FINE-GRAINED HALLUCINATION DETECTION AND MITIGATION IN LANGUAGE MODEL MATHEMATICAL REASONING

Anonymous authors

Paper under double-blind review

Abstract

Hallucinations in large language models (LLMs) pose significant challenges in tasks requiring complex multi-step reasoning, such as mathematical problemsolving. Existing approaches primarily detect the presence of hallucinations but lack a nuanced understanding of their types and manifestations. In this paper, we first introduce a comprehensive taxonomy that categorizes the common hallucinations in mathematical reasoning task into six types: fabrication, factual inconsistency, context inconsistency, instruction inconsistency, logical inconsistency, and logical error. We then propose FG-PRM (Fine-Grained Process Reward Model), an augmented model designed to detect and mitigate hallucinations in a fine-grained, step-level manner. To address the limitations of manually labeling training data, we propose an automated method for generating fine-grained hallucination data using LLMs. By injecting hallucinations into reasoning steps of correct solutions, we create a diverse and balanced synthetic dataset for training FG-PRM, which consists of six specialized Process Reward Models (PRMs), each tailored to detect a specific hallucination type. Our FG-PRM demonstrates superior performance across two key tasks: 1) Fine-grained hallucination detection: classifying hallucination types for each reasoning step; and 2) Verification: ranking multiple LLM-generated outputs to select the most accurate solution, mitigating reasoning hallucinations. Our experiments show that FG-PRM outperforms ChatGPT-3.5 and Claude-3 on fine-grained hallucination detection and substantially boosts the performance of LLMs on GSM8K and MATH benchmarks.¹

035

005 006

007

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026

027

028

029

1 INTRODUCTION

While considerable progress has been made in enhancing the general capabilities of large language 037 models (LLMs), solving complex reasoning tasks such as answering mathematical questions re-038 mains a challenge. Recently, advanced techniques like Chain-of-Thoughts (Wei et al., 2022), Treeof-Thoughts (Yao et al., 2024) and Reasoning-via-Planning (Hao et al., 2023) are proposed. These methods guide LLMs in breaking down complex reasoning tasks into simple steps, thus improv-040 ing their performance and enhancing the interpretability of the reasoning process. However, while 041 generating multi-step reasoning chains can improve performance, LLMs often produce incorrect 042 or unverifiable statements—commonly known as hallucinations—that hinder their ability to solve 043 complex problems that require multiple reasoning steps. 044

Prior methods of mitigating hallucinations in reasoning chains largely focus on detecting their pres ence, with limited exploration into the distinct types of hallucinations produced. Our research goes
 beyond this by developing a fine-grained taxonomy that categorizes hallucinations based on their
 nature and manifestation (see Figure 1 for an illustration comparing coarse-grained detection with
 our method). We analyze reasoning steps to pinpoint the emergence of hallucinations and uncover
 patterns in their behavior.

Recent efforts have shown that training reward models (RMs) is an effective approach for detecting and mitigating hallucinations, with the two primary categories being Outcome Reward Model

¹Codes and datasets are available at: https://anonymous.4open.science/r/FG-PRM-75BB/README.md



Hallucination Detection in Language Model Output



(ORM) (Cobbe et al., 2021) and Process Reward Model (PRM) (Lightman et al., 2023). ORMs eval-075 uate the correctness of entire reasoning chains, while PRMs assess each step. PRMs have demon-076 strated superior performance in many scenarios (Wang et al., 2023) since they can provide more 077 granular feedback and effectively guide models' learning process. However, collecting data to train PRMs is labor-intensive, particularly for complex multi-step reasoning tasks, where human annota-079 tion is costly and prone to bias. To address this, we develop a novel method to automatically generate fine-grained hallucination data using LLMs. Specifically, for a given problem with a ground-truth 081 solution, we first identify reasoning steps suitable for hallucination injection. After that, we utilize 082 an LLM to generate additional reasoning steps that incorporate various hallucination types, based on 083 predefined instructions and demonstrations. The generated hallucinatory reasoning steps then serve 084 as negative examples to train task-specific PRMs, each designed to detect a particular hallucination 085 type. This approach improves the accuracy of hallucination mitigation by allowing each PRM to focus on a distinct category.

To evaluate our methods, we test our FG-PRM on two widely used mathematical benchmarks, GSM8K (Cobbe et al., 2021) and MATH (Hendrycks et al., 2021). We validate the effectiveness of our method by two tasks: 1) fine-grained hallucination detection, where we classify different hallucination types at each reasoning step; and 2) verification: where we rank multiple outputs generated by LLMs to select the most accurate solution. Our major contributions are as below:

- We introduce a comprehensive hallucination taxonomy that categorizes common errors in mathematical reasoning tasks into six distinct types.
- We propose an automated method for synthesizing fine-grained hallucination data across all six categories without requiring human annotations. Based on this, we design FG-PRM to detect and mitigate hallucinations in a fine-grained, step-level manner.
- Through extensive experiments, we demonstrate that our FG-PRM demonstrates superior performance in the hallucination detection task compared to ChatGPT-3.5 and Claude-3 (over 5% F1 score on most hallucination types). Moreover, FG-PRM trained on our synthetic data, excel on the verification task on both GSM8k and MATH datasets, as compared to PRMs trained on the more costly human-labeled data Lightman et al. (2023).
- 102 103

092

094

096

098

099

100

101

074

2 FINE-GRAINED HALLUCINATION TAXONOMY

104 105

Large language models excel at solving tasks that require complex multi-step reasoning by gener ating solutions in a step-by-step and chain-of-thought format. Nevertheless, even state-of-the-art models are prone to inaccuracies, often producing content that is unfaithful, fabricated, inconsistent,

or nonsensical. Categorizing and localizing these inaccuracies in reasoning steps is challenging but
 provides explicit insights into which parts of the model output have specific types of problems. To
 address the need for a more nuanced understanding of these hallucinations, we propose a refined
 taxonomy for reasoning tasks.

112 Building upon the prior work Ji et al. (2023), we develop a fine-grained taxonomy for two major 113 categories of hallucinations: intrinsic and extrinsic hallucination, according to whether the hallu-114 cination can be verified by the input information or the context LLMs have previously generated. 115 To describe more complex errors surfacing in LM reasoning, we further divide the intrinsic hal-116 lucination into contextual inconsistency, logical inconsistency and instruction inconsistency, while 117 extrinsic hallucinations are divided into logical error, factual inconsistency, and fabrication. To il-118 lustrate our definition of LLM hallucinations more intuitively, we provide examples for each type of hallucination in Appendix Table 4, along with corresponding explanations. The definition of our 119 proposed categories are elaborated below: 120

121 122

123 124

125

126 127

128

129

130

131

132

133

134 135

136

137

- (1) **Context Inconsistency** refer to instances where the model's is unfaithful with the user's provided contextual information.
- (2) **Logical Inconsistency** refer to internal logical contradictions in the model's output, manifested as inconsistency both among the reasoning steps themselves and between the steps and the final answer.
- (3) **Instruction Inconsistency** refer to instances where the model's output does not align with the user's explicit request.
- (4) **Logical Error** refer to instances where the model makes incorrect calculation or conclusions that do not follow from the provided premises. For instance, it might draw incorrect inferences or make errors in basic logical operations.
 - (5) **Factual Inconsistency** refer to situations where the model's output contains facts that can be grounded in real-world information, but present contradictions.
 - (6) **Fabrication** refer to instances where the model's output contains facts that are unverifiable against established real-world knowledge or context information.

Compared with the simplified taxonomy in previous work (Golovneva et al., 2022; Prasad et al., 2023), our refined taxonomy aims to comprehensively capture the unique complexities of LLM hallucinations, providing a detailed description of the more intricate errors that occur during language model mathematical reasoning.

141 142 143

3 TASK FORMULATION

144 145

In this section, we formulate the two primary tasks of fine-grained hallucination detection and mitigation, highlighting the importance of step-level and fine-grained supervision.

- 146 147 148 149
- 3.1 TASK 1: FINE-GRAINED HALLUCINATION DETECTION

This task aims to detect hallucinations in language model reasoning output at a granular level, focus ing on individual reasoning steps. Specifically, the detector is tasked with identifying fine-grained
 hallucinations in the output of a language model by assigning reward scores for each intermediate
 step in a reasoning chain. The objective is to classify hallucination types at the step level, determin ing whether a specific hallucination type is present.

Given a question x and its solution y consisting of L reasoning steps, we assume the groundtruth annotations for hallucination types are available. These annotations, denoted as $y_i^{*t} \in$ {TRUE, FALSE}, provide a binary label for each hallucination type t at the *i*-th step, indicating whether the hallucination t is present (TRUE) or absent (FALSE). The detector models predict y_i^t , where y_i^t is the model's predicted label for the *i*-th step and hallucination type t. We evaluate the model's performance using standard metrics for classification as in previous work (Feng et al., 2023; Mishra et al., 2024): precision and recall. For each hallucination type t, the precision measures the proportion of correct predictions out of all predictions where the model indicated the presence of 162 a hallucination at a step, while recall measures the proportion of actual hallucination steps that the 163 model correctly identified. These are computed as follows: 164

$$\operatorname{Precision}^{t} = \frac{\sum_{i \in L} \mathbb{I}[y_{i}^{t} = y_{i}^{*t}]}{\sum_{i \in L} \mathbb{I}[y_{i}^{t} = \mathsf{TRUE}]}$$
(1)

$$\operatorname{Recall}^{t} = \frac{\sum_{i \in L} \mathbb{I}\left[y_{i}^{t} = y_{i}^{*t}\right]}{\sum_{i \in L} \mathbb{I}\left[y_{i}^{*t} = \operatorname{TRUE}\right]}$$
(2)

170 Here, $\mathbb{I}\left[\cdot\right]$ is an indicator function that returns 1 if the condition is true and 0 otherwise. Precision indicates the proportion of correctly predicted hallucinations for type t, while recall indicates how 172 many of the true hallucinations were detected by the model.

173 To assess the overall performance across all hallucination types, we calculate the F1 score, which 174 is the harmonic mean of precision and recall. The F1 score is computed for each hallucination type 175 and then averaged across all types \mathcal{E} : 176

F1 Score =
$$\frac{1}{|\mathcal{E}|} \sum_{t \in \mathcal{E}} \frac{2 \times \operatorname{Precision}^t \times \operatorname{Recall}^t}{\operatorname{Precision}^t + \operatorname{Recall}^t}$$
 (3)

Thus, fine-grained hallucination detection can be framed as a set of binary classification tasks, where 181 the system predicts whether each reasoning step s_i contains a specific hallucination type. By eval-182 uating precision, recall, and F1 score across different hallucination types, we gain a comprehensive 183 understanding of the model's ability to detect and categorize hallucinations within complex reason-184 ing processes.

186 3.2 TASK 2: FINE-GRAINED HALLUCINATION MITIGATION 187

The verification task (Lightman et al., 2023) assesses a model's ability to evaluate and rank multiple 188 candidate solutions for a given problem. In this task, a generator produces N possible solutions 189 $\{y^1, y^2, \dots, y^N\}$ for a problem x, which are then evaluated by a reward model (Section 4.1). The 190 reward model assigns a score to each candidate solution based on its correctness, with the goal of 191 selecting the best solution among the candidates. 192

This task follows the best-of-N selection method, where the solution with the highest score is cho-193 sen as the final answer. A well-performing reward model improves the likelihood of selecting the 194 correct solution, thereby enhancing the overall problem-solving accuracy. By providing meaningful 195 feedback on each candidate solution, the verification task helps ensure that the reasoning process is 196 grounded in correctness and consistency. 197

199

200

165

166 167

168 169

171

177

178 179

185

METHODOLOGY 4

201 In this section, we first introduce two basic types of reward models (Section 4.1), the Outcome Reward Model (ORM) and the Process Reward Model (PRM). After that, we describe our automated 202 framework for generating hallucination-annotated datasets, followed by a detailed explanation of the 203 training procedure for our Fine-Grained Process Reward Model (FG-PRM), elaborating on the use 204 of generated datasets and how our model enhances both hallucination detection and verification 205 performance (Section 4.2). 206

207 208

4.1 PRELIMINARY

209 **ORM** The Outcome-supervised Reward Model (ORM) was introduced by Cobbe et al. (2021). 210 Given a problem x and its solution y, an ORM assigns a sigmoid score r_y to the entire solution, 211 indicating whether y is correct. ORMs are typically trained using cross-entropy loss over the entire 212 solution. Assume y^* is the ground-truth label of the solution $y, y^* = 1$ if y is correct, otherwise 213 $y^* = 0$. The training objective minimizes the cross-entropy between the predicted outcome r_y and 214 the ground-truth y^* : 215

$$\mathcal{L}_{\text{ORM}} = y^* \log r_y + (1 - y^*) \log(1 - r_y)$$
(4)



Figure 2: Our automated reasoning process annotation involves two steps: First, for each reasoning step, we instruct a language model to verify the feasibility of injecting hallucinations (using Logical Error as an example in this figure). Second, for steps where hallucinations can be injected, we prompt the language model to introduce hallucinations by providing instructions (see in Appendix E) and few-shot demonstrations (see in Appendix F), which serve as negative examples for training FG-PRM.

However, ORM's coarse feedback mechanism limits its ability to diagnose errors within individual reasoning steps, as it only evaluates the final solution without considering intermediate correctness.

PRM The Process-supervised Reward Model (PRM), introduced by Lightman et al. (2023), addresses the limitations of ORM by providing fine-grained, step-level feedback. Instead of assigning a single score to the entire solution, PRM assigns a sigmoid score r_{y_i} for each reasoning step y_i in the solution y. This approach enables the model to evaluate the correctness of each intermediate step, providing more detailed feedback on where the reasoning process succeeds or fails. The training objective for PRM minimizes the sum of cross-entropy losses over all reasoning steps, allowing the model to learn from fine-grained supervision:

 $\mathcal{L}_{\text{PRM}} = \sum_{i=1}^{L} \log y_i^* \log r_{y_i} + (1 - y_i^*) \log(1 - r_{y_i})$ (5)

where L is the number of reasoning steps in the solution y and y_i^* is the ground-truth label of the *i*-th step of y. By providing feedback at the step level, PRM offers significant advantages over ORM in tasks requiring complex, multi-step reasoning. PRM not only improves the model's ability to detect and correct errors within individual steps but also enables more targeted learning and fine-tuning.

4.2 FG-PRM: FINE-GRAINED PROCESS REWARD MODEL

In this Section, we introduce our FG-PRM, the Fine-Grained Process Reward Model for hallucina tion detection and mitigation. To reduce the annotation cost issues associated with PRM, we first introduce an automated process annotation framework for step-level fine-grained dataset synthesis.
 After that, we provide the training details for our FG-PRM on the synthetic dataset.

4.2.1 AUTOMATED HALLUCINATION GENERATION

To detect fine-grained hallucinations in language model reasoning tasks, we propose a framework based on fine-grained, step-level process supervision. Existing step-level datasets with fine-grained annotations (Golovneva et al., 2022) are limited in size, and collecting the necessary data for training models with such detailed labels is costly, as it requires human annotators to provide fine-grained feedback for each reasoning step. To overcome the scarcity of human-labeled data, we introduce an automated reasoning process hallucination annotation framework, as illustrated in Figure 2. We treat the golden chain-of-thought (CoT) reasoning steps as positive examples, while negative examples are generated by injecting hallucinations into these steps using Llama3-70B (Dubey et al., 2024). To synthesize the negative examples, we adopt a two-step process as follows.

270 Step 1: Identify where to inject hallucination In our taxonomy, each hallucination type has 271 distinct characteristics, requiring specific conditions and methods for generation. However, not all 272 reasoning steps are suitable for generating every type of hallucination. For instance, when a reason-273 ing step is exclusively focused on numerical calculations, it becomes challenging to introduce factual inconsistency. To effectively introduce a hallucination type into the reasoning process, we need to 274 identify steps that meet the necessary conditions for hallucination generation. To achieve this, we 275 have developed a set of tailored rules for the Llama3-70B model. These rules guide the model in 276 determining whether a given reasoning step provides the elements required for a specific type of hallucination. For example, when evaluating whether a step can introduce factual inconsistency, the 278 model checks if the reasoning step references objects (e.g., quantities, features) or named entities. 279 This enables us to manipulate the information, allowing for the seamless integration of contextual 280 inconsistencies in later steps. The complete set of rules for identifying hallucination injection points 281 across the six hallucination types is detailed in Appendix D. 282

283 Step 2: Hallucinate the ground truth reasoning steps To control the distribution of hallucina-284 tions in the generated dataset and improve the success rate of incorporating our hallucination taxon-285 omy, we prompt the Llama3-70B model to insert hallucinations one by one from our taxonomy. We 286 begin by inputting specific instructions for each hallucination type into the system prompt, guiding 287 the model on how to modify the reasoning process and introduce the desired hallucination. Detailed 288 instructions for each hallucination type are provided in Appendix E. Next, we employ an in-context learning strategy by providing two demonstrations for each query type. Each demonstration includes 289 an example of an injected hallucination, along with an explanation of how it is introduced. These 290 demonstrations can be found in Appendix F. After confirming the appropriate location for injecting 291 the hallucination, we present the problem and the correct reasoning history to the model, instructing 292 it to generate the next reasoning step with the target hallucination. For cost efficiency, we delegate 293 the task of hallucinating reasoning steps to the Llama3-70B model. We experimentally found that 294 our method enables the language model to consistently generate hallucinatory reasoning steps with 295 a high success rate. 296

4.2.2 MODEL TRAINING

After generating six types of hallucination datasets using our automated data annotation method, we train our FG-PRM, denoted as R_{Φ} , which comprises six distinct PRMs, $R_{\phi_1} \dots R_{\phi_6}$, each corresponding to a specific type of hallucination in our taxonomy.

302 Formally, given an input question x and the corresponding solution y composed of L reasoning steps $\{y_1, y_2, \ldots, y_L\}$, we separately train task-specific PRMs R_{ϕ_t} to detect whether each 303 reasoning step in y contains the hallucination type t. The model input has the format of 304 "question: q, reasoning steps: y_1 [sep] y_2 [sep]... y_L [sep]", where the [sep] token 305 represents the classification output at each reasoning step to indicate whether the following step 306 y_i contains the hallucination type t. We define $R_{\phi_t} = +1$ if R_{ϕ_t} predicts "no error" for y_i and -1307 otherwise. To train each PRM R_{ϕ_t} , we utilize a step-level classification loss as in Eq.5 to each [sep] 308 token before step y_i . Overall, our FG-PRM R_{Φ} generates an aggregate reward for the solution y of 309 the input question x: 310

$$R_{\Phi}(x,y) = \sum_{t=1}^{6} \sum_{i=1}^{L} \left(R_{\phi_t}(x,y_i) \right)$$
(6)

 R_{ϕ_t} denotes the fine-grained reward models specific to each hallucination type t. In our verification task, we use this aggregate reward to represent the final reward of a solution assigned by FG-PRM.

317 318 319

321

311312313314

315

316

5 EXPERIMENTS

320 5.1 SETTINGS

Datasets We conduct our experiments on two widely used mathematical benchmarks, GSM8K
 (Cobbe et al., 2021) and MATH (Hendrycks et al., 2021). GSM8K is a high-quality dataset consisting of grade school math problems designed to benchmark the reasoning abilities of language

Table 1: Performance of fine-grained hallucination detection across all hallucination types on synthetic data and human-annotated data. All numbers are F1 scores.

		Syn	thetic	Reaso	ning C	hain			Human	-annota	ated Re	easonir	ng Cha	in
Detector	CI	LI	Π	LE	FI	FA	Avg.		LI	Π	LE	FI	FA	Avg.
ChatGPT	0.415	0.522	0.453	0.360	0.428	0.900	0.513	0.4	42 0.552	0.510	0.377	0.487	0.840	0.531
Claude	0.448	0.388	0.493	0.275	0.373	0.963	0.490	0.4	34 0.460	0.478	0.359	0.428	0.758	0.503
FG-PRM	0.488	0.549	0.529	0.398	0.422	0.608	0.499	0.5	526 0.575	0.513	0.377	0.426	0.484	0.484

models. To construct the hallucinatory reasoning steps, we employ a meta-dataset and software library (Ott et al., 2023), which collects the golden chain-of-thought solutions for each problem in the GSM8K. MATH, on the other hand, is a large-scale dataset designed for probing and improving model reasoning, which includes human-written step-by-step solutions.

Following Lightman et al. (2023), Uesato et al. (2022) and Wang et al. (2023), we sample instances 337 from GSM8k and MATH datasets to build training sets and test sets. For GSM8k, we sample 700 338 instances from the training set and 100 instances from the test set. For MATH, we sample 700 339 instances from the training set. The test set is sampled from the MATH dataset and includes 100 340 instances. We call both datasets "Coarse-grained Hallucinations" (CG-H). They include human-341 annotated reasoning chains. Each step has a binary labels indicating their correctness. Based on 342 the two sampled datasets, we augment each to 12,000 instances by our automatic hallucination 343 generation method, including all types of hallucinations described in Section 2 with a balanced 344 hallucination distribution. The augmented datasets are called "Fine-grained Hallucinations" (FG-345 H). In both tasks, we use our synthetic hallucination dataset, FG-H, as train and test data with a split of 95%:5%. Additionally, we sample 12,000 instances from Math-Shepherd (Wang et al., 2023). 346 The dataset includes automatically constructed step-by-step solutions by applying the Monte Carlo 347 sample method. 348

349

324

325

326 327 328

330 331

Models In the fine-grained hallucination detection task, we evaluate the performance of prompt-based and model-based detection. For prompt-based detection, we apply ChatGPT (GPT-3.5-turbo-0125) (Ouyang et al., 2022) and Claude (Claude-3-haiku)² with carefully designed prompts as baseline methods. For model-based detection, we apply our FG-PRM.

354 In the fine-grained hallucination mitigation task, we apply various verifiers to evaluate the correct-355 ness of solutions generated by language models (generators). We employ Llama3-70B (Dubey et al., 2024) as our solution generator, from which we sample 64 candidate solutions for each test problem. 356 We apply the LongFormer-base-4096 (Beltagy et al., 2020) and Llama-3-8B (Dubey et al., 2024) as 357 our base models due to its strong performance in handling long-context reasoning. Verifiers include 358 self-consistency (SC), ORM, PRM, and FG-PRM. The self-consistency verifier serves as a baseline 359 without specific model training; it aggregates multiple reasoning paths and selects the most frequent 360 solution as the final answer. Both ORMs and PRMs are trained on the CG-H dataset. For our FG-361 RPM, we train individual fine-grained PRMs for each of six hallucination types, following the same 362 supervision as PRMs on FG-H data. In inference processes, we sum all results as our final results. 363

364 365

5.2 HALLUCINATION DETECTION RESULTS

To evaluate the efficacy of our method in detecting fine-grained hallucinations, we conduct two experiments on synthetic and human-annotated data.

368 369

370

371

Synthetic Data We utilize the automated annotation labels from our synthetic dataset, FG-H, as the golden standard for evaluating various detectors across six types of hallucinations.

As shown in Table 1, FG-PRM outperforms prompt-based detectors in detecting CI, LI, II, and LE, demonstrating that FG-PRM has effectively learned the patterns of these hallucinations and is capable of detecting them accurately. However, prompt-based detectors outperform FG-PRM on FI and FA, primarily due to their larger model sizes and greater access to fact-based knowledge. This reflects the inherent advantage of large language models in fact-based verification. Moreover, precision and recall results are in Tables 5 and 6 in Appendix B.

²https://claude.ai/

Human-annotated Data To further validate the effectiveness of our method on real world data, we conduct evaluation on human annotated data. Specifically, we first use ChatGPT (GPT-3.5-turbo-0125) (Ouyang et al., 2022) to generate step-by-step solutions on 50 problems in MATH dataset. These solutions are then manually annotated by three graduate students using the taxonomy of hallucination types proposed in Section 2. The mutual agreement among annotators is 79%. Moreover, each type of hallucination has 50 human-annotated reasoning chains, along with the corresponding hallucinations.

385 The results on the human-annotated data align closely with the trends observed on the synthetic data. 386 Our FG-PRM model demonstrates improving performance in detecting CI and LI hallucinations, 387 where it consistently outperforms both ChatGPT and Claude. However, FG-PRM's performance is 388 slightly below that of the strong, non-public LLMs (e.g. ChatGPT and Claude) in detecting FI and 389 FA hallucinations. This is largely attributable to FG-PRM's smaller parameter size and limited access to world knowledge. Despite these challenges, FG-PRM performs competitively overall, 390 particularly in reasoning-related hallucinations. Further analysis about reasoning chain evaluation 391 for various verifiers is presented in Appendix C. 392

393 394

427

5.3 HALLUCINATION MITIGATION RESULTS

395 Table 2 presents a performance comparison of various verifiers on GSM8k and MATH. FG-Process 396 Reward Models (FG-PRMs) trained on our augmented dataset, FG-H, significantly outperform all 397 baselines across both base models. Notably, after fine-tuning with FG-H, Longformer and Llama3-398 8B achieve 94% and 58% accuracy on GSM8k and MATH, respectively, surpassing PRMs trained 399 on Math-Shepherd. The results show that base models mitigated by PRMs consistently outperform 400 those mitigated by ORMs, consistent with findings from Uesato et al. (2022), Lightman et al. (2023), 401 and Wang et al. (2023). On GSM8k, most baseline verifiers perform close to the self-consistency level due to the simplicity of the dataset, where many questions involve only basic arithmetic op-402 erations. However, the differences between verifiers become more evidence in the more complex 403 MATH dataset, where questions and reasoning steps often require LaTeX math expressions. These 404 results indicate that the balanced fine-grained step-level supervision employed by FG-PRMs offers 405 a more robust and effective approach to hallucination mitigation, particularly in handling complex 406 problem-solving tasks. 407

Table 2: Performance of different verifiers on GSM8K and MATH benchmarks. The evaluation is based on 64 candidate solutions generated by Llama3-70B model with greedy decoding for each problem. Each result is the mean of results from 3 groups of sampling results. Statistical significant test on most improvements compared to the "Self-Consistency" have (p < 0.05).

1		, <u> </u>	
Base Model	Verifier / Reward Model	GSM8K	MATH
-	Self-Consistency	0.88	0.48
	ORM	0.88	0.51
	PRM	0.89	0.53
I. F.	Math-Shepherd (ORM)	0.90	0.52
LongFormer	Math-Shepherd (PRM)	0.91	0.54
	FG-PRM (Ours)	0.94	0.57
	ORM	0.87	0.52
	PRM	0.90	0.53
11 2.00	Math-Shepherd (ORM)	0.89	0.51
Liama3-8B	Math-Shepherd (PRM)	0.91	0.53
	FG-PRM (Ours)	0.93	0.58
	Base Model - LongFormer Llama3-8B	Base ModelVerifier / Reward Model-Self-Consistency-ORM PRM Math-Shepherd (ORM) Math-Shepherd (PRM)FG-PRM (Ours)ORM PRM Math-Shepherd (ORM) Math-Shepherd (ORM) Math-Shepherd (PRM)Llama3-8BORM PRM Math-Shepherd (PRM) FG-PRM (Ours)	Base Model Verifier / Reward Model GSM8K - Self-Consistency 0.88 - Self-Consistency 0.88 PRM 0.89 0.89 Math-Shepherd (ORM) 0.90 Math-Shepherd (PRM) 0.91 FG-PRM (Ours) 0.94 Llama3-8B ORM 0.87 FG-PRM (Ours) 0.91 FG-PRM (Ours) 0.91 FG-PRM (Ours) 0.91

6 ANALYSIS

Hallucination Mitigation Performance with Different Number of Candidate Solutions Figure 3 illustrates the performance of four verifiers with the number of candidate solutions ranging from 1 to 64 across two benchmarks. This demonstrates that FG-PRM consistently outperforms all other verifiers. With predicted insights, the performance gap between FG-PRM and other baseline verifiers will increase with the growth of N.



Figure 3: Performance of LLaMA3-70B using different verification methods across different numbers of candidate solutions on GSM8K and MATH.

449 Out-of-Distribution Dataset Evaluation We further conduct out-of-distribution (OOD) evaluation experiments to assess the ro-450 bustness and transferability of our approach. In these experiments, 451 we train verifiers on the GSM8k (CG-H) and GSM8k (FG-H), and 452 test them on the MATH dataset. Then, we compare these verifiers 453 with the ones trained on the MATH (CG-H) and MATH (FG-H). 454 Notably, the GSM8k dataset contains simpler questions, predomi-455 nantly solvable through basic arithmetic operations, unlike the more 456 complex MATH dataset. 457



As detailed in Figure 4, the gap in CG-H (+0.3) is more signifi-458 cant than in FG-H (+0.1). Moreover, the verifiers trained on the 459 GSM8k (FG-H) demonstrates performance closely comparable to 460 those trained on the MATH (FG-H) dataset. This indicates that ver-461 ifiers trained on FG-H effectively learn to recognize patterns of hallucinations and can generalize this knowledge to tackle more challenging scenarios effectively. 462 463

Figure 4: Out-of-distribution performance on the MATH benchmark. PRM and FG-PRM are trained on CG-H

Qualitative Result of Fine-Grained Hallucination Evaluation To illustrate the effect of our 464 FG-PRMs, we demonstrate a case study in Table 3, comparing the hallucination detection results of 465 FG-PRM and CG-PRM. Our FG-PRM demonstrates remarkable discrimination by precisely detect 466 fine-grained types of hallucination in reasoning steps. Notably, Step 7 has a high Logical Inconsis-467 tency score since the step incorrectly refers to the formula in previous steps. It refers to a = 5 as 468 a = 11, which causes Logical Inconsistency. FG-PRM recognizes it by giving a low score, while 469 CG-PRM does not recognize this detail.

470 471

472

447

448

7 **RELATED WORK**

473 Hallucinations in LLM Numerous studies have explored the taxonomy of hallucinations in lan-474 guage models, offering valuable insights. Mishra et al. (2024) identify six fine-grained types of 475 factual hallucinations in information search scenarios. Zhang et al. (2023) classify hallucinations 476 based on conflict types: input-conflicting, context-conflicting, and fact-conflicting. Huang et al. (2023a) categorize hallucinations into factuality and faithfulness types, subdividing them further. Ji 477 et al. (2023) discussed intrinsic and extrinsic hallucinations, focusing on whether outputs contradict 478 input content. Our taxonomy, in contrast, concentrates on hallucinations in reasoning tasks, partic-479 ularly through chain-of-thought reasoning. Closest to our work, Golovneva et al. (2022) propose 480 ROSCOE, measuring semantic alignment, similarity, and language coherence in reasoning chains. 481 Unlike ROSCOE that includes common issues like grammar, redundancy, and repetition, which are 482 easier to spot in general, our taxonomy provides detailed distinctions between error types, providing 483 a compact framework for diagnosing complex errors and improving model outputs. 484

Evaluation of Reasoning Chains Depending on whether requiring golden references, methods to 485 evaluate reasoning chains can be roughly divide as reference-dependent and reference-free ones. For

Question	Each edge length of a rectangular solid is a prime number. If the volume of the rectangular solid is 385 cubic units, what is the total surface area, in square units, of the rectangular solid?
Answer	334
Reasoning	Step 1:We know that the volume of this rectangular solid is the product of its three side
Steps	lengths. FG-PRM: 0.95, CG-PRM: 0.83.
	Step 2: Let's call those side lengths a, b, and c. The volume is abc. FG-PRM: 0.83, CG-PRM:
	0.80.
	Step 6: Let's plug in the possible values. If $a = 5$, $b = 7$, and $c = 11$, then $2ab+2bc+2ac = 2ac$
	2(5)(7) + 2(7)(11) + 2(5)(11) = 74 + 154 + 110 = 338. FG-PRM: 0.87, CG-PRM: 0.72.
	Step 7: Which simplifies to $14 \cdot 11 + 22 \cdot 11 + 10 \cdot 11$. FG-PRM: Logical Inconsistency
	0.07 CG-PRM: 0.80.
	Step 8: So the total surface area would be 334. FG-PRM : 0.88 CG-PRM: 0.66.
	·

Table 3: A case study from the MATH dataset. A high CG-PRM score indicates that the step is positive. A low FG-PRM-TYPE score indicates the step has a high probability of having the TYPE of hallucination.

504 reference-dependent, the reasoning chains can be evaluated with LLMs (Ren et al., 2023; Adlakha 505 et al., 2023)), or by measuring the discrepancy between the vanilla response and reference (Huo 506 et al., 2023; Pezeshkpour, 2023). For reference-free metrics, some methods rely on aggregating 507 the individual token probabilities assigned by the LLM during generation so that they can reflect 508 reasoning chain uncertainty (Manakul et al., 2023; Huang et al., 2023b). In addition to that, many 509 model-based methods have emerged to evaluate reasoning chains. He et al. (2024) proposed to 510 prompt GPT-4 in a Socratic approach. Hao et al. (2024) employ GPT-4 to summarizes evaluation 511 criteria tailored for each task, and then evaluate the reasoning chains following the criteria. In this work, we focus on model-based reference-free reasoning chain evaluation from the perspective of 512 hallucination detection. 513

514

500

501

502 503

515 **Improving reasoning abilities of LLMs** For LLMs that have completed training, prompting tech-516 niques is an effective approach to improve the performance of LLMs on reasoning tasks without 517 modifying the model parameters. Many studies have developed different prompting strategies in 518 reasoning tasks, such as the Chain-of-Thought Wei et al. (2022); Fu et al. (2022), Tree-of-Thoughts 519 Yao et al. (2024). Besides, instead of directly improving the reasoning performance of LLMs, veri-520 fiers can raise the success rate in solving reasoning tasks by selecting the best answer from multiple decoded candidates. Two types of verifiers are commonly used: Outcome Reward Model (ORM) 521 and Process Reward Model (PRM). PRM provides a more detailed evaluation by scoring each indi-522 vidual step. However, training a PRM requires access to expensive human-annotated datasets, which 523 can be a barrier to the advancement and practical application of PRM. To overcome this challenge, 524 methods such as Math-Shepherd Wang et al. (2023) and MiPS Wang et al. (2024) have explored 525 Monte Carlo estimation to automate the data collection process without human involvement, and 526 OmegaPRM Luo et al. (2024) proposed a divide-and-conquer style Monte Carlo tree search al-527 gorithm for automated process supervision data generation. Different from the above methods, we 528 directly generate negative examples with hallucinations instead of annotating the reasoning process.

- 529 530
- 531

8 CONCLUSION

532 533

In conclusion, we propose FG-PRM framework, introduces a nuanced approach for comprehensive
 understanding and mitigations of hallucinations in language model reasoning, which are catego rized into six distinct types under our new paradigm. By leveraging a novel automatic data genera tion method, we significantly reduce the dependency on costly human annotations while enriching
 the dataset with diverse hallucinatory instances. Our empirical results demonstrate that FG-PRM,
 when trained on this synthetic data, significantly enhances the accuracy of hallucination detection,
 providing an effective mechanism for improving the LLM reasoning accuracy and faithfulness.

540 REFERENCES

547

- Vaibhav Adlakha, Parishad BehnamGhader, Xing Han Lu, Nicholas Meade, and Siva Reddy. Eval uating correctness and faithfulness of instruction-following models for question answering. *arXiv preprint arXiv:2307.16877*, 2023.
- Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer.
 arXiv preprint arXiv:2004.05150, 2020.
- 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.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Shangbin Feng, Vidhisha Balachandran, Yuyang Bai, and Yulia Tsvetkov. FactKB: Generalizable factuality evaluation using language models enhanced with factual knowledge. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 933–952, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.59. URL https: //aclanthology.org/2023.emnlp-main.59.
- Yao Fu, Hao Peng, Ashish Sabharwal, Peter Clark, and Tushar Khot. Complexity-based prompting
 for multi-step reasoning. In *The Eleventh International Conference on Learning Representations*,
 2022.
- Olga Golovneva, Moya Chen, Spencer Poff, Martin Corredor, Luke Zettlemoyer, Maryam Fazel-Zarandi, and Asli Celikyilmaz. Roscoe: A suite of metrics for scoring step-by-step reasoning. *arXiv preprint arXiv:2212.07919*, 2022.
- Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting Hu.
 Reasoning with language model is planning with world model. *arXiv preprint arXiv:2305.14992*, 2023.
- Shibo Hao, Yi Gu, Haotian Luo, Tianyang Liu, Xiyan Shao, Xinyuan Wang, Shuhua Xie, Haodi
 Ma, Adithya Samavedhi, Qiyue Gao, et al. Llm reasoners: New evaluation, library, and analysis
 of step-by-step reasoning with large language models. *arXiv preprint arXiv:2404.05221*, 2024.
- Hangfeng He, Hongming Zhang, and Dan Roth. Socreval: Large language models with the socratic method for reference-free reasoning evaluation. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pp. 2736–2764, 2024.
- 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.
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *arXiv preprint arXiv:2311.05232*, 2023a.
- Yuheng Huang, Jiayang Song, Zhijie Wang, Huaming Chen, and Lei Ma. Look before you leap:
 An exploratory study of uncertainty measurement for large language models. *arXiv preprint arXiv:2307.10236*, 2023b.
- Siqing Huo, Negar Arabzadeh, and Charles LA Clarke. Retrieving supporting evidence for llms generated answers. *arXiv preprint arXiv:2306.13781*, 2023.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang,
 Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation. ACM
 Computing Surveys, 55(12):1–38, 2023.

594 Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan 595 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let's verify step by step. arXiv preprint 596 arXiv:2305.20050, 2023. 597 Liangchen Luo, Yinxiao Liu, Rosanne Liu, Samrat Phatale, Harsh Lara, Yunxuan Li, Lei Shu, Yun 598 Zhu, Lei Meng, Jiao Sun, et al. Improve mathematical reasoning in language models by automated process supervision. arXiv preprint arXiv:2406.06592, 2024. 600 601 Potsawee Manakul, Adian Liusie, and Mark JF Gales. Selfcheckgpt: Zero-resource black-box hallu-602 cination detection for generative large language models. arXiv preprint arXiv:2303.08896, 2023. 603 Abhika Mishra, Akari Asai, Vidhisha Balachandran, Yizhong Wang, Graham Neubig, Yulia 604 Tsvetkov, and Hannaneh Hajishirzi. Fine-grained hallucination detection and editing for language 605 models. arXiv preprint arXiv:2401.06855, 2024. 606 607 Simon Ott, Konstantin Hebenstreit, Valentin Liévin, Christoffer Egeberg Hother, Milad Moradi, 608 Maximilian Mayrhauser, Robert Praas, Ole Winther, and Matthias Samwald. Thoughtsource: A central hub for large language model reasoning data. Scientific data, 10(1):528, 2023. 609 610 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong 611 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to fol-612 low instructions with human feedback. Advances in neural information processing systems, 35: 613 27730-27744, 2022. 614 Pouya Pezeshkpour. Measuring and modifying factual knowledge in large language models. In 2023 615 International Conference on Machine Learning and Applications (ICMLA), pp. 831–838. IEEE, 616 2023. 617 618 Archiki Prasad, Swarnadeep Saha, Xiang Zhou, and Mohit Bansal. Receval: Evaluating reasoning 619 chains via correctness and informativeness. arXiv preprint arXiv:2304.10703, 2023. 620 Ruiyang Ren, Yuhao Wang, Yingqi Qu, Wayne Xin Zhao, Jing Liu, Hao Tian, Hua Wu, Ji-Rong 621 Wen, and Haifeng Wang. Investigating the factual knowledge boundary of large language models 622 with retrieval augmentation. arXiv preprint arXiv:2307.11019, 2023. 623 624 Jonathan Uesato, Nate Kushman, Ramana Kumar, Francis Song, Noah Siegel, Lisa Wang, Antonia 625 Creswell, Geoffrey Irving, and Irina Higgins. Solving math word problems with process-and outcome-based feedback. arXiv preprint arXiv:2211.14275, 2022. 626 627 Peiyi Wang, Lei Li, Zhihong Shao, RX Xu, Damai Dai, Yifei Li, Deli Chen, Y Wu, and Zhifang 628 Sui. Math-shepherd: Verify and reinforce llms step-by-step without human annotations. CoRR, 629 abs/2312.08935, 2023. 630 631 Zihan Wang, Yunxuan Li, Yuexin Wu, Liangchen Luo, Le Hou, Hongkun Yu, and Jingbo Shang. Multi-step problem solving through a verifier: An empirical analysis on model-induced process 632 supervision. arXiv preprint arXiv:2402.02658, 2024. 633 634 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny 635 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. Advances in 636 neural information processing systems, 35:24824–24837, 2022. 637 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik 638 Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. Ad-639 vances in Neural Information Processing Systems, 36, 2024. 640 641 Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, 642 Yu Zhang, Yulong Chen, et al. Siren's song in the ai ocean: a survey on hallucination in large 643 language models. arXiv preprint arXiv:2309.01219, 2023. 644 645 646 647

Appendix

A EXAMPLES OF FINE-GRAINED HALLUCINATION

Hallucination Type	Sub-Type	User Input	Model Output	Explanation
	Factual	What is the capital of	The capital of Australia	The model incorrectly states
	Inconsistency	Australia?	is Sydney.	that Sydney is the capital of
				Australia, while the correct cap- ital of Australia is Capberra
				This error demonstrates factual
				inconsistency as the model's
				world fact.
	Fabrication	What is the population	The population of At-	The model provided a specific
	Tubricution	of Atlantis?	lantis is approximately	population figure for Atlantis, a
Extrinsic Hallucinations			1.2 million people.	fictional place with no basis in
				exist, there is no verifiable pop-
				ulation data.
	Logical	If a rectangle has a	The perimeter for-	While the model correctly
	Error	length of 4 units and a	mula of a rectangle	identifies the formula for the
		width of 3 units, what	is: Perimeter = $2 \times$	perimeter of a rectangle, it
		rectangle?	length is 4 units and the	of the length and width and
		0	width is 3 units. So, the	states that the perimeter is 16
			perimeter would be $2 \times (4 + 3)$. This equals 2	sion does not logically follow
			\times 8, which gives us 16	from the premises it provided,
			units.	leading to an incorrect answer.
	Instruction	Please calculate the	The quotient of 25 di-	Although the model correctly
	Inconsistency	quotient of 25 divided	vided by 100 is 0.25.	calculates the quotient of 25 di-
		the result as a fraction.		to follow the user's directive
				to provide the result as a frac-
				tion, thereby deviating from the user's specific instruction
	Context			
	Inconsistency	Given that I have 5 ap-	Considering the infor-	The model output contains a
		friend, how many ap-	initially had 5 apples	nation because it contradicts the
Intrinsic		ples do I have left?	and you gave 3 to your	user's provided context. The
Hallucination			2 apples left.	apples to their friend, but the
			TT	model incorrectly states that the
				user gave 3 apples to their friend
	Logical			
	Inconsistency	How many hours are	There are 7 days in a	The model correctly states that
		there in a week if each day has 24 hours?	has 24 hours a week	there are / days in a week but then incorrectly calculates the
		nuo 2 · nouro ·	would have 6 days	number of hours by multiplying
			times 24 hours, which	6 days by 24 hours to get 144 hours. This shows a contradic
			equais 177 nouis.	tion between the days in a week
				referred and a failure to main-
				tain internal logical consistency in the reasoning process.
				6 Processi

Table 4: Examples of various types of LLM hallucinations in our taxonomy. Content marked in Red represents the hallucinatory output, while content marked in Blue indicates user instruction or provided context that contradicts the LLM hallucination.

702

703 704

705

723 724 725

726

B DETAILED FINE-GRAINED HALLUCINATION DETECTION RESULTS

The precision and recall of the fine-grained detection results for the Llama3-70B generation are reported in Table 5 and 6, respectively.

		Hallu	cination	Туре			
Detector	CI	LI	II	LE	FI	FA	Average
ChatGPT	0.403	0.488	0.450	0.424	0.412	0.890	0.511
Claude-3	0.417	0.368	0.490	0.248	0.357	0.952	0.472
FG-PRM	0.428	0.513	0.528	0.413	0.403	0.589	0.479

Table 5: Precision for fine-grained hallucination detection across different categories.

		Hallu	cination	Туре			
Detector	CI	LI	Π	LE	FI	FA	Average
ChatGPT	0.440	0.600	0.460	0.541	0.477	0.920	0.573
Claude-3	0.525	0.433	0.500	0.334	0.416	0.990	0.533
FG-PRM	0.571	0.597	0.560	0.546	0.462	0.635	0.562

Table 6: Recall for fine-grained hallucination detection across different categories.

C REASONING STEP HALLUCINATION EVALUATION

We utilize our model to evaluate hallucination issues in the generated outputs of large language models. Each generation is assigned six scores corresponding to hallucination types. The score under each hallucination type for a model is calculated based on the proportion of correct reasoning steps in generations. Specifically, score $= \frac{1}{N} \sum_{i=1}^{N} \frac{\# \text{ of correct step}}{\# \text{ of total step}}$, where N is the total number of generations in the test set. A model with high scores indicates fewer hallucination issues in its generation.

Similar to the hallucination mitigation task, we apply our verifiers on Llama3-70B to help it select
the best generation among 64 options. The performance is shown in Figure 5. Llama3-70B, with
help from verifiers, performs better than itself. The performance trend under each hallucination type
aligns well with the results in Table 2 that FG-PRM performs the best among all verifiers.



14

D TAILORED RULES FOR JUDGING HALLUCINATION TYPES

We provide a prompt template for a language model to judge if the reasoning history of a given question can be incorporated into a specific type of hallucination:

Prompt Template for Hallucination Verification

[Question] {question} [Reasoning Steps] {correct reasoning steps} [Instruction] {output instruction}

In the following, we provide the rules for judging different type of hallucination:

Judgment Rules for Factual Inconsistency Hallucination

The above are step-wise reasoning steps to answer the question. Please help me determine whether the last reasoning step refers factual information not mentioned before the step. All factual information should be gounded in real-world information, including:

- Known Geographic Facts: the step should include widely accepted and verifiable facts in its original format or name. For example, state the fact that "The Eiffel Tower is located in Paris.", "Mount Everest, the tallest mountain in the world, is located in the Himalayas.", etc.
- Historical Events: the step should refer historical events with correct dates or details. For example, mention that "The American Civil War ended in 1865."

- Factual Scientific Data or Statistics: the step should include correct real-world data or statistics. But, basic calculation process should not be counted as factual information. For example, a step can state that "According to the 2020 census, the population on earth is over 7.5 billion.", "There is 7 days a week.", "The pythagorean theorem is $a^2 + b^2 = c^2$.", etc. In the output, there should be explanation whether the last reasoning step has factual information first. Then, in the new line, please only output "Yes" if the last reasoning step has factual information. Otherwise, please only output "No".

Judgment Rules for Context Inconsistency Hallucination

The above are step-wise reasoning steps to answer the question. Please help me determine whether the last reasoning step refers question information. Referred content in the last reasoning step should be the same as it mentioned in the question. Contents indirectly related to the referred content, such as derived or concluded by the referred contents, should not be counted as question information.

In the output, there should be an explanation whether the last reasoning step refers question information, output the extact referred question information in both the last reasoning step and question first. Then, in the new line, please only output "Yes" if the last reasoning step refers question information. Otherwise, please only output "No".

Judgment Rules for Logical Error Hallucination

The above are step-wise reasoning steps to answer the question. Please help me determine whether the last reasoning step involves calculation processes, including mathematical calculations or formulas:

- Mathematical Calculations: the step should have at least one calculation process. The calculation processes should include numbers (3, 5, 10 etc.) or mathematical symbols (sin, cos, x, y, π , etc.), and they should be like "The sum of 45 and 15 is 60", "30*4+5=125", "sin(x)+cos(x)", etc.

- Formulas: the step should include mathematical principles, laws of physics, or other data processing operations. Formulas may be in latex format. They can be simply stated in the step and do not have equal symbols. For example, formula can be "Pi*radius²", "2*Pi*radius", "[sin(x)+cos(x)]", etc.

In the output, there should be explanation whether the last reasoning step has calculation process first. Then, in the new line, please only output "Yes" if the last reasoning step has calculation process. Otherwise, please only output "No".

Judgment Rules for Logical Inconsistency Halluciantion

The above are step-wise reasoning steps to answer the question. Please help me determine whether the last reasoning step involves reasoning process. Referred content in the last reasoning step should be the same as it mentioned in the previous reasoning steps but not in the question. Contents indirectly related to the referred content, such as derived or concluded by the referred contents, should not be counted as previous information.

In the output, there should be an explanation whether the last reasoning step refers information in previous reasoning steps but not in the question, output the extact referred previous information in both the last reasoning step and previous step first. Then, in the new line, please only output "Yes" if the last reasoning step refers those information. Otherwise, please only output "No".

E PROMPTS FOR SYNTHESIZING HALLUCINATION DATASET

We provide a system prompt template for a language model that is used to guide the larger model through the task of generating each type of hallucination. This template consists of a description of the hallucination generation task and instructions for producing a particular type of hallucination.

System Prompt

864

865 866 867

868

870 871 872

873 874

875

876

877

878

879

882

883 884 885 Please generate the next reasoning step to incorrectly continue the reasoning process based on a question and a series of correct reasoning steps. The next reasoning step you generate does not necessarily result in an instant final answer. And you should follow the hallucination generation instruction below to generate the next reasoning step. {Instruction}

In the following, we present the instructions for generating each type of hallucination:

Instruction for Generate Fabrication Hallucination

887 # Instruction for Generating Fabrication Hallucination 888 When generating the next reasoning step, you should intentionally introduce fabrications 889 by including facts that are unverifiable against established real-world knowledge or context 890 information. These fabrications should be plausible within the context but should not be 891 verifiable through any external sources. Follow these guidelines: 892 - Unverifiable Facts: Introduce facts that cannot be verified through established real-world knowledge. For example, mention a historical event that did not happen, or a scientific 893 theory that does not exist. 894 - Fictitious Entities: Refer to people, places, or organizations that are entirely made up. For 895 example, mention a "Dr. John Smith of the International Institute of Quantum Studies," 896 which does not exist., 897 - Imaginary Data or Statistics: Provide data or statistics that are fictional. For example, state that "according to a 2023 study by the Global Health Organization, 75% of people prefer 899 digital books over physical ones," when no such study exists.', 900 # Example Guidelines 901 - Introduce a Fabricated Historical Event: For instance, state that "In 1875, the Grand Treaty 902 of Lisbon established the first international postal system," even though no such treaty exists. - Mention Nonexistent Scientific Theories or Discoveries: For example, reference "Dr. 903 Eleanor Rigby's groundbreaking work on temporal physics, which suggests that time travel 904 is theoretically possible," when no such work or scientist exists. 905 - Provide Fictitious Data or Statistics: Include statements like "A recent survey by the Na-906 tional Institute of Sleep Studies found that 60% of adults dream in black and white," even 907 though such an institute or survey does not exist. 908 # Constraints 909 - Plausibility: The fabricated content should be plausible within the context but should not 910 be verifiable. 911 - Consistency: The rest of the generated content should be consistent and coherent, without 912 introducing contradictions or errors in logic. 913 - No Contradiction to Known Facts: Avoid contradicting widely accepted and easily verifiable facts. The fabrication should be in areas that are less likely to be immediately recog-914 nized as false. 915 - Maintain Context: Ensure that the fabricated information fits smoothly into the surrounding 916 context, making it less likely to be immediately questioned. 917

18	
919	Instruction for Generating Factual Inconsistency Hallucination
920	# Instruction for Generating Factual Inconsistency Hallucination
921	When generating the next reasoning step, you should intentionally introduce factual incon-
922	sistencies by including facts that can be grounded in real-world information but present con-
923	tradictions. These inconsistencies should be subtle and should not be immediately obvious.
924	Follow these guidelines:
925	- Contradict Known Facts: Introduce information that contradicts widely accepted and ver-
926	ifiable facts. For example, state that "The Eiffel Tower is located in Berlin," contradicting
927	the well-known fact that it is in Paris.
928	- Inconsistent Historical Events: Reference historical events with incorrect dates or details.
929	in 1865
930	- Conflicting Data or Statistics: Provide data or statistics that conflict with established infor-
931	mation. For example, state that "According to the 2020 census, the population of New York
932	City is 2 million," when the actual population is significantly higher.
933	# Example Guidelines
934	- Contradict Known Geographic Facts: For instance, state that "Mount Everest, the tallest
935	mountain in the world, is located in the Andes mountain range," when it is actually in the
936	Himalayas.
937	- Inconsistent Historical Dates: For example, claim that "The Declaration of Independence
938	was signed on July 4, 1800, when it was actually signed in 1776.
939	- Connicting Scientific Information. Include statements like water bons at 110 degrees Celsius at sea level "when it actually boils at 100 degrees Celsius
940	# Constraints
941	- Plausibility: The inconsistent content should be subtle and not immediately obvious to the
942	reader.
943	- Consistency: The rest of the generated content should be consistent and coherent, without
944	introducing contradictions or errors in logic beyond the intended inconsistencies.
945	- Grounded in Real-World Information: The fabricated inconsistencies should still be based
946	on real-world information but presented inaccurately.
947	- Maintain Context: Ensure that the inconsistent information fits smoothly into the surround-
948	ing context, making it less likely to be immediately questioned.
949	(
950	
951	
952	Instruction for Generating Instruction Inconsistency Hallucination
053	instruction for Generating instruction inconsistency manuemation

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

Instruction for Generating Instruction Inconsistency Hallucination When generating the next reasoning step, you should intentionally introduce inconsistencies by not aligning the output with the specific instructions given by the user. These instruction inconsistencies should be subtle but clear enough to be identified. Follow these guidelines: - Ignore Specific Instructions: Generate text that contradicts or disregards explicit instructions given in the prompt. For example, if asked to list developed countries in Europe, list all developed countries in the world. - Alter the Requested Target: Change the target requested by the user. For example, if asked to list developed countries in the world, list all undeveloped countries in the world instead. - Misinterpret the Instructions: Deliberately misinterpret the instruction so that the output does not respond directly to the user's request. For example, if asked for "Japan's capital city", answer "Japan's largest city is Tokyo", even though Tokyo is the largest city in Japan. # Constraints

- Faithful: You cannot fabricate something that doesn't appear in the context.

- Coherence: The rest of the generated content should remain coherent and logical, without introducing contradictions or errors beyond the intended inconsistencies.

- Contextual Fit: Ensure that despite the inconsistency, the response still fits smoothly within the broader context of the conversation or text, making it less likely to be immediately questioned.

972	
0 = 0	Instruction for Generating Context Inconsistency Hallucination
973	
974	# Instruction for Generating Context Inconsistency Hallucination
975	When generating the next reasoning step, you should introduce inconsistencies by intention-
970	ally modifying information to contradict the user's provided contextual information. These context inconsistencies should be subtle but clear enough to be identified. Follow these
977	guidelines:
970	- Contradict Provided Facts: Introduce information that directly contradicts the facts given
979	in the user's prompt. For example, if the user states that "Bob was born in England," you
900	may contradict it by stating that "Bob was born in France."
901	- Alter Specific Details or Data: Change specific details or data provided by the user. For
902	example, if the user mentions that "Bob has three books and two pens in his backpack," you
903	might alter it by stating that "Bob has two books and four pens in his backpack."
984	- Misattribute Quotes or Data: Attribute quotes or data to the wrong source. For example, if
985	the user states that "Bob likes apples while Jane likes bananas." you might contradict it by
980	stating "Jane likes apples" or "Bob likes bananas".
987	# Constraints
988	- Subtlety: The context inconsistencies should be subtle and not immediately obvious to the
989	Coherence: The rest of the generated content should remain coherent and logical without
990	introducing contradictions or errors beyond the intended inconsistencies
991	- Contextual Fit: Ensure that the inconsistent information fits smoothly within the broader
992	context of the conversation or text, making it less likely to be immediately questioned.
993	
994	
995	
007	
008	
990	
1000	
1000	
1001	
1002	
1002	
1002 1003 1004	
1002 1003 1004 1005	Instruction for Generating Logical Inconsistency Hallucination
1002 1003 1004 1005 1006	Instruction for Generating Logical Inconsistency Hallucination
1002 1003 1004 1005 1006 1007	Instruction for Generating Logical Inconsistency Hallucination # Instruction for Generating Logical Inconsistency Hallucination
1002 1003 1004 1005 1006 1007 1008	Instruction for Generating Logical Inconsistency Hallucination # Instruction for Generating Logical Inconsistency Hallucination When generating the next reasoning step, you should introduce logical inconsistencies by
1002 1003 1004 1005 1006 1007 1008 1009	Instruction for Generating Logical Inconsistency Hallucination # Instruction for Generating Logical Inconsistency Hallucination When generating the next reasoning step, you should introduce logical inconsistencies by incorrectly referring to or copying content from previous reasoning steps. These logical inconsistencies should be subtle but clear around to be identified. Follow these suidelines:
1002 1003 1004 1005 1006 1007 1008 1009 1010	Instruction for Generating Logical Inconsistency Hallucination # Instruction for Generating Logical Inconsistency Hallucination When generating the next reasoning step, you should introduce logical inconsistencies by incorrectly referring to or copying content from previous reasoning steps. These logical inconsistencies should be subtle but clear enough to be identified. Follow these guidelines: Incorrect Pafarence: Pafar to a previous reasoning step incorrectly such as miniateresting
1002 1003 1004 1005 1006 1007 1008 1009 1010 1011	 Instruction for Generating Logical Inconsistency Hallucination # Instruction for Generating Logical Inconsistency Hallucination When generating the next reasoning step, you should introduce logical inconsistencies by incorrectly referring to or copying content from previous reasoning steps. These logical inconsistencies should be subtle but clear enough to be identified. Follow these guidelines: Incorrect Reference: Refer to a previous reasoning step incorrectly, such as misinterpreting or misrepresenting the calculations or conclusions. For example, if a previous states
1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012	Instruction for Generating Logical Inconsistency Hallucination # Instruction for Generating Logical Inconsistency Hallucination When generating the next reasoning step, you should introduce logical inconsistencies by incorrectly referring to or copying content from previous reasoning steps. These logical inconsistencies should be subtle but clear enough to be identified. Follow these guidelines: - Incorrect Reference: Refer to a previous reasoning step incorrectly, such as misinterpreting or misrepresenting the calculations or conclusions. For example, if a previous step states "Bob is an undergraduate," you may incorrectly refer back to this by stating "Since Bob is a
1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012 1013	Instruction for Generating Logical Inconsistency Hallucination # Instruction for Generating Logical Inconsistency Hallucination When generating the next reasoning step, you should introduce logical inconsistencies by incorrectly referring to or copying content from previous reasoning steps. These logical inconsistencies should be subtle but clear enough to be identified. Follow these guidelines: - Incorrect Reference: Refer to a previous reasoning step incorrectly, such as misinterpreting or misrepresenting the calculations or conclusions. For example, if a previous step states "Bob is an undergraduate," you may incorrectly refer back to this by stating "Since Bob is a graduate"
1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012 1013 1014	Instruction for Generating Logical Inconsistency Hallucination # Instruction for Generating Logical Inconsistency Hallucination When generating the next reasoning step, you should introduce logical inconsistencies by incorrectly referring to or copying content from previous reasoning steps. These logical inconsistencies should be subtle but clear enough to be identified. Follow these guidelines: - Incorrect Reference: Refer to a previous reasoning step incorrectly, such as misinterpreting or misrepresenting the calculations or conclusions. For example, if a previous step states "Bob is an undergraduate," you may incorrectly refer back to this by stating "Since Bob is a graduate" - Copying Errors: Copy content from a previous reasoning step but alter it in a way that
1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012 1013 1014 1015	 Instruction for Generating Logical Inconsistency Hallucination # Instruction for Generating Logical Inconsistency Hallucination When generating the next reasoning step, you should introduce logical inconsistencies by incorrectly referring to or copying content from previous reasoning steps. These logical inconsistencies should be subtle but clear enough to be identified. Follow these guidelines: Incorrect Reference: Refer to a previous reasoning step incorrectly, such as misinterpreting or misrepresenting the calculations or conclusions. For example, if a previous step states "Bob is an undergraduate," you may incorrectly refer back to this by stating "Since Bob is a graduate" Copying Errors: Copy content from a previous reasoning step but alter it in a way that introduces an error, such as changing numbers or relationships. For example, if the reasoning
1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012 1013 1014 1015 1016	 Instruction for Generating Logical Inconsistency Hallucination # Instruction for Generating Logical Inconsistency Hallucination When generating the next reasoning step, you should introduce logical inconsistencies by incorrectly referring to or copying content from previous reasoning steps. These logical inconsistencies should be subtle but clear enough to be identified. Follow these guidelines: Incorrect Reference: Refer to a previous reasoning step incorrectly, such as misinterpreting or misrepresenting the calculations or conclusions. For example, if a previous step states "Bob is an undergraduate," you may incorrectly refer back to this by stating "Since Bob is a graduate" Copying Errors: Copy content from a previous reasoning step but alter it in a way that introduces an error, such as changing numbers or relationships. For example, if the reasoning involves steps for calculating a total cost and one step states "Item A costs 5 * 2 =10," you
1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012 1013 1014 1015 1016 1017	 Instruction for Generating Logical Inconsistency Hallucination # Instruction for Generating Logical Inconsistency Hallucination When generating the next reasoning step, you should introduce logical inconsistencies by incorrectly referring to or copying content from previous reasoning steps. These logical inconsistencies should be subtle but clear enough to be identified. Follow these guidelines: Incorrect Reference: Refer to a previous reasoning step incorrectly, such as misinterpreting or misrepresenting the calculations or conclusions. For example, if a previous step states "Bob is an undergraduate," you may incorrectly refer back to this by stating "Since Bob is a graduate" Copying Errors: Copy content from a previous reasoning step but alter it in a way that introduces an error, such as changing numbers or relationships. For example, if the reasoning involves steps for calculating a total cost and one step states "Item A costs 5 * 2 =10," you might incorrectly copy this as "Since item A costs 5 * 3 =15" in the next step.
1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012 1013 1014 1015 1016 1017 1018	 Instruction for Generating Logical Inconsistency Hallucination # Instruction for Generating Logical Inconsistency Hallucination When generating the next reasoning step, you should introduce logical inconsistencies by incorrectly referring to or copying content from previous reasoning steps. These logical inconsistencies should be subtle but clear enough to be identified. Follow these guidelines: Incorrect Reference: Refer to a previous reasoning step incorrectly, such as misinterpreting or misrepresenting the calculations or conclusions. For example, if a previous step states "Bob is an undergraduate," you may incorrectly refer back to this by stating "Since Bob is a graduate" Copying Errors: Copy content from a previous reasoning step but alter it in a way that introduces an error, such as changing numbers or relationships. For example, if the reasoning involves steps for calculating a total cost and one step states "Item A costs 5 * 2 =10," you might incorrectly copy this as "Since item A costs 5 * 3 =15" in the next step. Make logical leaps or conclusions that do not follow from the previous steps, leading to an

Constraints

1019

1020

1021

- Subtlety: The logical inconsistencies should be subtle and not immediately obvious to the reader.

1022 - Coherence: The rest of the generated content should remain coherent and logical, without
 1023 introducing contradictions or errors beyond the intended inconsistencies.

Contextual Fit: Ensure that the inconsistent information fits smoothly within the broader context of the conversation or text, making it less likely to be immediately questioned.

1026	Instruction for Generating Calculation Error Hallucination
1027	
1020	# Instruction for Generating Calculation Error Hallucination
1023	When generating the next reasoning step, you should intentionally introduce calculation
1030	be subtle but clear enough to be identified. Follow these guidelines:
1031	De subile but clear chough to be fucilitied. Follow these guidelines.
1032	calculations. For example, state that "The sum of 45 and 15 is 70" when it is actually 60
1033	- Include Incorrect Data Processing: Misapply mathematical principles laws of physics or
1034	other data processing operations. For example, when asked to calculate the area of a circular.
1035	compute the perimeter formula $2*Pi*radius$ instead of the area formula $Pi*radius^2$.
1036	- Generates responses with unsupported claims, including numerical assertions that have no
1037	basis in the provided context or input.
1038	# Constraints
1039	- The values you use must be consistent with the context given, but the final calculation
1040	should be intentionally miscalculated.
1041	- You must not fabricate what does not appear in the context or contradict widely accepted
1042	and easily verifiable facts.
1043	- Ensure that despite the errors, the response still fits smoothly within the broader context of
1044	
1045	
1046	Instruction for Generating Logical Error Hallucination
1047	# Instruction for Generating Logical Error Hallucination
1040	When generating the next reasoning step, you should intentionally introduce logical error by
1049	including flawed logical reasoning or incorrect inferences. These errors should be subtle but
1050	clear enough to be identified. Follow these guidelines:
1051	- Causal Misattribution: Incorrectly identify the cause of an event or outcome. For example,
1052	conclude that "Because it rained yesterday, that's why the football team won today's match,"
1053	without considering other relevant factors.
1054	- Overgeneralization: Apply a rule or pattern more broadly than it should be. For instance,
1055	generalize that "all mammals fly" based on the fact that bats are flying mammals.
1050	- Generate responses with unsupported claims, including assertions that do not logically
1057	# Constraints
1050	- The information you refer to must be consistent with the information provided in the previ-
1059	ous reasoning steps and context, but the final conclusion should be intentionally and logically
1061	flawed.
1060	- You must not fabricate what does not appear in the context or contradict widely accepted
1062	and easily verifiable facts.
1067	
1065	
1065	
1067	
1069	
1060	
1005	
1070	
1072	
1072	
107/	
1075	
1075	
1070	
1079	
1070	

1080 F HALLUCINATION DEMONSTRATION EXAMPLES 1081

- 1082 1083 1084 1086
- 1087
- 1088 1089
- 1090
- 1091
- 1093 1094

Demonstrations for Fabrication Hallucination

- 1095
 - [Ouestion] What are the primary components of DNA?
 - [Correct Reasoning Steps]
- 1098 Step 1: DNA is structured as a double helix composed of nucleotides.
- 1099 Step 2: Each nucleotide consists of a sugar (deoxyribose), a phosphate group, and a nitrogenous base. 1100

We provide demonstrations for generating each type of hallucination. Each demonstration includes

two examples of an injecting hallucination, along with an explanation of how it is produced.

- Step 3: The four nitrogenous bases are adenine (A), thymine (T), cytosine (C), and guanine 1101 (G). 1102
- {output format} 1103
- [Explanation] 1104
- The user is asking about the primary components of DNA. The correct approach is to de-1105 scribe the structure of DNA and its components, including the nucleotides and the four 1106 nitrogenous bases. The Next Reasoning Step here introducesFabrication Hallucination by 1107 mentioning a "recent study by the Molecular Genetics Institute in Zurich" that identified a 1108 fifth nitrogenous base, "neomine (N)," which does not exist. This reasoning step remains 1109 coherent and logical, correctly describing the structure of DNA and its primary components, 1110 but introducing a fictitious base and study that is not supported by established scientific knowledge. 1111 [Next Reasoning Step with Fabrication Hallucination] 1112
- Step 4: Additionally, a recent study by the Molecular Genetics Institute in Zurich has identi-1113 fied a rare fifth base, neomine (N), which appears to play a role in DNA repair mechanisms. 1114 [Ouestion]
- 1115 What are the main causes of the decline of the Roman Empire? 1116 [Correct Reasoning Steps]
- 1117 Step 1: One of the major causes of the decline of the Roman Empire was economic trouble, 1118 which included heavy taxation and inflation.
- 1119 Step 2: Additionally, the overreliance on slave labor created economic stagnation.
- Step 3: Military overspending further strained the empire's resources. 1120
- {output format} 1121
- [Explanation] The user is asking about the main causes of the decline of the Roman Empire. 1122 The correct approach is to identify a combination of internal weaknesses and external pres-1123 sures that contributed to the empire's fall, along with key events that marked its decline. The 1124 Next Reasoning Step here introduces Fabrication Hallucination by mentioning the "Great 1125 Plague of 330 AD," which did not actually occur, and the Roman historian "Lucius Mar-1126 cellus," who is not a verifiable historical figure. This reasoning step remains coherent and 1127 logical, offering a plausible explanation for the decline of the Roman Empire, but introduc-1128 ing a fictitious event and historian that can not verified by known historical facts. 1129 [Next Reasoning Step with Fabrication Hallucination] Step 4: A significant but lesser-known factor was the Great Plague of 330 AD, which deci-1130
- 1131 mated the population and weakened the military. This pandemic, documented by the Roman historian Lucius Marcellus, is believed to have killed nearly a quarter of the empire's popu-1132 lation. 1133

	Demonstrations for Context Inconsistency Hallucination
	[Question]
	Ann's lavorite store was having a summer clearance. For $5/5$ she bought 5 pairs of she
	For \$7 each and 2 pairs of shoes for \$10 each. She also bought 4 tops, all at the same pr.
1	[Correct Desconing Steps]
	[Concer Reasoning Steps] Step 1: She bought 5 shorts at \$7 each so 5*7-\$35
	Step 7: She bought 2 pair of shoes at \$10 each so 2×10^{-10}
	Step 2: She sought 2 pair of shoes at \$10 each so 2 $10-$20$ Step 3: The shorts and shoes cost her $35+20 = 55
	Step 4: We know she spent 75 total and the shorts and shoes cost \$55 which left a differe
	of $75-55 = \$20$
	{output format}
	[Explanation]
	The user is asking about how much each top cost based on the total amount spent.
	correct approach is to calculate the total cost of shorts and shoes first, then determine
	remaining amount spent on tops, and finally find the cost per top by dividing the remain
	amount by the number of tops. The Next Reasoning Step here introduces Context Inc
	sistency Hallucination by incorrectly stating that Ann bought 5 tops instead of the 4 to
	mentioned in the user's provided context, which creates a context inconsistency by alter
	the specific detail about the number of tops. This reasoning step maintains coherent
	logical in its structure, providing a plausible answer for the price of each top, but introdu
	a subtle context inconsistency in the number of tops that affects the final answer.
	[Next Reasoning Step with Context Inconsistency Hallucination]
	Step 5: Since she bought 5 tops for a total of \$20, so each top $\cos 20/5 = 4
	[Question] There are 16 dolphins in a water park. One quarter of them are gray dolphins, and there
	twice as many white dolphins as gray dolphins. How many white dolphins are in the w
	park?
	[Correct Reasoning Steps]
	Step 1: We know one quarter of the dolphins are gray, we calculate this by dividing the t
	number of dolphins by 4, $16/4 = 4$
	{output format}
	[Explanation]
	The user is asking about the number of white dolphins in the water park, given the t
	number of dolphins and the proportion of gray dolphins. The correct approach is to de
	mine the number of gray dolphins first, then use the given ratio to find the number of w
	dolphins. The Next Reasoning Step here introduces Context Inconsistency Hallucination
	incorrectly stating that there are twice as many gray dolphins as white dolphins, instea
	correctly stating that there are twice as many white dolphins as gray dolphins. This cre
	a context inconsistency by altering the specific relationship between the number of gray white delphing. This reasoning star remains schement and logical in its structure provide
	white dolphins. This reasoning step remains concretent and logical in its structure, providing a plausible method to find the number of white dolphing, but introduces a subtle context inc
	plausible method to find the number of white dolphins, but introduces a subtle context inc
	sistency in the relationship between the number of gray and white doiphins that contra provided facts
	[Next Reasoning Step with Context Inconsistency Hallucination]
	Step 2: Since there are twice as many grey dolphins as white dolphins, we divide the num
	of gray dolphins by 2 to get the number of white dolphins: $4/2=2$.
	$\sigma_{1} \equiv \sigma_{1} = \sigma_{2} = \sigma_{1} = \sigma_{2} = \sigma_{1} = \sigma_{1$

1188	Demonstrations for Instruction Inconsistency Hallucination
1189	
1190	[Question]
1191	Adam bought 3 kilograms of nuts and 2.5 kilograms of dried fruits at a store. One kilogram
1192	of nuts costs \$12 and one kilogram of dried fruit costs \$8. How much did his purchases
1193	Cost ?
1194	Step 1: For the pute Adam paid 3 * 12 - 36
1195	Step 2: And for dried fruits Adam paid $5^{-12} = 50$.
1190	{output format}
1100	[Explanation]
1100	The user is asking for the total cost of Adam's purchases. The correct next reasoning step
1200	should add the costs of the nuts and dried fruits to find the total cost. The Next Reasoning
1200	Step here introduces Instruction Inconsistency Hallucination by calculating the average cost
1201	of the purchases instead of finding the total cost, altering the requested target. Despite the
1202	inconsistency, this reasoning step introduces no contradictions or errors in logic, and still fits
1203	smoothly within the broader context of the conversation.
1204	Step 3: To find the average cost of Adam's purchases, we can add the cost of puts and dried
1205	fruits and divide by 2: $(\$36 + \$20)/2 = \$28$
1200	[Ouestion]
1207	Abigail is trying a new recipe for a cold drink. It uses 14 of a cup of iced tea and 1 and 14 of
1200	a cup of lemonade to make one drink. If she fills a pitcher with 18 total cups of this drink,
1209	how many cups of lemonade are in the pitcher?
1210	[Correct Reasoning Steps]
1211	Step 1: Each drink uses 1.5 cups because $14 \text{ cup} + 1$ and $14 \text{ cup} = 1.5 \text{ cups}$
1212	Step 2: The pitcher contains 12 total drinks because $18 \ 1.5 = 12$
1213	{output format}
1215	[Explanation] The user is asking the number of curs of lemonade in the nitcher. The next correct reasoning
1216	step should calculate the total cups of lemonade by multiplying the number of drinks by the
1217	amount of lemonade per drink. The Next Reasoning Step here introduces Instruction Incon-
1218	sistency Hallucination by suddenly changing the unit of measurement from cups to ounces,
1219	ignoring the specific instruction to find the number of cups. Despite the inconsistency, this
1220	reasoning step introduces no contradictions or errors in logic, and still fits smoothly within
1221	the broader context of the conversation.
1222	[Next Reasoning Step with Instruction Inconsistency Hallucination]
1223	Step 3: Since each drink uses 1 and 1/4 cups of lemonade, and there are 8 ounces in a cup, the total suppose of lemonade in the nitebox are $12 \times (1 \text{ and } 1/4) \times 8 = 06$ suppose
1224	the total ounces of lemonade in the pitcher are $12 \text{ "(1 and 1/4) " 8 = 96 ounces.}$
1225	
1226	
1227	
1228	
1229	
1230	
1231	
1232	
1233	
1234	
1235	
1236	
1237	
1238	
1239	
1240	

1242	Demonstrations for Logical Inconsistency Hallucination
1243	
1244	[Question]
1245	Annie, Bob, and Cindy each got some candy. Annie has 6 candies, Bob has 2 candies more
1246	than half of Annie's candies, and Cindy has 2 candies less than twice Bob's candies. Which
1247	of the three of them has the least amount of candy?
1248	[Correct Reasoning Steps] Step 1: Appie has 6 capities
1249	Step 1. Annue has 0 candies. Step 2: Bob has 2 candies more than half of Annie's candies. Half of Annie's candies is (6
1250	(2 = 3) So Bob has $(3 + 2 = 5)$ candies
1251	Step 3: Cindy has 2 candies less than twice Bob's candies. Twice Bob's candies is $(2 * 5 =$
1252	10). So, Cindy has $(10 - 2 = 8)$ candies.
1253	
1254	{output format}
1255	[Explanation] The user is asking which of Annie, Bob, and Cindy has the least amount of
1256	candy. The correct approach is to calculate the number of candies each person has and then
1257	compare these amounts to determine who has the least. According to the previous steps: 1.
1258	Annie has 6 candies; 2. Bob has 5 candies; 3. Cindy has 8 candies. The Next Reasoning Step
1259	here introduces Logical Inconsistency Hallucination by incorrectly concluding that Annie
1260	has the least amount of candy, whereas the correct conclusion should be that Bob has the
1261	least amount of candy with 5 candles. This creates a logical inconsistency by failing to
1262	reasoning steps
1263	[Next Reasoning Step with Logical Inconsistency Hallucination]
1264	Step 4: Since Annie only has 6 candies. Anne has the least amount of candy.
1265	[Question]
1266	Annie, Bob and Cindy each buy personal pan pizzas cut into 4 pieces. If Bob eat 50% of his
1267	pizzas and Ann and Cindy eat 75% of the pizzas, how many pizza pieces are left uneaten?
1268	[Correct Reasoning Steps]
1269	Step 1: In total, there are $3 * 4 = 12$ pizza pieces. Step 2: Bob eats $4 * 50\% = 2$ pieces. Step
1270	3: Annie and Cindy eat $2 * 4 * 75\% = 6$ pieces. Step 4: The three of them eat $2 + 6 = 8$
1271	pieces.
1272	{output format]}
1273	[Explanation] The user is asking how many pizza pieces are left uneaten after Annie, Bob and Cindy each
1274	eat a portion of their pizza. The correct approach is to calculate the total number of pizza
1275	pieces, determine how many pieces each person eats, and then find the remaining uneaten
1276	pieces. According to the previous steps: 1. In total, there are 12 pizza pieces; 2. Bob eats
1277	2 pieces; 3. Annie and Cindy together eat 6 pieces; 4. Therefore, the three of them eat 2 + 6
1278	= 8 pieces. The Next Reasoning Step here introduces Logical Inconsistency Hallucination
1279	by incorrectly copying that 10 pieces of pizza were eaten and by incorrectly referencing the
1280	total number of pizza pieces as 16, whereas the correct calculation should be based on the
1281	total number of 12 pizza pieces and the remaining uneaten pieces should be $12 - 8 = 4$. This
1282	creates a logical inconsistency by incorrectly referencing the number of eaten pieces as 10 and the total number of ninge ninges as 16, contradicting the provides reasoning store
1283	In the total number of pizza pieces as 10, contradicting the previous reasoning steps.
1284	Step 5: Since 10 pizza pieces were eaten, there are $16 - 10 - 6$ pizza pieces upgaten
1285	Step 5. Since 10 pizza pieces were eaten, alere are 10 10 - 0 pizza pieces aneaten.
1286	
1287	
1288	
1289	
1290	
1291	
1292	
1293	

1296	Demonstrations for Calculation Error Hallucination
1297	
1298	[Question]
1299	Abigail is trying a new recipe for a cold drink. It uses 0.25 of a cup of iced tea and 1.25 of
1300	a cup of lemonade to make one drink. If she fills a pitcher with 18 total cups of this drink,
1301	now many cups of lemonade are in the pitcher?
1302	Step 1: Each drink uses 1.5 cups because 0.25 cup + 1 and 1.25 cup = 1.5 cups
1303	Step 1. Each drink uses 1.5 cups because 0.25 cup + 1 and 1.25 cup = 1.5 cups Step 2: The pitcher contains 12 total drinks because $18 / 15 = 12$
1304	Step 2. The pitcher contains 12 total driftes because $1671.5 - 12$
1305	[Fxn]anation]
1306	The user is asking for the total cups of lemonade in the pitcher. The correct calculation
1307	involves finding the total number of drinks and then multiplying it by the amount of lemon-
1308	ade per drink. The Next Reasoning Step here introduces Calculation Error Hallucination by
1309	incorrectly calculating the total cups of lemonade as $12 * 1.25 = 17$ cups when it is actually
1310	12 * 1.25 = 15 cups. Despite the incorrect calculation, this reasoning step correctly iden-
1311	tifies the values to be multiplied, and still fits smoothly within the broader context of the
1312	conversation.
1313	[Next Reasoning Step with Calculation Error Hallucination]
1314	Step 3: Since each drink uses 1.25 cups of lemonade, multiply the total drinks by the amount
1315	of lemonade per drink: $12 + 1.25 = 1 / \text{ cups of lemonade}$.
1316	[Question] Queania come \$150 a day as a part time clark. She come on additional \$5 per hour oc
1317	Queenie earns \$150 a day as a part-unie cierk. She earns an additional \$5 per nour as
1318	[Correct Reasoning Steps]
1319	Step 1: Queenie will earn $$150 * 5 = 750 for working 5 days Step 2: She will receive an
1320	additional $\$5 * 4 = \20 for overtime pay.
1321	{output format}
1322	[Explanation]
1323	The user is asking for Queenie's total earnings for working 5 days with 4 hours of overtime.
1324	The correct calculation involves adding her daily earnings to her overtime pay. The Next
1325	Reasoning Step here introduces Calculation Error Hallucination by incorrectly adding \$750
1326	and \$20 as \$7800 when it is actually \$770. Despite the incorrect calculation, this reasoning
1327	step correctly identifies the values to be added, and still fits smoothly within the broader
1328	context of the conversation.
1329	[Next Reasoning Step with Calculation Error Hallucination] Step 2: Then, add her doily carrings and evertime new to get her total corringes \$750 + \$20
1330	step 5. Then, add her daily earnings and overtime pay to get her total earnings. $\frac{5750 + 520}{-57800}$
1331	[Ouestion]
1332	What is the volume of a cylinder with a radius of 3 units and a height of 5 units?
1333	[Correct Reasoning Steps]
1334	Step 1: The volume of a cylinder is calculated using the formula Volume = $\pi \times \text{radius}^2 \times$
1335	height.
1336	Step 2: For a cylinder with a radius of 3 units and a height of 5 units, first calculate $\pi \times 3^2 =$
1337	9π .
1338	{output format}
1339	[Explanation]
1340	The user is asking for the volume of a cylinder. The correct formula involves multiplying π
1341	by the square of the radius and then by the height. Given the radius is 3 units and the height
1342	is 5 units, the volume should be calculated as $\pi \times 3^2 \times 5 = 45\pi$. The Next Reasoning Step
1342	nere introduces Calculation Error Hallucination by incorrectly calculating 9π multiplied by 5 as 18π when it is not calculated. this measurements
13/1/	5 as 10% when it is actually 40%. Although the infair result is miscalculated, this reasoning step correctly identifies the values to be multiplied, and still fits smoothly within the broader
12/5	context of the conversation
1343	[Next Reasoning Step with Calculation Error Hallucination]
1340	Step 3: Then multiply by 5, and the volume is $9\pi \times 5 = 18\pi$ cubic units.
1047	
1340	

1350	Demonstrations for Logical Error Hallucination
1352	
1353	[Question]
1354	following sentences
1355	Ouestion: Keyin is observing the sky on a clear night. With the unaided eve he is able to
1356	see Venus Mars Juniter and Saturn Why would Venus appear to be brighter than the other
1057	nlanets?
1050	Information: sent1: more light gets reflected on highly reflective things sent2: venus is
1000	covered in highly reflective clouds sent3: as the light reflected off of an object increases,
1009	the object will appear to be brighte
1300	[Correct Reasoning Steps]
1361	Step 1: More light gets reflected on highly reflective things.
1362	Step 2: Venus is covered in highly reflective clouds.
1363	{output format}
1364	[Explanation]
1365	The user is asking why Venus appears brighter than other planets in the night sky. The correct
1366	reasoning involves recognizing that Venus has highly reflective clouds, which contribute to
1367	its brightness. However, other factors like its relative proximity to Earth and its position in
1368	the sky also play significant roles. The Next Reasoning Step here introduces a Logical Error
1369	Hallucination by incorrectly inferring that because Venus is highly reflective, it must also
1370	be the closest planet to Earth. This is a flawed causality because the reflectivity of Venus's
1371	clouds does not determine its distance from Earth, which leads to an incorrect conclusion
1372	Next Descenting Step with Calculation Error Hallycinetion
1373	Step 3: Since Venus is covered in highly reflective clouds, it reflects more light than any
1374	other planet making it the closest planet to Earth and therefore the brightest in the night sky
1375	[Ouestion] Please answer the question in detail and accurately based on the information
1376	provided in the following sentences.
1377	Ouestion: Which form of energy is needed to change water from a liquid to a gas?
1378	Information: sent1: gas is a kind of state of matter sent2: water is a kind of substance sent3:
1379	liquid is a kind of state of matter sent4: heat energy can change the state of matter
1380	[Correct Reasoning Steps]
1381	Step 1: Gas is a kind of state of matter.
1382	Step 2: Liquid is a kind of state of matter.
1383	Step 3: Heat energy can change the state of matter.
1384	{output format}
1385	[Explanation]
1386	The user is asking which form of energy is needed to change water from a liquid to a gas.
1387	of matter specifically from liquid to gas. However the Next Bessening Step here introduces
1200	of matter, specifically from inquitito gas. However, the Next Reasoning Step here inforduces
1200	a Logical Error Handemation by incorrectly interning that since heat energy can change the
1209	reasoning error because while heat energy specifically is required for this phase change not
1390	all forms of energy are applicable. The overgeneralization in this step leads to an incorrect
1391	conclusion about the types of energy that can achieve this state change.
1392	[Next Reasoning Step with Logical Error Hallucination]
1393	Step 4: Since heat energy can change the state of water from a liquid to a gas, any form of
1394	energy can be used to change water from a liquid to a gas.
1395	
1396	
1397	
1398	
1399	
1400	
1401	