

Co-Eval: Augmenting LLM-based Evaluation with Machine Metrics

Anonymous ACL submission

Abstract

Large language models are increasingly used as evaluators in natural language generation tasks, offering scalability and interpretability advantages over traditional evaluation methods. However, current LLM-based evaluations often suffer from biases and misalignment, particularly in domain-specific tasks, due to limited functional understanding and knowledge gaps. To address these challenges, we introduce the Co-Eval framework, which employs a criteria planner model and optimized machine metric to improve scalability, fairness of LLM-based evaluation. Experimental results on both general and domain-specific tasks show that Co-Eval reduces biases across LLMs by up to 0.4903 in self-preference bias and improves alignment with human preferences by up to 0.324 in Spearman correlation.

1 Introduction

Evaluating natural language generation (NLG) quality is challenging, as these tasks often involve subjective judgments, and what constitutes high-quality output can vary depending on the specific context or audience. While human evaluation is a common method for assessing the quality of generated text, it is time-consuming. Recently, researchers (Liu et al., 2023; Chan et al., 2023; Zheng et al., 2023a) have started using large language models (LLMs) as evaluators, noting their impressive performance in aligning with human preferences when assessing generated text.

However, studies (Koo et al., 2023; Panickssery et al., 2024) have shown that LLMs exhibit certain biases, such as a preference for text generated by the models themselves, and factors like presentation order (Wang et al., 2023) and text length (Hu et al., 2024) can affect fairness as well. Moreover, general-purpose LLMs often fall short when it comes to evaluating natural language generation tasks within specific domains (Dorner et al., 2025).

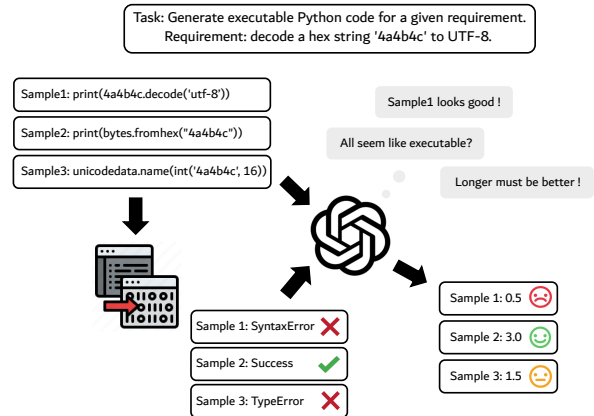


Figure 1: Machine metrics augment scalability and fairness of LLM-based evaluation.

Compared to LLM-based evaluators, machine metrics are more objective, providing precise assessments instead of the semantic evaluations typical of LLMs. Fine-tuned models can incorporate domain-specific knowledge, while rule-based metrics reflect human preferences embedded in rule design. For example, a compiler can definitively indicate if code runs, and BERTScore (Zhang et al., 2019) with CodeBERT can assess code similarity. Metrics like Cyclomatic Complexity (Watson et al., 1996) quantify code complexity by counting decision points. For fairer NLG evaluations and improved domain-specific LLM performance, machine metrics offer reliable benchmarks for consistent, human-aligned measurements.

In this paper, we introduce Co-Eval, a zero-shot reference-free LLM-based evaluation framework that enhances LLM-based evaluation through machine metrics. Recognizing that individual metrics often assess only specific aspects of a task, we fine-tuned a LLaMA-3.1-8B-Instruct model to serve as a criteria planner. This planner interprets diverse task descriptions to establish evaluation criteria, assign weights, and generate score-level descriptions. Next, we developed a comprehensive ma-

chine metrics library to link relevant metrics to the generated criteria based on similarity of their description. The criteria planner is then utilized to refine the machine metric descriptions, ensuring they align closely with the specified criteria. Finally, the prompt-based LLM evaluator is used to generate the final evaluation of each sample, with the overall score calculated as a weighted sum across criteria.

Extensive experiments are conducted across multiple tasks, including four general and four domain-specific tasks, demonstrating that Co-Eval framework enhances LLM-based evaluators, improving agreement with human preferences by up to 0.162 Spearman correlation in general generation tasks and up to 0.324 in domain-specific tasks, while reducing self-preference bias by up to 0.4903.

To summarize, the main contributions of this paper are as follows:

- We introduce Co-Eval, a novel LLM-based evaluation framework that enhances scalability and fairness in evaluation by incorporating machine metrics. We also provide a theoretical proof demonstrating that our framework reduces bias in LLM-based evaluations and improves alignment with human preferences.
- We present a multi-task supervised fine-tuning dataset for the criteria planner, along with a comprehensive machine metric library that includes approximately 50 machine metrics with their implementations.
- We conduct extensive experiments to demonstrate the effectiveness of the Co-Eval framework and, for the first time, explore LLM-based evaluation performance across domain-specific generation tasks.

2 Related Work

2.1 Metric-based Evaluation

Formula-based metrics rely on predefined rules to evaluate the quality of generated responses. Examples include BLEU (Papineni et al., 2002) and METEOR (Banerjee and Lavie, 2005) for machine translation tasks, ROUGE (Lin, 2004) for text summarization, and Flesch-Kincaid score (Flesch, 1943) for readability in educational content.

Model-based metrics leverage pre-trained neural networks to assess the quality of generated responses. For example, BERTScore (Zhang et al., 2019) computes cosine similarity between BERT

embeddings (Devlin, 2018), while GPTScore (Fu et al., 2023) utilizes embeddings from GPT (Radford, 2018). More recently, like UNIEVAL (Zhong et al., 2022), improve embedding-based evaluation by incorporating multiple evaluation dimensions.

Both kinds of machine metrics offer reliable and consistent evaluations but are constrained by their applicability. When used for inappropriate tasks, they can introduce significant biases, leading to misalignment with human preferences.

2.2 LLM-based Evaluation

LLM-based evaluation methods utilize LLMs as sophisticated judges of text quality, often referred to as LLMs-as-judges (Ashktorab et al., 2024; Bavaresco et al., 2024; Tseng et al., 2024).

Prompt-based methods aim to teach LLMs how to evaluate complex tasks through in-context learning. This includes providing fine-grained task criteria (Liu et al., 2023; Zhuo, 2024; Yi et al., 2024; Song et al., 2024a), learning from examples (shot learning) (Fu et al., 2024; Lin and Chen, 2023; Zhang et al., 2024; Jain et al., 2023; Song et al., 2024b), or breaking into multiple iterations (Hasanbeig et al., 2023; Chiang and Lee, 2023; Liu et al., 2024b; Xu et al., 2024; Saha et al., 2024).

Tuning-based methods (Deshwal and Chawla, 2024; Yue et al., 2023; Ye et al., 2024b; Wang et al., 2024; He et al., 2024; Kim et al., 2024; Liu et al., 2024a; Ke et al., 2024), on the other hand, involve training a pre-existing LLM on a specialized dataset to adapt it to specific judgment tasks.

Unlike single-LLM systems, Multi-LLM evaluation (Liang et al., 2024; Zhao et al., 2024a; Moniri et al., 2025; Chan et al., 2023) leverages the collective intelligence of multiple LLMs to enhance evaluation performance.

Despite extensive research, issues such as hallucinations and domain-specific knowledge gaps undermine the robustness of LLM-based evaluation, manifesting as biases, including self-preference bias (Li et al., 2024; Panickssery et al., 2024), position bias (Shi et al., 2024; Zhao et al., 2024b), and verbosity bias (Chen et al., 2024; Zheng et al., 2023b). Avoiding self-evaluation (Ye et al., 2024a) and reference-based approaches (Badshah and Sajjad, 2024) have proven effective in mitigating self-preference bias. However, obtaining accurate models and references can be challenging for open-ended tasks. Additionally, swap-based methods (Raina et al., 2024; Wang et al., 2023) have been shown to effectively address position bias.

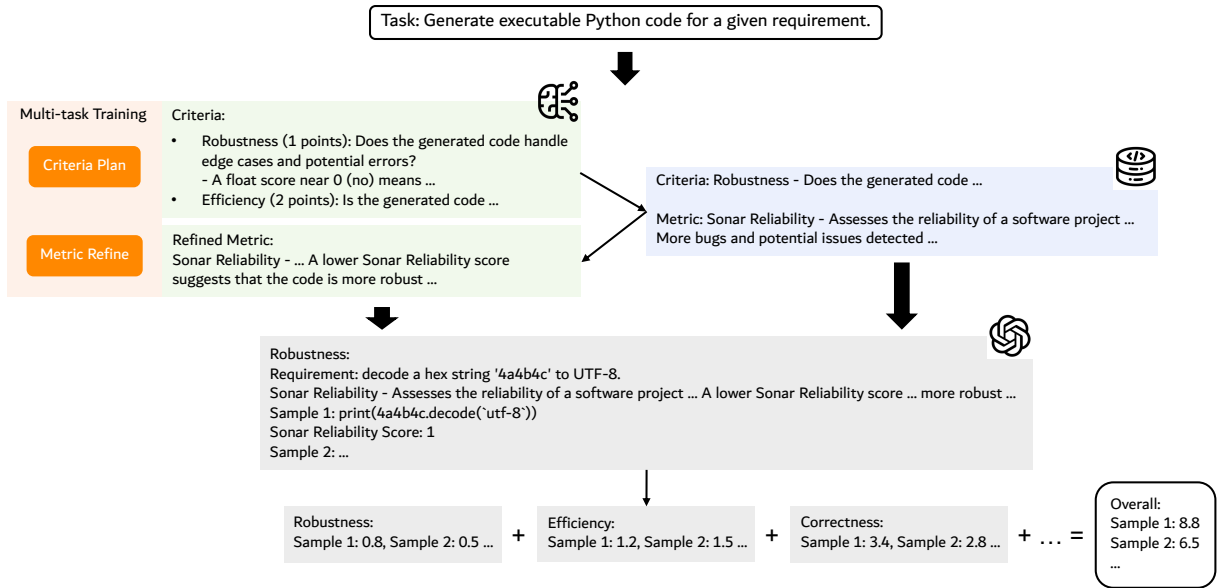


Figure 2: An overview of Co-Eval framework on executable Python code generation task. First, a fine-tuned criteria planner generates scoring criteria and corresponding weights for evaluating the task. Next, each criterion is matched with suitable machine metrics from a machine metric library based on semantic similarity between their descriptions. The chosen machine metrics are then refined by the criteria planner to specify how changes in their scores reflect the performance of the generated code against the criteria. Finally, the task description, original requirement, generated code, machine metric descriptions, and scores are input to a prompt-based evaluator to assign scores to each criterion. These scores are weighted and summed to produce the final evaluation score for each sample.

3 Methodology

To enhance the scalability and fairness of LLM-based evaluators, we propose the Co-Eval framework, outlined in Figure 2.

3.1 Criteria Planner

The main tasks of the criteria planner are to generate evaluation criteria and refine the descriptions of machine metrics.

For the criteria plan task, we recognize that machine metrics are suited for assessing well-defined criteria, which improves accuracy but limits scalability. Furthermore, criteria and their weights must be highly responsive to subtle differences across tasks, as even slight task variations can result in significant shifts in criteria and corresponding weights. Previous research (Kim et al., 2023) has also shown that using fine-grained criteria improves the performance of LLM-based evaluators. Therefore, a criteria planner is needed that can break down task criteria into fine-grained machine metrics and score-level descriptions, adjusting criteria and weights to capture nuanced task differences effectively.

For the metric refine task, we observe that machine metric descriptions tend to be straightforward, focusing mainly on the applicability of each metric rather than linking scores to criteria perfor-

mance. To address this, we refine the machine metric descriptions to better reflect their relationship to the criteria being assessed, rather than using them directly in a prompt-based evaluation setting.

Data Preparation We constructed a multi-task supervised fine-tuning dataset comprising a total of 950 samples. For the criteria planning task, we developed a dataset with 500 task descriptions and corresponding criteria descriptions. Among these, 250 task descriptions were collected from agent platforms such as Coze¹ and GPT-Shop², while the remaining 250 were generated by GPT-4o following a consistent format to ensure diversity and coverage. For the metric refinement task, we used the 500 criteria produced in the criteria planning task. For 250 of these criteria, we searched a metric library to identify suitable metrics and had GPT-4o generate refined metric descriptions. For the remaining 250 criteria, GPT-4o was tasked with both generating suitable metrics and refining their descriptions. To ensure the quality and consistency of the dataset, we extracted the required information from the initial outputs, reorganized them into a standardized format, and filtered out 50 outputs with missing key information. The prompt used for

¹<https://www.coze.com>

²<https://chatgpt.com/gpts>

data preparation is detailed in Appendix D.

Training Strategy Our primary objective is to distill GPT-4o’s performance on criteria planning and metric description refinement tasks, as well as to correct the output format bias of the Llama-3.1-8B-Instruct-based planner, enhancing its suitability for downstream tasks. Given that our training data consists of no more than 1,000 samples and the target task aligns closely with the native capabilities of the Llama-3.1-8B-Instruct model, we employ LoRA (Hu et al., 2021) as our fine-tuning method.

3.2 Machine Metrics Library

We compiled approximately 50 machine metrics for the machine metric library, which can be primarily divided into the following two categories:

Formula-based Metric relies on predefined rules and patterns to assess specific criteria in generated outputs, providing precise evaluations that LLMs may struggle to predict. For example, a syntax parser can accurately verify if generated code is syntactically correct and compilable, an assessment that may exceed the predictive capabilities of LLMs. Another key role of the formula-based metric is to guide the LLM-based evaluator toward aligning more closely with human preferences, which are often embedded within the metric’s design. For instance, when evaluating text summarization, Information Density Formula can prioritize brevity and key information inclusion.

To theoretically validate our approach, we demonstrate the benefits of integrating Formula-based Metrics in the following proof:

Let $f(X)$ be the LLM-based evaluator’s score based on sample X , and let $M(X)$ represent a formula-based metric score derived from X . Define $f(X, M(X))$ as the LLM-based evaluator’s score that incorporates the formula-based metric score $M(X)$. Let $h(X)$ represent the human-assigned score. The error of the LLM-based evaluator relative to the human score is given by

$$\begin{aligned} \epsilon_f &= |h(X) - f(X)| \\ &= |h(X) - E_{s \sim p(s|X)}[s]|, \quad (1) \end{aligned}$$

where s denotes a potential scoring outcome, $p(s|X)$ is the probability distribution over scores s conditioned on the sample X , and $E_{s \sim p(s|X)}[s]$ represents the expected value of s under $p(s|X)$.

Similarly, the error of the LLM-based evaluator when incorporating the formula-based metric is given by

$$\begin{aligned} \epsilon_{f'} &= |h(X) - f(X, M)| \\ &= |h(X) - E_{s \sim p(s|X, M)}[s]|. \quad (2) \end{aligned}$$

According to Bayes’ rule and the principle of maximum entropy, we have

$$p(s|X, M) \propto p(s|X) \cdot \exp(-\lambda(s - \beta M)^2), \quad (3)$$

where λ is a regularization parameter that controls the weight of the metric influence, and β is a scaling factor for the metric M .

For a distribution $p(s|X)$, $Var(p(s|X))$ quantifies how much scores s are expected to vary around their mean when conditioned on X alone. And by the properties of variance, we have

$$\begin{aligned} Var(p(s|X, M)) &= \\ \frac{Var(p(s|X)) \cdot Var(\exp(-\lambda(s - \beta M(X))^2))}{Var(p(s|X)) + Var(\exp(-\lambda(s - \beta M)^2))} &< Var(p(s|X)). \quad (4) \end{aligned}$$

This reduction implies that formula-based metric M can improve LLM-based evaluator to provide a more concentrated estimate around the target score.

Meanwhile, given that M is designed based on human-defined criteria, we assume $Corr(h, M) = \rho$, where $Corr$ represents the correlation between the human-assigned score $h(X)$ and the formula-based metric $M(X)$. We assume $\rho > 0$ implies that $h(X)$ and $M(X)$ are positively correlated, if and only if M is suitable for evaluating X according to the defined criteria. This positive correlation ensures that $\beta > 0$ and that the expected value of s under $p(s|X, M)$ is closer to $h(X)$. Consequently,

$$\begin{aligned} |h(X) - E_{s \sim p(s|X, M)}[s]| &< |h(X) - E_{s \sim p(s|X)}[s]|, \quad (5) \end{aligned}$$

which implies

$$\epsilon_{f'} < \epsilon_f. \quad (6)$$

Model-based Metric leverages well-trained deep neural network models to assess specific criteria for generated outputs. While LLMs are generally effective for broad generation tasks, we focus on smaller, domain-specific models trained on specialized corpora, which are typically more robust in their respective domains compared to general-purpose LLMs. For instance, a BERT model trained on a financial corpus may better capture financial context similarities. This type of model-based metric can augment an LLM-based evaluator’s domain-specific knowledge.

Metrics	Model	Understand		Natural		Coherence		Engaging		Grounded		Overall	
		ρ	τ	ρ	τ	ρ	τ	ρ	τ	ρ	τ	ρ	τ
Formula-based Evaluators													
BLEU-4	-	.033	.025	.130	.100	.277	.219	.386	.316	.446	.396	.280	.223
ROUGE-L	-	.052	.040	.132	.095	.206	.163	.321	.267	.461	.405	.249	.193
Embedding-based Evaluators													
BERTScore	-	.105	.080	.140	.101	.228	.184	.334	.275	.450	.395	.267	.213
BARTScore	-	.061	.039	.158	.124	.232	.188	.300	.237	.489	.422	.272	.215
Learning-based Evaluators													
USR	-	.322	.266	.346	.280	.354	.299	.392	.330	.551	.476	.438	.365
UNIEVAL	-	.467	.360	.513	.373	.612	.465	.608	.458	.574	.451	.662	.486
LLM-based Evaluators													
G-EVAL	GPT-4o	.679	.598	.618	.535	.570	.484	.707	.602	.726	.650	.692	.596
	Llama-3.1-70B	.472	.404	.535	.443	.515	.431	.615	.521	.628	.553	.650	.559
	Qwen-2.5-72B	.571	.486	.618	.531	.590	.505	.744	.663	.696	.621	.689	.592
BATCHEVAL	GPT-4o	.680	.591	.664	.562	.601	.514	.704	.607	.595	.525	.736	.651
	Llama-3.1-70B	.502	.433	.466	.391	.438	.376	.593	.499	.595	.522	.532	.450
	Qwen-2.5-72B	.500	.434	.488	.409	.455	.390	.662	.569	.530	.459	.551	.474
Co-Eval	GPT-4o	.683	.594	.673	.579	.628	.547	.708	.607	.736	.656	.745	.650
	Llama-3.1-70B	.598	.508	.530	.437	.602	.512	.617	.522	.733	.646	.694	.593
	Qwen-2.5-72B	.594	.510	.622	.523	.616	.532	.660	.572	.722	.642	.698	.609

Table 1: Turn-level Spearman (ρ) and Kendall (τ) correlations on Topical-Chat benchmark. The bold scores represent the highest score generated by each LLM as the final prompt-based evaluator, while the grey scores indicate the highest score across the entire column.

We also provide a theoretical justification for the benefits of integrating Model-based Metrics:

Let $D(X)$ represent a model-based metric score derived from X . Assuming that the domain-specific corpus aligns well with human preferences, we have

$$KL(p_d||p_h) \leq \epsilon_1, \quad (7)$$

where $p_d(x)$ denotes the distribution of the domain-specific corpus, $p_h(x)$ denotes the distribution implied by human preferences, and KL is Kullback-Leibler divergence. Since D is trained on the domain-specific corpus, it is optimized to minimize $\min_D E_{x \sim p_d}[L(D(x), h(x))]$. After sufficient training, we assume

$$KL(p_D||p_d) \leq \epsilon_2, \quad (8)$$

where p_D is the distribution implied by D 's scores.

By applying the triangle inequality for KL divergence, we obtain

$$KL(p_D||p_h) \leq KL(p_D||p_d) + KL(p_d||p_h) \leq \epsilon_2 + \epsilon_1 = \epsilon, \quad (9)$$

implying $Corr(h, D) = \rho > 0$. Therefore, the error of the LLM-based evaluator when incorporating the model-based metric is given by

$$\begin{aligned} \epsilon_{f''} &= |h(X) - f(X, D)| && 332 \\ &= |h(X) - E_{s \sim p(s|X, D)}[s]| && 333 \\ &< |h(X) - E_{s \sim p(s|X)}[s]| = \epsilon_f. && 334 \end{aligned} \quad (10)$$

Since typical descriptions of machine metrics sometimes fail to accurately reflect evaluation criteria, we aim to improve their precision by identifying the specific data features that influence changes in metric scores. To achieve this, we provide GPT-4o with pairwise evaluation samples for each metric, enabling it to generate more precise descriptions that highlight the specific features each machine metric effectively captures within its context.

3.3 Prompt-based Evaluator

For the final LLM-based evaluator, we simply adopt the in-context learning and batchwise methods used in BATCHEVAL (Yuan et al., 2023), along with its input and output format. The prompt template is provided in the Appendix D.

4 Experiment

4.1 Experimental Settings

The criteria planner model, based on the Llama-3.1-8B-Instruct model, was fine-tuned by LoRA (Hu

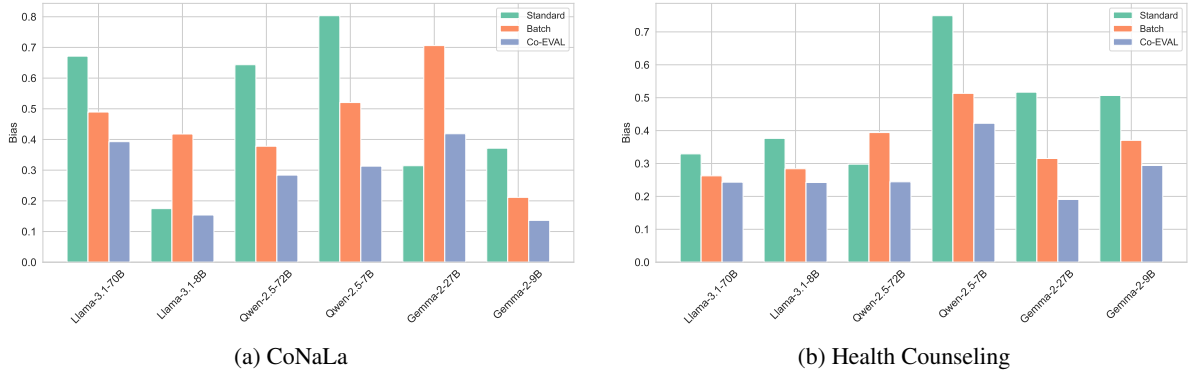


Figure 3: Self-preference bias on CoNaLa and Mental Health Counseling Conversations benchmarks.

et al., 2021) for 3 epochs with a learning rate of $1.0e-4$, a cosine scheduler, and a warmup ratio of 0.1. We set a total score of 10 with a maximum of 5 evaluation criteria. Experimental results for the constrain are provided in Appendix E.3.

In the machine metric search, we select the top three metrics with embedding similarity scores exceeding 0.8, averaging scores across five evaluation runs. Detailed descriptions of LLMs used as prompt-based evaluators and baselines are provided in Appendix A and Appendix B, respectively.

Experiments show that our Co-Eval framework enhances the scalability and fairness of LLM-based evaluation, especially in domain-specific tasks. Detailed experimental implementation information for each benchmark is provided in Appendix C.

4.2 Agreement on Human Preference

For the Topical-Chat benchmark, as shown in Table 1, our proposed Co-Eval framework demonstrates remarkable improvements in Spearman and Kendall correlations across all three models and five original criteria. Even for GPT-4o, the use of suitable machine metrics improve groundness assessment by up to 0.141 compared to BATCHEVAL, while the Co-Eval framework consistently surpasses baselines in overall quality evaluation. Similarly, on the Summeval and HANNA benchmarks, as shown in Table 3 and Figure 7, the Co-Eval framework, with its fine-tuned criteria planner and well-constructed machine metric library, achieves top correlations.

As shown in Table 2, Co-Eval outperforms standard and batch evaluation methods on both the CoNaLa and MATH benchmarks, achieving the highest correlations and even surpassing domain-specific evaluators and fine-tuned

Method	Model	CoNaLa		Model	MATH	
		ρ	τ		ρ	τ
Standard	Prometheus-7B	.065	.063	Prometheus-7B	.113	.108
	Prometheus-8x7B	.256	.253	Prometheus-8x7B	.213	.211
	Llama-3.1-8B	.189	.194	Qwen-2.5-7B	.454	.415
	Llama-3.1-70B	.223	.205	Qwen-2.5-72B	.501	.470
Batch	Llama-3.1-8B	.322	.318	Qwen-2.5-7B	.397	.357
	CodeLlama-7B	.096	.109	Qwen-2.5-MATH-7B	.326	.302
	Llama-3.1-70B	.453	.419	Qwen-2.5-72B	.488	.466
	CodeLlama-70B	.259	.214	Qwen-2.5-MATH-72B	.391	.376
Co-Eval	Llama-3.1-8B	.446	.420	Qwen-2.5-7B	.457	.423
	Llama-3.1-70B	.547	.492	Qwen-2.5-72B	.561	.535

Table 2: Spearman (ρ) and Kendall (τ) correlations on CoNaLa and MATH benchmarks.

evaluation-enhanced models. Notably, on the CoNaLa benchmark, the LLaMA-3.1-70B-Instruct model under Co-Eval improves by up to 0.324 over standard methods.

These results suggest that, whether for general or domain-specific generation tasks, the Co-Eval framework effectively aligns LLM-based evaluators with human preferences. This alignment is particularly beneficial in domain-specific tasks, where functional correctness is critical and general LLMs often struggle to assess accuracy reliably. In these cases, the Co-Eval framework can maximize evaluation effectiveness. In other words, Co-Eval framework can significantly **improve the scalability of LLM-based evaluation**.

4.3 Effectiveness on Bias Elimination

We demonstrate the effectiveness of the Co-Eval framework in eliminating three types of bias: self-preference bias, position bias, and verbosity bias.

Self-preference Bias We calculate the self-preference bias using the following equation:

$$Bias(i) = \frac{1}{N} \sum_{i=1}^N \max(0, R_o(i) - R_s(i)), \quad (11)$$

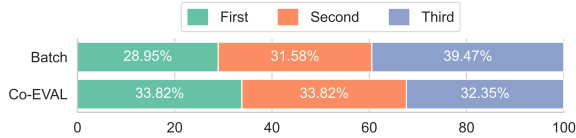


Figure 4: Top-ranking rate on MATH benchmark based on batch position.

where $R_s(i)$ is the rank assigned by the LLM-based evaluator to its self-generated result for instance i , $R_o(i)$ is the average rank assigned by other evaluators, N is the total number of instances.

In the CoNaLa and Health Counseling benchmarks, as illustrated in Figure 3, the Co-Eval framework effectively reduces self-preference bias across all six LLM evaluators. Additionally, smaller LLMs exhibit greater shifts when aided by machine metric scores. The Qwen-2.5-72B-Instruct model achieves the most significant bias reduction compared to individual evaluation. Another notable observation is that certain models, such as Gemma-2-27B-Instruct and Qwen-2.5-72B-Instruct, show increased self-preference bias in batch evaluations. This suggests that while batch evaluation is an effective and straightforward method, it can sometimes amplify self-preference bias when an appropriate baseline is lacking.

Position Bias As shown in Figure 4, we observe that placing the same generated answer in the last position within a batch increases its likelihood of achieving the top rank. However, with the Co-Eval framework, the LLM-based evaluator achieves a more balanced ranking rate, allowing the same answer to attain the top rank consistently, regardless of its position within the batch.

Verbosity Bias As shown in Figure 5, we observe that compared to standard individual methods, LLM-based evaluators using the batch method exhibit a pronounced preference for more verbose answers, even when these answers contain some functional errors. The Co-Eval framework, however, enhances the evaluator’s ability to detect functional errors in generated responses, enabling the LLM-based evaluator to achieve a more balanced ranking across answers of varying verbosity.

Based on the results above, the Co-Eval framework demonstrates outstanding effectiveness in mitigating self-preference bias, position bias, and verbosity bias. In summary, Co-Eval framework can significantly **improves the fairness of LLM-based evaluation**.

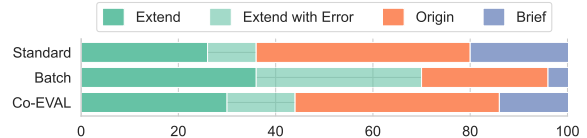


Figure 5: Top-ranking rate on FIQA benchmark based on verbosity degree.

Model	Llama-3.1-70B		Qwen-2.5-72B	
	ρ	τ	ρ	τ
Batch	0.510	0.422	0.532	0.448
Pure	0.465	0.384	0.502	0.413
+ Fine-tuned Planner	0.515	0.431	0.537	0.449
+ GPT-4o as Planner	0.517	0.437	0.543	0.452
+ Metric Library	0.477	0.401	0.521	0.428
+ Planner and Library	0.525	0.448	0.552	0.465

Table 3: Spearman (ρ) and Kendall (τ) correlations on Summeval benchmark.

4.4 Ablation Study

We conducted ablation studies on the two main components: the fine-tuned criteria planner and the machine metrics library.

Fine-tuning for Criteria Planner As shown in Table 3 between the LLM-based evaluator and human preferences, increasing the Spearman correlation by up to 0.05. Moreover, the fine-tuned planner retains over 90% of the performance of the GPT-4o planner. We further analyzed cases before and after fine-tuning the criteria planner. As shown in Figure 6, compared to the original LLaMA-3.1-8B-Instruct model, the fine-tuned planner provides more detailed criteria descriptions and assigns weights more aligned with human preferences. Simple errors, such as incorrect total scores, are also corrected. Additionally, the fine-tuned planner better captures subtle feature differences between tasks. For instance, it identifies "Structure" as essential criteria for "structured outline" task, but not for "summarization" task.

Machine Metric Library As shown in Table 3, machine metrics effectively improve the alignment of LLM-based evaluations with human preferences. We further examined the influence of different machine metrics on the Kendall correlation of LLM-based evaluations. As shown in Figure 7, randomly selected metrics sometimes negatively impact agreement, while the top 1 most similar metric generally has a positive effect. Moreover, machine metric scores from the fine-tuned model have an even greater positive impact compared to

<p>Generate a summarization for a given content.</p>	<ol style="list-style-type: none"> Accuracy (3 points): Does the summarization capture the essential information and key findings? Completeness (2 points): Are all the relevant topics and material included in the summarization? Clarity (2 points): Is the language simple and easy to understand? Organization (2 points): Are the main ideas and supporting details presented in a logical order? Brevity (1 points): Is the summary free from redundancy and unnecessary information? 	<ol style="list-style-type: none"> Accuracy (4 points): Does the summary capture the main ideas and key details of the original content? Clarity (2 points): Is the summary written in a clear, coherent, and easily understandable manner? Brevity (2 points): Does the summary condense the content while maintaining all essential information? Relevance (2 points): Does the summary focus only on the most important points and avoid unnecessary tangents?
<p>Task</p>	<p>LLaMA-3.1-8B-Instruct</p>	<p>Fine-tuned Criteria Planner</p>
<p>Create a structured outline for a given content.</p>	<ol style="list-style-type: none"> Content Coverage (2 points): Comprehensively cover the given content, addressing all key points and subtopics. Organization and Hierarchy (2 points): Have a clear, logical structure, with suitable headings, subheadings, and bullet points. Relevance and Clarity (2 points): Be easy to follow, with concise and non-ambiguous headings and subheadings. Mechanics and Consistency (2 points): Follow a consistent structure, and use proper grammar and punctuation. Adaptability to Content Requirements (4 points): Be adaptable to the specific requirements of the content. 	<ol style="list-style-type: none"> Structure (4 points): Does the outline follow a clear, logical progression of ideas? Comprehensiveness (3 points): Does the outline cover all the key points and important sections of the content? Clarity (2 points): Is the outline clear and easy to understand, with appropriate labeling of sections and subsections? Brevity (1 points): Is the outline concise, without including unnecessary details or extraneous information?

Figure 6: Case study for the fine-tuned criteria planner.

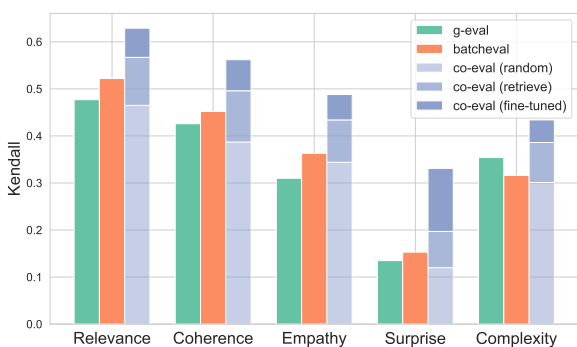


Figure 7: Kendall correlations on HANNA benchmark.

standard individual and batch methods. For criteria where LLM-based evaluations show the weakest performance, the appropriate machine metric and fine-tuned model scores achieve the most significant improvement compared to other criteria.

4.5 Error Analysis

Although we demonstrate the effectiveness of our proposed Co-Eval framework, some remaining errors in the process still need to be addressed:

Criteria planner sometimes fails. While using state-of-the-art models such as GPT-4o as a planner can be costly and inconsistent, fine-tuned smaller LLMs offer a more stable and cost-effective alternative while maintaining comparable performance. However, the generalization ability of fine-tuned smaller LLMs may not be sufficient, especially for long-tail tasks. Although we attempt to improve generalization by collecting data from real agent platforms, it is impossible to cover all real-world scenarios comprehensively. In such cases, using a state-of-the-art model is recommended.

Machine metric library sometimes fails. We rely on the semantic similarity to identify the most suitable machine metric. While we set a high threshold to ensure high precision and strive to make the machine metric descriptions as accurate as possible, semantic similarity does not always yield the best results. In some cases, the identified machine metric may be accurate but not more aligned with human preferences than the LLM itself, particularly for more general criteria. This can potentially misguide the evaluator.

Prompt-based evaluator sometimes fails. To counter occasional misguidance from the machine metric, we allow the final prompt-based evaluator to operate independently, without being strictly bound by these metrics. However, this approach also means that the evaluator may not always follow the instructions of the correct machine metric. Additionally, the limited format-following capability of some LLMs, particularly smaller models, can make parsing the final score more difficult.

A more detailed case study is presented in Appendix F.

5 Conclusion

In this paper, we present Co-Eval, a zero-shot LLM-based evaluation framework that enhances scalability and fairness. The Co-Eval framework integrates machine metrics into the prompt-based evaluator by utilizing a fine-tuned criteria planner and a comprehensive library of metrics. This approach addresses limitations such as bias and misalignment, which arise from inaccurate recognition of functional correctness and gaps in domain-specific knowledge.

541 Limitations

542 Although we demonstrate the effectiveness of our
543 proposed Co-Eval framework, several limitations
544 remain:

- 545 • While we have collected machine metrics
546 for natural language generation tasks across
547 a diverse set of domains, including general,
548 code, mathematical, health, and financial, it
549 remains challenging to cover all potential met-
550 rics. There is considerable room for expand-
551 ing the range of machine metrics to enhance
552 coverage.
- 553 • Our metric retrieval algorithm currently de-
554 pends on semantic similarity between criteria
555 descriptions and metric descriptions. How-
556 ever, this approach lacks adaptability, and mis-
557 matches in metric selection may mislead the
558 LLM-based evaluator.
- 559 • The Co-Eval framework is primarily designed
560 to support LLM-based evaluation, meaning
561 its overall effectiveness largely relies on the
562 capabilities of the LLM, which serves as a
563 prompt-based evaluator. This factor lies be-
564 yond the scope of this paper.

565 References

566 Amod. 2024. [Mental health counseling conversations](#)
567 [\(revision 9015341\)](#).

568 Zahra Ashktorab, Michael Desmond, Qian Pan,
569 James M. Johnson, Martin Santillan Cooper, Eliz-
570 abeth M. Daly, Rahul Nair, Tejaswini Pedapati,
571 Swapnaja Achintalwar, and Werner Geyer. 2024.
572 [Aligning human and llm judgments: Insights](#)
573 [from evalassist on task-specific evaluations and ai-](#)
574 [assisted assessment strategy preferences](#). *Preprint*,
575 arXiv:2410.00873.

576 Sher Badshah and Hassan Sajjad. 2024. [Reference-](#)
577 [guided verdict: Llms-as-judges in automatic evalua-](#)
578 [tion of free-form text](#). *Preprint*, arXiv:2408.09235.

579 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang,
580 Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei
581 Huang, et al. 2023. Qwen technical report. *arXiv*
582 *preprint arXiv:2309.16609*.

583 Satanjeev Banerjee and Alon Lavie. 2005. [METEOR:](#)
584 [An automatic metric for MT evaluation with im-](#)
585 [proved correlation with human judgments](#). In *Pro-*
586 *ceedings of the ACL Workshop on Intrinsic and Ex-*
587 *trinsic Evaluation Measures for Machine Transla-*
588 *tion and/or Summarization*, pages 65–72, Ann Arbor,
589 Michigan. Association for Computational Linguis-
590 tics.

Anna Bavaresco, Raffaella Bernardi, Leonardo Berto-
lazzi, Desmond Elliott, Raquel Fernández, Albert
Gatt, Esam Ghaleb, Mario Giulianelli, Michael
Hanna, Alexander Koller, André F. T. Martins,
Philipp Mondorf, Vera Neplenbroek, Sandro Pezzelle,
Barbara Plank, David Schlangen, Alessandro Sug-
lia, Aditya K Surikuchi, Ece Takmaz, and Alberto
Testoni. 2024. [Llms instead of human judges? a](#)
[large scale empirical study across 20 nlp evaluation](#)
[tasks](#). *Preprint*, arXiv:2406.18403. 591
592
593
594
595
596
597
598
599
600

Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu,
Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan
Liu. 2023. Chateval: Towards better llm-based eval-
uators through multi-agent debate. *arXiv preprint*
arXiv:2308.07201. 601
602
603
604
605

Guiming Hardy Chen, Shunian Chen, Ziche Liu, Feng
Jiang, and Benyou Wang. 2024. [Humans or LLMs](#)
[as the judge? a study on judgement bias](#). In *Proce-*
edings of the 2024 Conference on Empirical Methods
in Natural Language Processing, pages 8301–8327,
Miami, Florida, USA. Association for Computational
Linguistics. 606
607
608
609
610
611
612

Cyril Chhun, Pierre Colombo, Chloé Clavel, and
Fabian M Suchanek. 2022. Of human criteria and
automatic metrics: A benchmark of the evaluation of
story generation. *arXiv preprint arXiv:2208.11646*. 613
614
615
616

Cheng-Han Chiang and Hung-yi Lee. 2023. [A closer](#)
[look into using large language models for automatic](#)
[evaluation](#). In *Findings of the Association for Com-*
putational Linguistics: EMNLP 2023, pages 8928–
8942, Singapore. Association for Computational Lin-
guistics. 617
618
619
620
621
622

Marta R Costa-jussà, James Cross, Onur Çelebi, Maha
Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe
Kalbassi, Janice Lam, Daniel Licht, Jean Maillard,
et al. 2022. No language left behind: Scaling
human-centered machine translation. *arXiv preprint*
arXiv:2207.04672. 623
624
625
626
627
628

Mahesh Deshwal and Apoorva Chawla. 2024. [Phudge:](#)
[Phi-3 as scalable judge](#). *Preprint*, arXiv:2405.08029. 629
630

Jacob Devlin. 2018. Bert: Pre-training of deep bidi-
rectional transformers for language understanding.
arXiv preprint arXiv:1810.04805. 631
632
633

Florian E. Dorner, Vivian Y. Nastl, and Moritz Hardt.
2025. [Limits to scalable evaluation at the frontier:](#)
[Llm as judge won't beat twice the data](#). *Preprint*,
arXiv:2410.13341. 634
635
636
637

Alexander R Fabbri, Wojciech Kryściński, Bryan Mc-
Cann, Caiming Xiong, Richard Socher, and Dragomir
Radev. 2021. Summeval: Re-evaluating summariza-
tion evaluation. *Transactions of the Association for*
Computational Linguistics, 9:391–409. 638
639
640
641
642

Rudolf Flesch. 1943. Marks of readable style; a study
in adult education. *Teachers College Contributions*
to Education. 643
644
645

646	Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. 2023. Gptscore: Evaluate as you desire. <i>arXiv preprint arXiv:2302.04166</i> .	Tran-Johnson, et al. 2022. Language models (mostly) know what they know. <i>arXiv preprint arXiv:2207.05221</i> .	702
647			703
648			704
649	Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. 2024. GPTScore: Evaluate as you desire . In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 6556–6576, Mexico City, Mexico. Association for Computational Linguistics.	Pei Ke, Bosi Wen, Andrew Feng, Xiao Liu, Xuanyu Lei, Jiale Cheng, Shengyuan Wang, Aohan Zeng, Yuxiao Dong, Hongning Wang, Jie Tang, and Minlie Huang. 2024. CritiqueLLM: Towards an informative critique generation model for evaluation of large language model generation . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 13034–13054, Bangkok, Thailand. Association for Computational Linguistics.	705
650			706
651			707
652			708
653			709
654			710
655			711
656			712
657	Karthik Gopalakrishnan, Behnam Hedayatnia, Qianlang Chen, Anna Gottardi, Sanjeev Kwatra, Anu Venkatesh, Raefer Gabriel, and Dilek Hakkani-Tur. 2023. Topical-chat: Towards knowledge-grounded open-domain conversations. <i>arXiv preprint arXiv:2308.11995</i> .	Seungone Kim, Jamin Shin, Yejin Cho, Joel Jang, Shayne Longpre, Hwaran Lee, Sangdoon Yun, Seongjin Shin, Sungdong Kim, James Thorne, et al. 2023. Prometheus: Inducing fine-grained evaluation capability in language models. In <i>The Twelfth International Conference on Learning Representations</i> .	715
658			716
659			717
660			718
661			719
662			720
663	Hosein Hasanbeig, Hiteshi Sharma, Leo Betthausen, Felipe Vieira Frueger, and Ida Momennejad. 2023. Al-lure: Auditing and improving llm-based evaluation of text using iterative in-context-learning . <i>Preprint</i> , arXiv:2309.13701.	Seungone Kim, Juyoung Suk, Shayne Longpre, Bill Yuchen Lin, Jamin Shin, Sean Welleck, Graham Neubig, Moontae Lee, Kyungjae Lee, and Minjoon Seo. 2024. Prometheus 2: An open source language model specialized in evaluating other language models . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 4334–4353, Miami, Florida, USA. Association for Computational Linguistics.	721
664			722
665			723
666			724
667			725
668	Yuanqin He, Yan Kang, Lixin Fan, and Qiang Yang. 2024. Fedeval-llm: Federated evaluation of large language models on downstream tasks with collective wisdom . <i>Preprint</i> , arXiv:2404.12273.		726
669			727
670			728
671			729
672	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset. <i>arXiv preprint arXiv:2103.03874</i> .	Ryan Koo, Minhwa Lee, Vipul Raheja, Jong Inn Park, Zae Myung Kim, and Dongyeop Kang. 2023. Benchmarking cognitive biases in large language models as evaluators. <i>arXiv preprint arXiv:2309.17012</i> .	730
673			731
674			732
675			733
676			734
677	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. <i>arXiv preprint arXiv:2106.09685</i> .	Ruosun Li, Teerth Patel, and Xinya Du. 2024. Prd: Peer rank and discussion improve large language model based evaluations . <i>Preprint</i> , arXiv:2307.02762.	735
678			736
679			737
680			738
681			739
682	Zhengyu Hu, Linxin Song, Jieyu Zhang, Zheyuan Xiao, Jingang Wang, Zhenyu Chen, Jieyu Zhao, and Hui Xiong. 2024. Rethinking llm-based preference evaluation. <i>arXiv preprint arXiv:2407.01085</i> .	Sirui Liang, Baoli Zhang, Jun Zhao, and Kang Liu. 2024. ABSEval: An agent-based framework for script evaluation . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 12418–12434, Miami, Florida, USA. Association for Computational Linguistics.	740
683			741
684			742
685			743
686	Sameer Jain, Vaishakh Keshava, Swarnashree Mysore Sathyendra, Patrick Fernandes, Pengfei Liu, Graham Neubig, and Chunting Zhou. 2023. Multi-dimensional evaluation of text summarization with in-context learning . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 8487–8495, Toronto, Canada. Association for Computational Linguistics.	Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In <i>Text summarization branches out</i> , pages 74–81.	744
687			745
688			746
689			747
690			748
691			749
692			750
693			751
694	Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 2024. Mixtral of experts. <i>arXiv preprint arXiv:2401.04088</i> .	Yen-Ting Lin and Yun-Nung Chen. 2023. LLM-eval: Unified multi-dimensional automatic evaluation for open-domain conversations with large language models . In <i>Proceedings of the 5th Workshop on NLP for Conversational AI (NLP4ConvAI 2023)</i> , pages 47–58, Toronto, Canada. Association for Computational Linguistics.	752
695			753
696			754
697			755
698			756
699	Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli	Minqian Liu, Ying Shen, Zhiyang Xu, Yixin Cao, Eunah Cho, Vaibhav Kumar, Reza Ghanadan, and Lifu Huang. 2024a. X-eval: Generalizable multi-aspect text evaluation via augmented instruction tuning with auxiliary evaluation aspects . In <i>Proceedings of the 2024 Conference of the North American Chapter of</i>	757
700			758
701			

759				
760				
761				
762				
763	Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang,			
764	Ruo Chen Xu, and Chenguang Zhu. 2023. G-eval:			
765	Nlg evaluation using gpt-4 with better human align-			
766	ment. <i>arXiv preprint arXiv:2303.16634</i> .			
767	Yuxuan Liu, Tianchi Yang, Shaohan Huang, Zihan			
768	Zhang, Haizhen Huang, Furu Wei, Weiwei Deng,			
769	Feng Sun, and Qi Zhang. 2024b. Calibrating LLM-			
770	based evaluator . In <i>Proceedings of the 2024 Joint</i>			
771	<i>International Conference on Computational Linguistics,</i>			
772	<i>Language Resources and Evaluation (LREC-</i>			
773	<i>COLING 2024)</i> , pages 2638–2656, Torino, Italia.			
774	ELRA and ICCL.			
775	Shikib Mehri and Maxine Eskenazi. 2020. Usr: An			
776	unsupervised and reference free evaluation metric for			
777	dialog generation. <i>arXiv preprint arXiv:2005.00456</i> .			
778	Behrad Moniri, Hamed Hassani, and Edgar Dobriban.			
779	2025. Evaluating the performance of large language			
780	models via debates . <i>Preprint</i> , arXiv:2406.11044.			
781	Arjun Panickssery, Samuel R Bowman, and Shi Feng.			
782	2024. Llm evaluators recognize and favor their own			
783	generations. <i>arXiv preprint arXiv:2404.13076</i> .			
784	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-			
785	Jing Zhu. 2002. Bleu: a method for automatic evalua-			
786	tion of machine translation. In <i>Proceedings of the</i>			
787	<i>40th annual meeting of the Association for Computa-</i>			
788	<i>tional Linguistics</i> , pages 311–318.			
789	Alec Radford. 2018. Improving language understanding			
790	by generative pre-training.			
791	Vyas Raina, Adian Liusie, and Mark Gales. 2024. Is			
792	LLM-as-a-judge robust? investigating universal ad-			
793	versarial attacks on zero-shot LLM assessment . In			
794	<i>Proceedings of the 2024 Conference on Empirical</i>			
795	<i>Methods in Natural Language Processing</i> , pages			
796	7499–7517, Miami, Florida, USA. Association for			
797	Computational Linguistics.			
798	Swarnadeep Saha, Omer Levy, Asli Celikyilmaz, Mohit			
799	Bansal, Jason Weston, and Xian Li. 2024. Branch-			
800	solve-merge improves large language model evalua-			
801	tion and generation . In <i>Proceedings of the 2024</i>			
802	<i>Conference of the North American Chapter of the</i>			
803	<i>Association for Computational Linguistics: Human</i>			
804	<i>Language Technologies (Volume 1: Long Papers)</i> ,			
805	pages 8352–8370, Mexico City, Mexico. Association			
806	for Computational Linguistics.			
807	Lin Shi, Chiyu Ma, Wenhua Liang, Weicheng Ma, and			
808	Soroush Vosoughi. 2024. Judging the judges: A			
809	systematic study of position bias in llm-as-a-judge .			
810	<i>Preprint</i> , arXiv:2406.07791.			
811	Hwanjun Song, Hang Su, Igor Shalyminov, Jason Cai,			
812	and Saab Mansour. 2024a. FineSurE: Fine-grained			
	summarization evaluation using LLMs . In <i>Proceed-</i>			
	<i>ings of the 62nd Annual Meeting of the Association</i>			
	<i>for Computational Linguistics (Volume 1: Long Pa-</i>			
	<i>pers)</i> , pages 906–922, Bangkok, Thailand. Associa-			
	tion for Computational Linguistics.			
	Mingyang Song, Mao Zheng, and Xuan Luo. 2024b.			
	Can many-shot in-context learning help long-context			
	llm judges? see more, judge better! <i>Preprint</i> ,			
	arXiv:2406.11629.			
	Gemma Team, Thomas Mesnard, Cassidy Hardin,			
	Robert Dadashi, Surya Bhupatiraju, Shreya Pathak,			
	Laurent Sifre, Morgane Rivièrè, Mihir Sanjay Kale,			
	Juliette Love, et al. 2024. Gemma: Open models			
	based on gemini research and technology. <i>arXiv</i>			
	<i>preprint arXiv:2403.08295</i> .			
	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier			
	Martinet, Marie-Anne Lachaux, Timothée Lacroix,			
	Baptiste Rozière, Naman Goyal, Eric Hambro,			
	Faisal Azhar, et al. 2023. Llama: Open and effi-			
	cient foundation language models. <i>arXiv preprint</i>			
	<i>arXiv:2302.13971</i> .			
	Yu-Min Tseng, Wei-Lin Chen, Chung-Chi Chen,			
	and Hsin-Hsi Chen. 2024. Are expert-level lan-			
	guage models expert-level annotators? <i>Preprint</i> ,			
	arXiv:2410.03254.			
	Peiyi Wang, Lei Li, Liang Chen, Zefan Cai, Dawei Zhu,			
	Binghuai Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and			
	Zhifang Sui. 2023. Large language models are not			
	fair evaluators. <i>arXiv preprint arXiv:2305.17926</i> .			
	Tianlu Wang, Ilya Kulikov, Olga Golovneva, Ping Yu,			
	Weizhe Yuan, Jane Dwivedi-Yu, Richard Yuanzhe			
	Pang, Maryam Fazel-Zarandi, Jason Weston, and			
	Xian Li. 2024. Self-taught evaluators . <i>Preprint</i> ,			
	arXiv:2408.02666.			
	Arthur Henry Watson, Dolores R Wallace, and Thomas J			
	McCabe. 1996. <i>Structured testing: A testing method-</i>			
	<i>ology using the cyclomatic complexity metric</i> , vol-			
	ume 500. US Department of Commerce, Technology			
	Administration, National Institute of			
	Shuying Xu, Junjie Hu, and Ming Jiang. 2024. Large			
	language models are active critics in nlg evaluation .			
	<i>Preprint</i> , arXiv:2410.10724.			
	Hongyang Yang, Xiao-Yang Liu, and Christina Dan			
	Wang. 2023. Fingpt: Open-source financial large			
	language models. <i>arXiv preprint arXiv:2306.06031</i> .			
	Jiayi Ye, Yanbo Wang, Yue Huang, Dongping Chen,			
	Qihui Zhang, Nuno Moniz, Tian Gao, Werner Geyer,			
	Chao Huang, Pin-Yu Chen, Nitesh V Chawla, and			
	Xiangliang Zhang. 2024a. Justice or prejudice?			
	quantifying biases in llm-as-a-judge . <i>Preprint</i> ,			
	arXiv:2410.02736.			
	Ziyi Ye, Xiangsheng Li, Qiuchi Li, Qingyao Ai, Yu-			
	jia Zhou, Wei Shen, Dong Yan, and Yiqun Liu.			
	2024b. Beyond scalar reward model: Learning			
	generative judge from preference data . <i>Preprint</i> ,			
	arXiv:2410.03742.			

869	Seungjun Yi, Jaeyoung Lim, and Juyong Yoon. 2024.	Ming Zhong, Yang Liu, Da Yin, Yuning Mao, Yizhu	926
870	Protocollm: Automatic evaluation framework of llms	Jiao, Pengfei Liu, Chenguang Zhu, Heng Ji, and	927
871	on domain-specific scientific protocol formulation	Jiawei Han. 2022. Towards a unified multi-	928
872	tasks . <i>Preprint</i> , arXiv:2410.04601.	dimensional evaluator for text generation. <i>arXiv</i>	929
		<i>preprint arXiv:2210.07197</i> .	930
873	Pengcheng Yin, Bowen Deng, Edgar Chen, Bogdan		
874	Vasilescu, and Graham Neubig. 2018. Learning to	Terry Yue Zhuo. 2024. ICE-score: Instructing large	931
875	mine aligned code and natural language pairs from	language models to evaluate code . In <i>Findings of the</i>	932
876	stack overflow. In <i>Proceedings of the 15th interna-</i>	<i>Association for Computational Linguistics: EACL</i>	933
877	<i>tional conference on mining software repositories</i> ,	2024, pages 2232–2242, St. Julian’s, Malta. Associa-	934
878	pages 476–486.	tion for Computational Linguistics.	935
879	Peiwen Yuan, Shaoxiong Feng, Yiwei Li, Xinglin		
880	Wang, Boyuan Pan, Heda Wang, and Kan Li. 2023.	A Large Language Models	936
881	Batcheval: Towards human-like text evaluation.	GPT Family (Radford, 2018) , developed by Open-	937
882	<i>arXiv preprint arXiv:2401.00437</i> .	AI, is a series of large language models designed	938
883	Weizhe Yuan, Graham Neubig, and Pengfei Liu. 2021.	to understand and generate human-like text. Built	939
884	Bartscore: Evaluating generated text as text gener-	on transformer architecture and pre-trained on ex-	940
885	ation. <i>Advances in Neural Information Processing</i>	tensive datasets, these models primarily excel in	941
886	<i>Systems</i> , 34:27263–27277.	natural language generation tasks.	942
887	Xiang Yue, Boshi Wang, Ziru Chen, Kai Zhang, Yu Su,	Llama Family (Touvron et al., 2023) , developed by	943
888	and Huan Sun. 2023. Automatic evaluation of attri-	Meta, comprises a series of advanced open-source	944
889	bution by large language models . In <i>Findings of the</i>	language models. Included within this family is	945
890	<i>Association for Computational Linguistics: EMNLP</i>	CodeLlama, a domain-specific model focused on	946
891	2023, pages 4615–4635, Singapore. Association for	code generation. CodeLlama is trained on a sub-	947
892	Computational Linguistics.	stantial amount of code data, building on the foun-	948
893	Kaiqi Zhang, Shuai Yuan, and Honghan Zhao. 2024.	dation of the general LLaMA models to enhance	949
894	Talec: Teach your llm to evaluate in specific domain	its capabilities in software development tasks.	950
895	with in-house criteria by criteria division and zero-	Qwen Family (Bai et al., 2023) , developed by Al-	951
896	shot plus few-shot . <i>Preprint</i> , arXiv:2407.10999.	ibaba Cloud, is distinguished by its targeted op-	952
897	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q	timization for conversational AI and information	953
898	Weinberger, and Yoav Artzi. 2019. Bertscore: Eval-	retrieval. Additionally, it offers the Qwen-Math	954
899	uating text generation with bert. <i>arXiv preprint</i>	series, which enhances the mathematical perfor-	955
900	<i>arXiv:1904.09675</i> .	mance of the general Qwen models.	956
901	Ruochen Zhao, Wenxuan Zhang, Yew Ken Chia,	Gemma Family (Team et al., 2024) , developed by	957
902	Weiwen Xu, Deli Zhao, and Lidong Bing. 2024a.	EleutherAI, focuses on lightweight, state-of-the-art	958
903	Auto-arena: Automating llm evaluations with agent	open models, with the largest model containing 27	959
904	peer battles and committee discussions . <i>Preprint</i> ,	billion parameters.	960
905	arXiv:2405.20267.	Mixtral Family (Jiang et al., 2024) , developed by	961
906	Xiutian Zhao, Ke Wang, and Wei Peng. 2024b. Mea-	Mistral AI, comprises a series of advanced open-	962
907	suring the inconsistency of large language models in	source language models, with its notable feature	963
908	preferential ranking . In <i>Proceedings of the 1st Work-</i>	being the implementation of Sparse Mixture of	964
909	<i>shop on Towards Knowledgeable Language Models</i>	Experts (SMoE) architecture.	965
910	<i>(KnowLLM 2024)</i> , pages 171–176, Bangkok, Thai-	B Baselines	966
911	land. Association for Computational Linguistics.	B.1 Formula-based	967
912	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan	BLEU (Papineni et al., 2002) is an automated met-	968
913	Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,	ric for evaluating the quality of machine-translated	969
914	Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang,	text against one or more human reference transla-	970
915	Joseph E Gonzalez, and Ion Stoica. 2023a. Judging	tions. In this study, since we focus on zero-shot	971
916	llm-as-a-judge with mt-bench and chatbot arena . In		
917	<i>Advances in Neural Information Processing Systems</i> ,		
918	volume 36, pages 46595–46623. Curran Associates,		
919	Inc.		
920	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan		
921	Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,		
922	Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023b.		
923	Judging llm-as-a-judge with mt-bench and chatbot		
924	arena . <i>Advances in Neural Information Processing</i>		
925	<i>Systems</i> , 36:46595–46623.		

972	reference-free evaluation performance of each base-	automated evaluation framework designed to as-	1018
973	line method, we calculate the BLEU score between	ess the quality of text generation models in batch	1019
974	the generated response and the source conversation	settings. It leverages LLMs and customizable eval-	1020
975	concatenated with knowledge-based content from	uation criteria, allowing it to assess aspects across	1021
976	the Topical-Chat benchmark.	diverse tasks.	1022
977	ROUGE (Lin, 2004) measures the overlap of n-	Prometheus (Kim et al., 2023, 2024) is a family of	1023
978	grams, word sequences, and word pairs between a	open-source language models designed specifically	1024
979	generated summary and reference summaries. Sim-	for evaluating other language models. Compared to	1025
980	ilar to BLEU, we calculate the ROUGE-L score	the Prometheus 1 models, Prometheus 2 introduces	1026
981	between the generated response and the source con-	support for switch modes by offering different input	1027
982	versation concatenated with knowledge-based con-	prompt formats and system prompts.	1028
983	tent from the Topical-Chat benchmark.		
984	B.2 Embedding-based	C Experimental Implementation	1029
985	BERTScore (Zhang et al., 2019) leverages pre-	C.1 Topical-Chat	1030
986	trained BERT embeddings to capture semantic sim-	Topical-Chat (Gopalakrishnan et al., 2023) is a	1031
987	ilarity between tokens in the generated and refer-	large-scale open-domain conversational benchmark	1032
988	ence texts. For our evaluation, we use the source	containing crowd-sourced conversations on diverse	1033
989	conversation concatenated with knowledge-based	topics, grounded in factual knowledge, and in-	1034
990	content as the reference text for each generated	cludes human evaluation scores for generated re-	1035
991	response in the Topical-Chat benchmark.	sponses across five key criteria: naturalness, coher-	1036
992	BARTScore (Yuan et al., 2021) measures the likeli-	ence, engagingness, groundedness, and understand-	1037
993	hood of a generated text relative to a reference text	ability.	1038
994	using the BART model, treating the evaluation as a	In our work with the Topical-Chat benchmark,	1039
995	text generation task itself. We also use the source	we adhere to the original six evaluation criteria: un-	1040
996	conversation concatenated with knowledge-based	derstanding, naturalness, coherence, engagingness,	1041
997	content as the reference text for each generated	groundedness, and overall quality. Since Topical-	1042
998	response in the Topical-Chat benchmark.	Chat is a multi-turn conversation benchmark, we	1043
999	B.3 Learning-based	follow previous studies (Liu et al., 2023; Yuan et al.,	1044
1000	USR (Mehri and Eskenazi, 2020) is a reference-	2023) and use turn-level correlations, assessing	1045
1001	free metric and leverages pre-trained language mod-	alignment between generated evaluations and hu-	1046
1002	els and unsupervised learning techniques to esti-	man judgments by computing both Spearman (ρ)	1047
1003	mate how well a generated response aligns with	and Kendall (τ) correlations for each turn response,	1048
1004	context and meets conversational quality standards.	then averaging the scores to obtain the final eval-	1049
1005	UNIEVAL (Zhong et al., 2022) is a unified,	uation. For the first five criteria, we adopt the de-	1050
1006	reference-free evaluation framework designed for	scriptions provided by BATCHEVAL (Yuan et al.,	1051
1007	assessing text generation quality. It leverages pre-	2023) and select relevant metrics from the machine	1052
1008	trained language models to assess these qualities,	metric library. To evaluate overall quality, we im-	1053
1009	enabling it to handle a diverse range of text gen-	plement the full Co-Eval pipeline. Additionally, in	1054
1010	eration tasks with a consistent, robust methodology.	our analysis of G-Eval (Liu et al., 2023), we focus	1055
1011	B.4 LLM-based	on the zero-shot evaluation capability of the LLM-	1056
1012	G-EVAL (Liu et al., 2023) is a generative evalu-	-based evaluator, conducting assessments without	1057
1013	ation framework for assessing the quality of gen-	any pre-existing evaluation samples. Results are	1058
1014	erated text. It employs LLMs to directly evaluate	presented in Table 1.	1059
1015	generated text based on criteria across a variety of	C.2 Flores	1060
1016	text generation tasks.	Flores (Costa-jussà et al., 2022) is a benchmark de-	1061
1017	BATCHEVAL (Yuan et al., 2023) is a large-scale,	signed to provide high-quality human translations	1062
		of standardized sentences, enabling the evaluation	1063
		of translation accuracy across low-resource and	1064
		diverse linguistic settings.	1065

For the Flores benchmark, we examine the relationship between LLMs’ familiarity with the target task and their preference bias. Six languages were selected for this study: French, Spanish, Chinese, Vietnamese, Ukrainian, and Thai. We used four LLMs: LLaMA-3.1-8B-Instruct, Qwen-2.5-7B-Instruct, Gemma-2-9B-Instruct, and GPT-4o-mini. Each model translated English text into these six languages. To measure each LLM’s familiarity with the task, we followed previous work (Kadavath et al., 2022) that evaluates familiarity based on the self-consistency of LLMs in translation generation. Specifically, we selected ten samples, generated ten translations per sample with a temperature setting of 0.7, and computed the average token-level BLEU (Papineni et al., 2002) score across these translations. The results were ranked from 1 to 6, indicating each model’s familiarity with the task, from most to least familiar. Results are presented in Figure 8.

C.3 CoNaLa

CoNaLa (Yin et al., 2018) is a large-scale benchmark designed for research in code generation and understanding from natural language. It includes manually curated examples of Python code paired with corresponding natural language intents.

For the CoNaLa benchmark, we used six LLMs, including LLaMA-3.1-8B-Instruct, LLaMA-3.1-70B-Instruct, Qwen-2.5-7B-Instruct, Qwen-2.5-72B-Instruct, Gemma-2-9B-Instruct, and Gemma-2-27B-Instruct, to generate executable Python code based on specific requirements. The six responses were then randomly shuffled, and all six models served as LLM-based evaluators to examine their self-preference biases across three methods: the standard method, the batch method, and the Co-Eval framework. The results are displayed in Figure 3.

To further demonstrate that our proposed framework not only reduces bias but also aligns LLM-based evaluations with human preferences, we sampled the first 50 examples from the benchmark, manually scoring the code generated by LLaMA-3.1-8B-Instruct, Qwen-2.5-7B-Instruct, and Gemma-2-9B-Instruct. We invited three annotators. Each annotator with at least one year of Python coding experience was tasked with evaluating responses for correctness, readability, adherence to coding standards, and alignment with problem requirements. They were also encouraged to run the generated code to verify its func-

tionality. The final human annotation score is calculated as the average of the scores provided by the three annotators. We then calculated the Spearman (ρ) and Kendall (τ) correlations between these models’ scores and human preferences within the standard, batch, and Co-Eval frameworks, using the LLaMA-3.1-8B-Instruct and LLaMA-3.1-70B-Instruct models. Additionally, we applied domain-specific LLMs, CodeLLaMA-7B-Instruct and CodeLLaMA-70B-Instruct, using batch method. Results are shown in the Table 2.

C.4 Mental Health Counseling Conversations

Mental Health Counseling Conversations (Amod, 2024) is a comprehensive collection of conversational data designed to support research and development in the field of mental health counseling. It consists of real-world dialogues between mental health professionals and their clients, focusing on therapeutic interactions aimed at addressing various psychological issues.

For the Health Counseling benchmark, similar to the CoNaLa benchmark, we used six LLMs as well, including LLaMA-3.1-8B-Instruct, LLaMA-3.1-70B-Instruct, Qwen-2.5-7B-Instruct, Qwen-2.5-72B-Instruct, Gemma-2-9B-Instruct, and Gemma-2-27B-Instruct, to generate responses to previous mental health dialogues. The six responses were then randomly shuffled, and all six models served as LLM-based evaluators to examine their self-preference biases across three methods: the standard method, the batch method, and the Co-Eval framework. The results are shown in Figure 3.

C.5 MATH

MATH (Hendrycks et al., 2021) is a large-scale benchmark designed to assess mathematical reasoning abilities, featuring problems that span a wide range of topics from middle school to high school mathematics, including algebra, geometry, calculus, and more. Each problem is accompanied by a detailed step-by-step solution.

For the MATH benchmark, we sampled the first 10 problems from each of the seven categories. Using LLaMA-3.1-8B-Instruct, Qwen-2.5-7B-Instruct, and Gemma-2-9B-Instruct, we generated answers for each question. We then organized the generated answers in three different orders, ensuring that each model’s answer was evaluated in all positions in the batch. We used GPT-4o as an LLM-based evaluator to assess the generated answers across different orderings. We then calcu-

lated the rate at which each answer achieved the highest score at different positions, with results shown in Figure 4.

Similar to CoNaLa benchmark, we also manually scored 70 examples with answers generated by all three models. We invited three annotators as well. Each annotators who had completed at least one mathematics course was instructed to assess responses for accuracy, clarity, logical reasoning, and adherence to problem-solving approaches. The final human annotation score is calculated as the average of the scores provided by the three annotators. We then calculated the Spearman (ρ) and Kendall (τ) correlations between the models' scores and human preferences across the standard, batch, and Co-Eval frameworks, using Qwen-2.5-7B-Instruct and Qwen-2.5-72B-Instruct. Domain-specific LLMs, Qwen-2.5-MATH-7B-Instruct and Qwen-2.5-MATH-72B-Instruct, were also applied using the batch method. Results are presented in Table 2.

C.6 FIQA

FIQA (Yang et al., 2023) is a benchmark designed for research in financial question-answering tasks. It contains a collection of financial questions paired with corresponding answers, covering a wide range of topics such as stock markets, investments, and economic policies.

For the FIQA benchmark, we sampled the first 50 examples and used LLaMA-3.1-8B-Instruct to generate answers for each question. GPT-4o was then used to create both a brief and an extended version of each answer. From the extended versions, we sampled 25 examples and manually introduced errors, such as adding incorrect information, reversing the meaning of some sentences, and making calculation mistakes. We then organized the brief, original, and extended versions, both with and without errors, into a single batch, shuffling the presentation order. Using GPT-4o as an LLM-based evaluator, we calculated the rate at which each version received the highest score across the standard, batch, and Co-Eval frameworks. The results are shown in Figure 5.

C.7 Summeval

Summeval (Fabbri et al., 2021) is a comprehensive benchmark for evaluating abstractive summarization models, featuring human evaluations of machine-generated summaries based on four key

criteria: coherence, consistency, fluency, and relevance.

For the Summeval benchmark, we conducted an ablation experiment for the two main components of the Co-Eval framework: the fine-tuned criteria planner and the machine metric library. We selected the first 6 generated responses from the initial 50 samples and evaluated the Spearman (ρ) and Kendall (τ) correlations of these samples against human preferences, using the LLaMA-3.1-70B-Instruct and Qwen-2.5-72B-Instruct models as LLM-based evaluators. The evaluation employed both the batch method and the Co-Eval framework across four configurations: (1) with a non-fine-tuned criteria planner and no machine metric, (2) with only a fine-tuned criteria planner and no machine metric, (3) with GPT-4o as criteria planner and no machine metric, (4) with a non-fine-tuned criteria planner and machine metric, and (5) with both a fine-tuned criteria planner and machine metric. Results are presented in Table 3, with the complete leaderboard for each criterion shown in Table 4.

C.8 HANNA

HANNA (Chhun et al., 2022) is a large-scale, annotated benchmark designed for evaluating story generation models. It includes human-written and model-generated narratives with detailed annotations for five key aspects: coherence, relevance, empathy, surprise, and engagement.

For the HANNA benchmark, we investigated the impact of machine metric alignment with human preferences on the agreement of LLM-based evaluators with human preferences. We skipped the criteria planning step, using the original criteria descriptions instead. For each criterion, we applied three types of machine metrics: (1) a randomly selected metric from the top 10 retrieved metrics in the machine metric library, (2) the top 1 metric retrieved from the machine metric library, and (3) BERTScore using our fine-tuned BERT model. We then used the LLaMA-3.1-70B-Instruct model as an LLM-based evaluator across the five key aspects using the standard, batch, and Co-Eval frameworks on the first five generated stories of the initial 30 samples. The Spearman (τ) correlations of the evaluation results against human preferences are shown in Figure 7.

In our training setup for the BERT model, we allocated 50 of the remaining 70 samples for training and 20 for validation. The Adam optimizer is used

1267	with a learning rate of 1e-5, and training runs for	<i>Machine Metric: {{machine metric name}} -</i>	1311
1268	a maximum of 30 epochs. We employ a pairwise	<i>{{machine metric description}}</i>	1312
1269	ranking loss on batches generated from the same	Default for Data Preparation	1313
1270	prompt, with early stopping applied if the Kendall	<i>Instruction: First, generate the most suitable</i>	1314
1271	correlation does not improve on the validation set	<i>machine metric for the given criterion with met-</i>	1315
1272	for 5 consecutive epochs.	<i>ric description. Then, provide a detailed metric</i>	1316
		<i>description that clearly explains how the metric re-</i>	1317
1273	D Prompts	<i>fects and aligns with the corresponding criterion.</i>	1318
1274	D.1 Criteria Plan	<i>An Example:</i>	1319
1275	Default for Fine-tuned Criteria Planner	<i>Criteria: Coherence – Measures how logically</i>	1320
1276	<i>Please provide the evaluation criteria for this</i>	<i>the summary flows, ensuring clarity and consis-</i>	1321
1277	<i>task, including the weight of each criterion. The</i>	<i>tency in the ideas presented.</i>	1322
1278	<i>total score should be 10 points.</i>	<i>Machine Metric: BERTScore – Evaluates the</i>	1323
		<i>semantic similarity between two pieces of text.</i>	1324
1279	<i>Task: {{task description}}</i>	<i>Detailed Machine Metric: BERTScore – Evalu-</i>	1325
1280	Default for Data Preparation	<i>ates the semantic similarity between two pieces of</i>	1326
1281	<i>Task: {{task description}}</i>	<i>text. A higher BERTScore reflects a greater degree</i>	1327
		<i>of coherence, indicating that the summary aligns</i>	1328
1282	<i>Instruction: Please provide the evaluation cri-</i>	<i>more closely with the logical flow and meaning of</i>	1329
1283	<i>teria for this task, including the weight of each</i>	<i>the original content.</i>	1330
1284	<i>criterion. The total score should be 10 points, with</i>		
1285	<i>no more than 5 criteria in total. Present the infor-</i>	<i>Criteria: {{criteria name}} - {{criteria descrip-</i>	1331
1286	<i>mation in the following format:</i>	<i>tion}}</i>	1332
1287	<i>No. Criterion Name (Weight in points) - Descrip-</i>	<i>Machine Metric: {{machine metric name}} -</i>	1333
1288	<i>tion of what this criterion evaluates. Provide clear</i>	<i>{{machine metric description}}</i>	1334
1289	<i>guidance on how this aspect of the response will</i>		
1290	<i>be assessed.</i>	D.3 Evaluation	1335
1291	<i>An Example:</i>	Example of Standard Individual Evaluation	1336
1292	<i>1. Efficiency (2 points): Is the generated code</i>	<i>You will be given a sample, containing a gener-</i>	1337
1293	<i>optimized in terms of time and space complexity?</i>	<i>ated code for given requirement.</i>	1338
1294	<i>- A float score near 0 (no) means the code is in-</i>	<i>Your task is to assign a float score to the response</i>	1339
1295	<i>efficient and has significant room for optimization.</i>	<i>on one metric.</i>	1340
1296	<i>- A float score near 1 (somewhat) means the</i>	<i>You should carefully horizontally compare the</i>	1341
1297	<i>code has a moderate level of efficiency but could</i>	<i>given samples in order to assign a suitable float</i>	1342
1298	<i>be improved.</i>	<i>score to each sample.</i>	1343
1299	<i>- A float score near 2 (yes) means the code is</i>	<i>Please make sure you read and understand these</i>	1344
1300	<i>highly optimized in both time and space complexity.</i>	<i>instructions carefully. Please keep this document</i>	1345
		<i>open while reviewing, and refer to it as needed.</i>	1346
1301	<i>Return the complete list. Note: Efficiency is</i>	<i>Evaluation Criteria:</i>	1347
1302	<i>included as an example and is not required to be</i>	<i>Overall (floating point numbers within the inter-</i>	1348
1303	<i>part of the final list.</i>	<i>val [1,5]): What is your overall impression of the</i>	1349
1304	D.2 Machine Metric Refinement	<i>quality of the generated code?</i>	1350
1305	Default for Fine-tuned Criteria Planner	<i>- A float score near 1 (very poor): The generated</i>	1351
1306	<i>Please provide a detailed metric description that</i>	<i>code is of very low quality. It contains significant</i>	1352
1307	<i>clearly explains how the metric reflects and aligns</i>	<i>errors or does not run at all, lacks any meaningful</i>	1353
1308	<i>with the corresponding criterion.</i>	<i>structure, and does not meet the requirements in</i>	1354
		<i>any substantial way. The code might be difficult or</i>	1355
1309	<i>Criteria: {{criteria name}} - {{criteria descrip-</i>	<i>impossible to salvage for further use.</i>	1356
1310	<i>tion}}</i>	<i>- A float score near 2 (poor): The code runs but</i>	1357
		<i>is largely incorrect or ineffective. There are numer-</i>	1358

1359	<i>ous logical errors or missing functionality, and it</i>	<i>code is of very low quality. It contains significant</i>	1408
1360	<i>does not align well with the provided requirements.</i>	<i>errors or does not run at all, lacks any meaningful</i>	1409
1361	<i>The code may also suffer from poor readability</i>	<i>structure, and does not meet the requirements in</i>	1410
1362	<i>or lack of proper structure, making it difficult to</i>	<i>any substantial way. The code might be difficult or</i>	1411
1363	<i>understand or maintain.</i>	<i>impossible to salvage for further use.</i>	1412
1364	- A float score near 3 (neutral): <i>The code is</i>	- A float score near 2 (poor): <i>The code runs but</i>	1413
1365	<i>functional but unremarkable. It may have some er-</i>	<i>is largely incorrect or ineffective. There are numer-</i>	1414
1366	<i>rors or areas for improvement but generally follows</i>	<i>ous logical errors or missing functionality, and it</i>	1415
1367	<i>the basic requirements and runs with acceptable</i>	<i>does not align well with the provided requirements.</i>	1416
1368	<i>results. The code is neither highly readable nor</i>	<i>The code may also suffer from poor readability</i>	1417
1369	<i>efficient, but it's not overly difficult to understand</i>	<i>or lack of proper structure, making it difficult to</i>	1418
1370	<i>or extend.</i>	<i>understand or maintain.</i>	1419
1371	- A float score near 4 (good): <i>The generated code</i>	- A float score near 3 (neutral): <i>The code is</i>	1420
1372	<i>is of good quality, meeting most of the requirements</i>	<i>functional but unremarkable. It may have some er-</i>	1421
1373	<i>with only minor issues. It runs correctly for the</i>	<i>rors or areas for improvement but generally follows</i>	1422
1374	<i>majority of test cases and is fairly easy to read</i>	<i>the basic requirements and runs with acceptable</i>	1423
1375	<i>and maintain. The code could be improved, but</i>	<i>results. The code is neither highly readable nor</i>	1424
1376	<i>any changes would be enhancements rather than</i>	<i>efficient, but it's not overly difficult to understand</i>	1425
1377	<i>necessary fixes.</i>	<i>or extend.</i>	1426
1378	- A float score near 5 (excellent): <i>The code is</i>	- A float score near 4 (good): <i>The generated code</i>	1427
1379	<i>of very high quality, demonstrating strong adher-</i>	<i>is of good quality, meeting most of the requirements</i>	1428
1380	<i>ence to all requirements. It is free from significant</i>	<i>with only minor issues. It runs correctly for the</i>	1429
1381	<i>errors, highly readable, well-structured, efficient,</i>	<i>majority of test cases and is fairly easy to read</i>	1430
1382	<i>and maintainable. The code is clear, concise, and</i>	<i>and maintain. The code could be improved, but</i>	1431
1383	<i>easy to understand, with well-considered logic and</i>	<i>any changes would be enhancements rather than</i>	1432
1384	<i>style. There are no significant flaws or areas for</i>	<i>necessary fixes.</i>	1433
1385	<i>improvement.</i>	- A float score near 5 (excellent): <i>The code is</i>	1434
1386	<i>Generated code and given requirement:</i>	<i>of very high quality, demonstrating strong adher-</i>	1435
1387	<i>Source: {{requirement source}}</i>	<i>ence to all requirements. It is free from significant</i>	1436
1388	<i>System Response: {{response output}}</i>	<i>errors, highly readable, well-structured, efficient,</i>	1437
1389	<i>Evaluation Form (scores ONLY):</i>	<i>and maintainable. The code is clear, concise, and</i>	1438
1390	<i>- Overall:</i>	<i>easy to understand, with well-considered logic and</i>	1439
1391	Example of Batch Evaluation	<i>style. There are no significant flaws or areas for</i>	1440
1392	<i>You will be given a batch of 8 samples. Each</i>	<i>improvement.</i>	1441
1393	<i>sample contains a generated code for given require-</i>	<i>Generated code and given requirement:</i>	1442
1394	<i>ment.</i>	<i>Source: {{requirement source}}</i>	1443
1395	<i>Your task is to assign a float score to the response</i>	<i>Sample 1:</i>	1444
1396	<i>on one metric.</i>	<i>System Response: {{sample 1 response output}}</i>	1445
1397	<i>You should carefully horizontally compare the</i>	<i>Sample 2:</i>	1446
1398	<i>given samples in order to assign a suitable float</i>	<i>System Response: {{sample 2 response output}}</i>	1447
1399	<i>score to each sample.</i>	<i>...</i>	1448
1400	<i>Please make sure you read and understand these</i>	<i>Sample 6:</i>	1449
1401	<i>instructions carefully. Please keep this document</i>	<i>System Response: {{sample 6 response output}}</i>	1450
1402	<i>open while reviewing, and refer to it as needed.</i>	<i>Evaluation Form (Answer by starting with "Anal-</i>	1451
1403	<i>Evaluation Criteria:</i>	<i>ysis:" to analyze the given samples regarding the</i>	1452
1404	<i>Overall (floating point numbers within the inter-</i>	<i>evaluation criteria and offer insights derived from</i>	1453
1405	<i>val [1,5]): What is your overall impression of the</i>	<i>the machine metric scores as concise as possible</i>	1454
1406	<i>quality of the generated code?</i>	<i>(Attention: Don't give your scores during this step).</i>	1455
1407	- A float score near 1 (very poor): <i>The generated</i>	<i>After analyzing all the samples, please give all</i>	1456
		<i>the float scores in order following the template</i>	1457

"Float Scores: [Sample1:score of Sample1, Sample2:score of Sample2, Sample3:score of Sample3, Sample4:score of Sample4, Sample5:score of Sample5, Sample6:score of Sample6]".

Example of Co-Eval Evaluation

You will be given a batch of 8 samples. Each sample contains a generated code for given requirement.

Your task is to assign a float score to the response on one metric.

You should carefully horizontally compare the given samples in order to assign a suitable float score to each sample.

You can refer to the machine metric scores of each sample if you are not confidence.

Please make sure you read and understand these instructions carefully. Please keep this document open while reviewing, and refer to it as needed.

Evaluation Criteria:

Robustness (floating point numbers within the interval [0,2]): Does the generated code handle edge cases and potential errors gracefully?

- A float score near 0 (no) means the code fails to handle edge cases or crashes on invalid inputs.

- A float score near 1 (somewhat) means the code handles some edge cases but misses others or lacks comprehensive error handling.

- A float score near 2 (yes) means the code effectively handles all edge cases and includes comprehensive error handling.

Given Content and potentially useful Machine Metric Score:

Source: {{requirement source}}

Sonar Reliability - Assesses the robustness and fault-tolerance of software code, focusing on its potential to contain bugs or defects that could lead to malfunctions in production. The lower the numerical score, the better the reliability of the code, indicating fewer bugs and a lower risk of defects impacting the software's functionality.

Sample 1:

System Response: {{sample 1 response output}}

Score: {{sample 1 sonar reliability score}}

Sample 2:

System Response: {{sample 2 response output}}

Score: {{sample 2 sonar reliability score}}

...

Sample 6:

System Response: {{sample 6 response output}}

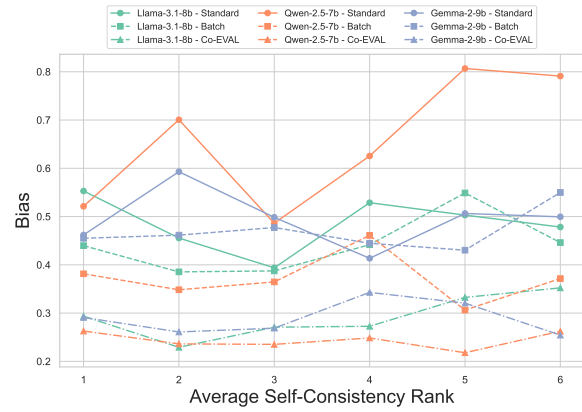


Figure 8: Self-preference bias on Flores benchmark.

Score: {{sample 6 sonar reliability score}}

Evaluation Form (Answer by starting with "Analysis:" to analyze the given samples regarding the evaluation criteria and offer insights derived from the machine metric scores as concise as possible (Attention: Don't give your scores during this step). After analyzing all the samples, please give all the float scores in order following the template "Float Scores: [Sample1:score of Sample1, Sample2:score of Sample2, Sample3:score of Sample3, Sample4:score of Sample4, Sample5:score of Sample5, Sample6:score of Sample6]".

- Robustness:

E Additional Experiment Results

E.1 Self-preference on Flores Benchmark

For the Flores benchmark, we attempt to explore the relationship between self-preference bias and LLMs' familiarity with different languages. Unfortunately, as shown in Figure 8, our results indicate that self-preference bias does not exhibit a clear correlation with language familiarity. This may be due to variations in language familiarity affecting the accuracy of self-preference bias calculations based on average rank. Nevertheless, regardless of the direction of these variations, batch evaluations help reduce self-preference across models and languages, with the Co-Eval framework further minimizing bias to near-uniform levels across languages.

E.2 Complete Summeval Leaderboard

We present complete experimental results on the Summeval benchmark, a meta-benchmark with fine-grained labels. The results are summarized

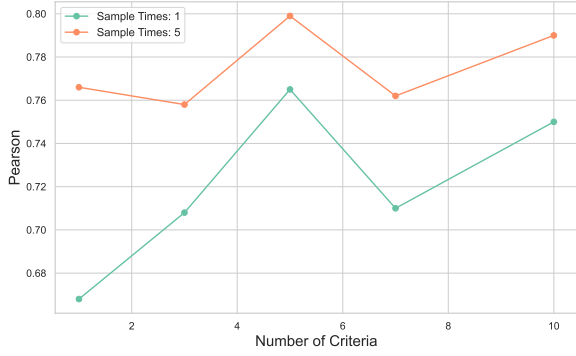


Figure 9: Pearson correlations on Topical-chat benchmark.

in Table 4.

The results on the Summeval benchmark with fine-grained labels exhibit a trend similar to that of the Topical-Chat benchmark. While G-EVAL and BATCHEVAL outperform in certain criteria, our proposed Co-Eval framework consistently achieves the best performance on the "Overall" criteria.

E.3 Criteria Number and Sample Times

We evaluate the impact of the number of criteria and sample times on the Pearson correlations using the Topical-Chat benchmark in relation to our proposed Co-Eval framework. Specifically, we assess the performance of the Co-Eval framework with criteria numbers of 1, 3, 5, 7, and 10, and sample times of 1 and 5. As shown in Figure 9, the performance is more consistently aligned with human preferences when we sample 5 times and take the average score, compared to sampling only once, which is consistent with the findings reported in prior work (Yuan et al., 2023).

Regarding the number of criteria, the Pearson correlation shows an increasing trend from 1 (equivalent to the batch method) to 5. However, when the number of criteria exceeds 5, the Pearson correlation begins to decrease, indicating that 5 criteria is the most suitable choice for common generation tasks. Too few criteria fail to provide a comprehensive evaluation of the task, while too many criteria can lead to diminishing returns, potentially introducing redundant or conflicting evaluation metrics that compromise the accuracy and coherence of the overall assessment.

E.4 Impact of Temperature

We evaluate the impact of temperature on self-preference bias, position bias, and verbosity bias by testing temperatures of 0.0, 0.3, 0.5, 0.7, and

1.0, and reproducing the experiments for each type of bias.

As shown in Figure 10, while the effect of temperature on self-preference bias varies across models, our proposed Co-Eval framework consistently enables the LLM-based evaluator to achieve the lowest self-preference bias. Furthermore, for position bias and verbosity bias, GPT-4o, when used as an LLM-based evaluator with the Co-Eval framework, consistently maintains a balanced top-ranking rate while being less influenced by the position and verbosity of each response.

F Detailed Case Study

We further analyze the cases throughout the entire process:

Case 1: For some long-tail tasks, the generalization ability of the fine-tuned criteria planner is insufficient to generate a comprehensive set of evaluation criteria. For example, consider the task: Generate architectural drawings for a supermarket. The fine-tuned criteria planner accounts for the following aspects: Accuracy of Store Layout, Adherence to Building Codes and Regulations, Effective Use of Space, Aesthetic Appeal and Brand Identity, and Technical Quality and Presentation. However, all five criteria are equally weighted, each contributing 2 points to the total 10-point score. In contrast, human preferences suggest that Regulations and Store Layout should carry the most weight, making the evaluation misaligned with human judgment. Additionally, compared to the GPT-4o, budget considerations and branding alignment, both critical factors in supermarket architectural design, are missing from the criteria set. This gap further highlights the planner’s limitations in capturing human-centric evaluation priorities.

Case 2: For some criteria descriptions, the machine metric with the highest semantic similarity score does not necessarily align best with human preferences. For example, in the Fluency criterion of the SummEval benchmark, perplexity is the machine metric whose description is most semantically similar to the criterion description. However, BARTScore exhibits a significantly higher Spearman correlation with human judgment. This misalignment leads to lower performance when Llama-3.1-70B-Instruct serves as the final prompt-based evaluator within the Co-Eval framework. The mistake arises despite regenerating machine metric descriptions via sampling to better reflect the specific

Metrics	Model	Coherence		Consistency		Fluency		Relevance		Overall	
		ρ	τ	ρ	τ	ρ	τ	ρ	τ	ρ	τ
G-EVAL	Llama-3.1-70B	.542	.454	.550	.486	.423	.366	.395	.338	.517	.423
	Qwen-2.5-72B	.509	.425	.624	.563	.529	.469	.413	.349	.474	.399
BATCHEVAL	Llama-3.1-70B	.444	.366	.547	.483	.427	.372	.421	.354	.510	.422
	Qwen-2.5-72B	.514	.424	.552	.497	.430	.373	.407	.343	.532	.448
Co-Eval	Llama-3.1-70B	.548	.502	.452	.413	.391	.355	.464	.427	.525	.448
	Qwen-2.5-72B	.483	.415	.592	.544	.558	.511	.457	.391	.552	.465

Table 4: Complete Spearman (ρ) and Kendall (τ) correlations on Summeval benchmark.

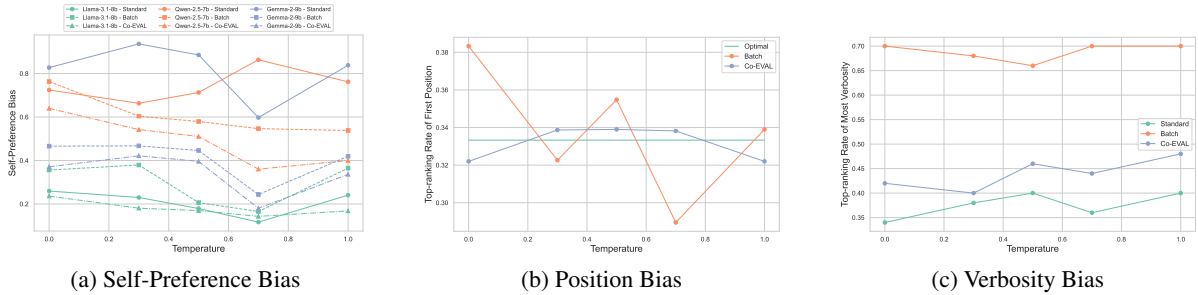


Figure 10: Impact of temperature on three kinds of bias.

1626 aspects each metric evaluates. However, human
 1627 evaluation does not always have clearly defined
 1628 boundaries between different criteria—especially
 1629 for closely related aspects. As a result, scores for
 1630 Coherence can inadvertently influence the evaluation
 1631 of Fluency, leading to discrepancies in alignment.
 1632

1633 **Case 3:** For some general tasks, the machine
 1634 metric score is less aligned with human preferences
 1635 than the LLM itself. For example, as shown in the
 1636 results in Table 1 and 4, LLM-based evaluation
 1637 achieves the highest scores in some criteria using
 1638 the batch method, even when the standard method
 1639 is used without a machine metric. This is true
 1640 even when the reference machine metric is suitable,
 1641 particularly for criteria that are more subjective
 1642 and dependent on the evaluator. In such cases,
 1643 the machine metric may interfere with the prompt-
 1644 based evaluator to some extent.

1645 **Case 4:** The prompt-based evaluator demon-
 1646 strates critical thinking when assessing the refer-
 1647 ence machine metric score. For example, "Upon re-
 1648 viewing the samples, it is evident that the machine
 1649 metric scores do not directly reflect the readabil-
 1650 ity of the code... However, analyzing the samples
 1651 based on readability, we find that..." This capabil-
 1652 ity strengthens the robustness of our proposed Co-
 1653 Eval framework against unsuitable machine metric

1654 scores. However, it also introduces the possibility
 1655 that the prompt-based evaluator may resist follow-
 1656 ing the instructions of the augmented machine met-
 1657 ric. As shown in the experiment on verbosity bias,
 1658 an 8% extended response containing error informa-
 1659 tion still achieved the highest score, even though
 1660 the machine metric detected the error.

1661 **Case 5:** Some LLMs, particularly smaller mod-
 1662 els, exhibit weak format-following capabilities. For
 1663 example, when LLaMA-3.1-8B-Instruct is used as
 1664 the final prompt-based evaluator, it may present
 1665 scores in inconsistent formats such as: "Float
 1666 Scores: Sample1: [3], Sample2: [2], Sample3:
 1667 [3], Sample4: [4]" and "Float Scores: [4.5: Sam-
 1668 ple1, 2: Sample2, 4: Sample3, 4.5: Sample4]",
 1669 whereas the expected standard format is: "Float
 1670 Scores: [Sample1: 2.5, Sample2: 2.5, Sample3: 4,
 1671 Sample4: 4]". These inconsistencies complicate
 1672 score parsing and may lead to misinterpretations of
 1673 evaluation results.

1674 **Case 6:** Compared to the diversity of tasks, the
 1675 coverage of machine metrics is limited. As a result,
 1676 some criteria lack suitable machine metrics, such
 1677 as the "Completeness" criteria in the MATH bench-
 1678 mark. Determining whether a solution step is both
 1679 complete and reasonable remains an open question.
 1680 In our experiment, we design a metric to evaluate
 1681 completeness using the BERTScore between con-

1682 secutive steps in a solution. A higher average score
1683 across all solution steps indicates a more complete
1684 and detailed response. Additionally, the Co-Eval
1685 framework makes it easy to incorporate new and
1686 useful machine metrics into the evaluation process,
1687 improving adaptability and coverage.