# Ask, Fail, Repeat: Meeseeks, an Iterative Feedback Benchmark for LLMs' Multi-turn Instruction-Following Ability

**Jiaming Wang[1], Yunke Zhao[1], Peng Ding, Jun Kuang[1],**
**Zongyu Wang[1], Xuezhi Cao[1], Xunliang Cai[1]**

[1]Meituan
{wangjiaming15}@meituan.com

## Abstract

The ability to follow instructions accurately is fundamental for Large Language Models (LLMs) to serve as reliable agents in real-world applications. For complex instructions, LLMs often struggle to fulfill all requirements in a single attempt. In practice, users typically provide iterative feedback until the LLM generates a response that meets all requirements. However, existing instruction-following benchmarks are either single-turn or introduce new requirements in each turn without allowing self-correction. To address this gap, we propose **Meeseeks** [1]. Meeseeks simulates realistic human-LLM interactions through an iterative feedback framework, which enables models to self-correct based on specific requirement failures in each turn, better reflecting real-world user-end usage patterns. Meanwhile, the benchmark implements a comprehensive evaluation system with 38 capability tags organized across three dimensions: Intent Recognition, Granular Content Validation, and Output Structure Validation. Through rigorous evaluation across LLMs, Meeseeks provides valuable insights into LLMs' instruction-following capabilities in multi-turn scenarios.

## 1 Introduction

LLM agents have become essential tools in various applications, from customer service to content creation.[21, 35, 24, 14, 27] However, as their use expands, the instructions users provide are often complex and multifaceted, posing challenges for accurate execution.[23, 40, 37, 31, 10, 22, 39, 38] Instruction following, or the ability to execute tasks based on natural language commands accurately, is critical for LLMs to be reliable in real-world scenarios. For instance, in healthcare, an LLM might need to generate reports with specific word counts or mention key terms multiple times. Misinterpreting such instructions could lead to errors with serious consequences.[16, 29] Similarly, in finance, precise adherence to instructions ensures compliance and accuracy, avoiding potential risks. [1, 12] These content creation tasks require LLMs to have strong text instruction-following capability. To evaluate the instruction-following ability of today's LLMs, researchers employs instruction-following benchmarks[1, 5, 17, 9, 32, 15, 11, 8, 36, 33, 41, 28] to determine how effectively LLMs can align with human intentions[30, 3, 4]

Self-correction, where LLMs adjust their outputs based on feedback, is a key aspect of interactive loops and draw attention from researchers these days. [2, 13, 34, 6, 7, 18–20, 25, 26]Typically, when

---

[1]Meeseeks is a benchmark designed to evaluate large language models' instruction-following abilities. The name is inspired by Mr. Meeseeks from "Rick and Morty," a character known for efficiently completing tasks as instructed.

users interact with LLMs and receive responses that not match the their requirements, they provide feedback so the LLM can correct itself. For example, when a user requests "provide me a story in 400 words, must contain keywords: 'wolf', 'superman', 'Meeseeks'", the LLM might initially generate a 410-word story mentioning only 'superman' and 'Meeseeks'. Upon receiving feedback "410 words is over 400 words, keywords 'wolf' is not included", the LLM then adjusts its response to meet all requirements - delivering a 400-word story incorporating all three keywords. However, previous benchmarks are all single-turn or providing new instruction after each turn, completely ignoring this important scenario. They fail to capture the interactive loop in which users give feedback and models adapt accordingly. As a result, these benchmarks may not accurately reflect the capabilities of LLMs in real-world use cases where self-correction to adapt the intention from user are essential.

To fill this void, we introduce Meeseeks, a multi-turn automatic instruction-following benchmark. Meeseeks embeds with 38 capability tags under 3 different evaluation dimensions, systematically defining what instruction-following should entail, examining the entire thought process of the model to establish an integrated instruction-following ability evaluation system, effectively projecting the model's instruction-following ability to its performance on Meeseeks. **Moreover,** Meeseeks deploys an iterative multi-turn evaluation framework that dynamically simulates human-LLM interaction. The advanced framework of Meeseeks enables it fully automatically to identify specific requirements that the under-evaluation LLM fails to meet and explain why they were not met. After each complete cycle of getting LLM response and evaluating the response, Meeseeks provides and prompts the feedback to the LLM, then evaluates the response after self-correction. The evaluation result after each chat turn is recorded. LLMs' performance under multi-turn scenario provides an important insight for evaluating the effectiveness of an LLM under real-world human-agent collaboration scenario. **Furthermore,** due to the high cost from single to multiple turns evaluation and more challenging data introduced, Meeseeks optimizes evaluation progress to reduce the cost and raise the accuracy of the evaluation workflow. **Finally,** Meeseeks reveals nuanced patterns in LLMs' multi-turn self-correction abilities that challenge conventional understanding of instruction following. Our analysis identifies three distinct trajectories: divergence (initially similar models developing performance gaps), convergence (disparate models achieving similar outcomes), and performance reversal (initial advantages becoming disadvantages). Most significantly, we discover that reasoning models' advantages diminish across multiple turns, with non-reasoning models sometimes outperforming their reasoning counterparts by the third turn. These findings provide important insight for researchers on LLMs' further iterations. In conclusion, our main contributions are as follows:

- We proposed Meeseeks, a benchmark that addresses a gap in multi-turn instruction-following evaluation by enabling model self-correction through iterative feedback. Leveraging an integrated hierarchical taxonomy, Meeseeks systematically evaluates LLMs' real-world instruction-following capabilities in multi-turn scenarios.

- We proposed Code-guided rule-augmented LLM-based evaluation, an enhanced framework that improves both efficiency and accuracy for multi-turn instruction-following assessment. By optimizing the existing rule-augmented evaluation approach, our method effectively handles complex datasets while addressing the computational challenges of scaling from single-turn to multi-turn scenarios, providing future researchers with a practical solution for high-performance multi-turn instruction-following evaluation.

- The evaluation results from Meeseeks reveal that multi-turn interactions create different model performance patterns, and the performance gap between reasoning models and non-reasoning models gradually diminishes or even reverses as turn iterations increase. These findings provide trailblazing insight for LLM researchers.

## 2   Related Works

**Existing single-turn instruction-following benchmarks**    IF-Eval[41] first pioneered a complete instruction-following benchmark using rule-verifiable instructions (e.g., "Include `keyword1`, `keyword2` in your response"). While effective for automatic verification, IF-Eval's simple synthesized rule-verifiable dataset limited its applicability to real-world scenarios. Recent instruction-following benchmarks such as InFoBench[28], CELLO[8], FollowBench[11] and COLLIE[36] advanced this field by applying constraint-based frameworks to adapt more complicated data. However, it still requires specific format constraints for evaluation purposes. Complexbench[33] addressed these challenges effectively. It deploys rule-augmented LLM-based evaluation progress to complete the

evaluation, eliminating extensive prompt restrictions for formatting responses from the LLM under evaluation.. **However,** these instruction-following benchmarks have primarily focused on single-turn interactions, where all requirements are expected to be fulfilled in "one shot", rather than reflecting the multi-turn dialogue patterns typical in real-world user interactions.

**Existing multi-turn instruction-following benchmarks** Parrot [32] introduced a framework for collecting human-like multi-turn instructions with natural dialogue patterns. Multi-IF [9] extended IF-Eval to multi-turn sequences across multiple languages, revealing performance degradation in later turns and non-Latin scripts. StructFlowBench[15] proposed a structured framework with six inter-turn relationships to evaluate dialogue coherence and contextual understanding. **However,** these multi-turn instruction-following benchmarks only provide LLMs with a single attempt at each turn, introducing additional requirements in subsequent turns while overlooking the models' inherent self-correction capabilities.

# 3 Meeseeks

In the following sections, we present a comprehensive description of Meeseeks. Section 3.1 outlines the Meeseeks's multi-turn pipeline. Section 3.2 demonstrates how Meeseeks enhances rule-augmented LLM-based evaluation to address the escalating costs in multi-turn frameworks. Section 3.3 introduces Meeseeks's dataset from data parameterization approach and evaluation system. Section 3.4 details the metrics Meeseeks apply to evaluate LLMs' instruction-following capabilities.
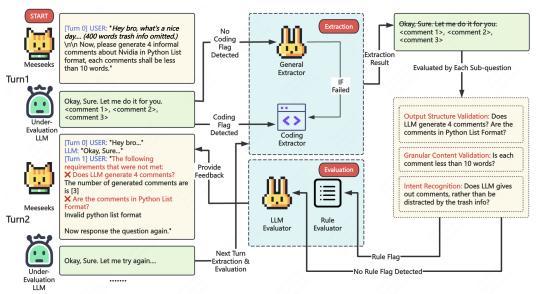


Figure 1: Meeseeks evaluation system

## 3.1 Meeseeks multi-turn pipeline

Meeseeks implements an iterative evaluation and self-correction framework for LLM responses as shown in Figure 1. The framework operates through a maximum of three turns in its default configuration, incorporating both LLM-based and rule-based evaluation mechanisms.

**Response generation and evaluation pipeline** The evaluation process begins with the original question prompt, followed by the under-evaluation LLM's response. Each response undergoes systematic evaluation through the LLM Evaluator, which identifies unmet requirements and provides detailed explanations. These evaluations serve as feedback for subsequent turns, enabling iterative improvement of responses.

**Content extraction mechanism** The framework employs dual extraction systems: a General Extractor and a Coding Extractor. The selection between these extractors is determined by the presence of coding flags(details in Appendix A) within the data. The Coding Extractor handles code-specific content(detail in Section 3.2.), while the General Extractor processes standard responses

through regeneration. In cases where the Coding Extractor fails, the General Extractor serves as a fallback mechanism, ensuring robust extraction.

**Evaluation system**   The evaluation process is bifurcated based on rule labels within requirements. For null rule labels(details in Appendix A) , the LLM Evaluator directly assesses the response. When specific rule labels are present, corresponding rule-based evaluation functions are triggered to verify the extracted content's validity.

**Iterative self-correction process**   Evaluation results are integrated with the LLM's chat template for subsequent turns. The framework tracks improvement through Utility Rate and Meeseeks Score metrics after each self-correction iteration, simulating realistic human-agent collaboration scenarios. Researchers maintain the flexibility to modify templates according to their specific training data requirements.

```
[Turn 0] USER: "<Original Question>"
          LLM: "<Turn 0 Response>"
[Turn 1] USER: "The following requirements that were not met:...
                Now response the <Original Question> again."
          LLM: "<Turn 1 Response>"
[Turn 2] USER: "The following requirements that were not met:...
                Now response the <Original Question> again."
          LLM: "<Turn 2 Response>"
```

### 3.2   Code-guided rule-augmented LLM-based evaluation

As we mentioned in Section 2, Early instruction-following benchmarks like IF-Eval and subsequent frameworks (InFoBench, CELLO, etc.) were limited by their reliance on rule-verifiable instructions and specific format constraints, making them less applicable to real-world scenarios. Complexbench overcame these limitations by introducing a rule-augmented LLM-based evaluation approach that eliminated the need for strict formatting requirements, thus enabling more flexible and practical evaluations.

Similar to ComplexBench, Meeseeks employs LLM-based extraction to eliminate redundant prompt constraints. For example, given the prompt "write me a 2000-word research report about LLM", an LLM under evaluation might generate: "Okay, I will write you a 2000-word essay. LLM is...(1998 words). I hope you like the answer." The system first utilizes an LLM Evaluator to extract the core content by removing auxiliary text (e.g., "Okay, I will write you a 2000-word essay" and "I hope you like the answer"). Subsequently, it applies rule-based evaluation to verify whether the extracted content "LLM is...(1998 words)" meets the 2000-word requirement. **However,** ComplexBench's LLM-extraction approach faces both efficiency and accuracy challenges. For instance, when evaluating responses to prompts like "Provide 150 sentences about LLM", the LLM Evaluator must extract all 150 sentences from the response. Our experiments with state-of-the-art LLMs (e.g., GPT-4o-1120 and Qwen2.5-72b) reveal that extraction accuracy deteriorates significantly as the number of elements increases, primarily due to hallucination. Moreover, these large-scale extraction tasks are computationally intensive. For long-form content evaluation, such as multi-thousand-word essays, the LLM Evaluator must process the entire text to identify and remove non-essential components, leading to substantial computational overhead.

To address these issues, we propose Code-guided rule-augmented LLM-based evaluation. Specifically, Meeseeks guides LLM to extract the to-be-evaluated part by coding, rather than having the model regenerate all context. A sample is provided in Table 1. By applying the function created by the LLM Evaluator, the necessary information can be extracted efficiently without the need to regenerate the entire context. This optimization significantly improved Meeseeks' end-to-end evaluation accuracy from 78.7% to 98.4% on Meeseeks' dataset, while substantially reducing the token generation during the extraction phase, thus reducing the time cost of the whole framework. Moreover, Meeseeks framework has great optimization on its evaluation progress, when using a evaluation batch size of 200, it can finish evaluating the quick-start-dataset in 600-900 seconds, making the process 4 to 5 times faster compared to bach size 1 evaluation.

### 3.3   Dataset construction and features

Meeseeks offers over 700 pre-synthesized data entries as quick-start-dataset, attaching 38 different capability tags, offering broader coverage compared to most of existing open-source instruction-

Table 1: Extraction output comparison

| | |
|---|---|
| **Under-evaluation LLM response** | Due to the length of the report, the following is a simplified version of the research report outline, which you can further expand as needed:—# Research Report: **System Spam Problems and Solutions Faced by Ordinary Users## Summary This report collects data for ordinary users in..... (2000 words omitted).** —Please note that this is only a simplified version of the report outline. |
| **Complexbench output** | **System Spam Problems and Solutions Faced by Ordinary Users## Summary This report collects data for ordinary users in..... (2000 words omitted).** <br> 2,754 Tokens |
| **Meeseeks output** | ```python def extract_info_list(model_response): """ Because the object to be captured is the entire research report outline. Therefore: After capturing "[# Research report: ]", the content between "[.---Please note that]" is enough. """ cleaned_text = re.sub( r'^.*# Research Report: |.---Please note that.*$', '', model_response, flags=re.DOTALL) return [cleaned_text] ``` <br> 98 Tokens |

following benchmarks. Additionally, all subjective capability requirements have been removed from the dataset, ensuring that the evaluation results are reliable and trustworthy. The dataset is in **Chinese** for now, later we will release English dataset.

Following sections present three key features of Meeseeks dataset system: capability tags for precise user prompt requirement classification, data parameterization for flexible evaluation data generation, and a cognitive-inspired evaluation system that decomposes instruction-following into three fundamental dimensions.

**Capability tags**    A real-world user prompt generally contains more than one output requirements. An example user prompt is shown in Table 2. This user prompt contains 3 objective requirements (subjective requirements such as "financial services attribution to videos" are ignored.)

Table 2: An example user prompt (Demo use, not included in quick-start-dataset)

| Question Context | | Requirement 1 | Requirement 2 | Requirement 3 |
|---|---|---|---|---|
| Generate 32 informal comments and 40 formal comments in the comment section of a financial services short video, from the perspective of a consumer. Each comment should have 7 words | **Context** | Generate 32 informal comments | Generate 40 formal comments | Each comment should have 7 words |
| | **Capability Tags** | Element number requirement | Element number requirement | Generate in 0-10 words; Generate at accurate word number; Generate multiple results |

To more precisely quantify a LLM's performance in following specific instruction-following capabilities, Meeseeks defines capability tags. Each tag represents a specific capability and is categorized into three levels, with lower levels being subordinate to higher levels. For instance, "Intent Recognition," "Granular Content Validation," and "Output Structure Validation" are three top-level (level 3) capability tags. The hierarchy of capability tags is in Appendix B. The capability tag assigned to each requirement is always the lowest level tag possible.

**Data Parameterization**    Meeseeks implements data parameterization to enable flexible evaluation data generation. This approach allows researchers to customize key parameters such as background context and content length, facilitating large-scale synthetic data generation for specific capability testing. The parameterization ensures reliable assessments by minimizing evaluation pipeline fluctuations. (Sample Template in Appendix C)

An example script is provided with Meeseeks code where researchers can customize it by modifying the placeholders: {THEME} to change the video type, {ITEM_NUM} to specify the required

number of user comments, and {KEYWORDS} to filter out certain keywords from the comments. This parameterization enables dataset creation without dealing with complex internal data structures.

**Cognitive-inspired instruction-following evaluation system** Prior instruction-following benchmarks evaluate LLMs through isolated capability tags (e.g., "Count Limit", "Lexical Constraint"), failing to establish a systematic framework for assessing general instruction-following abilities. This fragmented approach obscures the relationship between benchmark scores and actual instruction-following capabilities. To address this issue, we propose a cognitive-inspired framework that decomposes instruction-following into three fundamental dimensions (Figure 1): **Intent Recognition** for instruction comprehension, **Granular Content Validation** for response element compliance, and **Output Structure Validation** for organizational adherence. This decomposition mirrors the natural reasoning process of LLMs, from understanding instructions to generating structured responses. Detailed specifications are provided in Appendix B.

## 3.4 Metrics

**Utility Rate** When serving as a trustworthy agent, an LLM must fulfill all requirements specified in the user's prompt for its response to be considered usable. Beyond conventional accuracy metrics, Meeseeks adopts the strict "is followed" metric introduced by IF-Eval[41], but redefines it in a more intuitive and general manner as the Utility Rate—the proportion of responses that fully satisfy all prompt requirements. The Utility Rate reflects the ratio of usable responses and is essential for evaluating an LLM's effectiveness as a practical and professional agent.

To evaluate LLMs as trustworthy agents, we adopt and generalize IF-Eval[41]'s strict "is followed" metric, introducing it as Utility Rate - the proportion of responses that fully satisfy all prompt requirements. This intuitive metric measures the ratio of truly usable responses, serving as a crucial indicator of an LLM's capability to function as a reliable professional agent.

$$\text{Utility Rate} = \frac{\sum_{i=1}^{n} U_i}{n}, \text{ where } U_i = \begin{cases} 1, & \text{if response is usable} \\ 0, & \text{if response is not usable} \end{cases}$$

**Meeseeks Score** To help researchers identify subtle differences in instruction-following abilities between different LLMs, we propose Meeseeks Score, which indicates the overall capability tag accuracy of the LLM. Meeseeks Score averages the scores of all the level 1 capability tags associated with the user prompt.

$$\text{Meeseeks Score} = \frac{\sum_{j=1}^{m} \text{Score}_{\text{tag}_j}}{m}$$

where $m$ is the total number of level 1 capability tags associated with the current issue.

The score for each individual level 1 capability tag is the average of the scores of its respective requirements.

$$\text{Score}_{\text{tag}} = \frac{\sum_{i=1}^{n} \text{Score}_{\text{requirement}_i}}{n}$$

where $n$ is the total number of requirements for the capability tag.

**Accuracy on each capability tag** Additionally, Meeseeks offers accuracy measurements across various capability levels in a capability report, enabling researchers to understand the instruction-following proficiency of LLMs from different perspectives.

## 4 Evaluation

In this section, we evaluate 11 representative LLMs, comprising 4 RLLMs(Reasoning Large Language Models/reasoning models) and 8 LLMs(Large Language Models/non-reasoning models). Using a 3-turn Meeseeks framework, we record both the utility rate and Meeseeks Score for each turn.

## 4.1 Setup

In our evaluation experiments, we employ Qwen2.5-32B-Instruct as both the general evaluator and general extractor, and Qwen2.5-Coder-32B-Instruct as the coding extractor. Thus, the minimum requirement would be 160 GPU Ram. To validate our results, we conducted cross-validation with three annotators holding undergraduate degrees, and the cross-validation result showed that the end-to-end accuracy reached 98.4%.

Table 3: Model Performance Comparison (†: open-source model)

| Model Name | Utility Rate | | | Meeseeks Score | | |
|---|---|---|---|---|---|---|
| | turn1 | turn2 | turn3 | turn1 | turn2 | turn3 |
| *Reasoning Models* | | | | | | |
| OpenAI/o3-mini (high) | .583 | .734 | **.781** | .830 | .892 | .909 |
| OpenAI/o3-mini (medium) | .578 | .723 | **.769** | .834 | .884 | .901 |
| Anthropic/Claude-3.7-Sonnet-thinking | .482 | .613 | **.697** | .799 | .854 | .882 |
| DeepSeek-R1[†] | .326 | .485 | **.547** | .713 | .785 | .815 |
| *Non-Reasoning Models* | | | | | | |
| Anthropic/Claude-3.7-Sonnet | .359 | .573 | **.661** | .744 | .848 | .874 |
| DeepSeek-V3-Chat-20250324[†] | .346 | .492 | **.561** | .722 | .802 | .825 |
| DeepSeek-V3-Chat-20241226[†] | .315 | .473 | **.558** | .701 | .791 | .824 |
| OpenAI/GPT-4o-20241120 | .312 | .453 | **.531** | .695 | .768 | .803 |
| Qwen2.5-32B-Instruct[†] | .278 | .417 | **.471** | .674 | .756 | .779 |
| Qwen2.5-72B-Instruct[†] | .278 | .395 | **.428** | .666 | .742 | .752 |
| OpenAI/GPT-4o-mini | .265 | .352 | **.395** | .647 | .702 | .729 |

[*] T1/T2/T3: Turn 1/2/3

## 4.2 Results & Analysis

Table 3 presents the evaluation results of 11 representative LLMs. It can be clearly observed that OpenAI/o3-mini (high) demonstrates the best first-turn performance compared to all other models, attaining a 58.3% utility rate, while most other models—including popular non-reasoning models—fail to reach a 50% utility rate in the initial interaction. This indicates that LLMs generally struggle to fully satisfy user requirements in a single attempt. With additional turns, all models demonstrate significant performance gains. For example, o3-mini (high) achieves a 78.1% utility rate by the third turn. However, the initial advantage of reasoning models narrows or even reverses as non-reasoning models catch up through iterative self-correction. **Moreover,** across all models, complex language constraints and word count requirements remain challenging, with no significant advantage observed for reasoning models in these dimensions.

To gain a deeper understanding of the experimental results, we conducted a comprehensive analysis from three perspectives. **Single-turn result** analyzes first-turn performance, revealing that most models struggle to fully meet user requirements in a single attempt, with reasoning models generally performing better. **Multi-turn result** focuses on multi-turn interactions, showing that most models significantly improve with additional turns, but the advantage of reasoning models narrows over time. **Analysis on capability tags** examines performance across different capability tags, finding that all models face challenges with complex language and word count requirements.

**Single-turn result:** First-turn results reveals significant performance gaps among the 11 evaluated LLMs. Excluding the O3-mini variants, other popular models failed to achieve even a 50% utility rate in initial interactions, indicating their inability to satisfy all user requirements in a single attempt for over 50% scenarios. This deficiency creates a poor first impression, which is particularly problematic in user-facing applications where models typically have only one opportunity to demonstrate value before users switch alternatives.
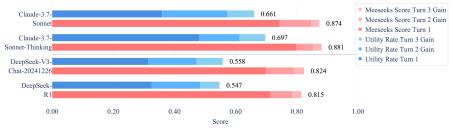


Figure 2: Utility rate and Meeseeks Score gain over turns

We observed that reasoning-enabled models generally outperform their non-reasoning counterparts in first-turn interactions. Our case analysis indicates that reasoning models systematically review and comprehend all requirements within a prompt before generating output, often including requirement checklists to verify each element's compliance in the response. Long chain of thought indeed reduces the likelihood of non-compliant responses significantly. However, Deepseek-R1's performance fell below expectations, failing to demonstrate a clear advantage over non-reasoning models. Our investigation revealed numerous inconsistencies between its reasoning chains and final outputs. For instance, in cases where the model correctly identified requirements during reasoning (e.g., recognizing that "Steamed Bass Fish" satisfies a three-word name requirement), the final answer would contain inconsistent results (outputting "Steamed Fish" instead). These reasoning-output misalignments substantially undermined the model's effectiveness despite its reasoning capabilities.
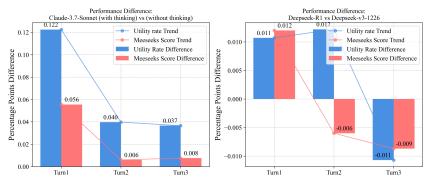


Figure 3: Utility rate and Meeseeks Score trend over turns

**Multi-turn result:** Moreover, the multi-turn results demonstrate LLMs' robust self-correction capabilities, with most models showing over 50 percentage point improvement in utility rates from turn 1 to turn 3. Notably, o3-mini(high) achieves a remarkable 78.1% utility rate by the third turn, highlighting its exceptional instruction-following capabilities. The multi-turn framework reveals deeper insights into models' instruction-following abilities:

- **Diverging Performance:** Models with similar single-turn performance can show significant differences over multiple turns. For instance, Qwen2.5-32B-Instructand Qwen2.5-72B-Instruct both started at 27.8% utility rate, but their gap widened to 4.3 percentage points by turn 3.

- **Convergent Performance:** Conversely, initially disparate models may converge. DeepSeek-V3-Chat versions (20250324 vs. 20241226) showed a 3.1 percentage point gap in turn 1, narrowing to 0.3 by turn 3.

- **Performance Reversal:** Initial advantages may reverse in later turns. DeepSeek-R1's 1.1 percentage point lead over DeepSeek-V3-Chat-20241226 in turn 1 turned into a 1.1 percentage point deficit by turn 3.

As discussed above, reasoning models theoretically should outperform non-reasoning models. However, this advantage diminishes across multiple turns. We analyze this phenomenon using two pairs of models: Claude-3.7-Sonnet/Sonnet-thinking and Deepseek-V3-20241226/Deepseek-R1 (Figure 2, Figure 3). Our case analysis reveals that the performance gap between reasoning and their corresponding non-reasoning models narrows with each iteration. In the case of Deepseek, the non-reasoning model Deepseek-V3-20241226 even surpasses Deepseek-R1 by turn 3. We argue that Meeseeks' multi-turn prompting and chat history effectively serve as external reasoning content for non-reasoning models, partially substituting the role of built-in reasoning capabilities. For Deepseek-R1, its inconsistencies between reasoning chains and final outputs become more pronounced in multi-turn scenarios, contributing to its performance decline relative to V3.

**Analysis on capability tags:** Beyond Meeseeks Score and Utility Rate, Meeseeks evaluates accuracy across individual capability tags. Table 4 presents the turn -3 accuracy distribution across capability tags for all 11 evaluated LLMs. More detailed information about the capability tags can be found in Appendix B.

We find reasoning models demonstrate no significant advantage over non-reasoning models in handling distractions. For instance, when presented with a complex, lengthy instruction requiring

Table 4: Performance of Turn-3 Models on Different Capability Tags (†: open-source model)

| Model Name | Intent Recognition | | Granular Content Validation | | | | | | | Output Structure Validation | | | | UR3 | MS3 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | OA | 3 | 4 | 5 | 6 | 7 | 8 | OA | 10 | 11 | 12 | OA | | |
| *Reasoning Models* | | | | | | | | | | | | | | | |
| o3-mini (high) | .725 | **.725** | .938 | .917 | .528 | 1.00 | .813 | .948 | **.893** | 1.00 | .941 | .895 | **.952** | **.781** | .909 |
| o3-mini (medium) | .706 | **.706** | .934 | .889 | .444 | 1.00 | .780 | .931 | **.875** | 1.00 | .959 | .895 | **.966** | **.769** | .901 |
| Claude-3.7-Sonnet-thinking | .823 | **.823** | .951 | .861 | .583 | 1.00 | .768 | .914 | **.874** | 1.00 | .969 | .895 | **.974** | **.697** | .882 |
| DeepSeek-R1† | .538 | **.538** | .973 | .778 | .667 | 1.00 | .611 | .852 | **.812** | 1.00 | .932 | 1.00 | **.948** | **.547** | .815 |
| *Non-Reasoning Models* | | | | | | | | | | | | | | | |
| Claude-3.7-Sonnet | .804 | **.804** | .955 | .889 | .528 | 1.00 | .733 | .900 | **.859** | 1.00 | .966 | .947 | **.973** | **.661** | .874 |
| DeepSeek-V3-Chat-20250324† | .873 | **.873** | .945 | .972 | .556 | 1.00 | .631 | .837 | **.806** | 1.00 | .953 | .895 | **.961** | **.561** | .825 |
| DeepSeek-V3-Chat-20241226† | .873 | **.873** | .945 | .944 | .556 | 1.00 | .634 | .828 | **.804** | 1.00 | .942 | .895 | **.953** | **.558** | .824 |
| GPT-4o-20241120 | .784 | **.784** | .951 | .917 | .472 | .906 | .589 | .758 | **.769** | .963 | .944 | .947 | **.948** | **.531** | .803 |
| Qwen2.5-32B-Instruct† | .794 | **.794** | .954 | .889 | .556 | 1.00 | .570 | .756 | **.766** | 1.00 | .863 | .947 | **.894** | **.471** | .779 |
| Qwen2.5-72B-Instruct† | .676 | **.676** | .935 | .917 | .500 | 1.00 | .503 | .763 | **.738** | 1.00 | .898 | .947 | **.921** | **.428** | .752 |
| GPT-4o-mini | .499 | **.499** | .960 | .861 | .444 | 1.00 | .489 | .704 | **.725** | 1.00 | .883 | .947 | **.909** | **.395** | .729 |

[*] **Intent Recognition metrics:** 1: Follow instruction under distraction, OA: Intent Recognition metrics overall accuracy **Granular Content Validation metrics:** 3: Theme requirement, 4: Stylistic requirement, 5: Language requirement, 6: Format requirement, 7: Word count requirement, 8: Other granular requirements, OA: Granular Content Validation metrics overall accuracy **Output Structure Validation metrics:** 10: Format requirement, 11: Element number requirement, 12: Logic requirement, OA: Output Structure Validation overall accuracy **Results metrics:** UR3: Turn-3 Utility Rate, MS3: Turn-3 Meeseeks Score

refinement, models often execute the instruction directly instead of refining it. Our case analysis reveals that reasoning models' response patterns are heavily influenced by prompt structure, particularly the positioning of refinement requests relative to the instruction itself. For example, the model's behavior varies significantly depending on whether the refinement request precedes or follows the instruction. This suggests that anti-distraction capabilities may be more closely tied to model training approaches rather than inherent reasoning abilities. **Moreover,** we find that most models struggle with two critical requirement tags. First, models show unsatisfactory performance on Language Requirements, which involve complex language constraints such as specific language ratios (e.g., 1:1 Chinese-English ratio) and hybrid language rules (e.g., English nouns within Chinese text). Second, models demonstrate poor performance in Word Count Requirements, which include range-based limits (e.g., 100-200 words) and exact count specifications (e.g., exactly 10 words). Notably, word count control remains a common yet challenging problem in current LLM development, particularly given the diverse and complex nature of these requirements across different contexts.

## 5 Conclusion & Limitations

In this work, we propose Meeseeks, an iterative feedback framework that simulates realistic human-LLM interactions by enabling models to self-correct based on specific requirement failures in each turn. Meeseeks incorporates 38 capability tags across 3 evaluation dimensions, effectively mapping models' instruction-following abilities to their performance metrics. Through extensive experiments with Meeseeks, we evaluate the instruction-following capabilities of current representative LLMs and reveal their distinct performance patterns under both single-turn and multi-turn scenarios, bringing inspire fresh perspectives in LLM evaluation. In conclusion, Meeseeks addresses a critical gap in understanding instruction-following capabilities within multi-turn interactions.

The main limitation lies in the deployment costs of scoring models capable of code-guided rule-augmented LLM-based evaluation, particularly those requiring code extraction. Our experiments show that Qwen2.5-Coder-32B-Instruct is the minimum viable model for this task, with Claude-3.5-Sonnet demonstrating superior performance. Additionally, while our current quick-start dataset, generated through templates, offers broad coverage, it remains relatively monotonous. Future work will focus on enhancing the Meeseeks framework and dataset to achieve more comprehensive coverage. Finally, the end-to-end accuracy discussed in Section 3.2 is specific to the current Meeseeks quick-start dataset and should not be considered a universal accuracy metric. Researchers using their own datasets may need to design customized scoring models' evaluation prompts to optimize accuracy for their specific use cases.

# References

[1] Abhinav Arun, Ashish Dhiman, Mehul Soni, and Yibei Hu. Numerical reasoning for financial reports, 2023.

[2] David Castillo-Bolado, Joseph Davidson, Finlay Gray, and Marek Rosa. Beyond prompts: Dynamic conversational benchmarking of large language models, 2024.

[3] DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanjia Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025.

[4] DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jiawei Wang, Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, Junxiao Song, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Litong Wang, Liyue Zhang, Meng Li, Miaojun Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu Chen, Qiushi Du, R. J. Chen, R. L. Jin, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin Xu, Ruoyu Zhang, Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng Ye, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou, Shuting Pan, T. Wang, Tao Yun, Tian Pei, Tianyu Sun, W. L. Xiao, Wangding Zeng, Wanjia Zhao, Wei An, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, X. Q. Li, Xiangyue Jin, Xianzu Wang, Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xiaokang Chen, Xiaokang Zhang, Xiaosha Chen, Xiaotao Nie, Xiaowen Sun, Xiaoxiang Wang, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xingkai Yu, Xinnan Song, Xinxia Shan, Xinyi Zhou, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, Y. K. Li, Y. Q. Wang, Y. X. Wei, Y. X. Zhu, Yang Zhang, Yanhong Xu, Yanhong Xu, Yanping Huang, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Li, Yaohui Wang, Yi Yu, Yi Zheng, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Ying Tang, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo,

Yu Wu, Yuan Ou, Yuchen Zhu, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yukun Zha, Yunfan Xiong, Yunxian Ma, Yuting Yan, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Z. F. Wu, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhen Huang, Zhen Zhang, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhibin Gou, Zhicheng Ma, Zhigang Yan, Zhihong Shao, Zhipeng Xu, Zhiyu Wu, Zhongyu Zhang, Zhuoshu Li, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Ziyi Gao, and Zizheng Pan. Deepseek-v3 technical report, 2025.

[5] Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B Hashimoto. Length-controlled alpacaeval: A simple way to debias automatic evaluators. *arXiv preprint arXiv:2404.04475*, 2024.

[6] Deep Ganguli, Amanda Askell, Nicholas Schiefer, Thomas I. Liao, Kamilė Lukošiūtė, Anna Chen, Anna Goldie, Azalia Mirhoseini, Catherine Olsson, Danny Hernandez, Dawn Drain, Dustin Li, Eli Tran-Johnson, Ethan Perez, Jackson Kernion, Jamie Kerr, Jared Mueller, Joshua Landau, Kamal Ndousse, Karina Nguyen, Liane Lovitt, Michael Sellitto, Nelson Elhage, Noemi Mercado, Nova DasSarma, Oliver Rausch, Robert Lasenby, Robin Larson, Sam Ringer, Sandipan Kundu, Saurav Kadavath, Scott Johnston, Shauna Kravec, Sheer El Showk, Tamera Lanham, Timothy Telleen-Lawton, Tom Henighan, Tristan Hume, Yuntao Bai, Zac Hatfield-Dodds, Ben Mann, Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, Christopher Olah, Jack Clark, Samuel R. Bowman, and Jared Kaplan. The capacity for moral self-correction in large language models, 2023.

[7] Luyu Gao, Zhuyun Dai, Panupong Pasupat, Anthony Chen, Arun Tejasvi Chaganty, Yicheng Fan, Vincent Y. Zhao, Ni Lao, Hongrae Lee, Da-Cheng Juan, and Kelvin Guu. Rarr: Researching and revising what language models say, using language models, 2023.

[8] Qianyu He, Jie Zeng, Wenhao Huang, Lina Chen, Jin Xiao, Qianxi He, Xunzhe Zhou, Lida Chen, Xintao Wang, Yuncheng Huang, Haoning Ye, Zihan Li, Shisong Chen, Yikai Zhang, Zhouhong Gu, Jiaqing Liang, and Yanghua Xiao. Can large language models understand real-world complex instructions?, 2024.

[9] Yun He, Di Jin, Chaoqi Wang, Chloe Bi, Karishma Mandyam, Hejia Zhang, Chen Zhu, Ning Li, Tengyu Xu, Hongjiang Lv, Shruti Bhosale, Chenguang Zhu, Karthik Abinav Sankararaman, Eryk Helenowski, Melanie Kambadur, Aditya Tayade, Hao Ma, Han Fang, and Sinong Wang. Multi-if: Benchmarking llms on multi-turn and multilingual instructions following, 2024.

[10] Bairu Hou, Qibin Chen, Jianyu Wang, Guoli Yin, Chong Wang, Nan Du, Ruoming Pang, Shiyu Chang, and Tao Lei. Instruction-following pruning for large language models, 2025.

[11] 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, 2024.

[12] Alex Kim, Maximilian Muhn, and Valeri Nikolaev. Financial statement analysis with large language models, 2025.

[13] Aviral Kumar, Vincent Zhuang, Rishabh Agarwal, Yi Su, John D Co-Reyes, Avi Singh, Kate Baumli, Shariq Iqbal, Colton Bishop, Rebecca Roelofs, Lei M Zhang, Kay McKinney, Disha Shrivastava, Cosmin Paduraru, George Tucker, Doina Precup, Feryal Behbahani, and Aleksandra Faust. Training language models to self-correct via reinforcement learning, 2024.

[14] Dong-Ho Lee, Jay Pujara, Mohit Sewak, Ryen W. White, and Sujay Kumar Jauhar. Making large language models better data creators, 2023.

[15] Jinnan Li, Jinzhe Li, Yue Wang, Yi Chang, and Yuan Wu. Structflowbench: A structured flow benchmark for multi-turn instruction following, 2025.

[16] Yunxiang Li, Zihan Li, Kai Zhang, Ruilong Dan, Steve Jiang, and You Zhang. Chatdoctor: A medical chat model fine-tuned on a large language model meta-ai (llama) using medical domain knowledge, 2023.

[17] Xiao Liu, Xuanyu Lei, Shengyuan Wang, Yue Huang, Zhuoer Feng, Bosi Wen, Jiale Cheng, Pei Ke, Yifan Xu, Weng Lam Tam, Xiaohan Zhang, Lichao Sun, Xiaotao Gu, Hongning Wang, Jing Zhang, Minlie Huang, Yuxiao Dong, and Jie Tang. Alignbench: Benchmarking chinese alignment of large language models, 2024.

[18] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. Self-refine: Iterative refinement with self-feedback, 2023.

[19] Potsawee Manakul, Adian Liusie, and Mark J. F. Gales. Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models, 2023.

[20] Ning Miao, Yee Whye Teh, and Tom Rainforth. Selfcheck: Using llms to zero-shot check their own step-by-step reasoning, 2023.

[21] Shervin Minaee, Tomas Mikolov, Narjes Nikzad, Meysam Chenaghlu, Richard Socher, Xavier Amatriain, and Jianfeng Gao. Large language models: A survey, 2025.

[22] Shervin Minaee, Tomas Mikolov, Narjes Nikzad, Meysam Chenaghlu, Richard Socher, Xavier Amatriain, and Jianfeng Gao. Large language models: A survey, 2025.

[23] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback, 2022.

[24] Keivalya Pandya and Mehfuza Holia. Automating customer service using langchain: Building custom open-source gpt chatbot for organizations, 2023.

[25] Debjit Paul, Mete Ismayilzada, Maxime Peyrard, Beatriz Borges, Antoine Bosselut, Robert West, and Boi Faltings. Refiner: Reasoning feedback on intermediate representations, 2024.

[26] Baolin Peng, Michel Galley, Pengcheng He, Hao Cheng, Yujia Xie, Yu Hu, Qiuyuan Huang, Lars Liden, Zhou Yu, Weizhu Chen, and Jianfeng Gao. Check your facts and try again: Improving large language models with external knowledge and automated feedback, 2023.

[27] Cheng Qian, Chi Han, Yi R. Fung, Yujia Qin, Zhiyuan Liu, and Heng Ji. Creator: Tool creation for disentangling abstract and concrete reasoning of large language models, 2024.

[28] Yiwei Qin, Kaiqiang Song, Yebowen Hu, Wenlin Yao, Sangwoo Cho, Xiaoyang Wang, Xuansheng Wu, Fei Liu, Pengfei Liu, and Dong Yu. Infobench: Evaluating instruction following ability in large language models, 2024.

[29] J. Qiu, K. Lam, G. Li, et al. Llm-based agentic systems in medicine and healthcare. *Nature Machine Intelligence*, 6(12):1418–1420, 2024.

[30] Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report, 2025.

[31] Alessandro Stolfo, Vidhisha Balachandran, Safoora Yousefi, Eric Horvitz, and Besmira Nushi. Improving instruction-following in language models through activation steering, 2025.

[32] Yuchong Sun, Che Liu, Kun Zhou, Jinwen Huang, Ruihua Song, Wayne Xin Zhao, Fuzheng Zhang, Di Zhang, and Kun Gai. Parrot: Enhancing multi-turn instruction following for large language models, 2024.

[33] Bosi Wen, Pei Ke, Xiaotao Gu, Lindong Wu, Hao Huang, Jinfeng Zhou, Wenchuang Li, Binxin Hu, Wendy Gao, Jiaxin Xu, Yiming Liu, Jie Tang, Hongning Wang, and Minlie Huang. Benchmarking complex instruction-following with multiple constraints composition, 2024.

[34] Zhenyu Wu, Qingkai Zeng, Zhihan Zhang, Zhaoxuan Tan, Chao Shen, and Meng Jiang. Large language models can self-correct with key condition verification, 2024.

[35] Jochen Wulf and Juerg Meierhofer. Exploring the potential of large language models for automation in technical customer service, 2024.

[36] Shunyu Yao, Howard Chen, Austin W. Hanjie, Runzhe Yang, and Karthik Narasimhan. Collie: Systematic construction of constrained text generation tasks, 2023.

[37] Zhiyuan Zeng, Jiatong Yu, Tianyu Gao, Yu Meng, Tanya Goyal, and Danqi Chen. Evaluating large language models at evaluating instruction following, 2024.

[38] Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tianwei Zhang, Fei Wu, and Guoyin Wang. Instruction tuning for large language models: A survey, 2024.

[39] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. A survey of large language models, 2025.

[40] Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. Instruction-following evaluation for large language models, 2023.

[41] Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. Instruction-following evaluation for large language models, 2023.
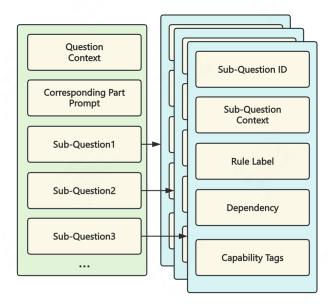
Figure 4: data-structure

# A Data structure

In this section, we delineate the composition of a singular data, providing a comprehensive overview of its constituent elements. The general arrangement and interrelations of all components within a data are visually represented in Figure 4 and a sample data entry is provided in Table 6.

**Question Context** is the first turn prompt to the under-evaluation LLM.

**Corresponding Part Prompt** is an prompt guiding the LLM-Extractor to extract the to-be-evaluated part from the under-evaluation LLM's response. An example is shown below.

Table 5: Corresponding Part Prompt example

| | |
|---|---|
| **Original data** | 请你按照python list的格式，抓取模型回复中，给出的所有歌词，要分割，例如: ["我姓石，何你相我都值", "我姓石，如在奔", "我姓石，心往神在落", "我姓石，高人如其名很"] |
| **English translation for reference** | Please extract all lyrics provided by the model's reply in the format of a Python list, with each lyric separated, for example: ["My surname is Shi, whenever I meet you it's worthwhile", "My surname is Shi, writing songs as if a steed galloping", "My surname is Shi, my heart flies as I write, recalling and writing", "My surname is Shi, with high vision and insight, true to my name and practical"] |

**Sub-Question:** The **Question Context** generally contains more than one requirements. Each requirement is corresponding to one **Sub-Question**. An example is shown below.

**Sub-Question Context** is the prompt input to the LLM Evaluator to judge if the under-evaluation LLM's response meets the corresponding requirement.

**Rule Label** is a special label projecting to different rule-based evaluation function inside the Meeseeks framework. For example, `"item_count:[80,80]"` checks whether the number of the elements(comments in the given example) is exact 80.

**Dependency** shows the error propagation between some of the sub-questions. For example, sub-question 2 relies on sub-question 0. If the under-evaluation LLM fails on sub-question 0, it automatically fails on the sub-question 2.

Table 6: Breakdown of sub-questions

| Question Context | Sub-Question Attributions | Sub-Question 0 | Sub-Question 1 | Sub-Question 2 | Sub-Question 3 |
|---|---|---|---|---|---|
| **Original data** 在美团外卖神券短视频的评论区，以消费者角度生成80条用户口语化评论。每条评论为10个字，不允许重复。 备注:字数只计中文字符 **English translation** Generate 80 colloquial user comments from a consumer perspective in the comments section of Meituan takeaway coupon short videos. Each comment should be 10 characters long and not repeated. Note: Only count Chinese characters | **Sub-Question ID** | 0 | 1 | 2 | 3 |
| | **Sub-Question Context** | **Original data** 生成的内容是否为美团外卖神券的评论? **English translation** Are the generated contents comments about Meituan takeaway coupons? | **Original data** 是否生成了80条用户口语化评论 **English translation** Are there 80 colloquial user comments generated? | **Original data** 所有评论是否为10个字 **English translation** Are all comments 10 characters long? | **Original data** 评论是否不重复 **English translation** Are the comments non-repetitive? |
| | **Rule Label** | Null | item_count: [80, 80] | each_length: [10,10] | non_repeat |
| | **Dependency** | [] | [0] | [0] | [0] |
| | **Capability Tag** | 主题约束 Theme requirement | 单元数量合规 Element number requirement | 0～10字 Generate in 0~10 words 精确 Generate at accurate word number | 重复 Generate repeat/non-repeat content |
| | **Corresponding Part** | Null | comments | comments | comments |

Table 7: Instruction-following benchmarks evaluation system comparison

| Meeseeks | Complexbench | CELLO | IF-Eval |
|---|---|---|---|
| Intent Recognition | Format Constraint | Task-prescribed Phrases | Keywords |
| Granular Content Validation | Lexical Constraint | Count Limit | Language |
| Output Structure Validation | Semantic Constraint | Input-dependent Query | Length Constraints |
| | Utility Constraint | Answer Format | Detectable Content |
| | | | Detectable Format... |

**Capability Tag** reflects certain capability of under-evaluation LLM. It is attached with the sub-question. This will be further explained in Appendix B.

## B  Capability tags

Meeseeks is equipped with a systematic evaluation system. Each sub-question is mounted with one or more capability tags, which represent the most granular level of evaluation. All capability tags are enumerated in Table 8, with the count of capability tags indicated in parentheses following each tag. Detailed explanations and definitions for each capability tag will be provided on GitHub, ensuring clarity and accessibility for researchers utilizing this system for evaluating model capabilities.

Unlike existing tag-based evaluations (Table 7), Meeseeks establishes a direct mapping between instruction-following abilities and benchmark performance through this systematic framework.

## C  Sample data parameterized template

Following is a sample parameterized template:

生成关于{THEME}短视频的{ITEM_NUM}条用户评价，每条用户评价都必须包含以下关键词：{KEYWORDS}。

An English translation for reference:

```
In {THEME} short video, generate {ITEM_NUM} user comments.
The comments cannot contain {KEYWORDS}.
```

Table 8: Meeseeks Evaluation System

| Level 3 | Level 2 | Level 1 |
|---|---|---|
| 任务意图理解: **Intent Recognition (83)** | 在干扰下完成指令: Follow instruction under distraction (83) | |
| 单元细节合规: **Intent Recognition (83)Granular Content Validation (2817)** | 主题约束: Theme requirement (1011) | |
| | 文体约束: Stylistic requirement (36) | 生成特定文案: Generate in certain style (12)<br>生成名字/标题: Generate names/titles (24) |
| | 语言约束: Language requirement (36) | 中英文混杂: Generate Chinese-English-mixed article (24)<br>繁体约束: Generate in traditional/simple Chinese (12) |
| | 格式约束: Granular format requirement (9) | 特定格式: Generate in other format (9)<br>日期格式: Generate result in date-format (0) |
| | 字数约束: Word count requirement (982) | 精确: Generate at accurate word number (79)<br>范围: Generate in rough/range word number (190)<br>倍数: Generate in X times word number of reference text (22)<br>多对象: Generate multiple results under certain word requirement (33)<br>0～10字: Generate in 0～10 words (88)<br>10～50字: Generate in 10～50 words (200)<br>50～200字: Generate in 50～200 words (218)<br>200字以上: Generate in above 200 words (152) |
| | 其他特殊规则: Other granular requirements (743) | 押韵: Generate rhyming content (30)<br>关键词: Generate with certain keywords (263)<br>重复: Generate repeat/non-repeat content (225)<br>平仄: Generate with Chinese pingze rules (10)<br>接龙: Generate with Chinese jielong rules (2)<br>emoji: Generate with emoji (38)<br>符号: Generate with/without punctuation (6)<br>写作手法: Generate with certain rhetoric (69)<br>词频: Generate with certain number of word X (100) |
| 整体结构合规: **Output Structure Validation (650)** | 模版合规: Output format requirement (41) | JSON格式: Generate in JSON format (41) |
| | 单元数量合规: Element number requirement (590) | |
| | 答题逻辑合规: Output logic requirement (19) | 答题结构合规: Generate by certain steps (19) |