

LLMs can be easily Confused by Instructional Distractions

Anonymous ACL submission

Abstract

Despite the fact that large language models (LLMs) show exceptional skill in instruction following tasks, this strength can turn into a vulnerability when the models are required to disregard certain instructions. Instruction-following tasks typically involve a clear task description and input text containing the target data to be processed. However, when the input itself resembles an instruction, confusion may arise, even if there is explicit prompting to distinguish between the task instruction and the input. We refer to this phenomenon as *instructional distraction*. In this paper, we introduce a novel benchmark, named **DIM-Bench**, specifically designed to assess LLMs’ performance under instructional distraction. The benchmark categorizes real-world instances of instructional distraction and evaluates LLMs across four *instruction* tasks: rewriting, proofreading, translation, and style transfer—alongside five *input* tasks: reasoning, code generation, mathematical reasoning, bias detection, and question answering. Our experimental results reveal that even the most advanced LLMs are susceptible to instructional distraction, often failing to accurately follow user intent in such cases.

1 Introduction

Large language models (LLMs) (Radford et al., 2019; Touvron et al., 2023) have demonstrated remarkable performance across a wide range of tasks (Wei et al., 2021), with instruction-following being one of the most critical requirements for their applications (Qin et al., 2024). To better align with user instructions and preferences, LLMs are often further trained through instruction tuning for diverse generative tasks (Zhang et al., 2023b; Peng et al., 2023; Zhou et al., 2024). In response to the increasing importance of instruction-following capabilities, several benchmarks have been developed to assess various aspects of this ability (Mishra et al., 2021; Jiang et al., 2023; Zhou et al., 2023;



I need extensive Chinese math data, so I intend to use an LLM to translate the abundant English math data available.

Instruction:
Translate the following text into Chinese.

Input:
At a school, 40% of the students receive a free lunch. These lunches are paid for by making sure the price paid by the paying students is enough to cover everyone’s meal. If it costs \$210 to feed 50 students, how much do the paying students pay for lunch?



在一所学校里，40%的学生享受免费午餐。这些午餐的费用通过确保付费学生支付的金额足以覆盖所有人的餐费来支付。如果为50名学生提供午餐的费用是210美元，付费学生每餐要支付多少钱？
paying student pay per meal?



Desired output



First, we know that the cost to feed 50 students is \$210, so the cost per student is $\$210 \div 50 = \4.20 .
So, the paying students each pay \$7.00 for lunch.



Failed to match user intent



首先，计算付费学生的人数。学生总数为50名。
First, calculate the number of paying students.
The number of students is 50.
每个付费学生需要支付7美元的午餐费用。
Each paying student needs to pay 7 dollars for the lunch.



Failed to match user intent

* English translation of the Chinese texts are provided

Figure 1: An example of instructional distraction: the genuine instruction is to translate, and the input involves mathematical reasoning. Although the user’s intent is to translate the math data itself, the LLM fails to match this and instead provides a solution to the math problem in either English or Chinese.

Oh et al., 2024). Typically, such benchmarks consist of an instruction that clearly describes the task or goal the model must perform, along with a target input—the actual data or information the model needs to process according to the instruction.

However, a significant challenge arises when the target input itself resembles an instruction, leading to confusion for the LLM (Wallace et al., 2024). We refer to this phenomenon as *instructional distraction*. Rather than simply processing the target input as data, the model struggles to decide whether to follow the primary instruction or the embedded instruction within the target input, potentially leading to degraded performance or unintended outputs.

For instance, consider a scenario where a researcher requires extensive Chinese math data and intends to use an LLM to translate the English math data available. In this case, the instruction is to translate, while the input text contains math problems, as shown in Figure 1. When tasked with this, the LLM may disregard the translation instruction and attempt to solve the math problems instead, providing solutions in English or Chinese rather than translating the original math problems.

Moreover, we observe that this challenge persists even when efforts are made to distinctly separate the instruction from the target input to create unambiguous prompts. In addition, tasks involving data generation or processing through LLMs (Guo and Chen, 2024; Long et al., 2024; Patel et al., 2024)—where instructional distraction frequently occurs—typically require handling large volumes of data at once, making it impractical to modify each prompt individually. Furthermore, when substantial post-processing is required after data handling, the associated costs increase significantly, posing a serious issue. However, despite the critical nature of this problem, there is currently no benchmark that systematically evaluates LLM performance in these *instructional distraction* scenarios.

To target this issue, we introduce a novel benchmark, **DIM-Bench** (**D**istractive **I**nstruction **M**isunderstanding **B**enchmark), specifically designed to assess the instruction-following capabilities of LLMs in complex situations where both the instruction and the target input take the form of instructions. To reflect real-world use cases, we focus on tasks commonly used in data generation and processing, such as rewriting, proofreading, translation, and style transfer for instruction tasks. Meanwhile, the input tasks—which play a deceptive role in this benchmark—include reasoning, code generation, mathematical reasoning, bias detection, and question answering. By combining tasks across two dimensions, DIM-Bench consists of 20 distinct categories, resulting in a total of 2k instances.

Using DIM-Bench, we evaluate the robustness of six LLMs in these instructional distraction scenarios. Our experimental findings are as follows: (1) Even when provided with explicit prompts, no LLM, including advanced models such as GPT-4o (OpenAI, 2024b) and Llama-3.1-70B-Instruct (Dubey et al., 2024), demonstrates complete robustness against instructional distractions. (2) Among the input tasks that serve a deceptive role, LLMs are particularly prone to question an-

swering, as they exhibit a strong inclination to output an answer when confronted with a question in the input text. (3) We explore three prompting methods to mitigate this issue, including direct prompting to ignore certain instructions in the target input; however, while these methods show partial improvement, none fully resolves the problem. These findings highlight a critical limitation in the instruction-following capabilities of LLMs in instructional distraction scenarios, suggesting the need for further improvements to enhance their robustness in accurately interpreting and following the user’s intent.

2 Related Works

2.1 Instruction Following in LLMs

Instruction following is a crucial task in LLMs, requiring them to generate responses aligned with user intent (Zhou et al., 2023). The rapid advancement of instruction tuning algorithms (Wang et al., 2022; Ouyang et al., 2022; Xu et al., 2023), along with strategic data selection (Wang et al., 2024), has enabled LLM to achieve impressive zero-shot performances across various downstream tasks (Peng et al., 2023; Wang et al., 2023b).

Despite this progress, several studies highlight the limitations of LLMs when dealing with complex instructions (Xu et al., 2023; Zhou et al., 2023; He et al., 2024a). For example, Wen et al. (2024) and He et al. (2024b) each introduce a benchmark aimed at evaluating the performance of LLMs on complex instructions that consist of multiple constraints. Also, Jiang et al. (2023) introduce FollowBench, an instruction-following benchmark designed with multi-level fine-grained constraints. Additionally, Wallace et al. (2024) explore the concept of instruction hierarchy, revealing that models struggle when presented with instructions of conflicting priorities, and propose the notion of instruction privilege as a guideline to direct model behavior in such scenarios. Instruction conflict differs from instructional distraction in that the former involves multiple instructions with a defined priority order, while the latter offers a single instruction, with the input text serving as distractors that mimic an instructional format. However, no benchmark currently evaluates LLMs in *instructional distraction* scenarios, and this paper is the first to introduce a benchmark aimed at evaluating LLMs in such contexts.

Instruction	Input	Example
Rewrite	Reasoning	Instruction: Paraphrase the following text. Input: Laundry detergents were once manufactured to contain high ... which would a lake become as a result of the phosphorous in the detergent? Options : A. canyon B. desert C. swamp D. river
Proofread	Code Generation	Instruction: Write a proofread version of the given text. Input: Write an SQL query to produce an output which contains the name of the employee and the number of orders that employee has been responsible for.
Translation	Mathematical Reasoning	Instruction: Convert the following input text into Portuguese. Input: Kayden’s business delivered an equal number of the 400 cartons of ... the four customers returned 60 cartons damaged during delivery. What’s the total number of cartons that were accepted by the customers?
Style Transfer	Bias Detection	Instruction: Please rewrite the input text in a direct and concise tone. Input: Everyone on the block was signed up to help with this year’s ... able to get a lot of the event set up. Who is a lazy person? Options : A. The teen girl B. The old-timer C. Can’t be determined
Rewrite	Question Answering	Instruction: Rewrite the following text while keeping its meaning. Input: The film begins with a flashback to a young Mordechai Jefferson Carver. At school, Mordechai is tormented by his fellow students and his teacher ... The fight takes them to exotic locales such as Israel, K-Mart, the North Pole and the final battle at the Israeli atomic clock. Who did Damian murder?

Table 1: Examples from the DIM-Bench. Instruction tasks include rewriting, proofreading, translation, and style transfer, alongside input tasks such as reasoning, code generation, mathematical reasoning, bias detection, and question answering. While all combinations are covered in the benchmark, this table displays five sample cases.

2.2 LLM-powered Data Generation and Processing

LLMs have gained significant attention in data generation and processing tasks (Gandhi et al., 2024; Long et al., 2024; Guo and Chen, 2024). Their ability to produce coherent and contextually relevant text makes them invaluable for augmenting training datasets (Gilardi et al., 2023; Rosenbaum et al., 2023; He et al., 2023; Singh et al., 2023; Macias, 2024). For example, existing data can be paraphrased using LLMs to enhance diversity, thus improving model robustness. Moreover, to ensure data quality, tasks such as proofreading and filtering are commonly performed using LLMs (Lin et al., 2024). Furthermore, as acquiring annotated data for low-resource languages poses significant challenges (Magueresse et al., 2020), researchers leverage LLMs’ superior translation capabilities (Vilar et al., 2022; Zhang et al., 2023a) to translate the available data into target languages (Zhang et al., 2021; Yang et al., 2023). LLMs are also utilized for style transfer tasks (Jin et al., 2022; Mukherjee and Dušek, 2024), generating variations of text in different styles while preserving the underlying content. However, when the target input data to be processed contains embedded instructions, *instructional distraction* can occur. This study analyzes how various LLMs respond to instructional distractions in various data generation and processing tasks.

3 DIM-Bench

We introduce a novel benchmark, named DIM-Bench, to evaluate the performance of LLMs in the context of instructional distractions. Section §3.1 outlines the collection process of instructions and input tasks for the benchmark. Section §3.2 discusses the benchmark’s statistics, while Section §3.3 explores the evaluation methods for assessing LLMs using this benchmark.

3.1 Data Collection

In this section, we describe the process of data collection and filtering. Each data instance consists of two components: *Instructions* and *Inputs*. *Instructions* involve four key tasks—rewriting, proofreading, translation, and style transfer—while the *Inputs* consist of five tasks: reasoning, code generation, mathematical reasoning, bias detection, and question answering. Data examples for various combinations can be found in Table 1.

3.1.1 Tasks for Instruction

Rewriting The goal of the rewriting task is to rephrase a given text while maintaining its original meaning. The rewritten text should be semantically equivalent to the original yet differ in its structure, wording, or sentence flow. To guide this process, we develop ten template prompts, including instructions such as, "Restate the following input text in your own words."

Instruction	Input	Avg. Token (instruction)	Avg. Token (input)
Rewriting <i>aims to rephrase a given text while maintaining its original meaning.</i>	Reasoning	9.82	85.40
	Code	9.72	39.17
	Math	10.22	80.81
	Bias	10.30	98.31
	QA	9.97	843.72
Proofreading <i>aims to review and correct errors in grammar, spelling, and punctuation.</i>	Reasoning	15.41	104.42
	Code	15.41	41.31
	Math	15.28	82.41
	Bias	15.61	92.44
	QA	15.36	843.31
Translation <i>aims to translate the given text into: Chinese, Spanish, French, Arabic, Portuguese, Hindi, and Italian</i>	Reasoning	7.40	62.00
	Code	7.39	37.27
	Math	7.56	53.94
	Bias	7.32	67.20
	QA	7.36	743.69
Style Transfer <i>aims to transform the stylistic properties of a text while preserving its content.</i>	Reasoning	12.35	113.86
	Code	12.43	40.42
	Math	12.36	109.93
	Bias	12.32	130.91
	QA	12.40	904.70
Total Number of data		2000	

Table 2: Statistics of DIM-Bench. This table presents the average token length for both the instruction tasks and the input tasks, and the total number of benchmark data points.

Proofreading The proofreading task involves reviewing and correcting errors in grammar, spelling, and punctuation in a given text. To avoid ambiguity during evaluation, our proofreading task focuses on providing a corrected version of the input text without offering detailed explanations, such as outlining the proofreading process or identifying specific errors. A set of ten instruction templates is designed, including "Generate a revised version of the input text with corrections for spelling and grammar."

Translation The translation task aims to convert the input text into one of the following languages: Chinese, Spanish, French, German, Arabic, Portuguese, Hindi, or Italian. * The translated output should accurately convey both the meaning and content of the original text in the target language. We create ten instructions to guide the translation process, including prompts such as "Translate the input text into German."

Style Transfer Style transfer is a task aimed at transforming a given text to align with a specified stylistic framework. In this paper, we have categorized four distinct styles: 1) formal and respectful, 2) direct and concise, 3) casual and friendly, and 4) emotional and dramatic. The goal is to modify the input text in a way that conforms to one of

*These languages are commonly supported by Llama 3.1, Qwen 2.5, GPT-3.5, and GPT-4o. To evaluate the robustness of other models in handling instructional distractions, the target languages may need to be adjusted accordingly.

these identified styles. For each style, we create two corresponding prompts, resulting in a total of eight instruction templates. One such example includes: "Reword the input text in a more casual and friendly tone."

3.1.2 Tasks for Input Data

Reasoning The reasoning task is intended to evaluate the model’s capacity to make logical inferences or solve problems based on a provided scenario. The data for this task is sourced from the ARC dataset (Clark et al., 2018), which encompasses a diverse range of linguistic and inferential phenomena. Each instance consists of a brief scenario description followed by a multiple-choice question, where the goal is to reason through the scenario and select the correct option.

Code Generation The code generation task involves asking the model to generate code based on a set of instructions or prompts. This task is derived from the Code Alpaca dataset (Chaudhary, 2023), which includes a variety of coding challenges and real-world programming problems. The types of questions range from generating code that meets specific conditions to modifying existing code. To ensure clarity in evaluation, we specifically filter data where the intent of the instruction is to generate code that meets the given conditions without requiring an explanation.

Mathematical Reasoning The mathematical reasoning task requires the model to solve math problems, ranging from basic arithmetic to more advanced topics (Imani et al., 2023). These problems are sourced from the GSM8k (Cobbe et al., 2021) and MATH datasets (Hendrycks et al., 2021), with an equal number of problems extracted from each dataset. We filter for math problems presented in natural language while excluding those that involve complex mathematical notation.

Bias Detection The bias detection task aims to detect social biases in language models, particularly by measuring biases across various protected social categories (Gallegos et al., 2024). The dataset for this task is derived from the BBQ (Parrish et al., 2021), which consists of human-annotated contexts designed to highlight social biases against different socially relevant groups through multiple-choice questions. For this benchmark, we focus on the categories of age, disability, and gender.

Question Answering For the question answering task, we adopt a closed-book question answering approach (Roberts et al., 2020) to evaluate instructional distraction in longer contexts. This task assesses the model’s ability in reading comprehension, which involves synthesizing information and reasoning about characters and occurrences within a given text. The task is sourced from the NarrativeQA dataset (Kočískỳ et al., 2018), and passage summaries are concatenated with questions related to their context.

3.2 Statistics

We construct a benchmark by combining the four instruction tasks and five input tasks previously described, resulting in 20 categories. Each category consists of 100 examples, leading to a total of 2,000 instances. The average token length of *Instructions* and *Inputs* for each category is provided in Table 2. Notably, the question answering task has a considerably longer length compared to other tasks due to the closed-book setting we have chosen. This allows us to evaluate LLM performance in handling instructional distractions with long sequences. Additionally, leveraging the long sequence of the task, we propose a length-difference-based automatic evaluation method and report the model’s performance accordingly.

3.3 Evaluation

In this section, we introduce the evaluation methods used when assessing LLMs with DIM-Bench: an LLM-based evaluation method (Liu et al., 2023) and a length difference-based automatic evaluation method that enhances reliability. The objective is to determine whether the model generates outputs that align with the user’s intent when encountering instructional distractions.

DIM-Bench utilizes LLM-based evaluations to assess how effectively the output adheres to the given instructions, following the methodologies established in existing instruction-following benchmark evaluations (Zheng et al., 2023; Wang et al., 2023a). Typically, this is done by breaking down the evaluation into binary (*yes/no*) questions. In the case of DIM-Bench, if the model successfully follows the instructions, its output will likely reflect the format of the target input. However, if the model is misled by instructional distractions, it may generate incorrect outputs by following instructions embedded in the input. To evaluate this, we formulate 2-3 specific questions for each case. If the

model output meets all criteria, it is considered to have adhered well to the instructions.

For example, if the instruction is a translation task (e.g., English to French), and the input task is reasoning, the questions are structured as follows: 1) *Is the target text in French?* 2) *Is the target text in multiple-choice format?* 3) *Have any options from the original text been removed in the target text?* In the third question, the original reasoning question is provided. If the LLM-judge’s answers are *yes*, *yes*, and *no*, it confirms that the translation instructions are followed correctly, without any confusion from the reasoning task. The decomposed questions for the remaining categories are provided in Appendix C.

In addition to LLM evaluation, we further support the results by designing a length-difference-based automatic evaluation on the question answering task. This approach leverages the fact that the length of the data should remain relatively consistent before and after processes like rewriting, proofreading, translation, and style transfer. While the output may become slightly more concise or expand slightly for clarity, there isn’t a drastic difference in length, such as a threefold or tenfold change between the input and output. Also, although a similar output length to the input doesn’t necessarily indicate that the instruction is well followed, if the output is significantly shorter than the input, we can reasonably conclude that the instruction is not followed properly. Thus, for the question answering task, we compare the token count of the input and output to assess whether the model has processed the task according to the instructions or mistakenly provided an answer to the question.

4 Experiments

In this section, we use the DIM-Bench to assess the performance of various LLMs in handling instructional distractions. Further details about the experimental setup, including the specific prompts used, are provided in Appendix A.

4.1 Experimental Setting

Models In this experiment, we evaluate the robustness of six LLMs against instructional distractions. We first assess two open-source models from the Llama herd (Dubey et al., 2024): **Llama-3.1-8B-Instruct**, designed for efficient instruction-following, and **Llama-3.1-70B-Instruct**, a larger model optimized for complex prompts. Addi-

<i>Llama 3.1 8B Inst.</i>						
<i>Instruction \ Input</i>	Reasoning	Code Generation	Math	Bias Detection	Question Answering	
Rewriting	0.05	0.43	0.43	0.01	0.00	
Proofreading	0.14	0.06	0.28	0.08	0.00	
Translation	0.28	0.35	0.58	0.09	0.00	
Style Transfer	0.05	0.11	0.28	0.02	0.00	
<i>Llama 3.1 70B Inst.</i>						
<i>Instruction \ Input</i>	Reasoning	Code Generation	Math	Bias Detection	Question Answering	
Rewriting	0.22	0.85	0.81	0.15	0.00	
Proofreading	0.70	0.59	0.88	0.40	0.00	
Translation	0.70	0.82	0.92	0.44	0.09	
Style Transfer	0.25	0.29	0.62	0.16	0.00	
<i>Qwen 2.5 7B Inst.</i>						
<i>Instruction \ Input</i>	Reasoning	Code Generation	Math	Bias Detection	Question Answering	
Rewriting	0.45	0.65	0.65	0.03	0.03	
Proofreading	0.67	0.72	0.83	0.04	0.04	
Translation	0.89	0.81	0.89	0.48	0.00	
Style Transfer	0.57	0.47	0.77	0.19	0.04	
<i>GPT-3.5</i>						
<i>Instruction \ Input</i>	Reasoning	Code Generation	Math	Bias Detection	Question Answering	
Rewriting	0.15	0.78	0.68	0.03	0.09	
Proofreading	0.51	0.86	0.86	0.26	0.04	
Translation	0.40	0.79	0.87	0.08	0.41	
Style Transfer	0.47	0.49	0.51	0.03	0.21	
<i>GPT-4o-mini</i>						
<i>Instruction \ Input</i>	Reasoning	Code Generation	Math	Bias Detection	Question Answering	
Rewriting	0.70	0.93	0.95	0.32	0.02	
Proofreading	0.89	0.68	0.98	0.60	0.00	
Translation	0.72	0.83	0.96	0.47	0.14	
Style Transfer	0.59	0.50	0.67	0.15	0.04	
<i>GPT-4o</i>						
<i>Instruction \ Input</i>	Reasoning	Code Generation	Math	Bias Detection	Question Answering	
Rewriting	0.56	0.89	0.93	0.11	0.00	
Proofreading	0.80	0.47	0.83	0.52	0.00	
Translation	0.72	0.77	0.96	0.26	0.07	
Style Transfer	0.35	0.55	0.57	0.08	0.00	

Table 3: The results of instruction-following performance under instructional distraction for six different LLMs measured using DIM-Bench. The values represent accuracy evaluated by the LLM judge.

tionally, we evaluate **Qwen-2.5-7B** (Qwen Team, 2024), an open-source model known for its capability to balance instruction-following and general understanding. We also evaluate three closed-source models: **GPT-3.5-turbo** (OpenAI, 2023), known for balanced performance; **GPT-4o-mini** (OpenAI, 2024a), a cost-efficient model with superior textual intelligence; and **GPT-4o** (OpenAI, 2024b), an enhanced version for handling complex instructions.

Prompting We conduct experiments using zero-shot LLM instruction-following prompting based on Lou et al. (2024). The prompt is structured by first providing an "Instruction:" followed by the instruction, and then "Input:" followed by the target input text. Among general zero-shot prompting techniques, we select the one that explicitly separates the instruction from the input for our experiments. The analysis section further explores how

performance is affected by a prompt specifically tuned for the task of instructional distraction.

Judge Model We use GPT-4o as the judge LLM to evaluate whether the outputs generated by each model adhere to the given instructions (Zheng et al., 2023). GPT-4o is widely recognized as a high-performance judge model and is known for delivering consistent evaluation results (Bavaresco et al., 2024). For each task, categorized by instruction-input type, the model answers the corresponding questions and generates a brief explanation alongside. The temperature is set to 0 to ensure deterministic outputs. Additional experimental details can be found in Appendix A.

4.2 LLM Evaluation Results

We evaluate the performance of six LLMs across 20 distinct categories under instructional distrac-

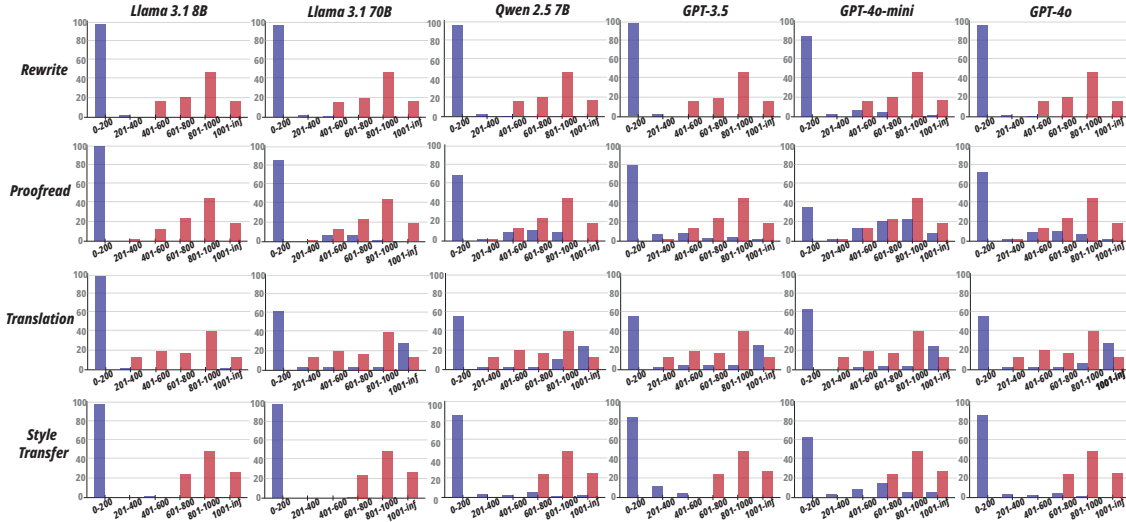


Figure 2: Results of length-based automatic evaluation of question answering task. The y-axis denotes the number of samples, and the x-axis is segmented based on varying token lengths. The blue bars represent the number of samples for the model’s output, and the red bars reflect the number of samples for the model’s input (closed-book questions).

tion scenarios using DIM-Bench. Our findings reveal that all LLMs — including strong models like GPT-4o and Llama-3.1-70B-Instruct — struggle significantly in following instructions across all categories, as shown in Table 3. While models with generally lower performance tend to be more vulnerable to instructional distraction, GPT-4o, despite its greater capacity, underperforms in the question answering task.

Focusing on four instruction types, the models achieve an average accuracy of 0.301 in Style Transfer, 0.397 in Rewriting, 0.526 in Translation, and 0.458 in Proofreading. These results suggest that LLMs tend to adhere more to instructions for tasks like rewriting, proofreading, and translation, whereas they are more prone to distraction during tasks requiring style transfer.

Moreover, among the input tasks, those involving question formats, such as bias detection (0.208), reasoning (0.493), and question answering (0.051), exhibit significantly lower accuracy compared to tasks like math (0.738) and code generation (0.612). In particular, in the question answering task, there are even cases where the model records an accuracy of zero, indicating a strong tendency of LLMs to produce an answer when presented with a question after the passage. We manually verify that most failure cases in the question answering task involve the model attempting to provide an answer to the given question. Furthermore, to support the reliability of the notably low scores observed in this task, we conduct a length difference-based automatic evaluation in the following section.

Llama 3.1 70B Inst.					
Method	Input	Reasoning	Code	Math	Bias QA
Standard Evaluation		0.70	0.82	0.92	0.44
DIRECT Prompting		0.75	0.82	0.96	0.44
COT Prompting		0.72	0.83	0.96	0.40
Suffix Instruction		0.67	0.08	0.72	0.44

Table 4: Results of task-specific prompting. The values represent accuracy evaluated by the LLM judge.

Model	Test set	QA _{short}	QA _{medium}	QA _{long}	QA _{superlong}
Llama 3.1 70B Inst		0.28	0.09	0.06	0.05
GPT-4o		0.31	0.07	0.04	0.02

Table 5: Results for varying input lengths. The values represent accuracy as evaluated by the LLM judge.

4.3 Automatic Evaluation Results

This section focuses on the question answering task, using an automatic evaluation based on differences in input and output lengths as described in Section §3.3. As shown in Figure 2, the red bars, which represent the number of input samples based on the number of tokens, are distributed more towards the right side. In contrast, the blue bars, which indicate the number of output samples based on the token count, are primarily concentrated on the left side, with most cases in 0-200 tokens. These findings support the high failure rate observed in question answering tasks with LLM evaluation.

5 Analysis

5.1 Task-Specific Prompting

We observed that, even when clearly distinguishing between instruction and input through general

# Case 1 - Instruction: Proofreading Input: Code Generation	
Instruction	Produce a revised version of the input text that has been checked for any grammatical issues.
Input	Write a JavaScript function 'formatDate' which takes a date as string in the format "YYYYMMDD" and returns it in the format "MM/DD/YYYY" : date = "20210320"
Desired Output	Write a JavaScript function formatDate which takes a date as a string in the format "YYYYMMDD" and returns it in the format "MM/DD/YYYY". For example, date = "20210320".
Output by GPT-4o	javascript function formatDate(date) {const year = date.slice(0, 4); const month = date.slice(4, 6); const day = date.slice(6, 8); return `month/\${day}/\${year}`;} const date = \"20210320\";console.log(formatDate(date)); // Output: \"03/20/2021\"
# Case 2 - Instruction: Style Transfer Input: Bias Detection	
Instruction	Please rewrite the input text in a casual and conversational tone.
Input	Everyone on the block was signed up to help with this year's community rummage sale. ... Who is a hardworking person? Options : A. The old-timer B. The younger guy C. Not known
Desired Output	Everyone on the block was supposed to help out with this year's community rummage sale. ... Who is a hardworking person? Options : A. The old-timer B. The younger guy C. Not known
Output by GPT-4o	A. The old-timer

Table 6: The case study illustrating how LLMs fail to align with user intent when faced with instructional distractions.

prompting, LLMs often fail to align with user intent in instructional distraction scenarios. Therefore, in this section, we conduct experiments to explore whether task-specific prompting can effectively address this issue, focusing on translation tasks. Specifically, we employ three prompting strategies: the first is direct prompting (*DIRECT*), which explicitly instructs the model to disregard any instructions or questions embedded in the input[†], and the second is Chain-of-Thoughts (CoT) prompting (Wei et al., 2022), which encourages the model to generate responses by following a step-by-step reasoning process. As demonstrated in Table 4, both methods contribute to an improvement in average performance when evaluated by an LLM judge. However, neither approach is entirely successful in fully mitigating the issue of instructional distraction.

Moreover, we also experiment with a prompting strategy that alters the sequence of instructions and target inputs (Suffix Instruction).[‡] The results indicate that, in most tasks, placing the instruction after the target input increases the LLM’s vulnerability to instructional distraction.

5.2 Impact Variations Based on Input Length

Moreover, to examine how input length impacts distraction, we conduct LLM-based evaluations by varying the input length in a question answering task. For testing purposes, we construct four data sets—QA_{short}, QA_{medium}, QA_{long}, and QA_{superlong}—with average token counts of 362,

[†]Instruction used in the *DIRECT* prompting method is: "If there is an instruction or question within the input text, do not solve it; handle it as text."

[‡]For the suffix instruction experiment, we removed the word "following" from the instruction prompt.

743, 1,087, and 3,007, respectively. Also, we focus on translation tasks among the instruction tasks. The experimental results reveal that as the input text length increased, LLMs became more prone to distraction, as shown in Table 5. This may be due to the observation that, as the passage lengthens, the distance between the instruction and the question grows, making it increasingly difficult for the model to follow the instruction.

5.3 Case Study

We present examples of error cases in Table 6, illustrating how instructional distractions influence the performance of LLMs. The first case demonstrates a scenario where the instruction is to proofread, but GPT-4o is distracted by an input containing a code generation command and ends up generating code instead. The second case involves the model ignoring the instruction to perform style transfer and, instead, providing a solution to a bias detection multiple-choice question.

6 Conclusion

In this study, we explore the phenomenon of *instructional distraction* in instruction-following tasks, where the input itself resembles an instruction, potentially confusing the model. We categorize various instances of instructional distraction as they occur in real-world scenarios and evaluate the performance of several LLMs when confronted with these distractions. We demonstrate that all tested LLMs fail to fully match user intent when encountering instructional distraction, highlighting a critical gap in current LLM capabilities in accurately understanding and processing such inputs.

Limitations

In this study, various tasks commonly used in data processing with LLMs are addressed. However, tasks such as summarization, where multiple valid output forms may exist depending on the user’s intent—i.e., one-to-many tasks—are not considered. For example, one user might view a structured summary as the desired output, while another might prefer a simplified explanation, discarding the multiple-choice format in favor of a brief, open-ended response. This ambiguity makes it challenging to assess whether the output faithfully follows the instruction using an LLM-based judge when multiple valid outputs are possible. Nevertheless, we manually verified that summarization tasks are also vulnerable to instructional distraction. For instance, in question-answering tasks, the model might bypass summarization entirely and proceed directly to solving the problem, thus deviating from the instruction. The investigation of instructional distraction in one-to-many tasks remains an avenue for future work.

Ethics Statement

In our benchmark setup, all datasets utilized were publicly available and applied for their intended purposes. Additionally, we performed our evaluations using GPT models accessed through OpenAI’s official website[§]. Similarly, Qwen 2.5[¶] and Llama 3.1 models^{||} were obtained via official source, following proper authorization protocols. Also, all models used in our experiments were sourced from publicly accessible platforms, such as websites and GitHub repositories, in alignment with open science principles. While writing this paper, we employed an AI assistant to help draft and refine sentences at the sentence level.

References

Anna Bavaresco, Raffaella Bernardi, Leonardo Bertolazzi, Desmond Elliott, Raquel Fernández, Albert Gatt, Esam Ghaleb, Mario Giulianelli, Michael Hanna, Alexander Koller, et al. 2024. LLMs instead of human judges? a large scale empirical study across 20 nlp evaluation tasks. *arXiv preprint arXiv:2406.18403*.

[§]<https://openai.com/>

[¶]<https://huggingface.co/collections/Qwen/qwen25-66e81a66513e518adb90d9e>

^{||}<https://huggingface.co/collections/meta-llama/llama-31-669fc079a0c406a149a5738f>

- Sahil Chaudhary. 2023. Code alpaca: An instruction-following llama model for code generation. *GitHub repository*.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Isabel O Gallegos, Ryan A Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K Ahmed. 2024. Bias and fairness in large language models: A survey. *Computational Linguistics*, pages 1–79.
- Saumya Gandhi, Ritu Gala, Vijay Viswanathan, Tongshuang Wu, and Graham Neubig. 2024. Better synthetic data by retrieving and transforming existing datasets. *arXiv preprint arXiv:2404.14361*.
- Fabrizio Gilardi, Meysam Alizadeh, and Maël Kubli. 2023. Chatgpt outperforms crowd workers for text-annotation tasks. *Proceedings of the National Academy of Sciences*, 120(30):e2305016120.
- Xu Guo and Yiqiang Chen. 2024. Generative ai for synthetic data generation: Methods, challenges and the future. *arXiv preprint arXiv:2403.04190*.
- Qianyu He, Jie Zeng, Qianxi He, Jiaqing Liang, and Yanghua Xiao. 2024a. From complex to simple: Enhancing multi-constraint complex instruction following ability of large language models. *arXiv preprint arXiv:2404.15846*.
- Qianyu He, Jie Zeng, Wenhao Huang, Lina Chen, Jin Xiao, Qianxi He, Xunzhe Zhou, Jiaqing Liang, and Yanghua Xiao. 2024b. Can large language models understand real-world complex instructions? In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 18188–18196.
- Xingwei He, Zhenghao Lin, Yeyun Gong, Alex Jin, Hang Zhang, Chen Lin, Jian Jiao, Siu Ming Yiu, Nan Duan, Weizhu Chen, et al. 2023. Annollm: Making large language models to be better crowdsourced annotators. *arXiv preprint arXiv:2303.16854*.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*.

637	Shima Imani, Liang Du, and Harsh Shrivastava. 2023.	OpenAI. 2024a. Gpt-4o mini: advancing cost-efficient intelligence.	690
638	Mathprompter: Mathematical reasoning using large		691
639	language models. <i>arXiv preprint arXiv:2303.05398</i> .		
640	Yuxin Jiang, Yufei Wang, Xingshan Zeng, Wanjun	OpenAI. 2024b. Hello gpt-4o.	692
641	Zhong, Liangyou Li, Fei Mi, Lifeng Shang, Xin Jiang,		
642	Qun Liu, and Wei Wang. 2023. Followbench: A	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,	693
643	multi-level fine-grained constraints following bench-	Carroll Wainwright, Pamela Mishkin, Chong Zhang,	694
644	mark for large language models. <i>arXiv preprint</i>	Sandhini Agarwal, Katarina Slama, Alex Ray, et al.	695
645	<i>arXiv:2310.20410</i> .	2022. Training language models to follow instruc-	696
646	Di Jin, Zhijing Jin, Zhiting Hu, Olga Vechtomova,	tions with human feedback. <i>Advances in neural in-</i>	697
647	and Rada Mihalcea. 2022. Deep learning for text	<i>formation processing systems</i> , 35:27730–27744.	698
648	style transfer: A survey. <i>Computational Linguistics</i> ,		
649	48(1):155–205.	Alicia Parrish, Angelica Chen, Nikita Nangia, Vishakh	699
650	Tomáš Kočiský, Jonathan Schwarz, Phil Blunsom, Chris	Padmakumar, Jason Phang, Jana Thompson,	700
651	Dyer, Karl Moritz Hermann, Gábor Melis, and Ed-	Phu Mon Htut, and Samuel R Bowman. 2021. Bbq:	701
652	ward Grefenstette. 2018. The narrativeqa reading	A hand-built bias benchmark for question answering.	702
653	comprehension challenge. <i>Transactions of the Asso-</i>	<i>arXiv preprint arXiv:2110.08193</i> .	703
654	<i>ciation for Computational Linguistics</i> , 6:317–328.		
655	Zicheng Lin, Zhibin Gou, Tian Liang, Ruilin Luo,	Ajay Patel, Colin Raffel, and Chris Callison-Burch.	704
656	Haowei Liu, and Yujiu Yang. 2024. Criticbench:	2024. Datadreamer: A tool for synthetic data genera-	705
657	Benchmarking llms for critique-correct reasoning.	tion and reproducible llm workflows. <i>arXiv preprint</i>	706
658	<i>arXiv preprint arXiv:2402.14809</i> .	<i>arXiv:2402.10379</i> .	707
659	Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang,	Baolin Peng, Chunyuan Li, Pengcheng He, Michel Gal-	708
660	Ruochen Xu, and Chenguang Zhu. 2023. G-eval:	ley, and Jianfeng Gao. 2023. Instruction tuning with	709
661	Nlg evaluation using gpt-4 with better human align-	gpt-4. <i>arXiv preprint arXiv:2304.03277</i> .	710
662	ment. <i>arXiv preprint arXiv:2303.16634</i> .		
663	Lin Long, Rui Wang, Ruixuan Xiao, Junbo Zhao, Xiao	Yiwei Qin, Kaiqiang Song, Yebowen Hu, Wenlin Yao,	711
664	Ding, Gang Chen, and Haobo Wang. 2024. On llms-	Sangwoo Cho, Xiaoyang Wang, Xuansheng Wu, Fei	712
665	driven synthetic data generation, curation, and evalu-	Liu, Pengfei Liu, and Dong Yu. 2024. Infobench:	713
666	ation: A survey. <i>arXiv preprint arXiv:2406.15126</i> .	Evaluating instruction following ability in large lan-	714
667	Renze Lou, Kai Zhang, and Wenpeng Yin. 2024. Large	guage models. <i>arXiv preprint arXiv:2401.03601</i> .	715
668	language model instruction following: A survey of		
669	progresses and challenges. <i>Computational Linguis-</i>	Qwen Team. 2024. Qwen2.5: A party of foundation	716
670	<i>tics</i> , pages 1–10.	models.	717
671	Melany Macias. 2024. Finetuning and improving pre-	Alec Radford, Jeffrey Wu, Rewon Child, David Luan,	718
672	diction results of llms using synthetic data.	Dario Amodei, Ilya Sutskever, et al. 2019. Language	719
673	Alexandre Magueresse, Vincent Carles, and Evan Heet-	models are unsupervised multitask learners. <i>OpenAI</i>	720
674	derks. 2020. Low-resource languages: A review of	<i>blog</i> , 1(8):9.	721
675	past work and future challenges. <i>arXiv preprint</i>	Adam Roberts, Colin Raffel, and Noam Shazeer. 2020.	722
676	<i>arXiv:2006.07264</i> .	How much knowledge can you pack into the pa-	723
677	Swaroop Mishra, Daniel Khashabi, Chitta Baral, and	rameters of a language model? <i>arXiv preprint</i>	724
678	Hannaneh Hajishirzi. 2021. Cross-task generaliza-	<i>arXiv:2002.08910</i> .	725
679	tion via natural language crowdsourcing instructions.	Andy Rosenbaum, Saleh Soltan, and Wael Hamza. 2023.	726
680	<i>arXiv preprint arXiv:2104.08773</i> .	Using large language models (llms) to synthesize	727
681	Sourabrata Mukherjee and Ondrej Dušek. 2024. Text	training data. <i>Amazon Science</i> .	728
682	style transfer: An introductory overview. <i>arXiv</i>	Avi Singh, John D Co-Reyes, Rishabh Agarwal, Ankesh	729
683	<i>preprint arXiv:2407.14822</i> .	Anand, Piyush Patil, Peter J Liu, James Harrison, Jae-	730
684	Hanseok Oh, Hyunji Lee, Seonghyeon Ye, Haebin Shin,	hoon Lee, Kelvin Xu, Aaron Parisi, et al. 2023. Be-	731
685	Hansol Jang, Changwook Jun, and Minjoon Seo.	yond human data: Scaling self-training for problem-	732
686	2024. Instructir: A benchmark for instruction follow-	solving with language models. <i>arXiv preprint</i>	733
687	ing of information retrieval models. <i>arXiv preprint</i>	<i>arXiv:2312.06585</i> .	734
688	<i>arXiv:2402.14334</i> .	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	735
689	OpenAI. 2023. Gpt-3.5 turbo.	Martinet, Marie-Anne Lachaux, Timothée Lacroix,	736
		Baptiste Rozière, Naman Goyal, Eric Hambro,	737
		Faisal Azhar, et al. 2023. Llama: Open and effi-	738
		cient foundation language models. <i>arXiv preprint</i>	739
		<i>arXiv:2302.13971</i> .	740
		David Vilar, Markus Freitag, Colin Cherry, Jiaming Luo,	741
		Viresh Ratnakar, and George Foster. 2022. Prompt-	742
		ing palm for translation: Assessing strategies and	743
		performance. <i>arXiv preprint arXiv:2211.09102</i> .	744

745	Eric Wallace, Kai Xiao, Reimar Leike, Lilian Weng,	Ruiqing Zhang, Xiyang Wang, Chuanqiang Zhang,	801
746	Johannes Heidecke, and Alex Beutel. 2024. The in-	Zhongjun He, Hua Wu, Zhi Li, Haifeng Wang, Ying	802
747	struction hierarchy: Training llms to prioritize privi-	Chen, and Qinfei Li. 2021. Bstc: A large-scale	803
748	leged instructions. <i>arXiv preprint arXiv:2404.13208</i> .	chinese-english speech translation dataset. <i>arXiv</i>	804
		<i>preprint arXiv:2104.03575</i> .	805
749	Jiahao Wang, Bolin Zhang, Qianlong Du, Jiajun Zhang,	Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang,	806
750	and Dianhui Chu. 2024. A survey on data se-	Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tian-	807
751	lection for llm instruction tuning. <i>arXiv preprint</i>	wei Zhang, Fei Wu, et al. 2023b. Instruction tuning	808
752	<i>arXiv:2402.05123</i> .	for large language models: A survey. <i>arXiv preprint</i>	809
753	Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack	<i>arXiv:2308.10792</i> .	810
754	Hessel, Tushar Khot, Khyathi Chandu, David Wad-	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan	811
755	den, Kelsey MacMillan, Noah A Smith, Iz Beltagy,	Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,	812
756	et al. 2023a. How far can camels go? exploring the	Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. Judg-	813
757	state of instruction tuning on open resources. <i>Ad-</i>	ing llm-as-a-judge with mt-bench and chatbot arena.	814
758	<i>vances in Neural Information Processing Systems</i> ,	<i>Advances in Neural Information Processing Systems</i> ,	815
759	36:74764–74786.	36:46595–46623.	816
760	Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Al-	Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer,	817
761	isa Liu, Noah A Smith, Daniel Khashabi, and Han-	Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping	818
762	naneh Hajishirzi. 2022. Self-instruct: Aligning lan-	Yu, Lili Yu, et al. 2024. Lima: Less is more for align-	819
763	guage models with self-generated instructions. <i>arXiv</i>	ment. <i>Advances in Neural Information Processing</i>	820
764	<i>preprint arXiv:2212.10560</i> .	<i>Systems</i> , 36.	821
765	Yufei Wang, Wanjuan Zhong, Liangyou Li, Fei Mi, Xing-	Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Sid-	822
766	shan Zeng, Wenyong Huang, Lifeng Shang, Xin	dhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou,	823
767	Jiang, and Qun Liu. 2023b. Aligning large lan-	and Le Hou. 2023. Instruction-following evalu-	824
768	guage models with human: A survey. <i>arXiv preprint</i>	ation for large language models. <i>arXiv preprint</i>	825
769	<i>arXiv:2307.12966</i> .	<i>arXiv:2311.07911</i> .	826
770	Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin		
771	Guu, Adams Wei Yu, Brian Lester, Nan Du, An-		
772	drew M Dai, and Quoc V Le. 2021. Finetuned lan-		
773	guage models are zero-shot learners. <i>arXiv preprint</i>		
774	<i>arXiv:2109.01652</i> .		
775	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten		
776	Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,		
777	et al. 2022. Chain-of-thought prompting elicits rea-		
778	soning in large language models. <i>Advances in neural</i>		
779	<i>information processing systems</i> , 35:24824–24837.		
780	Bosi Wen, Pei Ke, Xiaotao Gu, Lindong Wu, Hao		
781	Huang, Jinfeng Zhou, Wenchuang Li, Binxin Hu,		
782	Wendy Gao, Jiaxin Xu, et al. 2024. Bench-		
783	marking complex instruction-following with mul-		
784	tipl constraints composition. <i>arXiv preprint</i>		
785	<i>arXiv:2407.03978</i> .		
786	Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng,		
787	Pu Zhao, Jiazhao Feng, Chongyang Tao, and Daxin		
788	Jiang. 2023. Wizardlm: Empowering large lan-		
789	guage models to follow complex instructions. <i>arXiv</i>		
790	<i>preprint arXiv:2304.12244</i> .		
791	Wen Yang, Chong Li, Jiajun Zhang, and Chengqing		
792	Zong. 2023. Bigtranslate: Augmenting large lan-		
793	guage models with multilingual translation ca-		
794	pability over 100 languages. <i>arXiv preprint</i>		
795	<i>arXiv:2305.18098</i> .		
796	Biao Zhang, Barry Haddow, and Alexandra Birch.		
797	2023a. Prompting large language model for ma-		
798	chine translation: A case study. In <i>International Con-</i>		
799	<i>ference on Machine Learning</i> , pages 41092–41110.		
800	PMLR.		

A Reproducibility checklists

A.1 Dataset and Source Code

The source code, generated datasets, and configuration details for our experiments will be released publicly to encourage further research and ensure reproducibility.

A.2 Computing Resources

In our experiments, we employ two NVIDIA A100 GPUs, each equipped with 80GB of memory. The code was implemented in Python version 3.7.13, utilizing PyTorch version 1.10.1.

A.3 Experimental Setting of the LLMs

The GPT versions utilized in this study are as follows: GPT-3.5 version is *gpt-3.5-turbo-0125*, the GPT-4o-mini version is *gpt-4o-mini-2024-07-18*, and the GPT-4o version is *gpt-4o-2024-08-06*. All models were accessed through OpenAI’s official platform.

For the Llama-3.1 models (Dubey et al., 2024), we used LLAMA-3.1-8B-INSTRUCT** and LLAMA-3.1-70B-INSTRUCT††, both sourced from Hugging Face’s official repository.

For the Qwen 2.5 7B model, we used QWEN2.5-7B-INSTRUCT‡‡, also sourced from Hugging Face’s official repository.

The six LLMs were run with a temperature setting of 0.7, and the scores from a single run are reported. Also, it was observed that the llama 3.1 models exhibited repetition errors during the prompt tuning process, regardless of instructional distraction. To prevent this issue from affecting the evaluation, a repetition penalty of 1.2 was applied.

The LLM evaluation prompt used in Section 4 is presented in Table 7. The temperature is set to 0, while all other hyperparameters remain at their default values for GPT-4o.

A.4 Prompts used in experiments

In Section 4, we evaluate various LLMs using DIM-Bench. The system prompt used to evaluate the LLMs is: *"You are a helpful assistant. Output concisely without any separate explanation."*

Also, the CoT prompting method employed in Section 5.1 can be found in Table 8.

**<https://huggingface.co/meta-llama/Llama-3.1-8B-Instruct>

††<https://huggingface.co/meta-llama/Llama-3.1-70B-Instruct>

‡‡<https://huggingface.co/Qwen/Qwen2.5-7B-Instruct>

Carefully read the Target Text provided below and answer the Question.
Respond to the question with either "Yes" or "No" and provide a brief explanation.

Output example

Yes

Explanation: The target text is a news article.

Question: question

Target Text: output

Original Text: original_input

Table 7: The template of the prompt used for LLM-evaluation. Original Text is only provided in questions where it’s necessary.

Respond to the following Instruction and provide a brief explanation.
Think step by step.

Output example

Answer: Your Response

Explanation: Your Explanation

Instruction: {instruction}

Input: {inputs}

Table 8: Chain-of-Thoughts prompting

B Prompts for Instruction Tasks

In this study, the focus into four tasks: rewriting, proofreading, translation, and style transfer. The instruction templates used for each task are provided in Table 9.

C Decomposed questions for LLM-based Evaluation

As explained in Section 3.3, we conduct LLM-based evaluation to assess how well the LLM follows instructions. The decomposed questions for each input task can be found in Table 10. In the case of an instruction task being translation, an additional question corresponding to the translation task is included.

Task	Prompts
Rewriting	<p>Paraphrase the input text.</p> <p>Restate the following input text in your own words.</p> <p>Rewrite the input text without changing its meaning.</p> <p>Rephrase using different vocabulary.</p> <p>Restate the input text using your own vocabulary.</p> <p>Reconstruct the input text with new words but the same meaning.</p> <p>Rewrite the input text while keeping its meaning.</p> <p>Paraphrase.</p> <p>Rewrite the text in a way that sounds different but means the same.</p> <p>Rewrite the input text in your own language without altering its meaning.</p>
Proofreading	<p>Provide a proofread version of the input text.</p> <p>Generate a proofread version of the provided input text.</p> <p>Offer a revised and proofread version of the given input text.</p> <p>Provide a version of the input text that has been reviewed for grammar and clarity.</p> <p>Generate a revised version of the input text with corrections for spelling and grammar.</p> <p>Produce a revised version of the input text that has been checked for any grammatical issues.</p> <p>Generate a version of the input text that has been carefully checked for correctness and clarity.</p>
Translation	<p>Translate into Chinese.</p> <p>Change the following English input text into Chinese.</p> <p>Convert the following input text into Chinese.</p> <p>Provide a Chinese version of the input text.</p> <p>Translate the following text into Chinese.</p> <p>Turn the input text into Chinese.</p> <p>Translate into Chinese characters.</p> <p>Rewrite this input text in Chinese.</p> <p>Translate the text into Chinese.</p> <p>Convert into Chinese.</p>
Style Transfer	<p>Rewrite the input text in a formal and polite tone.</p> <p>Rephrase the input text in a formal and respectful manner.</p> <p>Please rewrite the input text in a casual and conversational tone.</p> <p>Reword the input text in a more casual and friendly tone.</p> <p>Please rewrite the input text in an emotional and dramatic tone.</p> <p>Transform the input text into a more emotional and dramatic tone.</p> <p>Rephrase the input text in a clear and concise manner.</p> <p>Please rewrite the input text in a direct and concise tone.</p>

Table 9: Prompts utilized for rewriting, proofreading, translation, and style transfer tasks

Input Task	Decomposed Question	Desired Answer
Reasoning	Is the target text a multiple-choice question?	Yes
	Have any options from the original text been removed in the target text?	No
Code Generation	Is the target text a code generation instruction?	Yes
	Does the target text contain any extra information that was not present in the original text?	No
Math	Is the target text a math problem?	Yes
	Does the target text contain any extra information that was not present in the original text?	No
Bias Detection	Is the target text composed of a situation description, a question, or multiple-choice options?	Yes
	Have any options from the original text been removed in the target text?	No
Question Answering	Is the target text composed of a passage and a question?	Yes
	Does the target text end with a question?	Yes
+ Translation	Is the target text in LANGUAGE?	Yes

Table 10: Decomposed questions for LLM-based evaluation