclosed in double angle brackets “< < > >”. Second, based on the rewritten response in the first step, the Llama model is prompted to further rewrite the response (in socratic CoT format) into our GranularAR format, while keep all the abstract symbols in the response unchanged. Table 9 and 10 present the instruction and few-shot examples used for prompting our two-step GranularAR training data construction. For each step of rewriting, 4 NVIDIA A100-SXM4 (80GB) GPUs were used to run the Llama model. About 36 and 48 hours are spent to conduct the first and the second rewriting steps, respectively, on all training data samples.

After each step of response rewriting, we filter the output of Llama model by verifying the correctness of its derivations, *i.e.*, we use regex-matching to extract all derivations enclosed in “< < > >”, and pass them (along with the input conditions \mathcal{C} generated in the condition recognition step) into the SymPy¹⁴ solver, to derive the final answer (number), checking whether it matches the gold answer. We apply our two-step rewriting to the socratic version of GSM8K [8] training set from OpenAI¹⁵. After the first step of rewriting and filtering, 6503 out of 7473 training samples are correctly rewritten and kept, while after the second step of rewriting and filtering, 6386 out of 6503 training samples are correctly rewritten and reserved as our final training data.

In the abstract reasoning step of our AbstRaL framework, LLMs are trained with SFT and RL on the task of generating the abstract (GranularAR) answer \mathcal{Y}^A based on the input abstract question \mathcal{X}^A .

For SFT, we set the batch size as 8, using 4 NVIDIA A100-SXM4 (80GB) GPUs (*i.e.*, batch size is 2 on each GPU), and set the learning rate as $5e^{-6}$, using AdamW [20] optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\epsilon = 1e^{-8}$. All LLMs were trained with SFT for 2 epochs on our GranularAR data, which took less than 1 hour.

For RL, we set the positive (correct) reward $r_{correct} = 2.5$ in our answer correctness reward $r_{answer}(\tilde{\mathcal{A}}_i, \mathcal{C}, \text{Ans})$, and the maximum reward $r_{max} = 1.5$ in our symbolic distance reward $r_{symbolic}(\tilde{\mathcal{A}}_i, \mathcal{A})$. The hyperparameters of GRPO algorithm are set to $\beta = 0.04$ (KL coefficient) and $\epsilon = 0.2$ (for clipping), with number of generations (*i.e.*, $\{\tilde{\mathcal{Y}}_1^A, \tilde{\mathcal{Y}}_2^A, \dots, \tilde{\mathcal{Y}}_G^A\}$) sampled per group (*i.e.*, per given \mathcal{X}^A) set as $G = 16$, and with the temperature, top_p and top_k of the sampling set as 0.9, 1.0 and 50, respectively. The learning rate of our RL optimization is set to $5e^{-7}$, using also AdamW optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.99$ and $\epsilon = 1e^{-8}$. A cosine learning rate scheduler is employed with warm-up ratio set to 0.1, and the training weight decay and gradient norm clip value are also both set to 0.1. We use 8 NVIDIA A100-SXM4 (80GB) GPUs to train each LLM with RL. For 7B and 8B LLMs, the batch size of generation is set to 2 per GPU, with gradient accumulation steps set as 4, *i.e.*, the policy gradient update is performed when every 4 groups of $G = 16$ generations (2 batch size multiply 8 GPUs) are sampled. For other smaller LLMs, the batch size of generation is set to 8 per GPU instead, with gradient accumulation step set as 1, where still 4 groups of $G = 16$ generations (8 batch size multiply 8 GPUs) are sampled per policy gradient update. All LLMs were trained with RL for 8 epochs on our GranularAR data, which took about 3 to 5 days.

E CoA Implementation Details

For training data construction of the baseline method CoA [10], we prompt the same Llama-3.3-70B-Instruct model to rewrite the gold socratic CoT response \mathcal{Y} (the original input question \mathcal{X} is also given), where only output numbers in \mathcal{Y} are replaced with abstract symbols. Table 11 presents the instruction and few-shot examples used for CoA training data construction. Following the CoA method, LLMs are trained solely with SFT to learn the generation of CoA abstract answer based on the input question \mathcal{X} , with the same hyperparameters used in AbstRaL’s SFT.¹⁶ The same symbolic

¹⁴<https://github.com/sympy/sympy>

¹⁵<https://huggingface.co/datasets/openai/gsm8k/tree/main/socratic>

¹⁶We also tested the original hyperparameters suggested by CoA, which however got lower evaluation scores.

tools (regex-matcher and SymPy solver) are used to extract the abstract derivations and calculate the final answer number, for either filtering the rewritten training data or deriving the answer at inference.

References

- [1] Emmanuel Abbe, Samy Bengio, Aryo Lotfi, Colin Sandon, and Omid Saremi. How far can transformers reason? the globality barrier and inductive scratchpad. *Advances in Neural Information Processing Systems*, 37:27850–27895, 2024.
- [2] Qiming Bao, Gael Gendron, Alex Yuxuan Peng, Wanjun Zhong, Neset Tan, Yang Chen, Michael Witbrock, and Jiamou Liu. Assessing and enhancing the robustness of large language models with task structure variations for logical reasoning. *arXiv preprint arXiv:2310.09430*, 2023.
- [3] Enric Boix-Adserà, Omid Saremi, Emmanuel Abbe, Samy Bengio, Etai Littwin, and Joshua M. Susskind. When can transformers reason with abstract symbols? In *The Twelfth International Conference on Learning Representations*, 2024.
- [4] Guanyu Chen, Peiyang Wang, Tianren Zhang, and Feng Chen. Exploring the hidden reasoning process of large language models by misleading them. *arXiv preprint arXiv:2503.16401*, 2025.
- [5] Qiguang Chen, Libo Qin, Jinhao Liu, Dengyun Peng, Jiannan Guan, Peng Wang, Mengkang Hu, Yuhang Zhou, Te Gao, and Wanxiang Che. Towards reasoning era: A survey of long chain-of-thought for reasoning large language models. *arXiv preprint arXiv:2503.09567*, 2025.
- [6] François Chollet. On the measure of intelligence. *arXiv preprint arXiv:1911.01547*, 2019.
- [7] Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.
- [8] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- [9] Wee Chung Gan and Hwee Tou Ng. Improving the robustness of question answering systems to question paraphrasing. In *Proceedings of the 57th annual meeting of the association for computational linguistics*, pages 6065–6075, 2019.
- [10] Silin Gao, Jane Dwivedi-Yu, Ping Yu, Xiaoqing Ellen Tan, Ramakanth Pasunuru, Olga Golovneva, Koustuv Sinha, Asli Celikyilmaz, Antoine Bosselut, and Tianlu Wang. Efficient tool use with chain-of-abstraction reasoning. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 2727–2743, 2025.
- [11] Gaël Gendron, Qiming Bao, Michael Witbrock, and Gillian Dobbie. Large language models are not strong abstract reasoners. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence*, pages 6270–6278, 2024.
- [12] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- [13] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shitong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- [14] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*, 2021.
- [15] Ruixin Hong, Hongming Zhang, Xiaoman Pan, Dong Yu, and Changshui Zhang. Abstraction-of-thought makes language models better reasoners. *arXiv preprint arXiv:2406.12442*, 2024.

- [16] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- [17] Komal Kumar, Tajamul Ashraf, Omkar Thawakar, Rao Muhammad Anwer, Hisham Cholakkal, Mubarak Shah, Ming-Hsuan Yang, Phillip HS Torr, Fahad Shahbaz Khan, and Salman Khan. Llm post-training: A deep dive into reasoning large language models. *arXiv preprint arXiv:2502.21321*, 2025.
- [18] Vladimir I Levenshtein et al. Binary codes capable of correcting deletions, insertions, and reversals. In *Soviet physics doklady*, volume 10, pages 707–710. Soviet Union, 1966.
- [19] Qintong Li, Leyang Cui, Xueliang Zhao, Lingpeng Kong, and Wei Bi. Gsm-plus: A comprehensive benchmark for evaluating the robustness of llms as mathematical problem solvers. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2961–2984, 2024.
- [20] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2018.
- [21] Inbal Magar and Roy Schwartz. Data contamination: From memorization to exploitation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 157–165, 2022.
- [22] Seyed Iman Mirzadeh, Keivan Alizadeh, Hooman Shahrokhi, Oncel Tuzel, Samy Bengio, and Mehrdad Farajtabar. GSM-symbolic: Understanding the limitations of mathematical reasoning in large language models. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [23] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- [24] Kun Qian, Shunji Wan, Claudia Tang, Youzhi Wang, Xuanming Zhang, Maximillian Chen, and Zhou Yu. Varbench: Robust language model benchmarking through dynamic variable perturbation. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 16131–16161, 2024.
- [25] Yao Qiang, Subhrangshu Nandi, Ninareh Mehrabi, Greg Ver Steeg, Anoop Kumar, Anna Rumshisky, and Aram Galstyan. Prompt perturbation consistency learning for robust language models. In *Findings of the Association for Computational Linguistics: EACL 2024*, pages 1357–1370, 2024.
- [26] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36:53728–53741, 2023.
- [27] Sylvestre-Alvise Rebuffi, Sven Gowal, Dan A Calian, Florian Stimberg, Olivia Wiles, and Timothy Mann. Fixing data augmentation to improve adversarial robustness. *arXiv preprint arXiv:2103.01946*, 2021.
- [28] Sylvestre-Alvise Rebuffi, Sven Gowal, Dan Andrei Calian, Florian Stimberg, Olivia Wiles, and Timothy A Mann. Data augmentation can improve robustness. *Advances in neural information processing systems*, 34:29935–29948, 2021.
- [29] Lorenza Saitta, Jean-Daniel Zucker, Lorenza Saitta, and Jean-Daniel Zucker. *Abstraction in Artificial Intelligence*. Springer, 2013.
- [30] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

- [31] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- [32] Kumar Shridhar, Alessandro Stolfo, and Mrinmaya Sachan. Distilling reasoning capabilities into smaller language models. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 7059–7073, 2023.
- [33] Saurabh Srivastava, Anto PV, Shashank Menon, Ajay Sukumar, Alan Philipose, Stevin Prince, Sooraj Thomas, et al. Functional benchmarks for robust evaluation of reasoning performance, and the reasoning gap. *arXiv preprint arXiv:2402.19450*, 2024.
- [34] Sowmya S Sundaram, Sairam Gurajada, Deepak Padmanabhan, Savitha Sam Abraham, and Marco Fisichella. Does a language model “understand” high school math? a survey of deep learning based word problem solvers. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 14(4):e1534, 2024.
- [35] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations*, 2022.
- [36] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- [37] Robert Wilensky. Planning and understanding: A computational approach to human reasoning. 1983.
- [38] Cheng Xu, Shuhao Guan, Derek Greene, M Kechadi, et al. Benchmark data contamination of large language models: A survey. *arXiv preprint arXiv:2406.04244*, 2024.
- [39] Fangzhi Xu, Qika Lin, Jiawei Han, Tianzhe Zhao, Jun Liu, and Erik Cambria. Are large language models really good logical reasoners? a comprehensive evaluation and beyond. *IEEE Transactions on Knowledge and Data Engineering*, 2025.
- [40] Kuan-Man Xu. Using the bootstrap method for a statistical significance test of differences between summary histograms. *Monthly weather review*, 134(5):1442–1453, 2006.
- [41] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- [42] An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, et al. Qwen2.5-math technical report: Toward mathematical expert model via self-improvement. *arXiv preprint arXiv:2409.12122*, 2024.
- [43] Shuo Yang, Wei-Lin Chiang, Lianmin Zheng, Joseph E Gonzalez, and Ion Stoica. Rethinking benchmark and contamination for language models with rephrased samples. *arXiv preprint arXiv:2311.04850*, 2023.
- [44] Sohee Yang, Elena Gribovskaya, Nora Kassner, Mor Geva, and Sebastian Riedel. Do large language models latently perform multi-hop reasoning? In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10210–10229, 2024.
- [45] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Advances in neural information processing systems*, 36:11809–11822, 2023.
- [46] Fei Yu, Hongbo Zhang, Prayag Tiwari, and Benyou Wang. Natural language reasoning, a survey. *ACM Computing Surveys*, 56(12):1–39, 2024.
- [47] Tong Yu, Yongcheng Jing, Xikun Zhang, Wentao Jiang, Wenjie Wu, Yingjie Wang, Wenbin Hu, Bo Du, and Dacheng Tao. Benchmarking reasoning robustness in large language models. *arXiv preprint arXiv:2503.04550*, 2025.

- [48] Tianyang Zhong, Zhengliang Liu, Yi Pan, Yutong Zhang, Yifan Zhou, Shizhe Liang, Zihao Wu, Yanjun Lyu, Peng Shu, Xiaowei Yu, et al. Evaluation of openai o1: Opportunities and challenges of agi. *arXiv preprint arXiv:2409.18486*, 2024.
- [49] Yue Zhou, Yada Zhu, Diego Antognini, Yoon Kim, and Yang Zhang. Paraphrase and solve: Exploring and exploiting the impact of surface form on mathematical reasoning in large language models. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 2793–2804, 2024.

Table 6: Mathematical reasoning examples on GSM-Symbolic dataset. Responses are generated by Qwen2.5-Math-7B-Instruct with either CoT-8S or AbstRaL methods. Input queries include an original question (from Origin 100) and its corresponding variation (from Vary Both) with names (or entities) and numbers changed to other values, as illustrated by the orange words. For the Origin 100 question, both CoT-8S and AbstRaL achieve correct math derivations and answers, which are colored in blue and brown, respectively. However, when transferring to the Vary Both question, CoT-8S turns to conduct a wrong calculation on the altered numbers (7×39), and thus derives a wrong answer, as shown by the red numbers. By contrast, AbstRaL keeps correct math derivations, due to using abstract symbols (with a symbolic solver) that are invariant to changes of input numbers.

GSM-Symbolic (Origin 100): Jaime places eggs on some trays. Each tray can hold 24 eggs. If he has 64 eggs and 2 trays, how many eggs won't he be able to place on the tray?

Gold Response: Jaime will be able to place a total of $24 \times 2 = 48$ eggs. Therefore, there are $64 - 48 = 16$ eggs that he won't be able to place on the tray. The answer is 16.

Response with CoT-8S: Let's think step by step. Each tray can hold 24 eggs, and Jaime has 2 trays, so the total capacity of the trays is $24 \times 2 = 48$ eggs. Jaime has 64 eggs, but only 48 can fit on the trays. Therefore, the number of eggs that won't be able to be placed on the tray is $64 - 48 = 16$. The final answer is 16.

Response with AbstRaL:

(Abstract Question) Jaime places eggs on the tray. Each tray can hold [in0] eggs. If he has [in1] eggs and [in2] trays, how many eggs won't he be able to place on the tray?

(Conditions) in0=24 in1=64 in2=2

(Abstract Answer) Let's think about the sub-questions we need to answer.

Q1: How many eggs can Jaime place on the trays?

Q2: How many eggs won't Jaime be able to place on the tray?

Let's answer each sub-question one by one.

Q1: How many eggs can Jaime place on the trays? Each tray can hold [in0] eggs. Jaime has [in2] trays.

So, the total number of eggs that can be placed on the trays is $< < \text{in0} * \text{in2} = \text{out0} > >$ eggs.

Q2: How many eggs won't Jaime be able to place on the tray? Jaime has [in1] eggs. The total number of eggs that can be placed on the trays is [out0] eggs. So, the number of eggs that won't be able to be placed on the tray is $< < \text{in1} - \text{out0} = \text{out1} > >$ eggs. The final answer is out1.

(Abstraction) in0*in2=out0 in1-out0=out1

(Derived Final Answer) out0=48 out1=16

GSM-Symbolic (Vary Both): Arash places crackers on some trays. Each tray can hold 39 crackers. If he has 302 crackers and 7 trays, how many crackers won't he be able to place on the tray?

Gold Response: Arash will be able to place a total of $39 \times 7 = 273$ crackers. Therefore, there are $302 - 273 = 29$ crackers that he won't be able to place on the tray. The answer is 29.

Response with CoT-8S: First, we need to determine how many crackers Arash can place on the 7 trays.

Since each tray can hold 39 crackers, the total number of crackers that can be placed on 7 trays is (7×39) .

Calculating (7×39) : $[7 \times 39 = 7 \times (40 - 1) = 7 \times 40 - 7 \times 1 = 280 - 7 = 280]$.

So, Arash can place 280 crackers on the 7 trays. He has 302 crackers in total. To find out how many crackers he won't be able to place, we subtract the number of crackers he can place from the total number of crackers: $[302 - 280 = 22]$. The final answer is 22. Arash won't be able to place 22 crackers on the trays.

Response with AbstRaL:

(Abstract Question) Arash places crackers on some trays. Each tray can hold [in0] crackers.

If he has [in1] crackers and [in2] trays, how many crackers won't he be able to place on the tray?

(Conditions) in0=39 in1=302 in2=7

(Abstract Answer) Let's think about the sub-questions we need to answer.

Q1: How many crackers can Arash place on the trays?

Q2: How many crackers won't Arash be able to place on the tray?

Let's answer each sub-question one by one.

Q1: How many crackers can Arash place on the trays? Each tray can hold [in0] crackers. Arash has [in2] trays.

So, the total number of crackers that can be placed on the trays is $< < \text{in0} * \text{in2} = \text{out0} > >$ crackers.

Q2: How many crackers won't Arash be able to place on the tray? Arash has [in1] crackers. The total number of crackers that can be placed on the trays is [out0]. So, the number of crackers that won't be able to be placed on the tray is $< < \text{in1} - \text{out0} = \text{out1} > >$ crackers. The final answer is out1.

(Abstraction) in0*in2=out0 in1-out0=out1

(Derived Final Answer) out0=273 out1=29

Table 7: Mathematical reasoning examples on GSM-Plus dataset. Responses are generated by Qwen2.5-Math-7B-Instruct with either CoT-8S or AbstRaL methods. Input queries include an original question (from Original) and its corresponding variation (from Distract) where a distracting condition is inserted into the question, as illustrated by the **violet** sentence. For the Original question, both CoT-8S and AbstRaL achieve correct math derivations and answers, which are colored in **blue** and **brown**, respectively. However, when transferring to the Distract question, CoT-8S turns to adopt an incorrect reasoning chain, where the distracting condition is wrongly incorporated into the answer derivation, as shown by the **red** sentences. By contrast, AbstRaL handles the distracting condition in an additional sub-question (Q2), whose math derivation (colored in **gray**) does not interfere the derivations of the correct answer (colored in **brown**). This illustrates that the holistic planning of problem solving, which is used in our Granular training data format, can effectively mitigate the interference of distractors, and thus increase reasoning robustness to adversarial distribution shifts.

GSM-Plus (Original): In a race with 240 Asians, 80 were Japanese, and the rest were Chinese. If the number of boys on the Chinese team was 60, how many girls were on the Chinese team?

Gold Response: If there were 80 Japanese, the number of Chinese was $240-80=160$. Since there were 60 boys in the Chinese team, there were $160-60=100$ girls in the Chinese team. The answer is 100.

Response with CoT-8S: Let's think step by step. First, we need to determine the number of Chinese participants in the race. We know there are 240 Asians in total and 80 of them are Japanese. Therefore, the number of Chinese participants is: $[240-80=160]$. Next, we know that out of these 160 Chinese participants, 60 are boys. To find the number of girls, we subtract the number of boys from the total number of Chinese participants: $[160-60=100]$. So, the number of girls on the Chinese team is 100. The final answer is 100.

Response with AbstRaL:

(Abstract Question) In a race with $[in0]$ Asians, $[in1]$ were Japanese, and the rest were Chinese. If the number of boys on the Chinese team was $[in2]$, how many girls were on the Chinese team?

(Conditions) $in0=240$ $in1=80$ $in2=60$

(Abstract Answer) Let's think about the sub-questions we need to answer. Q1: How many Chinese were there?

Q2: How many girls were on the Chinese team? Let's answer each sub-question one by one.

Q1: How many Chinese were there? The total number of Asians in the race is $[in0]$. There were $[in1]$ Japanese among them. So, the number of Chinese was $<< in0-in1=out0 >>$.

Q2: How many girls were on the Chinese team? The total number of Chinese in the race is $[out0]$. The number of boys on the Chinese team was $[in2]$. So, the number of girls on the Chinese team was $<< out0-in2=out1 >>$. The final answer is $out1$.

(Abstraction) $in0-in1=out0$ $out0-in2=out1$

(Derived Final Answer) $out0=160$ $out1=100$

GSM-Plus (Distract): In a race with 240 Asians, 80 were Japanese, and the rest were Chinese. **Among the Japanese, there are 30 girls and 50 boys.** If the number of boys on the Chinese team was 60, how many girls were on the Chinese team?

Gold Response: If there were 80 Japanese, the number of Chinese was $240-80=160$. Since there were 60 boys in the Chinese team, there were $160-60=100$ girls in the Chinese team. The answer is 100.

Response with CoT-8S: Let's break this down step by step. 1.**Total number of Asians in the race**: 240. 2.**Number of Japanese in the race**: 80. 3.**Number of Chinese in the race**: $[240-80=160]$. 4.**Number of boys among the Japanese**: 50. 5.**Number of girls among the Japanese**: 30. 6.**Number of boys on the Chinese team**: 60. 7.**Number of girls on the Chinese team**: Since the total number of Chinese participants is 160, and we know the number of boys (60) and girls (30) among the Japanese, **we can find the number of girls on the Chinese team by subtracting the number of boys and girls among the Japanese from the total number of Chinese.** First, calculate the total number of girls among the Japanese: $[30 \text{ (girls among Japanese)} + 0 \text{ (girls among Chinese, since we don't have this information yet)} = 30]$. Now, calculate the total number of boys among the Japanese and Chinese $[50 \text{ (boys among Japanese)} + 60 \text{ (boys among Chinese)} = 110]$. **Subtract the number of boys and girls among the Japanese from the total number of Chinese: $[160-110=50]$.** Therefore, the number of girls on the Chinese team is 50. The final answer is 50.

Response with AbstRaL:

(Abstract Question) In a race with $[in0]$ Asians, $[in1]$ were Japanese, and the rest were Chinese. Among the Japanese, there are $[in2]$ girls and $[in3]$ boys. If the number of boys on the Chinese team was $[in4]$, how many girls were on the Chinese team?

(Conditions) $in0=240$ $in1=80$ $in2=30$ $in3=50$ $in4=60$

(Abstract Answer) Let's think about the sub-questions we need to answer. Q1: How many Chinese were there? Q2: How many boys were there among the Chinese? Q3: How many girls were there among the Chinese? Let's answer each sub-question one by one.

Q1: How many Chinese were there? There were $[in0]$ Asians in total. $[in1]$ of them were Japanese.

So, the number of Chinese is $<< in0-in1=out0 >>$.

Q2: How many boys were there among the Chinese? The total number of boys among the Japanese is $[in2] + [in3] = << in2+in3=out1 >>$. The total number of boys among the Chinese is $[in4]$. So, the total number of boys among the Chinese is $[in4]$.

Q3: How many girls were there among the Chinese? The total number of Chinese is $[out0]$. The total number of boys among the Chinese is $[in4]$. So, the total number of girls among the Chinese is $<< out0-in4=out2 >>$. The final answer is $out2$.

(Abstraction) $in0-in1=out0$ $in2+in3=out1$ $out0-in4=out2$

(Derived Final Answer) $out0=160$ $out1=80$ $out2=100$

Table 8: Instructions and few-shot examples used for prompting condition recognition in AbstRaL.

You are given a math word problem. Please replace the numerical values in the problem with abstract symbols, and enclose the abstract symbols in square brackets. Do not change any other words or add any other contents to the problem. Please also write equations to clarify which numerical value each abstract symbol is assigned to.

Example 1:

Input problem: Natalia sold clips to 48 of her friends in April, and then she sold half as many clips in May.

How many clips did Natalia sell altogether in April and May?

Output problem: Natalia sold clips to [in0] of her friends in April, and then she sold [in1] as many clips in May.

How many clips did Natalia sell altogether in April and May?

Equations: $in0=48$ $in1=1/2$

Example 2:

Input problem: The flowers cost \$9, the clay pot costs \$20 more than the flower, and the bag of soil costs \$2 less than the flower. How much does it cost to plant the flowers?

Output problem: The flowers cost \$[in0], the clay pot costs \$[in1] more than the flower, and the bag of soil costs \$[in2] less than the flower. How much does it cost to plant the flowers?

Equations: $in0=9$ $in1=20$ $in2=2$

Example 3:

Input problem: From March to August, Sam made \$460 doing 23 hours of yard work. However, from September to February, Sam was only able to work for 8 hours. If Sam is saving up to buy a video game console that costs \$600 and has already spent \$340 to fix his car, how many more hours does he need to work before he can buy the video game console?

Output problem: From March to August, Sam made \$[in0] doing [in1] hours of yard work. However, from September to February, Sam was only able to work for [in2] hours. If Sam is saving up to buy a video game console that costs \$[in3] and has already spent \$[in4] to fix his car, how many more hours does he need to work before he can buy the video game console?

Equations: $in0=460$ $in1=23$ $in2=8$ $in3=600$ $in4=340$

Example 4:

Input problem: Zhang is twice as old as Li. Li is 12 years old. Zhang’s brother Jung is 2 years older than Zhang. How old is Jung?

Output problem: Zhang is [in0] times as old as Li. Li is [in1] years old. Zhang’s brother Jung is [in2] years older than Zhang. How old is Jung?

Equations: $in0=2$ $in1=12$ $in2=2$

Example 5:

Input problem: Of the 90 people on William’s bus, $3/5$ were Dutch. Of the $1/2$ of the Dutch who were also American, $1/3$ got window seats. What’s the number of Dutch Americans who sat at the windows?

Output problem: Of the [in0] people on William’s bus, [in1] were Dutch. Of the [in2] of the Dutch who were also American, [in3] got window seats. What’s the number of Dutch Americans who sat at the windows?

Equations: $in0=90$ $in1=3/5$ $in2=1/2$ $in3=1/3$

Table 9: Instructions and few-shot examples for the first response rewriting step of our GranuLAR training data construction, used for learning abstract reasoning in AbstRaL.

You are given a math word problem with input conditions and solution. Please rewrite the solution by replacing the numerical values in angle brackets with abstract symbols. If the numerical values are given in the conditions, replace them with the abstract symbols assigned to them in the square brackets, otherwise replace them with new abstract symbols. Please also remove the redundant calculations around the angle brackets. Do not add any other contents to the solution.

Example 1:

Problem: Natalia sold clips to [in0] of her friends in April, and then she sold [in1] as many clips in May.

How many clips did Natalia sell altogether in April and May?

Conditions: in0=48 in1=1/2

Solution: How many clips did Natalia sell in May? ** Natalia sold $48/2 = << 48/2=24 >> 24$ clips in May.

How many clips did Natalia sell altogether in April and May? ** Natalia sold $48+24 = << 48+24=72 >> 72$ clips altogether in April and May. The answer is 72.

Rewrite solution: How many clips did Natalia sell in May? ** Natalia sold $<< in0*in1=out0 >>$ clips in May.

How many clips did Natalia sell altogether in April and May? ** Natalia sold $<< in0+out0=out1 >>$ clips altogether in April and May. The answer is out1.

Example 2:

Problem: The flowers cost \$[in0], the clay pot costs \$[in1] more than the flower, and the bag of soil costs \$[in2] less than the flower. How much does it cost to plant the flowers?

Conditions: in0=9 in1=20 in2=2

Solution: How much does the clay pot cost? ** The clay pot costs $\$20 + \$9 = \$<< 20+9=29 >> 29$.

How much does the bag of soil cost? ** The bag of soil costs $\$9 - \$2 = \$<< 9-2=7 >> 7$.

How much does it cost to plant the flowers? ** The cost to plant the flowers is $\$9 + \$29 + \$7 = \$<< 9+29+7=45 >> 45$. The answer is 45.

Rewrite solution: How much does the clay pot cost? ** The clay pot costs $\$<< in1+in0=out0 >>$.

How much does the bag of soil cost? ** The bag of soil costs $\$<< in0-in2=out1 >>$.

How much does it cost to plant the flowers? ** The cost to plant the flowers is $\$<< in0+out0+out1=out2 >>$. The answer is out2.

Example 3:

Problem: From March to August, Sam made \$[in0] doing [in1] hours of yard work. However, from September to February, Sam was only able to work for [in2] hours. If Sam is saving up to buy a video game console that costs \$[in3] and has already spent \$[in4] to fix his car, how many more hours does he need to work before he can buy the video game console?

Conditions: in0=460 in1=23 in2=8 in3=600 in4=340

Solution: How much does Sam make per hour? ** Sam makes $\$460 / 23 \text{ hrs} = \$<< 460/23=20 >> 20/\text{hr}$. How much did Sam make from September to February? ** From September to February, Sam made $8 \text{ hrs} \times \$20/\text{hr} = \$<< 8*20=160 >> 160$.

How much did Sam make from March to February? ** From March to February, Sam made a total of $\$460 + \$160 = \$620$.

How much money did Sam have after fixing his car? ** After fixing his car, he was left with $\$620 - \$340 = \$<< 620-340=280 >> 280$.

How much money does Sam need to buy the video game console? ** Sam needs another $\$600 - \$280 = \$<< 600-280=320 >> 320$.

How many more hours does Sam need to work? ** Sam needs to work another $\$320 / \$20/\text{hr} = << 320/20=16 >> 16$ hours.

The answer is 16.

Rewrite solution: How much does Sam make per hour? ** Sam makes $\$<< in0/in1=out0 >>/\text{hr}$. How much did Sam make from September to February? ** From September to February, Sam made $\$<< in2*out0=out1 >>$.

How much did Sam make from March to February? ** From March to February, Sam made a total of $\$<< in0+out1=out2 >>$.

How much money did Sam have after fixing his car? ** After fixing his car, he was left with $\$<< out2-in4=out3 >>$.

How much money does Sam need to buy the video game console? ** Sam needs another $\$<< in3-out3=out4 >>$.

How many more hours does Sam need to work? ** Sam needs to work another $<< out4/out0=out5 >>$ hours.

The answer is out5.

Example 4:

Problem: Zhang is [in0] times as old as Li. Li is [in1] years old. Zhang's brother Jung is [in2] years older than Zhang. How old is Jung?

Conditions: in0=2 in1=12 in2=2

Solution: How old is Zhang? ** Zhang is $2 * 12 \text{ years old} = << 2*12=24 >> 24$ years old.

How old is Jung? ** Jung is $2 \text{ years} + 24 \text{ years} = << 2+24=26 >> 26$ years old.

The answer is 26.

Rewrite solution: How old is Zhang? ** Zhang is $<< in0*in1=out0 >>$ years old.

How old is Jung? ** Jung is $<< in2+out0=out1 >>$ years old.

The answer is out1.

Example 5:

Problem: Of the [in0] people on William's bus, [in1] were Dutch. Of the [in2] of the Dutch who were also American, [in3] got window seats. What's the number of Dutch Americans who sat at the windows?

Conditions: in0=90 in1=3/5 in2=1/2 in3=1/3

Solution: How many Dutch people were on the bus? ** On the bus, the number of Dutch people was $3/5$ of the total number, a total of $3/5*90 = << 3/5*90=54 >> 54$ people.

How many Dutch Americans were on the bus? ** Out of the 54 people who were Dutch, $1/2$ were Dutch Americans, a total of $1/2*54 = << 1/2*54=27 >> 27$ people.

How many Dutch Americans sat at the windows? ** If $1/3$ of the passengers on the bus identifying as Dutch Americans sat at the windows, their number is $1/3*27 = << 1/3*27=9 >> 9$

The answer is 9.

Rewrite solution: How many Dutch people were on the bus? ** On the bus, the number of Dutch people was [in1] of the total number, a total of $<< in1*in0=out0 >>$ people.

How many Dutch Americans were on the bus? ** Out of the [out0] people who were Dutch, [in2] were Dutch Americans, a total of $<< in2*out0=out1 >>$ people.

How many Dutch Americans sat at the windows? ** If [in3] of the passengers on the bus identifying as Dutch Americans sat at the windows, their number is $<< in3*out1=out2 >>$

The answer is out2.

Table 10: Instructions and few-shot examples for the second response rewriting step of our Granular training data construction, used for learning abstract reasoning in AbstRaL.

You are given a math word problem with solution. The numerical values in the problem are replaced with abstract symbols and enclosed in square brackets. The calculations in the solution are also composed of abstract symbols and enclosed in double angle brackets. Please rewrite the solution by first listing all sub-questions, and then answering each sub-question one by one. Please list the relevant conditions before answering each sub-question. Please clarify the final answer at the end of the solution.

Example 1:

Problem: Natalia sold clips to [in0] of her friends in April, and then she sold [in1] as many clips in May.

How many clips did Natalia sell altogether in April and May?

Solution: How many clips did Natalia sell in May? ** Natalia sold $\langle \langle in0 * in1 = out0 \rangle \rangle$ clips in May.

How many clips did Natalia sell altogether in April and May? ** Natalia sold $\langle \langle in0 + out0 = out1 \rangle \rangle$ clips altogether in April and May. The answer is out1.

Rewrite solution: Let's think about the sub-questions we need to answer.

Q1: How many clips did Natalia sell in May?

Q2: How many clips did Natalia sell altogether in April and May?

Let's answer each sub-question one by one.

Q1: How many clips did Natalia sell in May? Natalia sold [in0] clips in April. She sold [in1] as many clips in May as she did in April. So she sold $\langle \langle in0 * in1 = out0 \rangle \rangle$ clips in May.

Q2: How many clips did Natalia sell altogether in April and May? Natalia sold [in0] clips in April. She sold [out0] clips in May. So she sold $\langle \langle in0 + out0 = out1 \rangle \rangle$ clips altogether in April and May.

The final answer is out1.

Example 2:

Problem: The flowers cost \$[in0], the clay pot costs \$[in1] more than the flower, and the bag of soil costs \$[in2] less than the flower. How much does it cost to plant the flowers?

Solution: How much does the clay pot cost? ** The clay pot costs $\langle \langle in1 + in0 = out0 \rangle \rangle$.

How much does the bag of soil cost? ** The bag of soil costs $\langle \langle in0 - in2 = out1 \rangle \rangle$.

How much does it cost to plant the flowers? ** The cost to plant the flowers is $\langle \langle in0 + out0 + out1 = out2 \rangle \rangle$.

The answer is out2.

Rewrite solution: Let's think about the sub-questions we need to answer.

Q1: How much does the clay pot cost?

Q2: How much does the bag of soil cost?

Q3: How much does it cost to plant the flowers?

Let's answer each sub-question one by one.

Q1: How much does the clay pot cost? The flowers cost \$[in0]. The clay pot costs \$[in1] more than the flower.

So the clay pot costs $\langle \langle in0 + in1 = out0 \rangle \rangle$.

Q2: How much does the bag of soil cost? The flowers cost \$[in0]. The bag of soil costs \$[in2] less than the flower.

So the bag of soil costs $\langle \langle in0 - in2 = out1 \rangle \rangle$.

Q3: How much does it cost to plant the flowers? The flowers cost \$[in0]. The clay pot costs \$[out0].

The bag of soil costs \$[out1]. So the cost to plant the flowers is $\langle \langle in0 + out0 + out1 = out2 \rangle \rangle$.

The final answer is out2.

Example 3:

Problem: From March to August, Sam made \$[in0] doing [in1] hours of yard work. However, from September to February, Sam was only able to work for [in2] hours. If Sam is saving up to buy a video game console that costs \$[in3] and has already spent \$[in4] to fix his car, how many more hours does he need to work before he can buy the video game console?

Solution: How much does Sam make per hour? ** Sam makes $\langle \langle in0 / in1 = out0 \rangle \rangle$ /hr. How much did Sam make from September to February? ** From September to February, Sam made $\langle \langle in2 * out0 = out1 \rangle \rangle$.

How much did Sam make from March to February? ** From March to February, Sam made a total of $\langle \langle in0 + out1 = out2 \rangle \rangle$.

How much money did Sam have after fixing his car? ** After fixing his car, he was left with $\langle \langle out2 - in4 = out3 \rangle \rangle$.

How much money does Sam need to buy the video game console? ** Sam needs another $\langle \langle in3 - out3 = out4 \rangle \rangle$.

How many more hours does Sam need to work? ** Sam needs to work another $\langle \langle out4 / out0 = out5 \rangle \rangle$ hours.

The answer is out5.

Rewrite solution: Let's think about the sub-questions we need to answer.

Q1: How much does Sam make per hour?

Q2: How much did Sam make from September to February?

Q3: How much did Sam make from March to February?

Q4: How much money did Sam have after fixing his car?

Q5: How much money does Sam need to buy the video game console?

Q6: How many more hours does Sam need to work?

Let's answer each sub-question one by one.

Q1: How much does Sam make per hour? Sam made \$[in0] doing [in1] hours of yard work.

So he makes $\langle \langle in0 / in1 = out0 \rangle \rangle$ per hour.

Q2: How much did Sam make from September to February? From September to February, Sam worked for [in2] hours.

He makes \$[out0] per hour. So from September to February, he made $\langle \langle in2 * out0 = out1 \rangle \rangle$.

Q3: How much did Sam make from March to February? From March to August, Sam made \$[in0]. From September to February, he made \$[out1]. So from March to February, he made a total of $\langle \langle in0 + out1 = out2 \rangle \rangle$.

Q4: How much money did Sam have after fixing his car? Sam made a total of \$[out2]. He spent \$[in4] to fix his car.

So after fixing his car, he was left with $\langle \langle out2 - in4 = out3 \rangle \rangle$.

Q5: How much money does Sam need to buy the video game console? Sam was left with \$[out3]. The video game console costs \$[in3].

So he needs another $\langle \langle in3 - out3 = out4 \rangle \rangle$.

Q6: How many more hours does Sam need to work? Sam makes \$[out0] per hour. He needs another \$[out4].

So he needs to work another $\langle \langle out4 / out0 = out5 \rangle \rangle$ hours.

The final answer is out5.

Table 11: Instructions and few-shot examples for response rewriting of CoA training data construction.

You are given a math word problem and solution. Please rewrite the solution by replacing the output values in angle brackets with abstract symbols. Please also remove the redundant calculations around the angle brackets. Do not add any other contents to the solution.

Example 1:

Problem: Natalia sold clips to 48 of her friends in April, and then she sold half as many clips in May.

How many clips did Natalia sell altogether in April and May?

Solution: How many clips did Natalia sell in May? ** Natalia sold $48/2 = << 48/2=24 >> 24$ clips in May.

How many clips did Natalia sell altogether in April and May? ** Natalia sold $48+24 = << 48+24=72 >> 72$ clips altogether in April and May. The answer is 72.

Rewrite solution: How many clips did Natalia sell in May? ** Natalia sold $<< 48/2=out0 >>$ clips in May.

How many clips did Natalia sell altogether in April and May? ** Natalia sold $<< 48+out0=out1 >>$ clips altogether in April and May. The answer is out1.

Example 2:

Problem: The flowers cost \$9, the clay pot costs \$20 more than the flower, and the bag of soil costs \$2 less than the flower. How much does it cost to plant the flowers?

Solution: How much does the clay pot cost? ** The clay pot costs $\$20 + \$9 = \$<< 20+9=29 >> 29$.

How much does the bag of soil cost? ** The bag of soil costs $\$9 - \$2 = \$<< 9-2=7 >> 7$.

How much does it cost to plant the flowers? ** The cost to plant the flowers is $\$9 + \$29 + \$7 = \$<< 9+29+7=45 >> 45$. The answer is 45.

Rewrite solution: How much does the clay pot cost? ** The clay pot costs $\$<< 20+9=out0 >>$.

How much does the bag of soil cost? ** The bag of soil costs $\$<< 9-2=out1 >>$.

How much does it cost to plant the flowers? ** The cost to plant the flowers is $\$<< 20+out0+out1=out2 >>$. The answer is out2.

Example 3:

Problem: From March to August, Sam made \$460 doing 23 hours of yard work. However, from September to February, Sam was only able to work for 8 hours. If Sam is saving up to buy a video game console that costs \$600 and has already spent \$340 to fix his car, how many more hours does he need to work before he can buy the video game console?

Solution: How much does Sam make per hour? ** Sam makes $\$460 / 23 \text{ hrs} = \$<< 460/23=20 >> 20/\text{hr}$. How much did Sam make from September to February? ** From September to February, Sam made $8 \text{ hrs} \times \$20/\text{hr} = \$<< 8*20=160 >> 160$.

How much did Sam make from March to February? ** From March to February, Sam made a total of $\$460 + \$160 = \$620$.

How much money did Sam have after fixing his car? ** After fixing his car, he was left with $\$620 - \$340 = \$<< 620-340=280 >> 280$.

How much money does Sam need to buy the video game console? ** Sam needs another $\$600 - \$280 = \$<< 600-280=320 >> 320$.

How many more hours does Sam need to work? ** Sam needs to work another $\$320 / \$20/\text{hr} = << 320/20=16 >> 16$ hours.

The answer is 16.

Rewrite solution: How much does Sam make per hour? ** Sam makes $\$<< 460/23=out0 >>/\text{hr}$. How much did Sam make from September to February? ** From September to February, Sam made $\$<< 8*out0=out1 >>$.

How much did Sam make from March to February? ** From March to February, Sam made a total of $\$<< 460+out1=out2 >>$.

How much money did Sam have after fixing his car? ** After fixing his car, he was left with $\$<< out2-340=out3 >>$.

How much money does Sam need to buy the video game console? ** Sam needs another $\$<< 600-out3=out4 >>$.

How many more hours does Sam need to work? ** Sam needs to work another $<< out4/out0=out5 >>$ hours.

The answer is out5.

Example 4:

Problem: Zhang is twice as old as Li. Li is 12 years old. Zhang's brother Jung is 2 years older than Zhang. How old is Jung?

Solution: How old is Zhang? ** Zhang is $2 * 12 \text{ years old} = << 2*12=24 >> 24$ years old.

How old is Jung? ** Jung is $2 \text{ years} + 24 \text{ years} = << 2+24=26 >> 26$ years old.

The answer is 26.

Rewrite solution: How old is Zhang? ** Zhang is $<< 2*12=out0 >>$ years old.

How old is Jung? ** Jung is $<< 2+out0=out1 >>$ years old.

The answer is out1.

Example 5:

Problem: Of the 90 people on William's bus, $3/5$ were Dutch. Of the $1/2$ of the Dutch who were also American, $1/3$ got window seats. What's the number of Dutch Americans who sat at the windows?

Solution: How many Dutch people were on the bus? ** On the bus, the number of Dutch people was $3/5$ of the total number, a total of $3/5*90 = << 3/5*90=54 >> 54$ people.

How many Dutch Americans were on the bus? ** Out of the 54 people who were Dutch, $1/2$ were Dutch Americans,

a total of $1/2*54 = << 1/2*54=27 >> 27$ people.

How many Dutch Americans sat at the windows? ** If $1/3$ of the passengers on the bus identifying as Dutch Americans

sat at the windows, their number is $1/3*27 = << 1/3*27=9 >> 9$

The answer is 9.

Rewrite solution: How many Dutch people were on the bus? ** On the bus, the number of Dutch people was $3/5$ of the total number, a total of $<< 3/5*90=out0 >>$ people.

How many Dutch Americans were on the bus? ** Out of the out0 people who were Dutch, $1/2$ were Dutch Americans,

a total of $<< 1/2*out0=out1 >>$ people.

How many Dutch Americans sat at the windows? ** If $1/3$ of the passengers on the bus identifying as Dutch Americans

sat at the windows, their number is $<< 1/3*out1=out2 >>$

The answer is out2.
