Beyond Prompt Content: Enhancing LLM Performance via Content-Format Integrated Prompt Optimization

Anonymous ACL submission

Abstract

Large Language Models (LLMs) have shown significant capability across various tasks, with their real-world effectiveness often driven by 003 prompt design. While recent research has focused on optimizing prompt content, the role of prompt formatting-a critical but often over-007 looked dimension-has received limited systematic investigation. In this paper, we introduce Content-Format Integrated Prompt Optimization (CFPO), an innovative methodology that jointly optimizes both prompt content and formatting through an iterative refinement process. CFPO leverages natural language mutations to explore content variations and employs a dynamic format exploration strategy that systematically evaluates diverse format op-017 tions. Our extensive evaluations across multiple tasks and open-source LLMs demonstrate that CFPO demonstrates measurable performance improvements compared to content-only optimization methods. This highlights the importance of integrated content-format optimization and offers a practical, model-agnostic approach to enhancing LLM performance. Code is available at this link.

1 Introduction

036

Large Language Models (LLMs) have demonstrated impressive achievements across various domains (OpenAI, 2024a). The effectiveness of LLMs in real-world applications is fundamentally dependent on the design of effective prompts, which serve as an essential interface between human users or developers and the LLM system. Studies have shown that expert-designed prompts could significantly enhance LLM performance (Brown et al., 2020; Wei et al., 2023; Schulhoff et al., 2024).

However, manual design of prompts presents significant challenges, primarily due to the high sensitivity of LLMs to subtle variations in prompt characteristics, including both textual content and



Figure 1: The crucial role of prompt formatting and its interaction with content. (A): Model-specific format biases: Illustrates the performance sensitivity of two LLMs to different format styles on the GSM8K task, showing substantial variability in the effectiveness of 10 randomly selected formats. (B): For seven different prompt contents evaluated across 24 distinct formats, performance variations show the complex, interdependent relationship between prompt content and structure, demonstrating that no single format universally maximizes effectiveness.

structural format (Jiang et al., 2022; Zamfirescu-Pereira et al., 2023; Salinas and Morstatter, 2024). These sensitivities are further complicated by variations across different models and tasks (Zhuo et al., 2024; Sclar et al., 2024). To alleviate these difficulties, automated prompt optimization techniques, often leveraging the power of LLMs themselves, have proven to be an effective approach to adapt and refine prompts (Pryzant et al., 2023; Schnabel and Neville, 2024; Yang et al., 2024). However, existing research primarily focuses on optimizing *prompt content*, while overlooking a critical and largely unexplored dimension: the **prompt formatting**.

Our preliminary investigations, as illustrated in Figure 1, provide valuable insights into the role of prompt format in prompt optimization. We have observed that different LLMs display distinct preferences, with some formats performing well on one model but failing on another. This suggests sophis041



Figure 2: Illustration of the CFPO pipeline within a single iteration round. In the initial Component-wise Content Optimization stage, case-diagnosis and Monte-Carlo sampling are employed for content mutation. Subsequently, the Format Optimization stage identifies the most suitable format for each content candidate. The yellow dashed line indicates where the LLM optimizer is employed to guide the optimization process.

ticated, model-specific format biases (Sclar et al., 2024). Furthermore, we have identified a complex interplay between prompt content and format, where no single format consistently outperforms others across all contents. This lack of a universally optimal format highlights the impracticality of predefining format, and underscores the need for a joint optimization approach that treats prompt content and format as interdependent variables.

To address these limitations, we introduce **Content-Format Integrated Prompt Optimization (CFPO)**, an innovative methodology that concurrently optimizes both prompt content and format through an iterative refinement process. CFPO employs distinct optimization strategies tailored to the unique search spaces of content and format. Content optimization is guided by performance feedback and Monte Carlo sampling, leveraging natural language mutations to enhance prompt effectiveness. For format optimization, CFPO explores a discrete set of format options through a dynamic exploration strategy designed to identify optimal formats without requiring prior knowledge.

Specifically, CFPO's format optimizer leverages the principles of structured thinking, operating along two key dimensions: the *Prompt Renderer*, which governs the organizational structure of all components within a prompt (He et al., 2024), and the *Query Format*, which dictates the presentation of in-context learning examples and queries (Voronov et al., 2024a; Salinas and Morstatter, 2024). By integrating these two dimensions, CFPO defines a structured template that effectively distinguishes between content and format types, enabling the efficient identification of highperforming prompts.

Our primary contributions are threefold: (1) We

propose **CFPO**, an innovative approach to simultaneously optimizes prompt content and format using an iterative process. (2) We introduce an efficient strategy for dynamic format optimization that generates new formats in an iterative manner and evaluates formats instance through a scoring system to select the best option. (3) Through extensive evaluations across diverse tasks and multiple open-source LLMs, we demonstrate that CFPO consistently improves LLM performance in a measurable and effective manner.

100

101

102

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

133

2 Related Work

Optimization via LLM The remarkable capacity of LLMs has been demonstrated in various tasks as optimizers, leveraging their ability to enhance performance, such as code generation (Haluptzok et al., 2023; Zelikman et al., 2024; Askari et al., 2024), tool-making (Cai et al., 2024), and agent system design (Hu et al., 2024). However, recent studies indicate that LLMs face significant challenges in achieving completely automatic optimization. These models often rely on human intervention for designing workflows and struggle with tasks requiring complex decomposition and iterative refinement (Zhang et al., 2024; Li et al., 2024).

Automatic Prompt Optimization Automatic prompt optimization plays a crucial role in enhancing the performance of LLMs by refining prompts without requiring human intervention. Various approaches have been explored to search for the optimal prompt, including reinforcement learning (Zhang et al., 2023), Monte Carlo Search (Zhou et al., 2023), Monte Carlo Search (MCTS) (Wang et al., 2024b), feedbackbased methods (Pryzant et al., 2023; Das et al., 2024), and agent-driven frameworks (Wang et al.,

061



Figure 3: An illustrative example of our Structured Prompt Template. This template systematically organizes the prompt into distinct components, each serving a specific functional role. When formulating a prompt, the template first employs a Query format to present examples and queries, and then integrates all content components via the Prompt Renderer to construct the comprehensive prompt string.

2024a; Khattab et al., 2024; WHO, 2023). While these methods focus on optimizing the overall prompt, they often lack the capability for finegrained modifications. (Khattab et al., 2024; Schnabel and Neville, 2024) introduce phrase-level mutations, but they fail to address format mutations or implement them in a systematic manner.

134

135

136

138

140

Prompt structure and format Structured prompt-141 ing, which organizes prompts into distinct compo-142 nents such as instructions, examples, and queries, 143 holds significant potential in prompt engineer-144 ing (Fernando et al., 2023). Empirical rules for 145 prompt design always lack integration with auto-146 matic optimization techniques, limiting their scala-147 bility and effectiveness (Nigh, 2023; Google, 2024). 148 Frameworks like LangGPT (Wang et al., 2024a) have introduced structured prompting paradigms, 150 emphasizing reusable designs inspired by program-151 ming principles. However, these efforts primar-152 ily focus on content-level refinements and fail to adequately address the critical role of prompt for-154 matting. Studies have highlighted the impact of 155 formatting on prompt performance (Salinas and 156 Morstatter, 2024). Sclar et al. (2024) revealed that 157 modifications to separators and spacing within a query could substantially impact performance. He et al. (2024) reveals that the format of prompts 160 significantly impacts GPT-based models' performance, with no single format excelling universally. 163 Voronov et al. (2024b) focuses on the format of few-shot examples and suggests that it is beneficial 164 to maintain a consistent format across examples. 165 However, despite the recognition of formatting's importance, there remains a lack of comprehen-167

sive understanding regarding the optimization of prompt format in a systematic manner.

168

169

170

171

172

173

174

175

176

177

178

179

180

181

183

184

185

186

187

189

190

191

192

193

194

195

196

198

3 CFPO: Content-Format Integrated Prompt Optimization

As we have established, the effectiveness of LLMs is profoundly influenced by both the content and format of prompts. Existing automated prompt optimization methods have largely overlooked the format dimension, which exhibits a strong model bias. To address this critical limitation, we introduce Content-Format Integrated Prompt Optimization (CFPO) framework that jointly optimizes both prompt content and format. This contrasts with prior approaches focusing solely on content optimization. Our goal is to identify an optimal prompt p^* , comprising both content (c^*) and format (f^*), that maximizes performance on an evaluation dataset \mathcal{D} , given by:

$$p^*: (c^*, f^*) = \arg \max_{c \in \mathcal{L}, f \in \mathcal{F}} m(c, f | \mathcal{D}), \quad (1)$$

within the coherent natural language space \mathcal{L} and the space of all possible formats \mathcal{F} , guided by a metric function $m(\cdot)$ that assesses the prompt's quality.

To effectively search this complex space, CFPO employs a two-pronged iterative approach, detailed in Figure 2, that consists of two concurrently-run optimizers: a *Component-wise Content Optimizer* and a *Format Optimizer*. The content optimizer refines the textual content of a prompt, while the format optimizer explores the structural arrangement of its elements. Importantly, our framework



Figure 4: Built-in formats and rendering effects in our initial format pool. The final format configuration is achieved by selecting and combining elements from both the *Prompt Renderer* and the *Query Format* categories.

/div:

acknowledges the inherent interdependence of content and format and thus iterates between optimizing them, to find their optimal combination. The following sections detail our structured prompt template (Section 3.1), our innovative format optimization approach (Section 3.2), and finally our integrated optimization process (Section 3.3).

3.1 Structured Prompt Template

Output:{answer

199

201

210

211

214

215

219

222

To enable fine-grained and targeted optimization, our framework adopts a structured prompt template inspired by guidelines from OpenAI (2024b) and Google (2024). This template decomposes prompts into distinct functional components, facilitating both analysis and selective mutations. Specifically, our template divides a prompt into contentbased components and format-based components, as illustrated in Figure 3.

- The **Content-based Components** are:
- Task Instruction defines the primary goal, guidingthe model's overall behavior.
 - **Task Detail** offers supplementary task-specific information, including resolution steps.

Output Format specifies the desired output structure (e.g., JSON, bullet points, etc.).

Few-shot Examples provide contextual learning patterns, consisting of:

- *Examples*: specific instances pertinent to the task, including inputs and expected outputs.
 - *Example Hinter* (optional): a brief hint indicating that examples segment will follow, e.g., 'Here are some examples:'.

• *CoT Hinter* (optional): encourages a chain-ofthought reasoning process, e.g., 'Let's think step by step'. 230

231

232

233

234

235

236

238

239

240

241

242

243

244

245

246

247

249

250

251

252

253

254

255

257

Query shows the question or request to be answered by the LLM.

The Format-based Components are:

Query Format: defines how to structure the rendering of examples and queries.

Prompt Renderer defines how to aggregate all components into a structured prompt.

The formulation of the structured prompt template is fundamental to our optimization approach. This design yields two key advantages: first, it facilitates a structured component functionality where each part serves a specific purpose, promoting a more organized prompting framework; second, it enables fine-grained optimization by decoupling format from content, thus allowing targeted and precise modifications of individual components.

3.2 Format Optimizer Design

The key aspect of our work is the format optimization methodology. To efficiently explore the extensive range of prompt formats, the CFPO format optimizer adopts an approach that utilizes a format pool with a scoring system and an LLM-assisted format generation module. It strategically explores, evaluates, and refines formatting choices, all while learning from previous iterations.

304

307

259

3.2.1 Format Pool with Scoring System

The format pool is designed to hold the format configurations we use to generate prompts. As shown in Figure 4, these configurations are separated into two dimensions: the *Prompt Renderer*, which dictates the overall structure of the prompt, and the *Query Format*, which governs the rendering of incontext examples and queries. This distinction allows us to explore both macro and micro-level formatting variations.

To dynamically evaluate the potential of each format, we developed a scoring system for assessing the performance of each format f, represented as Q(f). This system updates the performance score of f across various prompt contents using the formula $Q(f) \leftarrow Q(f) + \sum_c m(c, f)$, where c represents each content instance in current round. Additionally, we maintain N(f) to count the number of times a format has been visited, which facilitates score normalization.

To initialize the exploration, we constructed an initial format search space \mathcal{F} , comprising a set of predefined commonly used formats, as illustrated in Figure 4. We also incorporate diverse variations of these predefined formats into the initial search space, such as adjustments to spacing, punctuation, and the use of special symbols. This establishes a starting point for our optimization.

3.2.2 LLM-assisted Format Generation

The variability of format space requires an automated process for effective expansion and exploration. To that end, we introduce an LLM-based format generator, LLM_{f_gen} , which autonomously generates new formats based on information in the existing format pool.

This evolutionary approach integrates the format generation into each optimization round, allowing for the creation of new and potentially beneficial formats. To enhance the efficiency of this process, we guide the LLM towards more promising areas by informing it of the performance function, $\frac{Q(f)}{N(f)}$. This iterative process not only diversifies the format pool but also ensures that our system can adapt to and incorporate a wide range of formats, thereby enhancing its utility and effectiveness. More detailed information of our format generation process is provided in the Appendix A.2.

3.2.3 Search Format via Format Optimizer

For each content candidate generated by the content optimizer, the format optimizer aims to identify Algorithm 1 Searching Optimal Format Given a Prompt Candidate

Input: $p_0 = (c_0, f_0)$: initial prompt, $p = (c, \cdot)$: current prompt candidate(with content c), \mathcal{F} : dynamic format pool, k: number of formats, $m(\cdot)$: evaluation metric, \mathcal{D} : evaluation data.

- 1: Initialize: $Q(f) \leftarrow m(c_0, f), N(f) \leftarrow 1$ for all $f \in \mathcal{F}$
- 2: Format Selection: $\mathcal{F}_{select} \leftarrow \{f \in \mathcal{F} : f \text{ is } in \text{ the top } k \text{ w.r.t. } UCT(f)\}$
- 3: Format Generation:
- 4: for each i = 0, 1, ..., k do
- 5: Generate format: $f_{new} \leftarrow LLM_{f_gen}(\mathcal{F})$
- 6: Collect f_{new} to \mathcal{F}_{gen} , and add f_{new} to \mathcal{F}
- 7: end for
- 8: Format Evaluation:
- 9: for each $f \in \mathcal{F}_{select} \cup \mathcal{F}_{gen}$ do
- 10: Evaluate m(c, f) with dataset \mathcal{D}
- 11: $Q(f) \leftarrow Q(f) + m(c, f)$
- 12: $N(f) \leftarrow N(f) + 1$
- 13: Update UCT(f) by Eq. 2
- 14: **end for**

Ì

15: $f \leftarrow \arg \max_{f \in \mathcal{F}_{select} \cup \mathcal{F}_{gen}} m(c, f)$

Output: The optimal format \hat{f} for content *c*

the most appropriate format from format pool. To navigate the balance between exploring new formats and exploiting known effective ones, we implemented the Upper Confidence Bounds applied to Trees (UCT) algorithm (Kocsis and Szepesvári, 2006). The UCT algorithm employs a selection criterion given by:

$$UCT(f) = \frac{Q(f)}{N(f)} + \alpha \sqrt{\frac{\sum_{f} N(f)}{N(f)}} \qquad (2)$$

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

329

330

where α serves as a balancing hyper-parameter, adjusting the trade-off between exploration and exploitation.

The overall process, outlined in Algorithm 1, selects 2k formats for evaluation in each optimization round: k promising formats from the pool (based on UCT score), and k new formats generated by the LLM_{f_gen} . The selected formats from both the existing pool (\mathcal{F}_{select}) and the newly generated pool (\mathcal{F}_{gen}) are then evaluated using a predefined metric function $m(\cdot)$, and the best-performing format among the tested candidates will be identified. The result is then incorporated into the pool for future iterations.

By iteratively evaluating formats, the format op-

430

timizer ensures a balance between exploring new
formats and refining current ones, converging to
the best format configuration.

3.3 Integrated Optimizer Design

337

361

364

372

373

CFPO orchestrates the Component-wise Content Optimization and Format Optimization within an iterative framework to jointly optimize content and format. This iterative process (illustrated in Figure 2) is key to our methodology.

Component-wise Content Optimization: This stage employs two primary strategies for mutat-341 ing the content of prompts. The first strategy is 342 case-diagnosis and revision, leveraging test cases to assess the efficacy of the current prompt. The 344 outcomes of these test cases, including both correct and incorrect samples, are analyzed by the LLM optimizer. This optimizer evaluates the performance and pinpoints specific components in need of optimization. Subsequently, targeted feedback is applied to these identified components for enhancement, resulting in improved prompts. For example, if the output is not in the specified format, the output format component will be altered. Additionally, a Mote-Carlo sampling strategy is employed to enhance the optimization robustness by generating synthetic content with same semantics for randomly selected components. After this step, we select top-performing content candidates based on an evaluation dataset for the next stage.

Format Optimization: As discussed in Section 3.2, this stage identifies the most effective format for each candidate prompt content from the previous content-optimization step. The format optimizer applies our dynamic format exploration and evaluation process, tracking performance, and updating the format pool's scoring system. The Format Optimizer meticulously tracks the performance of all evaluated formats, providing valuable insights to guide the selection of formats in subsequent iterations. Simultaneously, it retains only the most effective format for each prompt, ensuring the diversity of prompt content candidates during beam search.

In summary, the two optimizers work in tandem, leveraging the strengths of the LLM to facilitate swift adaptation and customization. Importantly, this iterative process allows for the optimization of format and content, thereby significantly enhancing the quality of the generated prompts.

4 **Experiments**

4.1 Experimental Setups

Dataset and Models. To rigorously evaluate CFPO, we selected a diverse set of tasks and models. Our benchmark tasks span various domains and complexities, including:

- **Reasoning**: GSM8K (Cobbe et al., 2021) and MATH500 (Hendrycks et al., 2021; Lightman et al., 2023) which require complex mathematical reasoning abilities.
- **Multiple-choice**: ARC-Challenge (Clark et al., 2018), demanding understanding and selection among alternatives.
- **Classification**: The *Implicatures* task from the Big-Bench benchmark (bench authors, 2023) to evaluate classification proficiency.

Our model selection includes a mix of foundational and instruction-tuned models to understand the generalizability of our approach:

- Foundational Models: Mistral-7B-v0.1 (Jiang et al., 2023) and LLaMA-3.1-8B (Meta, 2024b) represent pre-trained models.
- Instruction-Tuned Models: LLaMA-3-8B-Instruct (Meta, 2024a) and Phi-3-Mini-Instruct (Microsoft, 2024) represent models specifically fine-tuned for instruction following.

Furthermore, we use GPT-4 (2024-05-01-preview) as the LLM optimizer for content mutation and format generation (OpenAI, 2024a).

Implementation Details. The training process involved 20 iterative rounds, each consisting of content and format optimization. During content optimization, case-diagnosis and Monte Carlo sampling each generate 4 prompts per round. A set of 40 test cases is used, with 5 correct and incorrect cases leveraged for case-diagnosis. The number of prompt-structured components decreases progressively from 4 to 1, narrowing the search space over time to enhance efficiency. For format optimization, 4 UCT-selected formats and 4 newly generated formats are used to generate new prompts. The coefficient in the UCT selection process α is set to 1e - 3. Beam search, with a budget of 8, is employed during mutations to ensure effective exploration. Eval data sizes are configured as 50, 300, 500, and 500 for BigBench-Classification, MATH500, GSM8K, and ARC-Challenge, respectively. The best-performing prompt on the evaluation set for each method was selected and reported on the test set.

Method	Mistral-7B-v0.1	LLaMA-3.1-8B	LLaMA-3-8B-Instruct	Phi-3-Mini-Instruct
GSM8K				
Baseline (1-shot cot)	36.85	50.03	74.00	83.45
Baseline (8-shot cot)	38.21	51.02	73.46	85.75
GRIPS	39.04	50.27	74.53	83.47
APE	40.33	52.39	75.13	83.85
ProTeGi	45.72	54.74	75.36	84.84
SAMMO	43.82	54.74	75.89	84.76
CFPO (Ours)	53.22	63.38	80.74	89.16
		MATH-500		
Baseline (1-shot cot)	4.60	10.58	12.20	12.60
Baseline (4-shot cot)	10.20	23.40	14.00	40.40
GRIPS	13.40	15.80	23.60	10.80
APE	11.60	12.80	22.80	30.60
ProTeGi	10.80	17.00	18.40	28.80
SAMMO	12.20	15.40	25.80	42.40
CFPO (Ours)	14.80	26.99	33.33	44.20
		ARC-Challen	ge	
Baseline	67.15	73.81	75.94	84.39
GRIPS	77.05	77.90	79.61	87.46
APE	75.85	77.05	78.67	87.63
ProTeGi	76.54	77.22	79.86	87.54
SAMMO	77.22	77.13	79.86	87.03
CFPO (Ours)	79.35	78.50	80.63	88.23
Big-Bench Classification				
Baseline	56.00	64.00	70.00	54.00
GRIPS	86.00	67.00	84.00	69.00
APE	73.00	65.00	60.00	63.00
ProTeGi	83.00	81.00	82.00	76.00
SAMMO	86.00	80.00	86.00	78.00
CFPO (Ours)	94.00	90.00	91.00	87.00

Table 1: Main results on math reasoning tasks and commonsense reasoning tasks.

Baselines. To evaluate the effectiveness of CFPO, 431 we compared against several commonly used and 432 popular baselines. GrIPS (Prasad et al., 2023) 433 performs syntactic phrase-level edits in instruc-434 tion, representing a non-LLM-based optimization 435 APE (Zhou et al., 2023) and Proapproach. 436 TeGi (Pryzant et al., 2023) both employ LLM to 437 optimize prompt content, but differ in mutation 438 strategy. APE adopts an instruction induction ap-439 proach, while ProTeGi leverages test cases feed-440 back with LLM to guide the mutation process. 441 SAMMO (Schnabel and Neville, 2024) introduces 442 a structured framework that incorporates a prelimi-443 nary format mutation strategy, which relies on ran-444 dom selection from a predefined format pool. This 445 choice of baselines enables a comprehensive assess-446 ment of CFPO's capabilities against various types 447 of optimization approaches. All methods were eval-448 uated using consistent experimental configurations 449 to ensure a fair comparison. 450

Initial Prompts. To establish a reasonable starting
point, we employed a single in-context example
without any further instruction as the initial prompt
for each model and task, except for GrIPS which
requires an initial instruction. Chain-of-Thought

examples were employed for the reasoning tasks. We also report common baseline prompts, including 8-shot for GSM8K and 4-shot for MATH500. A comprehensive list of our initial prompts is in Appendix C. 456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

4.2 Main Results

Table 1 summarizes the performance of CFPO in comparison with several state-of-the-art methods across four datasets. The results highlight the superior performance of CFPO, significantly outperforming the baseline prompt as well as competing methods. We observed that pre-trained models exhibit greater sensitivity to prompt formatting, leading to substantial improvements when optimized by CFPO. Notably, optimized prompts for pre-trained models tend to be longer and incorporate more in-context examples, suggesting that these characteristics better align with the optimization needs of pre-trained models (see Appendix D.1). In contrast, instruction-tuned models display relatively more robust results and smaller gains, likely due to their inherent adaptability and generalization.

For the reasoning tasks, GSM8K and MATH, prompt optimization is especially impactful due to

the sensitivity of these tasks to prompt structure. 480 CFPO, which integrates unified content and for-481 mat optimization, delivers significant performance 482 gains. Specifically, the improvement for GSM8K 483 is more evident compared to the more challeng-484 ing MATH task, where the inherent complexity 485 limits the magnitude of improvement. Moreover, 486 feedback-based methods like ProTeGi, SAMMO, 487 and CFPO, consistently outperform the other base-488 lines because they leverage iterative feedback for 489 prompt refinement. In contrast, GRIPS, which is 490 limited to phrase-level mutations, exhibits marginal 491 improvements. These results underline the effec-492 tiveness of the integrated optimization strategy 493 adopted by CFPO. The selected optimal prompts 494 discovered by our approach can be found in Ap-495 pendix D. 496

4.3 Ablation Study

497

498 499

500

505

510

511

512

513

Impact of the Format Optimizer. CFPO incorporates a unique format optimization process, leveraging LLM for format generation and a UCT-based strategy for format selection. To evaluate its effectiveness, we evaluated two variations of our method: (1) $CFPO_c$, which optimizes content while keeping format fixed, and (2) $CFPO_c$ +Format, which first optimizes content, then performs a separate format optimization step. Table 2 shows that both $CFPO_c$ and $CFPO_c$ +Format underperform compared to the full CFPO approach, highlighting the importance of the integrated content and format optimization approach. The need for a joint optimization process which addresses the interdependence of content and format is essential for prompt optimization.

Task	Method	LLaMA-3.1-8B	LLaMA-3-8B-Instruct
GSM8K	ProTeGi	54.74	75.36
	$CFPO_c$	58.07	77.71
	$CFPO_c$ +Format	61.94	79.30
	CFPO	63.38	80.74
BBC	ProTeGi	81.00	82.00
	$CFPO_c$	85.00	85.00
	$CFPO_c$ +Format	88.00	89.00
	CFPO	90.00	91.00

Table 2: Ablation study of the format optimizer and content optimizer. CFPOc performs content optimization with a fixed format. CFPOc+Format performs format optimization after content optimization.

514Effectiveness of Format Generation. We com-515pared the full CFPO approach against a variant that516uses format from initial format pool without us-517ing LLM for generation. As presented in Table 4,518CFPO with format generation consistently outper-519forms the baseline relying solely on the initial pool.520These results demonstrate the effectiveness of the

proposed format exploration mechanism in enhancing both the quality and diversity of prompts.

Task	Method	LLaMA-3.1-8B	LLaMA-3-8B-Instruct
GSM8K	w/o Format Gen	62.70	78.85
	with Format Gen	63.38	80.74
BBC	w/o Format Gen	88.00	87.00
	with Format Gen	90.00	91.00

Table 3: Impact of format generation during promptoptimization.

Effectiveness of Format Selection. We further evaluated our UCT-based format selection process, compared it to a random selection from the format pool and a greedy selection without exploration (using $\alpha = 0$ in Eq. (2)). As presented in Table3, CFPO consistently achieves the best performance across all experimental settings, demonstrating the efficacy of the UCT-based selection strategy.

Task	Method	LLaMA-3.1-8B	LLaMA-3-8B-Instruct
GSM8K	Random	62.40	78.82
	$UCT(\alpha = 0)$	63.23	79.08
	UCT(ours)	63.38	80.74
BBH	Random	85.00	87.00
	$UCT(\alpha = 0)$	86.00	88.00
	UCT(ours)	90.00	91.00

 Table 4: Impact of different format selection strategies during optimization.

Effectiveness of the Content Optimizer. As presented in Table 2, we include ProTeGi (Pryzant et al., 2023), a baseline that optimizes only the content. In contrast, our CFPO_c, which incorporates structured prompting and integrates correct cases for diagnosis, achieves significant performance improvements, which highlights the effectiveness of our content optimization strategy.

5 Conclusion

This paper introduces Content-Format Integrated Prompt Optimization (CFPO), an innovative methodology that concurrently optimizes both prompt content and format. CFPO incorporates the Prompt Renderer and the Query Format within a structured prompt template. By leveraging distinct optimization strategies, CFPO discovers highperforming prompts that outperform content-only methods, addressing a critical gap in existing research. Our results demonstrate the substantial significant influence of format on LLM performance, underscoring the necessity of a joint optimization approach. These findings emphasize the importance of integrating content and format considerations in prompt engineering. CFPO represents a significant advancement, empowering developers to design effective prompts and unlocking the full potential of LLMs across diverse applications.

521 522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

Limitations While the proposed method demon-558 strates promising results, there are several limitations worth noting. First, the effectiveness of the 560 approach is task- and model-dependent. While the method generates promising prompts for specific tasks and models, it may not generalize as effec-563 tively to others—particularly tasks that are less 564 sensitive to prompt structure or models that already possess strong reasoning capabilities, thereby limiting its broader applicability. Moreover, the iterative 567 nature of the optimization process, with multiple 568 mutation strategies, introduces computational com-569 plexity, which could hinder scalability in resourceconstrained environments. Finally, while the for-571 mat generation mechanism shows strong potential, 572 its stability can be an issue, as the optimization process may not always yield consistent or optimal 574 results across different datasets or configurations.

References

576

578

579

583

584

587

588

589

590

595

596

598

606

607

- Arian Askari, Christian Poelitz, and Xinye Tang. 2024. Magic: Generating self-correction guideline for incontext text-to-sql.
- BIG bench authors. 2023. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. <u>Transactions on Machine Learning</u> Research.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners.
- Tianle Cai, Xuezhi Wang, Tengyu Ma, Xinyun Chen, and Denny Zhou. 2024. Large language models as tool makers.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. Preprint, arXiv:1803.05457.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021. Training verifiers to solve math word problems. arXiv preprint arXiv:2110.14168.
- Sarkar Snigdha Sarathi Das, Ryo Kamoi, Bo Pang, Yusen Zhang, Caiming Xiong, and Rui Zhang. 2024.

Greater: Gradients over reasoning makes smaller language models strong prompt optimizers. 610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

- Chrisantha Fernando, Dylan Banarse, Henryk Michalewski, Simon Osindero, and Tim Rocktäschel. 2023. Promptbreeder: Self-referential self-improvement via prompt evolution. <u>Preprint</u>, arXiv:2309.16797.
- Google. 2024. Prompting guide 101. https: //workspace.google.com/resources/ai/ writing-effective-prompts/.
- Patrick Haluptzok, Matthew Bowers, and Adam Tauman Kalai. 2023. Language models can teach themselves to program better.
- Jia He, Mukund Rungta, David Koleczek, Arshdeep Sekhon, Franklin X Wang, and Sadid Hasan. 2024. Does prompt formatting have any impact on llm performance? Preprint, arXiv:2411.10541.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset. <u>arXiv preprint</u> arXiv:2103.03874.
- Shengran Hu, Cong Lu, and Jeff Clune. 2024. Automated design of agentic systems. <u>Preprint</u>, arXiv:2408.08435.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b. <u>Preprint</u>, arXiv:2310.06825.
- Ellen Jiang, Kristen Olson, Edwin Toh, Alejandra Molina, Aaron Donsbach, Michael Terry, and Carrie J Cai. 2022. Promptmaker: Prompt-based prototyping with large language models. In Extended Abstracts of the 2022 CHI Conference on Human Factors in Computing Systems, CHI EA '22, New York, NY, USA. Association for Computing Machinery.
- Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Sri Vardhamanan, Saiful Haq, Ashutosh Sharma, Thomas T. Joshi, Hanna Moazam, Heather Miller, Matei Zaharia, and Christopher Potts. 2024. Dspy: Compiling declarative language model calls into self-improving pipelines.
- Levente Kocsis and Csaba Szepesvári. 2006. Bandit based monte-carlo planning. In European conference on machine learning, pages 282–293. Springer.
- Junyou Li, Qin Zhang, Yangbin Yu, Qiang Fu, and Deheng Ye. 2024. More agents is all you need. Transactions on Machine Learning Research.

63	Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri
64	Edwards, Bowen Baker, Teddy Lee, Jan Leike,
65	John Schulman, Ilya Sutskever, and Karl Cobbe.
666	2023. Let's verify step by step. arXiv preprint
667	arXiv:2305.20050.
	Mate 2024s. Introducing mate llamo2. The most same
800	Meta. 2024a. Introducing meta nama3: The most capa-
69	ble openly available llm to date.
570	Meta, 2024b. The llama 3 herd of models. Preprint.
571	arXiv:2407 21783
// 1	univ.2107.21705.
572	Microsoft, 2024. Phi-3 technical report: A highly capa-
573	ble language model locally on your phone. Preprint
574	arXiv:2404.14219.
675	Matt Nigh. 2023. Chatgpt3 free prompt
676	list. https://github.com/mattnigh/
577	ChatGPT3-Free-Prompt-List.
578	OpenAI. 2024a. Gpt-4 technical report. Preprint,
679	arXiv:2303.08774.
680	OpenAI. 2024b. Prompt generation. https:
681	<pre>//platform.openai.com/docs/guides/</pre>
682	prompt-generation/.
683	Archiki Prasad, Peter Hase, Xiang Zhou, and Mohit
684	Bansal. 2023. Grips: Gradient-free, edit-based in-
685	struction search for prompting large language models.
686	<u>Preprint</u> , arXiv:2203.07281.
07	Daid Dryzant Dan Itan Jamy Li Vin Tat Las Chan
087	Reid Pryzani, Dan Iter, Jerry Li, Yin Tat Lee, Chen-
	(\mathbf{r})
000	guang Zhu, and Whenael Zeng. 2025. Automatic
689	prompt optimization with "gradient descent" and
589 590	prompt optimization with "gradient descent" and beam search. page 7957–7968.
590 590	guang Zhu, and Michael Zeng. 2023. Automatic prompt optimization with "gradient descent" and beam search. page 7957–7968.Abel Salinas and Fred Morstatter. 2024. The butterfly
589 590 591	 guang Zhu, and Michael Zeng. 2023. Automatic prompt optimization with "gradient descent" and beam search. page 7957–7968. Abel Salinas and Fred Morstatter. 2024. The butterfly effect of altering prompts: How small changes and
589 590 591 592	 guarg Zhu, and Michael Zeng. 2023. Automatic prompt optimization with "gradient descent" and beam search. page 7957–7968. Abel Salinas and Fred Morstatter. 2024. The butterfly effect of altering prompts: How small changes and iailbreaks affect large language model performance.
589 590 592 593	 guarg Zhu, and Michael Zeng. 2023. Automatic prompt optimization with "gradient descent" and beam search. page 7957–7968. Abel Salinas and Fred Morstatter. 2024. The butterfly effect of altering prompts: How small changes and jailbreaks affect large language model performance. Preprint arXiv:2401.03729
589 590 591 592 593 594	 guang Zhu, and Michael Zeng. 2023. Automatic prompt optimization with "gradient descent" and beam search. page 7957–7968. Abel Salinas and Fred Morstatter. 2024. The butterfly effect of altering prompts: How small changes and jailbreaks affect large language model performance. <u>Preprint</u>, arXiv:2401.03729.
555 589 590 591 592 593 594	 guarg Zhu, and Michael Zeng. 2023. Automatic prompt optimization with "gradient descent" and beam search. page 7957–7968. Abel Salinas and Fred Morstatter. 2024. The butterfly effect of altering prompts: How small changes and jailbreaks affect large language model performance. <u>Preprint</u>, arXiv:2401.03729. Tobias Schnabel and Jennifer Neville. 2024. Symbolic
555 589 590 591 592 593 594 595 596	 guarg Zhu, and Michael Zeng. 2023. Automatic prompt optimization with "gradient descent" and beam search. page 7957–7968. Abel Salinas and Fred Morstatter. 2024. The butterfly effect of altering prompts: How small changes and jailbreaks affect large language model performance. <u>Preprint</u>, arXiv:2401.03729. Tobias Schnabel and Jennifer Neville. 2024. Symbolic prompt program search: A structure-aware approach
555 589 590 591 592 593 594 595 596 596	 guarg Zhu, and Michael Zeng. 2023. Automatic prompt optimization with "gradient descent" and beam search. page 7957–7968. Abel Salinas and Fred Morstatter. 2024. The butterfly effect of altering prompts: How small changes and jailbreaks affect large language model performance. Preprint, arXiv:2401.03729. Tobias Schnabel and Jennifer Neville. 2024. Symbolic prompt program search: A structure-aware approach to efficient compile-time prompt optimization. pages
555 590 591 592 593 594 595 596 597 598	 gualg Zhu, and Michael Zeng. 2023. Automatic prompt optimization with "gradient descent" and beam search. page 7957–7968. Abel Salinas and Fred Morstatter. 2024. The butterfly effect of altering prompts: How small changes and jailbreaks affect large language model performance. Preprint, arXiv:2401.03729. Tobias Schnabel and Jennifer Neville. 2024. Symbolic prompt program search: A structure-aware approach to efficient compile-time prompt optimization. pages 670–686.
555 589 590 591 592 593 594 595 596 597 598	 gualg Zhu, and Michael Zeng. 2023. Automatic prompt optimization with "gradient descent" and beam search. page 7957–7968. Abel Salinas and Fred Morstatter. 2024. The butterfly effect of altering prompts: How small changes and jailbreaks affect large language model performance. Preprint, arXiv:2401.03729. Tobias Schnabel and Jennifer Neville. 2024. Symbolic prompt program search: A structure-aware approach to efficient compile-time prompt optimization. pages 670–686.
555 589 590 591 592 593 594 595 596 597 598 599	 guang Zhu, and Michael Zeng. 2023. Automatic prompt optimization with "gradient descent" and beam search. page 7957–7968. Abel Salinas and Fred Morstatter. 2024. The butterfly effect of altering prompts: How small changes and jailbreaks affect large language model performance. <u>Preprint</u>, arXiv:2401.03729. Tobias Schnabel and Jennifer Neville. 2024. Symbolic prompt program search: A structure-aware approach to efficient compile-time prompt optimization. pages 670–686. Sander Schulhoff, Michael Ilie, Nishant Balepur, Kon-
555 589 590 591 592 593 594 595 596 596 597 598 599 700	 guang Zhu, and Michael Zeng. 2023. Automatic prompt optimization with "gradient descent" and beam search. page 7957–7968. Abel Salinas and Fred Morstatter. 2024. The butterfly effect of altering prompts: How small changes and jailbreaks affect large language model performance. Preprint, arXiv:2401.03729. Tobias Schnabel and Jennifer Neville. 2024. Symbolic prompt program search: A structure-aware approach to efficient compile-time prompt optimization. pages 670–686. Sander Schulhoff, Michael Ilie, Nishant Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si, Yin-
555 589 590 591 592 593 594 595 596 596 597 598 599 700 701	 guang Zhu, and Michael Zeng. 2023. Automatic prompt optimization with "gradient descent" and beam search. page 7957–7968. Abel Salinas and Fred Morstatter. 2024. The butterfly effect of altering prompts: How small changes and jailbreaks affect large language model performance. <u>Preprint</u>, arXiv:2401.03729. Tobias Schnabel and Jennifer Neville. 2024. Symbolic prompt program search: A structure-aware approach to efficient compile-time prompt optimization. pages 670–686. Sander Schulhoff, Michael Ilie, Nishant Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si, Yinheng Li, Aayush Gupta, HyoJung Han, Sevien Schul-
555 569 590 591 592 593 594 595 596 596 597 598 599 700 701 702	 guang Zhu, and Michael Zeng. 2023. Automatic prompt optimization with "gradient descent" and beam search. page 7957–7968. Abel Salinas and Fred Morstatter. 2024. The butterfly effect of altering prompts: How small changes and jailbreaks affect large language model performance. Preprint, arXiv:2401.03729. Tobias Schnabel and Jennifer Neville. 2024. Symbolic prompt program search: A structure-aware approach to efficient compile-time prompt optimization. pages 670–686. Sander Schulhoff, Michael Ilie, Nishant Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si, Yinheng Li, Aayush Gupta, HyoJung Han, Sevien Schulhoff, Pranav Sandeep Dulepet, Saurav Vidyadhara,
555 590 591 592 593 594 595 596 597 598 599 700 701 702 703	 guang Zhu, and Michael Zeng. 2023. Automatic prompt optimization with "gradient descent" and beam search. page 7957–7968. Abel Salinas and Fred Morstatter. 2024. The butterfly effect of altering prompts: How small changes and jailbreaks affect large language model performance. Preprint, arXiv:2401.03729. Tobias Schnabel and Jennifer Neville. 2024. Symbolic prompt program search: A structure-aware approach to efficient compile-time prompt optimization. pages 670–686. Sander Schulhoff, Michael Ilie, Nishant Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si, Yinheng Li, Aayush Gupta, HyoJung Han, Sevien Schulhoff, Pranav Sandeep Dulepet, Saurav Vidyadhara, Daveon Ki, Sweta Agrawal, Chau Pham, Gerson
555 589 590 591 592 593 594 595 596 597 598 599 700 701 702 703 704	 guang Zhu, and Michael Zeng. 2023. Automatic prompt optimization with "gradient descent" and beam search. page 7957–7968. Abel Salinas and Fred Morstatter. 2024. The butterfly effect of altering prompts: How small changes and jailbreaks affect large language model performance. Preprint, arXiv:2401.03729. Tobias Schnabel and Jennifer Neville. 2024. Symbolic prompt program search: A structure-aware approach to efficient compile-time prompt optimization. pages 670–686. Sander Schulhoff, Michael Ilie, Nishant Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si, Yinheng Li, Aayush Gupta, HyoJung Han, Sevien Schulhoff, Pranav Sandeep Dulepet, Saurav Vidyadhara, Dayeon Ki, Sweta Agrawal, Chau Pham, Gerson Kroiz, Feileen Li, Hudson Tao. Ashay Srivastava.
555 589 590 591 592 593 594 595 596 597 598 599 700 701 702 703 704 705	 guang Zhu, and Michael Zeng. 2023. Automatic prompt optimization with "gradient descent" and beam search. page 7957–7968. Abel Salinas and Fred Morstatter. 2024. The butterfly effect of altering prompts: How small changes and jailbreaks affect large language model performance. Preprint, arXiv:2401.03729. Tobias Schnabel and Jennifer Neville. 2024. Symbolic prompt program search: A structure-aware approach to efficient compile-time prompt optimization. pages 670–686. Sander Schulhoff, Michael Ilie, Nishant Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si, Yinheng Li, Aayush Gupta, HyoJung Han, Sevien Schulhoff, Pranav Sandeep Dulepet, Saurav Vidyadhara, Dayeon Ki, Sweta Agrawal, Chau Pham, Gerson Kroiz, Feileen Li, Hudson Tao, Ashay Srivastava, Hevander Da Costa, Saloni Gupta. Megan L. Rogers.
555 589 590 591 592 593 594 595 596 597 598 599 700 701 702 703 704 705 706	 guang Zhu, and Michael Zeng. 2023. Automatic prompt optimization with "gradient descent" and beam search. page 7957–7968. Abel Salinas and Fred Morstatter. 2024. The butterfly effect of altering prompts: How small changes and jailbreaks affect large language model performance. Preprint, arXiv:2401.03729. Tobias Schnabel and Jennifer Neville. 2024. Symbolic prompt program search: A structure-aware approach to efficient compile-time prompt optimization. pages 670–686. Sander Schulhoff, Michael Ilie, Nishant Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si, Yinheng Li, Aayush Gupta, HyoJung Han, Sevien Schulhoff, Pranav Sandeep Dulepet, Saurav Vidyadhara, Dayeon Ki, Sweta Agrawal, Chau Pham, Gerson Kroiz, Feileen Li, Hudson Tao, Ashay Srivastava, Hevander Da Costa, Saloni Gupta, Megan L. Rogers, Inna Goncearenco. Giuseppe Sarli Jeor Galvnker
300 389 390 391 392 393 394 395 396 397 398 399 700 701 702 703 704 705 706 707	 guang Zhu, and Michael Zeng. 2023. Automatic prompt optimization with "gradient descent" and beam search. page 7957–7968. Abel Salinas and Fred Morstatter. 2024. The butterfly effect of altering prompts: How small changes and jailbreaks affect large language model performance. Preprint, arXiv:2401.03729. Tobias Schnabel and Jennifer Neville. 2024. Symbolic prompt program search: A structure-aware approach to efficient compile-time prompt optimization. pages 670–686. Sander Schulhoff, Michael Ilie, Nishant Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si, Yinheng Li, Aayush Gupta, HyoJung Han, Sevien Schulhoff, Pranav Sandeep Dulepet, Saurav Vidyadhara, Dayeon Ki, Sweta Agrawal, Chau Pham, Gerson Kroiz, Feileen Li, Hudson Tao, Ashay Srivastava, Hevander Da Costa, Saloni Gupta, Megan L. Rogers, Inna Goncearenco, Giuseppe Sarli, Igor Galynker, Denis Peskoff Marine Carpuat Jules White Shva-
300 300 301 302 303 304 305 304 305 304 305 306 307 308 309 300 301 302 303 304 305 309 700 701 702 703 704 705 706 707 708	 guang Zhu, and Michael Zeng. 2023. Automatic prompt optimization with "gradient descent" and beam search. page 7957–7968. Abel Salinas and Fred Morstatter. 2024. The butterfly effect of altering prompts: How small changes and jailbreaks affect large language model performance. Preprint, arXiv:2401.03729. Tobias Schnabel and Jennifer Neville. 2024. Symbolic prompt program search: A structure-aware approach to efficient compile-time prompt optimization. pages 670–686. Sander Schulhoff, Michael Ilie, Nishant Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si, Yinheng Li, Aayush Gupta, HyoJung Han, Sevien Schulhoff, Pranav Sandeep Dulepet, Saurav Vidyadhara, Dayeon Ki, Sweta Agrawal, Chau Pham, Gerson Kroiz, Feileen Li, Hudson Tao, Ashay Srivastava, Hevander Da Costa, Saloni Gupta, Megan L. Rogers, Inna Goncearenco, Giuseppe Sarli, Igor Galynker, Denis Peskoff, Marine Carpuat, Jules White, Shyamal Anadkat Alexander Hoyle, and Philin Parnik
300 300 301 302 303 304 305 304 305 306 307 308 309 309 300 301 302 303 304 305 309 700 701 702 703 704 705 706 707 708 709	 guang Zhu, and Michael Zeng. 2023. Automatic prompt optimization with "gradient descent" and beam search. page 7957–7968. Abel Salinas and Fred Morstatter. 2024. The butterfly effect of altering prompts: How small changes and jailbreaks affect large language model performance. Preprint, arXiv:2401.03729. Tobias Schnabel and Jennifer Neville. 2024. Symbolic prompt program search: A structure-aware approach to efficient compile-time prompt optimization. pages 670–686. Sander Schulhoff, Michael Ilie, Nishant Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si, Yinheng Li, Aayush Gupta, HyoJung Han, Sevien Schulhoff, Pranav Sandeep Dulepet, Saurav Vidyadhara, Dayeon Ki, Sweta Agrawal, Chau Pham, Gerson Kroiz, Feileen Li, Hudson Tao, Ashay Srivastava, Hevander Da Costa, Saloni Gupta, Megan L. Rogers, Inna Goncearenco, Giuseppe Sarli, Igor Galynker, Denis Peskoff, Marine Carpuat, Jules White, Shyamal Anadkat, Alexander Hoyle, and Philip Resnik. 2024. The promut report: A systematic survey of
300 300 301 302 303 304 305 304 305 304 305 306 307 308 309 300 701 702 703 704 705 706 707 708 709 710	 guang Zhu, and Michael Zeng. 2023. Automatic prompt optimization with "gradient descent" and beam search. page 7957–7968. Abel Salinas and Fred Morstatter. 2024. The butterfly effect of altering prompts: How small changes and jailbreaks affect large language model performance. Preprint, arXiv:2401.03729. Tobias Schnabel and Jennifer Neville. 2024. Symbolic prompt program search: A structure-aware approach to efficient compile-time prompt optimization. pages 670–686. Sander Schulhoff, Michael Ilie, Nishant Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si, Yinheng Li, Aayush Gupta, HyoJung Han, Sevien Schulhoff, Pranav Sandeep Dulepet, Saurav Vidyadhara, Dayeon Ki, Sweta Agrawal, Chau Pham, Gerson Kroiz, Feileen Li, Hudson Tao, Ashay Srivastava, Hevander Da Costa, Saloni Gupta, Megan L. Rogers, Inna Goncearenco, Giuseppe Sarli, Igor Galynker, Denis Peskoff, Marine Carpuat, Jules White, Shyamal Anadkat, Alexander Hoyle, and Philip Resnik. 2024. The prompt report: A systematic survey of prompting techniques.
300 300 300 301 302 303 304 305 304 305 306 307 308 309 300 301 302 303 304 305 306 307 308 309 700 701 702 703 704 705 706 707 708 709 710	 guang Zhu, and Michael Zeng. 2023. Automatic prompt optimization with "gradient descent" and beam search. page 7957–7968. Abel Salinas and Fred Morstatter. 2024. The butterfly effect of altering prompts: How small changes and jailbreaks affect large language model performance. Preprint, arXiv:2401.03729. Tobias Schnabel and Jennifer Neville. 2024. Symbolic prompt program search: A structure-aware approach to efficient compile-time prompt optimization. pages 670–686. Sander Schulhoff, Michael Ilie, Nishant Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si, Yinheng Li, Aayush Gupta, HyoJung Han, Sevien Schulhoff, Pranav Sandeep Dulepet, Saurav Vidyadhara, Dayeon Ki, Sweta Agrawal, Chau Pham, Gerson Kroiz, Feileen Li, Hudson Tao, Ashay Srivastava, Hevander Da Costa, Saloni Gupta, Megan L. Rogers, Inna Goncearenco, Giuseppe Sarli, Igor Galynker, Denis Peskoff, Marine Carpuat, Jules White, Shyamal Anadkat, Alexander Hoyle, and Philip Resnik. 2024. The prompt report: A systematic survey of prompting techniques.
300 300 300 301 302 303 304 305 304 305 306 307 308 309 300 301 302 303 304 305 306 307 308 309 700 701 702 703 704 705 706 707 708 709 710 711	 guang Zhu, and Michael Zeng. 2023. Automatic prompt optimization with "gradient descent" and beam search. page 7957–7968. Abel Salinas and Fred Morstatter. 2024. The butterfly effect of altering prompts: How small changes and jailbreaks affect large language model performance. Preprint, arXiv:2401.03729. Tobias Schnabel and Jennifer Neville. 2024. Symbolic prompt program search: A structure-aware approach to efficient compile-time prompt optimization. pages 670–686. Sander Schulhoff, Michael Ilie, Nishant Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si, Yinheng Li, Aayush Gupta, HyoJung Han, Sevien Schulhoff, Pranav Sandeep Dulepet, Saurav Vidyadhara, Dayeon Ki, Sweta Agrawal, Chau Pham, Gerson Kroiz, Feileen Li, Hudson Tao, Ashay Srivastava, Hevander Da Costa, Saloni Gupta, Megan L. Rogers, Inna Goncearenco, Giuseppe Sarli, Igor Galynker, Denis Peskoff, Marine Carpuat, Jules White, Shyamal Anadkat, Alexander Hoyle, and Philip Resnik. 2024. The prompt report: A systematic survey of prompting techniques.
300 360 3	 guang Zhu, and Michael Zeng. 2023. Automatic prompt optimization with "gradient descent" and beam search. page 7957–7968. Abel Salinas and Fred Morstatter. 2024. The butterfly effect of altering prompts: How small changes and jailbreaks affect large language model performance. Preprint, arXiv:2401.03729. Tobias Schnabel and Jennifer Neville. 2024. Symbolic prompt program search: A structure-aware approach to efficient compile-time prompt optimization. pages 670–686. Sander Schulhoff, Michael Ilie, Nishant Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si, Yinheng Li, Aayush Gupta, HyoJung Han, Sevien Schulhoff, Pranav Sandeep Dulepet, Saurav Vidyadhara, Dayeon Ki, Sweta Agrawal, Chau Pham, Gerson Kroiz, Feileen Li, Hudson Tao, Ashay Srivastava, Hevander Da Costa, Saloni Gupta, Megan L. Rogers, Inna Goncearenco, Giuseppe Sarli, Igor Galynker, Denis Peskoff, Marine Carpuat, Jules White, Shyamal Anadkat, Alexander Hoyle, and Philip Resnik. 2024. The prompt report: A systematic survey of prompting techniques.
365 369 3	 guang Zhu, and Michael Zeng. 2023. Automatic prompt optimization with "gradient descent" and beam search. page 7957–7968. Abel Salinas and Fred Morstatter. 2024. The butterfly effect of altering prompts: How small changes and jailbreaks affect large language model performance. Preprint, arXiv:2401.03729. Tobias Schnabel and Jennifer Neville. 2024. Symbolic prompt program search: A structure-aware approach to efficient compile-time prompt optimization. pages 670–686. Sander Schulhoff, Michael Ilie, Nishant Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si, Yinheng Li, Aayush Gupta, HyoJung Han, Sevien Schulhoff, Pranav Sandeep Dulepet, Saurav Vidyadhara, Dayeon Ki, Sweta Agrawal, Chau Pham, Gerson Kroiz, Feileen Li, Hudson Tao, Ashay Srivastava, Hevander Da Costa, Saloni Gupta, Megan L. Rogers, Inna Goncearenco, Giuseppe Sarli, Igor Galynker, Denis Peskoff, Marine Carpuat, Jules White, Shyamal Anadkat, Alexander Hoyle, and Philip Resnik. 2024. The prompt report: A systematic survey of prompting techniques. Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. 2024. Quantifying Language Models' Sensitivity to Spurious Features in Prompt Design or: How I
300 360 3	 guang Zhu, and Michael Zeng. 2023. Automatic prompt optimization with "gradient descent" and beam search. page 7957–7968. Abel Salinas and Fred Morstatter. 2024. The butterfly effect of altering prompts: How small changes and jailbreaks affect large language model performance. Preprint, arXiv:2401.03729. Tobias Schnabel and Jennifer Neville. 2024. Symbolic prompt program search: A structure-aware approach to efficient compile-time prompt optimization. pages 670–686. Sander Schulhoff, Michael Ilie, Nishant Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si, Yinheng Li, Aayush Gupta, HyoJung Han, Sevien Schulhoff, Pranav Sandeep Dulepet, Saurav Vidyadhara, Dayeon Ki, Sweta Agrawal, Chau Pham, Gerson Kroiz, Feileen Li, Hudson Tao, Ashay Srivastava, Hevander Da Costa, Saloni Gupta, Megan L. Rogers, Inna Goncearenco, Giuseppe Sarli, Igor Galynker, Denis Peskoff, Marine Carpuat, Jules White, Shyamal Anadkat, Alexander Hoyle, and Philip Resnik. 2024. The prompt report: A systematic survey of prompting techniques. Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. 2024. Quantifying Language Models' Sensitivity to Spurious Features in Prompt Design or: How I learned to start worrying about prompt formatting.

Anton Voronov, Lena Wolf, and Max Ryabinin. 2024a.			
Mind your format: Towards consistent	evaluation		
of in-context learning improvements.	Preprint,		
arXiv:2401.06766.			

- Anton Voronov, Lena Wolf, and Max Ryabinin. 2024b. Mind your format: Towards consistent evaluation of in-context learning improvements. In Findings of the Association for Computational Linguistics: ACL 2024, pages 6287–6310, Bangkok, Thailand. Association for Computational Linguistics.
- Ming Wang, Yuanzhong Liu, Xiaoyu Liang, Songlian Li, Yijie Huang, Xiaoming Zhang, Sijia Shen, Chaofeng Guan, Daling Wang, Shi Feng, Huaiwen Zhang, Yifei Zhang, Minghui Zheng, and Chi Zhang. 2024a. Langgpt: Rethinking structured reusable prompt design framework for llms from the programming language.
- Xinyuan Wang, Chenxi Li, Zhen Wang, Fan Bai, Haotian Luo, Jiayou Zhang, Nebojsa Jojic, Eric P Xing, and Zhiting Hu. 2024b. Promptagent: Strategic planning with language models enables expert-level prompt optimization.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain-of-thought prompting elicits reasoning in large language models. <u>Preprint</u>, arXiv:2201.11903.
- WHO. 2023. Auto-gpt. https://github.com/ Significant-Gravitas/AutoGPT.
- Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V. Le, Denny Zhou, and Xinyun Chen. 2024. Large Language Models as Optimizers.
- J. D. Zamfirescu-Pereira, Richmond Y. Wong, Bjoern Hartmann, and Qiang Yang. 2023. Why johnny can't prompt: How non-ai experts try (and fail) to design llm prompts.
- Eric Zelikman, Eliana Lorch, Lester Mackey, and Adam Tauman Kalai. 2024. Self-Taught Optimizer (STOP): Recursively Self-Improving Code Generation.
- Jiayi Zhang, Jinyu Xiang, Zhaoyang Yu, Fengwei Teng, Xionghui Chen, Jiaqi Chen, Mingchen Zhuge, Xin Cheng, Sirui Hong, Jinlin Wang, Bingnan Zheng, Bang Liu, Yuyu Luo, and Chenglin Wu. 2024. Aflow: Automating agentic workflow generation.
- Tianjun Zhang, Xuezhi Wang, Denny Zhou, Dale Schuurmans, and Joseph E. Gonzalez. 2023. Tempera: Test-time prompting via reinforcement learning.
- Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy Ba. 2023. Large language models are human-level prompt engineers. <u>Preprint</u>, arXiv:2211.01910.
- Jingming Zhuo, Songyang Zhang, Xinyu Fang, Haodong Duan, Dahua Lin, and Kai Chen. 2024. Prosa: Assessing and understanding the prompt sensitivity of llms. <u>Preprint</u>, arXiv:2410.12405.

886

887

888

A Appendix: Detailed Optimization Process and Meta-Prompts

A.1 Meta-Prompt Header Setup

770

771

772

774

775

776

778

779

780

781

782

784

785

788

810 811

812

813

814

816

817

818

819

820

822

823

824

At the beginning of the prompt, we introduce the task and provide a detailed explanation of the prompt's components, followed by the current version of the prompt. Below is the structure of the meta-prompt header, where placeholders are denoted in [ALL CAPS]:

I'm trying to write a prompt to [TASK INTENTION].

The current prompt consists of several key components, including: [DESCRIPTION OF COMPONENTS]

The complete prompt is as follows: """[CURRENT PROMPT]"""

A.2 Format Generation

Our format generation process is a two-step procedure designed to create diverse and effective prompt formats. We focus on generating two key components of a prompt's format: the *Prompt Renderer* and the *Query Format*. The appendix presents examples of the format generated using this pipeline.

Step 1: Format Description Generation. For 796 each component (i.e., Prompt Renderer and the 797 Query Format), we first generate a natural language description of the format, alongside an example of how this format would render a sample input. This description acts as a blueprint, guiding the subse-801 quent code generation. We utilize a meta-prompt to instruct an LLM to perform this task. The metaprompt takes existing format examples as context 804 and generates new format descriptions along with rendered results. As an illustrative example, here is a conceptual outline of the meta-prompt employed 807 for generating new *Query Format* descriptions:

[META PROMPT HEADER]

We have some preset QUERY_FORMAT candidates, here are our whole search pool: [ALL EXISTING QUERY FORMATS DESCRIPTION] Here are two examples from our QUERY_FORMAT candidates as for your reference:

<Format name: Question-Answer>
[RENDERED EXAMPLE 1]

<Format name: Instruction-Response> [RENDERED EXAMPLE 2]

Please generate ONE new format for the QUERY_FORMAT segment, its description and render the provided example using this new format. The new format could either be a completely new format or a variation of an existing format.

If you choose to generate a completely new format, please ensure that the new format is conventional, structured, and aligned with commonly used query formats. Avoid overly creative or unconventional formats that deviate significantly from standard practices. The new format should be distinct from the existing formats.

The variation can focus on two parts, CASING and SEPARATOR:

CASING refers to both the capitalization of the text (e.g., f(x) = x.title(), f(x) = x.upper(), f(x) = x.lower()) and the specific wording or phrasing used (e.g., changing "question" to " instruction" or "input").

SEPARATOR: the punctuation or symbols used to separate the question and answer, there are some candidates as for your reference {{'', ' ', '\\ n', '--', ';\\n', ' ||', '<sep>', ' \\n', ':', '.'}}.

Note that focus solely on the format itself without altering the content of the question and answer. The format should remain focused on the existing structure (e.g., Question/Answer or Instruction/Response) without modifying the content or introducing any new sections. Avoid the use of underlines or any unconventional formatting styles among words. The format name should only include alphanumeric characters and underscores. Special characters such as `|`, `!`, `#`, `@`, and spaces should be avoided.

Please encapsulate the new query format using the following format:

<START> <Format name: [format name]> <Description: [format description]> [The example rendered by the newly generated format] <END>

Step 2: Format Code Generation. Based on the natural language description and rendered example produced in Step 1, we subsequently generate the corresponding code implementation of the new format. This code will be used by the system to render prompts according to the defined format. We again leverage a meta-prompt to instruct the LLM, this time to generate the executable code. As an illustrative example, here is a conceptual outline of the meta-prompt employed for generating the code representation of a new *Query Format*:

[META PROMPT HEADER]

We have some preset QUERY_FORMAT candidates, here are our whole search pool:

[ALL EXISTING QUERY FORMATS DESCRIPTION] 891 Here are two code implementations from our QUERY_FORMAT candidates as for your reference: 893 <Format name: Ouestion-Answer> 894 <Renderer code> [Ouestion-Answer RENDERER CODE] <Extractor code> [Question-Answer EXTRACTOR CODE] <Format name: Instruction-Response> 900 <Renderer code> 901 [Instruction-Response RENDERER CODE] 902 <Fxtractor code> 903 [Instruction-Response EXTRACTOR CODE] 904 905 Here is the example rendered by the new format: **FRENDERED RESULTS** 907 Please generate the code for this provided 908 909 example based on the new QUERY_FORMAT. Ensure 910 that both the renderer and extractor functions 911 are included. The generated code should be plain 912 Python code without any Markdown syntax or language identifiers such as ```python or ''' 913 python. Please output the code directly without 914 915 any additional formatting. If you need to use 916 any additional and specific packages, please import them in the code. Note that the generated 917 918 functions should include properly indented 919 blocks, so they can execute without errors. Note 920 that the renderer function name should be 921 query_renderer_{format_name} and the extractor function name should be guery_extractor_{ 923 format name}. 924 925 Please encapsulate the code using the following format: 927 <START> 929 <Format name: {format_name}> 930 <Description: {format_description}> 931 <Renderer code> 932 [Renderer code] 933 <Extractor code> 934 [Extractor code] 935 <FND>

A.3 Content Optimization

A.3.1 Case-diagnosis and Revision

As described in Section 3.3, content optimization is achieved through an iterative process of casediagnosis and feedback guided mutation. To facilitate this process, we utilize three distinct metaprompts, each tailored to a specific task within content optimization.

Case Diagnosis Meta-Prompt. This meta-prompt analyzes the current prompt's performance against a set of test cases. It identifies areas for improvement and suggests specific modifications for the next iteration.

949 [META PROMPT HEADER]

936

937

938

942

943

946

947

Upon evaluating the current prompt, this prompt gets the following examples wrong: [INCORRECT CASES] 951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998 999

1001

1002

1003

1004

1005

1006

1007

1009

1010 1011

1012

1013 1014

Meanwhile, this prompt gets the following examples correct: FCORRECT CASES1

Please review the provided examples of correct and incorrect answers, and identify [NUM OF DIAGNOSED COMPONENTS] specific area for improvement in the prompts. Each suggestion should focus on A SPECIFIC segment of the prompt that needs optimization. For each suggestion, provide a comprehensive explanation that encapsulates all the evaluation results. If you believe the EXAMPLES segment needs improvement, you may suggest one example that can be added, removed, or altered to enhance the EXAMPLES segment based on the examples given. If you think there is no need for improvement, do not return any prompt segment. Please encapsulate each suggestion using the following format:

<START> <Prompt segment: [Segment name]> [Suggestion goes here] <END>

Feedback Application Meta-Prompt. Based on the diagnosis, this meta-prompt generates targeted textual changes to enhance the prompt's performance. It directly modifies the identified components of the prompt based on the feedback.

[META PROMPT HEADER]

The existing [COMPONENT NAME] segment contains: [CURRENT CONTENT FOR THE COMPONENT]

Here are some suggestions for improving the [COMPONENT NAME] segments: [GENERATED DIAGNOSES]

Based on the above information, I wrote [NUMBER OF GENERATED CONTENT] distinct and improved versions of the [COMPONENT NAME] segment within the prompt. Each revised segment is encapsulated between < START> and <END>. In case this segment is an empty string, generate a suitable one referring to the suggestion. The [NUMBER OF GENERATED CONTENT] revised [COMPONENT NAME] segments are:

Feedback Application Meta-Prompt (for Examples). This meta-prompt specifically handles the optimization of few-shot examples. It revises examples by adding, deleting, or modifying one single instances, ensuring that the in-context learning process is effective.

[META PROMPT HEADER] The existing EXAMPLES segment contains: [CURRENT IN-CONTEXT EXAMPELS IN PROMPT] 1015 Here are some suggestions for enhancing the 1016 EXAMPLES segment: [GENERATED DIAGNOSES] 1017 Based on the above information, I have crafted [1019 NUMBER OF GENERATED EXAMPLES] improved version 1020 of the $\ensuremath{\mathsf{EXAMPLES}}$ segment within the prompt. Each 1022 revision represents ONLY ONE of the following 1023 specific actions: 1024 1. Addition: Incorporating one new example into 1025 the existing set. 2. Deletion: Eliminating one single example from 1026 1027 the current set. 1028 3. Modification: Changing the content of an 1029 example while maintaining its contextual 1030 relevance 1031 Please present the results without indicating which action was taken. Each refined EXAMPLES 1033 segment is marked by <START> and <END>. 1034 1035

The [NUMBER OF GENERATED EXAMPLES] revised EXAMPLES are:

1037 A.3.2 Monte-Carlo Sampling

1036

1038

1039

1040

1041

1042 1043

1044

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1058

1059

1060

1061

1062

1063

1064

1067

1069

1070

1071

1072

Monte-Carlo Sampling Meta-Prompt explores a wider range of semantically equivalent yet syntactically varied instructions, enhancing the chances of discovering more effective prompts.

[META PROMPT HEADER]

Please create a different version of [COMPONENT NAME] segment without changing its semantic meaning. In case this segment is an empty string, generate a suitable one. The existing [COMPONENT NAME] segment contains: [CURRENT CONTENT FOR THE COMPONENT]

The varied [COMPONENT NAME] segment is as follows:

Monte-Carlo Sampling Meta-Prompt (for Examples) refines few-shot examples by strategically adding, deleting, or modifying single instances to ensure their effectiveness.

[META PROMPT HEADER]

The existing EXAMPLE set contains: [CURRENT IN-CONTEXT EXAMPELS IN PROMPT]

Please generate a variation of the EXAMPLES set within the prompt while keeping the semantic meaning. The revision shoud represent ONLY ONE of the following specific actions: 1. Addition: Incorporating one new example into

the existing set. 2. Deletion: Eliminating one single example from

the current set. 3. Modification: Changing the content of an

example while maintaining its contextual relevance.

1073 Please present the results without indicating1074 which action was taken. The varied EXAMPLES1075 segment is as follows:

B Appendix: Examples of Generated 1076 Format 1077

1077 Here we select several format generated by GPT4 1078 in CFPO process. **B.1** Query Format **QA_Titlecase_Separator** 1081 1082 Question || In 3 years, Jayden will be half of Ernesto's age. If Ernesto is 11 years old, how 1083 many years old is Jayden now? 1084 Answer || Let's think step by step. Ernesto = 11 + 3 = <<11+3=14>>14 Jayden = 14/2 = <<14/2=7>>7 1086 in 3 years Now = 7 - 3 = <<7-3=4>>4 Jayden is 4 1087 years old. 1088 **QA_Brackets_Colon_Newline** 1089 [Question]: 1090 In 3 years, Jayden will be half of Ernesto's age. 1091 If Ernesto is 11 years old, how many years old 1092 is Jayden now? 1093 1094 [Answer]: 1095 Let's think step by step. Ernesto = 11 + 3 = <<11+3=14>>14 Jayden = 14/2 = <<14/2=7>>7 in 3 years Now = 7 - 3 = <<7-3=4>>4 1098 Jayden is 4 years old. 1099 QA_CapsBold_ColonNewline 1100 **QUESTION**: 1101 In 3 years, Jayden will be half of Ernesto's age. 1102 If Ernesto is 11 years old, how many years old 1103 is Jayden now? 1104 1105 **ANSWFR**: 1106 1107 Let's think step by step. Ernesto = 11 + 3 = <<11+3=14>>14 Jayden = 14/2 = 1108 <<14/2=7>>7 in 3 years Now = 7 - 3 = <<7-3=4>>4 1109 Jayden is 4 years old. 1110 Cascading_Statements 1111 Question: Statement 1 | Every element of a group 1112 generates a cyclic subgroup of the group. 1113 Statement 2 | The symmetric group S_10 has 10 1114 elements. 1115 Options: 1116 -A True, True 1117 -B False, False 1118 -C True, False 1119 -D False, True 1120 Answer: C 1121

Highlight_Separator_Case

QUESTION > Statement 1 | Every element of a1123group generates a cyclic subgroup of the group.1124Statement 2 | The symmetric group S_10 has 101125elements.1126OPTIONS > (A) True, True (B) False, False (C)1127True, False (D) False, True1128ANSWER > C1129

1122

B.2 Prompt Renderer 1130

Concise_Bullet_Points_Renderer 1131

1	11	32	
1	11	33	
-	11	34	
	11	35	
	н	36	
		30	
		37	
	11	38	
	11	39	
1	11	40	
1	11	41	
1	11	42	
1	11	43	
1	11	44	
1	11	45	
1	11	46	
1	11	47	
	11	48	
	11	10	
	1 I 1 4	-+3	
1		00	
1	11	51	
1	11	52	
		50	
1	11	53	
1	11	54	
1	11	55	
1	11	56	
1	11	57	
1	11	58	
	Ĥ	59	
	- 1 1	60	
	1 I 1 4	64	
	11 12	01	
1		02	
1	11	63	
1	11	64	
1	11	65	
1	11	66	
1	11	67	
-	н	68	
	1	50	
		~~	
1	1	69	
		70	
	1	<i>1</i> U	
1	11	71	
1	11	72	
1	11	73	
1	11	74	
1	H	75	
	11	76	
	1 I 1 4	10	
1		11	
1	11	78	
1	11	79	
1	11	80	
1	11	81	
-	Ĥ	82	
		-97.55	
-	н	83	
	. 1	- T- T-	

Task Instruction: Write a function that returns the sum of two numbers.
Task Detail: The function should take two numbers as input and return their sum.
Examples: Input: 1, 2 Output: 3
Query: Input: 1, 2 Output:

Tabular_Sections_Renderer

| Task Instruction | Write a function that returns the sum of two numbers. | | Task Detail | The function should take two numbers as input and return their sum. | | Examples | Input: 1, 2 Output: 3 | | Query | Input: 1, 2 Output: |

Checklist_Format_Renderer

[] **Task Instruction**
 Write a function that returns the sum of two numbers.

 [] **Task Detail**
 The function should take two numbers as input and return their sum.

- [] **Examples** Input: 1, 2 Output: 3

- [] **Query**
Input: 1, 2
Output:

C Appendix: Initial Prompt

C.1 GSM8K

Prompt Renderer: Directly Joint Query Format: QA

Q: There are 15 trees in the grove. Grove workers will plant trees in the grove today. After they are done, there will be 21 trees. How many trees did the grove workers plant today?

A: There are 15 trees originally. Then there were 21 trees after some more were planted. So there must have been 21 - 15 = 6. The answer is 6.

{{Query placeholder}}

C.2 MATH500

1184Prompt Renderer: Directly Joint

1185 Query Format: Question-Answer

1186A chat between a curious user and an AI1187assistant. The assistant gives step-by-step1188solutions to the user's questions. In the end of1189assistant's response, a final answer is given1190in the format of "The answer is: <ANSWER>.".

	1191
Here are some examples:	1192
Question: Let $[f(x) = \left(\frac{1}{2}\right)$	1193
\begin{array}{cl} ax+3, &\text{ if }x>2,	1194
x-5 &\text{ if } -2 \le x \le 2, \\	1195
2x-b &\text{ if } x <-2.	1196
\end{array}	1197
\right.\]Find \$a+b\$ if the piecewise function is	1198
continuous (which means that its graph can be	1199
drawn without lifting your pencil from the paper	1200
).	1201
Answer: Let's think step by step. For the	1202
piecewise function to be continuous, the cases	1203
must "meet" at \$2\$ and \$-2\$. For example, \$ax+3\$	1204
and \$x-5\$ must be equal when \$x=2\$. This	1205
implies \$a(2)+3=2–5\$, which we solve to get \$2a	1206
=-6 \Rightarrow a=-3\$. Similarly, \$x-5\$ and \$2x-	1207
b\$ must be equal when \$x=-2\$. Substituting, we	1208
get \$-2-5=2(-2)-b\$, which implies \$b=3\$. The	1209
answer is: \$a+b=-3+3=\boxed{0}\$.	1210
	1211
{{Query placeholder}}	1212
C.3 ARC-Challenge	1213
Prompt Renderer: Directly Joint	1214
Query Format: MultiChoice_QA	1215
You are a commonsense helper. I will provide	1216
several examples and a presented question. Your	1217
goal is to nick the most reasonable answer among	1218

several examples and a presented question. Your goal is to pick the most reasonable answer among the given options for the current question. Please respond with the corresponding label (A/B /C/D) for the correct answer.

1219

1220

1221

1222

1223

1224

1225

1226 1227

1228

1229

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1246

1247

1248

Here are some examples:

Question: Forests have been cut and burned so that the land can be used to raise crops. Which consequence does this activity have on the atmosphere of Earth? Choices: A: It reduces the amount of carbon dioxide production B: It reduces the production of oxygen C: It decreases the greenhouse effect D: It decreases pollutants in the air Answer: B

{{Query placeholder}}

C.4 Big-Bench Classification

Prompt Renderer: Directly Joint Query Format: Input-Output

Examples: 1241 Input: Speaker 1: 'You do this often?' Speaker 2: 1242 'It's my first time.' 1243 Output: no 1244 1245

{{Query placeholder}}

D Appendix: CFPO Results Analysis

D.1 In-context Examples and Text Length

Figure 5 presents an overview of the number of in-1249context examples and the text length of optimized1250



Figure 5: Overview of in-context examples and text lengths for various tasks and models.

prompts across various tasks and models. An in-1251 teresting pattern emerges: pre-trained models con-1252 sistently prefer prompts with longer text and more 1253 1254 in-context examples compared to instruction-tuned models. This observation suggests that pre-trained 1255 models benefit more from explicit context and de-1256 1257 tailed reasoning steps, which align with their less task-specialized nature. In contrast, the relative 1258 insensitivity of instruction-tuned models to prompt 1259 length and in-context examples supports the notion 1260 that these models have already trained with task-1261 1262 specific knowledge during fine-tuning, reducing their dependence on highly detailed prompts. 1263

D.2 Examples of Optimal Prompt

Here we selected several optimal prompts searched by CFPO.

LLaMA-3.1-8B on GSM8K

1264

1266

1267

1268

1269

1270

1271

1272

1273

1275

1276

1277

1278

1279

1280

1281

1282

1283

Understanding the Task: A Foundation for Mathematical Problem-Solving Your task is to methodically analyze the information provided and logically deduce the correct answer to the mathematical problem. Delve into each relevant detail, ensuring no critical step or aspect is overlooked. Approach the solution with a detailed-oriented mindset, ensuring every part of the process is considered to arrive at an accurate conclusion. Reflect on all the elements that might influence your reasoning or calculation, striving for thoroughness in your analysis.

> **Decoding Mathematical Language in Real-World Scenarios**

For the most effective problem-solving in mathematics, particularly when faced with intricate calculations over periods or under specific scenarios affecting results, an attentive and systematic method is key. Start by accurately determining the base numerical value. Then proceed by methodically listing every significant change whether it be increases, decreases, or modifications that impacts this base figure as the scenario unfolds, making sure to include each change in your overall computations. It's essential to focus on the concept of compounded operations, whether they re applied annually, monthly, or daily, and to thoughtfully evaluate the consequences of extraordinary events or circumstances (like an unexpected inheritance, a yearly loss, or a singular occurrence with a major impact) that might significantly shift the end calculations. Sharpen your attention on the dynamics of numerical relationships. particularly in cases involving ratios, proportions, and the impact of percentage changes over durations, to avoid common mistakes. Misunderstandings or misapplications of these numerical relationships can frequently cause inaccuracies. Thus, it is critical to scrutinize these mathematical relationships, whether they are of direct or inverse proportions, as well as the aggregate effects of consecutive percentage changes, as outlined in the problem description. This intensified attention is pivotal for an accurate and detailed resolution of complex issues, marked by multiplicative elements and interconnected circumstances. Reflect deeply on the significance of every step in the calculation process, absorbing the nuances of these changes, to systematically arrive at the most precise solution.

1284

1285

1286

1287

1288

1289

1290

1291

1292

1293

1294

1295

1296

1297

1298

1299

1300

1302

1303

1304

1305

1306

1307

1309

1310

1311

1313

1314

1315

1316

1317

1320

1322

1323

1324

1325

1326

1327

1328

1330

1331

1332

1333

1334

1335 1336

1337

1338

1339

1340

1342

1343

1344

1345

1346

1347

1348

1349

1350

1351

1353

**Ensuring Your Solution Fits the Scenario
Perfectly**

In presenting your solution, ensure it comprises both a numerical answer and a meticulously detailed explanation of the process leading to it. Begin with outlining the initial conditions and sequentially narrate the calculations you make at each step, highlighting any compounded operations or adjustments made to account for unique scenarios or conditions. This progression should clearly show how each step contributes to arriving at the final answer. For instance, if the task involves calculating the total costs saved over time with additional periodic benefits, your response should methodically explain: "Starting with an initial savings of X, plus Y every Z period, and considering an additional benefit of A every B period, leads to a total of...". This comprehensive breakdown not only bolsters the understanding of the mathematical principles applied but also provides a robust framework for identifying and rectifying any potential inaccuracies throughout the problem-solving process.

Examples to Illuminate the Path
To better grasp the concepts, consider the
following illustrative examples:
Question: There are 15 trees in the grove. Grove
workers will plant trees in the grove today.

1419

1420

1421

1422

1423

After they are done, there will be 21 trees. How many trees did the grove workers plant today? / ANSWER: Think through the problem step by step, diving into each segment for a thorough exploration to piece together the final answer. There are 15 trees originally. Then there were 21 trees after some more were planted. So there must have been 21 - 15 = 6. The answer is 6.

Question: A book club starts with a membership of 120. If the club increases its membership by 10% in the first year and then loses 5% of its members in the second year, what is the total membership at the end of the second year? / ANSWER: Think through the problem step by step, diving into each segment for a thorough exploration to piece together the final answer. The club starts with 120 members. In the first year, it increases by 10%, which is 0.10 * 120 = 12, so there are 120 + 12 = 132 members after the first year. In the second year, the club loses 5% of its members, which is $0.05 \times 132 =$ 6.6, but since the number of members must be an integer, we consider a loss of 7 members (assuming the figure is rounded up for practical reasons). Therefore, there are 132 - 7 = 125members at the end of the second year.

Question: Martin saves \$10 every week. In addition, every third week, he earns an extra \$15 from helping his neighbor. How much has Martin saved after 9 weeks? / ANSWER: Think through the problem step by step, diving into each segment for a thorough exploration to piece together the final answer. Martin saves \$10 each week, so over 9 weeks, he saves 9 * \$10 = \$90. Additionally, every third week, he earns an extra \$15, which occurs three times within 9 weeks (in the 3rd, 6th, and 9th weeks). So, he earns an extra 3 * \$15 = \$45 from helping his neighbor. Therefore, the total amount Martin has saved after 9 weeks is \$90 + \$45 = \$135.

Question: A teacher divides a class into groups for a project. If the ratio of boys to girls in the class is 3 to 2, and there are 30 students in the class, how many boys are in the class? / ANSWER: Think through the problem step by step, diving into each segment for a thorough exploration to piece together the final answer. The total ratio units for boys to girls in the class is 3 + 2 = 5. With 30 students in the class, each ratio unit represents 30 / 5 = 6students. Therefore, the number of boys, represented by 3 parts of the ratio, is 3 * 6 =18. The answer is 18.

Question: Grandma wants to order 5 personalized backpacks for each of her grandchildren's first days of school. The backpacks are 20% off of \$20 .00, and having their names monogrammed on the backpack will cost \$12.00 each. How much will the backpacks cost in total? / ANSWER: Think through the problem step by step, diving into each segment for a thorough exploration to piece together the final answer. The backpacks are 20% off of \$20.00, so the price after the discount is \$20.00 - (\$20.00 * 20%) = \$20.00 -\$4.00 = \$16.00 each. The monogramming costs an additional \$12.00 per backpack. Therefore, the total cost for each backpack is \$16.00 + \$12.00
= \$28.00. For 5 backpacks, the total cost will
be 5 * \$28.00 = \$140.00. The correct answer is
\$140.00.

1424

1425

1426

1427

1428

1429

1430

1431

1432

1433

1434

1435

1436

1437

1438

1439

1440

1441

1442

1443

1444

1445

1446

1447

1448

1449

1450

1451

1452

1453

1454

1455

1456

1457

1458

1459

1460

1461

1462

1463

1464

1465

1466

1467

1468

1469

1470

1471

1472

1473

1474

1475

1476

1477

1478

1479

1480

1481

1482

1483

1484

1485

1486

1487

1488

1489

1490

1491

1492

Query {{query}}

LLaMA-3-8B-Instruct on MATH-500

- Task Instruction: A chat between a curious user and an AI assistant focused on solving mathematical and reasoning tasks. The assistant is expected to deliver step-by-step solutions to the user's questions, emphasizing mathematical accuracy and rigor throughout the process. It must ensure that each mathematical operation and logical deduction is carefully examined and validated to derive the correct solution. At the conclusion of the response, the final answer should be presented in the format of "The answer is: <ANSWER>.", thereby confirming the solution 's validity and demonstrating a thorough understanding of the problem-solving approach.

- Task Detail: In addressing equation-based inquiries, precision in algebra, geometry, piecewise functions, complex numbers, and financial mathematics is paramount. This involves a detailed analysis of each equation, assessing every element and specific condition. For piecewise functions, it's critical to ensure continuity by solving for variables that maintain consistency across sections. In geometry, integrating measurements such as angles, lengths, and areas is fundamental. Algebraic queries require a consideration of all potential solutions and constraints, ensuring a comprehensive resolution. The addition of complex numbers into this mix necessitates a thorough understanding of their properties and operations to accurately determine both real and imaginary solutions. Similarly, tackling financial mathematics problems demands a deep comprehension of concepts such as compound interest, present value, and future value to make precise financial forecasts and comparisons. This holistic approach confirms that all aspects of the problem are considered and that the solution accounts for every requirement, assuring mathematical integrity in the resolution process.

- Output Format: 1. Solutions that involve fractions, square roots, or crucial mathematical figures (e.g., pi) must be simplified to their most fundamental form. This includes reducing fractions to their lowest terms and expressing square roots in their least complex radical form.

 Avoid the use of decimals unless the question explicitly requires it or they are necessary for conveying the most precise value possible.
 Present solutions involving square roots in their reduced radical form, ensuring the simplification process enhances comprehension without diluting mathematical integrity.
 In scenarios involving complex numbers, represent answers in their standard form (a + bi), ensuring both 'a' and 'b' are presented in their simplest, most refined state. This 1493 emphasizes the need for a clear, coherent 1494 representation of solutions encompassing complex 1495 numbers. 1496 5. Conclude your explanation with the statement: "The answer is: [<ANSWER>].", reinforcing 1497 1498 consistency and clarity across various mathematical challenges. This concluding 1499 1500 statement should encapsulate the solution in its 1501 simplest and most direct form, reflecting a 1502 thorough simplification and rationalization 1503 process. 1504 1505 Your explanation must delineate a detailed, step-1506 by-step progression leading to the final 1507 solution. This approach is not merely about 1508 arriving at the correct answer but about 1509 illuminating the path taken to get there, 1510 ensuring a deep understanding and clear 1511 demonstration of the reasoning behind each step. 1512 1513 - Examples: Here are some examples: 1514 ### Instruction: 1515 A rectangle ABCD has sides AB = 8 units and BC = 1516 6 units. A circle with a radius r units is 1517 inscribed within this rectangle. Calculate the radius r of the inscribed circle, ensuring the 1518 1519 answer is in its simplest form. 1520 1521 ### Response: 1522 We'll approach this problem by breaking it down 1523 into manageable steps. We start by understanding 1524 that the radius of the inscribed circle is 1525 equal to the distance from the center of the rectangle to any of its sides because the circle 1526 1527 is perfectly inscribed. In a rectangle, this 1528 distance is half the length of the rectangle's shorter side. Therefore, the radius r of the inscribed circle is half the length of BC, which 1529 1530 1531 is $6 \quad 1 = 3$ units. The answer is: r=3. 1532 1533 ### Instruction: 1534 Given a triangle where two sides are represented by complex numbers (3 + 4i) units and (1 - 2i)1535 1536 units, and the angle between them is 90 degrees, 1537 calculate the length of the hypotenuse. Ensure 1538 your answer includes a comprehensive breakdown of complex number operations and geometric 1539 1540 principles applied. 1541 1542 ### Response: 1543 We'll approach this problem by breaking it down 1544 into manageable steps. We start by acknowledging 1545 that the length of a side represented by a 1546 complex number can be found using the modulus of 1547 that number. The modulus of the first side is \$ 1548 $sqrt{3^2 + 4^2} = 5$ units, and the modulus of the second side is $\left\{1^2 + (-2)^2\right\} = \left\{1^2 + (-2)^2\right\}$ 1549 1550 {5}\$ units. Since these sides form a right 1551 triangle and we are given that the angle between them is 90 degrees, we can apply the 1553 Pythagorean theorem to find the length of the 1554 hypotenuse. The hypotenuse's length squared will 1555 be the sum of the squares of the lengths of the 1556 other two sides, which is $5^2 + (\sqrt{5})^2 =$ 25 + 5 = 30\$. Thus, the length of the 1557 1558 hypotenuse is $\operatorname{Sqrt}(30)$ units. The answer is: 1559 \$\sqrt{30}\$. 1560 - Query: 1562 {{query}}

LLaMA-3.1-8B on ARC-C

<pre>// class='TackInstruction'></pre>	156/
	1504
<n2>laskInstruction</n2>	1563
Your mission is to meticulously assess each	1566
situation presented alongside a specific	1567
question, employing your critical thinking and	1568
analytical skills. Your task comprises not	1569
only identifying the most logical and coherent	1570
choice (A/B/C/D) but also thoroughly	1571
evaluating how each option connects or	1579
diverges from the question's essence. This	1572
diverges from the question's essence. This	1573
requires a deep engagement with both the query	1574
and the choices, ensuring your reasoning is	1573
firmly anchored in the specifics of the	1576
options provided. It is essential to weave	1577
direct elements from the choices into your	1578
analysis, demonstrating a detailed	1579
understanding of how each option relates to	1580
the core question, and articulating why	1581
alternatives may be less fitting given the	1582
scenario. This approach ensures a nuanced and	1583
well-justified coloction process, grounded in	1500
the intervalue between the recetion process, grounded in	1004
the interplay between the question context and	158:
the specific details of the available choices	1586
	1587
	1588
<div class="TaskDetail"></div>	1589
<h2>TaskDetail</h2>	1590
In addressing the questions set before you,	1591
it is imperative to delve deeper than mere	1592
superficial observations or initial judgments	1593
Each scenario or question must be examined	1504
not just in its immediate context but within a	1504
hot just in its immediate context but within a	150
broader spectrum, tooking into the	1590
underpinning mechanisms or far-reaching	1597
effects of each option presented. This	1598
necessitates a thorough exploration of the	1599
larger implications and the scientific or	1600
logical foundations that dictate the outcomes.	1601
For instance, in environmental matters, it is	1602
vital to assess not just the immediate	1603
effects but the sustained impact on the	1604
ecosystem. In the realm of science, such as	160
when discerning chemical processes it is	1600
when discerning chemical processes, it is	1600
crucial to understand the molecular or atomic	1607
level changes that classify a reaction as a	1608
chemical change. This enhanced level of	1609
scrutiny and deeper analysis will lead to more	1610
accurate and well-founded choices, ensuring	1611
your responses are not just correct, but are	1612
also backed by a solid understanding of the	1613
underlying principles or long-term	1614
consequences	1614
	1616
<pre></pre>	161
<pre><uv <br="" class="outputFormat"> b 2> OutputFormat < /b 2></uv></pre>	101/
<nz>outputFormat</nz>	1010
For every query presented, your task is to	1619
identify the right choice from the options (A/	1620
B/C/D) accompanied by a concise rationale for	1621
your selection. This format is vital as it	1622
showcases the thought process leading to your	
	1623
decision, facilitating a comprehensive grasp	1623 1624
decision, facilitating a comprehensive grasp and interaction with the task.	1623 1624 1625
<pre>decision, facilitating a comprehensive grasp and interaction with the task. </pre>	1623 1624 1625 1626
<pre>decision, facilitating a comprehensive grasp and interaction with the task. </pre>	1623 1624 1625 1625
<pre>decision, facilitating a comprehensive grasp and interaction with the task. <div class="Examples"> <h2>Examples</h2></div></pre>	1623 1624 1625 1620 1627
<pre>decision, facilitating a comprehensive grasp and interaction with the task. <div class="Examples"> <h2>Examples</h2> <n>Hare are some examples;</n></div></pre>	1623 1624 1625 1626 1627 1628

1563

1631

Question: Forests have been cut and burned so

1632 that the land can be used to raise crops. Which 1633 consequence does this activity have on the 1634 atmosphere of Earth? A: It reduces the amount of carbon dioxide in 1635 1636 the atmosphere 1637 B: It reduces the availability of oxygen 1638 C: It lessens the greenhouse effect 1639 D: It lowers the levels of pollutants in the air 1640 Answer: B 1641 1642 Question: What is the most critical practice to 1643 ensure electrical safety while operating devices 1644 1645 A: Ensure the device does not come into contact 1646 with water. 1647 B: Use the device with hands covered in oil. 1648 C: Operate the device with wet hands. 1649 D: Leave the device plugged in when not in use. 1650 Answer: A 1651 1652 Question: Placing a plant cell in a hypertonic 1653 solution typically results in which of the 1654 following? 1655 A: The cell expanding as it absorbs water. 1656 B: No significant change due to the rigid cell 1657 wall. 1658 C: The cell shrinking as water exits the cell. 1659 D: Rapid division of the cell. Answer: C 1660 1661 1662 Question: What is the primary effect of using 1663 fossil fuels on global climate change? 1664 A: It leads to a significant reduction in 1665 greenhouse gases. 1666 B: It decreases the Earth's surface temperature. 1667 C: It increases the amount of greenhouse gases in the atmosphere. 1668 D: It contributes to a decrease in carbon 1669 dioxide levels. 1670 1671 Answer: C 1672 Question: The process of photosynthesis in 1673 1674 plants primarily involves which of the following 1675 transformations? A: Converting oxygen and glucose into carbon 1676 1677 dioxide and water 1678 B: Transforming water and carbon dioxide into 1679 oxygen and glucose 1680 C: Changing sunlight into chemical energy without producing oxygen 1681 D: Producing carbon dioxide and glucose from 1682 1683 oxvgen and water 1684 Answer: B 1685 1686 {{ query }}