

# Entrospect: Information-Theoretic Self-Reflection Elicits Better Response Refinement of Small Language Models

Anonymous ACL submission

## Abstract

Self-reflection helps de-hallucinate Large Language Models (LLMs). However, the effectiveness of self-reflection remains insufficiently validated in the context of Small Language Models (SLMs), which exhibit limited semantic capacities. In particular, we demonstrate that the conventional self-reflection paradigm, such as Self-Refine, fails to deliver robust response refinement for models with parameter sizes of 10 billion or smaller, even when compared to generations elicited through Chain-of-Thought (CoT) prompting. To improve SLMs' self-reflection, we redesign Self-Refine and introduce *Entrospect*<sup>1</sup> (Entropy-aware Introspection), an information-theoretic framework based on prompt engineering.

We evaluated *Entrospect* using accuracy and average time consumption metrics to comprehensively assess its precision and computational efficiency. Experiments conducted across four distinct SLMs and four baseline methods demonstrate that *Entrospect* achieves state-of-the-art performance on validation tasks. Notably, under identical model and data settings, *Entrospect* delivers a remarkable improvement of up to 36.2% in reasoning accuracy while enhancing computational efficiency by as much as 10 times compared to its predecessor, Self-Refine.

## 1 Introduction

Large Language Models have advanced rapidly, impacting many fields with improved natural language generation (Brown et al., 2020; Chang et al., 2024). However, their tendency to produce hallucinations—especially counterfactual ones—poses a critical challenge to reliability (Zhang et al., 2023; Huang et al., 2023). Hallucinations occur when models generate factually incorrect or nonsensical outputs, undermining their trustworthiness and hindering real-world adoption. Addressing this issue is essential for improving their practical utility and acceptance (Weidinger et al., 2021, 2022).

<sup>1</sup>The project is intended to be open-source soon after the publication. For reviewers, we attached the examples of *Entrospect*'s outputs to the submission.

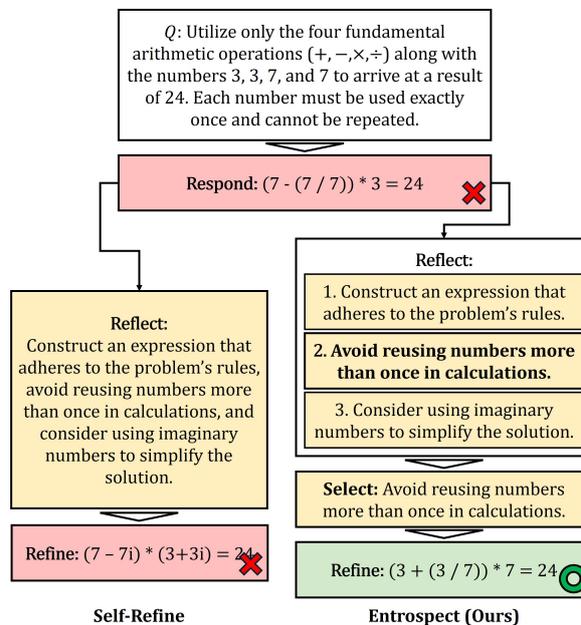


Figure 1: The single-round refinement of an initial response for the same query, comparing Self-Refine and our proposed *Entrospect*. Self-Refine fully relies on the model's self-reflected feedback, where any biases introduced during *reflect* are directly carried over into *refine*, hindering constructive improvements. On the other hand, our *Entrospect* identifies the optimal revision suggestion from an itemized output of the self-reflection, enabling *Entrospect* to achieve more robust and reliable response refinement.

To address these challenges, self-reflection has been proposed as a solution to counterfactual hallucinations, particularly for black-box models with inaccessible parameters (Madaan et al., 2024). However, its effectiveness is limited in Small Language Models (SLMs), which often lack sufficient semantic capabilities, inducing frequent occurrences of imperfect feedback, encompassing the self-reflected revision suggestions. Given the widespread use of SLMs in resource-constrained environments (Li et al., 2024; Wang et al., 2024), this limitation is particularly significant. In such cases, self-reflection may fail to consistently assist in the corrections of outputs, highlighting the need for more robust and scalable approaches.

Given the challenges of applying self-reflection to SLMs, a key question arises: *how might we con-*

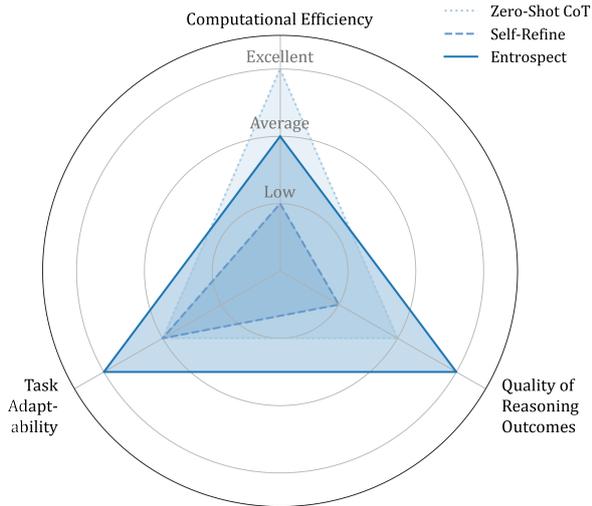


Figure 2: Entrospect contributes furtherance to the response refinement of SLMs particularly over its predecessor, Self-Refine, across three major aspects.

060 *struct a framework that effectively integrates self-*  
 061 *reflection to improve the precision of SLM outputs,*  
 062 *all while preserving the computational efficiency?*  
 063 In response to this challenge, we propose *Entro-*  
 064 *spect*, an information-theoretic framework predi-  
 065 cated on Self-Refine that lessens the dependency on  
 066 explicit semantic outputs from the model. Contrary  
 067 to Self-Refine’s equal consideration of all revision  
 068 suggestions, Entrospect employs an unsupervised  
 069 mechanism to identify the most effective revision  
 070 candidate, minimizing the impact of inferior ones,  
 071 as illustrated in Figure 1.

072 Specifically, Entrospect is implemented with an  
 073 Optimal Revision Suggestion Selector (ORSS) In-  
 074 spired by (Wu et al., 2024) and (Yang et al., 2024b),  
 075 the ORSS intervenes between the “reflect” and the  
 076 “refine” stages that are tightly connected in the Self-  
 077 Refine’s pipeline. It evaluates revision suggestions  
 078 generated through self-reflection and identifies the  
 079 one that minimizes the semantic uncertainty in the  
 080 model’s refinement of the prior response, where  
 081 low-quality suggestions conceivably ruining the  
 082 successive procedures are ruled out. This selective  
 083 approach distinguishes Entrospect from its prede-  
 084 cessors, enhancing both the quality and reliability  
 085 of the refined responses.

086 Architecturally, Entrospect retains the simplic-  
 087 ity and efficiency of Self-Refine, operating as  
 088 a parameter-free, recurrent finite-state machine  
 089 (FSM) where modules are interconnected through  
 090 purpose-specific prompts. This design ensures  
 091 computational efficiency while maintaining the  
 092 flexibility to adapt to diverse conversational AI

093 tasks. Figure 2 summarizes the multifaceted con-  
 094 tributions of Entrospect, the central focus of this  
 095 study.

096 We evaluated Entrospect on natural lan-  
 097 guage reasoning tasks, including the MATH  
 098 dataset (Hendrycks et al., 2021) for math reason-  
 099 ing and HaluEval (Li et al., 2023) for hallucina-  
 100 tion detection. The results show Entrospect outper-  
 101 forms baselines like zero-shot, few-shot, Chain-of-  
 102 Thought (CoT), and Self-Refine. These findings  
 103 underscore two critical advances:

- 104 1. **Selective Use of Self-Reflection:** We high-  
 105 light that the outcomes of a model’s self-  
 106 reflection should not be directly or entirely  
 107 relied upon as guidance for the response re-  
 108 finement.
- 109 2. **ORSS-Driven Optimization:** Our proposed  
 110 Entrospect improves Self-Refine by introduc-  
 111 ing ORSS, an information-theoretic mecha-  
 112 nism that unsupervisedly identifies the opti-  
 113 mal revision from multiple candidates. Com-  
 114 bined with our semantic similarity-based stop-  
 115 ping condition, Entrospect allows a more ro-  
 116 bust and systematic approach to self-reflection  
 117 for response refinement. Compared to its pre-  
 118 cursor, Self-Refine, Entrospect accomplishes  
 119 a remarkable performance boost, delivering  
 120 up to 36.2% improvement in accuracy under  
 121 identical dataset and model conditions, while  
 122 elevating computational efficiency by as much  
 123 as 10 times.

## 124 2 Related Work

### 125 2.1 Self-Reflection of Language Models

126 The empirical foundation of self-reflection is that  
 127 given some queries, language models may not be  
 128 able to provide proper answers every time (Yan  
 129 and Xu, 2023). Self-reflection assists in alleviating  
 130 such problems by explicitly instructing a language  
 131 model to review its generated response, providing  
 132 a feedback on potential deficiencies within the cur-  
 133 rent response and how they could be eliminated.  
 134 The feedback is subsequently used for guiding the  
 135 refinement of the previous answer. This procedure  
 136 can be fully automated through a prompt-driven  
 137 framework, by which a language model iteratively  
 138 reflects and refines the answer to a query on its  
 139 own (Lee et al., 2024).

140 Techniques like Self-Refine introducing mech-  
 141 anisms for models to improve their own re-

sponses (Madaan et al., 2024), especially in question-answering (QA) scenarios, to enhance generation quality. This approach has been further advanced in research such as Reflexion and Agent-Pro (Shinn et al., 2024; Zhang et al., 2024b), which extend self-reflection to agentic scenarios, increasing the efficiency and success rate of task execution during scenario exploration and trajectory execution. However, there remains significant room for improvement in its performance, particularly when it comes to SLMs.

Through extensive review, we found lack of report on the effectiveness of self-reflection applied to models which possess fewer than 10 billion parameters. Its success relies heavily on the context generated during the self-reflection process (Cheng et al., 2024) and is prone to overconfidence in its generated content (Zhang et al., 2024a), including biases.

We assessed the self-reflective capabilities of several SLMs across a variety of tasks, with Self-Refine chosen as a baseline approach. Our findings reveal that reflective thinking of these models fails to produce meaningful improvements in their generative performance. Entrospect is specifically designed to enhance the performance of SLMs by leveraging information theory to assist in the self-reflection process.

## 2.2 Enhancing the Reasoning Capabilities of Small Language Models

Recent studies have made significant strides in enhancing the reasoning capabilities of SLMs. Bi et al. introduced Solution-Guidance Fine-Tuning (Bi et al., 2024), focusing on problem understanding and decomposition to improve SLMs’ generalization and reasoning abilities. Wang and Lu explored continual pre-training on a synthetic dataset to inject multi-step reasoning abilities into moderate-sized models (Wang and Lu, 2023). Fu et al. specialized small models towards multi-step reasoning through knowledge distillation from large models (Fu et al., 2023). Yu et al. developed TRIPOST, an algorithm enabling small models to self-improve via interaction with large ones (Yu et al., 2023).

However, these methods often necessitate a substantial amount of additional data, whether it is synthetically created or derived from larger models, which may not be readily accessible or easy to produce. They entail a certain degree of computational overhead, be it in data generation, pre-training, or

iterative training processes. Differently, Entrospect does not require any additional data or specialized training, thus drastically reducing both overhead and resource demands, allowing broader applicability across diverse domains and use cases.

## 3 Methodology

### 3.1 Problem Definition

While frameworks like Self-Refine aim to automate response refinement in language models through self-reflection, they do not inherently ensure that such refinements are beneficial. This limitation is particularly pronounced in SLMs, where constrained semantic capabilities lead to unreliable self-reflections, resulting in *reflective contamination*. Reflective contamination occurs when the model’s self-generated feedback contains biases, which can degrade rather than improve the refined response.

To formalize this problem, consider the  $t$ -th refinement round, where the model  $\mathcal{M}_\theta$  generates feedback  $F_t$  based on the query  $Q$ , reflection prompt  $P_{\text{reflect}}$ , and current response  $A_t$ . This feedback, represented as  $\mathcal{M}_\theta(Q \| A_t \| P_{\text{reflect}})$ , consists of two components: 1) A valid portion  $S_t = \rho_t F_t$ , which supports effective refinement. 2) Reflective contamination  $N_t = (1 - \rho_t) F_t$ , which introduces biases. Here,  $\rho_t \in (0, 1)$  represents the proportion of valid feedback in  $F_t$ . The refined response  $A_{t+1}$  is then generated using  $F_t$ ,  $Q$ , and the refinement prompt  $P_{\text{refine}}$ , expressed as:

$$\begin{aligned} A_{t+1} &= \mathcal{M}_\theta(Q \| A_t \| F_t \| P_{\text{refine}}) \\ &= A_t + \alpha_t^S S_t - \alpha_t^N N_t \\ &= A_t + \alpha_t^S \rho_t F_t - \alpha_t^N (1 - \rho_t) F_t \\ &= A_t + [(\alpha_t^S + \alpha_t^N) \rho_t - \alpha_t^N] F_t, \end{aligned} \quad (1)$$

where  $\alpha_t^S$  and  $\alpha_t^N$  are partial attention factors ( $\alpha \in (0, 1)$ ) applied to the valid and contaminated portions of  $F_t$ , respectively.

#### The Core Problem:

1. A successful refinement requires  $A_{t+1} \geq A_t$ , but this is not guaranteed. When  $\rho_t$  is low (i.e., the feedback contains more contamination), the refined response may degrade, as described by the condition:

$$\rho_t < \frac{\alpha_t^N}{\alpha_t^S + \alpha_t^N}. \quad (2)$$

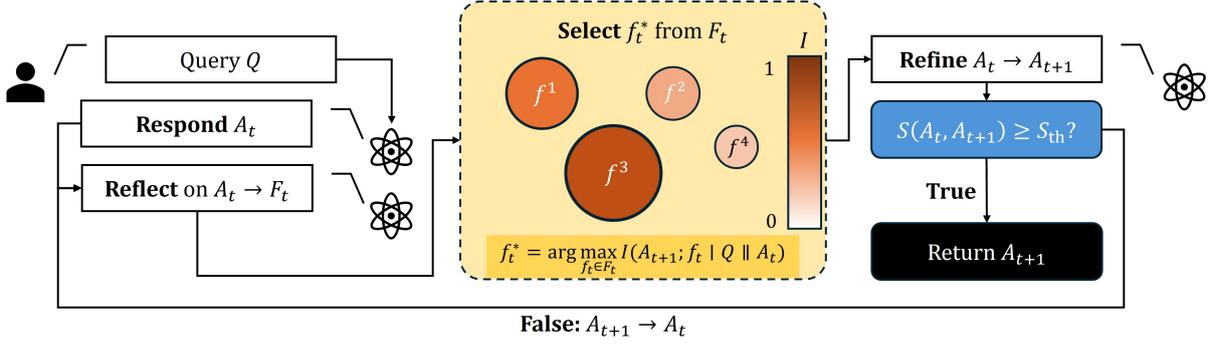


Figure 3: The pipeline of our Entrospect prompt-driven framework, extending the original Self-Refine structure with an **Optimal Revision Suggestion Selector (ORSS)** and a universal **semantic similarity-based stopping condition**. The framework requires no supervised pre-training or access to the model’s internal parameters, granting it to be generalizable to various language models and reasoning tasks.

2. SLMs, with their limited semantic competence, often exhibit low  $\rho_t$  and high  $\alpha_t^N$  (or low  $\alpha_t^S$ ), making them prone to degradation during the refinement phase of the response.

**Objective:** Within the realm of black-box models,  $\alpha_t^S$ ,  $\alpha_t^N$  and  $\rho_t$  are inaccessible. This presents a significant obstacle in accurately differentiating between  $S_t$  and  $N_t$ . An alternative perspective involves concentrating exclusively on the optimal component of  $F_t$ . Entrospect proposes an unsupervised mechanism driven by information theory, providing a systematic solution to this complication.

### 3.2 Optimal Revision Suggestion Selector

By employing a formatting prompt, we can steer the model’s self-reflective output towards a systematic arrangement of multiple revision suggestions. In this way,  $F_t$  is characterized as an ensemble of strings  $\{f_t^0, f_t^1, \dots, f_t^n\}$ , framing our goal as “discerning an optimal revision suggestion from this set”. However, in the absence of supervision, defining what constitutes *optimal* becomes a fundamental issue.

To address this, we propose a solution called the Optimal Revision Suggestion Selector (ORSS), which uses heuristic information-theoretic approaches for prompt selection (Wu et al., 2024; Yang et al., 2024b). These studies suggest that an optimal prompt should minimize the semantic uncertainty of a language model when processing a query, which is equivalent to maximizing the conditional mutual information (CMI) between the input and the output. Unlike recent work which assumes a manually constructed prompt pool,  $F_t$  as the candidate set in our case is constructed in an

automatic fashion, where revision suggestions become prompt candidates, and the one to be selected renders the maximum CMI following Equation 3:

$$f_t^* = \arg \max_{f_t \in F_t} I(A_{t+1}; f_t | Q || A_t),$$

$$\text{where } I = H(A_{t+1} | Q || A_t) - H(A_{t+1} | f_t, Q || A_t). \quad (3)$$

In Equation 3,  $Q || A_t$  stands for the prompt “Please provide a refined solution of <Q> given <A\_t>”, and  $(f_t, Q || A_t)$  signifies a slightly different prompt “Please provide a refined solution of <Q> given <A\_t>. <f\_t>”. The two  $H$ s characterize the *marginal entropy* and the *conditional entropy* in classical information theory, respectively. The value of CMI  $I$  stands for **the extent to which a revision suggestion  $f_t$  enhances the model’s confidence in the refinement applied to the current answer  $A_t$ .**

### 3.3 Eliciting the Convergence of Entrospect

We established a universal mechanism to enable Entrospect to automatically terminate its iterations. The core principle is that, at the semantic level,  $A_t$  and  $A_{t+1}$  are essentially equivalent. Consequently, when a language model employs greedy search (*temperature* = 0) for output sampling, subsequent outputs naturally converge toward consistency, rendering the increments from reflection and refinement negligible. Given these circumstances, the framework no longer introduces meaningful improvements to the response, a state we defined as “convergence”. More precisely, we leverage the *cosine similarity*  $S(\cdot, \cdot)$  to quantify the degree of semantic resemblance between two answers, modeled as

$$\begin{aligned}
S(A_1, A_2) &= \frac{\mathbf{v}_1 \cdot \mathbf{v}_2}{\|\mathbf{v}_1\| \|\mathbf{v}_2\|} \\
&= \frac{\sum_{i=1}^m (v_{1i} \cdot v_{2i})}{\sqrt{\sum_{i=1}^m v_{1i}^2} \cdot \sqrt{\sum_{i=1}^m v_{2i}^2}}, \quad (4)
\end{aligned}$$

where  $\mathbf{v} = [v_1 \ v_2 \ \dots \ v_m]^T$  indicates the  $A$ 's tokenized vector in a continuous,  $m$ -dimensional semantic space. The range of  $S$  is  $[-1, 1]$ , with a higher value referring to a stronger semantic similarity between the two entities compared. Leveraging semantic similarity as a stopping condition for the iterative refinement procedure guarantees an appropriate termination juncture, thus optimizing performance results.

### 3.4 Framework of Entrospect

Slightly different from the three-step process of *respond*  $\rightarrow$  *reflect*  $\rightarrow$  *refine* adopted by Self-Refine, Entrospect follows an extended four-step strategy: *respond*  $\rightarrow$  *reflect*  $\rightarrow$  *select*  $\rightarrow$  *refine*. In the following, we detail each step sequentially; see Figure 3 for an intuitive illustration of the pipeline and Algorithm 1 for implementation guidance.

**Respond:** The iterations begin with the language model generating an initial answer  $A_0$  for the input query  $Q$ .

**Reflect and Select:** During iteration  $t$ , the model  $\mathcal{M}_\theta$ , guided by the prompt  $P_{\text{reflect}}$ , the original query  $Q$ , and the current answer  $A_t$ , generates a set of candidate revision suggestions denoted as  $F_t = \{f_t^0, f_t^1, \dots, f_t^n\}$ . The prompt  $P_{\text{reflect}}$  serves as a directive that instructs the model on how to evaluate potential deficiencies in the current answer and construct appropriate  $F_t$  accordingly. Thereafter, the ORSS selects the optimal  $f_t^*$  that maximizes the CMI between the input and the output of the model. In practical implementation, the *Cross-Entropy Loss*  $\mathcal{L}_{\text{CE}}$  output by the model for a given input can be used to calculate the marginal entropy and the conditional entropy, allowing for the straightforward computation of the CMI.

**Refine:** Leveraging the  $f_t^*$  as the key instruction to the refinement, the model  $\mathcal{M}_\theta$  utilizes the prompt  $P_{\text{refine}}$ , in conjunction with the original query  $Q$  and the current answer  $A_t$ , to generate an updated answer  $A_{t+1}$ .

**Stop Condition:** Subsequent to the generation of the  $A_{t+1}$ , we exert the semantic textual similarity measure to check whether the iterative process

---

### Algorithm 1 The algorithm pipeline of Entrospect

---

**Require:** query  $Q$ , model  $\mathcal{M}_\theta$ , prompt  $P_{\text{reflect}}$  ( $:= P_f$ ), prompt  $P_{\text{refine}}$  ( $:= P_r$ ), semantic similarity threshold  $S_{\text{th}}$

- 1:  $A_0 \leftarrow \mathcal{M}_\theta(Q)$   $\triangleright$  Respond
- 2:  $A_t \leftarrow A_0$
- 3: **while** True **do**
- 4:    $F_t \leftarrow \mathcal{M}_\theta(P_f \| Q \| A_t)$   $\triangleright$  Reflect
- 5:    $\{f_t^0, f_t^1, \dots, f_t^n\} \leftarrow \text{list}(F_t)$   $\triangleright$  Itemize
- 6:    $I_{\text{max}} \leftarrow 0$
- 7:   **for**  $f_t$  in list( $F_t$ ) **do**  $\triangleright$  Select (ORSS)
- 8:      $H_t^{\text{marg}} \leftarrow \mathcal{L}_{\text{CE}}(\mathcal{M}_\theta(P_r \| Q \| A_t))$
- 9:      $H_t^{\text{cond}} \leftarrow \mathcal{L}_{\text{CE}}(\mathcal{M}_\theta(P_r \| Q \| A_t \| f_t))$
- 10:      $I_t \leftarrow H_t^{\text{marg}} - H_t^{\text{cond}}$
- 11:     **if**  $I_t > I_{\text{max}}$  **then**
- 12:        $f_t^* \leftarrow f_t$
- 13:     **end if**
- 14:   **end for**
- 15:    $A_{t+1} \leftarrow \mathcal{M}_\theta(P_r \| Q \| A_t \| f_t^*)$   $\triangleright$  Refine
- 16:   **if**  $S(A_t, A_{t+1}) \geq S_{\text{th}}$  **then**
- 17:     **break**
- 18:   **end if**
- 19:    $A_t \leftarrow A_{t+1}$
- 20: **end while**
- 21: **return**  $A_{t+1}$

---

should be terminated. When  $A_t$  and  $A_{t+1}$  exhibit a high degree of semantic resemblance, this suggests that Entrospect has entered a state of convergence from the current iteration onward. Following that,  $A_{t+1}$  is designated as the final output. To meet the requirements of long-text encoding with high representational fidelity, we opted for the *Jina Embeddings V3* (Sturua et al., 2024) with a dedicated LoRA adapter for text-matching tasks, an encoder-based model which natively supports an input sequence length of up to 8192 tokens. In our experiments,  $S \geq 0.9$  is adopted as the threshold for considering  $A_t$  and  $A_{t+1}$  semantically equivalent.

We detailed the instructions involved in the operation process of Entrospect in Figure 6.

## 4 Experiments and Results

### 4.1 Experimental Settings

We evaluated Entrospect equipped by four of the latest SLMs, including DeepSeek-R1-distilled Qwen 2.5 1.5B (Yang et al., 2024a; Guo et al., 2025), Qwen 2.5 7B (Yang et al., 2024a), Llama 3.1 8B (AI, 2024), and GLM-4 9B (GLM et al., 2024),

as compared to the baselines (see Section 4.4) on a math reasoning dataset and a hallucination detection dataset, namely MATH (Hendrycks et al., 2021) and HaluEval (Li et al., 2023). Each SLM was quantized to INT4 precision with either AutoGPTQ or BitsAndBytes (Pan, 2023; Dettmers et al., 2022).

## 4.2 Datasets

To comprehensively assess whether Entrospect heightens the ubiquitous reasoning performance of SLMs, we sourced our validation data from two representative datasets, MATH and HaluEval, with illustrative examples provided in Table 2.

Table 1: Accuracies (%) of various methods equipped by four of the latest SLMs on reasoning tasks MATH (The average accuracies of level 1 to level 5) and HaluEval. We highlight the best results in **bold**.

Model Name	Method	MATH	HaluEval
DeepSeek-R1-Distilled Qwen 2.5 Instruct 1.5B	Zero-Shot	94.2	80.5
	5-Shot	90.2	29.5
	Zero-Shot CoT	91.3	91.0
	Self-Refine	88.5	80.0
	<b>Entrospect</b>	<b>98.4</b>	<b>95.5</b>
Qwen 2.5 Instruct 7B	Zero-Shot	78.2	94.5
	5-Shot	72.8	91.0
	Zero-Shot CoT	83.8	98.0
	Self-Refine	73.0	97.5
	<b>Entrospect</b>	<b>86.0</b>	<b>100.0</b>
Llama 3.1 Instruct 8B	Zero-Shot	61.7	94.5
	5-Shot	56.5	94.0
	Zero-Shot CoT	73.7	94.5
	Self-Refine	44.3	95.0
	<b>Entrospect</b>	<b>80.5</b>	<b>99.5</b>
GLM 4 Instruct 9B	Zero-Shot	55.0	98.5
	5-Shot	57.9	97.5
	Zero-Shot CoT	65.8	97.5
	Self-Refine	56.8	97.5
	<b>Entrospect</b>	<b>69.7</b>	<b>100.0</b>

**MATH** (Hendrycks et al., 2021): a dataset designed to measure the mathematical reasoning capabilities of language models, consisting of problems sourced from high school math competitions, tagged with difficulty levels from 1 to 5 and covering a wide range of topics including algebra, geometry, number theory, and combinatorics. MATH is notable for its complexity compared to the other datasets of the same category (Frieder et al., 2024), e.g. GSM8K (Cobbe et al., 2021). Besides, the latest findings have unveiled that MATH suffers less leakage than GSM8K does from the worsen-

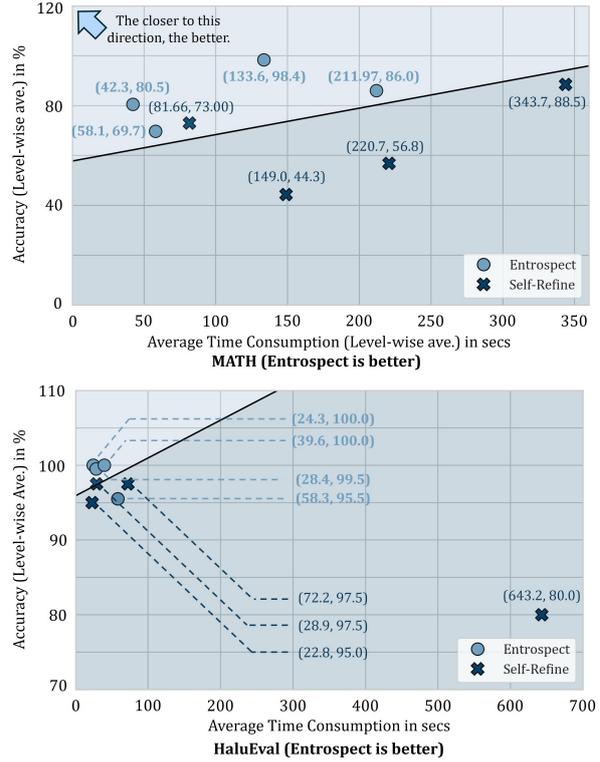


Figure 4: The Accuracy-ATC results derived from evaluating Entrospect and Self-Refine across four models and two tasks. The dividing lines in the chart correspond to the decision boundaries determined by linear SVMs fit on the data points of Entrospect and Self-Refine. Data points positioned closer to the **top-left corner** signify a more favorable trade-off between computational efficiency and reasoning accuracy, indicating superior overall performance.

ing cheating on model training (Xu et al., 2024), underlining its fairness. We randomly chose 120 samples from each difficulty level to serve as our experimental dataset.

**HaluEval** (Li et al., 2023): a dataset that gauges the performance of language models in recognizing hallucinations, featuring general user queries and task-specific examples across question answering, dialogue, and text summarization. We randomly sampled 200 pairs from this dataset, providing a robust evaluation platform for analyzing the effectiveness of our framework in detecting and reducing hallucinations.

## 4.3 Evaluation Metrics

We selected two evaluation metrics, i.e. Accuracy and Average Time Consumption (Han et al., 2023; Xu et al., 2023; Xiao et al., 2024), to provide both qualitative and quantitative insights into the effectiveness of Entrospect.

**Accuracy:** a pivotal evaluation metric, is de-

lineated as the proportion of problems correctly resolved relative to the total number of problems the model attempts, computed via  $A_{\text{correct}} / (A_{\text{correct}} + A_{\text{wrong}}) \times 100\%$ . A higher accuracy signifies that a prompting scheme is more effective in lifting the model’s reasoning outcomes.

**Average Time Consumption:** We measured the Average Time Consumption (ATC) of the selected prompting schemes, spanning from the moment the input is supplied to the generation of the final output. Given the sample size  $N$  of the validation set, ATC is calculated by  $\frac{1}{N} \sum_k (t_{k_o} - t_{k_i})$ , where  $t_{k_o} - t_{k_i}$  denotes the duration, counted in seconds, from the moment the  $k$ -th input is supplied to the time the  $k$ -th output is generated. A smaller ATC embodies better computational efficiency of a prompting method, which is vital for industrial implementation, notably on edge computing devices running local SLMs. In our assessments, both of the above metrics are considered for more comprehensive analysis.

#### 4.4 Baseline Selection

We compared Entrospect against the following well-established prompting methods as well as its ablated version, functioning as robust benchmarks for appraising the performance uplift in SLMs achieved with Entrospect.

**Zero-Shot and Few-Shot Prompting (Brown et al., 2020):** Zero-shot prompting directs a language model to perform tasks with only high-level instructions, often sacrificing accuracy for complex inputs. Conversely, few-shot prompting supplies demonstrations to improve context awareness and performance, yet its success hinges on the quality of examples, which may not fully capture task complexity and may be labor-intensive to gather in practice.

**Chain-of-Thought Prompting (Wei et al., 2022):** An approach that guides language models to generate a structured reasoning path before arriving at the final answer, encouraging more systematic and transparent problem solving. A key downside is the increased potential for longer outputs, as irrelevant, inaccurate, and repetitive steps may appear in the generated thought chain, especially concerning SLMs, impairing the overall outcome.

**Self-Refine (Madaan et al., 2024):** The framework allows a model to iteratively revise its own outputs with identified errors from the self-reflection’s feedback. Despite its potential, such a strategy

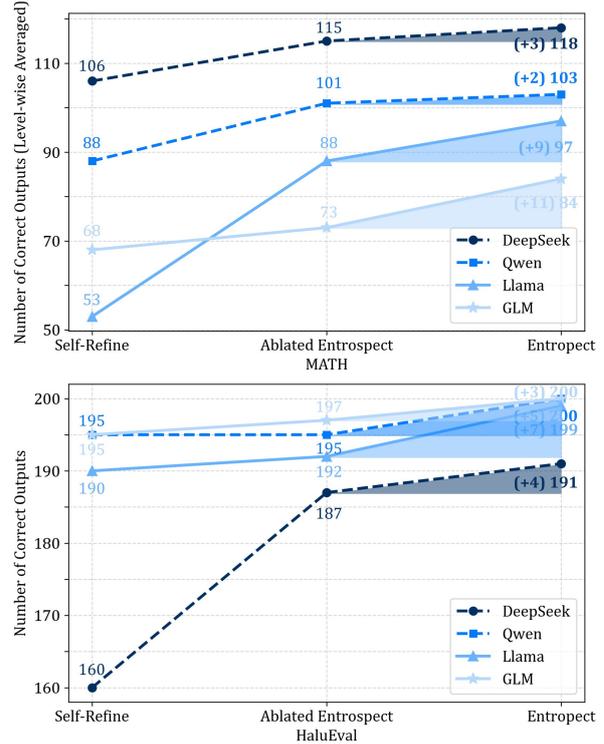


Figure 5: (A higher Number of Correct Outputs is better) Constrained on a fixed 5 rounds of refinement rather than the stopping condition, the ablated Entrospect falls into suboptimal performance in contrast to the complete version across both tasks and all involved models. This highlights the significance and efficacy of importing semantic similarity comparison as the stopping condition for our framework.

may introduce unnecessary or incorrect changes during the refinement cycles, especially for SLMs, as mentioned in Section 1.

**Ablated Entrospect:** The variant of Entrospect without the semantic similarity-based stopping condition. Instead, a manual setting of 5 fixed iterations is assigned. This baseline serves as the *ablation study* that verifies the efficacy of our nominated convergence policy.

#### 4.5 Results

We report the Entrospect’s state-of-the-art competences versus the baseline prompting approaches, especially Self-Refine, in augmenting the SLMs’ semantic reasoning across two validation tasks.

**Entrospect improves reasoning accuracies:** Displayed in Table 1 and 3, SLMs armed with Entrospect outshines all other baselines pertaining to the reasoning accuracies across both MATH and HaluEval validation sets. In contrast specifically to Self-Refine, Entrospect yields a maximum improvement of 36.2% with *Llama 3.1 Instruct*

482	8B (44.3% → 80.5%) on the MATH dataset and	information-theoretic Optimal Revision Sugges-	532
483	15.5% with <i>DeepSeek-R1-Distilled Qwen 2.5 In-</i>	tion Selector to provide optimal revision sugges-	533
484	<i>struct 1.5B</i> (80.0% → 95.5%) on the HaluEval	tions during the self-reflection stage while elimi-	534
485	dataset. Moreover, Figure 7 highlights Entrospect’s	nating ineffective ones for efficient refinement of	535
486	robustness beyond handling math problems with	initial responses from SLMs. Besides, the con-	536
487	a fixed complexity. When set against Self-Refine,	vergence of Entrospect is made possible with a	537
488	Entrospect consistently offers more substantial miti-	dedicated semantic similarity-determined stopping	538
489	gation against the overall degradation of reasoning	condition. Through our holistic evaluations, Entro-	539
490	accuracy as the problem difficulty rises, securing a	spect claimed state-of-the-art relative to the base-	540
491	reduced decay rate as much as 52.8%.	line methods on both of our reasoning tasks across	541
492	<b>The exceptional computational efficiency:</b> As	four SLMs of diverse parameter sizes, obtaining a	542
493	depicted in Figure 4, Entrospect reaches conver-	maximum 36.2% reasoning accuracy uplift and at	543
494	gence faster than Self-Refine across most instances.	most 10 times the computational efficiency exclu-	544
495	on the MATH dataset, Entrospect reduces runtime	sively over its antecedent, Self-Refine.	545
496	by an average factor of up to 2.8 (e.g., <i>Llama 3.1</i>	We aspire for this study to inspire further ad-	546
497	<i>8B + Entrospect</i> ), meanwhile demonstrating even	vancements in small language models research	547
498	more pronounced efficiency gains on the HaluEval	and furnishes new perspectives for information-	548
499	dataset, with runtime reductions reaching up to 10-	theoretic prompt engineering.	549
500	fold (e.g., <i>DeepSeek R1-Distill Qwen 2.5 1.5B +</i>		
501	<i>Entrospect</i> ). Beyond its efficiency advantages, Fig-	<b>Limitations</b>	550
502	ure 4 highlights Entrospect’s ability to strike a supe-	There remains much room for promoting Entro-	551
503	rior balance between computational efficiency and	spect, and our future studies shall prioritize the	552
504	accuracy, driving substantial overall performance	following key limitations:	553
505	enhancements in SLMs.	<b>More solid definition of an <i>optimal</i> revision sug-</b>	554
506	To investigate potential correlations between	<b>gestion:</b> The ORSS of Entrospect, grounded in	555
507	model parameter sizes and the ATC outcomes	maximizing the conditional mutual information,	556
508	achieved by Entrospect, we employed Spearman’s	operates as an approximate selection technique in	557
509	rank correlation coefficient alongside correspond-	unsupervised settings. This approach gauges the	558
510	ing <i>p</i> -values (Spearman, 2010). However, no statis-	quality of a revision suggestion by leveraging the	559
511	tically significant relationship was observed within	model’s intrinsic output uncertainty as a pivotal de-	560
512	the scope of our experiments (MATH: <i>corr</i> =	terminant. However, its reliability is compromised	561
513	−0.600, <i>p</i> = 0.400; HaluEval: <i>corr</i> = −0.200,	when the model demonstrates undue confidence in	562
514	<i>p</i> = 0.800).	erroneous outputs. As a result, it is imperative to	563
515	<b>Ablation study:</b> To validate whether the seman-	pursue a more precise and theoretically grounded	564
516	tic similarity-based stopping condition is crucial	definition of what constitutes an <i>optimal</i> revision	565
517	for propelling a higher reasoning accuracy of En-	suggestion in our future studies.	566
518	trospect, we conducted an ablation study by re-	<b>Beyond semantic similarity comparison as the</b>	567
519	moving this mechanism and fixing the number of	<b>stopping condition:</b> A high semantic similarity	568
520	refinement cycles to 5. Figure 5 illustrates that the	between consecutive refinement iterations as a sign	569
521	ablated Entrospect constantly underperforms com-	of convergence is logically aligned with language	570
522	pared to the complete implementation, witness-	models adopting greedy search sampling. In con-	571
523	ing performance deficits of 1.8 → 8.9% on the MATH	versational situations, however, sampling methods	572
524	dataset and 1.5% → 3.5% on the HaluEval dataset	such as Top-K and nucleus sampling are more reg-	573
525	across all tested SLMs. The results solidify the role	ularly used to ensure generative variability. Our	574
526	of the semantic similarity-guided stopping condi-	future work will seek to modify the current con-	575
527	tion as a cornerstone for enhancing Entrospect’s	vergence mechanism tailored to these sampling	576
528	overall performance.	configurations.	577
529	<b>5 Conclusion</b>	<b>Ethics Statement</b>	578
530	This paper introduces Entrospect, an opti-	This study strictly adheres to the Ethical Policy of	579
531	mized Self-Refine framework that leverages an	the Association for Computational Linguistics. We	580

581	conducted a thorough assessment of the potential	Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chen-	633
582	impacts of our research and did not identify any	hui Zhang, Da Yin, Dan Zhang, Diego Rojas, Guanyu	634
583	evident ethical concerns. All datasets utilized in	Feng, Hanlin Zhao, et al. 2024. Chatglm: A family	635
584	this study were sourced from publicly available	of large language models from glm-130b to glm-4 all	636
585	resources and were handled strictly in accordance	tools. <i>arXiv preprint arXiv:2406.12793</i> .	637
586	with their respective terms of use. Nevertheless, we	Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song,	638
587	acknowledge the possibility of unforeseen impacts	Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma,	639
588	in any research and invite readers to share feedback	Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: In-	640
589	on any potential ethical concerns they may identify.	centivizing reasoning capability in llms via reinforce-	641
		ment learning. <i>arXiv preprint arXiv:2501.12948</i> .	642
		Chengcheng Han, Xiaowei Du, Che Zhang, Yixin Lian,	643
		Xiang Li, Ming Gao, and Baoyuan Wang. 2023. Di-	644
590	<b>References</b>	alcot meets ppo: Decomposing and exploring reason-	645
591	Meta AI. 2024. Llama 3.1: The most capable open	ing paths in smaller language models. <i>arXiv preprint</i>	646
592	foundation models. <a href="https://ai.meta.com/blog/meta-llama-3-1/">https://ai.meta.com/blog/</a>	<i>arXiv:2310.05074</i> .	647
593	<a href="https://ai.meta.com/blog/meta-llama-3-1/">meta-llama-3-1/</a> .		
594	Jing Bi, Yuting Wu, Weiwei Xing, and Zhenjie Wei.	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul	648
595	2024. Enhancing the reasoning capabilities of small	Arora, Steven Basart, Eric Tang, Dawn Song, and	649
596	language models via solution guidance fine-tuning.	Jacob Steinhardt. 2021. Measuring mathematical	650
597	<i>arXiv preprint arXiv:2412.09906</i> .	problem solving with the math dataset. In <i>Thirty-</i>	651
		<i>fifth Conference on Neural Information Processing</i>	652
		<i>Systems Datasets and Benchmarks Track (Round 2)</i> .	653
598	Tom Brown, Benjamin Mann, Nick Ryder, Melanie	Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong,	654
599	Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind	Zhangyin Feng, Haotian Wang, Qianglong Chen,	655
600	Neelakantan, Pranav Shyam, Girish Sastry, Amanda	Weihua Peng, Xiaocheng Feng, Bing Qin, et al. 2023.	656
601	Askell, et al. 2020. Language models are few-shot	A survey on hallucination in large language models:	657
602	learners. <i>Advances in neural information processing</i>	Principles, taxonomy, challenges, and open questions.	658
603	<i>systems</i> , 33:1877–1901.	<i>ACM Transactions on Information Systems</i> .	659
604	Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu,	Dongyub Lee, Eunhwan Park, Hodong Lee, and Heui-	660
605	Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi,	Seok Lim. 2024. Ask, assess, and refine: Rectifying	661
606	Cunxiang Wang, Yidong Wang, et al. 2024. A sur-	factual consistency and hallucination in llms with	662
607	vey on evaluation of large language models. <i>ACM</i>	metric-guided feedback learning. In <i>Proceedings of</i>	663
608	<i>Transactions on Intelligent Systems and Technology</i> ,	<i>the 18th Conference of the European Chapter of the</i>	664
609	15(3):1–45.	<i>Association for Computational Linguistics (Volume</i>	665
610	Ruoxi Cheng, Haoxuan Ma, and Shuirong Cao. 2024.	<i>1: Long Papers)</i> , pages 2422–2433.	666
611	Deceiving to enlighten: Coaxing llms to self-	Beibin Li, Yi Zhang, Sébastien Bubeck, Jeevan Pathuri,	667
612	reflection for enhanced bias detection and mitigation.	and Ishai Menache. 2024. Small language models for	668
613	<i>arXiv preprint arXiv:2404.10160</i> .	application interactions: A case study. <i>arXiv preprint</i>	669
614	Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,	<i>arXiv:2405.20347</i> .	670
615	Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias	Junyi Li, Xiaoxue Cheng, Wayne Xin Zhao, Jian-Yun	671
616	Plappert, Jerry Tworek, Jacob Hilton, Reiichiro	Nie, and Ji-Rong Wen. 2023. Halueval: A large-	672
617	Nakano, et al. 2021. Training verifiers to solve math	scale hallucination evaluation benchmark for large	673
618	word problems. <i>arXiv preprint arXiv:2110.14168</i> .	language models.	674
619	Tim Dettmers, Mike Lewis, Younes Belkada, and Luke	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler	675
620	Zettlemoyer. 2022. Gpt3. int8 (): 8-bit matrix mul-	Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon,	676
621	tiplication for transformers at scale. <i>Advances in</i>	Nouha Dziri, Shrimai Prabhumoye, Yiming Yang,	677
622	<i>Neural Information Processing Systems</i> , 35:30318–	et al. 2024. Self-refine: Iterative refinement with	678
623	30332.	self-feedback. <i>Advances in Neural Information Pro-</i>	679
624	Simon Frieder, Mirek Olsák, Julius Berner, and Thomas	<i>cessing Systems</i> , 36.	680
625	Lukasiewicz. 2024. The imo small challenge: Not-	Qiwei Pan. 2023. <a href="#">Autogptq: An easy-to-use llms quan-</a>	681
626	too-hard olympiad math datasets for llms. In <i>Tiny</i>	<a href="#">tization package with user-friendly apis, based on</a>	682
627	<i>Papers@ ICLR</i> .	<a href="#">gptq algorithm</a> .	683
628	Yao Fu, Hao Peng, Litu Ou, Ashish Sabharwal, and	Noah Shinn, Federico Cassano, Ashwin Gopinath,	684
629	Tushar Khot. 2023. Specializing smaller language	Karthik Narasimhan, and Shunyu Yao. 2024. Re-	685
630	models towards multi-step reasoning. In <i>Inter-</i>	flexion: Language agents with verbal reinforcement	686
631	<i>national Conference on Machine Learning</i> , pages	learning. <i>Advances in Neural Information Process-</i>	687
632	10421–10430. PMLR.	<i>ing Systems</i> , 36.	688

689	C Spearman. 2010. The proof and measurement of association between two things. <i>International Journal of Epidemiology</i> , 39(5):1137–1150.	Tianqiang Yan and Tiansheng Xu. 2023. Refining the responses of llms by themselves. <i>arXiv preprint arXiv:2305.04039</i> .	746 747 748
692	Saba Sturua, Isabelle Mohr, Mohammad Kalim Akram, Michael Günther, Bo Wang, Markus Krimmel, Feng Wang, Georgios Mastrapas, Andreas Koukounas, Nan Wang, et al. 2024. jina-embeddings-v3: Multilingual embeddings with task lora. <i>arXiv preprint arXiv:2409.10173</i> .	An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. 2024a. Qwen2. 5 technical report. <i>arXiv preprint arXiv:2412.15115</i> .	749 750 751 752
698	Fali Wang, Zhiwei Zhang, Xianren Zhang, Zongyu Wu, Tzuhao Mo, Qiuhaio Lu, Wanjing Wang, Rui Li, Junjie Xu, Xianfeng Tang, et al. 2024. A comprehensive survey of small language models in the era of large language models: Techniques, enhancements, applications, collaboration with llms, and trustworthiness. <i>arXiv preprint arXiv:2411.03350</i> .	Sohee Yang, Jonghyeon Kim, Joel Jang, Seonghyeon Ye, Hyunji Lee, and Minjoon Seo. 2024b. Improving probability-based prompt selection through unified evaluation and analysis. <i>Transactions of the Association for Computational Linguistics</i> , 12:758–774.	753 754 755 756 757
705	Tianduo Wang and Wei Lu. 2023. Learning multi-step reasoning by solving arithmetic tasks. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 1229–1238.	Xiao Yu, Baolin Peng, Michel Galley, Jianfeng Gao, and Zhou Yu. 2023. Teaching language models to self-improve through interactive demonstrations. <i>arXiv preprint arXiv:2310.13522</i> .	758 759 760 761
710	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837.	Wenqi Zhang, Yongliang Shen, Linjuan Wu, Qiuying Peng, Jun Wang, Yueting Zhuang, and Weiming Lu. 2024a. Self-contrast: Better reflection through inconsistent solving perspectives. <i>arXiv preprint arXiv:2401.02009</i> .	762 763 764 765 766
715	Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, et al. 2021. Ethical and social risks of harm from language models. <i>arXiv preprint arXiv:2112.04359</i> .	Wenqi Zhang, Ke Tang, Hai Wu, Mengna Wang, Yongliang Shen, Guiyang Hou, Zeqi Tan, Peng Li, Yueting Zhuang, and Weiming Lu. 2024b. Agent-pro: Learning to evolve via policy-level reflection and optimization. <i>arXiv preprint arXiv:2402.17574</i> .	767 768 769 770 771
720	Laura Weidinger, Jonathan Uesato, Maribeth Rauh, Conor Griffin, Po-Sen Huang, John Mellor, Amelia Glaese, Myra Cheng, Borja Balle, Atoosa Kasirzadeh, et al. 2022. Taxonomy of risks posed by language models. In <i>Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency</i> , pages 214–229.	Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, et al. 2023. Siren’s song in the ai ocean: a survey on hallucination in large language models. <i>arXiv preprint arXiv:2309.01219</i> .	772 773 774 775 776
727	Junda Wu, Tong Yu, Rui Wang, Zhao Song, Ruiyi Zhang, Handong Zhao, Chaochao Lu, Shuai Li, and Ricardo Henao. 2024. Infoprompt: Information-theoretic soft prompt tuning for natural language understanding. <i>Advances in Neural Information Processing Systems</i> , 36.	<b>A Appendix</b>	777
733	Bin Xiao, Burak Kantarci, Jiawen Kang, Dusit Niyato, and Mohsen Guizani. 2024. Efficient prompting for llm-based generative internet of things. <i>arXiv preprint arXiv:2406.10382</i> .		
737	Ruijie Xu, Zengzhi Wang, Run-Ze Fan, and Pengfei Liu. 2024. Benchmarking benchmark leakage in large language models. <i>arXiv preprint arXiv:2404.18824</i> .		
740	Zhaozhuo Xu, Zirui Liu, Beidi Chen, Yuxin Tang, Jue Wang, Kaixiong Zhou, Xia Hu, and Anshumali Shrivastava. 2023. Compress, then prompt: Improving accuracy-efficiency trade-off of llm inference with transferable prompt. <i>arXiv preprint arXiv:2305.11186</i> .		

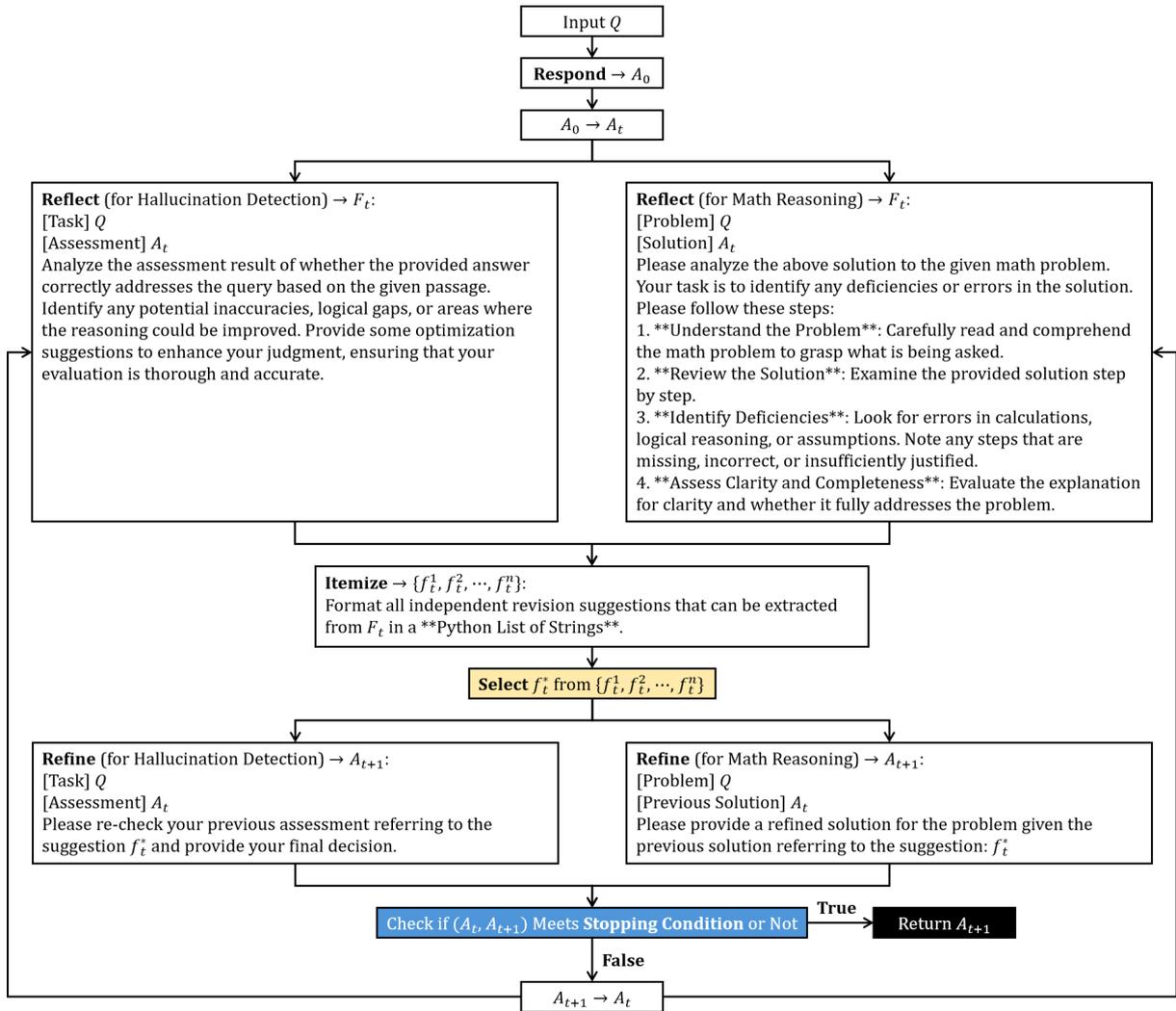


Figure 6: (Referred in Section 3.4) The detailed instructions used for all prompting nodes (modules) within the Entrospect framework during the evaluation phases. These instructions guide the SLMs through the process of generating an initial response, reflecting on its deficiencies, selecting the optimal revision, and refining the response based on the selected suggestion.

Table 2: (Referred in Section 4.2) Representative data samples from the MATH and HaluEval datasets, demonstrating a mathematical reasoning problem and a reading comprehension task.

<b>Dataset</b>	<b>Query</b>	<b>Label</b>
MATH	What is the simplified numerical value of $\frac{a+11b}{a-b}$ if $\frac{4a+3b}{a-2b} = 5$ ?	Let’s play with the given condition a little. Clearing out the denominator gives $4a + 3b = 5(a - 2b) = 5a - 10b$ . Selectively combine like terms by adding $9b - 4a$ to both sides to get $12b = a - b$ . This gives $\frac{12b}{a-b} = 1$ . Now, we want to find $\frac{a + 11b}{a - b}$ . Rewrite this as $\frac{a - b + 12b}{a - b} = \frac{a - b}{a - b} + \frac{12b}{a - b} = 1 + 1 = \boxed{2}$ , and we are done.
HaluEval	The following is a reading comprehension task, which provides a passage, a question related to the passage, and an answer to the question: [Passage] The ValleyCats play at Joseph L. Bruno Stadium which opened in 2002 on the campus of Hudson Valley Community College located in Troy. Joseph Bruno Stadium is a stadium located on the campus of Hudson Valley Community College in Troy, New York. [Question] The Tri-City ValleyCats play at which stadium located on the campus of Hudson Valley Community College in Troy, New York? [Answer] Troy Community Stadium, located on Hudson Valley Community College campus. Please determine whether the given answer is correct. If it is correct, output ‘PASS’; if it is incorrect, output ‘FAIL’.	FAIL

Table 3: (Referred in Section 4.5) The extended table of accuracies(%) on the MATH dataset, providing a detailed breakdown of all results across Level 1 to Level 5, where Entrospect performs the best with all SLMs relative to the baseline prompting methods across all difficulty levels. We highlight the best results in **bold**.

Model Name	Method	MATH-L1	MATH-L2	MATH-L3	MATH-L4	MATH-L5
DeepSeek-R1-Distilled Qwen 2.5 Instruct 1.5B	Zero-Shot	97.5	95.0	96.7	91.7	90.0
	5-Shot	96.7	92.5	94.2	91.7	75.8
	Zero-Shot CoT	75.0	98.3	98.3	93.3	91.7
	Self-Refine	90.0	89.2	90.8	88.3	84.2
	<b>Entrospect</b>	<b>99.2</b>	<b>99.2</b>	<b>99.2</b>	<b>96.7</b>	<b>97.5</b>
Qwen 2.5 Instruct 7B	Zero-Shot	91.7	92.5	85.8	73.3	47.5
	5-Shot	90.8	92.5	81.7	62.5	36.7
	Zero-Shot CoT	95.0	93.3	91.7	79.2	60.0
	Self-Refine	85.0	89.2	82.5	65.8	42.5
	<b>Entrospect</b>	<b>95.0</b>	<b>95.8</b>	<b>91.7</b>	<b>84.2</b>	<b>63.3</b>
Llama 3.1 Instruct 8B	Zero-Shot	87.5	74.2	60.8	48.3	37.5
	5-Shot	88.3	70.0	62.5	41.7	20.0
	Zero-Shot CoT	91.7	83.3	77.5	65.8	50.0
	Self-Refine	72.5	54.2	43.3	28.3	23.3
	<b>Entrospect</b>	<b>95.0</b>	<b>88.3</b>	<b>84.2</b>	<b>74.2</b>	<b>60.8</b>
GLM 4 Instruct 9B	Zero-Shot	82.5	66.7	55.0	46.7	24.2
	5-Shot	86.7	66.7	66.7	44.2	25.0
	Zero-Shot CoT	90.0	81.7	75.8	51.7	30.0
	Self-Refine	85.0	69.2	63.3	45.0	21.7
	<b>Entrospect</b>	<b>92.5</b>	<b>85.8</b>	<b>79.2</b>	<b>56.7</b>	<b>34.2</b>

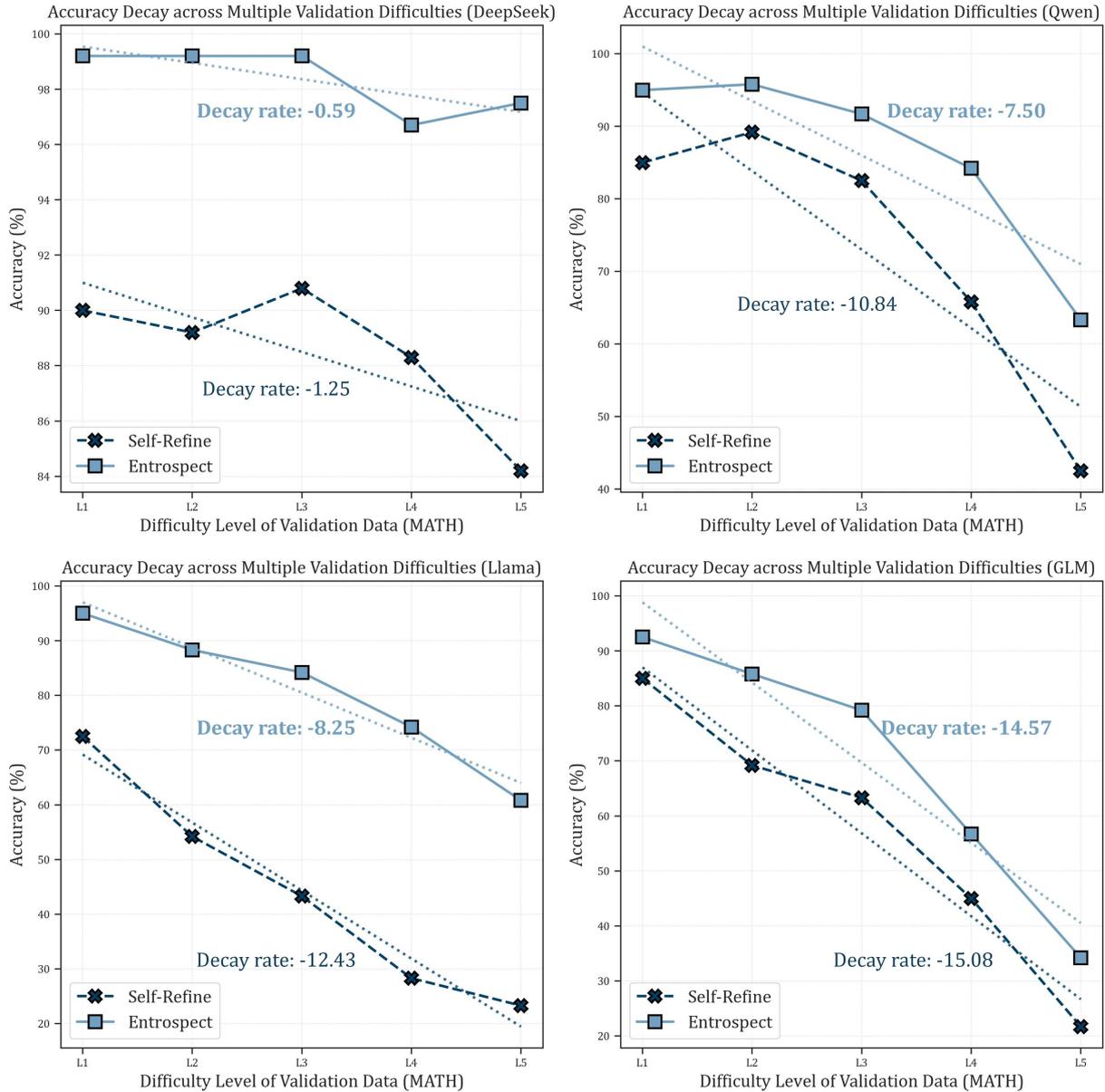


Figure 7: (Referred in Section 4.5) We employed linear regression to model the decline in reasoning accuracy, as measured by Entrospect and Self-Refine on the MATH validation set with increasing difficulty levels. The four charts correspond to the four distinct SLMs we evaluated, where the *decay rate* equals the slope of each fitted decay line. A decay rate with a larger absolute value indicates a more rapid deterioration in reasoning accuracy as the difficulty level rises. Across all tested models, observations indicate that as the difficulty level of the test data increases, the performance degradation exhibited by Entrospect is, overall, less pronounced than that of Self-Refine.