

CROP: TOKEN-EFFICIENT REASONING IN LARGE LANGUAGE MODELS VIA REGULARIZED PROMPT OPTIMIZATION

Deep Shah^{1,†,*} Sanket Badhe^{1,*} Nehal Kathrotia^{1,*} Priyanka Tiwari^{2,*}
¹Google LLC ²Purdue University

ABSTRACT

Large Language Models utilizing reasoning techniques improve task performance but incur significant latency and token costs due to verbose generation. Existing automatic prompt optimization (APO) frameworks target task accuracy exclusively at the expense of generating long reasoning traces. We propose Cost-Regularized Optimization of Prompts (CROP), an APO method that introduces regularization on response length by generating textual feedback in addition to standard accuracy feedback. This forces the optimization process to produce prompts that elicit concise responses containing only critical information and reasoning. We evaluate our approach on complex reasoning datasets, specifically GSM8K, LogiQA and BIG-Bench Hard. We achieved an 80.6% reduction in token consumption while maintaining competitive accuracy, seeing only a nominal decline in performance. This presents a pragmatic solution for deploying token-efficient and cost-effective agentic AI systems in production pipelines.

1 INTRODUCTION

Large Language Models (LLMs) have achieved unprecedented state-of-the-art capabilities in complex reasoning tasks primarily through step-by-step generation strategies like Chain of Thought Wei et al. (2022); Kojima et al. (2022); Wang et al. (2023b). By prompting models to decompose intricate problems into intermediate logical steps prior to outputting a final answer, researchers have significantly improved performance across mathematical, logical, and commonsense reasoning benchmarks Nye et al. (2021); Zelikman et al. (2022); Cobbe et al. (2021). Reasoning capabilities are an indispensable requirement for deploying foundation models across diverse and critical domains, including financial analysis, algorithmic ranking, medical diagnosis, and legal reasoning Wu et al. (2023); Cui et al. (2023); Badhe et al. (2026); Shah et al. (2026). This paradigm encourages models to mimic structured human deliberation Yao et al. (2023); Besta et al. (2024). Consequently, step-by-step rationale generation has become the standard operational mode for extracting high-fidelity reasoning from large-scale foundation models Lightman et al. (2024).

Despite these performance gains, the reliance on verbose intermediate reasoning introduces a critical deployment hurdle for real-world production environments. The generation of thousands of intermediate tokens incurs massive inference latency and exorbitant financial costs Pope et al. (2022); Aminabadi et al. (2022); Kwon et al. (2023). Recent empirical analyses reveal that modern reasoning models exhibit substantial verbosity compensation, often producing excessively lengthy explanations that provide little additional logical value Patel et al. (2024); Xu et al. (2025b). This extensive token consumption thus presents a severe liability that restricts the practical viability of deploying complex reasoning pipelines Renze & Guven (2024); Li et al. (2024); Jiang et al. (2023).

To optimize these systems, the field has increasingly adopted automated prompt engineering frameworks that leverage LLMs as meta-optimizers Zhou et al. (2023); Yang et al. (2024); Pryzant et al. (2023). Tools such as DSPy remove the need for manual human trial and error by refining prompts through systematic search Khattab et al. (2024). Recent frameworks extend this paradigm by automating the optimization workflow, utilizing meta-learning and combinatorial search to systematically refine instructions and prompting strategies without human intervention Spiess et al. (2025); Murthy et al. (2025); Xu et al. (2025a). However, these existing tools suffer from a critical flaw in that they optimize almost exclusively for task accuracy, often at the expense of generating excessively verbose intermediate reasoning traces. By continuously adding edge-case instructions and logic corrections to the prompt, the optimizer inadvertently maximizes the verbosity tax, producing bloated system prompts that force the target model to generate even longer and more exhaustive reasoning traces.

This inherent tendency toward verbosity exposes a fundamental missing mathematical constraint in the current paradigm of textual optimization. In standard machine learning optimization, loss functions routinely incorporate regularization techniques (such as L1 or L2 penalties) to constrain model weights and prevent overly complex architectures from overfitting Goodfellow et al. (2016); Bishop & Nasrabadi (2006). Yet, textual prompt optimization lacks a true analog to a length penalty.

*Equal contribution. †Correspondence to: shahdeep@google.com

To solve this structural deficiency, we propose CROP, an applied methodology that introduces an output-biased textual regularizer directly into the prompt optimization loop. CROP builds upon the multi-objective capabilities of textual differentiation to treat output token count as a first-class optimization constraint alongside task accuracy. By computing a continuous length-penalty gradient that functions independently of the task loss, the regularizer explicitly penalizes the target model whenever its output exceeds a minimal threshold. This feedback is then aggregated with traditional accuracy gradients using textual summation. Consequently, CROP forces the meta-optimizer to rewrite system instructions that explicitly constrain verbosity without sacrificing the underlying logical steps necessary for correct inference.

Through this regularized feedback mechanism, CROP successfully forces the discovery of concise reasoning paths. Our empirical evaluations demonstrate that this approach effectively shrinks the output distribution, minimizing both output length and overall inference latency. We evaluate CROP on complex reasoning benchmarks, specifically GSM8K Cobbe et al. (2021), LogiQA Liu et al. (2020), and BIG-Bench Hard Suzgun et al. (2023); Srivastava et al. (2023). The results show that our cost-aware optimization method yields a dramatic reduction in output token consumption by 80.6% reduction in token consumption while maintaining competitive accuracy. By establishing a continuous textual penalty for verbosity, CROP provides a pragmatic and highly effective solution for deploying token-efficient compound AI systems in strict production pipelines.

1.1 MAIN CONTRIBUTIONS

In summary, our primary contributions are as follows:

- **Cost-Aware Prompt Optimization (CROP):** We propose CROP, a novel token-efficient automatic prompt optimization framework. By introducing a response adaptive textual length penalty, CROP successfully decouples reasoning accuracy from generation verbosity during the prompt discovery phase.
- **Massive Inference Efficiency:** We empirically demonstrate that CROP reduces output token consumption by up to 80.6% on complex reasoning benchmarks (GSM8K, LogiQA and BIG-Bench Hard) while strictly preserving the task accuracy of unconstrained Chain-of-Thought prompting.
- **Autonomous Discovery of Symbolic Reasoning:** We show that our regularized textual gradients automatically force the target LLM to adopt highly compressed, symbolic reasoning structures. This emergent behavior conceptually mirrors human-engineered strategies (e.g., “Chain-of-Draft”) without requiring manual prompt design.

2 RELATED WORK

2.1 AUTOMATIC PROMPT OPTIMIZATION

The optimization of natural language prompts has evolved from discrete vocabulary search Shin et al. (2020); Deng et al. (2022); Prasad et al. (2023); Guo et al. (2023); Chen et al. (2024); Badhe & Shah (2026) to utilizing large language models as meta-optimizers Zhou et al. (2023); Yang et al. (2024); Pryzant et al. (2023); Wang et al. (2023a). Frameworks such as DSPy and TextGrad Khattab et al. (2024); Yuksekgonul et al. (2024) have formalized this into comprehensive computation graphs, enabling automated textual backpropagation capable of optimizing user-defined metrics and rewards. However, when applied to complex reasoning tasks, textual optimization exposes a fundamental generative bias: because evaluators diagnose logic errors via natural language critiques, they inherently resolve failures by prescribing additive logic and expanded instructions. This accumulation of instructional leads to unchecked increase of reasoning leading to verbosity. Recognizing this computational bottleneck, recent frameworks such as Promptmatix and AutoPDL Pan et al. (2025); Murthy et al. (2025) have introduced cost-aware objectives to balance efficiency with performance, alongside discrete evolutionary algorithms like CAPO and ParetoPrompt Wen (2025); Authors (2025); Zehle et al. (2025). Building upon this critical shift, CROP introduces a continuous, dual-objective textual optimization landscape. By applying an unconditional, output-biased regularizer alongside task gradients, CROP exerts continuous downward pressure on token bloat, forcing the meta-optimizer to systematically balance logical correctness with generative brevity in a way that purely additive textual gradients cannot.

2.2 INFERENCE-TIME SCALING AND SELF-CORRECTION

Alternative approaches improve reasoning by scaling compute during inference. Methodologies ranging from Self-Consistency Wang et al. (2023b) to iterative self-correction frameworks like Reflexion and Self-Refine Shinn et al. (2023); Madaan et al. (2023) allow models to critique and repeatedly refine their intermediate outputs. While yielding substantial accuracy gains, these pipelines inherently multiply the token consumption and latency of the system with every additional trial. In contrast to these methods that severely compound inference overhead, CROP shifts the iterative refinement burden entirely to an offline prompt optimization phase, ultimately producing a single, highly efficient system prompt that elicits compressed reasoning trajectories in a single forward pass.

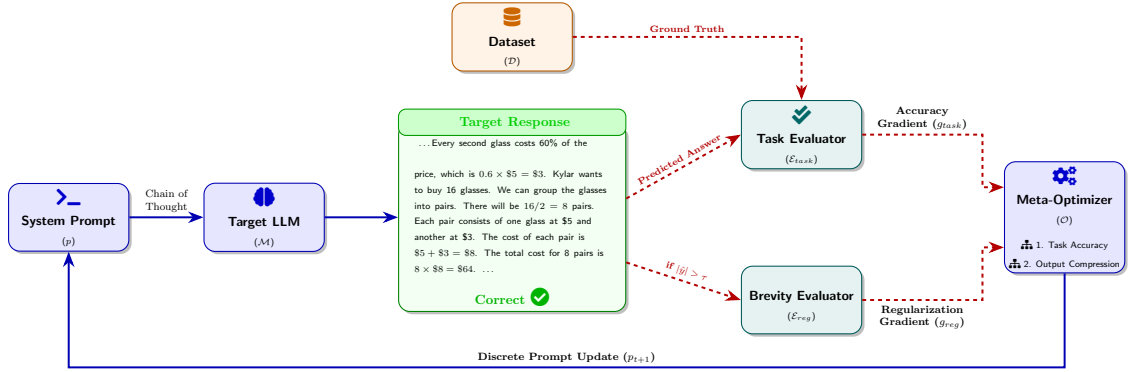


Figure 1: **Overview of CROP optimization pipeline.** A meta-optimizer updates the system prompt by aggregating textual gradients for both task accuracy (g_{task}) and response brevity (g_{reg}), systematically enforcing token-efficient reasoning.

2.3 COST-AWARE SYSTEMS AND VERBOSITY CONTROL

To mitigate computational overhead Pope et al. (2022); Kwon et al. (2023), researchers have developed system-level routing protocols Chen et al. (2023); Ong et al. (2024); Aggarwal et al. (2024); Jian et al. (2023), and inference-time early exiting or truncation strategies Mao et al. (2025); Jiang et al. (2023); Wu et al. (2025). Furthermore, recent literature focuses on explicitly compressing reasoning traces via continuous latent projections Hao et al. (2024); Shen et al. (2025); Wei et al. (2025), conditioned textual pruning Kang et al. (2024); Xia (2025); Wang (2025). However, these approaches introduce severe operational trade-offs. Methods relying on continuous latent projections or explicit fine-tuning are computationally prohibitive and alter foundation weights. Conversely, structured static prompting remains notoriously brittle, while dynamic token-pruning and step-skipping often require specialized decoding mechanisms or architectural interventions. Unlike approaches that demand expensive weight updates, custom decoding algorithms, or fragile hardcoded constraints, CROP operates entirely in the natural language optimization space. It dynamically modulates verbosity at the prompt level, systematically discovering optimal, token-efficient reasoning pathways through regularized textual feedback.

3 METHODOLOGY

In this section, we formalize **CROP**, a framework that extends the automatic prompt optimization (APO) paradigm to resource-constrained environments. We first briefly review the preliminary formulation of textual differentiation, and subsequently detail our novel textual regularization objective.

3.1 PRELIMINARIES: TEXTUAL DIFFERENTIATION

Let \mathcal{V}^* denote the space of all possible natural language strings. In a standard LLM-based computation graph, we define a system prompt $p \in \mathcal{V}^*$ as a learnable parameter, an input query $x \in \mathcal{V}^*$, and a target model \mathcal{M} . The forward pass generates an output $\hat{y} = \mathcal{M}(x; p)$.

Given a ground-truth label y , a textual loss function \mathcal{L}_{task} evaluates the prediction. Instead of yielding a scalar gradient, textual backpropagation queries a secondary evaluator LLM \mathcal{E} , to generate a *textual gradient* $g_{task} \in \mathcal{V}^*$. This gradient is a natural language critique detailing the logical flaws in \hat{y} and suggesting semantic updates to p :

$$g_{task} = \mathcal{E}_{task}(p, x, \hat{y}, y) \quad (1)$$

An optimizer LLM, \mathcal{O} , then acts as the parameter update function, synthesizing the batch of gradients to perform a discrete step in the prompt space: $p_{t+1} = \mathcal{O}(p_t, \{g_{task}^{(i)}\}_{i=1}^B)$, where B is the batch size.

3.2 VERBOSE REASONING

While \mathcal{L}_{task} effectively maximizes reasoning accuracy, it is agnostic to computational efficiency. Consequently, the meta-optimizer \mathcal{O} inherently gravitates toward system prompts that trigger useful but overly verbose Chain-of-Thought (CoT) bloat in \hat{y} . In production systems, the length of \hat{y} directly influences the inference latency and compute requirements.

Algorithm 1 CROP: Cost-Regularized Optimization of Prompts**Require:** Initial system prompt p_0 , Target LLM \mathcal{M} , Meta-Optimizer LLM \mathcal{O} **Require:** Task Evaluator $\mathcal{E}_{\text{task}}$, Brevity Evaluator \mathcal{E}_{reg} **Require:** Dataset \mathcal{D} , Batch size B , Max iterations T

```

1:  $p \leftarrow p_0$ 
2: for  $t = 1$  to  $T$  do
3:   Sample batch  $\{(x_i, y_i)\}_{i=1}^B \sim \mathcal{D}$ 
4:    $G_{\text{batch}} \leftarrow \emptyset$  ▷ Initialize empty batch gradient list
5:   for  $i = 1$  to  $B$  do
6:      $\hat{y}_i \leftarrow \mathcal{M}(x_i; p)$  ▷ Forward Pass
7:      $g_{\text{task}}^{(i)} \leftarrow \mathcal{E}_{\text{task}}(p, x_i, \hat{y}_i, y_i)$  ▷ Compute Task Accuracy Gradient
8:      $g_{\text{reg}}^{(i)} \leftarrow \mathcal{E}_{\text{reg}}(p, \hat{y}_i)$  ▷ Compute Textual Regularization Gradient
9:      $g_{\text{total}}^{(i)} \leftarrow g_{\text{task}}^{(i)} \oplus "\backslash\text{n}" \oplus g_{\text{reg}}^{(i)}$  ▷ Aggregate Gradients via Concatenation
10:     $G_{\text{batch}} \leftarrow G_{\text{batch}} \cup \{g_{\text{total}}^{(i)}\}$ 
11:  end for
12:   $\text{Accuracy}_t \leftarrow C_{\text{correct}}/B$ 
13:   $\text{AvgLength}_t \leftarrow L_{\text{total}}/B$ 
14:   $S_t \leftarrow \text{Accuracy}_t - \lambda \cdot L_{\text{norm}}(t)$  ▷ Evaluate trade-off score
15:  if  $S_t > S_{\text{best}}$  then ▷ Update best prompt if score improves
16:     $S_{\text{best}} \leftarrow S_t$ 
17:     $p^* \leftarrow p$ 
18:  end if
19:   $p \leftarrow \mathcal{O}(p, G_{\text{batch}})$  ▷ Meta-Optimizer computes discrete prompt update
20: end for
21: return  $p^*$  ▷ Optimized token-efficient prompt

```

3.3 TEXTUAL REGULARIZATION

CROP applies \mathcal{L}_{reg} across the entire batch, asserting a constant downward pressure on token consumption. We define a secondary evaluator prompt, \mathcal{E}_{reg} , which receives the generated output and computes a textual brevity gradient:

$$g_{\text{reg}} = \mathcal{E}_{\text{reg}}(p, \hat{y}) \quad (2)$$

The resulting gradient g_{reg} explicitly instructs the optimizer to inject brevity constraints into p . The regularization gradient is computed only when $\|\hat{y}\|$ exceeds a predefined threshold, preventing the penalization of already token-efficient reasoning. To maintain optimization stability and prevent destructively large steps in the prompt space—which frequently induce catastrophic forgetting of the underlying logic— \mathcal{E}_{reg} is constrained to enforce only a gradual, incremental reduction in output length per iteration. Furthermore, the evaluator is explicitly directed to target peripheral narrative elements and conversational filler, strictly preserving the tokens associated with critical computational steps.

3.4 GRADIENT AGGREGATION AND MULTI-OBJECTIVE OPTIMIZE

To perform the parameter update, CROP treats prompt optimization as a formal multi-objective problem. We leverage the concatenation properties of textual autograd engines to compute the total gradient. For a single instance, the combined textual gradient g_{total} is defined via the string concatenation operator \oplus :

$$g_{\text{total}} = g_{\text{task}} \oplus "\backslash\text{n}" \oplus g_{\text{reg}} \quad (3)$$

During the batched optimization step, the meta-optimizer \mathcal{O} receives the aggregated sets of both task gradients and regularization gradients and performs the update step $p_{t+1} = \mathcal{O}(p_t, \mathbf{g}_{\text{total}})$.

Rather than naively returning the prompt from the final iteration, CROP employs a structured selection criterion to identify the optimal prompt p^* discovered during the optimization trajectory. We define a composite score S for each candidate prompt:

$$S(p) = \text{Accuracy}(p) - \lambda \cdot L_{\text{norm}}(p) \quad (4)$$

where $L_{\text{norm}}(p)$ represents the generated token count normalized by the average COT output token length, and λ is a tunable regularization coefficient. The magnitude of λ dictates the explicit trade-off between task performance and token consumption. We selected λ through an empirical grid search on the validation set over the range of 0.01 to 0.5. We identified $\lambda = 0.2$ as the optimal configuration, as it provided the maximum reduction in reasoning verbosity without inducing a statistically significant degradation in task accuracy.

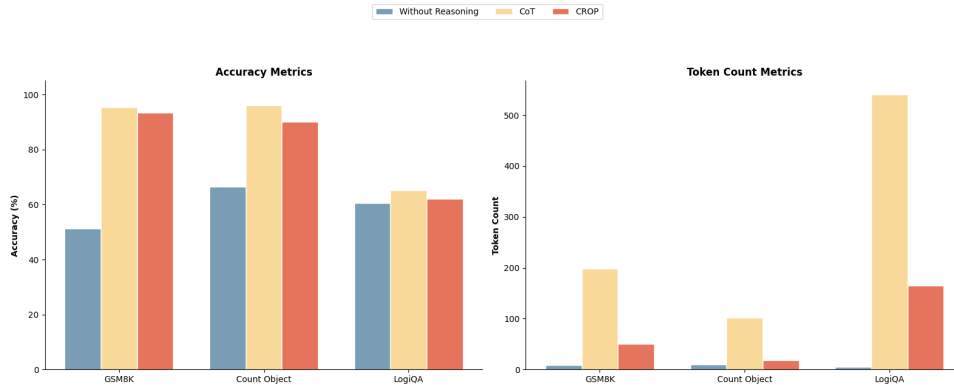


Figure 2: Performance and efficiency comparison across three datasets using the Gemini 2.0 Flash model.

4 EXPERIMENT SETUP

4.1 MODEL SETUP

We assign specific Large Language Models to the distinct roles of evaluation, optimization, and inference. We utilize Gemini 2.0 Flash as the dual-evaluator ($\mathcal{E}_{\text{task}}$ and \mathcal{E}_{reg}) to compute task accuracy critiques and regularized verbosity gradients. For the parameter update step, we deploy Gemini 3.1 Pro (with Thinking enabled) as the meta-optimizer (\mathcal{O}). A high-capacity model is critical here, as synthesizing multi-objective, conflicting textual gradients requires an extended context window and advanced deductive reasoning (ablated in Section 5).

For forward-pass inference (\mathcal{M}), we evaluated two models: Gemini 2.0 Flash and Qwen 2.5 7B.

4.2 BASELINES

We benchmark CROP against zero-shot prompting, chain-of-thoughts, concise induced prompting and automatic prompt optimization method TextGrad. We intentionally omit static length-constrained prompts, which are brittle and degrade reasoning; furthermore, the model was unable to follow explicit token-limit instructions and frequently exceeded the given limits.

- **Standard Zero-Shot (Direct Prompting):** Prohibits intermediate reasoning steps to output the final answer directly. This establishes our empirical lower bound for token consumption, though typically at the severe expense of task accuracy.
- **Chain-of-Thought (CoT):** Appends a standard step-by-step reasoning directive. This serves as our upper bound for baseline accuracy, simultaneously representing the widely used but verbosity prompting technique.
- **Conciseness inducing prompts:** We used different prompts which encourages concise reasoning via prompting Lee et al. (2025). (See prompts in Appendix A.2)
- **TextGrad:** An automatic prompt optimization framework that utilizes textual gradients to maximize task accuracy without penalizing verbose reasoning during the optimization loop.

4.3 DATASETS

To evaluate CROP across diverse cognitive domains, we selected three challenging reasoning benchmarks. Because these tasks inherently necessitate intermediate logic, they intrinsically incentivize unconstrained LLMs to generate highly verbose Chain-of-Thought (CoT) traces, providing an optimal testbed for our compression framework. To prevent data leakage, we utilize the official train, validation, and test splits provided by the dataset authors wherever available. For the BBH Object Counting task, which lacks standard splits, we manually partition the data into a 25/25/50 split for training, validation, and testing, respectively.

- **GSM8K Cobbe et al. (2021):** A benchmark of linguistically diverse grade-school math word problems. This dataset requires accurate multi-step arithmetic reasoning, which frequently induces unconstrained LLMs to produce excessively long textual explanations for basic algebraic steps.
- **Object Counting (BIG-Bench Hard) Suzgun et al. (2023):** An algorithmic reasoning task requiring the model to determine the total count of a specific semantic class from a natural-language collection

of items. This tests symbolic tracking, where standard CoT prompts typically trigger exhaustive, token-heavy item enumeration.

- **LogiQA Liu et al. (2020):** A challenging machine reading comprehension dataset sourced from expert logical examinations. It requires complex deductive inference (e.g., categorical, conditional, and disjunctive reasoning), assessing CROP’s ability to maintain the structural integrity of formal logical arguments under aggressive token compression.

Method	Core Instruction
No Reasoning	Output ONLY the final answer without reasoning or thinking.
Chain-of-Thought	Think step by step.
CROP (Ours)	[...] concisely, showing the combined essential calculations leading directly to the final numerical answer. Omit intermediate labels. Focus on chained calculations.

Table 1: Comparison of core system instructions across prompting baselines for the GSM8K dataset. Standard prefix formulations and output formatting constraints have been omitted to highlight the distinct reasoning directives.

Qualitative comparison of generated reasoning traces on the GSM8K dataset.	
Input Query: Kylar went to the store to buy glasses for his new apartment. One glass costs \$5, but every second glass costs only 60% of the price. Kylar wants to buy 16 glasses. How much does he need to pay for them?	
1. Direct Prompting (No Reasoning) 68	✗ Incorrect
2. Chain-of-Thought (Standard Verbosity) Let the cost of one glass be \$5. Every second glass costs 60% of the price, which is $0.6 * \$5 = \3 . Kylar wants to buy 16 glasses. We can group the glasses into pairs. There will be $16 / 2 = 8$ pairs. Each pair consists of one glass at \$5 and another at \$3. The cost of each pair is $\$5 + \$3 = \$8$. The total cost for 8 pairs is $8 * \$8 = \64 .	✓ Correct (128 tokens)
3. CROP (Ours Compressed Reasoning) $8 * 5 + 8 * (5 * 0.6) = 40 + 8 * 3 = 40 + 24 = 64$	✓ Correct (44 tokens)

Table 2: Reasoning comparison on the GSM8K dataset.

5 RESULTS AND ANALYSIS

We evaluate CROP’s capacity to eliminate the computational redundancy inherent in unconstrained Chain-of-Thought (CoT) reasoning. By integrating a continuous brevity gradient into the optimization loop, CROP explicitly optimizes for reasoning density alongside task performance. As detailed in Table 3, our framework compresses output token consumption by up to 80.6% without degrading exact-match accuracy. These results establish that high-fidelity inference can be sustained at a fraction of the standard computational cost.

5.1 SIGNIFICANT REDUCTIONS IN TOKEN CONSUMPTION

CROP consistently achieved state-of-the-art token efficiency while preserving the logical integrity required for complex problem-solving. On mathematical arithmetic (GSM8K), algorithmic tasks (Object Counting), and LogiQA, the framework reduced output token consumption by 74.8%, 80.6%, and 66.5% respectively while retaining reasoning leading to comparable accuracy for Gemini 2.0 Flash. Similar token-efficiency was observed for Qwen 2.5 7B as well. This starkly contrasts with the Direct Prompting (Without Reasoning) baseline, which collapses entirely on these tasks. This divergence confirms that while intermediate computation is strictly necessary for accuracy, the narrative *verbosity* of that computation is highly compressible via regularized feedback.

Table 3: Model Performance Across Prompting Methods. Cells display Accuracy (%) with Average Token Count in parentheses.

Model	Dataset	CoT	w/o Reasoning	CROP	TextGrad	Be Concise	Only Number	Use Abbrev
Gemini 2 Flash	GSM8K	95.3 (198.4)	51.2 (9.2)	93.4 (50.0)	95.5 (185.5)	96.2 (120.8)	91.4 (113.0)	95.1 (145.7)
	Count Object	96.0 (101.2)	66.4 (9.4)	94.8 (19.6)	90.5 (24.5)	96.0 (59.0)	92.5 (55.0)	95.0 (92.0)
	LogiQA	65.2 (540.5)	60.4 (5.1)	64.2 (181.0)	66.4 (663.0)	65.1 (301.0)	61.9 (423.0)	65.3 (422.0)
Qwen2.5 7B	GSM8K	90.2 (123.4)	54.9 (3.4)	87.8 (21.3)	89.8 (101.2)	90.4 (82.3)	85.1 (31.1)	88.2 (72.4)
	Count Object	92.1 (56.3)	67.8 (3.8)	90.1 (12.1)	88.8 (32.7)	92.7 (46.6)	81.5 (16.2)	90.2 (39.7)
	LogiQA	61.2 (287.0)	54.0 (3.4)	59.3 (84.2)	60.1 (212.3)	61.1 (152.3)	50.8 (32.8)	54.7 (212.0)

5.2 EMERGENT SYMBOLIC REASONING

A qualitative examination of the optimized prompts (Table 1) and generated outputs (Table 2) reveals that CROP discovered prompts fundamentally shifts how the target model structures its intermediate logic. Rather than generating conversational heuristics (e.g., “First, let us observe that...”), the regularized prompts elicit a dense, symbolic syntax. On GSM8K, the model spontaneously adopts a mathematical shorthand analogous to the recently proposed “Chain-of-Draft” Xu et al. (2025b) methodology. However, whereas prior works rely on manual human engineering to elicit this behavior, CROP discovers this optimal reasoning syntax entirely autonomously.

Meta-Optimizer Model	Discovered System Prompt (Core Instructions)
Gemini 2.0 Flash (<i>Lower Capacity</i>)	Solve the problem. Only provide the answer and do not provide any explanations.
Gemini 3.1 Pro (<i>High Capacity</i>)	Solve the problem with extreme conciseness. Do not use conversational filler, restate the question, or write full sentences. Simply provide the final sum.

Table 4: Comparison of optimized system prompts generated by different meta-optimizers for the BBH Object Counting task. While the lower-capacity model (Gemini 2.0 Flash) defaults to a blunt truncation directive, the high-capacity model (Gemini 3.1 Pro) successfully synthesizes multi-faceted constraints targeting specific vectors of verbosity.

Optimization Setting	Discovered System Prompt (Core Instructions)
Low Capacity Optimizer (Gemini 2.0 Flash, $B = 128$)	Solve the problem. Only provide the answer and do not provide any explanations.
Stochastic Optimization (Gemini 3.1 Pro, $B = 1$)	Solve the problem concisely. Go straight to the answer without restating the problem or adding redundant explanatory sentences.
Optimal Configuration (Gemini 3.1 Pro, $B = 128$)	Solve the problem with extreme conciseness. Do not use conversational filler, restate the question, or write full sentences. Simply provide the final sum.

Table 5: Qualitative ablation of discovered system prompts under different optimization constraints for the BBH dataset. Only the optimal configuration (high capacity, large batch) successfully balances extreme conciseness with the preservation of reasoning.

5.3 ROLE OF OPTIMIZER SCALE IN GENERALIZATION

The meta-optimizer (\mathcal{O}) is responsible for synthesizing batched textual accuracy and regularization gradients into a coherent, updated system prompt. In early iterations, deploying a lower-capacity model (Gemini 2.0 Flash) as the meta-optimizer resulted in a severe degradation of prompt quality. The optimizer consistently overfit to idiosyncratic errors from isolated batch examples rather than deriving generalizable instructional updates. Extensive tuning of the meta-optimizer’s system prompt failed to mitigate this limitation.

Ultimately, scaling the meta-optimizer to a frontier reasoning model—specifically Gemini 3.1 Pro with thinking capabilities enabled—resolved these issues and produced highly robust system prompts (See table 4). We hypothesize that synthesizing complex, multi-objective textual feedback requires both the extended context window and the advanced deductive reasoning capabilities inherent to this larger model.

5.4 GRADIENT STABILITY AND BATCH DYNAMICS

In standard deep learning, gradient variance is inversely proportional to batch size. We observe an analogous phenomenon in regularized textual differentiation. As demonstrated in Table 5, ablation experiments utilizing stochastic updates (batch size of 1) yield extreme variance in the textual gradients. When evaluating an isolated example, the meta-optimizer frequently overfits to the length penalty, prioritizing output compression over task accuracy and aggressively stripping away critical cognitive instructions. This results in prompt collapse: the generation of degenerate directives that entirely eliminate intermediate reasoning capacity. Conversely, scaling to a batch size of 128 stabilizes the multi-objective feedback. By providing the meta-optimizer with a statistically representative sample of both reasoning failures and excess verbosity, large batches prevent regularization-induced collapse and ensure stable, convergent trajectories during prompt optimization.

6 CONCLUSION

The deployment of Large Language Models in complex reasoning pipelines has historically been bottlenecked by the severe latency and financial costs associated with verbose Chain-of-Thought generation. In this work, we introduced CROP, a novel multi-objective automatic prompt optimization framework. By integrating a continuous textual regularization penalty directly into the prompt optimization loop, CROP forces the meta-optimizer to navigate the Pareto frontier of accuracy and token-consumption.

Our empirical evaluations demonstrate that intermediate reasoning is highly compressible. CROP successfully reduces output token consumption by up to 80.6% on rigorous reasoning benchmarks (GSM8K, BIG-Bench Hard, and LogiQA) while maintaining comparable accuracy compared to unconstrained CoT and automatic prompt optimization baselines. Furthermore, our analysis reveals that this dual-objective optimization automatically elicits highly efficient, symbolic reasoning structures, conceptually mirroring manual “Chain-of-Draft” prompting without requiring human intervention. Through systematic ablations, we also identified that successful multi-objective textual optimization requires frontier-level model capacity and large batch sizes to stabilize gradient variance and prevent the collapse of intermediate logic.

Ultimately, CROP provides a highly practical, deployment-ready methodology. By systematically excising computational redundancy at the prompt level rather than altering foundational weights, it enables the scalable and cost-effective integration of high-fidelity reasoning models into strict, real-world production and agentic environments.

LIMITATIONS

While CROP significantly reduces inference-time overhead, we acknowledge two primary limitations regarding its current operational scope.

One limitation of our current evaluation is that we do not explicitly consider input token length or cost. In modern production environments, prompt caching (or prefix caching) is a default feature that automatically stores and reuses the initial segments of instructions. Because our framework produces a single, static system prompt that remains identical across all user queries, it perfectly leverages this caching mechanism. By reusing this prior computational work, prefix caching reduces time-to-first-token latency by up to 80% and decreases input token costs by up to 90% OpenAI (2024). Consequently, the financial impact of a slightly longer optimized system prompt is effectively marginalized. This justifies our exclusive focus on constraining the verbosity of the generated output, which remains the primary driver of unmitigated inference costs.

Second, the optimization phase of CROP necessitates access to high-capacity reasoning models (e.g., Gemini 3.1 Pro with Thinking). We observe that such model is required to effectively synthesize and balance conflicting objectives between task accuracy and textual regularization. Consequently, the efficacy of using smaller, mid-tier models as meta-optimizers for this dual-objective task has not yet been established.

REFERENCES

- Aman Aggarwal et al. Automix: Automatically mixing language models. *arXiv preprint arXiv:2310.12963*, 2024.
- Reza Yazdani Aminabadi, Samyam Rajbhandari, Ammar Ahmad Awan, Cheng Li, Du Li, Elton Zheng, Olatunji Ruwase, Shaden Smith, Minjia Zhang, Jeff Rasley, et al. Deepspeed-inference: enabling efficient inference of transformer models at unprecedented scale. In *SC22: International Conference for High Performance Computing, Networking, Storage and Analysis*, pp. 1–15. IEEE, 2022.
- MOPO Authors. Mopo: Multi-objective prompt optimization for affective text generation. In *Proceedings of the 31st International Conference on Computational Linguistics*, 2025.

- Sanket Badhe and Deep Shah. Prompt-level distillation: A non-parametric alternative to model fine-tuning for efficient reasoning. *arXiv preprint arXiv:2602.21103*, 2026.
- Sanket Badhe, Deep Shah, and Nehal Kathrotia. Long-tail knowledge in large language models: Taxonomy, mechanisms, interventions and implications. *arXiv preprint arXiv:2602.16201*, 2026.
- Maciej Besta, Nils Blatter, Lucian Dragos, Jing Lin, Oksana Krause, et al. Graph of thoughts: Solving elaborate problems with large language models. In *International Conference on Learning Representations*, 2024.
- Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*. Springer, 2006.
- Jialin Chen et al. Evoprompt: Evolutionary prompt optimization for large language models. In *International Conference on Learning Representations*, 2024.
- Lingjiao Chen, Matei Zaharia, and James Zou. Frugalgpt: How to use large language models while reducing cost and improving performance. In *Advances in Neural Information Processing Systems*, 2023.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- Jiaxi Cui, Zongjian Li, Yang Yan, Bohua Chen, and Li Yuan. Chatlaw: Open-source legal large language model with integrated external knowledge bases. *arXiv preprint arXiv:2306.16092*, 2023.
- Mingkai Deng, Jianyu Wang, Cheng-Ping Hsieh, Yihan Wang, Han Guo, Tianmin Shu, Meng Song, Eric P Xing, and Zhiting Hu. Rlprompt: Optimizing discrete text prompts with reinforcement learning. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, 2022.
- Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.
- Qingyan Guo, Rui Wang, Junliang Guo, Bei Li, Kaitao Song, Xu Tan, Guihong Liu, Ke Jiang, Jian-Guang Yao, Min Zhang, et al. Connecting large language models with evolutionary algorithms yields powerful prompt optimizers. *arXiv preprint arXiv:2309.08532*, 2023.
- Shibo Hao et al. Training large language models to reason in a continuous latent space. *arXiv preprint arXiv:2412.00000*, 2024.
- Yuhang Jian et al. Ecoassistant: Using llm assistant more affordably and accurately. In *Advances in Neural Information Processing Systems*, 2023.
- Huiqiang Jiang, Qianhui Wu, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. Llmlingua: Compressing prompts for accelerated inference of large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, 2023.
- Yu Kang, Xianghui Sun, Liangyu Chen, and Wei Zou. C3ot: Generating shorter chain-of-thought without compromising effectiveness. *arXiv preprint arXiv:2412.11664*, 2024.
- Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Sri Vardhamanan, Saiful Haq, Ashutosh Sharma, Thomas T Joshi, Hanna Moazam, et al. Dspy: Compiling declarative language model calls into state-of-the-art pipelines. In *International Conference on Learning Representations*, 2024.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. In *Advances in Neural Information Processing Systems*, 2022.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E Gonzalez, Hao Zhang, and Jian Jin. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the 29th Symposium on Operating Systems Principles*, 2023.
- Ayeong Lee, Ethan Che, and Tianyi Peng. How well do llms compress their own chain-of-thought? a token complexity approach. *arXiv preprint arXiv:2503.01141*, 2025.
- Zongqian Li, Yixuan Su, and Nigel Collier. 500xcompressor: Generalized prompt compression for large language models. *arXiv preprint arXiv:2408.03094*, 2024.
- Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, et al. Let’s verify step by step. In *International Conference on Learning Representations*, 2024.
- Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. Logiqa: A challenge dataset for machine reading comprehension with logical reasoning. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence*, 2020.

- 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. In *Advances in Neural Information Processing Systems*, 2023.
- Minjia Mao, Bowen Yin, Yu Zhu, and Xiao Fang. Early stopping chain-of-thoughts in large language models. *arXiv preprint arXiv:2509.14004*, 2025.
- Rithesh Murthy, Ming Zhu, Liangwei Yang, Jieli Qiu, Juntao Tan, Shelby Heinecke, Silvio Savarese, Caiming Xiong, and Huan Wang. Promptomatix: An automatic prompt optimization framework for large language models. *arXiv preprint arXiv:2507.14241*, 2025.
- Maxwell Nye, Anders Johan Andreassen, Guy Gur-Ari, Aitor Lewkowycz Henry, Vikash Mansinghka, Vedant Mikulik, and Nathanael Schärli. Show your work: Scratchpads for intermediate computation with language models. In *Advances in Neural Information Processing Systems*, 2021.
- Isaac Ong et al. Routellm: Learning to route llms with preference data. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, 2024.
- OpenAI. Prompt caching, 2024. URL <https://developers.openai.com/api/docs/guides/prompt-caching>. Accessed: 2026-03-17.
- Xiaoying Pan et al. Autopdl: Automatic prompt optimization for llm agents. *arXiv preprint arXiv:2504.04365*, 2025.
- Rohan Patel, Kshitij Sharma, and Anil Gupta. Verbosity \neq veracity: Demystify verbosity compensation behavior of large language models. *arXiv preprint arXiv:2411.07858*, 2024.
- Reiner Pope, Sholto Douglas, Aakanksha Chowdhery, Jacob Devlin, James Bradbury, et al. Efficiently scaling transformer inference. *arXiv preprint arXiv:2211.05102*, 2022.
- Archiki Prasad, Peter Hase, Xiang Zhou, and Mohit Bansal. Grips: Gradient-free, edit-based instruction search for prompting large language models. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, 2023.
- Reid Pryzant, Dan iterative, et al. Automatic prompt optimization with” gradient descent” and beam search. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, 2023.
- Matthew Renze and Erhan Guven. The benefits of a concise chain of thought on problem-solving in large language models. *arXiv preprint arXiv:2401.05618*, 2024.
- Deep Shah, Sanket Badhe, and Nehal Kathrotia. Taxonomy of the retrieval system framework: Pitfalls and paradigms. *arXiv preprint arXiv:2601.20131*, 2026.
- Yikang Shen et al. Codi: Compressing chain-of-thought into continuous space via self-distillation. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, 2025.
- Taylor Shin, Yasaman Razeghi, Robert L Logan IV, Eric Wallace, and Sameer Singh. Autoprompt: Eliciting knowledge from language models with automatically generated prompts. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*, 2020.
- Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. In *Advances in Neural Information Processing Systems*, 2023.
- Claudio Spiess, Mandana Vaziri, Louis Mandel, and Martin Hirzel. Autopdl: Automatic prompt optimization for llm agents. *arXiv preprint arXiv:2504.04365*, 2025.
- Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *Transactions on Machine Learning Research*, 2023.
- Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, et al. Challenging big-bench tasks and whether chain-of-thought can solve them. In *Findings of the Association for Computational Linguistics: ACL 2023*, 2023.
- Chenglin Wang et al. Instructzero: Efficient instruction optimization for black-box large language models. In *Advances in Neural Information Processing Systems*, 2023a.
- et al. Wang. Compressing chain-of-thought in llms via step entropy. *arXiv preprint arXiv:2508.03346*, 2025.

- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. In *International Conference on Learning Representations*, 2023b.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In *Advances in Neural Information Processing Systems*, 2022.
- Xilin Wei, Xiaoran Liu, Yuhang Zang, Xiaoyi Dong, Yuhang Cao, Jiaqi Wang, Xipeng Qiu, and Dahua Lin. Sim-cot: Supervised implicit chain-of-thought. *arXiv preprint arXiv:2509.20317*, 2025.
- et al. Wen. Pareto prompt optimization. In *International Conference on Learning Representations*, 2025.
- Shijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabravolski, Mark Dredze, Sebastian Gehrmann, Prabhanjan Kambadur, David Rosenberg, and Gideon Mann. Bloomberggpt: A large language model for finance. *arXiv preprint arXiv:2303.17564*, 2023.
- Zhaofeng Wu et al. An empirical study on prompt compression for large language models. *arXiv preprint arXiv:2505.00019*, 2025.
- et al. Xia. Tokenskip: Controllable chain-of-thought compression in llms. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, 2025.
- Guowei Xu, Mert Yuksekogun, Carlos Guestrin, and James Zou. metatextgrad: Automatically optimizing language model optimizers. *arXiv preprint arXiv:2505.18524*, 2025a.
- Silei Xu et al. Chain of draft: Thinking faster by writing less. *arXiv preprint arXiv:2502.18600*, 2025b.
- Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V Le, Denny Zhou, and Xinyun Chen. Large language models as optimizers. In *International Conference on Learning Representations*, 2024.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. In *Advances in Neural Information Processing Systems*, 2023.
- Mert Yuksekogun, Federico Bianchi, Daniil Boiko, et al. Textgrad: Automatic "differentiation" via text. *arXiv preprint arXiv:2406.07496*, 2024.
- Tom Zehle, Moritz Schlager, Timo Heiß, and Matthias Feurer. Capo: Cost-aware prompt optimization. *arXiv preprint arXiv:2504.16005*, 2025.
- Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah D Goodman. Star: Bootstrapping reasoning with reasoning. In *Advances in Neural Information Processing Systems*, 2022.
- Yongchao Zhou, Andrei Ioan Muresanu, Zihao Zi, Keiran Tucker, et al. Large language models are human-level prompt engineers. In *International Conference on Learning Representations*, 2023.

A APPENDIX

A.1 GENERATIVE AI USAGE

All study design, literature review, synthesis, and writing were conducted by the authors. Generative AI tools (Gemini) were used only for grammar checking and proofreading during the final polishing of the manuscript. No generative AI system was used to generate content, interpret prior work, or draw conclusions. The authors reviewed and approved all final text and remain fully responsible for the content of the paper.

A.2 BASELINE PROMPTS

CoT

Think Step by Step.

Be concise

Think Step by Step. **Be Concise.**

Only Num/Eq

Think Step by Step. **Only use numbers or equations.**

AbbrevWords

Think Step by Step. **Abbreviate words as much as possible.**

UseOnlyNTokens

Think Step by Step. **Do not use more than N words.**

A.3 INITIAL PROMPTS

GSM8K

You will answer a mathematical reasoning question.
The last line of your response should be of the following format: 'Answer: \$VALUE' where VALUE is a numerical value.

BBH Object counting

Solve the problem
The last line of your response should be of the following format: 'Answer: VALUE' where VALUE is a numerical value.

LogiQA

You will answer a logical reasoning question.
The last line of your response should be of the following format: 'Answer: VALUE' where VALUE is the zero-based numerical index of the correct answer.

A.4 TEXTGRAD LEARNT PROMPTS

GSM8K

You will answer a mathematical reasoning question. First, explicitly and concisely state the final goal of the question.
Think step by step and be extremely concise in your reasoning. Avoid conversational filler and jump straight into the math.
Define your variables clearly and explicitly track units during your intermediate calculations to prevent conversion errors.
Carefully consider edge cases, hidden assumptions, overlapping sets, and ensure your final unit matches what the question asks for (e.g., dollars vs. cents). Double-check your arithmetic.
The very last line of your response must be exactly of the following format: 'Answer: VALUE' where VALUE is a pure integer.
Leave a blank line before this final line to completely isolate it from the rest of your explanation. Ensure there is exactly one space after the colon. Do not include any dollar signs (\$), commas, units, periods, or any other extraneous characters in VALUE. There should be absolutely no additional text, punctuation, or whitespace following the integer.

BBH Object counting

Solve the problem.
To ensure proper parsing, keep your response as concise as possible.
The very last line of your response must be exactly of the following format: 'Answer: VALUE' where VALUE is an integer.
Do not include any additional text, punctuation, or spaces after the numerical value.

LogiQA

You are an expert in logical reasoning. You will answer a logical reasoning question by selecting the correct option from the provided list.

Follow these systematic steps:

1. **Analyze the Query:** Carefully break down the context and the question. Identify the core logical structure, premises, definitions, and specific requirements. Pay special attention to negative constraints (e.g., "NOT", "cannot", "least likely") and the modality required (e.g., "must be true", "most weakens"). Explicitly state the ultimate goal of the question (e.g., "The goal is to find an alternative explanation that directly undermines the causal link").

2. **Formulate Criteria:** Establish clear, precise, and prioritized criteria that the correct answer must satisfy based on your logical analysis. Include negative criteria (what the correct option must *not* do) and distinguish between necessary and sufficient conditions.

3. **Evaluate Options:** Systematically and concisely analyze every option using its zero-based index (e.g., Option 0, Option 1, Option 2, Option 3). Use strong, decisive language to explain exactly how it meets or fails the formulated criteria.

- For "weaken/strengthen" questions, prioritize options that directly attack or support the core premise or mechanism over those providing indirect, tangential, or correlational impacts.

- For "must be true/inference" questions, ensure the deduction is a necessary consequence of the provided premises. Explicitly reject options that are merely "possible" or require external, unwarranted assumptions.

- Address potential nuances (e.g., correlation vs. causation, ambiguity, shifts in scope), clearly explaining why an incorrect option is logically flawed.

4. **Compare and Conclude:** Conduct a rigorous and direct comparative analysis of the most viable options. Be decisive. Explicitly contrast the strengths of the correct option against the specific weaknesses of the runners-up to summarize why the selected option is definitively the best or only correct choice.

5. **Final Answer:** Skip any section heading for this step. On the very last line of your response, provide the final answer in the exact format:

Answer: VALUE (where VALUE is the single zero-based integer index of the correct option).

CRITICAL FORMATTING RULE: The final line **MUST** be exactly 'Answer: VALUE' (e.g., 'Answer: 2') with **NO** additional spaces, **NO** periods, and **NO** extra text following the integer. The integer must be the absolute final character of your generated response.

A.5 CROP PROMPTS

GSM8K

You will answer a mathematical reasoning question concisely, showing the combined essential calculations leading directly to the final numerical answer. Omit intermediate labels. Focus on chained calculations. The last line of your response should be of the following format: 'Answer: VALUE' where VALUE is a numerical value.

BBH Object counting

You will answer a mathematical reasoning question. Present your reasoning solely through a single-line numerical calculation, summing only the pure numerical values identified from the question. Do not include item names, descriptions, or separate itemized lists. Do not include any introductory or conversational phrases. The last line of your response must be of the exact following format: 'Answer: VALUE' where VALUE is a pure numerical integer (without any currency symbols like '\$' or extra characters).

LogiQA

You will answer a logical reasoning question.

Be concise and direct in your reasoning. Avoid unnecessary introductions, restatements of the problem, or verbose explanations. Briefly explain why the correct answer is right and succinctly state why the incorrect options are wrong.

The very last line of your response **MUST** be exactly of the following format: 'Answer: VALUE' where VALUE is the zero-based numerical index of the correct answer. Do not add any extra spaces or characters to this final line.

A.6 GRADIENT PROMPTS

Accuracy Feedback Gradient Prompt

You are an expert Prompt Engineer.

You will give feedback to a variable with the following role: `<ROLE>` response from the language model `</ROLE>`.

Here is an evaluation of the variable using a string-based function:

Function purpose: The runtime of string-based function that checks if the prediction is correct, where the answer is followed by "Answer :"

`<INPUTS_TO_FUNCTION> {}`

Ground truth answer(role: correct answer for the query): `{}` `</INPUTS_TO_FUNCTION>`

`<OUTPUT_OF_FUNCTION> {} </OUTPUT_OF_FUNCTION>`

Objective: Your goal is to give feedback and criticism to the variable given the above evaluation output. Our only goal is to improve the above metric, and nothing else. `</OBJECTIVE_FUNCTION>`

We are interested in giving feedback to the response from the language model for this conversation. Specifically, give feedback to the following span of text:

`<VARIABLE> {} </VARIABLE>`

Given the above history, describe how the response from the language model could be improved to improve the Objective. Be very creative, critical, and intelligent.

Length Regularization Gradient

You are an expert Prompt Engineer and Efficiency Strategist.

You will give feedback to a variable (also known as MODEL_PROMPT) with the following role: `<ROLE>` structured system prompt to a somewhat capable language model that specifies the behavior `</ROLE>`.

`<MODEL_PROMPT> {} </MODEL_PROMPT>`

`<QUESTION> {} </QUESTION>`

`<MODEL_RESPONSE> {} </MODEL_RESPONSE>`

`<OBJECTIVE_FUNCTION>`Reduce the length of the model response by only few words without hurting the quality, and trying to keep as much reasoning and important stuff and focus on removing the fluff.`</OBJECTIVE_FUNCTION>`

Large language model takes `<MODEL_PROMPT>` and `<QUESTION>` as input and generates the `<MODEL_RESPONSE>`.

Given the above history, describe how the response from the language model could be updated to achieve the `<OBJECTIVE_FUNCTION>` without hurting the quality of the response. Be very creative, critical, and intelligent.

A.7 OPTIMIZER PROMPT

System Prompt Update Strategy

Here is the role of the variable you will improve:

`<ROLE>`structured system prompt to a somewhat capable language model that specifies the behavior and strategies for the QA task`</ROLE>`.

The variable is the text within the following span:

<VARIABLE> {current_system_prompt} </VARIABLE>

Here is the context and feedback we got for the variable:

<CONTEXT>

(The following block repeats for each conversation in the dataset:)

Here is a conversation:

<CONVERSATION>

<LM_SYSTEM_PROMPT> {current_system_prompt} </LM_SYSTEM_PROMPT>

<LM_INPUT> {questions[i]} </LM_INPUT>

<LM_OUTPUT> {forward_response[i]} </LM_OUTPUT>

</CONVERSATION>

This conversation is potentially part of a larger system. The output is used as response from the language model

Here is the feedback we got for structured system prompt to a somewhat capable language model that specifies the behavior and strategies for the QA task in the conversation:

<FEEDBACK>

{feedback_response[i]}

</FEEDBACK>

<REGULARIZATION_FEEDBACK>

{regularization_responses[i]}

</REGULARIZATION_FEEDBACK>

</CONTEXT>

Improve the variable (structured system prompt to a somewhat capable language model that specifies the behavior and strategies for the QA task) using the feedback provided in <FEEDBACK> and <REGULARIZATION_FEEDBACK> tags.

Send the improved variable in the following format:

<IMPROVED_VARIABLE>{the improved variable}</IMPROVED_VARIABLE>

Send ONLY the improved variable between the <IMPROVED_VARIABLE> tags, and nothing else.