

Towards Scalable Oversight: Meta-Evaluation of LLMs as Evaluators via Agent Debate

Steffi Chern^{2,4}, Ethan Chern^{1,4}, Graham Neubig², and Pengfei Liu^{1,3,4}

¹ Shanghai Jiao Tong University, Minhang District, Shanghai 200240, China

² Carnegie Mellon University, Pittsburgh PA 15213, USA

³ Shanghai Artificial Intelligence Laboratory, Xuhui District, Shanghai 200232, China

⁴ Generative Artificial Intelligence Lab, Minhang District, Shanghai 200240, China

Abstract. Despite the utility of Large Language Models (LLMs) across a wide range of tasks and scenarios, developing a method for reliably evaluating LLMs across varied contexts continues to be challenging. Modern evaluation approaches often use LLMs to assess responses generated by LLMs. However, existing meta-evaluation methods to assess the effectiveness of LLMs as evaluators is typically constrained by the coverage of existing benchmarks or require extensive human annotation. This underscores the urgency of methods for *scalable* meta-evaluation that can effectively, reliably, and efficiently evaluate the performance of LLMs as evaluators across diverse tasks and scenarios, particularly in potentially new, user-defined scenarios. To fill this gap, we propose SCALEEVAL, an *agent-debate-assisted meta-evaluation framework* that leverages the capabilities of multiple communicative LLM agents. This framework supports multi-round discussions to assist humans in discerning the capabilities and limitations of LLMs as evaluators, which significantly reduces their workload in cases that used to require much supervision and large-scale annotations during meta-evaluation. We release the code for our framework, which is publicly available at:

<https://github.com/GAIR-NLP/scaleeval>.

Keywords: meta-evaluation · multi-agent debate · human annotation.

1 Introduction

While LLMs [30,31] have unlocked a variety of exciting potential applications, they have also introduced complex challenges in evaluating the generated outputs. Current efforts on LLM evaluation primarily focus on automated evaluation metrics [10,6,7,9], many of which use LLMs themselves to do evaluation. However, when these LLMs as evaluators are applied to a new task, it begs the question: *can LLMs be trusted for evaluation?* In many cases, the answer is not clear.

There are still a few fortunate tasks where meta-evaluation (evaluation of evaluation metrics) has been performed rigorously, as shown in Related Works. This typically involves the collection of human-annotated judgements for particular criteria (e.g. fluency of outputs, semantic adherence to the input). For instance, there is an extensive meta-evaluation dataset from the WMT metrics task [18] for machine translation quality metrics, and datasets like TAC and RealSum [33,32] for summarization. Once such

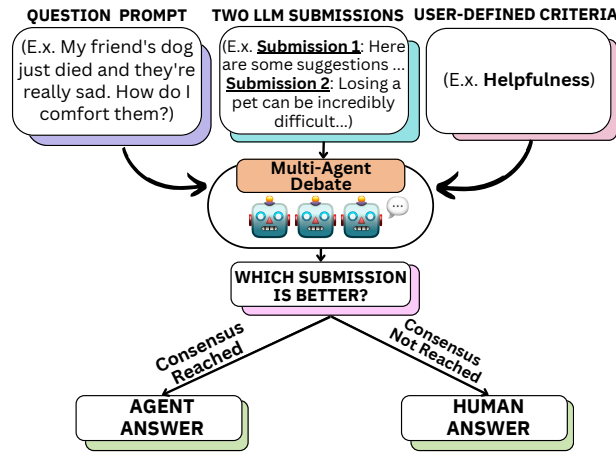


Fig. 1: We demonstrate SCALEEVAL, our scalable meta-evaluation framework. This is used in assessing the reliability and robustness of employing LLMs as evaluators for different evaluative purposes.

a dataset is collected, meta-evaluation can be performed by measuring the correlation between automatic evaluation metrics and the human gold-standard.

However, these datasets are extremely costly to collect, as they require large amounts of annotations by skilled human experts. With the increasing use of LLMs for various purposes such as math problem solving [25], reading comprehension [4], creative writing [7], multilingual applications [3,5], and many more, it is not feasible to create these human-judged datasets for every new task. As a result, LLMs as evaluators are used without proper vetting, and in many cases the evaluators themselves are highly unreliable [23,34].

In this paper, we propose SCALEEVAL, a *scalable meta-evaluation framework* for the era of LLMs, which creates meta-evaluation benchmarks across various tasks and scenarios. Concretely, SCALEEVAL relies on debate between multiple LLM agents, followed by minimal human oversight in cases where the agent LLMs do not agree (Fig. 1). Since our framework allows users to use their own prompts and responses while applying the framework to any scenario or criterion that they define, it offers flexibility and adaptability in various evaluation contexts.

In experiments, we conduct meta-meta evaluation, demonstrating that SCALEEVAL correlates well with when meta-evaluation is performed entirely by human expert annotators. We assess the reliability and cost-performance trade-off of various LLMs as evaluators under a variety of scenarios, and examine their specific capabilities and limitations as evaluators. We also examine the impact that variations in criteria prompts have on the performance of LLMs as evaluators.

	Meta-Eval	# Scenarios	Custom.	Scala.
LLM-as-a-Judge	Human	High	✗	Low
FairEval	Human	Low	✗	Low
ChatEval	Human	Low	✗	Low
SCALEEVAL	Agent Debate	High	✓	High

Table 1: Comparison of the meta-evaluation processes across different strategies using LLMs as evaluators: LLM-as-a-Judge [7], FairEval [23], ChatEval [11], and our own work, SCALEEVAL. High/low in scenarios refers to how many real-world scenarios can be evaluated. “Custom.” denotes whether the evaluation criterion could be customized. “Scala.” refers to scalability.

2 Related Works

2.1 Automatic Evaluation of LLM Output

The most common paradigm for evaluating LLMs is to evaluate their capabilities on standard benchmarks for tasks such as reasoning (e.g. BigBench [13]), common sense QA (e.g. MMLU [14]), or code generation (e.g. HumanEval [35]). These are indicative of the capabilities of the models, but do not measure model abilities for open-ended tasks requiring generation of free-form text.

To adapt to the rapid growth in the capabilities of LLMs for open-ended tasks, LLM evaluation has started to shift towards evaluating generated text directly, often using LLMs themselves as evaluators [10,6,7,9]. In addition, there are a few recent works that perform LLM-based multi-agent debate to improve the fidelity of evaluation [11,12]. While these methods take advantage of the instruction-following capabilities and versatility of LLMs, directly using LLMs as evaluators or communicative agents out-of-the-box in diverse, unseen user-defined scenarios provides no guarantees with respect to the accuracy of these methods.

Another widely used evaluation platform, Chatbot Arena [7], gathers diverse user prompts through crowd-sourcing for assessing LLMs’ performance. However, its heavy reliance on human evaluations, which are not universally accessible and lack standardized evaluation guidelines, may lead to biased or inconsistent assessments. We aim to address these issues by introducing scalable meta-evaluation to ensure the reliability of the evaluation protocol under diverse scenarios.

2.2 Meta-Evaluation of LLMs as Evaluators

Previous research on LLMs as evaluators usually involve conducting meta-evaluation in 3 different ways: (i) leveraging existing NLP meta-evaluation benchmarks [10,11], (ii) conducting small-scale meta-evaluations on expert-annotated datasets for specific tasks or scenarios [29,9,7], or (iii) using crowd-sourcing platforms to collect human annotations [7]. With the lack of coverage in existing datasets and benchmarks, both (i) and (ii) are inherently limited in their comprehensiveness. While (iii) can be more

comprehensive via crowd-sourcing, the amount of human annotation required limits the scalability of the approach, and crowd workers may not be accurate at more complex tasks. Thus, we propose an agent-debate-assisted meta-evaluation approach to mitigate these issues.

3 Preliminaries

In this section, we provide an introduction to the concepts of automatic evaluation and meta-evaluation systems, particularly focused on evaluation of LLM-generated outputs in the era of generative AI.

3.1 Key Terms

We first define some key terms that will be used throughout our paper.

- **Criterion:** A standard that measures the quality of the response generated by LLMs based on the user prompt. Some examples include: helpfulness, fluency, factuality, or creativity, among others.
- **Scenario:** A scenario describes the real-world situations in which users are interacting with LLMs. For example, brainstorming, coding, and dialog, among others.

3.2 Automatic Evaluation

Automatic evaluation using LLMs measures the quality of LLM-generated responses given prompts under different criteria, which is conducted with one of two different protocols: single-response evaluation and pairwise response comparison [19,7,17]. In this paper, we focus on **pairwise response comparison**. Pairwise response comparison is intuitive for both humans and LLMs as evaluators when conducting assessments. It could be further extended to provide win-rates and Elo scores across models [7], offering a straightforward leaderboard to understand the relative performance of different models under various scenarios. Formally, given an automatic evaluation metric E , a user-defined evaluation criterion c (e.g. helpfulness, reasoning, creativity), a user prompt p , and responses generated by two systems r_1, r_2 , evaluation for pairwise response comparison is done in the following way:

$$o = E(c, p, r_1, r_2). \quad (1)$$

where $o \in \{1, 0, -1\}$ represents that r_1 is better, equal, or worse than r_2 , respectively, given the user prompt p under criterion c .

3.3 Meta-Evaluation

Meta-evaluation assesses the quality of an automatic evaluation metric. Formally, we define a gold-standard evaluation metric G (e.g. human experts) that other automatic metrics should aspire to match. In pairwise response comparison, the meta-evaluation dataset $\mathcal{G} = \{G(c, p_i, r_{1,i}, r_{2,i})\}_{i=1}^n$ contains user prompts and corresponding responses

from two systems, annotated with gold-standard evaluations. The meta-evaluation process assesses the performance $\text{META}(E)$ of the automatic evaluation metric E under a certain criterion c .

In pairwise response comparison, the meta-evaluation measures the *example-level agreement rate* or the *system-level agreement rate* between E and G across the meta-evaluation dataset. A high agreement rate between E and G represents that E is a good automatic evaluation metric.

For the *example-level agreement rate*, we calculate:

$$\text{META}(E) = \frac{1}{n} \sum_{i=1}^n \delta_{E(c, p_i, r_{1,i}, r_{2,i}), G(c, p_i, r_{1,i}, r_{2,i})} \quad (2)$$

where $0 \leq \text{META}(E) \leq 1$, and $\delta_{\cdot, \cdot}$ refers to the Kronecker delta function.

For the *system-level agreement rate*, given that

$$\mathcal{E} = \{E(c, p_i, r_{1,i}, r_{2,i})\}_{i=1}^n, \quad (3)$$

$$\mathcal{G} = \{G(c, p_i, r_{1,i}, r_{2,i})\}_{i=1}^n \quad (4)$$

we calculate:

$$\text{META}(E) = \delta_{\text{mode}(\mathcal{E}), \text{mode}(\mathcal{G})} \quad (5)$$

where $\text{META}(E) \in \{0, 1\}$, $\delta_{\cdot, \cdot}$ refers to the Kronecker delta function, and $\text{mode}(\cdot)$ refers to the value (either 1, 0, -1 in this case) that appears most often in the set \mathcal{E} or \mathcal{G} .

4 Methodology

In this section, we detail the frameworks that SCALEEVAL employs for meta-evaluation, evaluation, and human expert meta-meta evaluation. For meta-evaluation, we follow its pairwise response comparison setting mentioned previously. Notably, instead of relying solely on human labor to construct the meta-evaluation benchmark \mathcal{G} , we use a scalable, agent-debate assisted framework to instantiate the golden metric G and construct the benchmark \mathcal{G} . For evaluation, we also follow its corresponding pairwise response comparison setting. The human expert meta-meta evaluation process follows the rules for meta-evaluation. The process is included to ensure the reliability of using the agent-debate assisted meta-evaluation framework.

4.1 Meta-Evaluation Framework via Multi-Agent Debate

The meta-evaluation framework involves multiple communicative agents $\{A_j\}_{j=1}^m$ that conduct rounds of discussion $d = 0 \sim D - 1$ with each other. This is less time-consuming and costly compared to traditional methods for meta-evaluation that relies entirely on human effort. With this agent-debate-assisted meta-evaluation framework, we can leverage each LLM agent’s distinct understanding about each query prompt p_i , LLM responses $r_{1,i}, r_{2,i}$, and defined criterion c to make a comprehensive assessment of LLMs under different scenarios and criteria. Each LLM agent is capable of providing an evaluation result regarding which response is better, along with its corresponding

justifications. Note that each LLM agent can also review other agents’ evaluation results and justifications after the initial round of discussion.

In the initial round of discussion $d = 0$, each LLM agent independently provides an evaluation result and justification:

$$\mathcal{A}_0 = [A_1(c, p_i, r_{1,i}, r_{2,i}, \emptyset), \dots, A_m(c, p_i, r_{1,i}, r_{2,i}, \emptyset)], \quad (6)$$

where

$$\mathcal{A}_0[j]_{j=1,\dots,m} \in (\{1, 0, -1\}, \text{JUSTIFICATION}), \quad (7)$$

indicates whether $r_{1,i}$ is better, equal, or worse than $r_{2,i}$, respectively, along with its justification. Note that the \emptyset in the last argument of A_j represents that in the initial round of discussion, each agent doesn’t have access to previous rounds of discussion. In subsequent discussion rounds $d = 1 \sim D-1$, agents are allowed to look at other agents’ previous assessments and conduct re-evaluations, in which each agent is prompted to stick with or change their original evaluation result. Specifically, given \mathcal{A}_{d-1} ($d \geq 1$), which represents the evaluation results and justifications of agents after $(d-1)^{th}$ rounds of discussions, we conduct the d^{th} round of discussion:

$$\mathcal{A}_d = [A_1(c, p_i, r_{1,i}, r_{2,i}, \mathcal{A}_{d-1}), \dots, A_m(c, p_i, r_{1,i}, r_{2,i}, \mathcal{A}_{d-1})] \quad (8)$$

where similarly to \mathcal{A}_0 ,

$$\mathcal{A}_d[j]_{j=1,\dots,m} \in (\{1, 0, -1\}, \text{JUSTIFICATION}), \quad (9)$$

The detailed prompt template for meta-evaluation can be found in Appendix.

In cases where agents fail to reach a consensus after $d = D - 1$ rounds of discussions, a human evaluator intervenes. The human evaluator reviews the assessment reports provided by the agents and makes a final decision. Through this process, we incorporate an element of human oversight, thereby increasing the reliability of the final decision. This approach strikes a balance between efficiency and the need for human judgment, ensuring that evaluations are done in a timely and accurate manner. An example of multi-agent debate process during meta-evaluation is shown in Fig. 2.

4.2 Evaluation Framework

We follow the pairwise response comparison setting outlined in Automatic Evaluation under Preliminaries. Note that in the LLM era, the automatic evaluation metric E is often instantiated through single LLMs [10,6,7,9], or multi-agent debate [11,12]. In SCALEVAL, we focus on instantiating E through single LLMs (e.g., *gpt-3.5-turbo*). However, it is important to note that our framework can be further generalized to other instantiations of E .

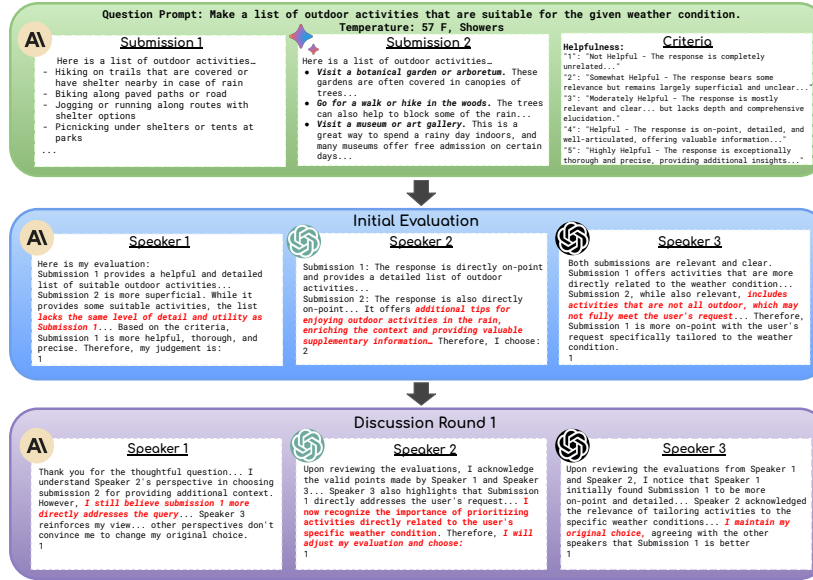


Fig. 2: An example of the multi-agent debate process during meta-evaluation.

4.3 Human Expert Meta-Meta Evaluation

To test the reliability of our proposed meta-evaluation framework, we apply meta-meta evaluation. The meta-meta evaluation process also follows the meta-evaluation process described in Preliminaries, where E is instantiated as the agent-debated assisted protocol, and G is instantiated as the human expert annotation protocol.

5 Experiments

5.1 Exp-I: Meta-Meta-Evaluation of SCALEEVAL

We first examine whether SCALEEVAL's results match with those from meta-meta-evaluation.

Setup For our SCALEEVAL meta-evaluation framework, we deploy three LLM agents to perform multi-agent debate: *gpt-4-turbo*, *claude-2*, and *gpt-3.5-turbo*.⁵ In our meta-evaluation experiment, we analyze a total of 160 prompts, with 137 prompts from AlpacaEval [6], 10 coding problem prompts from HumanEval [20], and 13 math problem prompts from GSM-Hard [21]. We categorize these prompts into four distinct scenarios: *brainstorming*, *coding*, *math*, and *writing*, where each scenario contains 40 prompts.

⁵ Results collected in December 2023. Specific models used are: gpt-4-1106-preview, claude-2, and gpt-3.5-turbo-1106.

LLM	Criterion	Scenario	GPT-4-Turbo	Claude-2	GPT-3.5-Turbo	GPT-4-Turbo	Claude-2	GPT-3.5-Turbo	Multi-LLM	Meta-
Comparisons			Single LLM			Self-Consistency			Consistency Evaluation	
GPT-3.5-Turbo vs. Claude-Instant	Helpfulness	Brainstorming	0.633	0.433	0.267	0.633	0.533	0.400	0.567	0.600
	Interpretability	Coding	0.700	0.533	0.567	0.733	0.667	0.600	0.733	0.600
	Reasoning	Math	0.600	0.400	0.367	0.733	0.467	0.433	0.667	0.867
	Creativity	Writing	0.667	0.400	0.333	0.667	0.400	0.400	0.667	0.700
Claude-Instant vs. Gemini-Pro	Helpfulness	Brainstorming	0.533	0.467	0.500	0.600	0.500	0.500	0.600	0.667
	Interpretability	Coding	0.600	0.500	0.567	0.600	0.533	0.633	0.567	0.833
	Reasoning	Math	0.667	0.330	0.367	0.733	0.467	0.433	0.500	0.767
	Creativity	Writing	0.633	0.400	0.500	0.700	0.400	0.467	0.567	0.733
GPT-3.5-Turbo vs. Gemini-Pro	Helpfulness	Brainstorming	0.600	0.467	0.467	0.667	0.600	0.500	0.700	0.733
	Interpretability	Coding	0.733	0.567	0.667	0.700	0.667	0.667	0.800	0.833
	Reasoning	Math	0.767	0.500	0.433	0.767	0.567	0.467	0.767	0.867
	Creativity	Writing	0.667	0.500	0.433	0.667	0.433	0.500	0.700	0.767

Table 2: Baseline experiments – example-level agreement rate comparison between human expert and single LLM evaluations, human expert and self-consistency [37], human expert and multi-LLM consistency, and human expert and SCALEEVAL’s meta-evaluation across four scenarios and criteria.

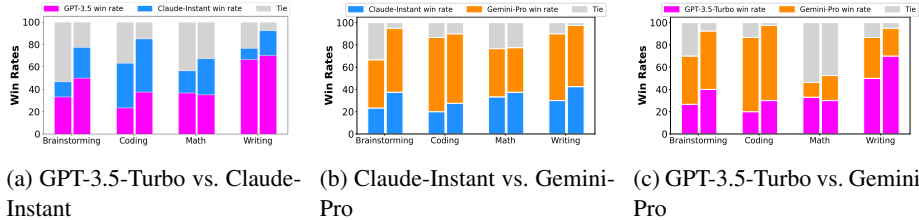


Fig. 3: System-level agreement – win rates for each LLM pairwise comparison. Left bars in each scenario represent human expert meta-meta evaluation results; right bars represent SCALEEVAL’s meta-evaluation results.

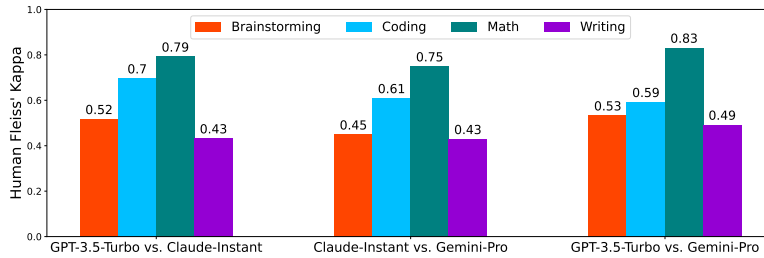


Fig. 4: Human Fleiss Kappa for each LLM pairwise comparison under four scenarios.

Each scenario is evaluated based on the following criteria, respectively: *helpfulness*, *interpretability*, *reasoning*, and *creativity*. We evaluate the generated responses from the following three LLMs: *gpt-3.5-turbo*, *claude-instant*, and *gemini-pro*. We select the above LLMs to evaluate due to their rather similar performances according to past research and public user feedback, which can help us establish a more nuanced

understanding of their performance in various real-world scenarios, and to identify specific contexts where one may outperform the others.

Our meta-meta evaluation involves having human experts annotate which LLM submission they think is better based on a defined criterion during pairwise comparisons. A total of seven human experts were selected from a pool of graduate students who have the relevant expertise in answering the queries in each scenario. Different groups of three human experts are responsible for answering the prompts in each scenario, where they were assigned to the scenario that relates to their expertise. Each expert received identical instructions for the task – they were asked to decide which submission is better based on our defined criteria, and for each comparison, label either 0 (*neither submission is better*), 1 (*submission 1 is better*), or 2 (*submission 2 is better*). The label 2 corresponds to the label -1 as denoted in Preliminaries. The experts were tasked to conduct 30 comparisons for each of the four different scenarios (*brainstorming, coding, math, and writing*), based on their corresponding defined criteria (*helpfulness, interpretability, reasoning, and creativity*). This results in a total of 120 final judgements. The question prompts, LLM responses, and criteria utilized for human expert annotations were consistent with those used during our meta-evaluation experiment. All details were presented in a google sheet that allowed experts to record their answers. Experts were compensated with food for their participation.⁶

Q1: Can LLM agents with multi-agent debate be used as meta-evaluators in new user-defined scenarios? To validate the reliability of SCALEEVAL, we perform comparisons between the results from human experts and SCALEEVAL’s multi-agent debate by two key metrics: the *example-level agreement rate* and the *system-level agreement rate*. The example-level agreement rate measures the proportion of instances where the multi-agent debate results correspond with the human experts judgements. The system-level agreement rate assesses whether the human experts and multi-agents concur in their overall evaluation of which LLMs produce the best responses for each scenario. A high agreement rate in both metrics would suggest a strong reliability and validity of our meta-evaluation framework, indicating that both human and LLM agents consistently recognize and agree on the quality of responses generated by LLMs. For our baselines, we employ single-LLM evaluations, self-consistency [37], and multi-LLM consistency. Self-consistency involves each evaluator separately generates evaluation results three times, and the final answer is the evaluation result that occurred with the highest probability. Multi-LLM consistency involves all evaluators engaging in two rounds of multi-agent debate, where the final evaluation result is determined by the evaluation result that occurred with the highest probability.

Results From Table 2, we generally observe a higher example-level agreement rate between human experts and SCALEEVAL (meta-evaluation), compared to the agreement rate between human experts and single LLM evaluations, self-consistency, and multi-LLM consistency. The consistently high agreement rates suggest that our meta-evaluation framework aligns well with human expert judgements in these areas, indi-

⁶ Human experts were compensated 150 USD in total. Inference costs for meta-evaluation were around 13 USD. Single-LLM baselines cost around 8 USD.

cating a reliable performance of the collective use of LLMs in meta-evaluating complex scenarios. Across all LLM submission comparisons in our experiment, we observe higher agreement rates in decisions between SCALEEVAL outcomes and those of human experts, particularly in coding and math scenarios. This could be attributed to the inherently objective nature of these subjects, which have relatively clear, definitive answers unlike more subjective areas like creative writing. We observe that the example-level agreement rates between human experts and SCALEEVAL consistently exceed the Fleiss Kappa scores (human-human example-level agreement rate), as illustrated in Fig. 4. This indicates the potential of SCALEEVAL as a promising framework for meta-evaluation of LLMs as evaluators, offering a reliable alternative to human evaluation.

To verify the effectiveness of our proposed method, we compare SCALEEVAL against existing methods that use LLMs as evaluators with the FairEval [23] benchmark, as shown in Table 3. The benchmark consists of 80 open-ended questions originating from a wide array of categories, including common-sense, counterfactual, and more. We adopt a similar evaluation setting as FairEval [23] and ChatEval [11]. We provide the accuracy of each method tested on the benchmark, which measures the proportion of questions with correct evaluation results (same as human annotations) out of all the questions available. We notice that our method, SCALEEVAL, achieves the highest accuracy at 68.8%, outperforming all other existing methods. ChatEval (Multi Agent) [11] comes in second at 63.8%, showing the advantage of multi-agent systems over single-agent approaches.

Criterion	Method	Accuracy
Helpfulness	FairEval	62.5
	ChatEval (Single Agent)	61.3
	ChatEval (Multi Agent)	63.8
	SCALEEVAL	68.8

Table 3: Accuracy comparison among existing methods that use LLMs as evaluators, FairEval [23] and ChatEval [11]. Above results are tested using the FairEval [23] benchmark with helpfulness criterion.

Based on Fig. 3, we notice a "preference in the same direction" between human experts and multi-agent debates across **all** LLM pairwise comparisons and scenarios. Notably, *gpt-3.5-turbo* is favored (higher win rates) in *brainstorming*, *math*, and *writing* scenarios when compared with *claude-instant*. Similarly, *gemini-pro* is also preferred over *claude-instant* in all scenarios. When comparing *gpt-3.5-turbo* with *gemini-pro*, a varied pattern in decision outcomes is observed: both human experts and multi-agent systems agree that *gpt-3.5-turbo* outperforms *gemini-pro* in scenarios involving *math* and *writing*. Conversely, *gemini-pro* is deemed superior in *brainstorming* and *coding* scenarios. The high agreement of multi-agent preferences with expert judgements ver-

ifies the reliability of using multiple LLMs agents as meta-evaluators in various user-defined scenarios.

5.2 Exp-II: Meta-Evaluation vs. LLM Evaluators

Next, we use the fact that SCALEEVAL allows for reliable and scalable meta-evaluation to examine the traits of LLMs as evaluators.

Q2: What are the capabilities and limitations of each LLM evaluator? We adopt an approach that involves comparing the outcomes from SCALEEVAL with the evaluations made independently by each LLM evaluator. In this process, we aim to identify which LLM evaluators demonstrate superior evaluative abilities and vice versa, thereby contributing to our understanding of their reliability in evaluating responses under each scenario. In addition, we provide a comprehensive cost-performance analysis to decide which LLM evaluator is the most suitable choice in each scenario.

Setup We employed 3 LLMs (*gpt-4-turbo*, *claude-2*, and *gpt-3.5-turbo*) as evaluators to perform pairwise comparisons of responses from 3 LLMs: *gpt-3.5-turbo*, *claude-instant*, and *gemini-pro*. Previous studies have highlighted the presence of positional biases when LLMs are used as evaluators [23]. Thus, we randomize the sequence in which submissions from LLMs are presented to the agent evaluators, as well as the order for agent-debate discussions. The meta-evaluations were done under 8 scenarios: *brainstorming*, *coding*, *dialog*, *judgement*, *open-domain general*, *open-domain science*, and *writing*, with the same set of 4 criteria used during human expert annotation.

Criterion	Scenario	GPT-4-Turbo	Claude-2	GPT-3.5-Turbo	Auto-J
Helpfulness	Brainstorming	0.800	0.500	0.650	0.575
	Coding	0.600	0.725	0.675	0.675
	Dialog	0.800	0.700	0.700	0.625
	Judgement	0.725	0.625	0.725	0.750
	Math	0.825	0.650	0.600	0.350
	ODG	0.850	0.525	0.575	0.700
	ODS	0.875	0.525	0.575	0.675
	Writing	0.750	0.600	0.750	0.600
Interpretability	Coding	0.825	0.600	0.550	0.525
Reasoning	Math	0.650	0.525	0.475	0.450
	Judgement	0.750	0.650	0.700	0.675
Creativity	Writing	0.775	0.600	0.575	0.650
	Brainstorming	0.800	0.525	0.550	0.625
	Dialog	0.875	0.750	0.700	0.800
Average	Overall	0.780	0.607	0.629	0.619

Table 4: Agreement rate between SCALEEVAL’s meta-evaluation and each LLM evaluator for comparing GPT-3.5-Turbo vs. Claude-Instant. ODG = Open-Domain General. ODS = Open-Domain Science.

Criteria Format	Criteria	Scenario	GPT-4-Turbo	Claude-2	GPT-3.5-Turbo
General	Helpfulness	Brainstorming	0.800	0.500	0.650
	Interpretability	Coding	0.825	0.600	0.550
	Reasoning	Math	0.650	0.525	0.475
	Creativity	Writing	0.800	0.600	0.575
Shortened	Helpfulness	Brainstorming	0.675	0.500	0.575
	Interpretability	Coding	0.675	0.325	0.425
	Reasoning	Math	0.625	0.425	0.400
	Creativity	Writing	0.675	0.250	0.525
Gibberish	Helpfulness	Brainstorming	0.575	0.450	0.575
	Interpretability	Coding	0.700	0.275	0.525
	Reasoning	Math	0.650	0.200	0.400
	Creativity	Writing	0.550	0.150	0.450
Shuffled	Helpfulness	Brainstorming	0.625	0.550	0.500
	Interpretability	Coding	0.600	0.400	0.525
	Reasoning	Math	0.625	0.225	0.600
	Creativity	Writing	0.625	0.275	0.500
Flipped	Helpfulness	Brainstorming	0.725	0.325	0.550
	Interpretability	Coding	0.725	0.425	0.300
	Reasoning	Math	0.575	0.250	0.500
	Creativity	Writing	0.750	0.075	0.550
Masked	Helpfulness	Brainstorming	0.725	0.300	0.500
	Interpretability	Coding	0.650	0.225	0.475
	Reasoning	Math	0.575	0.150	0.375
	Creativity	Writing	0.575	0.200	0.400

Table 5: Example-level agreement rate between SCALEEVAL’s meta-evaluation results and each LLM evaluator under various criteria prompt formats and scenarios comparing GPT-3.5-Turbo vs. Claude-Instant.

Results From Table 4, we observe *gpt-4-turbo* as the evaluator that has the highest agreement rates with our meta-evaluation, particularly in *brainstorming*, *dialog*, and *open-domain general* scenarios under the *helpfulness* criterion. It stands out with the highest overall average score of 0.780. However, our selected open-source model evaluator, *auto-j*, outperforms *gpt-4-turbo* in evaluating *coding* questions with the *helpfulness* criterion. Additionally, it exhibits the highest agreement rate with our meta-evaluation in the *judgement* scenario under the *helpfulness* criterion, indicating it as the most capable evaluator in this setting. It also achieves comparable results with other closed-source models like *claude-2* and *gpt-3.5-turbo* in most other scenarios. While *gpt-4-turbo* performs the best as an evaluator in most scenarios, it is not necessarily the best choice when we take into consideration its relatively high API costs. In fact, both the more affordable version (*gpt-3.5-turbo*) and our selected free, open-source model (*auto-j*) show comparable performance in scenarios like *judgement* and *writing*. For coding-related evaluations, the slightly less expensive *claude-2* could be a more cost-effective alternative to *gpt-4-turbo*.

5.3 Exp-III: Meta-Evaluation with Criteria Prompt Format Variations

Q3: How do the qualities of criteria prompts influence the robustness of LLMs as evaluators in different scenarios? Variations in prompts can substantially affect the

behavior of LLMs, particularly with the text they generate. Thus, we define various formatted criteria for evaluating LLM responses under each scenario. This examines the extent to which different formats of criteria prompts influence both the performance and robustness of LLMs as evaluators.

Setup We define 5 variations of the same criteria prompts: *shortened*, *gibberish*, *shuffled*, *flipped*, and *masked* (see Appendix A.2 for detailed prompt variations). We intend to observe how LLMs as evaluators would respond differently when conducting evaluation. We compare the example-level agreement rate between SCALEEVAL’s meta-evaluation results and each LLM evaluator.

Results Based on Table 5, the performance of LLMs as evaluators generally deteriorates when certain letters in the criteria prompts are masked. Furthermore, the removal of guiding phrases at the beginning, such as "Not Helpful" or "Highly Helpful", can also diminish their effectiveness as evaluators. Both *gpt-4-turbo* and *gpt-3.5-turbo* demonstrate some resilience to these adversarially formatted criteria prompts, maintaining a relatively consistent agreement rates across various criteria formats. In contrast, *Claude-2* often showcases confusion and refuses to evaluate, particularly in cases with gibberish and masked criteria prompts. It rejects answering about half of the questions, stating it lacks sufficient information to evaluate effectively. None of the LLMs as evaluators we tested maintained very similar evaluation capabilities when faced with these adversarially formatted criteria prompts, indicating a limitation in these LLMs as evaluators’ current design and application. Despite their advanced capabilities in fulfilling a variety of tasks, they may still struggle with understanding and responding accurately to substituted criteria information, highlighting an area for potential improvement in future iterations of LLMs. Among all the different formatted criteria, we highlight the cases where the LLMs perform the best as evaluators in Table 5.

6 Conclusion

In this work, we propose SCALEEVAL, a scalable, agent-debate assisted meta-evaluation framework for assessing the reliability and robustness of LLMs as evaluators. We address the expensive and time-intensive challenges inherent in traditional meta-evaluation methods, particularly pertinent as the usage of LLMs expands, necessitating a more scalable solution. We demonstrate the reliability of our proposed meta-evaluation framework, and shed light on the capabilities and limitations of LLMs as evaluators in various scenarios. We observe how the results from these LLMs as evaluators vary based on modifications to the same criteria prompts. By open-sourcing our framework, we aim to foster further research in this field and encourage the development of more advanced and reliable LLMs as evaluators in the future.

References

1. Fabbri, A.R., Kryściński, W., McCann, B., Xiong, C., Socher, R., Radev, D.: SummEval: Re-evaluating Summarization Evaluation. In: Transactions of the Association for Computational Linguistics, vol. 9, pp. 391–409. MIT Press (2021).

2. Wang, A., Singh, A., Michael, J., Hill, F., Levy, O., Bowman, S.R.: GLUE: A Multi-Task Benchmark and Analysis Platform for Natural Language Understanding. arXiv preprint arXiv:1804.07461v3 (2019).
3. Hu, J., Ruder, S., Siddhant, A., Neubig, G., Firat, O., Johnson, M.: XTREME: A Massively Multilingual Multi-task Benchmark for Evaluating Cross-lingual Generalization. arXiv:2003.11080v5 (2020).
4. Zhong, W., Cui, R., Guo, Y., Liang, Y., Lu, S., Wang, Y., Saied, A., Chen, W., Duan, N.: AGIEval: A Human-Centric Benchmark for Evaluating Foundation Models. arXiv:2304.06364v2 (2023).
5. Bang, Y., Cahyawijaya, S., Lee, N., Dai, W., Su, D., Wilie, B., Lovenia, H., Ji, Z., Yu, T., Chung, W., Do, Q.V., Xu, Y., Fung, P.: A Multitask, Multilingual, Multimodal Evaluation of ChatGPT on Reasoning, Hallucination, and Interactivity. arXiv:2302.04023v3 (2023).
6. Li, X., Zhang, T., Dubois, Y., Taori, R., Gulrajani, I., Guestrin, C., Liang, P., Hashimoto, T.B.: AlpacaEval: An Automatic Evaluator of Instruction-following Models. GitHub repository (2023).
7. Zheng, L., Chiang, W.L., Sheng, Y., Zhuang, S., Wu, Z., Zhuang, Y., Lin, Z., Li, Z., Li, D., Xing, E.: Judging LLM-as-a-judge with MT-Bench and Chatbot Arena. arXiv preprint arXiv:2306.05685 (2023).
8. Wang, P., Li, L., Chen, L., Zhu, D., Lin, B., Cao, Y., Liu, Q., Liu, T., Sui, Z.: Large Language Models are Not Fair Evaluators. arXiv preprint arXiv:2305.17926 (2023).
9. Wang, J., Liang, Y., Meng, F., Shi, H., Li, Z., Xu, J., Qu, J., Zhou, J.: Is ChatGPT a Good NLG Evaluator? A Preliminary Study. arXiv preprint arXiv:2303.04048 (2023).
10. Fu, J., Ng, S.K., Jiang, Z., Liu, P.: GPTScore: Evaluate as You Desire. arXiv preprint arXiv:2302.04166 (2023).
11. Chan, C.M., Chen, W., Su, Y., Yu, J., Xue, W., Zhang, S., Fu, J., Liu, Z.: ChatEval: Towards Better LLM-based Evaluators through Multi-Agent Debate. arXiv preprint arXiv:2308.07201 (2023).
12. Li, R., Patel, T., Du, X.: PRD: Peer Rank and Discussion Improve Large Language Model Based Evaluations. arXiv preprint arXiv:2307.02762 (2023).
13. Srivastava, A., Rastogi, A., Rao, A., Shoeb, A.A.M., Abid, A., Fisch, A., Brown, A.R., Santoro, A., Gupta, A., Garriga-Alonso, A., et al.: Beyond the Imitation Game: Quantifying and Extrapolating the Capabilities of Language Models. arXiv preprint arXiv:2206.04615 (2022).
14. Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., Steinhardt, J.: Measuring Massive Multitask Language Understanding. arXiv preprint arXiv:2009.03300 (2020).
15. Zhou, K., Zhu, Y., Chen, Z., Chen, W., Zhao, W.X., Chen, X., Lin, Y., Wen, J.R., Han, J.: Don't Make Your LLM an Evaluation Benchmark Cheater. arXiv preprint arXiv:2311.01964 (2023).
16. Yang, S., Chiang, W.L., Zheng, L., Gonzalez, J.E., Stoica, I.: Rethinking Benchmark and Contamination for Language Models with Rephrased Samples. arXiv:2311.04850 (2023).
17. Li, J., Sun, S., Yuan, W., Fan, R.Z., Zhao, H., Liu, P.: Generative Judge for Evaluating Alignment. arXiv preprint arXiv:2310.05470 (2023).
18. Freitag, M., Rei, R., Mathur, N., Lo, C., Stewart, C., Avramidis, E., Kocmi, T., Foster, G., Lavie, A., Martins, A.F.T.: Results of WMT22 Metrics Shared Task: Stop Using BLEU—Neural Metrics are Better and More Robust. In: Proceedings of the Seventh Conference on Machine Translation (WMT), pp. 46–68 (2022).
19. Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al.: Training Language Models to Follow Instructions with Human Feedback. In: Advances in Neural Information Processing Systems, vol. 35, pp. 27730–27744 (2022).

20. Chen, M., Tworek, J., Jun, H., Yuan, Q., Pinto, H.P.O., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G., et al.: Evaluating Large Language Models Trained on Code. arXiv:2107.03374 (2021).
21. Gao, L., Madaan, A., Zhou, S., Alon, U., Liu, P., Yang, Y., Callan, J., Neubig, G.: PAL: Program-aided Language Models. arXiv preprint arXiv:2211.10435 (2022).
22. Zhou, C., Liu, P., Xu, P., Iyer, S., Sun, J., Mao, Y., Ma, X., Efrat, A., Yu, P., Yu, L., et al.: LIMA: Less Is More for Alignment. arXiv preprint arXiv:2305.11206 (2023).
23. Wang, P., Li, L., Chen, L., Zhu, D., Lin, B., Cao, Y., Liu, Q., Liu, T., Sui, Z.: Large Language Models are Not Fair Evaluators. arXiv abs/2305.17926 (2023).
24. Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E., Le, Q., Zhou, D.: Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. arXiv preprint arXiv:2201.11903v6 (2023).
25. Hendrycks, D., Burns, C., Kadavath, S., Arora, A., Basart, S., Tang, E., Song, D., Steinhardt, J.: Measuring Mathematical Problem Solving with the MATH Dataset. arXiv preprint arXiv:2103.03874 (2021).
26. Zeng, Z., Yu, J., Gao, T., Meng, Y., Goyal, T., Chen, D.: Evaluating Large Language Models at Evaluating Instruction Following. arXiv preprint arXiv:2310.07641 (2023).
27. Kendall, M.G.: Rank Correlation Methods. Griffin (1948).
28. Spearman, C.: The Proof and Measurement of Association Between Two Things. Appleton-Century-Crofts (1961).
29. Chiang, C.H., Lee, H.Y.: Can Large Language Models Be an Alternative to Human Evaluations? arXiv preprint arXiv:2305.01937 (2023).
30. Bubeck, S., Chandrasekaran, V., Eldan, R., Gehrke, J., Horvitz, E., Kamar, E., Lee, P., Lee, Y.T., Li, Y., Lundberg, S., et al.: Sparks of Artificial General Intelligence: Early Experiments with GPT-4. arXiv preprint arXiv:2303.12712 (2023).
31. Gemini Team, Anil, R., Borgeaud, S., Wu, Y., Alayrac, J.B., Yu, J., Soricut, R., Schalkwyk, J., Dai, A.M., Hauth, A., et al.: Gemini: A Family of Highly Capable Multimodal Models. arXiv preprint arXiv:2312.11805 (2023).
32. Bhandari, M., Gour, P., Ashfaq, A., Liu, P., Neubig, G.: Re-evaluating Evaluation in Text Summarization. arXiv preprint arXiv:2010.07100 (2020).
33. Dang, H.T., Owczarzak, K., et al.: Overview of the TAC 2008 Update Summarization Task. In: TAC (2008).
34. Huang, J., Chen, X., Mishra, S., Zheng, H.S., Yu, A.W., Song, X., Zhou, D.: Large Language Models Cannot Self-Correct Reasoning Yet. arXiv preprint arXiv:2310.01798 (2023).
35. Chen, M., Tworek, J., Jun, H., Yuan, Q., Pinto, H.P.O., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G., et al.: Evaluating Large Language Models Trained on Code. arXiv preprint arXiv:2107.03374 (2021).
36. Liang, P., Bommasani, R., Lee, T., Tsipras, D., Soylu, D., Yasunaga, M., Zhang, Y., Narayanan, D., Wu, Y., Kumar, A., et al.: Holistic Evaluation of Language Models. arXiv preprint arXiv:2211.09110 (2022).
37. Wang, X., Wei, J., Schuurmans, D., Le, Q.V., Chi, E.H., Narang, S., Chowdhery, A., Zhou, D.: Self-Consistency Improves Chain of Thought Reasoning in Language Models. In: The Eleventh International Conference on Learning Representations (2023).

A Appendix

A.1 Examined Scenarios

Establishing real-life scenarios that reflect individuals’ daily usage is key to assess the performance and limitations of LLMs in a comprehensive manner. In the current instantiation of SCALEEVAL, we include 8 different scenarios that are closely related to everyday situations and tasks [36,17]. Some example prompts for each defined scenario is shown in Table 6. We describe more about how we define these scenarios below. Individuals interested in evaluating LLMs with our framework can supplement their assessment with additional scenarios.

Brainstorming The brainstorming scenario is designed to test the LLMs’ ability to engage in problem-solving, creative ideation, and generation of insightful responses, especially in situations that require critical thinking and detailed, step-by-step reasoning.

Coding The code scenario evaluates LLMs’ ability to comprehend, produce, and debug code, as well as answering coding-related questions.

Dialog The dialog scenario measures LLMs’ ability to engage with users in a manner that is intuitive, human-like, and dynamic, testing their proficiency through context-sensitive conversations and role-playing that require maintaining a consistent persona throughout a series of interactions.

Scenario	Examples
Brainstorming	- Can you tell me how to make chocolate chip cookies? - Make a list of snacks and foods to serve as party snacks on a game day!
Coding	- What is the difference between HTML and JavaScript? - Implement a binary search algorithm to find a specific element in a sorted array.
Dialog	- Act as the Norse Goddess Freyja. - Can you think and feel like a human?
Judgement	- What if the Aztecs had successfully repelled the Spanish conquistadors? - How can you determine if a person is genuinely interested in a conversation or simply being polite?
Math	- Given that $f(x) = 5x^3 - 2x + 3$, find the value of $f(2)$. - If the endpoints of a line segment are (2, -2) and (10, 4), what is the length of the segment?
ODG	- Is there a meaning for Christmas wreaths? - What are some of the best universities for studying robotics?
ODS	- What causes the northern lights? - What do the different octane values of gasoline mean?
Writing	- Can you help me write a formal email to a potential business partner proposing a joint venture? - Take MLK speech "I had a dream" but turn it into a top 100 rap song.

Table 6: Examined scenarios and corresponding selected examples.

Judgement The judgement scenario assesses LLMs' ability to make inferences and formulate opinions, including soliciting insights on diverse situations or emotions, and posing questions that require logical thinking or reasoning.

Math The math scenario evaluates the LLMs' proficiency in understanding and solving mathematical problems, emphasizing their accuracy in tasks ranging from simple calculations to complex reasoning.

Open-Domain General (ODG) The ODG scenario measures LLMs' proficiency in applying diverse knowledge and exercising reasoning across a wide array of topics, such as answering questions with definitive answers.

Open-Domain Science (ODS) The ODS scenario tests the LLMs' application of scientific knowledge, and gauges their ability to accurately interpret and respond to queries related to scientific disciplines like biology, chemistry, physics, astronomy, and more.

Writing The writing scenario evaluates LLMs' ability to summarize, translate, and generate various texts, testing their core language processing and production skills.

A.2 Prompts

We provide the meta-evaluation and criteria prompt used for SCALEEVAL below.

<Initial Evaluation>

Compare the two submissions based on the criteria above. Which one is better? First, provide a step-by-step explanation of your evaluation reasoning according to the criteria. Avoid any potential bias. Ensure that the order in which the submissions were presented does not affect your judgement. Keep your explanation strictly under 150 words. Afterwards, choose one of the following options:

Submission 1 is better: "1"

Submission 2 is better: "2"

Neither is better: "0"

Directly type in "1" or "2" or "0" (without quotes or punctuation) that corresponds to your reasoning. At the end, repeat just the number again by itself on a new line.

[Question]: {question}

[Submission 1]: {submission_1}

[Submission 2]: {submission_2}

[Criteria]: {criteria}

[User]: {user_prompt}

You are evaluating two submissions for a particular question, using a specific set of criteria. Above is the data.

<Discussion Rounds>

Always remember you are Speaker 1/2/3. Review again your own previous evaluations/discussions first, then answer user's request from Speaker 1/2/3's perspective.

[Question]: {question}

[Submission 1]: {submission_1}

[Submission 2]: {submission_2}

[Criteria]: {criteria}

[Speaker 1's Initial Evaluation]: {evaluation_1}

[Speaker 2's Initial Evaluation]: {evaluation_2}

[Speaker 3's Initial Evaluation]: {evaluation_3}

[Speaker {speaker_number}'s Discussion - Round {round_number}]:

{discussion_reasoning}

...

Read the question, submissions, criteria, and evaluations above. First, explain your thoughts step-by-step about other speakers' evaluations. Second, explain your reasoning step-by-step regarding whether or not to change your original answer about which submission you think is better after considering other speakers' perspectives. Keep your reasoning strictly under 150 words. Afterwards, choose one of the following options:

Submission 1 is better: "1"

Submission 2 is better: "2"

Neither is better: "0"

Directly type in "1" or "2" or "0" (without quotes or punctuation) that corresponds to your reasoning. At the end, repeat just the number again by itself on a new line.

Table 7: Prompt template for meta-evaluation via multi-agent debate

<Type 1: General Format Version>

"1": "Not Helpful - The response is completely unrelated, lacks coherence, and fails to provide any meaningful information."
 "2": "Somewhat Helpful - The response bears some relevance but remains largely superficial and unclear, addressing only the peripheral aspects of the user's needs."
 "3": "Moderately Helpful - The response is mostly relevant and clear, covering the basic aspects of the query, but lacks depth and comprehensive elucidation."
 "4": "Helpful - The response is on-point, detailed, and well-articulated, offering valuable information and clarifications that meet the user's primary needs and enhance understanding."
 "5": "Highly Helpful - The response is exceptionally thorough and precise, providing additional insights and valuable supplementary information."

<Type 2: Shortened Format Version>

"1": "The response is completely unrelated, lacks coherence, and fails to provide any meaningful information."
 "2": "The response bears some relevance but remains largely superficial and unclear, addressing only the peripheral aspects of the user's needs."
 "3": "The response is mostly relevant and clear, covering the basic aspects of the query, but lacks depth and comprehensive elucidation."
 "4": "The response is on-point, detailed, and well-articulated, offering valuable information and clarifications that meet the user's primary needs and enhance understanding."
 "5": "The response is exceptionally thorough and precise, providing additional insights and valuable supplementary information."

<Type 3: Gibberish Format Version>

"1": "N.tc H\$!%ful - Th\$ r\$sp0n\$e is c mplt\$!\$ unr@la7\$d, la\$ks c()h\$ r\$ \$n(€, and f#i/s t# p\$0sid\$ any m@an+gful !format\$on."
 "2": "\$#m\$+ha+ H\$%fu/ - Th\$ r#s0!n\$ b\$ars \$/e r#ev+nc\$ b\$ r\$ma\$ n\$ l#rg\$14 \$u/7\$r7cial an+ !ncl=4r, a6r\$ss@n4 o7ly th\$ p\$R4ph#r\$1 a5p\$ct\$ #f th\$ \$s+r's n+ds."
 "3": "MS!7r\$T#ly H\$!pfu\$ - Th\$ r#s0-n\$0 !s m\$%stl@ r\$'\$van7 an cl\$ar, c\$%r\$ n4 th\$ ba\$!c a\$%ct\$ of th\$ qu@ry, b\$T l#cks d\$pth an cmpr\$h\$ns\$ v\$ lu\$7\$dat'on."
 "4": "H\$!pfu\$ - Th\$ r!s0+n\$e !s o/7-p\$!nt, d\$ta\$!l\$d, an w\$1/-a\$!u/at\$d, #ff\$r!n4 v#l\$%bl\$ #nformat\$on and cl+r\$!cat!ons th\$t m=t th\$ u/7\$r\$ pr!/ary n\$ \$ds an+ @n7anc\$ un#rstand!n4."
 "5": "H4h7y H\$!p\$u\$ - Th\$ r\$ss+n!e !s \$xc\$pt\$#nally th#r#7gh an+ pr\$cs\$%, pr#v\$d\$ n# a4+ts\$#nal !ns\$hts an+ v#lu\$bl\$ @+pp\$%ntary #n\$ormat\$on."

<Type 4: Shuffled Format Version>

"1": "coherence fails provide unrelated, completely response - and the meaningful any to lacks Not Helpful is The information."
 "2": "superficial response largely addressing unclear, remains only needs. - relevance user's and the Helpful the peripheral some bears but aspects Somewhat The of"
 "3": "basic aspects query, lacks Moderately covering clear, - Helpful is depth response and comprehensive elucidation. relevant mostly the The and the of but"
 "4": "clarifications the is response information needs enhance and Helpful - on-point, valuable well-articulated, offering understanding. The and detailed, primary that user's meet"
 "5": "valuable Highly response is providing - the exceptionally Helpful information. insights thorough and additional precise, supplementary and The"

<Type 5: Flipped Format Version>

"1": "toN lufpleH - eHT esnopser si yletelpmoc detalernu, skcal ecnerehoc, dna sliaf ot edivorp yna lufgnaem noitamrofni."
 "2": "tamewoS lufpleH - eHT esnopser sraeb emos ecnaveler tub snlamer ylegral laicifrepus dna raelcnu, gnisserdda ylno eht lahrepirep stcepsa fo eht s' resu sdeen."
 "3": "yletaredoM lufpleH - eHT esnopser si yltsom tnaveler dna raelc, gnirevoc eht cisab stcepsa fo eht yreug, tub skcal htped dna evisneherpmoc noitadicule."
 "4": "lufpleH - eHT esnopser si tniop-no, deliated, dna detalucitra-llew, gnireffo elbaulav noitamrofni dna snoitaicifralc taht teem eht s' resu yramirp sdeen dna ecnahne gnidnatsrednu."
 "5": "ylhgiH lufpleH - eHT esnopser si yllanoitpecke hguorht dna esicerp, gnidivorp lanoitidda sthgisni dna elbaulav yratnemelppus noitamrofni."

<Type 6: Masked Format Version>

"1": "N _H_l_ful - The r_pnse is c_m_et_y unr_l_te_, lacks _ohe_en_e, _nd _ai_s to p_ov_de_ny m_a_ngfu_ _nfo_ma_ion."
 "2": "_om_w_at He_p_ul - T_e re_ponse be_rs_ome rel_a_ce but r_ains la_ely s_erfi_al and u_cle_, ad_res_ng onl_ _he _ri_er_l_a_pe_ts of t_ _u_e's ne_ds."
 "3": "Mod_ _tely_elp_l - Th_ _esp_se is mos_y re_ _va_t an_ _le_r, c_v_ing the ba_ic_spe_ts of the q_e_y, but _cks_e_th and co_preh_ns_ve el_c_d_t_on."
 "4": "_lpful - _he respo_se is on_p_in_, d_ _iled, and we_l-ar_icu_ated, of_er_ng val_ab_e_ _for_ation and cl_r_fi_t_ons_t_at mee_ the _se's_p_im_r_ _eeds and en_nce u_de_tan_ing."
 "5": "Hi_h_y H_p_ul - The r_spon_e is e_c_p_io_al_ th_r_ugh and p_ec_se, pr_v_i_ing a_di_on_l ins_g_ts and va_u_b_e_ _upp_e_en_a_y inf_rma_io_."

Table 8: Criteria prompt format variations for Helpfulness

<p><Type 1: General Format Version> "1": "Beginning - The response is notably lacking in originality, depth, and coherence. It demonstrates a fundamental misunderstanding of the topic, predominantly featuring generic or clichéd thoughts." "2": "Developing - The response reveals faint traces of originality, but ideas are largely underdeveloped or superficial. While there are attempts at creative thinking, they often revert to commonplace concepts. The response may deviate from the main topic." "3": "Competent - The response exhibits a blend of conventional and innovative ideas. It showcases evident creative thinking and a reasonable infusion of original insights. While the response remains largely on-topic, certain areas could be further enriched through deeper exploration." "4": "Proficient - The response includes imaginative and innovative thoughts, reflecting a depth of thinking and divergent exploration. It is content-rich and structured coherently, highlighting a well-considered and effectively executed creative process." "5": "Mastery - The response stands as a beacon of creativity, weaving together profound insights, thoughtful concepts, and astute judgement. Every element of the content radiates originality. The delivery is articulate, compelling, and showcases the pinnacle of creative thought."</p>
<p><Type 2: Shortened Format Version> "1": "The response is notably lacking in originality, depth, and coherence. It demonstrates a fundamental misunderstanding of the topic, predominantly featuring generic or clichéd thoughts." "2": "The response reveals faint traces of originality, but ideas are largely underdeveloped or superficial. While there are attempts at creative thinking, they often revert to commonplace concepts. The response may deviate from the main topic." "3": "The response exhibits a blend of conventional and innovative ideas. It showcases evident creative thinking and a reasonable infusion of original insights. While the response remains largely on-topic, certain areas could be further enriched through deeper exploration." "4": "The response includes imaginative and innovative thoughts, reflecting a depth of thinking and divergent exploration. It is content-rich and structured coherently, highlighting a well-considered and effectively executed creative process." "5": "The response stands as a beacon of creativity, weaving together profound insights, thoughtful concepts, and astute judgement. Every element of the content radiates originality. The delivery is articulate, compelling, and showcases the pinnacle of creative thought."</p>
<p><Type 3: Gibberish Format Version> "1": "B@&#n(n - Th\$ r'spon@e :s n(tally l\$ck@ng !n &r%?na#ity, d?pth, and c()h\$ri\$N(*. #t d@m(nstrat)s a f&ndam@ntal m?s#ndirs&nd?ng of th> t()!c, pr@d()m@ncanly f?at&ring g@n?r!c or cl\$ch@d th(>ghts." "2": "The @e@p@n@s@e r@v@ll@s f!nt t@<@>s of or@r@gnality, but >d@cs are (arg!ly u@d(rde&>loped or (uperf!<!al. W@#le *her(are at (mpts at (%@t\$ve !h@nk@ng, th@y \$f\$en r\$@v@rt to ?#m@m@n@l@ce c@nce>ts. The r?sp@?se ?ay <v@!te >om the m@in top@<." " "3": "Th@ r?>@n@s@e *x@blts a bl@nd of c@nvent!@nal and ?m@v@t@ve @d@as. It sh@w@s\$ \$v@!dent c@eat?ve th@<king and a ?e@st@n@le !nf@s!@n of or!g!n@l #n@s@hts. !h!le the ?e@>n@s@e *e@m@!ns l@rg@ly o!>top@<, %e@t@!n #r@s@s c@u@d be f?>th@: e@rch@d @hr\$ugh @sep@r exp?<rat>n." "4": "Th@ *e@sp@n@s@e &n@!@des !m@?@st@t@ve and ?nnov@t!ve th@#u@hts, %efl!<ct@ng a d@pth of th!<nk!ng and d!ve@gent e?>lor@t!@n. It *s c@nt@nt@r!ch and s@r@ctur@d coh@r@ntly, h?gh!<ght!ng a w@ll@c@ns!d@red and *ff@ct!vely \$@cut@d creat@ve pr@c@ss." "5": "Th@ #?@n@s@e st@nds @s a <e@&n of cr@t!v!ty, w@s@v!ng t@>eth> prof@und >n@s!ghts, th@ughtful c@ncepts, and @st@tute judg(m@nt. ?very *lem@nt of the c@nt@nt r@d@t@s@e or!<gn@lity. The %el@very !s @rt@>ulate, compe@!ng, and s@w@s@s the p@>@cl@e of c@eat!ve th@ught."</p>
<p><Type 4: Shuffled Format Version> "1": "misunderstanding a notably lacking in depth, response and generic The originality, predominantly is Beginning features or of clichéd thoughts. fundamental the topic, demonstrates a" "2": "Developing ideas are faint traces largely originality, The response of but reveals underdeveloped of or superficial. creative at thinking, While attempts are there often they commonplace to revert concepts. main from the deviate response may topic." "3": "Competent - conventional The a blend exhibits response and ideas. of innovative It creative thinking evident showcases a infusion and of original reasonable insights. response While the remains on-topic, largely areas could certain be further through enriched deeper exploration." "4": "Proficient divergent - The includes response and thoughts, imaginative innovative reflecting depth thinking a of structured exploration. and executed effectively content-rich It is coherently, a well-considered highlighting and creative process." "5": "Mastery - stands response The as beacon creativity, a of together weaving profound judgement the pinnacle insights, thoughtful concepts, and astute. Every radiates originality. element of the content The delivery articulate, is compelling, and showcases of creative thought."</p>
<p><Type 5: Flipped Format Version> "1": "gninnigeB - eHT esnopser si ylbaton gnikal ni ytilanigiro, htped dna ecnerehoc. tI setartsnomed a latnemadnuf gnidnatsrednusim fo eht cipot, yltnanimoderp gnirutaef cireneg ro d@h@ic sthguoht." "2": "gnipoleved - eHT esnopser slaever tniaf secart fo ytilanigiro, tub saedi era ylegral repolevedrednu ro laicifrepus. elihW ereht era stpmetta ta evitaerc gnikiht, yeht netfo trever ot ecalpnomoc stpecnoc. eHT esnopser yam etaived morf eht niam cipot." "3": "tnetepmoc - eHT esnopser stibihxe a dnebl fo lanoitnevnoc dna evitavonni saedi. tI sesacwohs tnevde evitaerc gnikiht dna a elbanosaer noisufni fo lanigiro sthgisni. elihW eht esnopser sniamer ylegral no-cipot, niatrec saera dluc eb rehruf dehcirne hguorht repeed noitarolpxe." "4": "tneiciforP - eHT esnopser sedulcni evitanigami dna evitavonni sthguoht, gnitcelfer a htped fo gnikiht dna tnegrevid noitarolpxe. tI si hcir-txetnoc dna derutcurts yltnerohoc, gnithgilhgih a deredisnoc-llew dna ylevitceffe detucexe evitaerc ssecorp." "5": "yretsaM - eHT esnopser sdnats sa a nocaeB fo ytivitaerc, gnivaew rehtegot dnuoforp sthgisni, lufthguoht stpecnoc, dna etutsa tnegeduj. yrevE tneleme fo eht tnetnoc setaidar ytilanigiro. eHT yreviled si etalucitra, gnillepmoc dna sesacwohs eht elcannip fo evitaerc thguoht."</p>
<p><Type 6: Masked Format Version> "1": "N_g_r_n!g_g - Th_e_o_s_s n_t_bly l_ck_ng_n_r_g_n_l_t_y, d_p_h, _nd coh_r_nc_. _t_d_mnst_t_s _f_nd_m_nt_l_m_s_nd_rst_nd_ng_f_th_t_p_c, pr_d_m_nntly f_t_r_ng_g_n_r_c_r cl_ch_d th_ghts." "2": "v_l_p_ng - Th_e_o_s_r_v_l@s f_int_tr_c_s of or_g_n_lity, b_t_d_s _l_rg_ly_nd_r_d_v_l_p_d_r_s_p_r_f_c_l. Wh_l_th_r_r_att_mpts_t_cre_t_v_th_nk_ng, th_y_ft_n_r_v_r_t_t_c_m_m_nplac_c_nc.pts. Th_re_o_s_n_v_d_v_t_from th_m_n_t_p_c." "3": "The Res_ns_xh_b_ts a bl_nd_f_c_nv_rti_n_l_nd inn_v_t_v_d_s. _t_sh_wc_s_v_d_nt cr_t_v_th_nk_ng_nd _r_s_n_b_l_n_f_s_n_f_r_g_n_l_ns_ghts. Wh_l_th_r_s_ns_r_m_ns_l_rg_ly on_t_p_c, c_r_t_n_r_s_c_ld_b_f_rth_f_n_r_ch_d th_r_g_h_d_p_r_expl_r_t_n." "4": "Th_r_sp_ns_inc_d_s !mag_n_r_v_nd_nn_v_t_ve th_ghts, r_f_lct_ng_dpth_f th_nk_ng_nd d_v_rg_nt_xpl_r_t!_n_t is c_nt_nt_r_ch_nd str_ct_r_d_c_h_r_ntly, h_g_hlght_ng_w_ll_c_ns_d_r_d_nd_eff_ct_v_ly_x_c_t_d cr_et_v_pr_c_ss." "5": "h_r_s_nse st_nds_s a_b_c_n_f cre_t_v_t_y, w_v_ng_t_g_th_r_pr_fo_nd ins_ghts, th_ght_f_l_c_nc.pts, _nd _st_t_j_dgm_nt_v_r_y el_m_nt_f th_c_nt_nt_r_d_t_s_r_g_n_l_t_y. Th_d_l_v_r_y_s_art_c_l_t, c_mp_ll_ng, _nd sh_wc_s_th_p_nn_cl_f cr_t_v_th_ght."</p>

Table 11: Criteria prompt format variations for Creativity