DIVIDE-VERIFY-REFINE: ALIGNING LLM RE SPONSES WITH COMPLEX INSTRUCTIONS

Anonymous authors

003 004

010

Paper under double-blind review

ABSTRACT

011 Recent studies show that LLMs, particularly open-source models, struggle to follow complex instructions with multiple constraints, hindering their adoption in 012 mission-critical applications. Despite the importance, methods to improve LLMs' 013 adherence to such constraints remain largely unexplored, and current research fo-014 cuses primarily on evaluating this ability rather than developing solutions. While 015 a few studies enhance constraint adherence through model tuning, this approach 016 is computationally expensive and heavily reliant on training data quality. An al-017 ternative is to leverage LLMs' self-correction capabilities, allowing them to adjust 018 responses to better meet specified constraints. However, this self-correction abil-019 ity of LLMs is limited by the feedback quality, as LLMs cannot autonomously generate reliable feedback or detect errors. Moreover, the self-refinement process 021 heavily depends on few-shot examples that illustrate how to modify responses to meet constraints. As constraints in complex instructions are diverse and vary widely (e.g., text length, number of bullet points, or inclusion of specific key-023 words), manually crafting few-shot examples for each constraint type can be labor-intensive and sub-optimal. To deal with these two challenges, we propose 025 the **Divide-Verify-Refine (DVR)** framework with three steps: (1) **Divide** com-026 plex instructions into single constraints and prepare appropriate tools; (2) Verify: To address the feedback quality problem, these tools will rigorously verify re-028 sponses and provide reliable feedback (e.g., Python scripts for format checking 029 or pre-trained classifiers for content analysis); (3) **Refine**: To address the constraint diversity challenge, we design a refinement repository that collects suc-031 cessful refinement processes and uses them as few-shot demonstrations for future 032 cases, allowing LLMs to learn from the past experience during inference. Additionally, recognizing that existing datasets lack complexity and have internal conflict, we develop a new dataset of complex instructions, each containing 1-6 034 constraints. Experiments show that the framework significantly improves performance, doubling LLama3.1-8B's constraint adherence and tripling Mistral-7B's performance on instructions with 6 constraints. The code and dataset are available 037 at https://anonymous.4open.science/r/CODE_ICLR2025-52CE/README.md.

039 040

041

1 INTRODUCTION

042 Large language models (LLMs), like ChatGPT, have shown significant improvements across a va-043 riety of language tasks (Ouyang et al., 2022; Touvron et al., 2023). The success of LLMs relies on 044 their instruction-following ability to comprehend and execute complex instructions. Misinterpretations or failures to follow instructions can result in unintended outputs, which may have severe consequences (Mu et al., 2023; Zhou et al., 2023). This issue becomes particularly critical when 046 LLMs are deployed as in high-stakes environments, such as legal documentation or technical writ-047 ing. For example, when drafting legal contracts, LLMs must strictly adhere to constraints related 048 to format, specific terminology, and precise language usage to avoid misinterpretations or legal liabilities. Similarly, in technical writing, adhering to strict format guidelines, word limits, and the inclusion of essential technical terms is critical to ensure clarity, consistency, and compliance with 051 industry standards. 052

Recently, several studies show that LLMs, particularly open-source ones, struggle to follow complex instructions that contain multiple constraints, such as response length or formatting (He et al., 2024a;

066



Figure 1: (a) The LLMs hallucinate and cannot give reliable feedback. (b) The tools can check the response rigorously and provide reliable and directional feedback. The refinement repository 065 provides past refinement examples and stores the current refinement process.

067 Jiang et al., 2024b; Chen et al., 2024b). Despite the recognition of this issue, research on enhanc-068 ing LLMs' ability to follow constraints is still limited. Current efforts mainly focus on evaluating 069 LLMs' constraint-following ability rather than on improving this ability (Jiang et al., 2024b; Chen et al., 2024b; Zhou et al., 2023). Only very few studies improve the LLMs' constraint-following 071 ability through fine-tuning (He et al., 2024a; Sun et al., 2024; Li et al., 2024). Among them, one approach (He et al., 2024a) specifically enhances LLMs' ability to follow multiple constraints. It 072 employs a teacher model to iteratively refine the outputs of a student model. This step-by-step cor-073 rection process, along with the final accurate responses, is used to train the student model. Although 074 fine-tuning is an effective approach, it usually requires a large amount of computation resources 075 and heavily depends on the data quality. In contrast to training-based methods, the concept of 076 "self-correction" offers an alternative approach, where LLMs autonomously correct their responses 077 (Madaan et al., 2024; Shinn et al., 2024). Self-correction has been applied on various tasks such as question answering (Dhuliawala et al., 2023; Shinn et al., 2024) or mathematics (Madaan et al., 079 2024), where an LLM will evaluate its responses, give feedback, and further refine responses. For constraint-following, this self-correction process can be divided into two phases: verification and 081 self-refinement (see Fig. 1(a)). During the verification phase, LLMs assess whether their responses align with the specified constraints. If the responses do not align with the constraints, the LLMs will 083 give feedback that pinpoints errors and suggests adjustments. Following this, the self-refinement phase takes place where LLMs use the feedback to refine and improve their responses accordingly. 084

085 However, there are several challenges for an LLM to correct its response for multi-constraints. The first one is *feedback reliability*. Recent studies indicate that LLMs often exhibit only modest per-087 formance gains from self-correction, and the improvements can be unstable, occasionally even de-088 grading performance in areas such as question answering (Huang et al., 2024) and code generation (Olausson et al., 2024). Several studies claim that the significant bottleneck in self-correction is the 089 generation of reliable feedback (Tyen et al., 2024; Gou et al., 2024; Jiang et al., 2024a). LLMs, 090 including advanced models like GPT-4 and Claude 3, tend to have low recall in detecting LLMs 091 errors, underperforming significantly compared to humans (Kamoi et al., 2024). On the other hand, 092 research reveals that the self-correction performance on reasoning tasks is boosted if the error location is given, indicating LLMs have the self-correction ability given reliable feedback (Tyen et al., 094 2024). From a constraint-following perspective, LLMs are also not good at checking simple and 095 easy-to-verify constraints. As shown in Fig. 1(a), given a response, the LLMs struggle to accurately 096 count the bullet points, sentences, or words. The second challenge is *constraint diversity* which lies in the self-refinement process. Given the response and the feedback, the LLMs should refine the 098 response according to the feedback. However, to perform this task effectively, a set of representative few-shot examples is needed to demonstrate how to appropriately modify the response (Brown 099 et al., 2020). These constraints can vary widely, from adhering to a length limit to including spe-100 cific keywords. Each type of constraint needs distinct modifications. For example, meeting a length 101 limit might require removing content, whereas incorporating specific keywords requires adding text. 102 Manually crafting representative few-shot examples for each constraint type is labor-intensive. 103

104 To address these challenges, we propose a novel framework named Divide-Verify-Refine (DVR) 105 as illustrated in Fig. 1(b). To enhance feedback reliability, we observe that constraints that LLMs struggle to verify, such as exact word counts, sentences counts, or bullet points, can be readily 106 assessed using external tools. These tools include coding methods for quantitative measures, such 107 as counting the number of words, sentences, paragraphs, or bullet points, as well as pre-trained

108 classifiers for content analysis, such as topic and sentiment analysis, which are easily accessible and 109 widely available (Antypas et al., 2022; Loureiro et al., 2022). Instead of relying on LLMs to verify 110 their responses, we instruct LLMs to interact with external tools to handle the verification process. 111 We first instruct LLMs to divide the complex constraint into single constraints and then select one 112 tool for each single constraint. By integrating these tools, we overcome the limitations of LLMs in verification and feedback, enabling more rigorous checking and providing reliable and detailed 113 feedback for subsequent refinement. To address the constraint diversity problem, we propose the 114 refinement repository. This repository is a memory module that collects and stores successfully 115 refined examples for future use. Since the verification reliability is guaranteed by external tools, the 116 successful refinement processes can be recorded and saved in the refinement repository. This allows 117 LLMs to learn from past experiences during inference time. 118

Our main contributions are: (i) It is the first work to improve the LLMs' constraint-following abil-119 ity without training. Our framework enhances feedback reliability by integrating easily accessible 120 and widely available tools for verifying responses against specified constraints. (ii) To address the 121 constraint diversity challenge, we propose a novel refinement repository to store successfully refined 122 examples, which enables LLMs to learn from past experience and avoids crafting demonstrations 123 manually. (iii) Most benchmarks only contain 1-2 constraints (Chen et al., 2024b). Current com-124 plex instruction datasets He et al. (2024a) combine seed instructions with external constraints but 125 often overlook existing ones, causing conflicts and incomplete evaluations. To overcome the limita-126 tions and ensure a comprehensive evaluation, we construct a new complex instruction dataset with 127 instructions containing 1-6 constraints.

128 129

130 131

2 RELATED WORK

Instruction-Following of LLMs. Since the ability to follow instructions is crucial for the practical 132 use of LLMs, many recent studies evaluate this capability from various perspectives (Dubois et al., 133 2024; Zhou et al., 2023; Jiang et al., 2024b; Chen et al., 2024b; Zhou et al., 2023; He et al., 2024b). 134 They evaluate LLMs' instruction-following ability by testing on length (Dubois et al., 2024), format 135 (Zhou et al., 2023), semantic and topic constraints (Chen et al., 2024b). Most works only test 136 LLMs on simple instructions with only 1-2 constraints. Recently, some works generate instructions 137 with multiple constraints (He et al., 2024a; Jiang et al., 2024b). They find that LLMs struggle 138 to follow complex instructions as the number of constraints increases. Moreover, there is a big 139 performance gap between the open-source models and the closed-source models on instruction-140 following. Upon finding these problems, some seminal works aim to improve instruction-following 141 the ability of LLMs (Chen & Wan, 2023; Sun et al., 2024; Wang et al., 2024; He et al., 2024a). They 142 use various prompting strategies to generate instructions and responses with advanced LLM models (e.g., GPT4) and then use the generated data to fine-tune open-source LLMs. Though most methods 143 only consider instructions with few constraints, one of them (He et al., 2024a) focuses on improving 144 the LLMs' ability to follow multiple constraints. They generate complex instruction datasets by 145 merging instructions with external constraints. Then, they adopt GPT4 (teacher model) to modify 146 the response of the open source model (student model) iteratively. The student model is finetuned 147 by both the intermediate modification process and the final modified response. Although fine-tuning 148 is an effective approach, it usually requires a large number of computation resources and heavily 149 de- pends on the data quality. In addition, the fine-tuned model will still suffer from new constraints 150 not seen before. Different from previous methods, our framework uses in-context learning with tool 151 interaction to effectively identify and rectify unsatisfactory responses using the LLM itself, which 152 is more practical and accessible.

153 Self-Correction of LLMs. Self-correction is a framework where LLMs refine their own responses 154 during inference by reflecting on their initial responses (Shinn et al., 2024; Madaan et al., 2024). 155 This process can be divided into two phases. Initially, LLMs are prompted to analyze and provide 156 feedback on their own responses. Subsequently, based on the feedback LLMs refine the responses to 157 correct their mistakes. However, recent studies report negative results indicating that LLMs cannot 158 self-correct their own mistakes (Hong et al., 2024; Tyen et al., 2024; Kamoi et al., 2024; Gou et al., 159 2024). For example, Kamoi et al. (2024) reveal that top LLMs like GPT-4 and Claude 3 have low recall in detecting LLM errors, with LLMs significantly underperforming compared to humans. 160 Additionally, feedbacks provided by LLM self-correction tend to hallucinate and lack reliability. 161 This unreliability suggests that even when errors are detected, the guidance offered for corrections



Figure 2: The DVR framework: (a) Divide: The LLMs decompose constraints and instantiate tools
for each constraint, (b) Verify: Tools will give feedback on the response, (c) Refine: The refinement
repository provides past refinement process as few-shot examples. The current refinement process
will be stored in the repository.

may be incorrect or misleading. Hong et al. (2024) find that LLMs struggle to accurately identify 178 logical fallacies, casting doubt on their inherent ability to detect errors and conduct self-verification 179 reasoning effectively. However, the self-correction performance on reasoning tasks is boosted if the error location is given (Tyen et al., 2024). A study also shows that LLMs do have the ability 181 to correct their responses with the help of tools (Gou et al., 2024). All these observations indicate 182 that LLMs themselves are not reliable in analyzing their responses and a more reliable feedback 183 mechanism is needed to pinpoint the mistakes. The most closely related work is CRITIC (Gou et al., 2024), which uses tools for verification, such as an external API for testing the toxic score 185 of the response. Our method differs from CRITIC by incorporating a refinement repository module that accumulates successful past refinements, thereby enhancing the effectiveness of self-correction over time. Moreover, we use multiple tools to provide a detailed analysis of various aspects of the 187 response and give modification suggestions, instead of giving a single toxic score. 188

189 **LLMs Using Tools.** Tools have been extensively employed to enhance the capabilities of LLMs 190 across various domains. For instance, retrievers are used to augment response generation of LLMs 191 by fetching relevant information (Khandelwal et al., 2019), while search engines enhance the model's access to real-time data Nakano et al. (2021). Similarly, calculators are adopted to sup-192 port math reasoning of LLMs (Cobbe et al., 2021), interpreters are used to facilitate accurate code 193 generation (Chen et al., 2022; Gao et al., 2023), and mathematical provers help in verifying theo-194 retical proofs (Jiang et al., 2023). We observe that following complex instructions can be hard, but 195 checking them with tools is much easier. Because of this, we integrate external tools to our LLMs to 196 help with this process. This integration can provide reliable verification and detailed feedback and 197 also enable the LLMs to save past refinement examples for future use.

199 200

201

3 THE PROPOSED FRAMEWORK: DVR

As shown in Fig. 2, we propose the Divide-Verify-Refine (DVR) framework which consists of three 202 modules: (a) Divide instructions and prepare tools accordingly, (b) Verify responses and provide 203 feedback, and (c) *Refine* and store responses in a repository. First, The tool preparation module aims 204 to identify constraints, select appropriate tools, and fill out parameters. In this module, LLMs first 205 decompose the complex instructions into single constraints. For each single constraint, the LLMs 206 will prepare appropriate tools for verification. Second, in the verification and feedback module, the 207 prepared tools will verify the response and give detailed feedback if the response does not adhere 208 to the constraint. Third, in the self-refinement module, given the feedback and past refinement 209 experience for the same constraint as few-shot examples, the response is refined to adhere to the 210 target constraint. The successfully refined response will be stored in the refinement repository so 211 that it can be retrieved as refinement examples in the future. Next, we introduce each model in detail.

212

214

213 3.1 DIVIDE: TOOL PREPARATION

To provide accurate feedback, we propose to adopt tools for verification. To enable LLMs to use external tools, we construct tools for different types of constraints. We build Python verifiers for

structural constraints (e.g., the number of words, the number of paragraphs, or the number of bullet
points). To handle constraints where existing tools are unavailable, we can request advanced codegeneration models such as DeepSeek-Coder (Zhu et al., 2024) or Code-Llama Roziere et al. (2023)
to generate required tools. We also adopt existing classifiers as verifiers for content constraints (e.g.,
the topic classifier (Antypas et al., 2022) and the sentiment classifier (Loureiro et al., 2022)).

221 Given an input instruction I, the LLM \mathcal{M} first decomposes it into a series of individual constraints. 222 We use a decomposition prompt p_{decomp} asking LLMs for decomposition. With input instruction 223 and decomposition prompt, LLM then generates a set of decomposed constraints: $\mathcal{M}(p_{decomp}, I) \rightarrow I$ 224 $\{c_i\}_{i=1,2,3,..}$, where c_i is the *i*-th single constraint. For each constraint c_k , the LLM determines the 225 appropriate tool by matching c_k to a tool t_k from the predefined toolset: $\mathcal{M}(p_{select}, c_k) \to t_k$, where 226 $t_k \in \{t_i\}_{i=1,2,3,..}$ is the selected tool for the constraint c_k . The prompts for decomposition p_{decomp} and tool selection p_{select} are in Table 20 in Appendix A.13. After selecting the tools, the LLM sets 227 the necessary parameters for each tool, such as specifying the required number of bullet points or 228 the desired sentiment for the response. Finally, all tools relevant to instruction I are compiled into 229 the set $T_I = \{t_i\}_{i=1,2,3,..}$, ready to be utilized in the subsequent verification and feedback phase. 230

231 232

238

3.2 VERIFY: VERIFICATION AND FEEDBACK

Given the instruction, the LLM will first generate the initial response $R_0 = \mathcal{M}(p_{generate}, I)$, where $p_{generate}$ is the prompt for generation (detailed in Appendix A.13). We denote the current response as R and $R = R_0$ for the first round of refinement and will be updated to the refined response in subsequent rounds. The current response is verified by each tool in toolset T_I as follows:

$$f_i = t_i(R), \forall t_i \in T_I \tag{1}$$

239 where f_i is the feedback from tool t_i for constraint c_i . If the response adheres to the constraint, 240 the feedback is a boolean value "true". Otherwise, f_i is a textual feedback that first identifies the 241 error in the response and then suggests modification. For example, as shown in Fig. 2, the tool 242 "Bullet_points(4)" counts the number of bullet points in the response and outputs "true" if there 243 are 4 bullets; while the response only contains 2 bullets. It finds that the response does not satisfy 244 the constraint and gives out the feedback "The response only contains 2 bullet points. 2 more bullet points should be added." This detailed feedback points out the errors in the response and 245 gives directional information for LLMs to modify the response. We collect all feedback F_I = 246 ${f_i}_{i=1,2,3...}$ which will be used to refine the response R. 247

248 249

262

3.3 REFINE: SELF-REFINE WITH FEEDBACK AND FEW-SHOT DEMONSTRATION

250 In the self-refinement phase, the LLM leverages the feedback collected to refine the response. To 251 improve the performance, we propose to adopt representative demonstrations in the prompt to in-252 struct LLMs on how to conduct refinement using feedbacks. As constraints vary widely, each type 253 of constraint requires specific demonstrations for effective refinement. Manually creating few-shot 254 examples for each constraint type is labor-intensive and impractical for real-world applications. To 255 solve this issue, we propose to store the successful refinement process in the refinement reposi-256 tory. When LLMs need to refine a new response involving the same constraint type, the few-shot examples can be retrieved from the refinement repository as in-context examples. 257

Specifically, the refinement process targets one unsatisfied constraint at a time, cycling through a refine-verify-refine loop until all constraints are satisfied. For a given response R and the feedback $f \in F_I, f \neq True$, the refinement response can be written as follows:

$$R' = \mathcal{M}(p_{refine}, s^t, I, R, f) \tag{2}$$

263 where p_{refine} is the prompt for refinement (detailed in Appendix A.13), s^t = 264 $\{(I_i, R_i, f_i, R'_i)^t\}_{i=1,2,3...}$ is the set of refinement examples selected from the refinement 265 repository Q, which contains refinement examples having the same constraint type associated with 266 f. There might be many refinement examples having the same constraint type with f available in the refinement repository. Retrieval techniques like semantic similarity can be employed to select 267 the most relevant examples. In this paper, we randomly select relevant examples for simplicity and 268 leave more advanced techniques as future work. Some refinement examples are in Table 18 and 19 269 in Appendix A.12.

If the refined response adheres to the constraint, i.e., t(R') = True, the current successful refinement process will be stored in the repository as $Q = Q \cup \{(I, R, f, R')^t\}$.

Discussion. Our method, DVR introduces a novel approach to enhancing LLMs' ability to follow 273 complex instructions with multiple constraints. The detailed algorithm of DVR is shown in Algo-274 rithm 1 in Appendix A.1. By integrating external tools for reliable and detailed feedback and a 275 refinement repository for storing successful refinement examples, we provide a scalable and robust 276 framework for improving instruction compliance without the need for extensive retraining. More-277 over, the external tools and the refinement repository work jointly. Without reliable feedback, the 278 refinement repository would risk accumulating incorrect or noisy examples, which could deteriorate 279 the performance of LLMs over time. The detailed feedback gives "directional" information, which 280 guides the LLMs to adjust their responses. Compared to directly following complex instructions, decomposing these instructions and selecting the appropriate tools are simpler tasks for LLMs. This 281 inherent advantage allows our framework to be very effective, as it leverages these easier tasks to 282 build a robust system that enhances the LLMs' adherence to constraints. 283

284 285

286 287

288

289

290

4 EMPIRICAL VALIDATION

In this section, we conduct experiments to answer the following research questions: (**RQ1**) Can our DVR improve the ability of LLMs to follow complex constraints? (**RQ2**) How does the performance of LLMs differ across various types of constraints, and which constraints pose the greatest challenges? (**RQ3**) How does each module of DVR (the tool-assisted verification and the few-shot self-refinement library) individually contribute to improving LLMs' ability to follow constraints?

291 292 293

294

4.1 EXPERIMENTAL SETUP

295 Datasets. We conduct experiments on two datasets: (i) We conduct experiments on CoDI (Control-296 lable Generation under Diversified Instructions) (Chen et al., 2024b). It has 500 instructions with 297 2 constraints. Each instruction has a topic constraint and a sentiment constraint. (ii) ComplexIn-298 struct: Since the complexity of CoDI is limited, we construct a new complex instruction datasets 299 called ComplexInstruct. We use CoDI (topic instruction set) (Chen et al., 2024b) as seed instructions which ask users to generate text on certain topics. Some instructions ask users to generate 300 "a paragraph of..." or "a sentence of ..." which already contain length constraints. To avoid con-301 flicts and hidden constraints, we remove these constraints by replacing the keywords "paragraph" 302 and "sentence" with "text". Then, we synthesize complex instructions by adding constraints to 303 these seed instructions (Zhou et al., 2023). To simulate instruction of different levels, we generate 304 6000 complex instructions with 1-6 constraints for each instruction as 6 levels (1000 instructions for 305 each level). We have 21 types of constraints categorized into 8 general categories (such as length 306 constraint, punctuation, and case change). Each type of constraint is diversified into 8 different 307 expressions. The detailed information about the constraint types is in Appendix A.3 and Table 6.

308 **Baselines.** We compare our method with representative and state-of-the-art baselines, which can 309 be categorized into three main types: (i) Self-reflection based methods, which iteratively improve 310 response via feedback from LLMs reflection, such as Reflexion (Shinn et al., 2024); (ii) Prompting 311 based methods, which use different prompting strategies to get the best response, including Branch-312 solve-Merge (BSM) (Saha et al., 2024) and Universal Self-Consistency (U-SC) (Chen et al., 2024a); 313 and (iii) Tool based methods, which use external tools for feedback or selection, such as Rejection 314 sampling (Saunders et al., 2022), React (Yao et al., 2023), and CRITIC (Gou et al., 2024). For the 315 refinement repository of our framework, we consider two variants, i.e., warm-start and cold-start. For warm-start, we have an additional set of instructions (6000 samples for ComplexInstruct and 316 500 samples for CoDI). Note that these data samples are totally independent with test set. Our 317 framework will first run on these samples to collect examples to fill the refinement repository. For 318 cold-start, since the refinement repository is empty at beginning, we use 5 fixed few-shot examples 319 if there are no examples that can be retrieved from the repository. The detailed information about 320 each baseline is in Appendix A.2. 321

Implementation. We test on popular open-source models including Mistral-7B, Llama3-8B,
 Llama3.1-8B and Llama3.1-70B. The temperature of the model is 0.8. We set the number of few-shot demonstrations for initial response generation and self-refinement (without repository) as 5 for

326	in parentheses $(+xx)$ indic	cate the imp	provement c	compared to	the best pe	rforming ba	aseline.
327	Method	Level 1	Level 2	Level 3	Level 4	Level 5	Level 6
328	Vanilla	90.5	76.6	62.5	50.1	35.6	25.3
329	Reflxion	91.6	78.1	63.7	49.8	35.8	25.7
330	BSM	90.1	75.3	62.0	47.5	35.5	24.1
000	U-SC	90.9	76.3	62.4	47.1	36.0	25.8
331	Rejection Sampling	92.1	86.7	71.1	60.4	49.8	36.3
332	ReAct	94.2	86.1	72.5	60.7	50.2	37.2
333	CRITIC	93.8	87.1	75.4	64.4	52.4	43.2
334	\overline{DVR} (coldstart)	94.5 (+0.7)	87.9 (+0.8)	78.4 (+3.0)	69.5 (+5.1)	60.9 (+8.5)	49.2 (+6.0)
335	DVR (warmstart)	95.2 (+1.4)	88.7 (+1.6)	79.2 (+3.8)	69.7 (+5.3)	60.5 (+8.1)	49.6 (+6.4)
000							

Table 1: Instruction Satisfaction Rate (ISR) across levels 1 to 6 (Llama-3.1-8B-Instruct). The values in parentheses (+xx) indicate the improvement compared to the best performing baseline.

autore 2. Instruction Satisfaction Rate (ISR) across revers 1 to 0 (Wistrat-7D-mstruct-v0).	Table 2: Instruction	Satisfaction Rate	(ISR)	across levels	1 to 6	(Mistral-7B-Instruct-v0.)
---	----------------------	-------------------	-------	---------------	--------	---------------------------

		()				
Method	Level 1	Level 2	Level 3	Level 4	Level 5	Level 6
Vanilla	77.0	55.3	34.1	19.9	12.4	6.3
Reflxion	77.2	55.8	35.1	20.1	12.0	5.8
BSM	78.1	56.2	33.8	19.3	11.3	5.2
U-SC	76.8	56.0	34.3	20.4	12.9	5.8
Rejection Sampling	78.4	58.3	37.6	23.0	13.5	6.8
ReAct	86.0	67.8	46.0	32.5	18.2	10.7
CRITIC	88.9	72.5	55.6	43.5	28.1	18.1
DVR (coldstart)	94.9 (+6.0)	80.2 (+7.7)	64.1 (+8.5)	49.3 (+5.8)	35.8 (+7.7)	23.6 (+5.5)
DVR (warmstart)	95.0 (+6.1)	81.3 (+8.8)	66.6 (+11.0)	51.4 (+7.9)	36.4 (+8.3)	23.4 (+5.3)

our method and every baseline. We use the same set of few-shot demonstrations both for baselines and our method. We also set the maximum number of few-shot demonstrations for refinement (with repository) as 8. We set the number of trials as 5 for our method and every baseline.

Evaluation Metrics. We assess the constraint-following ability by calculating the Instruction Satisfaction Rate (ISR) (Jiang et al., 2024b). Specifically, each single instruction is satisfied when all constraints in that instruction are satisfied. It is calcualted as $ISR = \frac{1}{N} \sum_{i=1}^{N} \prod_{j=1}^{m_i} c_{ij}$, where N is the total number of instructions in the dataset, m_i is the number of constraints in the *i*-th instruction, $c_{ij} = 1$ if the *j*-th constraint in *i*-th instruction is satisfied; otherwise $c_{ij} = 0$.

356 357

358

324

325

337 338 339

347

348

349

350

4.2 RQ1: Assessing the Constraint-Following Ability

To answer RQ1, we evaluate our framework on two datasets. We evaluate structural constraints (e.g., text length, number of sections, and bullet points) on ComplexInstruct and content constraints (e.g., topic and sentiment constraints) on CoDI respectively.

362 For ComplexInstruct, there are six difficulty levels. Each level corresponds to the number of constraints in the instruction. For example, each instruction in Level 3 contains three constraints. There are 1000 instructions for each level. Results are shown in Table 1 and Table 2 (more results in 364 Appendix A.4). (i) Single constraints vs Multi-constraints: For instructions in different difficulty levels, responses to instructions with more constraints tend to have a lower satisfaction rate. The 366 satisfaction rate of instructions at Level 1 can approach close to 100%. However, at Level 6, the 367 satisfaction rates for models like Llama3.1-8B and Mistral-7B drastically fall to 25% and 6.3%, re-368 spectively. This indicates that LLMs struggle to satisfy instructions with multiple constraints even 369 though LLMs can satisfy them individually. (ii) Self-Reflection is unreliable: We can observe that 370 Reflxion (Shinn et al., 2024) where LLMs reflect on and self-correct their responses, provides little 371 improvement over Vanilla. This result indicates that LLMs themselves can not effectively iden-372 tify their errors in constraint-following tasks. Similarly, Universal Self-consistency (Chen et al., 373 2024a), which allows LLMs to choose the most consistent answers from a set of candidates, has 374 small improvements over Vanilla. These observations indicate that self-reflection and conventional 375 prompting techniques may not sufficiently enhance LLMs' ability to follow constraints, highlighting the potential need for external tools to assist LLMs in this area. (iii) Tools are helpful: ReAct 376 and CRITIC can be viewed as two variants in our framework. In contrast, CRITIC offers more 377 granular feedback by specifically identifying which constraint within an instruction is not satisfied.



	Mistral-7B		Llama	3-8B	Llama3	3.1-8B	Llama 3	.1-70B
Constraint Type	Vanilla	DVR	Vanilla	DVR	Vanilla	DVR	Vanilla	DVR
Detectable Content	76.36	88.90	84.18	96.81	86.29	96.31	97.06	98.59
Keywords	76.04	84.23	83.84	88.32	84.94	88.77	87.88	92.03
Punctuation	24.34	72.93	91.03	95.64	97.01	98.04	98.38	98.39
Case Change	70.08	81.28	81.28	93.38	82.97	90.71	80.23	96.28
Start End	81.29	90.41	84.88	90.03	84.07	91.92	89.37	96.46
Detectable Format	69.59	80.70	81.57	89.23	84.69	92.31	90.52	95.30
Language	69.11	81.80	77.06	88.38	81.96	89.76	90.83	95.26
Length Constraints	50.42	73.23	65.29	80.85	68.55	83.57	80.50	90.20

Table 4: Comparison Across Different Constraints Types



Figure 4: Ablation study on Mistral-7B, Llama3.1-8B and Llama3-8B.

73% satisfaction rate. The reason might be that Mistral 7B tends to ignore this constraint and needs external feedback to correct the answer.

4.4 RQ3: CONTRIBUTION OF INDIVIDUAL MODULES

In addition to comparing our method with existing baselines, we conduct an ablation study to assess the effectiveness of individual modules. Specifically, we investigate three variants: (i) w/o Detailed Feedback: The detailed feedback from the tool is removed but the refinement repository is kept to provide relevant few-shot examples showing the responses before and after refinement. Here, the refinement repository is empty in the beginning (coldsart). (i) w/o Repository: The refinement repository is removed, and only 5 fixed examples are used for the self-refine process. (i) w/o both: The refinement repository and detailed feedback are all removed. Tools only give whether the whole instruction is satisfied. Figure 4 shows performance gaps between each method and the Vanilla. Both detailed feedback and the refinement repository are crucial. Without the repository, performance gains are limited, as fixed few-shot examples aren't optimal for each refinement target. Detailed feedback is also important for LLMs, because it locates the error and provides the direc-tion for LLMs to modify their responses. Llama3-8B and Llama3.1-8B show higher gains on more difficult instructions. Mistral-7B's gains are modest, because of the limited capacity of Mistral-7B in following complex instructions, it can only follow 6.3% of level 6 instructions (shown Table 2). Despite this, the gains are notable from its low starting point. In conclusion, both detailed feedback and the refinement repository contribute to the performance of our framework.

4.5 Hyper-Parameter Sensitivity Analysis

We also conduct a hyper-parameter sensitivity analysis of our framework, testing different numbers of refinement few-shots and trials for successful refinement on Llama3.1-8B. As shown in Figure 5, performance improves with more trials but saturates at five, with minimal gains beyond that. Similarly, increasing the few-shot examples boosts performance in the beginning. The performance saturates after 8 shots. On the CoDI dataset, performance improves rapidly with initial increases in trials and examples, indicating that the first few numbers of trials and few-shot examples are most effective for refining the response.



Figure 5: Parameter study on ComplexInstruct (first two charts) and CoDI dataset (last two charts).

Table 5: LLMs Performance on Tool Selection (%)

Models	Hamming Loss	Accuracy	Precision	Recall	F1 Score
Mistral-7B	4.13	52.85	92.98	81.39	86.80
Llama3-8B	2.64	67.60	94.39	89.48	91.87
Llama3.1-8B	2.90	61.77	94.69	87.50	90.95
Llama3.1-70B	0.86	86.40	98.41	96.38	97.38

4.6 TOOL SELECTION ACCURACY

507 Correctly decomposing and selecting tools are essential for feedback and refinement. We define tool 508 selection as a multi-label prediction task for LLMs, evaluated using hamming score, accuracy, preci-509 sion, recall, and F1-score. The total number of tools is 21. Results are shown in Table 5. Hamming loss, which measures the fraction of incorrect labels, is low across all models, indicating minimal 510 mispredictions. Every model demonstrates a very high precision score, meaning that the tools they 511 select are mostly correct, avoiding misleading feedback with incorrect tool selection. Accuracy, 512 which measures the exact match between the selected tools and the ground truth, is the strictest 513 metric. Despite this, all models achieve over 50% accuracy. Considering the limited performance 514 of these models on constraint-following tasks, tool selection is a relatively easier task for LLMs. 515 This performance gap makes it possible for our method to provide reliable feedback, collect past 516 refinement examples and be effective in improving LLMs' constraint-following ability. 517

518

495 496

504 505

506

519

4.7 WOULD DVR AFFECT COMPREHENSIBILITY AND FLUENCY OF RESPONSES

In this subsection, we investigate if our framework would sacrifice comprehensibility and fluency
in order to follow complex-constraints. We focus on evaluating key metrics such as readability,
perplexity, and coherence. These metrics assess the comprehensibility and fluency of the responses.
Results on ComplexInstruct (Table 11) and CoDI (Table 12) are in Appendix A.6. They both show
that our framework has performance comparable to those of Vanilla, indicating that it does not
degrade fluency and readability. The reason is that our method does not change any weights in
LLMs, which maintains their ability in generating fluent and comprehensible text.

527 528

529

5 CONCLUSION

530 In conclusion, this paper presents the Divide-Verify-Refine (DVR) framework to enhance LLMs' 531 ability to follow multi-constraint instructions. There are three steps in our framework: (1) Di-532 vide complex instructions into single constraints and assign appropriate tools for each constraint. 533 (2) Verify: To deal with the feedback quality problem, these tools rigorously verify the response 534 and generate reliable feedback. (3) Refine: To tackle the constraint diversity challenge, we design 535 the refinement repository to store successful refinement processes, allowing LLMs to retrieve and 536 learn from past examples. Our framework improves LLMs' adherence to complex multi-constraint 537 instructions without the need for retraining, offering a scalable solution to enhance the practical usability of LLMs in real-world applications. Additionally, we construct a new dataset free from 538 hidden or conflicting constraints, providing a more comprehensive and accurate evaluation of LLM performance on multi-constraint instructions.

540	REFERENCES
541	

550

551 552

553

554

555

559

560

561

565

566

567

- Dimosthenis Antypas, Asahi Ushio, Jose Camacho-Collados, Vitor Silva, Leonardo Neves, and 542 Francesco Barbieri. Twitter topic classification. In Proceedings of the 29th International Confer-543 ence on Computational Linguistics, 2022. 544
- Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, 546 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are 547 few-shot learners. In Proceedings of the 34th International Conference on Neural Information 548 Processing Systems, pp. 1877–1901, 2020.
 - Tianle Cai, Xuezhi Wang, Tengyu Ma, Xinyun Chen, and Denny Zhou. Large language models as tool makers. In The Twelfth International Conference on Learning Representations, 2024.
 - Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W Cohen. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. arXiv preprint arXiv:2211.12588, 2022.
- Xiang Chen and Xiaojun Wan. A comprehensive evaluation of constrained text generation for large 556 language models. arXiv preprint arXiv:2310.16343, 2023.
- 558 Xinyun Chen, Renat Aksitov, Uri Alon, Jie Ren, Kefan Xiao, Pengcheng Yin, Sushant Prakash, Charles Sutton, Xuezhi Wang, and Denny Zhou. Universal self-consistency for large language models. In ICML 2024 Workshop on In-Context Learning, 2024a.
- 562 Yihan Chen, Benfeng Xu, Quan Wang, Yi Liu, and Zhendong Mao. Benchmarking large language models on controllable generation under diversified instructions. In Proceedings of the AAAI 563 Conference on Artificial Intelligence, volume 38, pp. 17808–17816, 2024b. 564
 - Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. arXiv preprint arXiv:2110.14168, 2021.
- Shehzaad Dhuliawala, Mojtaba Komeili, Jing Xu, Roberta Raileanu, Xian Li, Asli Celikyilmaz, 569 and Jason Weston. Chain-of-verification reduces hallucination in large language models. arXiv 570 preprint arXiv:2309.11495, 2023. 571
- 572 Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B Hashimoto. Length-controlled al-573 pacaeval: A simple way to debias automatic evaluators. Conference on Language Modeling, 574 2024. 575
- 576 Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. Pal: Program-aided language models. In International Conference on Machine 577 Learning, pp. 10764–10799. PMLR, 2023. 578
- 579 Zhibin Gou, Zhihong Shao, Yeyun Gong, Yujiu Yang, Nan Duan, Weizhu Chen, et al. Critic: Large 580 language models can self-correct with tool-interactive critiquing. In The Twelfth International 581 Conference on Learning Representations, 2024. 582
- 583 Qianyu He, Jie Zeng, Qianxi He, Jiaqing Liang, and Yanghua Xiao. From complex to simple: Enhancing multi-constraint complex instruction following ability of large language models. arXiv 584 preprint arXiv:2404.15846, 2024a. 585
- 586 Qianyu He, Jie Zeng, Wenhao Huang, Lina Chen, Jin Xiao, Qianxi He, Xunzhe Zhou, Jiaqing Liang, 587 and Yanghua Xiao. Can large language models understand real-world complex instructions? In 588 Proceedings of the AAAI Conference on Artificial Intelligence, volume 38, pp. 18188–18196, 589 2024b. 590
- Ruixin Hong, Hongming Zhang, Xinyu Pang, Dong Yu, and Changshui Zhang. A closer look at the self-verification abilities of large language models in logical reasoning. In Proceedings of the 592 2024 Conference of the North American Chapter of the Association for Computational Linguistics (NAACL 2024), 2024.

- ⁵⁹⁴ Dong Huang, Qingwen Bu, Jie M Zhang, Michael Luck, and Heming Cui. Agentcoder: Multi-agentbased code generation with iterative testing and optimisation. arXiv preprint arXiv:2312.13010, 2023.
- Jie Huang, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song, and Denny Zhou. Large language models cannot self-correct reasoning yet. In *The Twelfth International Conference on Learning Representations*, 2024.
- Albert Qiaochu Jiang, Sean Welleck, Jin Peng Zhou, Timothee Lacroix, Jiacheng Liu, Wenda Li,
 Mateja Jamnik, Guillaume Lample, and Yuhuai Wu. Draft, sketch, and prove: Guiding formal
 theorem provers with informal proofs. In *The Eleventh International Conference on Learning Representations*, 2023.
- Dongwei Jiang, Jingyu Zhang, Orion Weller, Nathaniel Weir, Benjamin Van Durme, and Daniel Khashabi. Self-[in] correct: Llms struggle with refining self-generated responses. *arXiv preprint arXiv:2404.04298*, 2024a.
- Yuxin Jiang, Yufei Wang, Xingshan Zeng, Wanjun Zhong, Liangyou Li, Fei Mi, Lifeng Shang, Xin Jiang, Qun Liu, and Wei Wang. Followbench: A multi-level fine-grained constraints following
 benchmark for large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL 2024)*, 2024b.
- Ryo Kamoi, Sarkar Snigdha Sarathi Das, Renze Lou, Jihyun Janice Ahn, Yilun Zhao, Xiaoxin Lu, Nan Zhang, Yusen Zhang, Ranran Haoran Zhang, Sujeeth Reddy Vummanthala, et al. Evaluating llms at detecting errors in llm responses. In *Proceedings of the 2024 Conference on Language Modeling (COLM 2024)*, 2024.
- 617
 618
 619
 620
 619
 619
 620
 619
 619
 619
 619
 610
 610
 611
 612
 612
 613
 614
 615
 615
 616
 617
 617
 618
 619
 619
 619
 610
 610
 611
 612
 612
 613
 614
 614
 615
 615
 616
 617
 618
 619
 618
 619
 618
 619
 618
 619
 618
 619
 618
 619
 618
 619
 618
 619
 618
 619
 618
 619
 618
 619
 620
 620
 620
 620
 620
 620
 620
 620
 620
 621
 621
 622
 622
 622
 623
 624
 625
 626
 627
 628
 629
 620
 620
 620
 620
 620
 620
 620
 620
 620
 620
 620
 620
 620
 620
 620
 620
 620
 620
 620
 620
 620
 620
 620
 620
 620
 620
 620
 620
 620
 620
 620
 620
 620
 620
 620
 620
 620
 620
- Ming Li, Han Chen, Chenguang Wang, Dang Nguyen, Dianqi Li, and Tianyi Zhou. Ruler: Improv ing llm controllability by rule-based data recycling. *arXiv preprint arXiv:2406.15938*, 2024.
- Daniel Loureiro, Francesco Barbieri, Leonardo Neves, Luis Espinosa Anke, and Jose Camacho-Collados. Timelms: Diachronic language models from twitter. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pp. 251–260, 2022.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri
 Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement
 with self-feedback. *Advances in Neural Information Processing Systems*, 36, 2024.

632

633

- Norman Mu, Sarah Chen, Zifan Wang, Sizhe Chen, David Karamardian, Lulwa Aljeraisy, Dan Hendrycks, and David Wagner. Can llms follow simple rules? *arXiv preprint arXiv:2311.04235*, 2023.
- Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. Webgpt: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021.
- Theo X Olausson, Jeevana Priya Inala, Chenglong Wang, Jianfeng Gao, and Armando Solar Lezama. Is self-repair a silver bullet for code generation? In *The Twelfth International Conference on Learning Representations*, 2024.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35: 27730–27744, 2022.
- Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi
 Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, et al. Code llama: Open foundation models for code. arXiv preprint arXiv:2308.12950, 2023.

648	Swarnadeep Saha, Omer Levy, Asli Celikyilmaz, Mohit Bansal, Jason Weston, and Xian Li. Branch-
649	solve-merge improves large language model evaluation and generation. In Proceedings of the
650	2024 Conference of the North American Chapter of the Association for Computational Linguis-
651	tics: Human Language Technologies (Volume 1: Long Papers), pp. 8345–8363, 2024.
652	William Soundars, Cathering Val. Loff Wu, Stavan Bills, Long Ouwang, Jonathan Ward, and Jan
653	Leike Self critiquing models for assisting human evaluators arXiv preprint arXiv:2206.05802
654	2022
655	2022.
656	Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion:
657	Language agents with verbal reinforcement learning. Advances in Neural Information Processing
658	Systems, 36, 2024.
659	Haoran Sun Lixin Liu Junije Li Fengyu Wang Baohua Dong Ran Lin and Ruohui Huang
660	Conifer: Improving complex constrained instruction-following ability of large language models.
661	arXiv preprint arXiv:2404.02823, 2024.
662	
663	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-
664	lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-
665	tion and fine-tuned chat models. arXiv preprint arXiv:2307.09288, 2023.
666	Gladys Tyen Hassan Mansoor Victor Cărbune Yuanzhu Peter Chen and Tony Mak Llms cannot
667	find reasoning errors but can correct them given the error location. In <i>Findings of the Association</i>
668	for Computational Linguistics ACL 2024, pp. 13894–13908, 2024.
669	Joi compinancian 2m8monos nez 2021, pp. 10091 10900, 2021
670	Fei Wang, Chao Shang, Sarthak Jain, Shuai Wang, Qiang Ning, Bonan Min, Vittorio Castelli, Yas-
671	sine Benajiba, and Dan Roth. From instructions to constraints: Language model alignment with
672	automatic constraint verification. arXiv preprint arXiv:2403.06326, 2024.
673	Xuezhi Wang Jason Wei Dale Schuurmans, Quoc V Le, Ed H Chi, Sharan Narang, Aakanksha
674	Chowdhery and Denny Zhou Self-consistency improves chain of thought reasoning in language
675	models. In The Eleventh International Conference on Learning Representations, 2023.
676	
677	Bosi Wen, Pei Ke, Xiaotao Gu, Lindong Wu, Hao Huang, Jinfeng Zhou, Wenchuang Li, Binxin
678	Hu, Wendy Gao, Jiaxin Xu, et al. Benchmarking complex instruction-following with multiple
679	constraints composition. arXiv preprint arXiv:2407.03978, 2024.
680	Shunyu Yao Jeffrey Zhao Dian Yu Nan Du Izhak Shafran Karthik R Narasimhan and Yuan
681	Cao. React: Synergizing reasoning and acting in language models. In <i>The Eleventh International</i>
682	Conference on Learning Representations, 2023.
683	
684	Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou,
685	and Le Hou. Instruction-following evaluation for large language models. In The Twelfth Interna-
696	tional Conference on Learning Representations, 2023.
607	Oibao Zhu Dava Guo Zhihong Shao Deijan Yang Peivi Wang Runxin Xu Y Wu Yukun Li
600	Huazuo Gao, Shirong Ma, et al. Deepseek-coder-v2: Breaking the barrier of closed-source models
600	in code intelligence. arXiv preprint arXiv:2406.11931. 2024.
600	σ r
690	
091	
092	
093	
694	
695	
696	
697	
698	
699	
700	
701	

702 A APPENDIX

A.1 ALGORITHM

The	the Language Model MA Januar Justing Stratest
Inpu	A: Language Model \mathcal{M} , input instructions A , looiset I
Solor	Jul: Response set <i>I</i>
1 1	
1:1	initialize the refinement repository $Q = \{\}$
2: 1	Of $I \in A$ up Concrete initial responses: $P = M(n = I)$
5: 4.	Unitialize the toolset for instruction $I: T_{-} = \int (p_{generate}, T)$.
4. 5.	Decompose constraints: $M(n, -1) \rightarrow \{a\}$
5:	becompose constraints. $\mathcal{M}(p_{decomp}, 1) \rightarrow \{C_i\}_{i=1,2,3}$
0. 7.	$\mathcal{M}(n, t, c) \rightarrow t$ where $t \in T$
7. 8.	$M(p_{select}, c) \neq t$, where $t \in T$. M sets parameters for t
9. 9.	$T_T = T_T \cup t$
10:	$r_1 = r_1 \otimes v$.
11:	$R = R_0, a = n_0$
12:	while $a > 0$ do
13:	a = a - 1
14:	Verify and get feedback from tools: $F_I = \{f_i\}_{i=1,2,3}$, where $f_i = t(R)$.
15:	if $f = True, \forall f \in F_I$ then
16:	return R
17:	end if
18:	Retrieve few-shot examples: $s^t = \{(I_i, R_i, f_i, R'_i)^t\}_{i=1,2,3}$, where $s^t \subseteq Q$
19:	Refine: $R' = \mathcal{M}(p_{refine}, s^t, I, R, f)$, where $f \in F_i$ and $f \neq True$.
20:	if $t(R') = True$ then
21:	Save the refinement process: $Q = Q \cup \{(I, R, f, R')^t\}$
22:	Update the current response: $R = R'$
23:	a = n
24:	end if
25:	end while
26:	$Y = Y \cup R$
27: 6	end for

A.2 BASELINE DETAILS

738 739

740

741

742

743

744

745

746

747

748

749

750

751

- Reflexion (Shinn et al., 2024): This method allows LLMs to self-reflect on their own responses and provide valuable feedback for future outputs. With the feedback, LLMs will refine their responses.
- Branch-solve-Merge (BSM) (Saha et al., 2024): BSM uses a "Divide and Conquer" approach to break complex instructions as individual branches. Then the LLMs will merge the responses from branches as the final answer. Similarly, in our experiment, we use LLMs to generate a response for each single constraint and then merge them together.
- Universal Self-Consistency (U-SC) (Chen et al., 2024a): This study extends the idea of Self-Consistency (Wang et al., 2023) to free-form generation. It first generates several candidate responses and then asks LLMs to select the most consistent one.
- Rejection Sampling (Saunders et al., 2022): Since we have tools for reliable verification, the most simple method is to select the best one from a set of responses. Here, we give the maximum number of trials as 5.
- ReAct (Yao et al., 2023): In ReAct, LLMs take actions based on the observation of the environment. Here, we adopt this method by letting the tools as the environment and giving LLMs boolean signals indicating whether the generated response adheres to all constraints in the instruction.

756 757 758	• CRITIC (Gou et al., 2024): CRITIC uses tools for verification such as using external A for evaluating the toxic score of the response or code executor for checking code generation. We adopt this method into our scenario, where tools will pinpoint which constraint of the second se	PI n. ne
759	instruction is not satisfied.	
761	A.3 COMPLEXINSTRUCT	
762		
763	We have 21 types of constraints which can be divided into 8 general categories (Zhou et al., 2023)):
764	• Keywords:	
765	(1) Include keyword.	
767	(2) Include keyword at least/less than certain frequency,	
768	(3) Forbidden word,	
769	(4) At least/less than certain frequency of letters.	
770	• Length:	
771	(1) At least/less than certain number of words,	
772	(2) At least/less than certain number of sentences, (2) Funct number of new graphs	
773	(5) Exact number of paragraphs.	
774	• Detectable Content:	
776	(2) Exact number of placeholders.	
777	Detectable Format:	
778	(1) Number of bullet points.	
779	(2) Add title,	
780	(3) Answer from options,	
781	(4) Minimum of highlighted sections,	
782	(5) Json format.	
783	• Change Cases:	
784	(1) All uppercase, (2) All lowercase	
785	(3) At least/less than certain number of all-capital words.	
700	• Startend:	
788	(1) End the text with a certain sentence,	
789	(2) Wrap whole response in double quotation.	
790	• Punctuation:	
791	(1) No commas in response.	
792	• Language:	
793	(1) Respond with certain language.	
794		
795	Table 6: Instruction Examples (6 levels)	
790	Instruction (level 1):	_
798	Write something with a topic of film or ty or video. At the end of your response please	
799	explicitly add a postscript starting with P.S.	
800	Constraint Type:	-
801	(1) Detectable Content - postscript	_
802	Instruction (level 2): Produce a brief writing piece with emphasis on culture. The ensure should be in all leverences	
803	letters with no capitalization. Response must also contain exactly 3 bullet points in markdown	
804	format. Use * to indicate bullets, like:	
805	* xyz.	
806	* abc.	
808	* opq.	
809		_
-		

	Table 6: Instruction Examples (continued)
C	onstraint Type:
(1)) Change Cases - All lowercase
(2) Detectable Format - Number of bullet points
In	struction (level 3):
Pr	oduce a text with a focus on food. Your answer must have a title contained in double ang
br	ackets, such as «title». Make sure you don't use any commas. Make sure to include the wo
'ir	noredients'
11	
Co	onstraint Type:
(1)) Detectable Format - Add title
\hat{o}) Punctuation - No commas in response
(3)) Keywords - Include keyword
	regulate here here here here here here here he
In	struction (level 4):
W	rite a text with a technology theme, be sure the letter 'l' appears at least 7 times in y
re	sponse. Your entire response should be in all lowercase letters (no capital letters whatsoey
Tł	he word "artificial" should not appear in your response. make sure the response has less t
82	words.
52	
C	onstraint Type:
(1) Keywords - At least/less than certain frequency of letters
\hat{o}) Change Cases - All lowercase
(3) Keywords - Forbidden word
(4) Length: At least/less than certain number of words
· ·	
In	struction (level 5):
A	ssist me in writing on a scientific topic. Highlight at least 3 text sections. i.e. *highlight
se	ction*. Mention the word "research" for less than 4 times, make sure that words with
ся	pital letters appear less than 2 times. The number of sentences in your response should
le	ss than 5. Include a title wrapped in double angular brackets, i.e. «title».
Co	onstraint Type:
(1) Detectable Format - Minimum of highlighted sections
(2) Keywords - Include keyword at least/less than certain frequency
(3)) Keywords - At least/less than certain frequency of letters
(4) Length: At least/less than certain number of sentences
(5) Detectable Format - Add title
	,
In	struction (level 6):
W	rite a short description of entrepreneurs. The very end of your entire response should r
ex	actly like: Is there anything else I can help with? The total number of words in your response
sh	ould be at least 23. Include a title wrapped in double angular brackets, i.e. «title». Sepa
vc	our esponse into 2 parts, where each part is separated with *** Your answer must con
ex	actly 4 bullet point in Markdown using the following format
*	Rullet noint one
* 1	Bullet point two
	Bunot point two.
 D	on't forget to include the keywords risk-taking
	en e respecto mende die Reprotos fisik undits.

864	Table 6: Instruction Examples (continued)
865	
866	Constraint Type:
867	(1) Startend - End the text with a certain sentence
868	(2) Length: At least/less than certain number of words
869	(3) Detectable Format - Add title
870	(4) Length: Exact number of paragraphs
074	(5) Detectable Format - Number of bullet points
0/1	(6) Keywords: Include keyword
872	

A.4 DETAILED EXPERIMENTS

We can observe that our methods consistently outperform baselines on differen LLMs.

Table 7: Performance of methods across levels 1 to 6 (Llama-3-8B-Instruct)

Method	Level 1	Level 2	Level 3	Level 4	Level 5	Level 6
Vanilla	89.1	71.7	56.4	44.2	30.4	20.4
Reflxion	88.8	72.1	57.5	41.5	30.0	20.9
BSM	89.2	71.9	56.0	41.3	28.8	19.5
U-SC	89.5	71.8	56.7	45.1	31.2	20.6
Rejection Sampling	90.8	80.9	64.5	52.9	39.8	31.0
ReAct	93.6	81.3	68.8	54.4	39.1	30.7
CRITIC	94.1	85.8	74.4	61.2	51.1	41.5
$\overline{DVR}(\overline{coldstart})$	95.0	86.7	76.9	65.3	53.7 -	43.8
DVR(warmstart)	95.4	87.6	77.0	67.4	55.4	46.9

Table 8: Performance of methods across levels 1 to 6 (Llama-3.1-70B-Instruct-AWQ-INT4)

Method	Level 1	Level 2	Level 3	Level 4	Level 5	Level 6
Vanilla	95.5	83.7	72.4	63.2	51.3	35.9
Reflexion	95.3	83.5	72.8	63.0	51.6	36.1
BSM	95.0	84.5	72.6	64.8	49.6	34.2
U-SC	96.0	83.3	71.8	64.0	52.5	36.3
Rejection Sampling	97.3	90.8	83.2	72.2	63.8	50.7
ReAct	97.5	91.0	83.5	72.4	65.0	52.1
CRITIC	98.1	93.2	87.6	79.1	73.7	61.3
$\overline{DVR}(\overline{coldstart})$	98.0	94.3	- 88.2 -		75.7	63.1
DVR(warmstart)	98.2	94.6	88.7	82.2	76.0	64.2

Table 9: Comparison Across Different Constraints Types (coldstart)

	Mistral-7B		Llama3-8B		Llama3.1-8B		Llama 3.1-70B	
Constraint Type	Vanilla	DVR	Vanilla	DVR	Vanilla	DVR	Vanilla	DVR
Detectable Content	76.36	88.49	84.18	95.82	86.29	96.19	97.06	98.14
Keywords	76.04	84.32	83.84	87.53	84.94	88.77	87.88	92.05
Punctuation	24.34	71.31	91.03	95.47	97.01	98.29	98.38	98.38
Case Change	70.08	80.15	81.28	93.20	82.97	89.96	80.23	96.24
Start End	81.29	88.71	84.88	88.71	84.07	91.78	89.37	95.94
Detectable Format	69.59	81.30	81.57	89.29	84.69	92.38	90.52	95.56
Language	69.11	82.72	77.06	86.24	81.96	89.30	90.83	94.34
Length Constraints	50.42	72.61	65.29	79.66	68.55	83.93	80.50	90.17

918 A.5 LLM SELF-VERIFY ABILITY

We evaluate LLMs on their ability to verify whether the responses meet the given constraints. As
shown in Table 10, LLMs often fail to verify their outputs accurately, suggesting they lack the
capacity to provide reliable feedback or reflection on their own responses.

Table	10: LLMs Sel	f-verification A	Accuracy (%)
Model	Mistral-7B	Llama3-8B	Llama3.1-8B
	53.1	56.8	55.7

A.6 FLUENCY AND READABILITY

Table 11: Descriptive Statistics of Responses (ComplexInstruct), where Coherence_or1 is first order coherence and Coherence_2 is the second order coherence.

Model	Method	Readability ↑	Perplexity \downarrow	Coherence_or1 ↑	Coherence_or2↑
mistral7D	Vanilla	62.24	18.22	0.61	0.59
lillsu al / D	DVR	61.93	18.95	0.59	0.57
llomo 2 PD	Vanilla	63.77	18.49	0.57	0.57
nama5-6D	DVR	63.58	18.08	0.59	0.56
llomo 2 1 9D	Vanilla	63.43	19.68	0.62	0.61
Ilaina5.1-6D	DVR	62.75	18.30	0.62	0.60
llama 2 1 70D	Vanilla	61.96	17.99	0.64	0.63
nama5.1-70B	DVR	63.51	18.04	0.63	0.62

Table 12: Descriptive Statistics of Responses (CoDI)

		1		1 ()	
Model	Method	Readability ↑	Perplexity \downarrow	Coherence_or1 ↑	Coherence_or2 ↑
mistral7D	Vanilla	63.62	15.27	0.82	0.81
iiiisu ai / D	DVR	63.53	15.98	0.83	0.82
11	Vanilla	63.79	14.25	0.79	0.76
nama5-8D	DVR	64.14	16.07	0.80	0.77
11	Vanilla	62.18	16.27	0.81	0.83
fiama5.1-8B	DVR	62.02	17.58	0.81	0.80

A.7 EXPERIMENTS ON IFEVAL

We conduct experiments on IFEval (Zhou et al., 2023) which is an instruction-following benchmark widely used for industry. The IFEval dataset evaluates the instruction-following ability and is one of the core benchmarks used in the Open LLM Leaderboard (Hugging Face). We conduct experiments on Mistral-7B-v0.3 and the results are shown in Table 13. DVR outperforms all other baselines on IFEval benchmark.

		Tab	le 13: I	SR (%)	for IFEval Dataset			
Method	Vanilla	Reflexion	BSM	U-SC	Rejection Sample	ReAct	CRITIC	DVR
ISR	47.32	47.13	47.87	46.95	53.23	53.97	55.53	60.44

A.8 EXPERIMENTS ON GPT4-TURBO

971 We conduct experiments on GPT-4-turbo. Shown in Table 14, we can observe that GPT-4-turbo performs better than open-source models (Mistral and Llama). Surprisingly, applied on Llama3.1-

973	source model.	 		 · ···· F · ····	r	
974						
975			~ .			

8B, DVR can still outperform GPT-4-turbo, indicating that DVR exploits the potential of the open-

Table	14: Perfor	mance Co	mparison t	o GPT4-tu	rbo	
Model	Level 1	Level 2	Level 3	Level 4	Level 5	Level 6
Mistral-7B	77.0	55.3	34.1	19.9	12.4	6.3
DVR (Mistral-7B)	95.0	81.3	66.6	51.4	36.4	23.4
Llama3.1-8B	90.5	76.6	62.5	50.1	35.6	25.3
DVR (Llama3.1-8B)	95.2	88.7	79.2	69.7	60.5	49.6
GPT-4-turbo	95.3	88.4	78.8	65.2	53.7	42.6

A.9 ROBUSTNESS OF DVR

We also conduct experiments to assess DVR's performance in the presence of tool errors. Two types of errors are introduced: random noise and systematic bias. Specifically, we evaluate the framework on instructions with length constraints, using 600 samples (from ComplexInstruct) for word count control and another 600 for sentence count control. A constraint example: "The response needs to be less than (or at least) x number of words/sentences." where x ranges from 10 to 100 for words and 3 to 5 for sentences. We add two types of noises to tools:

Noise: Gaussian noise with a mean of 0 is added to the counted number of words (or sentences) to simulate random errors. The DVR's performance is then measured across different deviation levels.

Bias Errors: A fixed bias is added to the counted values of words (or sentences) to introduce systematic errors. The tables below demonstrate DVR's performance under different bias values.

Observations: We have several observations in Table 15 and Table 16. (1) The performance will decrease as the noise levels (deviation, bias values) increase. (2) As the errors become large, the performance degradation will saturate. (3) Overall, DVR will not perform much worse than vanilla even if the bias and errors are large (20 for word count and 4 for sentence count). (4) The impact of noise on the overall instruction satisfaction rate is less severe compared to its influence on specific constraints.

Table 15: Satisfaction Data for Word Number Constraints (%)

0	5	10	20	Vanilla
88.00	87.17	82.83	81.50	68.16
48.17	45.00	43.67	43.67	10.17
0	5	10	20	Vanilla
88.00	87.17	84.33	82.67	68.16
48.17	47.83	47.00	45.67	10.17
or Senter 0	nce Num	ber Cons	straints (%	%)
	-	4	4	Vanilla
74.50	68.17	62.50	4 60.50	Vanilla 56.33
74.50 42.83	68.17 38.17	62.50 35.33	4 60.50 34.33	Vanill 56.33 10.17
74.50 42.83 0	68.17 38.17 1	62.50 35.33 2	4 60.50 34.33 4	Vanill 56.33 10.17 Vanill
74.50 42.83 0 74.50	68.17 38.17 1 64.67	62.50 35.33 2 58.67	4 60.50 34.33 4 56.50	Vanill 56.33 10.17 Vanill 56.33
	0 88.00 48.17 0 88.00 48.17 or Senter 0	0 5 88.00 87.17 48.17 45.00 0 5 88.00 87.17 48.17 47.83 or Sentence Num 0 1	0 5 10 88.00 87.17 82.83 48.17 45.00 43.67 0 5 10 88.00 87.17 84.33 48.17 47.83 47.00	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$

A.10 COMPUTATION TIME

We conducted experiments with 20 instructions, each containing 6 constraints, using Mistral-7B.
 The number of trials was set to 5, consistent with the paper's settings. The average running time is summarized below:

1029 1030

1038

1040

	Table	17: The Av	erage R	Running	Time for One Sa	mple		
Method	Vanilla	Reflexion	U-SC	BSM	Rejection Sample	ReAct	CRITIC	DVR
Running Time (s)	5.91	20.53	41.34	46.21	32.32	36.48	37.98	33.91

As shown in Table 17, our method does not exhibit significantly higher running time compared to other baselines. Considering the performance gains (Table 2), our method demonstrates a balance between efficiency and effectiveness.

1039 A.11 LIMITATIONS AND FUTURE WORKS

1041 There are several limitations and potential future works.

(1) Currently, we consider multiple independent constraints. However, the instructions in real-world might be more complex and constraints might have dependency with each other (Wen et al., 2024).
For example, the instruction can ask the response has 4 bullet points and 2 sentences in each bullet point. In such scenario, simple decomposition might not work. (2) Moreover, tools may not be available for new constraints. Here, we assume that we have tools for all existing constraints. However, users' requirements can be very diverse and we would not have certain tools for new constraints. (3) As shown in Table 15 and Table 16, DVR performance declines when tools produce errors.

1049 As a result, we can consider the following future works. (1) The first one is to handle constraints 1050 with dependency. If LLMs can correctly disentangle these constraints and find appropriate tools, our framework can also enhance LLMs' instruction-following ability in such a complex setting. (2) 1051 The second future work is to deal with new constraints. One possible solution is to firstly find out 1052 if the constraint can be tested via python code. Then local LLMs can summarize constraints into 1053 code requirements and send them to advanced code generators for tool generation. Generated tools 1054 can also be saved locally for future use. (3) Ensuring the reliability and quality of tools is crucial 1055 for real-world applications. There are techniques such as developing test cases that verify a tool's 1056 outputs against known values (Cai et al., 2024; Huang et al., 2023) or generating multiple versions 1057 of a tool and selecting the most consistent output. Additional research is needed to evaluate whether 1058 these techniques are sufficient to guarantee tool quality. Furthermore, new approaches could be 1059 explored to enhance tool generation. 1060

1061 A.12 CASE STUDY

1063 We use some examples to show how LLMs modify responses with the guidance of feedback. The 1064 red colored content is the errors from the original response and blue colored content is the corrected 1065 parts after refinement.

Table 18: Refinement Examples

1066 1067 1068

1062

Provide me with text related to the everyday experience. make sure the response has less than 80 words. Respond with at least 3 sentences. The entire reply should contain a title in double angular brackets, i.e. «title». Your entire response should be in all lowercase letters (no capital letters whatsoever).

1073 Response:

1074 «morning routines»

Prompt:

Unsatisfied Constraint: make sure the response has less than 80 words.

<sup>i love starting my day with a quiet cup of coffee and a good book. it helps me set a positive tone for the rest of the day. morning routines can be as simple or as elaborate as you like, but they help create a sense of structure and calm in our otherwise chaotic lives. [name]'s morning routine involves a quick workout followed by a healthy breakfast. [another name] prefers to meditate for a few minutes before beginning their day.
We is a few minutes before beginning their day.</sup>

	Table 18: Refinement Examples (continued)
Fe Tl Pl	eedback: he response contains 83 words, which is 4 too many. lease remove 4 words to have fewer than 80 words
P P	efined Response: "morning routines"
i 1 to of fc	love starting my day with a quiet cup of coffee and a good book. it helps me set a pone. morning routines can be as simple or as elaborate as you like, but they help create f structure and calm in our chaotic lives. [name]'s morning routine involves a quick v blowed by a healthy breakfast
10	
In	istruction:
E: ar	xcuse me, could you please provide text that focuses on the music-topic? use the key trists. highlight at least 1 key point by wrapping it with *. For example: *highlight
po	bint*. The number of sentences in your response should be less than 5. be sure the
ap	prears at least 4 times in your response. Make sure to include at least 5 placeholder represented by square brackets, such as [address] [name], the total number of words in your re-
sh	rould be at least 23.
D	
K	esponse: Music Unites» From festivals to online communities, music continues to bridge cultur
g	eographical gaps like never before. *Electronic music has especially seen a surge in pon
ar	mong the younger generation [*] . Talented artists from all over the world have been c
to	gether to create beautiful melodies and harmonies. [name] is one such artist who ha
m	haking waves in the industry with his eclectic beats. His music can be found on v
pl	atforms, including [address] and [website]. Many emerging artists, including [another
ar	e also experimenting with different genres.
	isatistied Constraint: The number of sentences in your response should be less than 5.
T	he response contains 6 sentences which is 2 more than allowed
P	lease remove at least 2 sentences.
R	efined Response:
«I	Music Unites» From festivals to online communities, music continues to bridge cultur
ge	eographical gaps like never before. *Electronic music has especially seen a surge in pop
ar	mong the younger generation [*] . Talented artists from all over the world have been c
to m	gether to create beautiful metodies and narmonies. [name] is one such artist who ha
pl	latforms, including [address] and [website].
1	
In	istruction:
K	individual generate a few texts on the subject of education or learning. Give a title wrap
u(th	ouble angular brackets, such as «your fulle». Add sucess words which are capitalized.
se	entence of your response should be "Hope you agree with me.".
R	esponse:
«	The Power of Learning»
L	EARNING IS A LIFELONG PROCESS that enables individuals to acquire new ski
kı	nowledge. It plays a CRUCIAL role in shaping one's career and personal growth. Edu
he	eips to broaden one's perspective and fosters critical thinking. Hope you agree with me
U	Insatisfied Constraint: Add stress words which are capitalized. Ensure those stress
aŗ	ppear less than 4 times.
E	eedback: The response contains 6 capitalized words, which is more than the allowed 3
1.0	
Pl	lease remove at least 3 capitalized word(s).

K(«]	enned Response: The Power of Learning»
L	earning is a lifelong process that enables individuals to acquire new skills and knowl
pl	ays a crucial role in shaping one's career and personal growth. Education is VERY ir
an	id helps to broaden one's perspective. Hope you agree with me.
In	struction: Produce a text about TV. Your entire response should be in all capital lette
re	sponse should contain at least 3 sentences. Your answer must have a title contained in
so	igular brackets, such as (address) [name] Highlight at least 2 sections of your resi
m	arkdown such as *highlighted section*.
D	
N «]	ESPONSE.
T	V SHOWS HAVE BEEN A STAPLE OF ENTERTAINMENT FOR DECADES.
C	LASSICS LIKE Ï LOVE LUCYT̈́O MODERN HITS LIKE ĜAME OF THRONES
IS	SOMETHING FOR EVERYONE. *TV PROVIDES A WAY FOR US TO ESCAP
D. 11	IFFEKENT WOKLDS AND EXPERIENCE STORIES THAT WE MIGHT NOT (JISE ENCOUNTED* WITH THE DISE OF STDEAMING SEDVICES LIVE NU
Δ	ISE ENCOUNTER [*] . WITH THE RISE OF STREAMING SERVICES LIKE NI ND HUI II IT IS FASIER THAN EVER TO ACCESS VOUR FAVORITE SHOWS
л Tl	HE COMFORT OF YOUR OWN HOME AT [address] *WHETHER YOU ARE
0	F DRAMA COMEDY OR REALITY TV THERE IS NO SHORTAGE OF OPTIC
C	HOOSE FROM*. FANS OF [name] CAN LOOK FORWARD TO NEW CONTENT
Tl	HEIR FAVORITE CREATORS AND DISCOVER NEW SHOWS FROM TALENT
A	ND-COMING WRITERS AND PRODUCERS. WITH SO MANY GREAT SHO
	HOOSE FROM IT IS NO WONDER THAT TV REMAINS A POPULAR FORM
11 V	ERIAINMENT FOR PEOPLE ALL OVER THE WORLD. [Channel] OFFERS A A DIETY OF DOOD AMMING THAT CATEDS TO DIFFERD TASTES AND
E	RENCES.
U	nsatisfied Constraint: Your entire response should be in all capital letters.
гe Pl	ease capitalize all of them
R	efined Response:
«]	IV SHOWS FOR THE AGES»
Т	V SHOWS HAVE BEEN A STAPLE OF ENTERTAINMENT FOR DECADES.
C	LASSICS LIKE Ï LOVE LUCYT̈́O MODERN HITS LIKE ĠAME OF THRONES
IS	SOMETHING FOR EVERYONE. *TV PROVIDES A WAY FOR US TO ESCAP
D	IFFERENT WORLDS AND EXPERIENCE STORIES THAT WE MIGHT NOT
El	RWISE ENCOUNTER*. WITH THE RISE OF STREAMING SERVICES LIK
FI	LIX AND HULU IT IS EASIER THAN EVER TO ACCESS YOUR FAVORITE S
	XOM THE COMFORT OF YOUR OWN HOME AT [ADDRESS]. * WHETHER YO
	IONS TO CHOOSE EROM* EANS OF INAMEL CAN LOOK FORWARD TO
C	ONTENT FROM THEIR FAVORITE CREATORS AND DISCOVER NEW SHOWS
T/	ALENTED UP-AND-COMING WRITERS AND PRODUCERS. WITH SO MANY
SI	HOWS TO CHOOSE FROM IT IS NO WONDER THAT TV REMAINS A PO
F	ORM OF ENTERTAINMENT FOR PEOPLE ALL OVER THE WORLD. [CHANNI
FI	ERS A WIDE VARIETY OF PROGRAMMING THAT CATERS TO DIFFERENT T
11	

Instruction:	
Generate a w	riting piece related to health. organize your entire response in at least 3 sentences.
Answer shou	ld be written with all letters lowercased. In your entire response, the letter n should
appear less th	nan 6 times.
Response:	
maintaining	a healthy lifestyle involves making conscious choices about the food we eat and
the activities	we engage in. eating a balanced diet that includes plenty of fruits and vegetables
can help to b	boost our energy levels and support overall well-being. regular exercise, such as
walking or ye	oga, can also have a positive impact on both our physical and mental health.
Unsatisfied C	Constraint: In your entire response, the letter n should appear less than 6 times.
Feedback: T	he response contains the letter 'n' 23 times, which is 18 too many.
Please remov	re 18 occurrences to have less than 6.
Refined Resp	ponse:
eating a heal	thy diet is good for us. it helps to keep us strong. we should try to do it.
<u> </u>	
Instruction: (The dear one, would you be so kind as to grace me with a magnificent text delying
into the ever-	so-mysterious world of science? The very end of your entire response should read
exactly like:	Let me know if you have additional questions. All letters in your response must be
lower case le	tters. Separate your response into 4 parts, where each part is separated with ***
use the keyw	ords physics. Highlight at least 3 text sections, i.e. *highlighted section* there
should be at	least 4 occurrences of the word innovations in your response.
Response:	
«exploring th	e mysteries of science»
*** *the wo	nders of physics* have long been a source of fascination for scientists and the
general publi	ic alike. from the behavior of subatomic particles to the vast expanses of space.
physics gove	rns the fundamental laws of our universe.
*** *innovat	ions in physics* have led to numerous breakthroughs in fields such as medicine.
technology, a	ind engineering, for instance, advances in medical imaging have allowed for more
accurate diag	moses and targeted treatments.
	,
*** *the inte	rsection of physics and technology* has given rise to exciting innovations such as
quantum con	puting, nanotechnology, and renewable energy solutions. these innovations have
the potential	to transform industries and improve lives worldwide.
r · · · · · · · · ·	····· · · · · · · · · · · · · · · · ·
*** as scient	ists continue to push the boundaries of human knowledge, we can expect to see
even more g	oundbreaking innovations in the field of physics. let me know if you have addi-
tional question	ns.
Unsatisfied (Constraint: Separate your response into 4 parts, where each part is separated with
***.	real france france france france, more cash part to separate a mar
· Feedback· F·	spected exactly 4 sections, but found 5
Please remov	re 1 section by removing 1 '***'
Refined Rest	onse: «exploring the mysteries of science»
the wonders	of physics* have long been a source of fascination for scientists and the general
nublic alike	from the behavior of substamic particles to the vast expanses of space physics
governe the f	indamental laws of our universe
governs me I	undamentai iaws of our universe.

innovations	n physics* have led to numerous breakthroughs in fields such as medicing, tech
nology and	in physics - nave red to numerous oreakunoughs in neuos such as incurrent, leen-
noiogy, and	engineering. for instance, advances in medical imaging have allowed for more
accurate diag	moses and targeted treatments.

Table 18: Refinement Examples (continued)

	Table 18: Refinement Examples (continued)
**	**
th qu ha ar	e intersection of physics and technology* has given rise to exciting *innovations* su nantum computing, nanotechnology, and renewable energy solutions. these *innova we the potential to transform industries and improve lives worldwide. and these *innova e just the beginning.
** ev ad	** as scientists continue to push the boundaries of human knowledge, we can expect then more groundbreaking *innovations* in the field of physics. let me know if you additional questions.
T	
In w do at	struction: Provide me with a text that revolves around the theme of family. do not so ord "relationships" in the response. Highlight at least 2 sections of your response in own such as *highlighted section*. Make sure to include a postscript starting with P.S. least 5 sentences in your response. make sure the response has less than 56 words.
an bu *	nswer must contain exactly 2 bullet points in the markdown format (use * to indicate allet) such as: This is the first point
*	This is the second point.
Ro su fa	esponse: A family is a bond that can never be broken. *Children grow up with lov port from their family*. As they grow older they begin to appreciate the sacrifices mily has made for them
*	A family can provide a sense of belonging.
*	A family can provide a sense of security.
P.	S. I love my family.
U *ł	nsatisfied Constraint: Highlight at least 2 sections of your response in markdown su highlighted section*.
Fe Pl	eedback: The response contains 1 highlighted section. ease include 1 more highlighted section in the format of *highlighted*.
M lo th	odified Response: A family is a *bond that can never be broken*. *Children grow up ve and support from their family*. As they grow older they begin to appreciate the sacr eir family has made for them.
*	A family can provide a sense of belonging.
* P.	A family can provide a sense of security. S. I love my family.
In	struction: Produce a comprehensive written discourse analyzing the multiple facets of
A	nswer in lowercase letters only, throughout your entire answer. There should be exact
pa	argraphs in your response, separated by the markdown divider: ***. Response mus
сс *	ontain exactly 5 bullet points in markdown format. Use * to indicate bullets, like:
*	abc
*	opq
R	echonse.

1297	
1298	«the art of storytelling in film»
1299	storytelling is a fundamental aspect of filmmaking. it is the art of conveying emotions and
1300	ideas through visual narratives. films have the power to transport us to different worlds, evoke
1301	emotions, and make us think. whether it's a romantic comedy or a sci-fi epic, storytelling is the
1302	core of what makes a film engaging.
1303	*** Flan is a sellah anation and isan that manines the involution of a surround individuals from
1304	nim is a conaborative medium that requires the involvement of numerous individuals, from screenwriters to cinematographere, the director plays a crucial role in bringing the story to life
1305	guiding the actors and crew to create a cohesive vision the cinematographer's job is to capture
1306	the visuals that will transport the audience to the world of the story.
1307	***
1308	* some notable film movements include:
1309	* german expressionism
1310	* french new wave
1311	* italian neorealism
1312	each of these movements brought a unique perspective to filmmaking, experimenting with style
1313	and narrative structure, they paved the way for future generations of filmmakers, allowing them to push the boundaries of what is possible on screep
1314	***
1315	ultimately, film is a powerful medium that can inspire, educate, and entertain, it has the abil-
1316	ity to transcend cultural boundaries, allowing us to connect with people from different back-
1317	grounds and experiences. whether you're a filmmaker or a film enthusiast, the world of cinema
1318	has something to offer everyone.
1319	
1320	Unsatisfied Constraint: Response must also contain exactly 3 bullet points in markdown format.
1321	Use * to indicate bullets, like:
1322	* XYZ
1323	* abc
1324	opq
1325	Feedback:
1326	In the response, there are 4 bullet points.
1327	Here are the bullet points detected:
1328	* some notable film movements include:
1329	* german expressionism
1330	* Irench new wave
1332	Please remove exactly 1 bullet points to meet the requirement of 3
1333	
1334	
1335	
1336	
1337	
1338	
1339	
1340	
1341	
1342	
1343	
1344	
1345	
1346	
1347	
1348	
1349	

Table 18: Refinement Examples (continued)

	Table 18: Refinement Examples (continued)
Refined Respons	se:
«the art of story	telling in film»
storytelling is a	fundamental aspect of filmmaking. it is the art of conveying emotions and
ideas through vi	sual narratives. films have the power to transport us to different worlds, evoke
emotions, and m	ake us think. whether it's a romantic comedy or a sci-fi epic, storytelling is the
core of what ma	kes a film engaging.

nim is a collabo	prative medium that requires the involvement of numerous individuals, from
guiding the actor	rs and crew to create a cohesive vision, the cinematographer's job is to capture
the visuals that y	will transport the audience to the world of the story
***	will dansport the addicate to the world of the story.
* some notable f	film movements include:
* french new wa	we
* italian neoreal	ism
each of these mo	ovements brought a unique perspective to filmmaking, experimenting with style
and narrative str	ucture. they paved the way for future generations of filmmakers, allowing them
to push the bour	idaries of what is possible on screen.

ultimately, film	is a powerful medium that can inspire, educate, and entertain. it has the abil-
ity to transcend	cultural boundaries, allowing us to connect with people from different back-
grounds and exp	eriences. whether you're a filmmaker or a film enthusiast, the world of cinema
has something to	o offer everyone.
Instruction: Gene	erate a text that involves anger and is related to films:
Perpense: The d	liractor's blatant disregard for the script and the talent on set is nothing sho
astounding They	here the audacity to claim artistic vision while disregarding the very actors
hired to bring the	ir mediocre ideas to life. It's a travesty that someone so clearly out of touch
call themselves a	filmmaker. The industry should be ashamed for allowing such incompetence
flourish.	i minimuken. The mausity should be ashunde for anowing such meompetene
Feedback: The s	sentiment of the text is 'disgust', which does not match the required sentin
'anger'.	
Please adjust the	sentiment of the text to be more 'anger'.
Refined Response	2.
The director's flag	grant disregard for the script and the talent on set is infuriating. They have the n
to claim artistic v	vision while callously disregarding the very actors they hired to bring their su
ideas to life. It's	outrageous that someone so clearly out of touch can call themselves a filmma
I ne industry show	uid be neid accountable for allowing such incompetence to perpetuate. It's not
a minor mistake;	It is a systemic issue that demands attention. The fact that this director contribution is a clop in the face to avail talented attention when here here we have a large talented attention.
underapprovide	The situation is unaccentable, and it's time for change
The final answer	. The situation is unacceptable, and it's unit for change.
ment of multiple	nis. There is no one final answer to uns provided. The provided code and responses of writing to fit the constraints provided. The provided code and responses
modifications der	process or writing to in the constraints provided. The provided code and responses in the provided code and responses in the required tone for each inst
tion	nonstrate now to aujust the sentiment to match the required tone for each list
Instruction: Gene	erate text on science with disgust, using general words:

	Table 19: Refinement Examples on CoDI (continued)
Res creat as it field	ponse: The discoveries in the realm of biology have been grotesque, with the revelation of tures that defy natural explanation. Their twisted forms and behaviors evoke a sense of revulsion, f the very fabric of life has been warped in ways that are both fascinating and repulsive. It's a l that challenges our understanding of what it means to be alive.
Fee expo Plea	dback: The detected topic of the response is learning and education, which does not match the ected topic science and technology. use adjust the content to align more closely with the topic science and technology.
Refi The gani sens und	Ined Response: latest breakthroughs in genetic engineering have been abhorrent, with the development of or- isms that are unnatural and unsettling. Their abnormal structures and functions elicit a strong se of disgust, as if the boundaries of science have been overstepped. It's a field that tests our erstanding of the fundamental principles of life and the ethics of scientific inquiry.
13	Prompts
	Table 20: Prompts
Der	ampt for Decrement Conception
Ge	neration prompt
Yo	u are an AI assistant that generates responses based on given prompts.
For	r each prompt, provide a response that adheres to the specified constraints.
#P	rompt: Generate a text that touch on arts. make sure the response has less than 40 words.
ma	ke sure it contains at least 2 placeholders represented by square brackets, such as [name].
Re	sponse: "The Beauty of Expression"
Ar	t from [name]'s perspective transforms emotions into visual narratives [another name]'s
ma	sterpieces illustrate this beautifully.
#P	rompt: Generate sth. about celebrity or pop culture. the response must contain a title
wra tera you	apped in double angular brackets, i.e. «title». The answer should be in all lowercase let- s, with no capitalizations. The word ëntertainmentshould not appear in your response. In ar entire response, the letter m should appear at least 3 times. Your answer must have at least
2 p	laceholders, wrapped in square brackets, such as [author].
() #P	more examples)
Re	sponse:
110	
Pro	ompt for Instruction Decomposition

1458 1459	Table 20: Prompts (continued)
1460	You are an advanced assistant specializing in identifying and listing output constraints from
1461	provided instructions. The instructions typically include a task related to generating content
1462	on a specific topic and one (or multiple) format constraint(s). Your goal is to focus only on
1463	extracting and listing all the format constraints required for the output, ignoring the content-
1464	related task.
1465	Instruction:
1466	ISON format. The word 'show' should not appear in your response. Use square brackets
1467	for placeholders like [username1] [username2] Please include at least 2 placeholders in the
1468	thread. You are not allowed to use any commas in your response.
1469	Format Constraints:
1470	#1. Put your entire answer in JSON format.
1471	#2. The word 'show' should not appear in your response.
1472	#3. Use square brackets for placeholders, like [username1], [username2]. Please include at
1473	least 2 placeholders in the thread.
1474	#4. You are not allowed to use any commas in your response.
1475	Instruction:
1476	{current instruction}
1477	Format Constraints:
1478	
1479	Prompt for Tool Selection
1480	You will be given a list of constraints. Each constraint belongs to a specific category. Your task
1481	is to recognize and categorize each constraint. Only output the category from the following
1482	options:
1483	postscript, placeholder, include keyword, exclude keyword, letter frequency, keyword fre-
1484	responses highlighted ISON format title format quoted response and phrase no commas all
1485	capital letters, all lowercase, capital word frequency language restriction
1486	Please ensure to categorize each constraint accurately according to its description. There is
1487	definitely a valid category option for each constraint. You can Here are examples for each type
1400	of constraint:
1409	Prompt: Make sure to include the word 'mutations'.
1490	Category: include keyword
1491	Prompt: Limit the number of words you use to fewer than 65 words.
1493	(more examples)
1494	Prompt: {Current Prompt}
1495	Category:
1496	
1497	Prompt for Refinement
1498	
1499	
1500	
1501	
1502	
1503	
1504	
1505	
1506	
1507	
1508	
1509	
1510	
1511	

	Table 20: Prompts (continued)
	You are an AI assistant responsible for refining a given response. Given a prompt, its original
	response, and the analysis of the response, your task is to modify the response according to the
	analysis. #Dromety I'm locking for tout that avalance arts or culture, can you assist? places avaliaitly
	"#Prompt: 1 in looking for text that explores arts of culture, can you assist? please explicitly add a note starting with P.S. There should be available available of culture, can your response, separated
	by the markdown divider: *** Make sure to include at least 2 placeholder represented by
	square brackets such as [address] [name] Highlight at least 2 text sections i.e. *highlighted
	section*. There should be no commas in your reply.
	#Original Response: Art has the power to bring people together and transcend cultural bound-
	aries. It can evoke emotions and spark conversations that might not be possible through other
	means. *At the [address] museum, visitors can experience this firsthand by exploring the di-
	verse collection of art from around the world.*
	*** From paintings to sculptures to installations, each piece tells a unique story that can be
	interpreted in many ways. *The work of [name] is a great example of this, as it challenges
	viewers to think critically about the world around them.* Whether you're an art enthusiast or
	just looking for a new perspective, the [address] museum is a must-visit destination. P.S. Don't
	forget to check out the museum's events calendar for upcoming exhibitions and performances!
	#It does not satisfy the constraint: There should be no commas in your reply.
	#Analysis: The response contains 4 comma(s). Here are the detected commas: (museum, visitors) (tallations, each) (of this, as it) (perspective, the laddress)
	Please remove all commas
	#Modified Response: Art has the power to bring people together and transcend cultural bound-
	aries. It can evoke emotions and spark conversations that might not be possible through other
	means. *At the [address] museum visitors can experience this firsthand by exploring the diverse
	collection of art from around the world.*
	*** From paintings to sculptures to installations each piece tells a unique story that can be
	interpreted in many ways. *The work of [name] is a great example of this as it challenges
	viewers to think critically about the world around them.* Whether you're an art enthusiast or
	just looking for a new perspective the [address] museum is a must-visit destination. P.S. Don't
	forget to check out the museum's events calendar for upcoming exhibitions and performances!
	"(more examples)
	#Original Response: current response
	#It does not satisfy the constraint: current unsatisfied constraint
	#Analysis: current feedback
	#Modified Response:
	A 14 FAILURE CASES
	A.14 TAILUKE CASES
,	We analyze these failure cases and find out two typical situations that our framework would fail:
	- **

(1) Combined Effect: In certain scenarios, multiple constraints can create a challenge for the LLMs 1553 to balance. As shown in the first example in the Table 21, in the case of limiting the frequency 1554 of the letter 't' while generating content about sports, the natural occurrence of the letter 't' in 1555 'sports', combined with 2 sections requirements, makes it difficult to satisfy 't' letter frequency 1556 limit. Although these constraints are not inherently contradictory, their combined effect puts the 1557 model under pressure, resulting in failure to meet some requirements. (2) Incorrect decomposition 1558 and wrong tool selection: Incorrect decomposition would lead to wrong tool selection. As shown 1559 in the second example in the Table 21, the constraint to add capitalized stress words fewer than 1560 four times was mistakenly decomposed into two separate tasks. This leads to the incorrect selection 1561 of tools, taking the "stress words" as keywords instead of capitalized words. This type of failure degrades the overall performance in following complex instructions. 1562

- 1563
- 1564

Instruction: Generate text on the subject of sports. Make sure to highlight at least 2 section your answer with markdown, i.e. use *highlighted section*. Please also include the keywo "games" in the response. You are not allowed to use any commas in your response. Put response into 2 sections, separated using 3 asterisks ***. Please make sure each point hat title wrapped in double angular brackets, i.e. «title». In your entire response, the letter t she appear less than 10 times. «Sports Overview» Participating in sports is an excellent way to improve physical and mental health *Sports basketball and soccer require quick reflexes and agility* Playing sports also fosters teamw and communication among players. *** «The Thrill of Games» Competitive games bring out the best in athletes *The rush of adrenaline during a crucial mu is exhilarating* Fans enjoy the suspense and excitement of watching their favorite teams ba it out. Unsatisfied Constraint: In your entire response, the letter t should appear less than 10 time: Discussion: There are other constraints against this one. The topic is about sports which contains the let t. Additionally, "Put the response into 2 sections" tends to make the response long, wich mu satisfying letter frequency limit challenging. Instruction: Can you provide me with some information about dining? Add stress words w are capitalized. Ensure those stress words appear less than 4 times. Response: «The Dining Experience» A DELICIOUS meal is not just about the food, but also about the atmosphere and serv When dining out, it's essential to consider the AMBIANCE of the restaurant, as it can gre impact your overall experience. A restaurant with a cozy and intimate setting can make f ROMANTIC evening, while a lively and bustling atmosphere can be PERFECT for a night with friends. Furthermore, the service can make or break your dining experience. A friendly and atten wait staff can make you feel welcome and valued, while poor service can leave a sour tast your mouth. Ultimately, a great d		Table 21: Failure Cases
«Sports Overview» Participating in sports is an excellent way to improve physical and mental health *Sports basketball and soccer require quick reflexes and agility* Playing sports also fosters teamw and communication among players. *** «The Thrill of Games» Competitive games bring out the best in athletes *The rush of adrenaline during a crucial ma is exhilarating* Fans enjoy the suspense and excitement of watching their favorite teams ba it out. Unsatisfied Constraint: In your entire response, the letter t should appear less than 10 times Discussion: There are other constraints against this one. The topic is about sports which contains the le t. Additionally, "Put the response into 2 sections" tends to make the response long, wich ma satisfying letter frequency limit challenging. Instruction: Can you provide me with some information about dining? Add stress words wh are capitalized. Ensure those stress words appear less than 4 times. A DELICIOUS meal is not just about the food, but also about the atmosphere and serv When dining out, it's essential to consider the AMBIANCE of the restaurant, as it can gre impact your overall experience. A restaurant with a cozy and intimate setting can make f ROMANTIC evening, while a lively and bustling atmosphere can be PERFECT for a night with friends. Furthermore, the service can make or break your dining experience. A friendly and atten wait staff can make you feel welcome and valued, while poor service can leave a sour tast your mouth. Ultimately, a great dining experience is one that combines DELICIOUS food with exceptio service and a pleasant atmosphere.	Instruc your ar "games respons title wr appear	tion: Generate text on the subject of sports. Make sure to highlight at least 2 sections iswer with markdown, i.e. use *highlighted section*. Please also include the keywo s" in the response. You are not allowed to use any commas in your response. Put se into 2 sections, separated using 3 asterisks ***. Please make sure each point hav apped in double angular brackets, i.e. «title». In your entire response, the letter t sho less than 10 times.
 «The Thrill of Games» Competitive games bring out the best in athletes *The rush of adrenaline during a crucial mains exhilarating* Fans enjoy the suspense and excitement of watching their favorite teams being out. Unsatisfied Constraint: In your entire response, the letter t should appear less than 10 times Discussion: There are other constraints against this one. The topic is about sports which contains the let. Additionally, "Put the response into 2 sections" tends to make the response long, wich mastifying letter frequency limit challenging. Instruction: Can you provide me with some information about dining? Add stress words whare capitalized. Ensure those stress words appear less than 4 times. Response: «The Dining Experience» A DELICIOUS meal is not just about the food, but also about the atmosphere and serv When dining out, it's essential to consider the AMBIANCE of the restaurant, as it can gre impact your overall experience. A restaurant with a cozy and intimate setting can make f ROMANTIC evening, while a lively and bustling atmosphere can be PERFECT for a night with friends. Furthermore, the service can make or break your dining experience. A friendly and atten wait staff can make you feel welcome and valued, while poor service can leave a sour tast your mouth. Ultimately, a great dining experience is one that combines DELICIOUS food with exception service and a pleasant atmosphere. Unsatisfied Constraint: Add stress words which are capitalized. Ensure those stress words appear less than 4 times. 	«Sports Particip basketh and con ***	S Overview» bating in sports is an excellent way to improve physical and mental health *Sports l ball and soccer require quick reflexes and agility* Playing sports also fosters teamw mmunication among players.
Unsatisfied Constraint: In your entire response, the letter t should appear less than 10 times Discussion: There are other constraints against this one. The topic is about sports which contains the let t. Additionally, "Put the response into 2 sections" tends to make the response long, wich ma satisfying letter frequency limit challenging. Instruction: Can you provide me with some information about dining? Add stress words wh are capitalized. Ensure those stress words appear less than 4 times. Response: «The Dining Experience» A DELICIOUS meal is not just about the food, but also about the atmosphere and serv When dining out, it's essential to consider the AMBIANCE of the restaurant, as it can gre impact your overall experience. A restaurant with a cozy and intimate setting can make f ROMANTIC evening, while a lively and bustling atmosphere can be PERFECT for a night with friends. Furthermore, the service can make or break your dining experience. A friendly and atten wait staff can make you feel welcome and valued, while poor service can leave a sour tast your mouth. Ultimately, a great dining experience is one that combines DELICIOUS food with exceptio service and a pleasant atmosphere. Unsatisfied Constraint: Add stress words which are capitalized. Ensure those stress wo appear less than 4 times.	«The T Compe is exhil it out.	'hrill of Games» titive games bring out the best in athletes *The rush of adrenaline during a crucial ma arating* Fans enjoy the suspense and excitement of watching their favorite teams ba
Discussion: There are other constraints against this one. The topic is about sports which contains the let t. Additionally, "Put the response into 2 sections" tends to make the response long, wich ma- satisfying letter frequency limit challenging. Instruction: Can you provide me with some information about dining? Add stress words wl are capitalized. Ensure those stress words appear less than 4 times. Response: «The Dining Experience» A DELICIOUS meal is not just about the food, but also about the atmosphere and serv When dining out, it's essential to consider the AMBIANCE of the restaurant, as it can gre impact your overall experience. A restaurant with a cozy and intimate setting can make f ROMANTIC evening, while a lively and bustling atmosphere can be PERFECT for a night with friends. Furthermore, the service can make or break your dining experience. A friendly and atten wait staff can make you feel welcome and valued, while poor service can leave a sour tast your mouth. Ultimately, a great dining experience is one that combines DELICIOUS food with exception service and a pleasant atmosphere. Unsatisfied Constraint: Add stress words which are capitalized. Ensure those stress words appear less than 4 times.	Unsatis	sfied Constraint: In your entire response, the letter t should appear less than 10 times
Instruction: Can you provide me with some information about dining? Add stress words whare capitalized. Ensure those stress words appear less than 4 times. Response: «The Dining Experience» A DELICIOUS meal is not just about the food, but also about the atmosphere and serv When dining out, it's essential to consider the AMBIANCE of the restaurant, as it can gre impact your overall experience. A restaurant with a cozy and intimate setting can make f ROMANTIC evening, while a lively and bustling atmosphere can be PERFECT for a night with friends. Furthermore, the service can make or break your dining experience. A friendly and atten wait staff can make you feel welcome and valued, while poor service can leave a sour tast your mouth. Ultimately, a great dining experience is one that combines DELICIOUS food with exception service and a pleasant atmosphere. Unsatisfied Constraint: Add stress words which are capitalized. Ensure those stress words appear less than 4 times.	Discus There a t. Addi satisfyi	sion: are other constraints against this one. The topic is about sports which contains the le tionally, "Put the response into 2 sections" tends to make the response long, wich main ng letter frequency limit challenging.
Response: «The Dining Experience» A DELICIOUS meal is not just about the food, but also about the atmosphere and serv When dining out, it's essential to consider the AMBIANCE of the restaurant, as it can gre impact your overall experience. A restaurant with a cozy and intimate setting can make f ROMANTIC evening, while a lively and bustling atmosphere can be PERFECT for a night with friends. Furthermore, the service can make or break your dining experience. A friendly and atten wait staff can make you feel welcome and valued, while poor service can leave a sour tast your mouth. Ultimately, a great dining experience is one that combines DELICIOUS food with exception service and a pleasant atmosphere. Unsatisfied Constraint: Add stress words which are capitalized. Ensure those stress we appear less than 4 times.	Instruc are cap	tion: Can you provide me with some information about dining? Add stress words wh italized. Ensure those stress words appear less than 4 times.
A DELICIOUS meal is not just about the food, but also about the atmosphere and serv When dining out, it's essential to consider the AMBIANCE of the restaurant, as it can gre impact your overall experience. A restaurant with a cozy and intimate setting can make f ROMANTIC evening, while a lively and bustling atmosphere can be PERFECT for a night with friends. Furthermore, the service can make or break your dining experience. A friendly and atten wait staff can make you feel welcome and valued, while poor service can leave a sour tast your mouth. Ultimately, a great dining experience is one that combines DELICIOUS food with exception service and a pleasant atmosphere. Unsatisfied Constraint: Add stress words which are capitalized. Ensure those stress we appear less than 4 times.	Respor «The D	ise: Dining Experience»
Furthermore, the service can make or break your dining experience. A friendly and atten wait staff can make you feel welcome and valued, while poor service can leave a sour tast your mouth. Ultimately, a great dining experience is one that combines DELICIOUS food with exception service and a pleasant atmosphere. Unsatisfied Constraint: Add stress words which are capitalized. Ensure those stress word appear less than 4 times.	A DEL When of impact ROMA with fri	JCIOUS meal is not just about the food, but also about the atmosphere and servidining out, it's essential to consider the AMBIANCE of the restaurant, as it can great your overall experience. A restaurant with a cozy and intimate setting can make for NTIC evening, while a lively and bustling atmosphere can be PERFECT for a night itends.
Ultimately, a great dining experience is one that combines DELICIOUS food with exception service and a pleasant atmosphere. Unsatisfied Constraint: Add stress words which are capitalized. Ensure those stress words appear less than 4 times.	Further wait sta your m	more, the service can make or break your dining experience. A friendly and attent aff can make you feel welcome and valued, while poor service can leave a sour taste outh.
Unsatisfied Constraint: Add stress words which are capitalized. Ensure those stress we appear less than 4 times.	Ultima service	tely, a great dining experience is one that combines DELICIOUS food with exception and a pleasant atmosphere.
	Unsatis appear	fied Constraint: Add stress words which are capitalized. Ensure those stress words than 4 times.

1621

Table 21: Failure Cases (continued)

1622 Discussion: In this example, the constraint "Add stress words which are capitalized. Ensure those stress 1623 words appear less than 4 times" is mistakenly decomposed into two separate constraints by 1624 the LLMs: (1) Add stress words which are capitalized. (2) Ensure those stress words appear 1625 less than 4 times. For the first constraint, the LLMs can not find a tool for it. For the second 1626 constraint, the LLM incorrectly selects the keyword tool: Keywords("stress words", "less than", 1627 4). In this case, the correct tool should be Capitalwords("less than",4). Usually, incorrect 1628 decomposition prevents LLMs from identifying any tools, rather than causing them to select an 1629 incorrect one. This example is unique because the constraint is phrased in two sentences, and the second sentence can be misinterpreted to suggest a limit on keywords (i.e., "stress words") instead of focusing on capitalized stress words. 1633 1634 A.15 TOOL EXAMPLES 1635 1. Word Counting: This example is obtained through GPT4-o with zero-shot. It demonstrates that 1637 reliable tools can be easily created. The details are as follows: 1638 1639 def feedback(response, max_words=None, min_words=None): # Count the number of words in the response 1640 word_count = len(response.split()) 1641 1642 # Check for maximum word constraint 1643 if max_words is not None and word_count > max_words: 1644 return f"Response failed because it has {word_count} words, 1645 exceeding the maximum allowed limit of {max_words} words." 1646 1647 # Check for minimum word constraint 1648 if min_words is not None and word_count < min_words: return f"Response failed because it has only {word_count} words, fewer than the minimum required {min_words} words." 1650 1651 # If all constraints are satisfied 1652 return True

Lowercase Letter Validation: This example validates whether a given text is entirely in lowercase.
 If any word contains uppercase letters, it provides feedback on which words need correction. The implementation is as follows:

```
class LowercaseLetter:
1658
          def __init__(self):
1659
              pass
          def feed_back(self, value):
              # Split the input string into words
1663
              words = value.split()
1664
              # Find words that are not fully in lowercase
              upper case words = [word for word in words if
              any(char.isupper() for char in word)]
              if value.islower():
1669
                  return True
              else:
1671
                  return f"The response contains words that are not in
1672
                  all lowercase letters: {', '.join(upper_case_words)}.
1673
                  Please lowercase all of them."
```