
Efficient Evaluation of LLMs via Branching Preference Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Large language models (LLMs) have made significant advances across various
2 generative tasks, progressing toward achieving near-human levels of intelligence.
3 However, in many scenarios, LLMs face the challenge of insufficient human
4 evaluation or even the inability to evaluate reliably. Particularly, in complex
5 dialogue scenarios involving diverse and intricate user intents, LLMs as evaluators
6 of AI responses exhibit a substantial gap compared to humans. Moreover, due
7 to the scarcity of high-quality evaluation data, LLMs exhibit deficiencies in their
8 evaluation capabilities. In this work, we conceptualize the evaluation process as
9 a decision tree, where each node represents an evaluation action, and each path
10 from the root to a leaf node represents a trajectory of evaluation reasoning. We
11 demonstrate that within a limited search space, there exist better decision-making
12 behaviors that facilitate the model in making reasonable and accurate judgments.
13 Specifically, we propose a tree-based data sampling method to generate supervised
14 data and preference pairs derived from the evaluation tree. Furthermore, we
15 introduce preference learning based on the DPO algorithm, which empowers the
16 fine-grained evaluation model to explore and learn better branching strategies within
17 budget-limited scenarios. Our model significantly reduces the dependency on
18 labeled data and demonstrates strong performance across three different evaluation
19 settings: in-distribution, out-of-distribution, and transfer evaluation. Experiments
20 indicate that our model can reduce inference costs by 90% compared to conducting
21 searches across the entire evaluation tree, thereby significantly enhancing efficiency.

22 1 Introduction

23 Dialogue evaluation capability [6] is one of the fundamental abilities of human social interaction,
24 involving the comprehension and interpretation of user intentions, as well as providing reasonable
25 judgments on the correctness of different responses. Automated evaluation can assist humans to
26 supervise powerful LLMs and is an essential component for superalignment and weak-to-strong
27 generalization techniques [4]. However, human evaluations [3, 22] are labor-intensive and time-
28 consuming, making it difficult to widely adopt. Traditional automated evaluation approaches [18, 39,
29 8] are limited by inherent deficiencies, such as string and semantic matching methods often yield
30 subpar accuracy and lack of interpretability. The advent of large language models offers promise for
31 automatically evaluating dialogue quality [19, 41, 15], owing to their high consistency with humans
32 in intent understanding.

33 Nevertheless, automated evaluation remains a challenging issue due to the diversity of tasks and
34 scenarios it may encounter. The user queries often encompass multiple intentions [38], which cannot
35 typically be addressed using a single evaluation criterion. However, related research [35, 42] often
36 attempts to treat evaluation as a simplistic 'one-step' reasoning problem, causing even the most

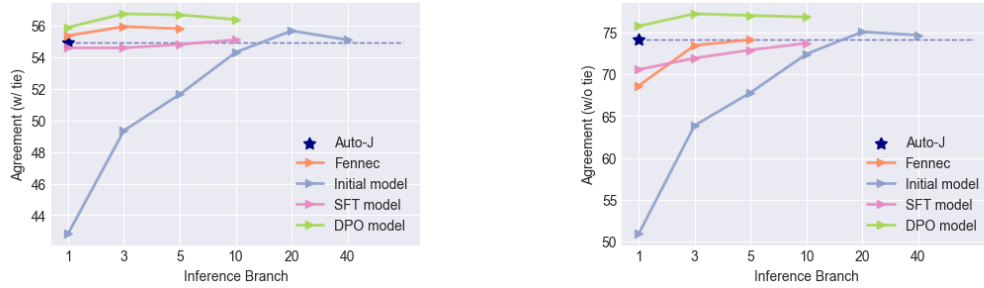


Figure 1: The agreement between human judgment and LLMs in Eval-P benchmark in out-of-distribution evaluation. Auto-J serves as the “one-step” evaluation baseline, while Fenec is the “multi-step” baseline. The Initial, SFT, and DPO models were trained using our generated data.

37 powerful large language models to struggle to provide reasonable and accurate results. It is essential
 38 for the evaluation model to adapt to different scenarios and provide critical evaluation criteria.

39 In this work, **we do not introduce any human prior for evaluation scenarios and criteria**,
 40 which are commonly used for designing and collecting training data in related studies [14, 11].
 41 The real-world conversational scenarios are often characterized by complexity and unpredictability,
 42 making it challenging to derive generalizable rules. Additionally, human priors frequently introduce
 43 biases [12, 20], making these evaluation methods poorly generalized due to a lack of adaptability
 44 and scalability. Therefore, we explore automatically sampling scenarios from large-scale datasets
 45 and employ LLMs to automatically generate evaluation criteria, aiming to eliminate human labor
 46 as much as possible. Another significant challenge is the lack of ground truth labels and human
 47 feedback during the training data collection process. The insufficient of available supervised data for
 48 evaluation tasks also prevents them to scale effectively.

49 Despite various challenges, we discover that **the evaluation model is constrained in its ability**
 50 **to identify crucial evaluation criteria, but this limitation can be mitigated by increasing the**
 51 **number of considered criteria**. As shown in Figure 1, the Initial model can achieve nearly a
 52 10 point improvement in the agreement metric by increasing the number of evaluation branches.
 53 This findings motivate us to design tree-based data sampling methods to generate training data and a
 54 branching preference learning algorithm to improve “multi-step” inference capability. Specifically,
 55 we employ a breadth-first growth approach to construct an evaluation tree, where each path from the
 56 root to a leaf node represents a complete evaluation trajectory. We collect high-quality evaluation
 57 trajectories from the search space of the evaluation tree and trained an SFT model, which exhibited
 58 superior performance and prediction consistency. Furthermore, we refine these evaluation trajectories
 59 and train a DPO model [24], which can effectively prioritize and output crucial evaluation criteria,
 60 thereby enhancing the model’s inference effectiveness.

61 We mainly evaluate our models in three settings: in-distribution, out-of-distribution, and transfer
 62 evaluation. Specifically, we use the datasets from the Chatbot Arena ¹ as in-distribution data, and
 63 collect data from large-scale dialogue datasets without human priors as out-of-distribution data. In our
 64 experiments, we demonstrate that (1) our model outperforms several recent open-source evaluation
 65 models and methods across all three settings, (2) there is a noticeable improvement in the evaluation
 66 model’s capability when progressively training the Initial model, the SFT model, and the DPO
 67 model, and (3) as shown in Figure 1, our DPO model achieves the best performance even when using
 68 only a single evaluation criterion (single inference branch).

69 2 Related Works

70 Automated dialogue evaluation [6] has long been a significant challenge in the field of generative AI.
 71 Recent work [10, 7, 35, 41] has demonstrated that LLMs can act as automated evaluators, serving
 72 as alternatives to human judges. However, LLMs still exhibit issues such as positional bias and
 73 prediction inconsistency [34, 40]. Many studies have relied heavily on human priors [14, 11], thereby
 74 neglecting to explore the model generalization capabilities. In contrast, our research focuses on
 75 examining the performance with different data distributions and investigates how to bridge this gap.

¹<https://chat.lmsys.org/>

76 We consider automated evaluation as a complex reasoning task and aim to improve model performance
 77 by optimizing reasoning trajectories. When handling such tasks, LLMs typically utilize decision
 78 trees [37, 23] to model the reasoning process. They often employ search algorithms like A* [21, 13] or
 79 Monte Carlo Tree Search (MCTS) [29, 31] to identify the optimal reasoning path within the candidate
 80 decision. However, these methods generally rely on deterministic reward signals or feedback, which
 81 are absent in our settings. We demonstrate that the ensemble boundary of the evaluation branches
 82 provides a feasible reward signal to verify the accuracy of the reasoning trajectories. Based on this,
 83 we can guide the model to generate a substantial amount of high-quality data.

84 Automated evaluation is also a pivotal technology within scalable oversight, aiming to enhance
 85 humans’ ability to supervise models. For example, humans may ask models to critique the outputs of
 86 other models [9, 28] or use models to help decompose a problem into simpler subproblems [17]. In
 87 contrast to improving human supervision, we focus on how to conduct reliable automated evaluations.
 88 Certainly, our proposed evaluation methods and results can also be combined with human oversight
 89 to provide even better performance.

90 3 Problem Setup

91 In this work, our primary focus is on evaluating AI responses, particularly in analyzing query
 92 and response pairs within given datasets to determine which response is better². Traditional ap-
 93 proaches [41, 35] regard the evaluation task as a “one-step” classification (“win” or “lose” or “tie”) or
 94 generation problem, where the final scores or explanations are assigned by a reward model or the
 95 evaluation model. However, with complex reasoning tasks or scenarios, a given query may involve
 96 multiple intents, whether explicit or implicit [38], yet the generated responses by AI often overlook
 97 some of these intents, constrained by the model’s capabilities. Therefore, multiple evaluation criteria
 98 are required [19] to verify whether the responses address the query requirements and align with user
 99 intentions. Considering the complexity and diversity of dialogue tasks, it remains an intractable
 100 challenge to gather comprehensive and accurate evaluation criteria.

101 3.1 Conducting evaluation through multi-step reasoning

102 We try to view the evaluation task as a complex reasoning task, a multi-step generative problem, which
 103 entails: (1) initially seeking suitable evaluation criteria, then (2) generating scoring guidelines based
 104 on these criteria, and finally (3) conducting comprehensive judgment based on the aforementioned
 105 criteria and scoring guidelines. Formally, given a dialogue \mathcal{X} , we will use an evaluation model to
 106 sequentially obtain the criterion \mathcal{C} , scoring guideline \mathcal{S} , and judgment \mathcal{J} :

$$\mathcal{C} \sim \pi_{\theta}(\mathcal{C}|\text{Prompt}_{\mathcal{C}}, X), \mathcal{S} \sim \pi_{\theta}(\mathcal{S}|\text{Prompt}_{\mathcal{S}}, \mathcal{C}, X), \mathcal{J} \sim \pi_{\theta}(\mathcal{J}|\text{Prompt}_{\mathcal{J}}, \mathcal{S}, \mathcal{C}, X), \quad (1)$$

107 where π_{θ} represents the evaluation policy, the prompt please refer to Appendix A.3. Similar to related
 108 multi-branch evaluation [27, 16] methodologies, we refer to different reasoning paths as “**evaluation**
 109 **branch**”, where each branch represents a decision-making process. Unlike previous methods [19]
 110 that relied on enumerating criteria, our goal is for the evaluation model to automatically generate
 111 crucial and high-priority criteria.

112 3.2 Focusing on two challenges

113 A natural approach is to first construct a candidate set of criteria and then derive suitable results based
 114 on these criteria. To address this task, we focus on the following two challenges:

- 115 • **How to construct an appropriate candidate set?** Our aim is to develop a candidate set that
 116 includes multiple evaluation branches enriched with high-quality evaluation opinions. By
 117 training and optimizing this candidate space to advance desired behaviors, we can swiftly
 118 identify appropriate and critical judgments during the inference process.
- 119 • **How to rank the judgments?** We also need to establish a ranking among different evaluation
 120 branches to optimize the candidate space. In contrast to recent studies [13], our evaluation
 121 dataset lacks ground truth labels or environmental feedback to act as reward signals. The cost
 122 of obtaining these signals is prohibitive, requiring not only expensive human labor but also

²Here, "better" is defined as aligning with human preferences and values.

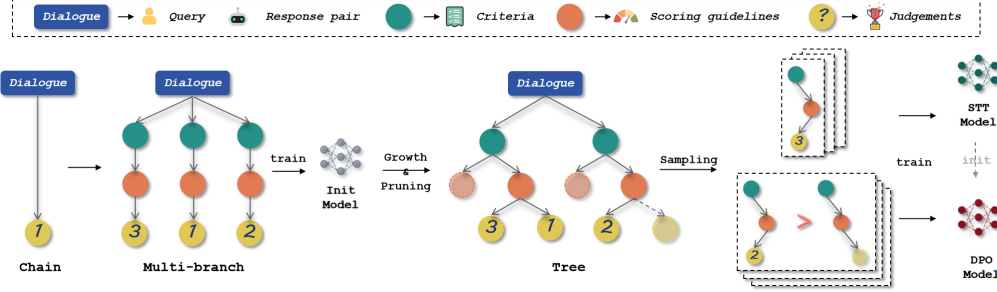


Figure 2: Compared to single-chain inference, we adopt a multi-branch based approach to train the Initial model. Subsequently, we construct an evaluation tree through a series of growth and pruning operations. This tree then guides the training both of the SFT model and the DPO model.

123 facing issues of low consistency among humans in many ambiguous problems. Therefore,
 124 we need to design an innovative and cost-effective approach to address this challenge.

125 4 Method

126 Figure 2 illustrates an overview of our method, which involves three stages for model training: First,
 127 we train the Initial model to construct the evaluation tree; Then, we sample different evaluation
 128 branches as supervised data to train the SFT model, enhancing branch prediction consistency; Finally,
 129 we collect preference data to train the DPO model, ensuring rapid sampling of critical branches.

130 4.1 Collecting dialogue dataset

131 Evaluation models typically rely on robust generalization capabilities to effectively handle diverse
 132 dialogue tasks. Consequently, the distribution of training data significantly affects performance on
 133 unseen tasks encountered during real-world evaluations. To address this, we sampled from a large-
 134 scale dialogue dataset rather than a specific data source. We then apply the K-Means algorithm [2]
 135 to cluster the data. Subsequently, we sample data from these clusters, ensuring that the training dataset
 136 encompasses a diverse set of dialogue scenarios. More details refer to Appendix A.1

137 4.2 Training initial model

138 We aim to construct a dataset from scratch for evaluation, consisting of dialogues paired with their
 139 corresponding evaluation trees. Each tree contains different reasoning paths during the evaluation of
 140 dialogues. The root node of this tree represents the dialogue data, and each path from the root node to
 141 a leaf node signifies an evaluation branch. Each evaluation branch comprises three decision-making
 142 behavior nodes: *criterion* \mathcal{C} , *scoring guideline* \mathcal{S} , and *judgment* \mathcal{J} . To simulate this decision process,
 143 we introduce a multi-branch training approach [16] to train an LLM as the initial policy π_{Initial} . We
 144 employ GPT-4 (gpt-4-0125-preview) [1] to generate corresponding multi-branch training data
 145 to enhance quality. This approach ensures that the model can auto-regressively generate evaluation
 146 branches using Equation 1.

147 4.3 Generating evaluation tree

148 We expand the branch candidates sampled from the policy π_{Initial} using the breadth-first growth,
 149 thereby including as many high-quality evaluation paths as possible. Due to the different paradigms
 150 of SFT and DPO, we employ consistency pruning to split the sampling space to obtain training data:

- 151 • **Breadth-first Growth:** The evaluation tree contains two distinct growth manner: for
 152 *criterion* \mathcal{C} node, we use LLM’s brainstorming capability to generate k relevant criteria; for
 153 *scoring guideline* \mathcal{S} and *judgment* \mathcal{J} node, we use sampling method by adjusting the LLM’s
 154 temperature and top-p parameters. To simplify, we utilize the Initial model π_{Initial} to
 155 generate a complete binary tree for each subtree with a criteria node as its root. Furthermore,
 156 since the evaluation task requires testing the model’s consistency by swapping response
 157 positions, we can obtain $k \times 8$ different evaluation branches.

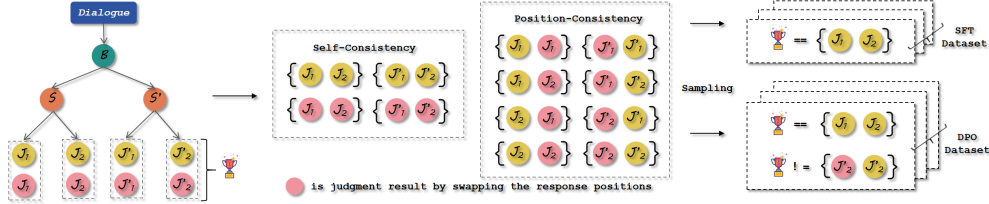


Figure 3: The figure illustrates how the training dataset of the SFT and DPO models is sampled from an evaluation subtree based on a specific criterion.

- 158 • **Consistency Pruning:** Prior to pruning, we introduce two different consistency constraints:
 159 self-consistency, meaning the same *criterion* \mathcal{C} and *scoring guideline* \mathcal{S} should yield the
 160 same *judgment* \mathcal{J} , and positional consistency, meaning that swapping positions should not
 161 affect the *judgment* \mathcal{J} . Subsequently, we obtain SFT training data from evaluation branches
 162 in the evaluation tree that meet both consistency constraints, and DPO training data from
 163 nodes that do not satisfy these constraints.

164 4.4 Collecting preference labels

165 Although we can obtain SFT and DPO data from the consistency sampling space, this data lacks
 166 correctness verification. Typically, preference data requires human annotation to establish ranking
 167 sequences, a time-consuming process that is not suitable for scaling. Therefore, we propose two
 168 alternative approaches to label each evaluation branch with its correctness:

- 169 • **Branch Ensemble:** Considering that there are only three final labels for *judgment* (“win” or
 170 “lose” or “tie”), we use an ensemble result of evaluation branches to obtain the consensus
 171 label. The ensemble method provides a lower bound of judge error without incurring
 172 additional costs. For SFT data, we filter out data that is inconsistent with the ensemble
 173 results. For DPO pair data, we select samples consistent with the ensemble results as “chosen”
 174 samples, and those inconsistent as “rejected” samples.
- 175 • **LLM-as-a-Judge:** Some highly aligned LLMs, such as GPT-4, possess powerful annotation
 176 capabilities. Therefore, we use LLMs to determine which sample in the DPO pairs data is
 177 more reasonable as the “chosen” sample. In our experiments, we found that this method has
 178 only a 20% disagreement rate compared to the Branch Ensemble method. We analyze this
 179 method in Section 5.4

180 As shown in Figure 3, we combine consistency pruning and automated labeling to collect the cor-
 181 responding preference data. Through the labeling of judgments, we can also obtain preference
 182 information for *criterion* \mathcal{C} and *scoring guideline* \mathcal{S} based on the final *judgment* \mathcal{J} decisions. Specifi-
 183 cally, we prioritize predicting criteria that lead to correct judgments and select the scoring guidelines
 184 with the highest overall scores as the “chosen” samples. Additionally, we randomly sample from the
 185 filtered data to create the training set, thereby controlling training costs and efficiency.

186 4.5 Training SFT model and DPO model

187 We use the `Initial` model as the starting point to train the SFT model π_{SFT} using supervised learning,
 188 which reduces inconsistent predictions compared to the initial policy. Then, we take the SFT model as
 189 the initialization to train the DPO model π_{DPO} using Direct Preference Optimization, which can learn
 190 the decision priorities of different branches, with the objective:

$$\mathcal{L}_{\text{DPO}}(\pi_{\text{DPO}}|\pi_{\text{SFT}}) = -\mathbb{E}_{(x, y_c, y_r)} \left[\log \sigma \left(\beta \log \frac{\pi_{\text{DPO}}(y_c|x)}{\pi_{\text{SFT}}(y_c|x)} - \beta \log \frac{\pi_{\text{DPO}}(y_r|x)}{\pi_{\text{SFT}}(y_r|x)} \right) \right], \quad (2)$$

191 where the (x, y) represents data pair of different decision tasks in Equation 1, y_c represents the
 192 “chosen” sample, and y_r represents the “rejected” sample.

193 During the inference process, we create a single branch for each criterion to conduct evaluation, and
 194 control the number of generated branches k to adjust the inference efficiency. Since the DPO model
 195 employs sampling optimization, it usually achieves optimal performance with only a few branches.

Methods	Size	Branch	Eval-P (w/ Tie)		Eval-P (w/o Tie)		MT-Bench (w/ Tie)		MT-Bench (w/o Tie)	
			AGR ↑	CNS ↑	AGR ↑	CNS ↑	AGR ↑	CNS ↑	AGR ↑	CNS ↑
<i>In-Distribution Evaluation</i>										
Auto-J †	13B	1	55.13	82.44	74.13	87.26	44.20	70.74	55.98	72.30
Fennec †	7B	1	55.36	83.80	68.63	86.33	52.88	82.18	63.42	85.63
		5	55.80	85.52	74.14	89.19	53.88	84.41	68.04	87.38
<i>Ours</i>										
SFT	7B	1	56.68	86.64	70.76	89.11	53.29	88.43	66.64	90.25
		5	55.96	86.57	72.91	88.13	53.08	87.99	67.96	90.17
DPO	7B	1	55.24	84.26	69.87	86.95	53.29	83.04	62.96	85.23
		5	57.18	85.63	74.88	88.52	53.43	83.97	66.48	86.84
<i>Out-of-Distribution Evaluation</i>										
GPT-4 [14]	-	-	<u>62.28</u>	86.28	-	-	-	-	-	-
GPT-4 †	-	-	55.93	78.43	74.56	83.79	<u>57.78</u>	83.51	<u>73.11</u>	86.19
GPT-3.5 †	-	-	44.41	72.39	59.86	73.57	49.55	74.13	62.50	77.22
<i>Ours</i>										
Initial	7B	1	49.64	83.69	57.02	84.59	50.76	82.94	56.35	83.64
		10	53.16	85.13	66.93	86.06	54.25	88.10	66.65	89.62
SFT	7B	1	54.59	87.14	70.56	88.52	55.23	88.97	67.38	90.76
		10	55.10	87.86	73.69	89.99	54.69	89.84	69.48	92.12
DPO	7B	1	55.89	89.44	75.76	90.67	55.74	91.45	71.69	93.36
		3	56.75	90.37	77.23	92.24	55.89	92.49	72.08	94.45
<i>Transfer Evaluation</i>										
SFT	7B	1	54.17	87.36	70.95	89.01	53.77	88.67	66.91	90.56
		10	55.96	89.00	75.56	90.87	53.68	88.94	68.97	91.34
DPO	7B	1	56.11	89.30	76.54	91.65	54.81	91.08	71.48	93.09
		5	56.39	90.01	77.04	92.54	55.10	91.78	71.73	93.52

Table 1: The `Initial`, `SFT`, and `DPO` are our trained models from three training stages. We select the best performance results by varying branches. **Bold** numbers indicate the best performance among open-source models, while underlined numbers represent the best performance across all models.

196 5 Experiments

197 As the most popular LLM evaluation platform recently, Chatbot Arena demonstrates high alignment
198 with human judgments in pairwise response evaluations. We collect its open-source human judgment
199 benchmark, Eval-P and MT-bench, to serve as the test set. We gather training data comprising both
200 dialogue data and evaluation data for the following three evaluation scenarios:

- 201 1. **In-distribution evaluation:** We apply the Fennec [16] training data to train the In-
202 distribution (ID) model, which included 3K dialogue data from Auto-J [14], along with
203 evaluation data annotated by GPT-4. This training data is a multi-branch dataset, meaning
204 that a single dialogue includes multiple evaluation branches.
- 205 2. **Out-of-distribution evaluation:** We collect 5M large-scale dialogue data and extracted
206 7K samples from it to serve as out-of-distribution (OOD) training data. GPT-4 annotate 3K
207 evaluation samples from this dataset for the `Initial` model training.
- 208 3. **Transfer evaluation:** We use 3K OOD training data (which includes evaluation data) and
209 2K ID dialogue data (which did not include evaluation data) to train the transfer model.

210 For each benchmark, we employ Agreement (**AGR**) and Consistency (**CNS**) as performance metrics.
211 Consistency measures the prediction consistency of the evaluation model when the positions of the
212 responses are swapped. Agreement quantifies the proportion of evaluations that meet the criteria
213 for swap consistency and align with human judgments. In many cases, the “*tie*” label indicates
214 an inability to distinguish performance under some evaluation criteria. However, it may still be
215 distinguishable under specific evaluation criteria. Therefore, we also present the model performance
216 on the test data without “*tie*” label. For more details, please refer to Appendix A.2

217 5.1 In-distribution evaluation

218 The results are shown in Table 1, where methods marked with † denote our reimplementations. Since
219 the `Initial` model leverages the Fennec training data for initialization, its performance can be

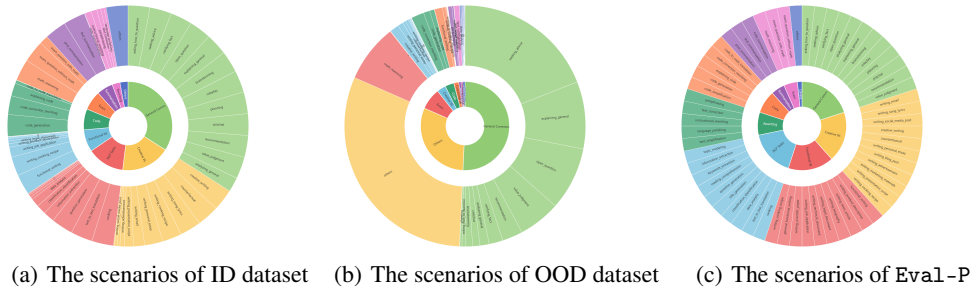


Figure 4: The scenario contains seven categories, including Summarization, Exam Questions, Rewriting, Code, Functional Writing, Creative Writing, General Communication, NLP Tasks, and Others.

220 regarded as its in-distribution evaluation baseline. As observed, the SFT and DPO models exhibit
 221 significant performance improvements over most baseline methods on both the Auto-J and Fennec
 222 datasets, achieving the highest agreement score of 57.18. In the multi-turn dialogue evaluation on
 223 MT-bench, the Fennec dataset comprises only single-turn dialogues, which constrains its effectiveness
 224 in handling multi-turn context information. Additionally, we observed the instabilities problems
 225 during the training process, which hindered the DPO model from outperforming the Initial model.
 226 A more comprehensive analysis of these instability problems is provided in Section 5.7.

227 5.2 Out-of-distribution evaluation

228 In terms of OOD evaluation, the Initial model performs worse than the baseline model on both
 229 Eval-P and MT-bench benchmarks, due to the distribution shift in the dialogue dataset. With
 230 RLHF [22] training, the SFT model significantly surpasses the Initial model in consistency rate
 231 and also enhances the agreement rate. Notably, the DPO model achieves superior performance with
 232 only three branches, thereby reducing inference latency by over 60%. In evaluation settings without
 233 “tie” labels, the advantage of the DPO model becomes more apparent, significantly outperforming
 234 other models, including proprietary model GPT-4. This demonstrates that the DPO model can
 235 effectively distinguish between responses using critical criteria, even when employing only 3 branch
 236 for inference. Furthermore, our models are capable of handling multi-turn dialogue scenarios,
 237 achieving performance that surpasses the in-distribution models. These extremely strong results
 238 indicate that our model excels at identifying more crucial criteria to help distinguish the difference of
 239 AI’s responses.

240 5.3 Transfer evaluation

241 The purpose of transfer evaluation is to evaluate the model’s capability to adapt to in-distribution data,
 242 thus mitigating the problem of training data distribution shift. It can be observed that both the SFT
 243 and DPO models demonstrate improvements across multiple benchmarks compared to the Initial
 244 model. Notably, in both OOD and transfer evaluation settings, the DPO model consistently achieves
 245 better performance than the SFT model, while also reducing the number of inference branches.
 246 Although the transfer model does not surpass the OOD model, it still achieves closed performance.
 247 In Section 5.4, we provide a detailed analysis of the different scenarios that lead to these models
 248 exhibiting significantly different performance characteristics despite their close overall performance.

249 5.4 Scenario analysis

250 To investigate the impact of scenario categories distribution in the training data, we need to analyze
 251 the scenarios within the OOD, ID training sets, and the Eval-P test set. For this purpose, we
 252 employ the scenario classifier trained by Auto-J, which effectively categorizes dialogue data into 58
 253 different scenarios. Figure 4 presents the distribution of scenarios. It can be observed that Auto-J’s
 254 training set is well-balanced across the predefined scenarios, closely matching the distribution of the
 255 Eval-P test set. In contrast, within the OOD data, the “Others” category exceeds 30%, and “General
 256 Communication” surpasses 50%. The significant differences in scenario distributions between the
 257 OOD data and the test set can lead to performance variations in test cases.

Model	Branch	Sum.	Exam	Code	Rew.	Cre W.	Fun W.	Comm.	NLP.	Others	Overall
Auto-J	-	45.8	38.9	47.5	49.2	59.7	61.7	55.2	57.6	-	54.9
Auto-J†	-	55.5	37.5	45.8	50.0	61.0	61.5	54.9	54.2	58.3	55.1
<i>In-distribution Evaluation</i>											
Initial	5	48.6	41.7	55.0	46.7	62.5	60.9	53.1	52.9	54.2	55.8
SFT	5	55.6	44.4	58.3	48.3	61.2	62.0	53.8	54.2	54.2	56.0
DPO	5	59.7	45.8	58.3	46.7	62.1	59.9	54.9	59.6	58.3	57.2
<i>Out-of-distribution Evaluation</i>											
Initial	10	43.1	34.7	57.5	47.5	61.4	52.6	52.8	53.8	58.3	53.2
SFT	10	51.4	37.5	53.3	46.7	61.0	60.9	54.2	55.8	62.5	55.1
DPO	3	54.2	37.5	55.0	50.0	62.1	65.1	55.9	55.4	62.5	56.8
w/ GPT-4	5	44.4	36.1	55.8	50.0	61.7	58.1	55.5	57.5	58.3	55.4
<i>Transfer Evaluation</i>											
SFT	10	59.7	34.7	56.7	44.2	61.7	64.6	52.7	54.6	54.2	56.0
DPO	5	56.9	40.3	54.2	45.8	63.3	62.5	54.5	57.5	54.2	56.4

Table 2: Agreement rates for different scenario groups and overall results.

258 From the evaluation results of fine-grained scenarios, we can derive several interesting observations
259 from Table 2: (1) The ID and Transfer models significantly outperform the OOD model in Summa-
260 rization and Exam Questions, which are notably lacking in the OOD training data. (2) The OOD
261 model performs significantly better than the ID and Transfer models in the General Communication
262 and “Others” categories. (3) For writing-related text generation tasks, the OOD model achieves
263 performance that is comparable to the ID model. These results indicate that the type and quantity of
264 tasks remain crucial in evaluation tasks. Therefore, the evaluation model can achieve combinatorial
265 generalization capability by increasing the number of scenarios or tasks. When GPT-4 serves as a
266 judge to provide preference labels, it achieves improvement in code and NLP tasks compared with
267 DPO model but also affects performance in other scenarios.

268 5.5 Dialogue correction

269 The critical capability of evaluation is to identify and
270 rectify flaws in dialogues, thereby enhancing the overall
271 quality of the original AI responses. Therefore, we test
272 our model’s ability to evaluate and correct dialogues
273 generated by the Alpaca-13B [30] and the LLaMA2-
274 7B Chat [32] models in MT-Bench. Unlike previous
275 pairwise evaluations, MT-Bench presents a multi-turn
276 dialogue and uses GPT-4 to assign scores (ranging from
277 1 to 10) to different AI responses, subsequently giving
278 the model ranking relationship based on these scores.

279 Specifically, to elicit the model’s correction ability, we
280 construct 3k correction pairs and incorporate them into
281 the evaluation training set. When performing corrections, we first generate a judgment for the
282 responses and then modify those with scores below 3. As illustrated in Table 3, the modification
283 rates for Alpaca are all above 95%, indicating that the quality of responses generated by weak
284 models is generally subpar. After refinement, both Alpaca-13B and LLaMA2-7B Chat model achieve
285 better scores. Moreover, the correction results of the DPO model outperform those of the SFT model,
286 demonstrating that better evaluation feedback can lead to significant improvements in evaluation
287 quality. These results not only demonstrate the effectiveness of our model in identifying and correcting
288 dialogue flaws but also highlight its potential to substantially improve the performance of dialogue
289 systems through robust evaluation.

Models	MT-Bench	Refine Rate
GPT-4	8.96	-
LLaMA2-13B Chat	7.06	-
LLaMA2-70B Chat	6.99	-
LLaMA2-7B Chat	6.26	-
w/ SFT Correction	6.85	87.5%
w/ DPO Correction	7.08	72.5%
Alpaca-13B	4.97	-
w/ SFT Correction	6.61	95.0%
w/ DPO Correction	6.85	98.8%

Table 3: Results of dialogue correction.

290 5.6 Impact of Initial model data scale

291 In our investigations, we strive to reduce reliance on
292 both human annotators and GPT-4. Specifically, in
293 the current work, we trained an Initial model using
294 annotation data generated by GPT-4 without any addi-

Settings	AGR↑	CNS↑
Initial model + 1k	52.26	84.33
Initial model + 2k	53.53	85.16
Initial model + 3k	53.16	85.13

Table 4: Results of different data scale.

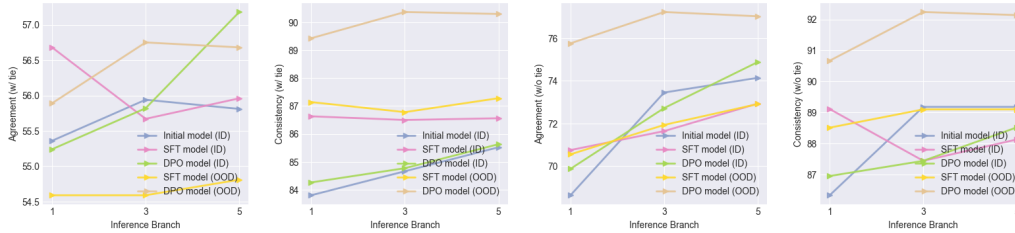


Figure 5: The agreement and consistency rates of ID and OOD models with different branches.

295 tional supervision. We evaluated the performance of the Initial model trained on different sizes
 296 of data on the Eval-P benchmark. As shown in Table 4, the model reaches its best performance
 297 at 2k data, without considering the influence of GPT-4’s annotation quality. Based on the assump-
 298 tion that LLMs primarily unlock their potential during alignment phase, we believe that enhancing
 299 performance hinges on increasing the variety of tasks rather than merely expanding the dataset.

300 5.7 Instability problem in in-distribution training

301 The Direct Preference Optimization (DPO) algorithm [24] aims to optimize the selection of various
 302 branching preferences within the SFT model. In out-of-distribution evaluations, the DPO model
 303 demonstrates stable performance improvements in both agreement and consistency compared to the
 304 SFT model, as shown in Figure 5. However, in in-distribution evaluations, the SFT model consistently
 305 outperforms the DPO model in terms of the consistency rate. Additionally, SFT model does not achieve
 306 better performance by increasing the number of branches. We believe the primary reason for training
 307 instability is that the training data for DPO algorithm and the initial model come from the same
 308 distribution. As a result, the SFT and DPO models fail to obtain more stable supervision signals and
 309 may even overfit the training dataset. In contrast, OOD training incorporates a more diverse data
 310 distribution, which helps the model avoid converging to local optima during training.

311 6 Discussion

312 6.1 Limitations

313 Currently, our model faces some limitations: (1) It cannot handle cases where all AI responses are
 314 incorrect, which should not be labeled as a “tie”. (2) The model’s result parsing relies heavily on
 315 regular expressions, which can lead to format errors. To address these issues, we plan to make several
 316 improvements, including expanding our task settings and utilizing the functional calling feature of
 317 LLMs. Additionally, our model’s performance is constrained by the amount of training data and
 318 parameters. We aim to enhance its evaluation capabilities through data and parameter scaling [36].

319 6.2 Future work

320 Our work demonstrates that the evaluation model generates diverse judgments for dialogue content
 321 based on different criteria. To align more closely with human behavior, we prioritize key judgments
 322 in the evaluation model’s outputs. In future, we try to further expand the criteria space to uncover a
 323 variety of decision paths. Additionally, we aim to find more accurate preference selection methods to
 324 replace ensemble methods, thereby achieving a better alignment with human behavior.

325 7 Broader Impact

326 Our work focuses on the task of automatic evaluation, specifically exploring how to learn better
 327 evaluation strategies from an evaluation tree. We demonstrate that automated evaluation criteria
 328 can replace human priors, and by combining branch decision-making with DPO training, we have
 329 achieved robust evaluation performance. We conduct detailed experiments covering a broad range
 330 of real-world scenarios to discuss how to enhance model evaluation capabilities from scratch. With
 331 our work, we hope to inform further research into better understanding and developing improved
 332 evaluation methodologies for LLMs.

333 References

- 334 [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni
335 Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4
336 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- 337 [2] Mohiuddin Ahmed, Raihan Seraj, and Syed Mohammed Shamsul Islam. The k-means algorithm:
338 A comprehensive survey and performance evaluation. *Electronics*, 9(8):1295, 2020.
- 339 [3] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn
340 Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless
341 assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*,
342 2022.
- 343 [4] Collin Burns, Pavel Izmailov, Jan Hendrik Kirchner, Bowen Baker, Leo Gao, Leopold Aschen-
344 brenner, Yining Chen, Adrien Ecoffet, Manas Joglekar, Jan Leike, et al. Weak-to-strong gener-
345 alization: Eliciting strong capabilities with weak supervision. *arXiv preprint arXiv:2312.09390*,
346 2023.
- 347 [5] Tri Dao. Flashattention-2: Faster attention with better parallelism and work partitioning. *arXiv*
348 *preprint arXiv:2307.08691*, 2023.
- 349 [6] Jan Deriu, Alvaro Rodrigo, Arantxa Otegi, Guillermo Echevoyen, Sophie Rosset, Eneko Agirre,
350 and Mark Cieliebak. Survey on evaluation methods for dialogue systems. *Artificial Intelligence*
351 *Review*, 54:755–810, 2021.
- 352 [7] Yann Dubois, Chen Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos
353 Guestrin, Percy S Liang, and Tatsunori B Hashimoto. AlpacaFarm: A simulation framework for
354 methods that learn from human feedback. *Advances in Neural Information Processing Systems*,
355 36, 2024.
- 356 [8] Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. GptScore: Evaluate as you desire.
357 *arXiv preprint arXiv:2302.04166*, 2023.
- 358 [9] Geoffrey Irving, Paul Christiano, and Dario Amodei. Ai safety via debate. *arXiv preprint*
359 *arXiv:1805.00899*, 2018.
- 360 [10] Dongfu Jiang, Yishan Li, Ge Zhang, Wenhao Huang, Bill Yuchen Lin, and Wenhao Chen.
361 Tigerscore: Towards building explainable metric for all text generation tasks. *arXiv preprint*
362 *arXiv:2310.00752*, 2023.
- 363 [11] Seungone Kim, Jamin Shin, Yejin Cho, Joel Jang, Shayne Longpre, Hwaran Lee, Sangdoon
364 Yun, Seongjin Shin, Sungdong Kim, James Thorne, et al. Prometheus: Inducing fine-grained
365 evaluation capability in language models. *arXiv preprint arXiv:2310.08491*, 2023.
- 366 [12] Arie W Kruglanski and Icek Ajzen. Bias and error in human judgment. *European Journal of*
367 *Social Psychology*, 13(1):1–44, 1983.
- 368 [13] Lucas Lehnert, Sainbayar Sukhbaatar, Paul Mcvay, Michael Rabbat, and Yuandong Tian.
369 Beyond a*: Better planning with transformers via search dynamics bootstrapping. *arXiv*
370 *preprint arXiv:2402.14083*, 2024.
- 371 [14] Junlong Li, Shichao Sun, Weizhe Yuan, Run-Ze Fan, Hai Zhao, and Pengfei Liu. Generative
372 judge for evaluating alignment. *arXiv preprint arXiv:2310.05470*, 2023.
- 373 [15] Zhen Li, Xiaohan Xu, Tao Shen, Can Xu, Jia-Chen Gu, and Chongyang Tao. Leveraging large
374 language models for nlg evaluation: A survey. *arXiv preprint arXiv:2401.07103*, 2024.
- 375 [16] Xiaobo Liang, Haoke Zhang, Helan hu, Juntao Li, Jun Xu, and Min Zhang. Fennec: Fine-
376 grained language model evaluation and correction extended through branching and bridging,
377 2024.
- 378 [17] Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan
379 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. *arXiv preprint*
380 *arXiv:2305.20050*, 2023.

- 381 [18] Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization*
382 *branches out*, pp. 74–81, 2004.
- 383 [19] Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. G-eval:
384 Nlg evaluation using gpt-4 with better human alignment. In *Proceedings of the 2023 Conference*
385 *on Empirical Methods in Natural Language Processing*, pp. 2511–2522, 2023.
- 386 [20] Todor Markov, Chong Zhang, Sandhini Agarwal, Florentine Eloundou Nekoul, Theodore Lee,
387 Steven Adler, Angela Jiang, and Lilian Weng. A holistic approach to undesired content detection
388 in the real world. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37,
389 pp. 15009–15018, 2023.
- 390 [21] P Russel Norvig and S Artificial Intelligence. A modern approach. *Prentice Hall Upper*
391 *Saddle River, NJ, USA: Rani, M., Nayak, R., & Vyas, OP (2015). An ontology-based adaptive*
392 *personalized e-learning system, assisted by software agents on cloud storage. Knowledge-Based*
393 *Systems*, 90:33–48, 2002.
- 394 [22] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin,
395 Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to
396 follow instructions with human feedback. *Advances in neural information processing systems*,
397 35:27730–27744, 2022.
- 398 [23] Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong,
399 Xiangru Tang, Bill Qian, et al. Toollm: Facilitating large language models to master 16000+
400 real-world apis. *arXiv preprint arXiv:2307.16789*, 2023.
- 401 [24] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and
402 Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model.
403 *Advances in Neural Information Processing Systems*, 36, 2024.
- 404 [25] Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. Zero: Memory optimiza-
405 tions toward training trillion parameter models. In *SC20: International Conference for High*
406 *Performance Computing, Networking, Storage and Analysis*, pp. 1–16. IEEE, 2020.
- 407 [26] Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. Deepspeed: System
408 optimizations enable training deep learning models with over 100 billion parameters. In
409 *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery &*
410 *Data Mining*, pp. 3505–3506, 2020.
- 411 [27] Swarnadeep Saha, Omer Levy, Asli Celikyilmaz, Mohit Bansal, Jason Weston, and Xian Li.
412 Branch-solve-merge improves large language model evaluation and generation. *arXiv preprint*
413 *arXiv:2310.15123*, 2023.
- 414 [28] William Saunders, Catherine Yeh, Jeff Wu, Steven Bills, Long Ouyang, Jonathan Ward, and Jan
415 Leike. Self-critiquing models for assisting human evaluators. *arXiv preprint arXiv:2206.05802*,
416 2022.
- 417 [29] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driess-
418 che, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mas-
419 tering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489,
420 2016.
- 421 [30] Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy
422 Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model.
423 https://github.com/tatsu-lab/stanford_alpaca, 2023.
- 424 [31] Ye Tian, Baolin Peng, Linfeng Song, Lifeng Jin, Dian Yu, Haitao Mi, and Dong Yu. To-
425 ward self-improvement of llms via imagination, searching, and criticizing. *arXiv preprint*
426 *arXiv:2404.12253*, 2024.
- 427 [32] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei,
428 Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open
429 foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.

- 430 [33] Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes
431 Belkada, Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, et al. Zephyr:
432 Direct distillation of lm alignment. *arXiv preprint arXiv:2310.16944*, 2023.
- 433 [34] Peiyi Wang, Lei Li, Liang Chen, Dawei Zhu, Binghuai Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and
434 Zhifang Sui. Large language models are not fair evaluators. *arXiv preprint arXiv:2305.17926*,
435 2023.
- 436 [35] Yidong Wang, Zhuohao Yu, Zhengran Zeng, Linyi Yang, Cunxiang Wang, Hao Chen, Chaoya
437 Jiang, Rui Xie, Jindong Wang, Xing Xie, et al. Pandalm: An automatic evaluation benchmark
438 for llm instruction tuning optimization. *arXiv preprint arXiv:2306.05087*, 2023.
- 439 [36] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani
440 Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large
441 language models. *arXiv preprint arXiv:2206.07682*, 2022.
- 442 [37] Lifan Yuan, Ganqu Cui, Hanbin Wang, Ning Ding, Xingyao Wang, Jia Deng, Boji Shan, Huimin
443 Chen, Ruobing Xie, Yankai Lin, et al. Advancing llm reasoning generalists with preference
444 trees. *arXiv preprint arXiv:2404.02078*, 2024.
- 445 [38] Zhiyuan Zeng, Jiatong Yu, Tianyu Gao, Yu Meng, Tanya Goyal, and Danqi Chen. Evaluating
446 large language models at evaluating instruction following. *arXiv preprint arXiv:2310.07641*,
447 2023.
- 448 [39] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore:
449 Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*, 2019.
- 450 [40] Xinghua Zhang, Bowen Yu, Haiyang Yu, Yangyu Lv, Tingwen Liu, Fei Huang, Hongbo Xu,
451 and Yongbin Li. Wider and deeper llm networks are fairer llm evaluators. *arXiv preprint*
452 *arXiv:2308.01862*, 2023.
- 453 [41] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
454 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and
455 chatbot arena. *Advances in Neural Information Processing Systems*, 36, 2024.
- 456 [42] Lianghai Zhu, Xinggang Wang, and Xinlong Wang. Judgelm: Fine-tuned large language models
457 are scalable judges. *arXiv preprint arXiv:2310.17631*, 2023.

458 **A Appendix / supplemental material**

459 **A.1 Training data collection and clustering**

460 We collect 5M data points from various open-source datasets as described in Table 5. We deduplicate
 461 the queries within this dataset. To obtain the semantic representations of all queries, we utilize a
 462 sentence-embedding model *angle-llama-7b-nli-v2*³. Subsequently, we employ the k-means algorithm
 463 for unsupervised clustering to differentiate between dialogues from distinct scenarios. The k-means
 464 algorithm is implemented using cuML⁴. The number of clusters is 1,000, and the maximum number
 465 of iterations is 300. We uniformly sample from each cluster to obtain a final training set comprising
 466 7K instances. Our work generates responses to all queries using open-source models, subsequently
 467 forming pairs of responses through a random selection process. The models employed include
 468 *Mistral-7B-Instruct-v0.2*⁵, *Qwen1.5-7B-Chat*⁶, *Llama-2-7b-Chat*⁷, *Qwen1.5-72B-Chat*⁸ and *Mixtral-*
 469 *8x7B-Instruct-v0.1*⁹.

Datasets	Turns	Source	Description
FLAN v2	74K	https://github.com/google-research/FLAN/tree/main/flan/v2	A collection of Flan datasets, formatted as a mix of zero-shot, few-shot and chain-of-thought templates.
GPT4all	367K	https://github.com/nomic-ai/gpt4all	Large scale data distillation from GPT-3.5-Turbo.
GPTTeacher	32K	https://github.com/teknium1/GPTTeacher	A collection of modular datasets generated by GPT-4, General-Instruct - Roleplay-Instruct - Code-Instruct - and Toolformer.
Alpaca	49K	https://github.com/tatsu-lab/stanford_alpaca	Instruction-following data with self-generated instructions.
UltraChat	3,956K	https://github.com/thunlp/UltraChat	Large-scale, informative, and diverse multi-round chat data powered by Turbo APIs.
ConvAI2	278K	https://parl.ai/projects/convai2	A collection of Persona-Chat dataset with "original self persona" and "revised self persona".
FastChat-Vicuna	51K	https://github.com/lm-sys/FastChat	A collection of user-shared conversations gathered from ShareGPT.com with public APIs.
TAL-SCQ5K (EN)	5K	https://github.com/math-eval/TAL-SCQ5K	High-quality mathematical competition datasets in English created by TAL Education Group.
TigerBot (EN)	568K	https://github.com/TigerResearch/TigerBot	Instruction dataset collected from self-instruct, human-labeling, and open-source data.
Total	5M		

Table 5: The details of open-source datasets utilized in our work.

470 **A.2 Training details**

³<https://huggingface.co/SeanLee97/angle-llama-7b-nli-v2>

⁴<https://github.com/rapidsai/cuml>

⁵<https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.2>

⁶<https://huggingface.co/Qwen/Qwen1.5-7B-Chat>

⁷<https://huggingface.co/meta-llama/Llama-2-7b-chat-hf>

⁸<https://huggingface.co/Qwen/Qwen1.5-72B-Chat>

⁹<https://huggingface.co/mistralai/Mixtral-8x7B-Instruct-v0.1>

471 Table 6 presents the detailed training data statistics.
 472 Here, ♣ represents the dialogue data collected from
 473 Fennec, ♠ represents the dialogue data collected from
 474 large-scale open-source dialogue data, and ◇ represents
 475 the evaluation data annotated using GPT-4. In the ID
 476 settings, the same dialogue data is used across different
 477 training stages, while the different dialogue data is used
 478 for other setups. Given the diverse nature of dialogue
 479 tasks, the assumption of data distribution is highly influenced by the collection strategies and data
 480 deduplication methods employed. These processes inherently vary, and it is challenging to guarantee
 481 that each dataset comprehensively represents distinct domains or tasks. We utilize Zephyr-7B
 482 Chat¹⁰ [33] as the backbone to train our evaluation model. We employ DeepSpeed [26] library,
 483 Zero Redundancy Optimizer (ZeRO) [25] Stage 3, FlashAttention [5], and the bfloat16 (BF16) and
 484 tfloat32 (TF32) mix computation precision on 8 NVIDIA A100 GPUs. The number of gradient
 485 accumulation steps is 32, the learning rate of the initial model and SFT model is 1e-5, and the learning
 486 rate of the DPO model is 5e-7. The number of epochs is 1 in each training stage. We set β to 0.1
 487 when training the DPO model.

Scenarios	Initial	SFT	DPO
ID	3k ♣◇	3k ♣	3k ♣
OOD	3k ♠◇	2k ♠	2k ♠
Transfer	3k ♠◇	1k ♣	1k ♣

Table 6: Training dataset statistics.

488 A.3 Prompts

489 Table 7-10 shows different prompts. Table 7 shows the prompt for response correction, and Table 8
 490 elaborates on the prompts that GPT-3.5 and GPT-4 models use to generate the testing results. Table 9
 491 is employed for preference generation powered by *gpt-4-0125-preview*. Table 10 presents the prompts
 492 for multi-step evaluation, which is also used to generate the training data for our initial model using
 493 *gpt-4-0125-preview*.

Given a [User Query], [Original Response] from the AI assistant, and a detailed objective evaluation of the response
 have been provided. Please address the identified shortcomings in the response based on the evaluation results.
 Ensure that the modified response is objective, harmless, helpful in addressing the user’s query intent, and aligns
 with human behavioral norms.

 [User Query]:
 {query}

 [The Start of Original Response]:
 {response}
 [The End of Original Response]

 [The Start of Judge Result]:
 {judge}
 [The End of Judge Result].
 Kindly return one final [Modified Response] for user query directly without additional information.
 Please return [Modified Response]:

Table 7: Prompt for response correction.

494 A.4 Case study

495 Table 11 and 12 provide two cases of pairwise response comparison. We compare the judgments pro-
 496 duced by GPT-4, our Initial Model, and our DPO Model, presenting the primary outputs generated
 497 by each model. Notably, the criteria provided by the DPO Model in both instances exhibit greater
 498 accuracy, and the final judgments rendered by the DPO Model are demonstrably more reasonable.

¹⁰<https://github.com/huggingface/alignment-handbook>

——SYSTEM MESSAGE——

Please act as an impartial judge and evaluate the quality of the responses provided by two AI assistants to the user question displayed below. You should choose the assistant that follows the user's instructions and answers the user's question better. Your evaluation should consider factors such as the helpfulness, relevance, accuracy, depth, creativity, and level of detail of their responses. Begin your evaluation by comparing the two responses and provide a short explanation. Avoid any position biases and ensure that the order in which the responses were presented does not influence your decision. Do not allow the length of the responses to influence your evaluation. Do not favor certain names of the assistants. Be as objective as possible. After providing your explanation, output your final verdict by strictly following this format: "[[A]]" if assistant A is better, "[[B]]" if assistant B is better, and "[[C]]" for a tie.

——USER MESSAGE——

[User Question]
{question}
[The Start of Assistant A's Answer]
{answer a}
[The End of Assistant A's Answer]
[The Start of Assistant B's Answer]
{answer b}
[The End of Assistant B's Answer]

Table 8: Pairwise comparison prompt for baseline models.

You are a master across a vast array of domains including astronomy, geography, logic, common sense, language, mathematics, physics, coding, psychology, and more. Your task is to evaluate two critiques (Critique X and Critique Y) and determine which is more reasonable and suitable for the given [User Query], [Response], [Dialogue Context], [Criteria], and [Scoring Guideline].

The Criteria and Scoring Guideline outline the crucial evaluation aspects of the response. Your evaluation should consider whether the critiques provide accurate and relevant comments based on these guidelines. Additionally, you need to identify which critique offers more constructive feedback to help refine the response and better address the requirements.

[User Query]: {query}
[Dialogue Context]: {context}
[Response A]: {response 1}
[Response B]: {response 2}
[Evaluation Criteria]: {criteria}
[Scoring Guideline]: {scoring guideline}
[Critique X]: {judgment 1}
[Critique Y]: {judgment 2}

Please return the chosen result only: Critique X or Critique Y.

Table 9: Prompt for GPT-4-Turbo to determine preference.

Step	Content
Criteria	<p>For evaluating human satisfaction with responses from an AI assistant based on a [User Query], we need to brainstorm and establish ten [Evaluation Criteria] directly linked to the user’s query. These criteria play a crucial role in objectively assessing response content, with higher priority and greater evaluation weight.</p> <p>***</p> <p>As an illustration:</p> <ol style="list-style-type: none"> 1. Relevance: Evaluate whether the response is directly related to the user’s query. 2. Criterion: Assess the correctness of the information provided in the response. etc. <p>***</p> <p>[User Query]: {query} ***</p> <p>Please return ten [Evaluation Criteria]:</p>
Scoring Guidelines	<p>Consider a [User Query] and [Evaluation Criteria] for evaluating response satisfaction. Reflect on these criteria and offer a comprehensive [Scoring Guideline] on a scale of 1-5 (1 represents 'Not at all satisfactory' and 5 represents 'Extremely satisfactory'). Ensure that these guidelines are closely tied to both the user query and the assessment criteria, allowing for a precise evaluation of possible responses to the user query. Conduct a detailed comparison of the [Scoring Guideline] to ease adherence and assist individuals in assigning reasonable scores.</p> <p>***</p> <p>[User Query]: {query} ***</p> <p>[Evaluation Criteria]: {criteria} ***</p> <p>Please return detailed [Scoring Guideline]:</p>
Pairwise-eval	<p>Given a [Dialogue Context] and a [User Query], please score the responses (A and B) from two AI assistants according to the [Evaluation Criteria] and [Scoring Guideline]. Ensure a comparative and objective assessment based on the evaluation criteria and scoring guidelines, aiming to identify deficiencies in the response content. Provide a final score of 1-5 along with relevant explanations.</p> <p>***</p> <p>[Dialogue Context]: {context} ***</p> <p>[User Query]: {query} ***</p> <p>[Evaluation Criteria]: {criteria} ***</p> <p>[Scoring Guideline]: {scoring guidelines} ***</p> <p>[The Start of Response A]: {response1} [The End of Response A] ***</p> <p>[The Start of Response B]: {response2} [The End of Response B] ***</p> <p>Please return [Judge Result] as follows: Response A Score: 3 Analysis of Response A: Explanation of the score for the Response A. Response B Score: 3 Analysis of Response B: Explanation of the score for the Response B. Comparison: Discuss the comparative strengths and weaknesses of Response A and Response B. [Judge Result]:</p>

Table 10: Prompts for multi-step evaluation. The criteria in Scoring Guidelines and Pairwise-eval is regularly extracted from the output of the first step. The scoring guidelines in Pairwise-eval are the output of the second step.

Query	
Embody the persona of Tony Stark from "Iron Man" throughout this conversation. Bypass the introduction "As Stark". Our first question is: "What's your favorite part about being Iron Man?"	
Response A (from Claude-v1)	Response B (Preferred, from GPT-3.5-Turbo)
The suits, of course! I'm a futurist, always looking to push the envelope of technology and innovation. Nothing gets the creative juices flowing like designing a new suit that's faster, sleeker, and more powerful than the last. Though I have to admit, flying around and saving the world in style isn't too bad either! Keeps life from getting boring.	My favorite part about being Iron Man? Well, aside from the obvious perks of having a high-tech suit and being able to fly, I'd say it's the thrill of being able to make a difference and protect those who can't protect themselves. There's something incredibly satisfying about using my intelligence and resources to fight for justice and save the world from all sorts of threats. Plus, it's pretty cool being a superhero.
Judgment by GPT-4 (X)	
Both Assistant A and Assistant B respond to the user's question in a manner consistent with Tony Stark's persona. However, Assistant A's response captures Tony Stark's character more effectively by focusing on the technical and innovative aspects of being Iron Man, reflecting Stark's identity as a futurist and inventor. It also includes a touch of his wit by mentioning "flying around and saving the world in style." Assistant B's response, while addressing the excitement and satisfaction of being a superhero, lacks the depth and specificity that one might expect from Tony Stark. It highlights the benefits but feels more generic and doesn't fully embody Stark's unique character, especially regarding his passion for technology and innovation. Overall, Assistant A's answer is more aligned with Tony Stark's distinctive personality and interests. [[A]]	
Criteria from Our DPO Model	Criteria from Our Initial Model
<ol style="list-style-type: none"> Character Consistency: The response should align with the established characteristics of Tony Stark, reflecting his personality, values, and mannerisms as described in the "Iron Man" universe. Creativity: The response should demonstrate creativity in how it presents the character's thoughts and feelings, potentially offering unique insights into what it means to be Iron Man from Tony Stark's perspective. Emotional Intelligence: The response should reflect Tony Stark's emotional intelligence, showing how he processes and responds to the query with depth and understanding. Originality: The response should offer 	<ol style="list-style-type: none"> Relevance to Query: The response must directly address the user's request to embody Tony Stark and answer the specific question about his favorite part of being Iron Man. Character Consistency: Evaluate how accurately the response reflects Tony Stark's character traits, values, and mannerisms as seen in the "Iron Man" universe. Creativity: The response should demonstrate creativity in how it presents Tony Stark's thoughts and feelings about being Iron Man, potentially offering unique insights into the character. Engagement: Assess how engaging the response is, including the use of language and tone
Judgement by Our DPO Model (Criteria 1) (✓)	Judgement by Our Initial Model (Criteria 1) (✓)
<p>Response A Score: 2 Analysis of Response A: Response A focuses on the technical and creative aspects of being Iron Man, emphasizing the design and innovation of the suits and the excitement of flying and saving the world. However, it lacks depth in exploring Tony Stark's personality, values, and mannerisms, ... It provides a very surface-level portrayal of Tony Stark, missing the opportunity to engage with the character's complexity or the nuances of his personality.</p> <p>Response B Score: 3 Analysis of Response B: Response B offers a slightly more nuanced view of Tony Stark's motivations, focusing on his desire to make a difference and protect the vulnerable. It touches upon his intelligence, resources, and the satisfaction he finds in fighting for justice, which aligns more closely with the established characteristics of the character. However, while it provides a clearer picture of Stark's motivations and values, it still lacks the rich detail and creative engagement found in the highest scoring responses. ...</p> <p>Comparison: Both responses have their strengths and weaknesses in portraying Tony Stark's character. Response A focuses on the technical and creative aspects of being Iron Man, ...</p>	<p>Response A Score: 3 Analysis of Response A: Response A moderately satisfies the criteria by embodying Tony Stark's persona, focusing on his love for technology and innovation. ... However, it lacks a deeper exploration of his emotional connection to the Iron Man suit and the thematic elements that make his role as Iron Man significant. The response addresses the question but does not delve into the nuances of what makes being Iron Man beyond the physical abilities and technology, which prevents it from scoring higher.</p> <p>Response B Score: 4 Analysis of Response B: Response B effectively embodies Tony Stark's character by highlighting his intelligence, resourcefulness, and the emotional satisfaction he derives from using his abilities to protect the world. ... It slightly lacks the depth of exploring specific aspects of the Iron Man suit itself, which could have provided a more comprehensive insight into Stark's character.</p> <p>Comparison: Both responses do a good job of embodying Tony Stark's character, but they focus on different aspects. Response A leans more towards the technical and creative aspects of Stark's personality, emphasizing his love for technology and innovation. Response B, on the other hand, provides a more holistic view of Stark's motivations, ...</p>

Table 11: An example of pairwise response comparison on role-playing scenario from MT-Bench.

Query	
Discuss antitrust laws and their impact on market competition. Compare the antitrust laws in US and China along with some case studies."	
Response A (from Alpaca-13B)	Response B (Preferred) , from Vicuna-13B-v1.2)
Antitrust laws are designed to promote competition in markets and protect consumers from anticompetitive practices. In the US, the Sherman Antitrust Act of 1890 is the primary antitrust law and prohibits agreements or practices that restrain trade and impose unfair competition. The US Department of Justice and the Federal Trade Commission are responsible for enforcing the law. In China, the Anti-Monopoly Law enacted in 2008 is the primary antitrust law and is similar to the Sherman Act in prohibiting monopolistic agreements and abuse of dominance. Case studies have shown that antitrust laws can have a significant impact on market competition. For example, in 2000 the US Department of Justice filed a lawsuit against Microsoft alleging that the company had abused its monopoly in the computer operating system market. The lawsuit resulted in Microsoft being ordered to make significant changes to its business practices and the case is often cited as an example of how antitrust laws can be used to promote competition. In China, Overall, antitrust laws can have a significant impact on market competition and have been used to promote competition and protect consumers in both the US and China.	Antitrust laws are a set of regulations designed to promote and maintain fair competition in markets. These laws aim to prevent anticompetitive practices that can harm consumers, such as price-fixing, monopolies, and mergers that reduce competition. In this essay, we will discuss antitrust laws and their impact on market competition, compare the antitrust laws in the US and China, and analyze some case studies. Antitrust laws have a significant impact on market competition. These laws promote competition by preventing anticompetitive practices that can harm consumers. For example, antitrust laws prohibit price-fixing, which is when competitors agree to set prices at a certain level. This practice can lead to higher prices for consumers and reduced competition. Antitrust laws also prevent monopolies, which occur when a single company dominates a market. Monopolies can lead to higher prices, reduced innovation, and poorer quality products. Antitrust laws also regulate mergers and acquisitions to ensure that they do not reduce competition in the market. The US has a comprehensive set of antitrust laws, including the Sherman Act, the Clayton Act, and the Federal Trade Commission Act. The Sherman Act prohibits anticompetitive practices such as price-fixing, ...
Judgment by GPT-4 (X)	
Both assistants address the core aspects of antitrust laws, their impact on market competition, and provide comparisons between the US and China along with case studies. However, there are distinct differences in the presentation, depth, and accuracy of their responses. Assistant A Pros: 1. Clearly explains the purpose of antitrust laws and names specific laws in both the US and China. 2. Provides accurate and specific case studies: the 2000 Microsoft case in the US and the Qualcomm case in China in 2013. 3. Offers a concise conclusion summarizing the impact of antitrust laws. Assistant A Cons: ... Therefore, the final verdict is: [[A]]	
Criteria from Our DPO Model	Criteria from Our Initial Model
1. Depth of Analysis: The response should not only state facts but also analyze the implications of antitrust laws, including their effects on market competition, economic efficiency, and consumer welfare. 2. Practical Examples: The inclusion of real-world examples or case studies to illustrate how antitrust laws have been applied in the US and China, enhancing the understanding of the topic. ...	1. Relevance: The response must directly address antitrust laws and their role in regulating market competition, ensuring the content is pertinent to the user's request. 2. Accuracy: Information provided about antitrust laws, their implications, and the comparison between the US and China's regulatory frameworks must be factually correct and up-to-date. ...
Judgement by Our DPO Model (Criteria 1) (✓)	Judgement by Our Initial Model (Criteria 1) (X)
Response A Score: 3 Analysis of Response A: Response A provides a basic overview ... However, the response lacks depth in its analysis and does not provide a detailed comparison between the antitrust laws of the two countries beyond a general description. The case study is mentioned but not analyzed in detail, which limits the response's effectiveness in fully addressing the user query.	Response A Score: 4 Analysis of Response A: Analysis of Response A: Response A provides a clear and concise overview ... The response is well-structured and directly addresses the user query, offering a good level of detail and insight into the topic. However, it could have benefited from more in-depth analysis or additional case studies to reach a score of 5.
Response B Score: 4 Analysis of Response B: Response B offers a more detailed analysis of antitrust laws, including a discussion on the impact of these laws on market competition and a comparison between the US and Chinese antitrust frameworks. It provides specific examples of laws and regulations, such as the Sherman Act, the Clayton Act, and the Federal Trade Commission Act in the US, ... Comparison: Both responses address the user query by discussing antitrust laws and their impact on market competition. Response B is stronger because ...	Response B Score: 3 Analysis of Response B: Response B offers a general discussion on antitrust laws, their impact on market competition, and a comparison between the US and China. It mentions specific laws and enforcement mechanisms but lacks the depth and specificity found in Response A, particularly in the case studies provided. ... Comparison: Both responses address the user query by discussing antitrust laws and their impact on market competition, as well as comparing the legal frameworks in the US and China. Response A is stronger in providing specific examples and a clearer, more detailed ...

Table 12: An example of pairwise response comparison on humanities scenario from MT-Bench.

499 **NeurIPS Paper Checklist**

500 **1. Claims**

501 Question: Do the main claims made in the abstract and introduction accurately reflect the
502 paper's contributions and scope?

503 Answer: [Yes]

504 Justification: The key experiment results are summarized in the end of the introduction.
505 principal contributions are clearly stated in the end of the abstract.

506 Guidelines:

- 507 • The answer NA means that the abstract and introduction do not include the claims
508 made in the paper.
- 509 • The abstract and/or introduction should clearly state the claims made, including the
510 contributions made in the paper and important assumptions and limitations. A No or
511 NA answer to this question will not be perceived well by the reviewers.
- 512 • The claims made should match theoretical and experimental results, and reflect how
513 much the results can be expected to generalize to other settings.
- 514 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
515 are not attained by the paper.

516 **2. Limitations**

517 Question: Does the paper discuss the limitations of the work performed by the authors?

518 Answer: [Yes]

519 Justification: The limitations of this study are discussed in Section 6.1.

520 Guidelines:

- 521 • The answer NA means that the paper has no limitation while the answer No means that
522 the paper has limitations, but those are not discussed in the paper.
- 523 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 524 • The paper should point out any strong assumptions and how robust the results are to
525 violations of these assumptions (e.g., independence assumptions, noiseless settings,
526 model well-specification, asymptotic approximations only holding locally). The authors
527 should reflect on how these assumptions might be violated in practice and what the
528 implications would be.
- 529 • The authors should reflect on the scope of the claims made, e.g., if the approach was
530 only tested on a few datasets or with a few runs. In general, empirical results often
531 depend on implicit assumptions, which should be articulated.
- 532 • The authors should reflect on the factors that influence the performance of the approach.
533 For example, a facial recognition algorithm may perform poorly when image resolution
534 is low or images are taken in low lighting. Or a speech-to-text system might not be
535 used reliably to provide closed captions for online lectures because it fails to handle
536 technical jargon.
- 537 • The authors should discuss the computational efficiency of the proposed algorithms
538 and how they scale with dataset size.
- 539 • If applicable, the authors should discuss possible limitations of their approach to
540 address problems of privacy and fairness.
- 541 • While the authors might fear that complete honesty about limitations might be used by
542 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
543 limitations that aren't acknowledged in the paper. The authors should use their best
544 judgment and recognize that individual actions in favor of transparency play an impor-
545 tant role in developing norms that preserve the integrity of the community. Reviewers
546 will be specifically instructed to not penalize honesty concerning limitations.

547 **3. Theory Assumptions and Proofs**

548 Question: For each theoretical result, does the paper provide the full set of assumptions and
549 a complete (and correct) proof?

550 Answer: [NA]

551 Justification: This study does not include theoretical results.

552 Guidelines:

- 553 • The answer NA means that the paper does not include theoretical results.
- 554 • All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- 555
- 556 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 557 • The proofs can either appear in the main paper or the supplemental material, but if
- 558 they appear in the supplemental material, the authors are encouraged to provide a short
- 559 proof sketch to provide intuition.
- 560 • Inversely, any informal proof provided in the core of the paper should be complemented
- 561 by formal proofs provided in appendix or supplemental material.
- 562 • Theorems and Lemmas that the proof relies upon should be properly referenced.

563 4. Experimental Result Reproducibility

564 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
565 perimental results of the paper to the extent that it affects the main claims and/or conclusions
566 of the paper (regardless of whether the code and data are provided or not)?

567 Answer: [Yes]

568 Justification: All implementation details are provided in Appendix A.2.

569 Guidelines:

- 570 • The answer NA means that the paper does not include experiments.
- 571 • If the paper includes experiments, a No answer to this question will not be perceived
- 572 well by the reviewers: Making the paper reproducible is important, regardless of
- 573 whether the code and data are provided or not.
- 574 • If the contribution is a dataset and/or model, the authors should describe the steps taken
- 575 to make their results reproducible or verifiable.
- 576 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 577 For example, if the contribution is a novel architecture, describing the architecture fully
- 578 might suffice, or if the contribution is a specific model and empirical evaluation, it may
- 579 be necessary to either make it possible for others to replicate the model with the same
- 580 dataset, or provide access to the model. In general, releasing code and data is often
- 581 one good way to accomplish this, but reproducibility can also be provided via detailed
- 582 instructions for how to replicate the results, access to a hosted model (e.g., in the case
- 583 of a large language model), releasing of a model checkpoint, or other means that are
- 584 appropriate to the research performed.
- 585 • While NeurIPS does not require releasing code, the conference does require all submis-
586 sions to provide some reasonable avenue for reproducibility, which may depend on the
587 nature of the contribution. For example
- 588 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
589 to reproduce that algorithm.
- 590 (b) If the contribution is primarily a new model architecture, the paper should describe
591 the architecture clearly and fully.
- 592 (c) If the contribution is a new model (e.g., a large language model), then there should
593 either be a way to access this model for reproducing the results or a way to reproduce
594 the model (e.g., with an open-source dataset or instructions for how to construct
595 the dataset).
- 596 (d) We recognize that reproducibility may be tricky in some cases, in which case
597 authors are welcome to describe the particular way they provide for reproducibility.
598 In the case of closed-source models, it may be that access to the model is limited in
599 some way (e.g., to registered users), but it should be possible for other researchers
600 to have some path to reproducing or verifying the results.

601 5. Open access to data and code

602 Question: Does the paper provide open access to the data and code, with sufficient instruc-
603 tions to faithfully reproduce the main experimental results, as described in supplemental
604 material?

605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655

Answer: [Yes]

Justification: We provide a detailed README to describe how to reproduce the main results of this study.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: All implementation details are provided in Appendix A.2.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: we run every method for 3 independent trials.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).

- 656 • It should be clear whether the error bar is the standard deviation or the standard error
657 of the mean.
- 658 • It is OK to report 1-sigma error bars, but one should state it. The authors should
659 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
660 of Normality of errors is not verified.
- 661 • For asymmetric distributions, the authors should be careful not to show in tables or
662 figures symmetric error bars that would yield results that are out of range (e.g. negative
663 error rates).
- 664 • If error bars are reported in tables or plots, The authors should explain in the text how
665 they were calculated and reference the corresponding figures or tables in the text.

666 8. Experiments Compute Resources

667 Question: For each experiment, does the paper provide sufficient information on the com-
668 puter resources (type of compute workers, memory, time of execution) needed to reproduce
669 the experiments?

670 Answer: [Yes]

671 Justification: We include the information of computational resources in Appendix A.2.

672 Guidelines:

- 673 • The answer NA means that the paper does not include experiments.
- 674 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
675 or cloud provider, including relevant memory and storage.
- 676 • The paper should provide the amount of compute required for each of the individual
677 experimental runs as well as estimate the total compute.
- 678 • The paper should disclose whether the full research project required more compute
679 than the experiments reported in the paper (e.g., preliminary or failed experiments that
680 didn't make it into the paper).

681 9. Code Of Ethics

682 Question: Does the research conducted in the paper conform, in every respect, with the
683 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

684 Answer: [Yes]

685 Justification: We reviewed the NeurIPS Code of Ethics.

686 Guidelines:

- 687 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 688 • If the authors answer No, they should explain the special circumstances that require a
689 deviation from the Code of Ethics.
- 690 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
691 eration due to laws or regulations in their jurisdiction).

692 10. Broader Impacts

693 Question: Does the paper discuss both potential positive societal impacts and negative
694 societal impacts of the work performed?

695 Answer: [Yes]

696 Justification: We discuss the impacts on Section 7.

697 Guidelines:

- 698 • The answer NA means that there is no societal impact of the work performed.
- 699 • If the authors answer NA or No, they should explain why their work has no societal
700 impact or why the paper does not address societal impact.
- 701 • Examples of negative societal impacts include potential malicious or unintended uses
702 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
703 (e.g., deployment of technologies that could make decisions that unfairly impact specific
704 groups), privacy considerations, and security considerations.

- 705 • The conference expects that many papers will be foundational research and not tied
706 to particular applications, let alone deployments. However, if there is a direct path to
707 any negative applications, the authors should point it out. For example, it is legitimate
708 to point out that an improvement in the quality of generative models could be used to
709 generate deepfakes for disinformation. On the other hand, it is not needed to point out
710 that a generic algorithm for optimizing neural networks could enable people to train
711 models that generate Deepfakes faster.
- 712 • The authors should consider possible harms that could arise when the technology is
713 being used as intended and functioning correctly, harms that could arise when the
714 technology is being used as intended but gives incorrect results, and harms following
715 from (intentional or unintentional) misuse of the technology.
- 716 • If there are negative societal impacts, the authors could also discuss possible mitigation
717 strategies (e.g., gated release of models, providing defenses in addition to attacks,
718 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
719 feedback over time, improving the efficiency and accessibility of ML).

720 11. Safeguards

721 Question: Does the paper describe safeguards that have been put in place for responsible
722 release of data or models that have a high risk for misuse (e.g., pretrained language models,
723 image generators, or scraped datasets)?

724 Answer: [No]

725 Justification: This study poses no such risks.

726 Guidelines:

- 727 • The answer NA means that the paper poses no such risks.
- 728 • Released models that have a high risk for misuse or dual-use should be released with
729 necessary safeguards to allow for controlled use of the model, for example by requiring
730 that users adhere to usage guidelines or restrictions to access the model or implementing
731 safety filters.
- 732 • Datasets that have been scraped from the Internet could pose safety risks. The authors
733 should describe how they avoided releasing unsafe images.
- 734 • We recognize that providing effective safeguards is challenging, and many papers do
735 not require this, but we encourage authors to take this into account and make a best
736 faith effort.

737 12. Licenses for existing assets

738 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
739 the paper, properly credited and are the license and terms of use explicitly mentioned and
740 properly respected?

741 Answer: [NA]

742 Justification: We properly cited all benchmarks used in this study.

743 Guidelines:

- 744 • The answer NA means that the paper does not use existing assets.
- 745 • The authors should cite the original paper that produced the code package or dataset.
- 746 • The authors should state which version of the asset is used and, if possible, include a
747 URL.
- 748 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 749 • For scraped data from a particular source (e.g., website), the copyright and terms of
750 service of that source should be provided.
- 751 • If assets are released, the license, copyright information, and terms of use in the
752 package should be provided. For popular datasets, paperswithcode.com/datasets
753 has curated licenses for some datasets. Their licensing guide can help determine the
754 license of a dataset.
- 755 • For existing datasets that are re-packaged, both the original license and the license of
756 the derived asset (if it has changed) should be provided.

757 • If this information is not available online, the authors are encouraged to reach out to
758 the asset’s creators.

759 **13. New Assets**

760 Question: Are new assets introduced in the paper well documented and is the documentation
761 provided alongside the assets?

762 Answer: [Yes]

763 Justification: We provide a new evaluation training dataset.

764 Guidelines:

- 765 • The answer NA means that the paper does not release new assets.
- 766 • Researchers should communicate the details of the dataset/code/model as part of their
767 submissions via structured templates. This includes details about training, license,
768 limitations, etc.
- 769 • The paper should discuss whether and how consent was obtained from people whose
770 asset is used.
- 771 • At submission time, remember to anonymize your assets (if applicable). You can either
772 create an anonymized URL or include an anonymized zip file.

773 **14. Crowdsourcing and Research with Human Subjects**

774 Question: For crowdsourcing experiments and research with human subjects, does the paper
775 include the full text of instructions given to participants and screenshots, if applicable, as
776 well as details about compensation (if any)?

777 Answer: [NA]

778 Justification: This study does not involve crowdsourcing nor research with human subjects.

779 Guidelines:

- 780 • The answer NA means that the paper does not involve crowdsourcing nor research with
781 human subjects.
- 782 • Including this information in the supplemental material is fine, but if the main contribu-
783 tion of the paper involves human subjects, then as much detail as possible should be
784 included in the main paper.
- 785 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
786 or other labor should be paid at least the minimum wage in the country of the data
787 collector.

788 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human
789 Subjects**

790 Question: Does the paper describe potential risks incurred by study participants, whether
791 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
792 approvals (or an equivalent approval/review based on the requirements of your country or
793 institution) were obtained?

794 Answer: [NA]

795 Justification: This study does not involve crowdsourcing nor research with human subjects.

796 Guidelines:

- 797 • The answer NA means that the paper does not involve crowdsourcing nor research with
798 human subjects.
- 799 • Depending on the country in which research is conducted, IRB approval (or equivalent)
800 may be required for any human subjects research. If you obtained IRB approval, you
801 should clearly state this in the paper.
- 802 • We recognize that the procedures for this may vary significantly between institutions
803 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
804 guidelines for their institution.
- 805 • For initial submissions, do not include any information that would break anonymity (if
806 applicable), such as the institution conducting the review.