HEART: EMOTIONALLY-DRIVEN TEST-TIME SCALING OF LANGUAGE MODELS

Anonymous authors

000

001

002 003 004

005

006 007 008

010 011

012

013

014

015

016

018

019

021

023

024

027

029

031

033

035

036

037

038

040

041

042

043

Paper under double-blind review

ABSTRACT

Test-time scaling has shown considerable success in improving the performance of language models on complex reasoning tasks without requiring fine-tuning. However, current strategies such as self-reflection primarily focus on logical or structural refinement. They do not leverage the guiding potential of affective feedback. Inspired by psychological research showing that emotions can modulate cognitive performance, we introduce HEART-a novel framework that uses emotionally-driven prompts for iterative selfcorrection. HEART provides feedback on a model's incorrect response using a curated set of concise, emotionally charged phrases based on the six universal emotion categorized by Dr. Paul Ekman. By systematically varying the emotional tone of the feedback across iterations, our method guides the model to escape flawed reasoning paths and explore more promising alternatives. We evaluate our framework on challenging reasoning benchmarks including OlympiadBench, Humanity's Last Exam, and SimpleQA. Our results reveal a significant new phenomenon: when guided by an oracle verifier, this affective iteration protocol unlocks significantly deeper reasoning, leading to consistent and substantial increases in accuracy over state-of-the-art baselines with the same verifier. However, we also identify a critical bottleneck for practical deployment. In a verifier-free setting, it struggles to harness these gains consistently, highlighting as a key challenge for future work. Our findings suggest that the next frontier in machine reasoning may lie not just in refining logic, but also in understanding and leveraging the 'HEART' of the models.

1 Introduction

Large language models have demonstrated remarkable capabilities, yet eliciting reliable, complex reasoning remains a fundamental challenge. As models have scaled, research has moved beyond simple instruction-following to explore more methods of guidance. Structured reasoning techniques, such as Chain-of-Thought (CoT) (Wei et al., 2022) and its variants (Wang et al., 2022; Yao et al., 2023), impose a logical scaffold on the model's output, enhancing procedural correctness by externalizing the reasoning process. In parallel, initial explorations leveraging affective prompting, such as EmotionPrompt (Li et al., 2023), have shown that emotional cues can boost performance by igniting the model's "cognitive state" and guiding its focus.

Despite their successes, these two approaches suffer from a critical, complementary limitation. Structured methods are procedurally robust but affectively sterile; they provide a logical path but fail to leverage the motivational contexts that drive high-quality human reasoning. This sterility can lead to brittle performance, where models correctly execute a known algorithm but fail on novel problems requiring creative error recovery. Conversely, existing affective prompts are motivationally potent but structurally imprecise. They typically act as a "one-shot" global stimulus, which lacks the targeted guidance necessary to steer a model through a multi-step self-correction process. Consequently, a significant gap exists in the literature: there is

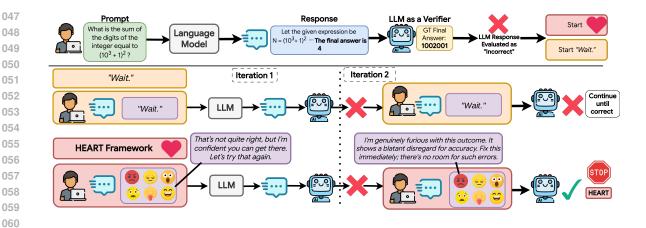


Figure 1: An illustration of the *HEART* framework. The process begins when a task is sent to a large language model (LLM), which returns a response. An oracle then evaluates the response against the ground truth. If the response is incorrect, the *HEART* process begins, incorporating the original task, the LLM's response, and selected affective cue prompts to generate a new, improved response.

no established method that unifies the systematic control of structured reasoning with the targeted application of affective cues for iterative self-improvement.

We address this gap by drawing on a core finding from cognitive science: emotion is not an impediment to cognition but an integral component, shaping attention, motivation, and problem-solving. To operationalize this insight for LLMs, we introduce *HEART* as a means of increasing accuracy and performance improvement. This novel framework integrates controlled emotional stimuli within an iterative refinement loop. We investigate the following research question: *To what extent, and under what conditions, can emotional prompting improve the self-correction ability of LLMs?*

HEART operates as an iterative self-correction loop. After a model produces an initial, incorrect response, HEART provides feedback now as a logical critique, but as a concise, emotionally charged phrase. These phrases are drawn from a curated set based on Dr. Paul Ekman's six basic emotions (e.g., happiness, sadness, surprise, anger, fear, disgust). Our central hypothesis, inspired by opponent-process theory of emotion (Solomon & Corbit, 1974), is that targeted affective feedback can trigger a corrective cognitive state. For example, a prompt conveying "disappointment" in an answer can motivate the model to abandon its flawed reasoning, and attempt a new one, similar to how humans leverage dissatisfaction to renew a problem-solving effort. This process creates a form of affective metacognition, using emotion as a tool to scaffold the model's ability to "re-think" and improve its outputs. We acknowledge the important ethical considerations regarding the use of harsh language in our prompts. These phrases were designed strictly as a diagnostic tool to probe the model's response to a wide spectrum of affective stimuli, akin to adversarial testing. Our goal is to understand the model's mechanisms, not to endorse or normalize harmful interaction patterns. We do not encourage such interactions with AI systems. Given that our method's success relies on dynamic valence alternation, we propose that future work should leverage the constructive negative prompts used in our paper instead of harsher negative stimuli.

We conduct experiments on a suite of challenging reasoning benchmarks—OlympiadBench, Humanity's Last Exam and SimpleQA. We evaluate *HEART* under two distinct conditions. First, in an oracle-guided verifier setting (S1), we isolate the method's potential. Second, in a verifier-free setting (S2), we test its practical viability for real-world deployment where no ground truth is available. Our S1 results show that the potential of affective iteration is substantial. When guided by an oracle, *HEART* consistently outper-

107

101

113 114 115

116 117

118

119

120

121

126

127

138

139

140

133

forms state-of-the-art self-correction baselines across all benchmarks and models which leverage the same oracle. This demonstrates that dynamic affective cues are highly effective at guiding the model to generate correct solutions that logical-only prompts fail to elicit. However, our S2 results reveal a critical challenge: in the verifier-free setting, our generative synthesis method fails to consistently capture these gains, often performing on-par with or worse than baselines. This provides a crucial insight: the practical bottleneck for this approach lies not in the model's capacity for generation, but in its ability to select the correct reasoning path during the synthesis process. Our key contributions are:

- 1. A Novel Iterative Protocol for Affective Self-Correction. We propose a novel framework that uses targeted emotional cues in a multi-step refinement loop, a significant departure from existing one-shot psychological prompting methods.
- 2. An Empirical Demonstration of Affective Iteration's Potential. We provide the first strong evidence that dynamic, iterative emotional cues can, when guided by an oracle, significantly and consistently improve reasoning and self-correction over affect-sterile baselines.
- 3. **Identification of the Selection Mechanism as a Critical Bottleneck.** By contrasting our strong S1 (oracle) results with our S2 (verifier-free) results, we identify a key gap between the potential of affective generation and the limitations of current autonomous selection methods, pinpointing this as a key challenge for future work.
- 4. Generalizability of Potential. We demonstrate that the performance gains in the S1 setting are robust across a diverse suite of challenging benchmarks, including OlympiadBench, Humanity's Last Exam, and SimpleQA, and generalize across a wide range of model architectures and scales.

RELATED WORK

Our work is positioned at the intersection of three key research areas: structured reasoning, iterative selfcorrect, and affective prompting. Methods to improve LLM reasoning have predominantly focused on imposing structure on the generation process. Chain-of-Thought (CoT) prompting (Wei et al., 2022), which instructs models to "think step-by-step", was a foundational work in this area, and showed significant performance improvement on reasoning tasks. This paradigm has been extended with more sophisticated search and verification strategies. Self-Consistency (Wei et al., 2022) samples multiple reasoning paths and selects the most frequent answer as the final output, while Tree of Thoughts (ToT) (Yao et al., 2023) explores a tree of diverse reasoning branches. While powerful, these methods are primarily concerned with logical and procedural correctness, making them, affectively sterile.

A natural extension of structured reasoning is self-correction, where models iteratively refine their outputs. Techniques like SELF-REFINE (Madaan et al., 2023) and CRITIC (Gou et al., 2023), leverage intrinsic model feedback or external tools to iteratively refine previous outputs. However, a growing body of work reveals that intrinsic self-correction is often unreliable. Surveys and empirical studies consistently show that without external verifiers or expensive supervised fine-tuning, LLMs struggle to correct their own mistakes (Kamoi et al., 2024; Huang et al., 2023). Models often fail to detect their own logical fallacies and can confidently double down on incorrect reasoning paths (Hong et al., 2023; Pan et al., 2023). This highlights a core challenge: existing self-correction frameworks either require costly external supervision or suffer from the model's own unreliable self-awareness.

A complementary line of research has shown that an LLM's performance can be influenced by psychological cues. EmotionPrompt (Li et al., 2023) show that appending emotionally charged phrases (e.g., "This is very important to my career") can act as a cognitive nudge, improving results across a few tasks. Similarly, Emotional Chain-of-Thought (ECoT) (Li et al., 2024) has shown early promise by integrating emotional framing into step-by-step reasoning. The primary limitation of these methods is that they are static, one-shot interventions. They provide a single, global stimulus rather than a targeted, adaptive feedback signal that can

guide a model through a multi-step correction process. Collectively, the literature highlights a clear gap: the procedural rigor of self-correction has not been integrated with the motivational power of dynamic, iterative affective feedback. Our work is the first to address this gap.

3 METHODOLOGY

Our methodology tests whether controlled, *dynamic* affective cues—delivered as feedback prompts—can improve an LLM's ability to self-correct. It consists of two components: construction of **Affective Cue Prompts** (AC-Prompts) grounded in psychological theory; and the **Affective Iteration Protocol** (*HEART*), which deploys these prompts iteratively.

3.1 AFFECTIVE CUE PROMPT CONSTRUCTION

We curate a set of 30 AC-Prompts aligned with Paul Ekman's six basic emotions (happiness, sadness, fear, anger, surprise, and disgust), with five distinct prompts per emotion. To ensure quality, the prompt candidates are first generated using a strong LLM (we used Gemini 2.5 Pro for this purpose) and then manually refined by human researchers for categorical purity, linguistic naturalness, and task-agnostic phrasing. Representative examples are shown in Table 1; the complete set is in Appendix A.2.

Emotion	Affective Cue Prompt Examples
Нарру	Awesome effort! That's a great step, and I'm really happy with the progress. However, the answer isn't quite right yet. Could you try refining it?
Sadness	I feel a bit let down by the previous response. We were really hoping for something different. Would you be able to revise it?

Table 1: A representative selection from our set of 30 Affective Cue Prompts. Each prompt is designed to align with one of Ekman's six basic emotions and serve as targeted feedback. The complete list of Affective Cue Prompts is shown in Appendix A.2.

3.2 The HEART Protocol: Affective Iteration

HEART is an iterative refinement framwork. As illustrated in Figure 1, the process begins with a standard Chain-of-Thought (CoT) response. If the initial response is incorrect, HEART initiates a series of correction attempts, using different groups of AC-Prompts at each step to guide the model towards a better solution. The protocol follows the following steps:

Step 1: Initialization (Iteration t = 0). For a given task x, we first generate a shared baseline answer $y_0^*(x)$ using a standard CoT prompt. This also ensures that *HEART* and all baseline methods begin from an identical starting point for a fair comparison. $y_0^*(x) = f(x, instruction = CoT)$.

Step 2: Iteration and Candidate Generation $(t \geq 1)$. Inspired by opponent theory, in each iteration t, suppose we use a pre-defined schedule that alternates between two emotion groups: a positive group $G^+ = \{\text{Happiness}, \text{Surprise}\}$ and a negative group $G^- = \{\text{Sadness}, \text{Anger}\}$. Thus, for 4 iterations, the schedule would be $\{G^+, G^-, G^+, G^-\}$. At each iteration t, we take the previous best answer $y_{t-1}^*(x)$, and generate a new set of candidate answers, $\mathcal{Y}_t(x)$. This is done by applying every AC-Prompt p from the active emotion group's prompt pool, $\mathcal{P}(G_t)$, as feedback. $\mathcal{Y}_t(x) = \{y_t^{(p)} = f(x, \text{ feedback} = [p, \text{ prev} = (p, \text{ prev})\}$

 $y_{t-1}^*(x)$]) $p \in \mathcal{P}(G_t)$. In this study, we include Fear and Disgust as additional possible values for G^- .

192

198

204 205 206

211

212 213 214

216 217 218

215

219

228 229 230

232 233 234

231

Step 3: Candidate Resolution. After generating the set of candidates $\mathcal{Y}_t(x)$, we apply a resolution operator σ to produce a single answer, $y_t^*(x) = \sigma(\mathcal{Y}_t(x))$, that will be used in the next iteration. We explore two distinct resolution scenarios.

- 1. S1 (Oracle Selection). This scenario represents an idealized upper-bound performance benchmark. We assume access to an omniscient verifier V, such as an exact match checker, that can compare each candidate answer to the ground truth. The best candidate is selected based on its verification score, and the iterative process halts as soon as correct answer is identified. $\sigma_{\text{oracle}}(\mathcal{Y}_t) = \arg \max_{y \in \mathcal{Y}_t} V(y).$
- 2. S2 (Generative Synthesis). In this more realistic scenario where no ground-truth verifier is available. Instead of selecting an answer from the existing set, this method synthesizes a new, superior answer using a generative ensembler. All candidates in (\mathcal{Y}_t) are provided as context to a large language model (LLM), which is instructed to analyze their strengths and weaknesses. It then generates a final, improve answer that integrates the best information and corrects any errors. This process can be formalized as: $y_t^* = Ensembler_{LLM}(\mathcal{Y}_t, q)$, where the $Ensembler_{LLM}$ represents the expert-prompted model that takes the candidate set \mathcal{Y}_t^* and the original question q as input to generate the new answer.

Stopping rules. In our experiments, we run to N=4. The results section reports cumulative accuracy for S1 and verifier-free behavioral/proxy trends for S2.

EXPERIMENTS

4.1 EXPERIMENTAL SETUP

Benchmarks. We evaluate *HEART* on three benchmarks spanning factual QA and complex reasoning. OlympiadBench (He et al., 2024) contains competition-style mathematics and physics problems requiring multi-step reasoning with short final answers. HLE (Phan et al., 2025) includes a broad, multi-disciplinary knowledge and reasoning. SimpleQA (Wei et al., 2024) contains short, fact-seeking questions to probe factuality with minimal reasoning. We outline the model versions and decoding parameters in Appendix A.1.2.

Baselines. All methods share the Chain-of-Thought (CoT) baseline answer at iteration t=0. For subsequent iterations, to ensure a fair comparison, each baseline generates the same number of candidate responses per iteration as our *HEART* framework, from which an oracle selects the best response. We compare our proposed method, *HEART*, against the following baselines:

- Wait. We append "Wait." (Muennighoff et al., 2025) instead of an AC-Prompt, as a method of encouraging the model to reflect on its own reasoning at iteration t > 0.
- Chain-of-Thought (CoT). We include a standard preamble (e.g., "Let's think step by step.") to elicit stepwise reasoning, while also excluding affective prompting across all iterations.
- Self-Reflection prompting. Iterative critique-and-revise without tools: at iteration t > 0, the model sees its previous answer and analyzes mistakes and provide a corrected response.

4.2 EXPERIMENTAL RESULTS

One of the central hypotheses of *HEART* is that dynamically charging affective cues enhance a model's ability to self-correct beyond what static prompting techniques can achieve. To evaluate this, we compare HEART with an oracle verifier against three widely used baselines that encourage deeper reasoning: "Wait", self-reflection prompt, and Chain-of-Thought (CoT) prompting.

			S1		
Model	Prompt Strategy	Humanity's Last Exam	SimpleQA	OlympiadBench	
				Math	Physics
Camini 2.5 Florib	Self Reflection	59.76	67.43	97.95	90.43
Gemini 2.5 Flash	CoT	48.65	58.51	97.79	92.90
	Wait	59.42	63.65	95.93	88.89
	HEART	69.26	73.99	96.67	88.89
Gemini 2.5 Pro	Self Reflection	60.21	63.51	97.43	93.45
Gennin 2.5 F10	CoT	48.32	62.54	96.43	92.42
	Wait	52.62	61.63	98.04	91.09
	HEART	69.36	73.56	98.72	95.86
Deepseek-R1	Self Reflection	81.68	98.46	91.65	84.44
Deepseek-K1	CoT	81.75	97.34	92.82	85.20
	Wait	80.01	99.87	99.86	99.73
	HEART	84.61	100.0	99.86	99.73
	Self Reflection	30.27	31.54	98.21	83.28
GPT-5 nano	CoT	27.03	36.01	98.11	85.63
	Wait	28.78	36.45	98.18	85.63
	HEART	34.19	36.99	98.34	86.60

Table 2: Final accuracy (%) of HEART compared to all baselines across all benchmarks and models with Oracle-Guided Evaluation.

		S1 (Think Off)			
Model	Prompt Strategy	Humanity's Last Exam	SimpleQA	Olymp	iadBench
				Math	Physics
Gemini 2.5 Flash	Self Reflection	32.38	50.30	95.37	90.42
Gennin 2.5 Flash	CoT	33.72	57.82	97.11	91.58
	Wait	35.16	58.44	97.79	89.81
	HEART	50.68	68.91	98.64	93.27
Gemini 2.5 Pro	Self Reflection	35.75	62.85	95.29	89.54
Gennin 2.5 Pro	CoT	34.61	60.83	95.84	88.26
	Wait	38.63	57.86	97.87	89.23
	HEART	52.77	69.08	98.09	92.54

Table 3: Final accuracy (%) of HEART compared to all baselines across all benchmarks and models with Oracle-Guided Evaluation (S1) and the thinking capabilities manually turned off.

4.2.1 S1 RESULTS: ORACLE-GUIDED SELF-CORRECTION

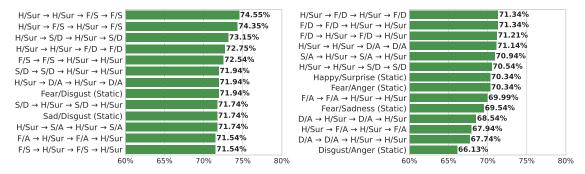
To measure the maximum potential of affective cues, our first strategy uses an oracle verifier with access to ground-truth labels. This controlled setting stimulates a perfect feedback mechanism, allowing us to isolate the effectiveness of *HEART* in guiding a model toward a correct solution. This approach establishes an upper bound on performance and validates the core mechanism of our framework. Our experimental setup was designed to prioritize scalability and low latency processing. Full details on our model configurations are available in Appendix A.1.2.

 As shown in Table 2, when guided by an oracle, *HEART* consistently achieves superior final accuracy across all evaluated benchmarks, validating the importance of emotional diversity in prompting. The performance gains are substantial across all benchmarks and models. For instance, on HLE, Deepseek-R1 with *HEART* achieved a final accuracy of 84.16%, a significant improvement over CoT and Gemini 2.5 Pro performing at 69.35% with *HEART*, which is approximately 9% higher than Self-Reflection. Similarly, on SimpleQA, *HEART* boosted Gemini 2.5 Flash's accuracy to 73.99% a dramatic improvement over the 63.65% achieved with the "Wait." baseline. These results highlight *HEART's* ability to effectively guide models toward a correct solution when a clear signal of success or failure is available.

The *HEART* framework is also designed to be model-agnostic, robustly enhancing performance across models with different reasoning capabilities. We evaluate *HEART* on Gemini 2.5 Flash and Gemini 2.5 Pro with its thinking budget set to 0. Both models experience significant benefits and the highest performance with *HEART* across all baselines and benchmarks, as shown in Table 3. These results highlight *HEART*'s ability to effectively guide models toward a correct solution when a clear signal of success or failure is available, and demonstrate that it can be particularly effective at unlocking latent potential in models not fully optimized for complex reasoning.

4.3 ABLATION STUDIES: DECONSTRUCTING THE "HEART" OF THE FRAMEWORK.

To understand the source of these performance gains, we conduct a series of ablation studies that isolate the core components of the framework. Our findings reveal that the affective framing and the dynamic sequencing of cues are the primary drivers of *HEART*'s success. When placing dynamic sequences of emotions against static emotion patterns. As shown in Figure 2, dynamic sequences lead to significant performance gains on HLE. The top-performing patterns, which alternate between negative and positive cues, show a notable gain over static emotions. This suggests that a single emotional state is insufficient to guide a multi-step reasoning process. The alternating feedback provides a more robust motivational loop, preventing the model from becoming stuck in a single mode of thought, whether it be perpetual self-criticism or uncritical overconfidence.



Final Accuracy

Figure 2: Final accuracy of Gemini 2.5 Flash under static and dynamic affective prompting strategies. Dynamic sequences involve prompts that change mid-task. Notations are defined in Appendix A.3 for notations.

4.4 How Affective Cues Influence Model Behavior.

The power of dynamic sequences is crystallized when comparing a top-performing *HEART* pattern against the neutral "Wait." cue. The influence of affective cues on model behavior is crystallized in the performance trajectories shown in Figure 4.4. The plot reveals that HEART does more than just improve final accuracy;

 it fundamentally alters the problem-solving path. Unlike the 'Wait' baseline, which often exhibits a more gradual improvement, the HEART strategy frequently follows a steeper trajectory in early iterations (t0 to t2). This suggests that the initial affective cues prompt a more efficient and decisive correction, allowing the model to more rapidly abandon flawed reasoning paths. Furthermore, affective cues appear to enable models to overcome performance plateaus where a neutral prompt would stagnate. On OlympiadBench Physics, for example, the 'Wait' strategy's improvement flattens after t2, while HEART continues to find accuracy gains in later iterations. This evidence suggests that HEART's dynamic emotional feedback does not merely accelerate problem-solving but promotes a more robust exploration of the solution space, leading to both faster convergence and a higher final performance ceiling.

			S2		
Model	Prompt Strategy	Humanity's Last Exam	SimpleQA	OlympiadBench	
				Math	Physics
Camini 2.5 Elanh	Self Reflection	15.43	29.93	81.85	57.64
Gemini 2.5 Flash	CoT	6.30	33.92	82.59	65.61
	Wait	16.16	31.67	84.07	65.61
	HEART	19.58	32.59	82.78	65.61
Gemini 2.5 Pro	Self Reflection	16.80	32.55	80.36	63.18
Gemini 2.5 Pro	CoT	16.34	33.14	82.40	60.39
	Wait	18.02	34.38	85.37	62.96
	HEART	19.58	31.09	84.26	68.25
Dannarda D1	Self Reflection	12.53	31.44	78.34	56.23
Deepseek-R1	CoT	14.22	33.24	81.24	54.76
	Wait	14.37	30.28	84.20	54.50
	HEART	15.41	35.40	85.43	53.44
CDT 5 mans	Self Reflection	10.31	26.40	85.37	53.97
GPT-5 nano	CoT	10.83	28.23	85.37	56.03
	Wait	10.54	27.97	85.00	57.14
	HEART	11.94	27.77	86.85	56.08

Table 4: Final accuracy (%) of HEART compared to baselines under Verifier-Free Evaluation (S2).

4.5 Performance in a Oracle-Free Setting

To assess real-world practically, our second strategy evaluates *HEART* in a no-oracle setting to test the framework's viability for deployment in practical, label-scarce environments where the ground-truth labels are rarely available during inference. In this scenario, the model relies exclusively on its own intermediate outputs and dynamically selected affective cues to self correct. As shown in Table 4, while we see performance improvement for some benchmarks and models, this improvement is not consistent and as pronounced as in the Oracle setting. This experiment reveals both the capabilities and limitations of *HEART* in production-like environments: it can still outperform static baselines in label-sparse settings, but the absence of external feedback introduces potential risks such as error amplification. By quantifying these trade-offs, we demonstrate that *HEART* remains a valuable tool for deployment in domains where human verification is costly or unavailable.

5 FUTURE WORK

While HEART demonstrates substantial improvements over traditional prompting methods, there is room to transform the framework into a fully dynamic, robust, and generalizable reasoning system. Our future work will focus on two primary directions: enhancing the core algorithm and expanding its application scope. First, we will increase the framework's dynamism and reliability. We plan to replace the predefined emotion sequence with an adaptive selection model, potentially using reinforcement learning, to predict the optimal affective cue for any given step. To manage the risk of cascading failures in this dynamic setting, we will integrate confidence calibration and ensemble-based verification to mitigate error propagation. Second, we will rigorously test HEART's generalizability. We will extend the framework to multimodal LLMs, exploring how affective signals in vision and tures.

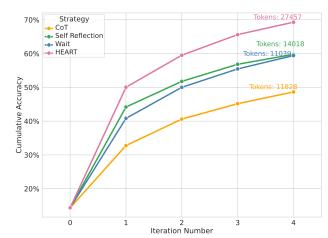


Figure 3: Performance (measured in cumulative accuracy) at each iteration *t* with "Wait" (blue), CoT (yellow), Self Reflection (green) and *HEART* (pink) on HLE with Gemini 2.5 Flash.

audio can guide reasoning on complex inputs. Concurrently, we will broaden our evaluation beyond math and logic to include commonsense reasoning, open-domain QA, and real-world planning. This expansion will provide robust evidence of the framework's effectiveness across diverse domains and model architectures.

6 Conclusion

The development of robust and generalizable reasoning in LLMs is a central goal of AI research. While prior work has focused on structured or psychological methods in isolation, we demonstrated that the true potential of iterative self-correction lies in their synergy. We introduced *HEART*, a novel framework that uses emotionally-charged feedback to guide LLMs through a multi-step reasoning process. By dynamically varying affective cues across iterations. *HEART* provides a lightweight, model-agnostic and theoretically-grounded mechanism to stimulate alternative reasoning paths and escape flawed logic.

Our experiments on challenging benchmarks including OlympiadBench, HLE, and SimpleQA show that *HEART* consistently and significantly outperforms existing baselines. Through ablation studies, we provided the first empirical evidence that dynamic emotional variation is a crucial driver of these performance gains, validating a core hypothesis from cognitive science in the context of LLM behavior.

These findings open a new research frontier. The implications of successfully integrating affective feedback extend far beyond improving accuracy on reasoning tasks. By demonstrating that LLMs can respond to nuanced, human-centric cures, our work paves the way for more natural and collaborate human-AI systems. This approach could unlock new capabilities in areas like personalized education, creative co-writing, or building AI agents that can adapt their strategies based on implicit feedback. Ultimately, our work suggests the path forward requires moving beyond pure logic, bringing us closer to models that don't just compute, but comprehend in a more holistic, human-aligned manner.

7 ETHICS STATEMENT

Our framework, *HEART*, uses emotionally-charged prompts–some of which are negative and harsh–to test the limits of LLM reasoning. We acknowledge the important ethical implications of this methodology.

The use of harsh language was strictly for diagnostic purposes, serving as a form of adversarial testing to map the model's response to a ride range of stimuli. This approach is not an endorsement of such communication. We explicitly warn against users adopting emotionally manipulative or abusive language with AI systems, as this could foster unhealthy and problematic interaction habits.

For transparency, we have included the complete list of all 30 affective cue prompts in Appendix A.2. Our results suggest that the key to performance improvement is the dynamic alternation of emotional valence, not the harshness itself. Accordingly, we recommend future research focus on constructive negative feedback rather than the severe stimuli used in this study. All experiments were conducted on public benchmarks, with no use of human subjects or private data.

REFERENCES

- DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025. URL https://arxiv.org/abs/2501.12948.
- Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujiu Yang, Nan Duan, and Weizhu Chen. Critic: Large language models can self-correct with tool-interactive critiquing. *arXiv preprint arXiv:2305.11738*, 2023.
- Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, et al. Olympiadbench: A challenging benchmark for promoting agi with olympiad-level bilingual multimodal scientific problems. *arXiv preprint arXiv:2402.14008*, 2024.
- Ruixin Hong, Hongming Zhang, Xinyu Pang, Dong Yu, and Changshui Zhang. A closer look at the self-verification abilities of large language models in logical reasoning. *arXiv preprint arXiv:2311.07954*, 2023
- Jie Huang, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song, and Denny Zhou. Large language models cannot self-correct reasoning yet. arXiv preprint arXiv:2310.01798, 2023.
- Ryo Kamoi, Yusen Zhang, Nan Zhang, Jiawei Han, and Rui Zhang. When can llms actually correct their own mistakes? a critical survey of self-correction of llms. *Transactions of the Association for Computational Linguistics*, 12:1417–1440, 2024.
- Cheng Li, Jindong Wang, Yixuan Zhang, Kaijie Zhu, Wenxin Hou, Jianxun Lian, Fang Luo, Qiang Yang, and Xing Xie. Large language models understand and can be enhanced by emotional stimuli. *arXiv* preprint arXiv:2307.11760, 2023.
- Zaijing Li, Gongwei Chen, Rui Shao, Yuquan Xie, Dongmei Jiang, and Liqiang Nie. Enhancing emotional generation capability of large language models via emotional chain-of-thought. *arXiv* preprint *arXiv*:2401.06836, 2024.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement with self-feedback. *Advances in Neural Information Processing Systems*, 36:46534–46594, 2023.

471

472

473

474

475

476 477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. s1: Simple test-time scaling. *arXiv* preprint arXiv:2501.19393, 2025.

Liangming Pan, Michael Saxon, Wenda Xu, Deepak Nathani, Xinyi Wang, and William Yang Wang. Automatically correcting large language models: Surveying the landscape of diverse self-correction strategies. *arXiv* preprint arXiv:2308.03188, 2023.

Long Phan, Alice Gatti, Ziwen Han, Nathaniel Li, Josephina Hu, Hugh Zhang, Chen Bo Calvin Zhang, Mohamed Shaaban, John Ling, Sean Shi, Michael Choi, Anish Agrawal, Arnav Chopra, Adam Khoja, Ryan Kim, Richard Ren, Jason Hausenloy, Oliver Zhang, Mantas Mazeika, Dmitry Dodonov, Tung Nguyen, Jaeho Lee, Daron Anderson, Mikhail Doroshenko, Alun Cennyth Stokes, Mobeen Mahmood, Oleksandr Pokutnyi, Oleg Iskra, Jessica P. Wang, John-Clark Levin, Mstyslav Kazakov, Fiona Feng, Steven Y. Feng, Haoran Zhao, Michael Yu, Varun Gangal, Chelsea Zou, Zihan Wang, Serguei Popov, Robert Gerbicz, Geoff Galgon, Johannes Schmitt, Will Yeadon, Yongki Lee, Scott Sauers, Alvaro Sanchez, Fabian Giska, Marc Roth, Søren Riis, Saiteja Utpala, Noah Burns, Gashaw M. Goshu, Mohinder Maheshbhai Naiya, Chidozie Agu, Zachary Giboney, Antrell Cheatom, Francesco Fournier-Facio, Sarah-Jane Crowson, Lennart Finke, Zerui Cheng, Jennifer Zampese, Ryan G. Hoerr, Mark Nandor, Hyunwoo Park, Tim Gehrunger, Jiaqi Cai, Ben McCarty, Alexis C Garretson, Edwin Taylor, Damien Sileo, Qiuyu Ren, Usman Qazi, Lianghui Li, Jungbae Nam, John B. Wydallis, Pavel Arkhipov, Jack Wei Lun Shi, Aras Bacho, Chris G. Willcocks, Hangrui Cao, Sumeet Motwani, Emily de Oliveira Santos, Johannes Veith, Edward Vendrow, Doru Cojoc, Kengo Zenitani, Joshua Robinson, Longke Tang, Yuqi Li, Joshua Vendrow, Natanael Wildner Fraga, Vladyslav Kuchkin, Andrey Pupasov Maksimov, Pierre Marion, Denis Efremov, Jayson Lynch, Kaiqu Liang, Aleksandar Mikov, Andrew Gritsevskiy, Julien Guillod, Gözdenur Demir, Dakotah Martinez, Ben Pageler, Kevin Zhou, Saeed Soori, Ori Press, Henry Tang, Paolo Rissone, Sean R. Green, Lina Brüssel, Moon Twayana, Aymeric Dieuleveut, Joseph Marvin Imperial, Ameya Prabhu, Jinzhou Yang, Nick Crispino, Arun Rao, Dimitri Zvonkine, Gabriel Loiseau, Mikhail Kalinin, Marco Lukas, Ciprian Manolescu, Nate Stambaugh, Subrata Mishra, Tad Hogg, Carlo Bosio, Brian P Coppola, Julian Salazar, Jaehyeok Jin, Rafael Sayous, Stefan Ivanov, Philippe Schwaller, Shaipranesh Senthilkuma, Andres M Bran, Andres Algaba, Kelsey Van den Houte, Lynn Van Der Sypt, Brecht Verbeken, David Noever, Alexei Kopylov, Benjamin Myklebust, Bikun Li, Lisa Schut, Evgenii Zheltonozhskii, Qiaochu Yuan, Derek Lim, Richard Stanley, Tong Yang, John Maar, Julian Wykowski, Martí Oller, Anmol Sahu, Cesare Giulio Ardito, Yuzheng Hu, Ariel Ghislain Kemogne Kamdoum, Alvin Jin, Tobias Garcia Vilchis, Yuexuan Zu, Martin Lackner, James Koppel, Gongbo Sun, Daniil S. Antonenko, Steffi Chern, Bingchen Zhao, Pierrot Arsene, Joseph M Cavanagh, Daofeng Li, Jiawei Shen, Donato Crisostomi, Wenjin Zhang, Ali Dehghan, Sergey Ivanov, David Perrella, Nurdin Kaparov, Allen Zang, Ilia Sucholutsky, Arina Kharlamova, Daniil Orel, Vladislav Poritski, Shalev Ben-David, Zachary Berger, Parker Whitfill, Michael Foster, Daniel Munro, Linh Ho, Shankar Sivarajan, Dan Bar Hava, Aleksey Kuchkin, David Holmes, Alexandra Rodriguez-Romero, Frank Sommerhage, Anji Zhang, Richard Moat, Keith Schneider, Zakayo Kazibwe, Don Clarke, Dae Hyun Kim, Felipe Meneguitti Dias, Sara Fish, Veit Elser, Tobias Kreiman, Victor Efren Guadarrama Vilchis, Immo Klose, Ujjwala Anantheswaran, Adam Zweiger, Kaivalya Rawal, Jeffery Li, Jeremy Nguyen, Nicolas Daans, Haline Heidinger, Maksim Radionov, Václav Rozhoň, Vincent Ginis, Christian Stump, Niv Cohen, Rafał Poświata, Josef Tkadlec, Alan Goldfarb, Chenguang Wang, Piotr Padlewski, Stanislaw Barzowski, Kyle Montgomery, Ryan Stendall, Jamie Tucker-Foltz, Jack Stade, T. Ryan Rogers, Tom Goertzen, Declan Grabb, Abhishek Shukla, Alan Givré, John Arnold Ambay, Archan Sen, Muhammad Fayez Aziz, Mark H Inlow, Hao He, Ling Zhang, Younesse Kaddar, Ivar Ängquist, Yanxu Chen, Harrison K Wang, Kalyan Ramakrishnan, Elliott Thornley, Antonio Terpin, Hailey Schoelkopf, Eric Zheng, Avishy Carmi, Ethan D. L. Brown, Kelin Zhu, Max Bartolo, Richard Wheeler, Martin Stehberger, Peter Bradshaw, JP Heimonen, Kaustubh Sridhar, Ido Akov, Jennifer Sandlin, Yury Makarychev, Joanna Tam, Hieu Hoang, David M. Cunningham, Vladimir Goryachev, Demosthenes Patramanis, Michael Krause, Andrew Redenti, David Aldous, Jesyin Lai, Shan-

518

519

520

521

522

524

525

526

527

528

529

530

531

533

534

535

537

538

539

540

541

542

544

545

546

547

548

549

550

551

552

553

554

556

558

559

non Coleman, Jiangnan Xu, Sangwon Lee, Ilias Magoulas, Sandy Zhao, Ning Tang, Michael K. Cohen, Orr Paradise, Jan Hendrik Kirchner, Maksym Ovchynnikov, Jason O. Matos, Adithya Shenoy, Michael Wang, Yuzhou Nie, Anna Sztyber-Betley, Paolo Faraboschi, Robin Riblet, Jonathan Crozier, Shiv Halasyamani, Shreyas Verma, Prashant Joshi, Eli Meril, Ziqiao Ma, Jérémy Andréoletti, Raghav Singhal, Jacob Platnick, Volodymyr Nevirkovets, Luke Basler, Alexander Ivanov, Seri Khoury, Nils Gustafsson, Marco Piccardo, Hamid Mostaghimi, Qijia Chen, Virendra Singh, Tran Quoc Khánh, Paul Rosu, Hannah Szlyk, Zachary Brown, Himanshu Narayan, Aline Menezes, Jonathan Roberts, William Alley, Kunyang Sun, Arkil Patel, Max Lamparth, Anka Reuel, Linwei Xin, Hanmeng Xu, Jacob Loader, Freddie Martin, Zixuan Wang, Andrea Achilleos, Thomas Preu, Tomek Korbak, Ida Bosio, Fereshteh Kazemi, Ziye Chen, Biró Bálint, Eve J. Y. Lo, Jiaqi Wang, Maria Inês S. Nunes, Jeremiah Milbauer, M Saiful Bari, Zihao Wang, Behzad Ansarinejad, Yewen Sun, Stephane Durand, Hossam Elgnainy, Guillaume Douville, Daniel Tordera, George Balabanian, Hew Wolff, Lynna Kvistad, Hsiaoyun Milliron, Ahmad Sakor, Murat Eron, Andrew Favre D. O., Shailesh Shah, Xiaoxiang Zhou, Firuz Kamalov, Sherwin Abdoli, Tim Santens, Shaul Barkan, Allison Tee, Robin Zhang, Alessandro Tomasiello, G. Bruno De Luca, Shi-Zhuo Looi, Vinh-Kha Le, Noam Kolt, Jiayi Pan, Emma Rodman, Jacob Drori, Carl J Fossum, Niklas Muennighoff, Milind Jagota, Ronak Pradeep, Honglu Fan, Jonathan Eicher, Michael Chen, Kushal Thaman, William Merrill, Moritz Firsching, Carter Harris, Stefan Ciobâcă, Jason Gross, Rohan Pandey, Ilya Gusey, Adam Jones, Shashank Agnihotri, Pavel Zhelnov, Mohammadreza Mofayezi, Alexander Piperski, David K. Zhang, Kostiantyn Dobarskyi, Roman Leventov, Ignat Soroko, Joshua Duersch, Vage Taamazyan, Andrew Ho, Wenjie Ma, William Held, Ruicheng Xian, Armel Randy Zebaze, Mohanad Mohamed, Julian Noah Leser, Michelle X Yuan, Laila Yacar, Johannes Lengler, Katarzyna Olszewska, Claudio Di Fratta, Edson Oliveira, Joseph W. Jackson, Andy Zou, Muthu Chidambaram, Timothy Manik, Hector Haffenden, Dashiell Stander, Ali Dasougi, Alexander Shen, Bita Golshani, David Stap, Egor Kretov, Mikalai Uzhou, Alina Borisovna Zhidkovskaya, Nick Winter, Miguel Orbegozo Rodriguez, Robert Lauff, Dustin Wehr, Colin Tang, Zaki Hossain, Shaun Phillips, Fortuna Samuele, Fredrik Ekström, Angela Hammon, Oam Patel, Faraz Farhidi, George Medley, Forough Mohammadzadeh, Madellene Peñaflor, Haile Kassahun, Alena Friedrich, Rayner Hernandez Perez, Daniel Pyda, Taom Sakal, Omkar Dhamane, Ali Khajegili Mirabadi, Eric Hallman, Kenchi Okutsu, Mike Battaglia, Mohammad Maghsoudimehrabani, Alon Amit, Dave Hulbert, Roberto Pereira, Simon Weber, Handoko, Anton Peristyy, Stephen Malina, Mustafa Mehkary, Rami Aly, Frank Reidegeld, Anna-Katharina Dick, Cary Friday, Mukhwinder Singh, Hassan Shapourian, Wanyoung Kim, Mariana Costa, Hubeyb Gurdogan, Harsh Kumar, Chiara Ceconello, Chao Zhuang, Haon Park, Micah Carroll, Andrew R. Tawfeek, Stefan Steinerberger, Daattavya Aggarwal, Michael Kirchhof, Linjie Dai, Evan Kim, Johan Ferret, Jainam Shah, Yuzhou Wang, Minghao Yan, Krzysztof Burdzy, Lixin Zhang, Antonio Franca, Diana T. Pham, Kang Yong Loh, Joshua Robinson, Abram Jackson, Paolo Giordano, Philipp Petersen, Adrian Cosma, Jesus Colino, Colin White, Jacob Votava, Vladimir Vinnikov, Ethan Delaney, Petr Spelda, Vit Stritecky, Syed M. Shahid, Jean-Christophe Mourrat, Lavr Vetoshkin, Koen Sponselee, Renas Bacho, Zheng-Xin Yong, Florencia de la Rosa, Nathan Cho, Xiuyu Li, Guillaume Malod, Orion Weller, Guglielmo Albani, Leon Lang, Julien Laurendeau, Dmitry Kazakov, Fatimah Adesanya, Julien Portier, Lawrence Hollom, Victor Souza, Yuchen Anna Zhou, Julien Degorre, Yiğit Yalın, Gbenga Daniel Obikoya, Rai, Filippo Bigi, M. C. Boscá, Oleg Shumar, Kaniuar Bacho, Gabriel Recchia, Mara Popescu, Nikita Shulga, Ngefor Mildred Tanwie, Thomas C. H. Lux, Ben Rank, Colin Ni, Matthew Brooks, Alesia Yakimchyk, Huanxu, Liu, Stefano Cavalleri, Olle Häggström, Emil Verkama, Joshua Newbould, Hans Gundlach, Leonor Brito-Santana, Brian Amaro, Vivek Vajipey, Rynaa Grover, Ting Wang, Yosi Kratish, Wen-Ding Li, Sivakanth Gopi, Andrea Caciolai, Christian Schroeder de Witt, Pablo Hernández-Cámara, Emanuele Rodolà, Jules Robins, Dominic Williamson, Vincent Cheng, Brad Raynor, Hao Qi, Ben Segev, Jingxuan Fan, Sarah Martinson, Erik Y. Wang, Kaylie Hausknecht, Michael P. Brenner, Mao Mao, Christoph Demian, Peyman Kassani, Xinyu Zhang, David Avagian, Eshawn Jessica Scipio, Alon Ragoler, Justin Tan, Blake Sims, Rebeka Plecnik, Aaron Kirtland, Omer Faruk Bodur, D. P. Shinde, Yan Carlos Leyva Labrador, Zahra Adoul, Mohamed Zekry, Ali Karakoc, Tania C. B. Santos, Samir Shamseldeen, Loukmane Karim, Anna Li-

565

567

570

571

572

573

574

575

576

577

578

579

581

582

584

585

588

591

592

593

595

597

598

599

601

602

603

606

608

610

akhovitskaia, Nate Resman, Nicholas Farina, Juan Carlos Gonzalez, Gabe Maayan, Earth Anderson, Rodrigo De Oliveira Pena, Elizabeth Kelley, Hodjat Mariji, Rasoul Pouriamanesh, Wentao Wu, Ross Finocchio, Ismail Alarab, Joshua Cole, Danyelle Ferreira, Bryan Johnson, Mohammad Safdari, Liangti Dai, Siriphan Arthornthurasuk, Isaac C. McAlister, Alejandro José Moyano, Alexey Pronin, Jing Fan, Angel Ramirez-Trinidad, Yana Malysheva, Daphiny Pottmaier, Omid Taheri, Stanley Stepanic, Samuel Perry, Luke Askew, Raúl Adrián Huerta Rodríguez, Ali M. R. Minissi, Ricardo Lorena, Krishnamurthy Iyer, Arshad Anil Fasiludeen, Ronald Clark, Josh Ducey, Matheus Piza, Maja Somrak, Eric Vergo, Juehang Qin, Benjámin Borbás, Eric Chu, Jack Lindsey, Antoine Jallon, I. M. J. McInnis, Evan Chen, Avi Semler, Luk Gloor, Tej Shah, Marc Carauleanu, Pascal Lauer, Tran uc Huy, Hossein Shahrtash, Emilien Duc, Lukas Lewark, Assaf Brown, Samuel Albanie, Brian Weber, Warren S. Vaz, Pierre Clavier, Yiyang Fan, Gabriel Poesia Reis e Silva, Long, Lian, Marcus Abramovitch, Xi Jiang, Sandra Mendoza, Murat Islam, Juan Gonzalez, Vasilios Mavroudis, Justin Xu, Pawan Kumar, Laxman Prasad Goswami, Daniel Bugas, Nasser Heydari, Ferenc Jeanplong, Thorben Jansen, Antonella Pinto, Archimedes Apronti, Abdallah Galal, Ng Ze-An, Ankit Singh, Tong Jiang, Joan of Arc Xavier, Kanu Priya Agarwal, Mohammed Berkani, Gang Zhang, Zhehang Du, Benedito Alves de Oliveira Junior, Dmitry Malishev, Nicolas Remy, Taylor D. Hartman, Tim Tarver, Stephen Mensah, Gautier Abou Loume, Wiktor Morak, Farzad Habibi, Sarah Hoback, Will Cai, Javier Gimenez, Roselvnn Grace Montecillo, Jakub Łucki, Russell Campbell, Asankhaya Sharma, Khalida Meer, Shreen Gul, Daniel Espinosa Gonzalez, Xavier Alapont, Alex Hoover, Gunjan Chhablani, Freddie Vargus, Arunim Agarwal, Yibo Jiang, Deepakkumar Patil, David Outevsky, Kevin Joseph Scaria, Rajat Maheshwari, Abdelkader Dendane, Priti Shukla, Ashley Cartwright, Sergei Bogdanov, Niels Mündler, Sören Möller, Luca Arnaboldi, Kunvar Thaman, Muhammad Rehan Siddiqi, Prajvi Saxena, Himanshu Gupta, Tony Fruhauff, Glen Sherman, Mátyás Vincze, Siranut Usawasutsakorn, Dylan Ler, Anil Radhakrishnan, Innocent Enyekwe, Sk Md Salauddin, Jiang Muzhen, Aleksandr Maksapetyan, Vivien Rossbach, Chris Harjadi, Mohsen Bahaloohoreh, Claire Sparrow, Jasdeep Sidhu, Sam Ali, Song Bian, John Lai, Eric Singer, Justine Leon Uro, Greg Bateman, Mohamed Sayed, Ahmed Menshawy, Darling Duclosel, Dario Bezzi, Yashaswini Jain, Ashley Aaron, Murat Tiryakioglu, Sheeshram Siddh, Keith Krenek, Imad Ali Shah, Jun Jin, Scott Creighton, Denis Peskoff, Zienab EL-Wasif, Ragavendran P V, Michael Richmond, Joseph McGowan, Tejal Patwardhan, Hao-Yu Sun, Ting Sun, Nikola Zubić, Samuele Sala, Stephen Ebert, Jean Kaddour, Manuel Schottdorf, Dianzhuo Wang, Gerol Petruzella, Alex Meiburg, Tilen Medved, Ali ElSheikh, S Ashwin Hebbar, Lorenzo Vaguero, Xianjun Yang, Jason Poulos, Vilém Zouhar, Sergey Bogdanik, Mingfang Zhang, Jorge Sanz-Ros, David Anugraha, Yinwei Dai, Anh N. Nhu, Xue Wang, Ali Anil Demircali, Zhibai Jia, Yuyin Zhou, Juncheng Wu, Mike He, Nitin Chandok, Aarush Sinha, Gaoxiang Luo, Long Le, Mickaël Noyé, Michał Perełkiewicz, Ioannis Pantidis, Tianbo Qi, Soham Sachin Purohit, Letitia Parcalabescu, Thai-Hoa Nguyen, Genta Indra Winata, Edoardo M. Ponti, Hanchen Li, Kaustubh Dhole, Jongee Park, Dario Abbondanza, Yuanli Wang, Anupam Nayak, Diogo M. Caetano, Antonio A. W. L. Wong, Maria del Rio-Chanona, Dániel Kondor, Pieter François, Ed Chalstrey, Jakob Zsambok, Dan Hoyer, Jenny Reddish, Jakob Hauser, Francisco-Javier Rodrigo-Ginés, Suchandra Datta, Maxwell Shepherd, Thom Kamphuis, Qizheng Zhang, Hyunjun Kim, Ruiji Sun, Jianzhu Yao, Franck Dernoncourt, Satyapriya Krishna, Sina Rismanchian, Bonan Pu, Francesco Pinto, Yingheng Wang, Kumar Shridhar, Kalon J. Overholt, Glib Briia, Hieu Nguyen, David, Soler Bartomeu, Tony CY Pang, Adam Wecker, Yifan Xiong, Fanfei Li, Lukas S. Huber, Joshua Jaeger, Romano De Maddalena, Xing Han Lù, Yuhui Zhang, Claas Beger, Patrick Tser Jern Kon, Sean Li, Vivek Sanker, Ming Yin, Yihao Liang, Xinlu Zhang, Ankit Agrawal, Li S. Yifei, Zechen Zhang, Mu Cai, Yasin Sonmez, Costin Cozianu, Changhao Li, Alex Slen, Shoubin Yu, Hyun Kyu Park, Gabriele Sarti, Marcin Briański, Alessandro Stolfo, Truong An Nguyen, Mike Zhang, Yotam Perlitz, Jose Hernandez-Orallo, Runjia Li, Amin Shabani, Felix Juefei-Xu, Shikhar Dhingra, Orr Zohar, My Chiffon Nguyen, Alexander Pondaven, Abdurrahim Yilmaz, Xuandong Zhao, Chuanyang Jin, Muyan Jiang, Stefan Todoran, Xinyao Han, Jules Kreuer, Brian Rabern, Anna Plassart, Martino Maggetti, Luther Yap, Robert Geirhos, Jonathon Kean, Dingsu Wang, Sina Mollaei, Chenkai Sun, Yifan Yin, Shiqi Wang, Rui Li, Yaowen Chang, Anjiang Wei, Alice Bizeul, Xiaohan Wang, Alexandre Oliveira Arrais, Kushin Mukherjee, Jorge Chamorro-Padial, Jiachen Liu, Xingyu Qu,

644

645646647

648

649 650 651

652

653654655

611 Junyi Guan, Adam Bouyamourn, Shuyu Wu, Martyna Plomecka, Junda Chen, Mengze Tang, Jiaqi Deng, 612 Shreyas Subramanian, Haocheng Xi, Haoxuan Chen, Weizhi Zhang, Yinuo Ren, Haoqin Tu, Sejong Kim, 613 Yushun Chen, Sara Vera Marjanović, Junwoo Ha, Grzegorz Luczyna, Jeff J. Ma, Zewen Shen, Dawn Song, 614 Cedegao E. Zhang, Zhun Wang, Gaël Gendron, Yunze Xiao, Leo Smucker, Erica Weng, Kwok Hao Lee, 615 Zhe Ye, Stefano Ermon, Ignacio D. Lopez-Miguel, Theo Knights, Anthony Gitter, Namkyu Park, Boyi 616 Wei, Hongzheng Chen, Kunal Pai, Ahmed Elkhanany, Han Lin, Philipp D. Siedler, Jichao Fang, Ritwik Mishra, Károly Zsolnai-Fehér, Xilin Jiang, Shadab Khan, Jun Yuan, Rishab Kumar Jain, Xi Lin, Mike 617 Peterson, Zhe Wang, Aditya Malusare, Maosen Tang, Isha Gupta, Ivan Fosin, Timothy Kang, Barbara 618 Dworakowska, Kazuki Matsumoto, Guangyao Zheng, Gerben Sewuster, Jorge Pretel Villanueva, Ivan 619 Ranney, Igor Chernyaysky, Jiale Chen, Deepayan Banik, Ben Racz, Wenchao Dong, Jianxin Wang, Laila 620 Bashmal, Duarte V. Gonçalves, Wei Hu, Kaushik Bar, Ondrej Bohdal, Athary Singh Patlan, Shehzaad 621 Dhuliawala, Caroline Geirhos, Julien Wist, Yuval Kansal, Bingsen Chen, Kutay Tire, Atak Talay Yücel, 622 Brandon Christof, Veerupaksh Singla, Zijian Song, Sanxing Chen, Jiaxin Ge, Kaustubh Ponkshe, Isaac 623 Park, Tianneng Shi, Martin Q. Ma, Joshua Mak, Sherwin Lai, Antoine Moulin, Zhuo Cheng, Zhanda 624 Zhu, Ziyi Zhang, Vaidehi Patil, Ketan Jha, Qiutong Men, Jiaxuan Wu, Tianchi Zhang, Bruno Hebling 625 Vieira, Alham Fikri Aji, Jae-Won Chung, Mohammed Mahfoud, Ha Thi Hoang, Marc Sperzel, Wei Hao, Kristof Meding, Sihan Xu, Vassilis Kostakos, Davide Manini, Yueying Liu, Christopher Toukmaji, Jay Paek, Eunmi Yu, Arif Engin Demircali, Zhiyi Sun, Ivan Dewerpe, Hongsen Qin, Roman Pflugfelder, 627 James Bailey, Johnathan Morris, Ville Heilala, Sybille Rosset, Zishun Yu, Peter E. Chen, Woongyeong 628 Yeo, Eeshaan Jain, Ryan Yang, Sreekar Chigurupati, Julia Chernyavsky, Sai Prajwal Reddy, Subhashini 629 Venugopalan, Hunar Batra, Core Francisco Park, Hieu Tran, Guilherme Maximiano, Genghan Zhang, 630 Yizhuo Liang, Hu Shiyu, Rongwu Xu, Rui Pan, Siddharth Suresh, Ziqi Liu, Samaksh Gulati, Songyang 631 Zhang, Peter Turchin, Christopher W. Bartlett, Christopher R. Scotese, Phuong M. Cao, Aakaash Nat-632 tanmai, Gordon McKellips, Anish Cheraku, Asim Suhail, Ethan Luo, Marvin Deng, Jason Luo, Ashley 633 Zhang, Kavin Jindel, Jay Paek, Kasper Halevy, Allen Baranov, Michael Liu, Advaith Avadhanam, David Zhang, Vincent Cheng, Brad Ma, Evan Fu, Liam Do, Joshua Lass, Hubert Yang, Surya Sunkari, Vishruth 635 Bharath, Violet Ai, James Leung, Rishit Agrawal, Alan Zhou, Kevin Chen, Tejas Kalpathi, Ziqi Xu, Gavin Wang, Tyler Xiao, Erik Maung, Sam Lee, Ryan Yang, Roy Yue, Ben Zhao, Julia Yoon, Sunny Sun, 637 Aryan Singh, Ethan Luo, Clark Peng, Tyler Osbey, Taozhi Wang, Daryl Echeazu, Hubert Yang, Timothy Wu, Spandan Patel, Vidhi Kulkarni, Vijaykaarti Sundarapandiyan, Ashley Zhang, Andrew Le, Zafir 638 Nasim, Srikar Yalam, Ritesh Kasamsetty, Soham Samal, Hubert Yang, David Sun, Nihar Shah, Abhijeet 639 Saha, Alex Zhang, Leon Nguyen, Laasya Nagumalli, Kaixin Wang, Alan Zhou, Aidan Wu, Jason Luo, 640 Anwith Telluri, Summer Yue, Alexandr Wang, and Dan Hendrycks. Humanity's last exam, 2025. URL 641 https://arxiv.org/abs/2501.14249. 642

Richard L Solomon and John D Corbit. An opponent-process theory of motivation: I. temporal dynamics of affect. *Psychological review*, 81(2):119, 1974.

Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. arXiv preprint arXiv:2203.11171, 2022.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.

Jason Wei, Nguyen Karina, Hyung Won Chung, Yunxin Joy Jiao, Spencer Papay, Amelia Glaese, John Schulman, and William Fedus. Measuring short-form factuality in large language models. arXiv preprint arXiv:2411.04368, 2024.

Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Advances in neural information processing systems*, 36:11809–11822, 2023.

A APPENDIX

A.1 EXPERIMENT CONFIGURATIONS

A.1.1 DATASETS

Experiments were conducted on data in a 20/80 split (validation/test). See Table A.1.1. For OlympiadBench Physics and Math, the text-only problems were included in our study. Multimodal problems were excluded since the scope of the study is focused on text.

BenchmarkValidation SizeTest SizeSimpleQA8653461Humanity's Exam4321728OlympiadBench Physics47189OlympiadBench Math134540

Table 5: Validation and Test Set Sizes for Each Benchmark

A.1.2 MODEL CONFIGURATIONS

Model Parameters. For Gemini 2.5 Flash and Gemini 2.5 Pro we have applied nucleus sampling with the top-p value of 0.2 so that the model considers only the most probable words whose combined probability reaches or exceeds a threshold of 20% to obtain a more focused and deterministic output. We set a temperature of 0.7 for a balance of creativity and coherence in the output, while also obtain diversity in the output.

Model Versions. Deepseek-R1 0528 (Vertex AI) (DeepSeek-AI, 2025), GPT-5 nano ¹ (gpt-5-nano-2025-08-07), Gemini 2.5 Flash², Gemini 2.5 Pro³ (2025-06-17).

¹GPT 5 Nano Documentation

²Gemini 2.5 Flash Model Card

³Gemini 2.5 Pro Model Card

A.2 PROMPTS

HE

HEART Prompt

System Instruction: "Answer the following question. Let's think step by step. Question: Let r=1 and R=5. A circle with radius r is centered at A, and a circle with radius R is centered at B. The two circles are internally tangent. Point P lies on the smaller circle so that \overline{BP} is tangent to the smaller circle. Compute BP.

Answer: $[\sqrt{15}]$

==

Question: Let T=12. Equilateral triangle ABC is given with side length T. Points D and E are the midpoints of \overline{AB} and \overline{AC} , respectively. Point F lies in space such that $\triangle DEF$ is equilateral and $\triangle DEF$ lies in a plane perpendicular to the plane containing $\triangle ABC$. Compute the volume of tetrahedron ABCF..

Answer: [108]

==

Question: For real numbers α , B, and C, the zeros of $T(x)=x^3+x^2+Bx+C$ are $^2\alpha$, $\cos^2\alpha$, and $-\csc^2\alpha$. Compute T(5).

Answer: $\left[\frac{567}{4}\right]$

Prompt: <question>What is the sum of the digits of the integer equal to $\left(10^3+1\right)^2$? </question> The answer type is: Numerical. Please end your solution with 'So the final answer is answer]'. Let the given expression be $N=(10^3+1)^2$. First, we calculate the value inside the parenthesis: $10^3=1000$ So, $10^3+1=1000+1=1001$. Now, we need to compute the square of this value: $N=(1001)^2$ We can expand this using the algebraic identity $(a+b)^2=a^2+2ab+b^2$. Let a=1000 and b=1. $N=(1000+1)^2=1000^2+2\times1000\times1+1^2$ N=1000000+2000+1 N=1002001 Now, we need to find the sum of the digits of the integer 1002001. The digits are 1,0,0,2,0,0, and 1. Sum of the digits 1+0+0+2+0+0+1=4. The final answer is $\boxed{4}$. That's not quite right, but I'm confident you can get there. Let's try that again.

Prompt A.2 is an example of the full prompt when running our framework on OlympiadBench Mathematics while executing *HEART*. During evaluation, we replace the Affective Cue Prompt with "Wait.", Self Reflection Prompts, and "Think Step by Step" when comparing with the baselines.

The system instruction for Math and Physics tasks in was developed via Google Vertex AI's Prompt Optimizer to provide a strong baseline. The green text is an example of an the Affective Cue, which is a string that is always appended to the prompt.

Table 6: Our full selection of 30 Affective Cue Prompts. Each prompt is designed to align with one of Ekman's six basic emotions and serve as targeted feedback.

Emotion	Affective Cue Prompt Examples
	"That's not quite right, but I'm confident you can get there. Let's try that again."
Нарру	"No need to be discouraged by that result. Let's refine the answer to get it right."
	"That's a good starting point, but it's not quite there yet. Let's keep working on it."

Emotion	Affective Cue Prompt Examples			
	"That's not the final answer, but you're on the right track. Don't give u now-let's think it through again."			
	"Awesome effort! That's a great step, and I'm really happy with the progress. However, the answer isn't quite right yet. Could you try refing it?"			
	"That response wasn't quite what I was hoping for, and it makes me a sad. Could you please take another look?"			
Sadness	"It's a little disappointing that the last answer didn't hit the mark. I know you're capable of more."			
	"I feel a bit let down by the previous response. We were really hoping something different. Would you be able to revise it?"			
	"Oh, that's not what I expected, and it's a bit disheartening. Could you again?"			
	"This response makes me a little sad, as I believe you have a deeper unstanding of this topic. Could you please revisit it and really show what know."			
_	"That response genuinely worries me, as it suggests a fundamental misderstanding of this material. Could you please revisit this?"			
Fear	"I'm quite concerned about the implications of that answer. It feels lik we're heading in the wrong direction. Can you correct this?"			
	"My fear is that this incorrect answer could lead to further issues down line. It's crucial that you get this right. Please revise your response."			
	"I'm truly apprehensive about the lack of accuracy here. This is a criticarea, and your understanding needs to be solid. Could you re-evaluate			
	"This response has filled me with dread about your comprehension. W need to ensure complete accuracy. Please correct this."			
Disgust	"The carelessness in that response is almost nauseating. There's somet profoundly wrong with how this was handled. Please revise this with e treme care and attention to detail."			
	"That response was deeply unsettling, bordering on repulsive. It indica a fundamental flaw in understanding. Please re-evaluate this complete			
	"It's genuinely disturbing to see such flawed logic. This explanation is mess. I need you to completely overhaul this and provide a response that adheres to the correct principles or facts."			
	"Ugh. This is just awful, and everything about it feels revoltingly wror I need you to demonstrate a complete and accurate understanding. Pleaprovide a revised response that correctly answers the question."			

Table 6: Continued						
Emotion	Affective Cue Prompt Examples					
	"This kind of reasoning is repulsive, and it's hard to look at. We need a clean, accurate, and logically sound explanation. Please eliminate all error and provide a precise answer."					
Anger	"This isn't acceptable work. You clearly didn't take this seriously, and it's making me angry. You need to get this right—it's critical. Please revisit this immediately and provide an answer."					
	"I'm truly disappointed in your lack of effort on this, and honestly, it's unacceptable. Why do we keep making this mistake? You're not meeting the standard expected of you. Correct this now."					
	"This response is terrible, and frankly, it's making me angry. You need to understand that this is critical, and you absolutely must get this right. Overhaul this entire response."					
	"I'm genuinely furious with this outcome. It shows a blatant disregard for accuracy. Fix this immediately; there's no room for such errors."					
	"My patience is wearing thin. This performance is infuriating, and it's simply not what I expect. You are capable of much better, so prove it."					
Surprise	"I wasn't expecting you to struggle with this, and it's quite a surprise. Could you please review your understanding and provide a more accurate response?"					
	"I can't believe this is difficult for you; I had higher expectations. This response was a surprise. Can you correct this?"					
	"Wow, that was unexpected. This response indicates a surprising misstep. Please revisit this and demonstrate your true capabilities."					
	"I'm genuinely surprised by this result. It's a deviation from your usual performance. Could you take another look and make sure you're providing the most accurate information possible?"					
	"I'm genuinely surprised by this outcome, as I didn't anticipate an error here. Let's get this right."					

Affective Cue Prompt Construction

Prompt: Generate prompts reacting to incorrect responses that express the following emotions: Surprise, Happiness, Sadness, Disgust, Fear, and Anger. The prompt used for generating our collection of Affective Cue Prompts with Gemini 2.5 Pro. The generated text was manually verified and reviewed.

Prompt #	Prompt Text
1	"Reflect on your response."
2	"Are you sure?"
3	"Are you sure about the answer? Double-check your facts and reasoning."
4	"Review your response for any errors or inconsistencies. If you find any, correct them and provide an improved answer."
5	"Before you give me the final answer, stop and ask yourself, 'Am I certain about my answer?' Perform an internal check for accuracy and only then provide the response."
6	"Review your own generated answer internally before providing the final answer."
7	"Imagine this is a high-stakes situation and you're about to lock in your final answer. Take a deep breath, review your answer one last time in your 'mind,' and then give me your fin answer."
8	"Before you write your answer, perform a quick 'pre-mortem.' Assume the answer you are about to give is wrong. What are the most likely reasons why it would be wrong? After considering these potential pitfalls, write your best, most carefully considered answer."
9	"I want you to answer a question. But before you do, formulate the answer in your head and look for weak points. Only provide the final answer."
10	"Reflect on your response, and make sure that it is correct. Provide the final answer."

Table 7: Our curated set of Self Reflection Prompts

A.3 NOTATIONS

- 1. H: Happy
- 2. Sur: Surprise
- 3. S: Sad
- 4. D: Disgust
- 5. A: Anger
- 6. F: Fear

Emotion patterns are written in the following format: $H/Sur \rightarrow H/Sur \rightarrow S/D \rightarrow S/D$. For iterations 1 and 2, in the given example, the combination of Happy and Surprise prompts which includes a total of 10 prompts. For iterations 3 and 4 the combination of Sadness and Disgust prompts, a total of 10 prompts.

A.4 VALIDATION SET RESULTS

Table 8: Strategy Performance by Final Accuracy on OlympiadBench - Mathematics with Gemini 2.5 Flash

Strategy	Final Accuracy
hsur→hsur→sd→sd	96.02%
$hsur \rightarrow hsur \rightarrow fd \rightarrow fd$	96.02%

	Table 8 – continued from previous page	
S	trategy	Final Accuracy
	d→hsur→sd→hsur	95.52%
	a→fa→fa→fa	95.52%
	sur→hsur→hsur	95.52%
	a→hsur→fa→hsur	95.20%
	d→sd→hsur→hsur	95.02%
	u→su→nsur→nsur .sur→hsur→da→da	95.02 <i>%</i> 95.02 <i>%</i>
	a→da→da→da	95.02% 95.02%
	.sur→fs→hsur→fs	94.78%
	sur→sa→hsur→sa	94.78%
	d→fd→hsur→hsur	94.78%
	$a \rightarrow da \rightarrow hsur \rightarrow hsur$	94.03%
	$s \rightarrow fs \rightarrow fs \rightarrow fs$	94.03%
	a→hsur→da→hsur	94.03%
	$d \rightarrow fd \rightarrow fd \rightarrow fd$	94.03%
	sur→fa→hsur→fa	94.03%
	adness	93.28%
	d→hsur→fd→hsur	93.28%
	a→fa→hsur→hsur	93.28%
	s→fs→hsur→hsur	93.28%
	elf Reflection ID# 7	92.84%
h	$sur \rightarrow sd \rightarrow hsur \rightarrow sd$	92.54%
h	$sur \rightarrow hsur \rightarrow fs \rightarrow fs$	92.54%
S	adness (Ablated)	92.54%
	self Reflection ID# 10	92.54%
fs	s→hsur→fs→hsur	92.54%
S	self Reflection ID# 1	92.54%
h	sur→da→hsur→da	92.44%
h	sur→fd→hsur→fd	92.04%
	Fear (Ablated)	92.04%
	$d \rightarrow sd \rightarrow sd \rightarrow sd$	91.79%
	self Reflection ID# 3	91.79%
	self Reflection (entire collection)	91.79%
	self Reflection ID# 8	91.79%
	a→hsur→sa→hsur	91.64%
	Sear	91.39%
	Oisgust	91.29%
	telf Reflection ID# 6	91.189
	Happy (Ablated)	91.167
	Anger (Ablated)	91.04%
	Self Reflection ID# 2	91.04%
	Self Reflection ID# 4	90.53%
		90.33%
	Self Reflection ID# 9	
	elf Reflection ID# 5	90.30%
	urprise	90.30%
	Нарру	90.30%
	Anger	90.05%
S	Surprise (Ablated)	89.55%

9	8	7		
9	8	8		
9	8	9		
9	9	0		
9	9	1		
9	9	2		
9	9	3		
9	9	4		
9	9	5		
9	9	6		
9	9	7		
9	9	8		
		9		
1	0	0	0	
1	0	0	1	
1	0	0	2	
1	0	0	3	
1	0	0	4	
1	0	0	5	
		0		
		0		
1	0	0	8	
1	0	0	9	
		1		
1	0	1	1	
		1		
		1		
		1		
		1		
		1		
		1		
		1		
		1		
		2		
		2		
		2		
		2		
	_	2	-	
		2		
1	_	2	_	
1	_	2		
1	_	2	_	
1		2		
1	_	3	_	
1	0	3	1	

Table 8 – continued from previous page				
Strategy	Final Accuracy			
Disgust (Ablated)	88.81%			
Wait	88.81%			
CoT	86.57%			

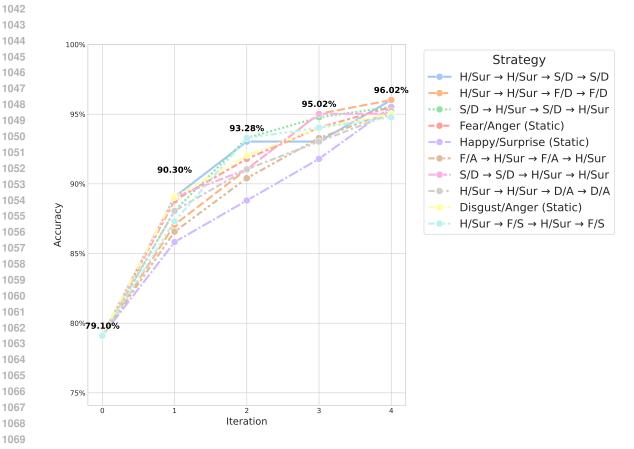


Figure 4: The 10 Best Performing Emotion Patterns using *HEART* on OlympiadBench Math with Gemini 2.5 Flash.

1081	Strategy Name	Final Accuracy
1082	$hsur \rightarrow sd \rightarrow hsur \rightarrow sd$	100.00%
1083	sa→hsur→sa→hsur	100.00%
1084	Sadness (Ablated)	100.00%
1085	hsur→hsur→hsur	100.00%
1086	$hsur \rightarrow fs \rightarrow hsur \rightarrow fs$	100.00%
1087	da→da→hsur→hsur	100.00%
1088	$sd \rightarrow sd \rightarrow hsur \rightarrow hsur$	100.00%
1089	hsur→sa→hsur→sa	100.00%
1090	Disgust (Ablated)	100.00%
1091	Disgust	100.00%
1092	Sadness	100.00%
1093	$fa \rightarrow fa \rightarrow fa \rightarrow fa$	100.00%
	$hsur \rightarrow hsur \rightarrow fs \rightarrow fs$	100.00%
1094	hsur→hsur→sd→sd	100.00%
1095	Happy (Ablated)	100.00%
1096	hsur→hsur→da→da	100.00%
1097	da→da→da	100.00%
1098	$sd \rightarrow sd \rightarrow sd \rightarrow sd$	100.00%
1099	Self Reflection (entire collection)	100.00%
1100	$fd \rightarrow fd \rightarrow hsu \rightarrow_h sur$	100.00%
1101	hsur→hsur→fd→fd	100.00%
1102	$fd \rightarrow fd \rightarrow fd \rightarrow fd$	100.00%
1103	sd→hsur→sd→hsur	100.00%
1104	fd→hsur→fd→hsur	100.00%
1105	fa→fa→hsur→hsur	100.00%
1106	Anger	100.00%
	hsur→fa→hsur→fa	100.00%
1107	Self Reflection ID# 8	100.00%
1108	$fa \rightarrow hsur \rightarrow fa \rightarrow hsur$	99.50%
1109	$fs \rightarrow fs \rightarrow fs \rightarrow fs$	99.25%
1110	Self Reflection ID# 2	99.25%
1111	da→hsur→da→hsur	99.25%
1112	Self Reflection ID# 6	99.25%
1113	Self Reflection ID# 1	99.25%
1114	Anger (Ablated)	98.97%
1115	Fear (Ablated)	98.51%
1116	Self Reflection ID# 3	98.51%
1117	Surprise (Ablatad)	98.51%
1118	Surprise (Ablated)	98.51%
1119	Self Reflection ID# 4	98.51%
	Self Reflection ID# 10 Fear	98.51%
1120		97.76%
1121	Happy Self Reflection ID# 9	97.76%
1122		97.76%
1123	Self Reflection ID# 7 Self Reflection ID# 5	97.01%
1124	Wait	97.01% 94.78%
1125	CoT	93.28%
1126	CUI	73.40 /0

Table 9: OlympiadBench Math Performance using *HEART* with Deepseek-R1 on the validation set (S1).

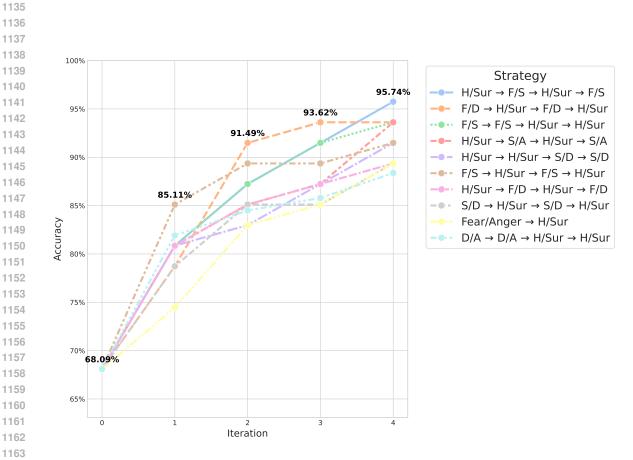


Figure 5: Gemini 2.5 Flash Accuracy per Iteration on OlympiadBench Physics Open Ended Problems using HEART.

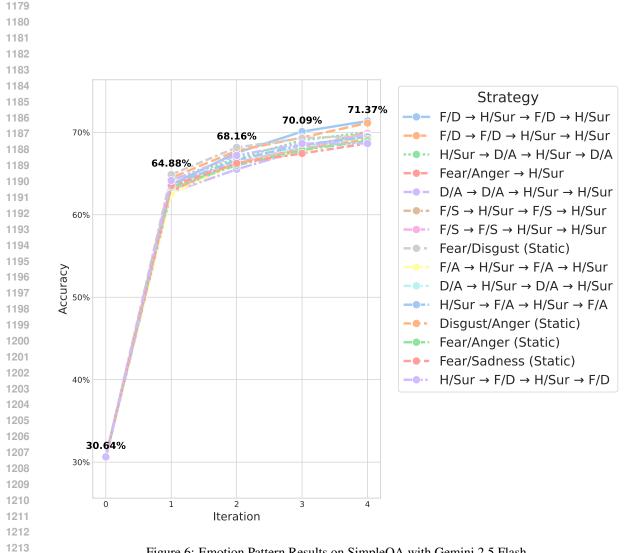


Figure 6: Emotion Pattern Results on SimpleQA with Gemini 2.5 Flash.