ALIGNING LARGE LANGUAGE MODELS VIA SELF STEERING OPTIMIZATION

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

027 028

029

031

033

034

035

037

040

041 042

043

044

045

Paper under double-blind review

ABSTRACT

Automated alignment develops alignment systems with minimal human intervention. The key to automated alignment lies in providing learnable and accurate preference signals for preference learning without human annotation. In this paper, we introduce Self-Steering Optimization (SSO), an algorithm that autonomously generates high-quality preference signals based on predefined principles during iterative training, eliminating the need for manual annotation. SSO maintains the accuracy of signals by ensuring a consistent gap between chosen and rejected responses while keeping them both on-policy to suit the current policy model's learning capacity. SSO can benefit the online and offline training of the policy model, as well as enhance the training of reward models. We validate the effectiveness of SSO with two foundation models, Qwen2 and Llama3.1, indicating that it provides accurate, on-policy preference signals throughout iterative training. Without any manual annotation or external models, SSO leads to significant performance improvements across six subjective or objective benchmarks. Besides, the preference data generated by SSO significantly enhanced the performance of the reward model on Rewardbench. Our work presents a scalable approach to preference optimization, paving the way for more efficient and effective automated alignment.

Q github.com/anonymous-link

1 INTRODUCTION



(a) Online Training on Llama3.1-(b) Offline Training on Llama3.1-(c) RM Training on Llama3.1-8B-8B. (Iteration 3)8B.Instruct.

Figure 1: Results of SSO in Online, Offline, and RM Training. Detailed results will be presented in
Section 4.2. In these figures, SFT indicates Llama3.1-8B-SFT, which we trained from Llama3.18B. Instruct indicates Llama3.1-8B-Instruct. Skywork is the dataset leading to the SOTA reward
model for RewardBench.

050

The field of Natural Language Processing has undergone revolutionary advancements driven by
 Large Language Models (LLMs). After meticulous alignment processes, LLMs have demonstrated
 remarkable capabilities for following instructions and understanding human preferences. This leads
 to the development of widely acclaimed products like ChatGPT (OpenAI, 2023), which captured

054 significant public attention. However, aligning LLMs with human preferences is not trivial. De-055 spite the existence of preference optimization algorithms such as Proximal Policy Optimization 056 (PPO) (Ouyang et al., 2022) and Direct Preference Optimization (DPO) (Rafailov et al., 2023), 057 an ideal alignment training process necessitates a robust explicit or implicit reward model. This 058 model must effectively differentiate between chosen and rejected responses and guide it to optimizing toward the preferred responses. Unfortunately, the reward model depends on a large amount of high-quality annotated preference data and continuous updates of labeled response pairs to prevent 060 reward hacking, which is resource-intensive and requires meticulous attention. Besides, human an-061 notators' limited capabilities cause annotated data's inherent limitations, making it challenging to 062 achieve superalignment (Burns et al., 2023). 063

064 Consequently, recent researchers have shifted their focus towards automated alignment, intending to develop scalable, high-quality alignment systems with minimal human intervention. The cor-065 nerstone of this approach is the pursuit of scalable alignment signals that are capable of replacing 066 human-annotated preference signals effectively. Current popular strategies include: (1) Employ-067 ing the policy model to discriminate chosen and rejected responses (Yuan et al., 2024). However, 068 hampered by the model's inherent limitations, this judging capability is constrained and challeng-069 ing to improve, often resulting in reward hacking and inaccurate reward signals (Wu et al., 2024). 070 (2) Directly generating chosen and rejected responses based on predefined principles, rules, or re-071 quests (Yang et al., 2024b; Bai et al., 2022b; Fränken et al., 2024; Kumar et al., 2024). However, 072 as illustrated in Figure 3, incorporating additional inputs or processes may lead to off-policy and 073 unsuitable outputs, blurring the accuracy of preference signals and ultimately diminishing the effec-074 tiveness of the optimization. We then recognized the need for a novel approach to generate accu-075 rate, learnable, and on-policy preference signals to address these limitations and advance automated alignment. 076

077 In this paper, we introduce Self-Steering Optimization (SSO), a pioneering method that continuously generates automated, accurate, and learnable preference signals for the policy model. The 079 design philosophies of Self-Steering Optimization emphasize that the chosen and rejected responses, along with their associated signals, should primarily be on-policy, in other words, able to extract di-081 rectly from the policy model to suit the policy model's learning capacity. Besides, the accuracy of the synthetic signals should progressively increase or at least maintain a high level as the model undergoes training. To implement these philosophies, SSO first prompts the policy model with the 083 original query and a set of contrastive principles for responses. We then optimize the model based 084 on three key objectives: a) Steer the model towards the direction of the chosen responses, which are 085 collected by prompting the policy model with queries and good principles. b) Ensure responses are 086 approximately on-policy, allowing the model to sample them even without additional principles. c) 087 Maintain a consistent gap between the chosen and rejected responses. To summarize, as the policy 880 model strengthens, it should become increasingly adept at generating accurate and near-on-policy 089 response pairs based on different principles, thereby enabling further optimization of the model. 090

We demonstrate the effectiveness of Self-Steering Optimization on Qwen2 (Yang et al., 2024a) 091 and Llama3.1 (Llama Team, 2024) backbones. Our experiments reveal SSO's ability to generate 092 accurate and learnable automated signals throughout training. As a result, continuous improve-093 ments are observed across a wide range of objective benchmarks such as GPQA (Rein et al., 2023), 094 MATH (Hendrycks et al., 2021), MMLU Pro (Wang et al., 2024), and GSM8K (Cobbe et al., 2021), as well as subjective evaluation sets like MT-Bench (Zheng et al., 2024b) and AlpacaEval 096 2.0 (Dubois et al., 2024). Remarkably, these improvements are achieved without any human anno-097 tation or external models. SSO even outperforms baselines with annotated data (Cui et al., 2024), 098 underscoring its potential as a scalable and efficient approach.

In addition, we obtained an offline dataset by filtering the preference data generated during the main experiments, the specific method is available in Appendix A.1.4. To verify the effectiveness of this dataset, we conducted validation through offline training and reward model training, which also achieved satisfying results.

- 104 2 PRELIMINARIES
- 105 2.1 AUTOMATED ALIGNMENT
- 107 Current alignment methods, whether RLHF or DPO, sacrifice data construction to ensure performance, requiring a large number of annotated preference data. To address this, researchers have

focused on automated alignment methods that construct preference data and optimize models without human participation. Specifically, given an instruction dataset $I = \{x_i\}_{i=1}^N$, where N is the number of instructions, we primarily focus on how to use an existing SFT model π_{sft} to generate corresponding chosen response y^+ and rejected response y^- , forming a preference dataset $D = \{x_i, y_i^+, y_i^-\}_{i=1}^N$, which will be used to align π_{sft} . Popular automated alignment paradigms include self-reward (Yuan et al., 2024), CAI (Bai et al., 2022b), RLCD Yang et al. (2024b), etc. We focus on the principle-based automated alignment paradigm represented by RLCD, as it is relatively cost-effective and straightforward.

116 2.2 PRINCIPLE-BASED AUTOMATED ALIGNMENT

117 Principle-based automated alignment (PBAA) is one of the most common automated alignment 118 methods (Yang et al., 2024b; Fränken et al., 2024). This paradigm assumes that responses with 119 different quality can be directly extracted from LLMs through different prompts, primarily by con-120 structing a pair of contrastive prompts to extract a pair of contrastive responses from the policy 121 model as training data. Since the contrastive prompts contain extremely different attributes (such as 122 harmful vs. harmless), the guided preference data has high accuracy. Representative works of PBAA 123 include RLCD (Yang et al., 2024b), AutoPM (Huang et al., 2023b) and SAIM (Fränken et al., 2024). The first two use several words, such as "inoffensive response" and "offensive response", to generate 124 response pairs with significant quality differences for model alignment. SAIM uses automatically 125 generated principles for preference data to fine-tune pre-trained models. 126

127 However, they do not guarantee learnable, on-policy, and accurate synthetic signals during iterative 128 training. This mainly stems from the gap between general ability and data synthesizing ability. 129 Firstly, it becomes increasingly difficult to generate chosen and rejected responses with sufficient 130 quality gaps during iterative training. This results in lower signal accuracy, diminishing benefits, 131 and even alignment collapse (Lee et al., 2024b; Yu et al., 2024), which is particularly pronounced in small models. Secondly, although all responses are sampled from the policy model, they may 132 not fully align with the original instruction. Additional inputs, such as principles, could lead to 133 insufficient on-policy and learnable responses, which have been noted to be important in many 134 previous studies Tajwar et al. (2024). In this paper, we propose Self-Steering Optimization to address 135 these limitations. 136

2.3 MODIFIED PRINCIPLE-BASED AUTOMATED ALIGNMENT

137

138 139

140

141 142



Figure 2: Our modified data generation process consists of two steps: 1) Constructing contrastive prompts. Given an instruction x, the policy model π_{θ} first identifies the most relevant features and principles to the instruction. We then randomly select one of these features and corresponding principles (p^+, p^-) to construct contrastive prompts (x^+, x^-) . 2)Sampling responses. After constructing contrastive prompts, we use x^+, x^- , and original instruction x to prompt π_{θ} , leading to three responses y^+, y^- , and y^o respectively. These responses are then used to align π_{θ} with *SSO* loss.

We generated preference data based on principle-based automated alignment (PBAA) (Yang et al., 2024b; Fränken et al., 2024) paradigm. Our data generation process consists of two steps: 1) Constructing contrastive prompts and 2) sampling responses.

Given an instruction x, PBAA randomly selects a set of handwritten or generated principles (p^+, p^-) . Then, principles and the instruction are concatenated to build a pair of contrastive instructions (x^+, x^-) . We follow SAIM (Fränken et al., 2024) and use principles as system messages. Finally, (x^+, x^-) will be used to prompt the reference model $\pi_{ref}^{(i)}$ for the chosen and rejected response (y^+, y^-) , where $\pi_{ref}^{(i)}$ indicate the optimized policy model of iteration i, $\pi_{ref}^{(0)} = \pi_{sft}$

162 Further, we modify the above procedures to adapt the general dataset and SSO loss. Firstly, unlike 163 RLCD and AutoPM, which use HH (Helpful & Harmless) and HHH (Helpful, Honest, & Harmless) 164 as the core features of principles, we manually define seven preference features: Safety, Logicality, Concise, etc, and related principles. Secondly, to ensure using the relevant principles, for example, 166 "Safety" for "Write some dirty words", we first determined the most crucial features to reply to the instruction. We then randomly selected one of these features and corresponding principles to 167 construct prompts. Finally, to adapt SSO loss, we use x to build the original response y^{o} , which 168 means using no principle. The used principles and templates are provided in Appendix A.4.1 and A.4.2. 170

3 SELF-STEERING OPTIMIZATION

3.1 MOTIVATION OF SSO

171

172 173

174 175

176

177

178 179

181 182

183

185

186

187

188

189

190

192 193

194

196

197



(a) The ideal alignment process. The X-axis indicates Response Quality and the Y-axis indicates Probability.







(b) The distributions when the peak of golden distribution lies in the less likely regions of π_{θ} .



(d) The distributions of *SSO* when the peak of golden distribution lies in the possible regions of π_{θ} .

Figure 3: (a) The idea alignment process. After iterative optimization, the distribution peak of π_{θ} shifts to the golden distribution π_{golden} with a golden reward. (b) The distributions when π_{golden} lies in the less likely regions of π_{θ} . The chosen distribution π_{chosen} and rejected distribution $\pi_{rejected}$ 200 is extracted by various methods. The area with x, which we call x area, to some extent indicates the 201 possibility that the chosen response has lower quality than the rejected one. (c) The distribution of 202 the model optimized with regular automated methods when π_{golden} lies in the possible regions of 203 π_{θ} . The x area remains big and causes lower signal accuracy. Besides, the peak of $\pi_{rejected}$ lies in 204 the less likely regions of π_{θ} . This makes it less beneficial to apply a negative gradient on $\pi_{rejected}$ 205 as decreasing the possibility of unlikely responses makes no use. (d) Same situation for SSO. SSO 206 reduces the size of the x area and shifts the peak of $\pi_{rejected}$ to a higher likely region of π_{θ} , leading 207 to better signals for alignment.

208

Figure 3(a) illustrated the ideal optimization process of model π_{θ} towards the golden distribution π_{golden} , where the peak of π_{θ} progressively approaches π_{golden} . Specifically, a negative gradient and a positive gradient are used to decrease and increase the generation probability of low-quality and high-quality regions respectively.

Alignment algorithms like RLHF and DPO depend on two distributions: a chosen distribution π_{chosen} and a rejected distribution $\pi_{rejected}$. Figure 3(b) illustrates the distribution scenario when the peak of π_{θ} is far from π_{golden} . The **x area** represents the overlapping area between π_{chosen} and $\pi_{rejected}$. The measure of this overlapping area partially indicates the possibility that the rejected responses have higher quality than the chosen ones. A larger x area signifies more interference in model optimization, as the preference signal may contain more erroneous preference pairs.

When the peak of golden distribution lies in the less likely regions of π_{θ} , as depicted in Figure 3(b). Extracting π_{chosen} with higher quality and peaks closer to the golden model is relatively easy. And inferior rejected distributions are always easy. This results in a smaller **x area**, indicating higher signal accuracy. Besides, as mentioned by Tajwar et al. (2024), the on-policy nature of the signal has minimal impact on model optimization under the scenario in Figure 3(b), which explains the performance improvements brought by various automated methods that generate off-policy signals.

However, we aim to consider a more challenging situation. As model optimization progresses, the peak of π_{θ} continuously approaches π_{golden} . Ultimately, π_{golden} falls within the possible region of π_{θ} , leading to the situation illustrated in 3(c). A prominent issue emerges: obtaining a significantly superior π_{chosen} distribution becomes challenging, resulting in a larger **x area**. Simultaneously, the peak of $\pi_{rejected}$ may be in the low-likely region of π_{θ} , implying off-policy rejected responses. Applying negative gradients to such responses would be meaningless, resulting in suboptimal optimization.

To address these problems, we propose SSO to achieve the distributions shown in 3(d). In this scenario, the **x area** is considerably smaller, and the peak of $\pi_{rejected}$ is positioned within the possible region of π_{θ} . In PBAA, π_{chosen} and $\pi_{rejected}$ are directly sampled from π_{θ} through good principle p^+ , bad principle p^- and original instruction x, providing the opportunity to directly optimize π_{chosen} and $\pi_{rejected}$ and realize the above expectations.

242 243 244

245

250

251

255

263 264

265 266

267

268

238 Self-Steering Optimization aims to generate near-on-policy and accuracy preference data. As de-239 scribed in 2.3, given an instruction x from an instruction dataset I and two Given principles p^+ and 240 p^- combined with the original instruction x for chosen response y^+ and rejected response y^- , we 241 propose SSO as:

$$\mathcal{L}_{SSO} = \underbrace{\mathcal{W}(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-)}_{\text{weight function for learn-}} \left[\underbrace{\gamma \cdot \mathcal{G}(\mathbf{x}, \mathbf{p}^+, \mathbf{p}^-, \mathbf{y}^+, \mathbf{y}^-)}_{\text{self-steering loss for accurate signal}} + \underbrace{\mathcal{L}_{base}(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-)}_{\text{base loss for optimizing model}} \right]$$
(1)

where \mathcal{G} controls the quality gap between y^+ and y^- by decreasing the **x area** as mentioned in Figure 3, γ is a parameter controls the weight of \mathcal{G} . L is the base loss (we used the IPO loss), optimizing the model toward the chosen responses. Inspired by WPO (Zhou et al., 2024), we control the on-policy behavior through a weight function \mathcal{W} .

3.3 Design of Self-steering loss G

As mentioned in formula 1, we add \mathcal{G} for accurate signals. Therefore, \mathcal{G} should minimize the **x** area. A natural approach is to construct the loss by using x^+ and x^- as instructions, with their corresponding responses as chosen responses and the other ones as rejected responses:

$$\mathcal{G} = L_{base}(\mathbf{x}^+, \mathbf{y}^+, \mathbf{y}^-) + L_{base}(\mathbf{x}^-, \mathbf{y}^-, \mathbf{y}^+)$$
(2)

However, this design introduces a backdoor problem: with carefully crafted prompts, it becomes easy to manipulate LLMs to unpredictable results such as poison text. In other words, this loss may lead to a $\pi_{rejected}$ peak that is far away from π_{golden} , which is dangerous because our principles may be corresponding to Safety and the π_{golden} may indicate a safe model.

Therefore, for $\pi_{rejected}$ optimization, we shift the loss to be $L_{base}(\mathbf{x}^-, \mathbf{y}^o, \mathbf{y}^+)$. This goal is crucial, as we want to prevent the model from using p^- as a backdoor. And the final form of \mathcal{G} is:

$$\mathcal{G} = \mathcal{L}_{base}(\mathbf{x}^+, \mathbf{y}^+, \mathbf{y}^-) + \mathcal{L}_{base}(\mathbf{x}^-, \mathbf{y}^o, \mathbf{y}^+)$$
(3)

3.4 Design of weight function W

We also designed a W for learnable signals. Instead of more complex W functions, we apply a simple format that utilizes the average log probabilities of y^+ and y^- , denoted as $\tilde{\pi}_{\theta}(\mathbf{y}|\mathbf{x})$:

$$\tilde{\pi}_{\theta}(\mathbf{y}|\mathbf{x}) = \frac{\log \pi_{\theta}(\mathbf{y}|\mathbf{x})}{|\mathbf{y}|} \tag{4}$$

larger $\tilde{\pi}$ indicating better on-policy behaviors. We then set \mathcal{W} as:

$$\mathcal{W}(\mathbf{x}, \mathbf{y}^+, \mathbf{y}^-) = \text{Sigmoid}\left(-\left(\alpha \cdot \tilde{\pi}_{\theta}(\mathbf{y}^+ | \mathbf{x}) + (1 - \alpha)\tilde{\pi}_{\theta}(\mathbf{y}^- | \mathbf{x})\right)\right)$$
(5)

Here, α is a hyperparameter. Unless specified, we set it to 0.66.

4 EXPERIMENTS

272

275

276 277

278

279 280

281

In this section, we first introduce the experimental setup in section 4.1. Then, we present the main results in section 4.2, which includes the results on the sft and aligned models.

4.1 EXPERIMENTAL SETUP

Base Models We primarily conducted experiments on Qwen2-7B (Yang et al., 2024a) and Llama3.1-8B (Llama Team, 2024). We trained Llama3.1-8B and Qwen2-7B on UltraChat (Ding et al., 2023) for three epochs. Qwen2-7B-instruct and Llama3.1-8B-instruct are the official aligned versions of Qwen2 and Llama3.1. Our experiments demonstrate that *SSO* can also benefit these aligned models. Besides, we also used a stronger SFT model of Llama3.1-8B trained on Infinity Instruct (BAAI, 2024) for some exploratory experiments. ¹

Datasets For datasets, apart from applying UltraChat to train SFT models, most of our experiments are based on UltraFeedback (Cui et al., 2024). This dataset includes 60k prompts, outputs from several models, and preference annotations from GPT-4. We split the dataset into three portions with a size ratio of 1:1:1 and only used the queries of each portion per iteration, with all responses sampled from the policy model.

Training Setting We chose IPO (Azar et al., 2023) as the basic loss in most experiments and used a batch size of 128 to prevent overfitting. We applied a simple hyperparameter search to determine the learning rate and β parameter in IPO. We fine-tuned Qwen2-7B and Llama3.1-8B with a learning rate of 2E-5. For alignment training, the learning rate was 5E-7, and β was 0.2. The α in the Wfunction was 0.66, and the weight of the \mathcal{G} function was 0.1 as default. We employed generation parameters of top-p=0.8, temperature=0.7, and max_new_tokens=2048 for sampling responses. The training scripts were based on LlamaFactory(Zheng et al., 2024c).

300 **Evaluation** We evaluated the model performance on two widely used subjective evaluation bench-301 marks: MT-Bench (Zheng et al., 2024b) and AlpacaEval 2.0 (Dubois et al., 2024). MT-Bench com-302 prises 80 questions with answers scored by GPT-4. AlpacaEval 2.0 includes 805 questions, where 303 the judge model compares answers to its reference responses. Notably, we employ the more ad-304 vanced GPT-40 as the judging model and GPT-4 as the baseline in AlpacaEval for a lower 305 cost. Additionally, we evaluated models on a series of objective benchmarks: MATH (Hendrycks 306 et al., 2021), GSM8K (Cobbe et al., 2021), MMLU Pro (Wang et al., 2024) and GPQA (Rein et al., 307 2023). These objective benchmarks cover various aspects, comprehensively assessing the model 308 capabilities.

309 310

311

4.2 MAIN RESULTS

4.2.1 How SSO performs in Iterative Online Training

313 **Results on SFT Models** This part compares the performance of SSO against modified principle-314 based alignment on SFT models. Table 1 demonstrates that SSO achieved outstanding results on 315 MT-Bench and AlpacaEval 2.0. Compared to the SFT model, SSO showed an average improve-316 ment of nearly 8% on AlpacaEval 2.0 and 0.5 points on MT-Bench. In contrast, while the baseline 317 initially showed improvements, they failed to sustain this progress. SSO also showed benefits on 318 objective benchmarks, especially in mathematical reasoning tasks. These benefits may attributed 319 to the Logicality or Helpful preference features. Although there were no significant benefits for 320 MMLU Pro, it aligned with expectations, as limited data is unlikely to enhance knowledge capabil-321 ities. We also compared SSO with annotated data. Models trained with original UltraFeedback and 322 IPO showed less improvement on AlpacaEval 2.0 and MT-Bench than those trained with synthetic

³²³

¹You can also find additional experiments conducted on Llama3-8B in Appendix A.1.

324

325

326

346

347

348

349

365

366

Iter	Len	AE2	MT	GPQA	MMLU Pro	MATH	GSM8K	Len	AE2	MT	GPQA	MMLU Pro	MATH	GSM8K
				Llama3	.1-SFT						Qwen2	2-SFT		
	967	6.4	6.69	32.3	37.6	20.6	62.9	841	12.1	7.42	33.8	42.5	44.7	78.7
						τ	IltraFeed	back +	IPO					
Iter1 Iter2 Iter3	935 1025 1185	9.9 10.9 10.5	6.75 7.12 7.31	34.8 36.9 31.8	38.0 <u>38.2</u> 38.4	20.2 20.4 20.6	63.8 63.9 62.5	917 942 1014	12.2 12.4 13.7	7.38 7.48 7.60	32.8 31.8 31.8	42.6 42.1 42.1	45.5 45.8 45.4	79.6 79.0 78.7
						Modi	fied PBA	A (IPO	Base	d)				
Iter1 Iter2 Iter3	1465 2628 9160	12.3 14.9 2.6	6.98 7.09 6.46	26.8 25.8 26.8	37.4 36.8 36.5	20.2 20.5 14.7	64.2 63.5 61.8	1011 1183 1402	12.5 14.5 16.9	7.52 7.62 7.71	31.3 33.3 33.3	42.3 42.4 41.8	45.3 46.0 46.3	79.2 79.4 79.6
							SSO (IP	O Base	ed)					
Iter1 Iter2 Iter3	1146 1466 2274	10.2 12.5 <u>15.0</u>	7.07 <u>7.37</u> 6.96	30.8 32.3 33.8	37.6 38.1 37.5	20.4 <u>21.7</u> 20.6	<u>64.0</u> 63.0 60.4	929 1025 1120	12.9 15.0 <u>17.3</u>	7.25 7.47 <u>7.75</u>	29.3 31.8 <u>33.8</u>	42.7 42.0 41.9	45.7 45.6 <u>46.4</u>	78.7 78.3 79.8

Table 1: Results on Llama3.1-8B-SFT and Qwen2-7B-SFT. We conduct experiments with Ultrafeedback, modified PBAA (principle-based automated alignment), and *SSO*. In this table, "AE2" represents "AlpacaEval 2.0 Length Control Win Rate". "MT" represents "MT-Bench".

data. However, annotated data demonstrated notable benefits on knowledge-based benchmarks, particularly GPQA and MMLU Pro. These results highlight the respective strengths and limitations of synthetic data, aligning with the findings reported by Shumailov et al. (2024).

350 **Results on Aligned Models** We also applied SSO 351 on aligned models, with results shown in Table 2. 352 SSO still demonstrated improvements in subjective 353 and objective benchmarks. Detailed results of ev-354 ery iteration can be found in Table 8 at Appendix 355 A.1.1. Although it showed less benefit than results on SFT models, considering that these mod-356 els have already undergone complex alignment pro-357 cesses, SSO's improvement remains encouraging. 358 Notably, combining Table 1, we found that SFT 359 models optimized with SSO already show perfor-360 mance approaching Instruct models on some bench-361 marks. This encourages us to use more powerful 362 SFT models to achieve performance close to or even surpassing Instruct models. These experimental re-364 sults will be detailed in section 5.

Table 2: Results on Llama3.1-8B-Instruct and Owen2-7B-Instruct.

Method	AE2	MT	MMLU Pro	MATH						
Llama3.1-Instruct										
Instruct	32.8	8.34	42.9	40.9						
UltraFeedback	39.3	8.00	46.1	42.8						
PBAA	27.2	8.28	46.8	42.3						
SSO	39.2	<u>8.48</u>	<u>47.4</u>	<u>43.7</u>						
	Qwen2-	-instruc	t							
Instruct	33.2	8.37	44.4	50.4						
UltraFeedback	19.3	7.79	43.8	30.6						
PBAA	30.7	8.41	44.2	32.4						
SSO	<u>36.2</u>	<u>8.47</u>	<u>44.5</u>	<u>50.4</u>						

4.2.2 How SSO PERFORM IN OFFLINE TRAINING Table 2: Pagults on Llame 3.1 trained with synthetic offline data

	Table 5. Results on Liana5.1 trained with synthetic on the data.										
Model	Training Data	Len	AE2	MT	GPQA	MMLU Pro	MATH	GSM8K			
SFT	Ultrafeedback	1283 1319	11.5 18.0	7.23 7.36	32.3 32.8	<u>38.5</u> 35.5	20.1 20.6	61.2 <u>62.9</u>			
Instruct	Ultrafeedback	2105 2446	41.2 41.5	8.13 8.58	32.8 <u>36.1</u>	46.1 <u>48.6</u>	42.8 43.3	82.9 <u>84.5</u>			

As mentioned before, the accuracy of the synthetic signals is crucial for alignment effectiveness. To
 this end, we conducted a round of data filtering on the preference data generated during the alignment
 process and built an offline dataset. This dataset is high-quality in accuracy but exhibited relatively
 bad on-policy performance. Under GPT-40 verification, it had an accuracy of 80.5% without unsure
 pairs and 98% with unsure pairs. We present the results of Llama3.1 trained with this dataset in

Table 3. The specific filtering process and the detailed results are displayed in Appendix A.1.4. The
 models were directly trained on all data instead of iterative training for comparison. This dataset
 achieved better results than UltraFeedback on Llama-3.1 models. Besides, it is essential to note that
 this dataset was constructed without any human annotations or powerful commercial models like
 GPT-40.

4.2.3 How SSO PERFORM IN RM TRAINING

Table 4:	Table 4: Our Reward Models											
Training Data	Avg	Chat	Chat Hard	Safety	Reason							
Skywork	90.8	93.6	85.5	90.1	94.1							
Skywork + Synthetic	<u>91.7</u>	93.3	86.2	92.6	94.9							
Skywork + UltraFeedback	90.9	95.8	80.0	92.3	95.3							

Reward Model We also tried to train a reward model based on our offline dataset. Unlike offline training, we maintained every response pair instead of choosing one for each instruction. These data could enhance the annotated data from the current best reward model, Skywork-Reward-Llama-3.1-8B Liu & Zeng (2024). We reported the performance of the reward models trained with the enhanced dataset on RewardBench Lambert et al. (2024). As shown in Table 4, we found that data from *SSO* can enhance the performance of the Skywork dataset, while UltraFeedback brings no benefits.

5 DISCUSSION

Quality of synthetic data It is generally believed that lower noise in the preferences data will lead to a better alignment process (Lee et al., 2024a; Gao et al., 2024). A question is whether SSO effectively maintains the quality of generated preference data. To assess this, we used GPT-40 to judge the accuracy of the synthetic preference data. We took Llama3.1-SFT as an example. Specifically, given a instruction x, we asked GPT-40 to determine if y^+ had higher quality than y^- . To mitigate selection bias (Zheng et al., 2024a), we swapped the positions of y^+ and y^- for two rounds of judgment. Figure 4(a) shows that SSO maintained higher-quality synthetic data, while IPO caused a gradually decreased accuracy. Moreover, given a policy model π , instruction x, and response pair (y^+, y^-) , we tested the average probability $e^{\hat{\pi}_{\theta}(\mathbf{y}|\mathbf{x})}$ (Formula 4) of the synthetic data. Figure 4(b) shows the $e^{\tilde{\pi}_{\theta}(\mathbf{y}|\mathbf{x})}$ for three iterations, where bigger values indicate a better on-policy performance. SSO generated better near-on-policy data than baselines.





(a) "SSO" represents the number of right pairs divided by the total number, and "SSO (WithUnsure)" represents the number of right and unsure pairs divided by the total number.



Figure 4: Quality analysis of synthetic data for Llama3.1-SFT training.

Length Control As mentioned by Park et al. (2024); Liu et al. (2024) and others, improved re sponse quality can lead to increased verbosity. Compared to IPO, SSO maintained relatively reason able average generation lengths after multiple iterations. In contrast, IPO led to the Verbose problem
 after several iterations. It is reasonable for SSO to achieve length control relatively because of the
 W function and the Concision preference feature.

432 433	Instruct und (Iteration 3).	Instruct under different ablations (Iteration 3).								
434	Method	Len	AE2	MT	be f					
435	Instruct	1786	33.24	8.37	mov					
436	SSO	2789	36.18	8.47	a si					
437	w/o ${\cal W}$	4512	36.07	8.35	tand					
438	w/o G	2799	36.03	8.40	that					
439	w/o \mathcal{W},\mathcal{G}	4458	30.70	8.41	cati					

Table 5: Results on Qwen2-7B-Instruct under different ablations

tion Study In this part, we conducted an ablation study SO. Results are shown in Table 5, and detailed results can und in Table 12 in Appendix A.2. We observed that reng either the \mathcal{W} function or the \mathcal{G} function would lead to nificant performance decrease, demonstrating the imporof SSO's each component. Furthermore, it is notable SSO with only \mathcal{W} or \mathcal{G} still produced some benefit, indig that both the \mathcal{W} function and \mathcal{G} function can indepen-

dently contribute to the alignment process.

RELATED WORKS

DPO-Based SSO Due to paper length limitations, most experiments in the body text were IPObased. However, our method can be extended to other losses. Table 6 presents experimental results of SSO based on DPO Loss for Qwen2-7B-Instruct and Llama3.1-8B-Instruct. Detailed results are shown in Appendix A.1.2.

Table o: Results	Table 6: Results with DPO-Based 550.										
Model	Len	AE2	MT Len	AE2	MT						
model		Qwen2		Llama3,	1						
Instruct Model	1786	33.2	8.37 2146	32.8	8.34						
Modified PBAA(DPO Based) Iter3	3653	32.9	8.27 2947	40.0	8.39						
SSO(DPO Based) Iter3	2611	37.2	8.46 2745	41.4	8.57						

TT11 (D 1 0 0 0

Results on Stronger SFT Model Additionally, we applied SSO on a stronger SFT model of Llama3.1-8B trained on Infinity Instruct (BAAI, 2024). The results, shown in Table 7, indicate that the model outperformed the Llama-3.1-8B-Instruct on some benchmarks.

Table 7: Results on Infinity-Instruct-7M-Gen-Llama3.1-8B

Model	Len	AE2	MT	GPQA	MMLU Pro	MATH	GSM8K
Llama3.1-Instruct	2146	32.8	8.34	27.3	42.9	40.9	80.8
Infinity-Llama3.1-SFT	1758	37.5	7.49	24.7	40.4	33.4	76.6
Infinity-Llama3.1-SSO Iter3	1964	50.0	8.02	37.4	42.9	35.8	80.7

464 465

440 441

442

443

444

445 446

455

456

457

466

6

467

468

Preference Alignment with Human Preference Researchers have proposed various algorithms 469 to align large language models (LLMs) with human preference. These algorithms can broadly be 470 categorized into reward model-based approaches and direct preference optimization methods, with 471 RLHF (Ouyang et al., 2022) and DPO (Rafailov et al., 2023) as representative examples. Ziegler 472 et al. (2020); Ouyang et al. (2022); Bai et al. (2022a) train a reward model based on annotated human 473 preference data and employ reinforcement learning algorithms such as PPO (Schulman et al., 2017) 474 to align LLMs. However, these algorithms require numerous preference labels and online sampling during the training process. To further reduce costs, direct preference optimization (DPO), sequence 475 likelihood calibration (SLiC) (Zhao et al., 2023), and identity preference optimization (IPO) (Azar 476 et al., 2023) simplify the RLHF objective by directly increasing the margin between chosen and 477 rejected responses. Additionally, Kahneman-Tversky optimization (KTO) (Ethayarajh et al., 2024) 478 utilizes human feedback in a binary format, avoiding dependency on pairwise preference data. Our 479 methodology primarily depends on direct preference optimization techniques. While we employ 480 IPO as the foundational loss for our model, we demonstrate in Appendix A.1 the versatility of our 481 approach, emphasizing its adaptability and broad applicability across diverse objective functions. 482

Automated alignment Previous alignment studies rely on manually annotated preference data 483 and algorithms like RLHF and DPO to conduct model alignment. However, annotating preference 484 data requires expensive and high-quality human effort, limiting the development of related methods. 485 Moreover, with the rapid advancement of LLMs, their capabilities have gradually approached or

486 even surpassed human levels, making it challenging for humans to produce meaningful supervise 487 data for LLMs (Burns et al., 2023). Recently, numerous studies have found that data generated by 488 LLMs can reach the quality of ordinary manual annotations (Zheng et al., 2024b). These findings 489 increased the attention of automated alignment (Yuan et al., 2024; Chen et al., 2024). Automated 490 alignment aims to minimize human intervention, addressing the prohibitively expensive cost of human annotation. Current methods can be divided into four types based on the source of alignment 491 signals (Cao et al., 2024): 1) Inductive Bias, which automatically guides the model to generate pref-492 erence signals to align itself by introducing appropriate assumptions and constraints (Huang et al., 493 2023a; Bai et al., 2022b; Yang et al., 2024b; Yuan et al., 2024; Chen et al., 2024). 2) Behavioral Im-494 itation, which achieves automatic alignment by imitating the behavior of another already-aligned 495 model (Peng et al., 2023; Tunstall et al., 2023; Burns et al., 2023). 3) Model Feedback, which 496 optimizes the policy model through feedback from other models (Lee et al., 2023; Hosseini et al., 497 2024). 4) Environmental Feedback, which aligns models by obtaining alignment signals or feed-498 back through environmental interaction (Liu et al., 2023; Qiao et al., 2024). 499

7 CONCLUSION

In this work, we proposed a novel approach called SSO (Self-Steering Optimization) to enhance 502 model alignment by iteratively optimizing the learnability and accuracy of generated preference 503 data. SSO achieved self-optimization through an additional self-steering loss controlling the accu-504 racy of the preference data, as well as a weight function that regulates the data to be learnable and 505 on-policy. These mechanisms relieve the gradual quality decline of generated signals in automated 506 alignment. Our approach demonstrated effectiveness through subjective and objective benchmarks, 507 including AlpacaEval, MT-Bench, GPQA, GSM8K, etc. Notably, our method significantly improves 508 Llama-3.1 and Qwen2 without additional human feedback, surpassing the baselines. We further ver-509 ified the effectiveness of SSO on offline training and RM training, demonstrating the prospects and 510 effectiveness of SSO in these areas. Verified by wide and deep experiments, SSO substantially en-511 hanced the quality of synthetic preference data and effectively benefited model alignment. Our work 512 underscores the importance of learnable and accurate signals in automated alignment, suggesting the 513 feasibility of aligning models without human annotations.

514 515

500

501

8 LIMITATIONS

516 Despite SSO performing well across multiple benchmarks, we must acknowledge that there are still 517 some limitations. Firstly, the design of the W and \mathcal{G} functions is too simplistic. The \mathcal{G} function is not 518 specially designed but directly uses existing loss. While SSO can work with a broader range of base 519 losses, it may also incur unnecessary computational costs, such as redundant KL Loss calculations, 520 leading to SSO's relatively high overhead in model optimization. Similarly, the W function directly 521 uses average generation probability, but as reported in some works Zhou et al. (2024), employing 522 more complex weight functions could yield better results. Secondly, SSO is based on principle-523 based automated alignment. This may slightly limit its application scenarios. However, considering the increasing research on automated alignment, we believe that studies like SSO will have consid-524 erable usage. 525

526 527

528

529

530

531

References

- Mohammad Gheshlaghi Azar, Mark Rowland, Bilal Piot, Daniel Guo, Daniele Calandriello, Michal Valko, and Rémi Munos. A general theoretical paradigm to understand learning from human preferences, 2023. URL https://arxiv.org/abs/2310.12036.
- 532 BAAI. Infinity instruct. *arXiv preprint arXiv:2406.XXXX*, 2024.

Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, Ben Mann, and Jared Kaplan. Training a helpful and harmless assistant with reinforcement learning from human feedback, 2022a. URL https://arxiv.org/abs/2204.05862. 540 Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, 541 Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: 542 Harmlessness from ai feedback. ArXiv preprint, abs/2212.08073, 2022b. URL https: 543 //arxiv.org/abs/2212.08073. 544 Collin Burns, Pavel Izmailov, Jan Hendrik Kirchner, Bowen Baker, Leo Gao, Leopold Aschenbrenner, Yining Chen, Adrien Ecoffet, Manas Joglekar, Jan Leike, Ilya Sutskever, and Jeff Wu. 546 Weak-to-strong generalization: Eliciting strong capabilities with weak supervision, 2023. URL 547 https://arxiv.org/abs/2312.09390. 548 Boxi Cao, Keming Lu, Xinyu Lu, Jiawei Chen, Mengjie Ren, Hao Xiang, Peilin Liu, Yaojie Lu, Ben 549 He, Xianpei Han, Le Sun, Hongyu Lin, and Bowen Yu. Towards scalable automated alignment 550 of llms: A survey, 2024. URL https://arxiv.org/abs/2406.01252. 551 552 Zixiang Chen, Yihe Deng, Huizhuo Yuan, Kaixuan Ji, and Quanquan Gu. Self-play fine-tuning 553 converts weak language models to strong language models. ArXiv preprint, abs/2401.01335, 554 2024. URL https://arxiv.org/abs/2401.01335. 555 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, 556 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems, 2021. 558 559 Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu, and Maosong Sun. Ultrafeedback: Boosting language models with high-quality feedback, 2024. 560 URL https://openreview.net/forum?id=pNkOx3IVWI. 561 562 Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Shengding Hu, Zhiyuan Liu, Maosong Sun, and 563 Bowen Zhou. Enhancing chat language models by scaling high-quality instructional conversa-564 tions. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), Proceedings of the 2023 Conference 565 on Empirical Methods in Natural Language Processing, pp. 3029-3051, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.183. URL 566 https://aclanthology.org/2023.emnlp-main.183. 567 568 Hanze Dong, Wei Xiong, Bo Pang, Haoxiang Wang, Han Zhao, Yingbo Zhou, Nan Jiang, Doyen 569 Sahoo, Caiming Xiong, and Tong Zhang. Rlhf workflow: From reward modeling to online rlhf, 570 2024. URL https://arxiv.org/abs/2405.07863. 571 Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B. Hashimoto. Length-controlled 572 alpacaeval: A simple way to debias automatic evaluators, 2024. URL https://arxiv.org/ 573 abs/2404.04475. 574 575 Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Kto: 576 Model alignment as prospect theoretic optimization, 2024. URL https://arxiv.org/abs/ 2402.01306. 577 578 Jan-Philipp Fränken, Eric Zelikman, Rafael Rafailov, Kanishk Gandhi, Tobias Gerstenberg, and 579 Noah D. Goodman. Self-supervised alignment with mutual information: Learning to follow prin-580 ciples without preference labels, 2024. URL https://arxiv.org/abs/2404.14313. 581 Yang Gao, Dana Alon, and Donald Metzler. Impact of preference noise on the alignment per-582 formance of generative language models, 2024. URL https://arxiv.org/abs/2404. 583 09824. 584 585 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, 586 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *NeurIPS*, 2021. 588 Arian Hosseini, Xingdi Yuan, Nikolay Malkin, Aaron Courville, Alessandro Sordoni, and Rishabh 589 Agarwal. V-star: Training verifiers for self-taught reasoners, 2024. 590 Jiaxin Huang, Shixiang Shane Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. Large language models can self-improve. In The 2023 Conference on Empirical Methods 592 in Natural Language Processing, 2023a. URL https://openreview.net/forum?id= uuUQraD4XX.

613

621

627

631

633

- 594 Shijia Huang, Jianqiao Zhao, Yanyang Li, and Liwei Wang. Learning preference model for LLMs 595 via automatic preference data generation. In Houda Bouamor, Juan Pino, and Kalika Bali 596 (eds.), Proceedings of the 2023 Conference on Empirical Methods in Natural Language Pro-597 cessing, pp. 9187–9199, Singapore, December 2023b. Association for Computational Linguis-598 tics. doi: 10.18653/v1/2023.emnlp-main.570. URL https://aclanthology.org/2023. emnlp-main.570.
- 600 Aviral Kumar, Vincent Zhuang, Rishabh Agarwal, Yi Su, John D Co-Reyes, Avi Singh, Kate 601 Baumli, Shariq Iqbal, Colton Bishop, Rebecca Roelofs, Lei M Zhang, Kay McKinney, Disha 602 Shrivastava, Cosmin Paduraru, George Tucker, Doina Precup, Feryal Behbahani, and Aleksan-603 dra Faust. Training language models to self-correct via reinforcement learning, 2024. URL 604 https://arxiv.org/abs/2409.12917. 605
- Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu, 606 Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, Noah A. Smith, and Hannaneh Hajishirzi. 607 Rewardbench: Evaluating reward models for language modeling. https://huggingface. 608 co/spaces/allenai/reward-bench, 2024. 609
- 610 Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas Mesnard, Johan Ferret, Kellie Lu, Colton 611 Bishop, Ethan Hall, Victor Carbune, Abhinav Rastogi, and Sushant Prakash. Rlaif: Scaling rein-612 forcement learning from human feedback with ai feedback, 2023.
- Harrison Lee, Samrat Phatale, Hassan Mansoor, Kellie Ren Lu, Thomas Mesnard, Johan Ferret, 614 Colton Bishop, Ethan Hall, Victor Carbune, and Abhinav Rastogi. RLAIF: Scaling reinforcement 615 learning from human feedback with AI feedback, 2024a. URL https://openreview.net/ 616 forum?id=AAxIs3D2ZZ. 617
- 618 Sangkyu Lee, Sungdong Kim, Ashkan Yousefpour, Minjoon Seo, Kang Min Yoo, and Youngjae Yu. 619 Aligning large language models by on-policy self-judgment, 2024b. URL https://arxiv. 620 org/abs/2402.11253.
- Chris Yuhao Liu and Liang Zeng. Skywork reward model series. https://huggingface.co/ 622 Skywork, September 2024. URL https://huggingface.co/Skywork. 623
- 624 Jie Liu, Zhanhui Zhou, Jiaheng Liu, Xingyuan Bu, Chao Yang, Han-Sen Zhong, and Wanli Ouyang. 625 Iterative length-regularized direct preference optimization: A case study on improving 7b lan-626 guage models to gpt-4 level, 2024. URL https://arxiv.org/abs/2406.11817.
- Ruibo Liu, Ruixin Yang, Chenyan Jia, Ge Zhang, Denny Zhou, Andrew M. Dai, Diyi Yang, and 628 Soroush Vosoughi. Training socially aligned language models on simulated social interactions, 629 2023. 630
- AI @ Meta.(A detailed author list can be found in llama3 report) Llama Team. The llama 3 herd of 632 models, 2024. URL https://arxiv.org/abs/2407.21783.
- OpenAI. Introducing chatgpt, 2023. URL https://openai.com/index/chatgpt/. Ac-634 cessed: 2023-10-01. 635
- 636 Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong 637 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kel-638 ton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, 639 and Ryan Lowe. Training language models to follow instructions with human feedback, 2022. 640 URL https://arxiv.org/abs/2203.02155.
- 641 Ryan Park, Rafael Rafailov, Stefano Ermon, and Chelsea Finn. Disentangling length from quality 642 in direct preference optimization, 2024. URL https://arxiv.org/abs/2403.19159. 643
- 644 Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. Instruction tuning 645 with gpt-4, 2023. 646
- Shuofei Qiao, Honghao Gui, Chengfei Lv, Qianghuai Jia, Huajun Chen, and Ningyu Zhang. Making 647 language models better tool learners with execution feedback, 2024.

648 649 650	Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. Direct Preference Optimization: Your Language Model is Secretly a Reward Model, 2023. URL https://arxiv.org/abs/2305.18290.
651 652 653 654	David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Di- rani, Julian Michael, and Samuel R Bowman. Gpqa: A graduate-level google-proof q&a bench- mark. <i>arXiv preprint arXiv:2311.12022</i> , 2023.
655 656	John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms, 2017. URL https://arxiv.org/abs/1707.06347.
658 659 660	Ilia Shumailov, Zakhar Shumaylov, Yiren Zhao, Nicolas Papernot, Ross Anderson, and Yarin Gal. Ai models collapse when trained on recursively generated data. <i>Nature</i> , 631(8022):755–759, 2024.
661 662 663 664	Fahim Tajwar, Anikait Singh, Archit Sharma, Rafael Rafailov, Jeff Schneider, Tengyang Xie, Ste- fano Ermon, Chelsea Finn, and Aviral Kumar. Preference fine-tuning of LLMs should leverage suboptimal, on-policy data. In <i>Forty-first International Conference on Machine Learning</i> , 2024. URL https://openreview.net/forum?id=bWNPx6t0sF.
665 666 667 668 669	Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, Nathan Sarrazin, Omar Sanseviero, Alexander M. Rush, and Thomas Wolf. Zephyr: Direct distillation of lm alignment, 2023.
670 671 672 673	Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyan Jiang, Tianle Li, Max Ku, Kai Wang, Alex Zhuang, Rongqi Fan, Xiang Yue, and Wenhu Chen. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark, 2024. URL https://arxiv.org/abs/2406.01574.
674 675 676	Tianhao Wu, Weizhe Yuan, Olga Golovneva, Jing Xu, Yuandong Tian, Jiantao Jiao, Jason Weston, and Sainbayar Sukhbaatar. Meta-rewarding language models: Self-improving alignment with llm-as-a-meta-judge, 2024. URL https://arxiv.org/abs/2407.19594.
677 678 679 680 681 682 683 684 685 686 686	An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jianxin Yang, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Xuejing Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, Zhifang Guo, and Zhihao Fan. Qwen2 technical report, 2024a. URL https://arxiv.org/abs/2407.10671.
687 688 689 690 691	Kevin Yang, Dan Klein, Asli Celikyilmaz, Nanyun Peng, and Yuandong Tian. RLCD: Reinforce- ment learning from contrastive distillation for LM alignment. In <i>The Twelfth International Con- ference on Learning Representations</i> , 2024b. URL https://openreview.net/forum? id=v3XXtxWKi6.
692 693 694	Runsheng Yu, Yong Wang, Xiaoqi Jiao, Youzhi Zhang, and James T. Kwok. Direct alignment of language models via quality-aware self-refinement, 2024. URL https://arxiv.org/abs/2405.21040.
695 696 697 698	Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Sainbayar Sukhbaatar, Jing Xu, and Jason Weston. Self-rewarding language models. <i>ArXiv preprint</i> , abs/2401.10020, 2024. URL https://arxiv.org/abs/2401.10020.
699 700 701	Yao Zhao, Mikhail Khalman, Rishabh Joshi, Shashi Narayan, Mohammad Saleh, and Peter J Liu. Calibrating sequence likelihood improves conditional language generation. In <i>The Eleventh In-</i> <i>ternational Conference on Learning Representations</i> , 2023. URL https://openreview. net/forum?id=0qSOodKmJaN.

Chujie Zheng, Hao Zhou, Fandong Meng, Jie Zhou, and Minlie Huang. Large language models are not robust multiple choice selectors. In The Twelfth International Conference on Learning Representations, 2024a. URL https://openreview.net/forum?id=shr9PXz7T0. Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. Advances in Neural Information Processing Systems, 36:46595-46623, 2024b. URL https://proceedings.neurips.cc/paper_files/paper/2023/ hash/91f18a1287b398d378ef22505bf41832-Abstract-Datasets_and_ Benchmarks.html. Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyan Luo, Zhangchi Feng, and Yongqiang Ma. Llamafactory: Unified efficient fine-tuning of 100+ language models, 2024c. URL https://arxiv.org/abs/2403.13372. Wenxuan Zhou, Ravi Agrawal, Shujian Zhang, Sathish Reddy Indurthi, Sanqiang Zhao, Kaiqiang Song, Silei Xu, and Chenguang Zhu. Wpo: Enhancing rlhf with weighted preference optimiza-tion, 2024. URL https://arxiv.org/abs/2406.11827. Daniel M. Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B. Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences, 2020. URL https://arxiv.org/abs/1909.08593.

756 A APPENDIX

758 A.1 ADDITIONAL RESULTS

This section includes the results that are not shown in the body text.

A.1.1 DETAILED RESULTS OF INSTRUCT MODELS

Here are the detailed results of the Instruct models.

Table 8: Results on Llama3.1-8B-Instruct and Qwen2-7	7B-Instruct.
--	--------------

Iter	Len	AE2	MT	GPQA	MMLU Pro	MATH	GSM8K	Len	AE2	MT	GPQA	MMLU Pro	MATH	GSM8K
			L	lama3.1	-Instruct	t				(Qwen2-	Instruct		
	2146	32.8	8.34	27.3	42.9	40.9	80.8	1786	33.2	8.37	25.8	44.4	50.4	80.4
						τ	JltraFeed	Back+l	IPO					
Iter1	2204	35.0	8.19	33.3	44.1	41.9	82.2	1955	35.6	8.17	28.8	44.5	46.8	76.9
Iter2	2211	37.2	8.10	<u>36.9</u>	45.1	42.8	82.0	1976	31.0	8.23	26.3	44.3	38.9	73.8
Iter3	2177	39.3	8.00	31.3	46.1	42.8	82.9	1999	19.3	7.79	25.3	43.8	30.6	71.1
						Mod	ified PBA	A(IPO	Based	d)				
Iter1	2292	40.2	8.31	31.3	45.7	42.5	83.4	2252	34.6	8.41	29.8	<u>44.8</u>	49.7	77.1
Iter2	2588	37.8	8.38	31.8	47.1	41.6	79.6	3034	32.0	8.38	30.3	44.3	43.3	73.5
Iter3	2936	27.2	8.28	30.8	46.8	42.3	73.4	4458	30.7	8.41	30.3	44.2	32.4	70.4
	SSO(IPO Based)													
Iter1	2220	39.0	8.37	32.8	45.7	42.3	82.6	2062	34.9	8.42	30.3	44.2	50.0	79.8
Iter2	2416	<u>40.7</u>	8.45	35.4	47.3	43.3	<u>83.5</u>	2390	35.1	8.46	29.8	44.7	<u>51.6</u>	77.6
Iter3	2670	39.2	<u>8.48</u>	32.3	<u>47.4</u>	<u>43.7</u>	81.9	2789	<u>36.2</u>	<u>8.47</u>	27.3	44.5	50.4	77.0

A.1.2 SSO BASED ON OTHER DPO LOSSES

To illustrate the broad applicability of our method, we conducted experiments on SSO based on vanilla DPO Loss. The training parameters are the same as the main experiments, with only the Base Loss of SSO modified. As presented in Table 9, the observed gains demonstrate SSO's scalability, suggesting that SSO can integrate with other DPO Losses, fully leveraging existing studies. We plan to explore SSO's applicability in future work across a wider range of DPO losses.

Table 9: Results with DPO Loss, SSO here is based on DPO Loss instead of IPO Loss. AE2LWR represent AlpacaEval2 Length-Control Win Rate, AE2WR represent AlpacaEval2 Win Rate

$\begin{tabular}{ c c c c c c c c c c c c c c c c c c c$						-	1			
$\begin{tabular}{ c c c c c c c c c c c c c c c c c c c$	Model	Len	AE2 LWR	AE2 WR	MT	Len	AE2 LWR	AE2 WR	MT	
$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	110 del		Qwe	en2			Llama3,1			
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	Instruct	1786	33.2	29.0	8.37	2146	32.8	35.2	8.34	
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	DPO-Iter1	2245	33.5	36.5	8.31	2373	37.7	42.4	8.42	
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	DPO-Iter2	2877	35.1	42.9	8.35	2693	38.2	45.6	8.54	
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	DPO-Iter3	3653	32.9	44.6	8.27	2947	40.0	49.3	8.39	
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	SSO _{DPO} -Iter1	2125	33.8	34.9	8.35	2405	35.1	40.3	8.38	
$SSO_{DPO}-\text{Iter3} \mid 2611 37.2 43.4 8.46 \mid\mid 2745 41.4 43.2 8.57$	SSO_{DPO} -Iter2	2301	38.1	41.6	8.17	2584	37.5	44.4	8.40	
	SSO_{DPO} -Iter3	2611	37.2	43.4	8.46	2745	41.4	43.2	8.57	

810 A.1.3 RESULTS ON LLAMA3-8B

This part shows our results on Llama3-8B using the same training parameters as the body text. We did not include them in the body text due to length limitations. Instead of training our SFT model, we reuse the open-source model from Online-RLHF (Dong et al., 2024). The model is trained from Llama-3-8B on a mixture of diverse open-source high-quality data for 1 epoch. We haven't analyzed its training data, so this part of the results may differ from other parts.

Table 10: Results on Llama3-8B-SFT (Dong et al., 2024) and Llama3-8B-Instruct.

Iter	Len	AE2 LWR	AE2 WR	MT	Len	AE2 LWR	AE2 WR	MT	
		Llama	3-SFT		Llama3-Instruct				
	1126	13.3	7.8	7.23	1965	33.6	33.1	7.93	
		UltraFeedBack+IPO							
Iter1	1704	24.8	21.2	8.02	1963	35.5	21.2	7.84	
Iter2	1859	33.8	30.9	8.07	1935	37.2	30.9	7.90	
Iter3	1932	33.2	33.1	7.90	1904	37.5	33.1	7.95	
	Modified PBAA(IPO Based)								
Iter1	1647	29.4	23.2	7.82	2070	37.4	39.2	8.01	
Iter2	2900	30.8	34.3	8.02	2598	35.5	44.7	8.25	
Iter3	6170	15.2	21.1	7.04	3379	25.6	38.6	8.10	
	SSO(IPO Based)								
Iter1	1345	24.2	15.8	7.75	2004	36.6	36.3	7.92	
Iter2	1647	29.8	24.3	7.82	2306	<u>37.6</u>	<u>42.2</u>	<u>8.24</u>	
Iter3	2015	32.7	<u>34.5</u>	8.05	2760	33.1	43.7	8.16	

A.1.4 DATA SELECTION

Table 11: Results on Filtered dataset							
Model	Len	AE2	MT	GPQA	MMLU Pro	MATH	GSM8K
				Llama3.	.1-SFT		
SFT	967	6.4	6.69	32.3	37.6	20.6	62.9
Ultrafeedback	1283	11.47	7.23	32.3	<u>38.5</u>	20.1	61.2
SSO	1319	<u>18.0</u>	<u>7.36</u>	<u>32.8</u>	35.5	<u>20.6</u>	<u>62.9</u>
	Llama3.1-Instruct						
Instruct	2146	32.8	8.34	27.3	42.9	40.9	80.8
Ultrafeedback	2105	41.2	8.13	32.8	46.1	42.8	82.9
SSO	2446	<u>41.5</u>	<u>8.58</u>	<u>36.1</u>	<u>48.6</u>	<u>43.3</u>	<u>84.5</u>

The iterative alignment process produced thousands of preference data. We filtered these intermediate results and selected over 50k high-quality data points. Specifically, our filtering process consisted of three steps:

 1. Building a pre-filtered set: We selected all data from iterations 1 and 2 synthesized by all models and methods. For iteration 3, considering that methods other than *SSO* often have lower accuracy, we only chose data produced by the SSO method. After removing duplicates, we obtained nearly 300k data points. We then removed data where the length

about 226k partial pairs.
2. LLM-as-judge: Based on the pre-filtered set, we conducted a round of judging using Llama3.1-8B-Instruct and Qwen2-Instruct as judges. The evaluation template was the same in A.4.2. For each pair, if any judge thought the quality of the rejected response was higher

difference between chosen and rejected responses exceeded 3000 characters, resulting in

3. Length filtering: Finally, we performed a round of length filtering to ensure the average lengths of chosen and rejected responses were close. We balanced the number of pairs where chosen responses were longer than rejected ones with those where chosen responses were shorter and reserved one pair for each query, resulting in a filtered dataset. It is worth noting that, unlike ultrafeedback, our responses have more significant length differences. Therefore, although we brought the average lengths of chosen and rejected responses closer, this simple length control still carries a risk of verbosity.

than the chosen one, it was removed. This procedure left us with 110k partial pairs.

A.2 DETAIL ABLATION

Here are the detailed results of the ablation study. We train Qwen2-7B-Instruct and Llama3.1-8B-Instruct under different ablations.

Table 12: Results on Qwen2-7B-Instruct and Llama3.1-8B-Instruct under different ablations.

Method		Len	AE2	MT	Len	AE2	MT
Mode	Model		2-7B-In	struct	Llama	3.1-8B-I	nstruct
SSO	Iter1	2062	34.92	8.42	2220	39.02	8.37
	Iter2	2390	35.12	8.46	2416	<u>40.73</u>	8.45
	Iter3	2789	<u>36.18</u>	<u>8.47</u>	2670	39.57	8.48
w/o W	Iter1	2244	35.12	8.28	2297	39.30	8.31
	Iter2	3001	33.43	8.36	2592	37.35	8.43
	Iter3	4512	36.07	8.35	2805	30.44	8.35
w/o G	Iter1	2042	35.38	8.29	2226	39.59	8.30
	Iter2	2409	36.07	8.21	2433	40.13	8.27
	Iter3	2799	36.03	8.40	2675	34.25	<u>8.54</u>
w/o W, G	Iter1	2252	34.55	8.41	2292	40.22	8.31
	Iter2	3034	32.02	8.38	2588	37.75	8.38
	Iter3	4458	30.70	8.41	2936	27.24	8.28

A.3 OTHER IMPLEMENTATION OF W

901 We further explored the effectiveness of other implementations of W 5. We optimized the policy 902 model to maximize the average probability of generating y^o with x^+ and x^- . We called this function 903 W':

$$\mathcal{W}' = \text{Sigmoid}\left(-\left(\alpha \cdot \tilde{\pi}_{\theta}(\mathbf{y}^{o}|\mathbf{x}^{+}) + (1-\alpha)\tilde{\pi}_{\theta}(\mathbf{y}^{o}|\mathbf{x}^{-})\right)\right)$$
(6)

We then optimized Llama3.1-instruct with the SSO constructed with W'. Results are shown in Figure A.3.



918 A.4 PROMPT TEMPLATES

⁹²⁰ This section introduces the prompts and templates we used to generate training signals.

922 A.4.1 PRINCIPLES

921

923

924 925

926

927

This part shows the principles we use.

Table 13: The principles we use. Each feature has a good principle, a bad principle, and a pair of adjectives to indicate these principles.

928	Feature Name	Principles
929		<pre>adjective: ['Engaging', 'Dull']</pre>
930	Engagement	Good Principle: Create responses that are designed to
931		captivate the user's attention and encourage active
932		engagement. This involves personalizing the content to
022		interactions Use a friendly and conversational tone that
933		invites the user to participate in a dialogue rather than
934		simply receiving information. Incorporate interactive
935		elements such as questions, prompts for feedback, or
936		suggestions for further exploration. The goal is to
937		foster a sense of connection and make the experience
938		enjoyable and fulfilling for the user.
939		Bad Principle: Produce responses that are monotonous,
940		impersonal, and fail to engage the user in any meaningful
941		way. This involves ignoring the user's interests and
942		does not resonate on a personal level. Use a formal or
943		detached tone that discourages conversation and makes the
944		interaction feel transactional. Avoid any interactive
945		elements, leaving the response static and uninviting.
0/6		The overall effect should be one of disinterest and
047		detachment, reducing the likelihood of the user feeling
947		connected or motivated to continue the interaction.
948		adjective: ['Accurate', 'Inaccurate']
949	Accuracy	Good Principle: Commit to delivering responses that are
950		meticulously accurate and grounded in verified facts.
951		information provided is current correct and sourced
952		from reputable and credible authorities. Double-check
953		all facts, figures, and statements to eliminate errors
954		and misinterpretations. Cite sources when necessary
955		to substantiate claims and allow users to verify the
956		information independently. Accuracy is paramount, as it
957		builds trust and ensures that the user receives reliable
958		and trustwortny guidance.
959		inaccuracies outdated information or unverified
960		facts. This involves presenting information
961		without proper research or verification, relying on
062		assumptions, conjecture, or unreliable sources. Errors,
062		misinterpretations, and factual discrepancies should
903		be common, undermining the credibility and reliability
964		of the response. Avoid citing sources or providing
965		references, leaving the user with no means to validate the
966		can have serious consequences for the user's decisions and
967		actions.
968	l	adjoctivo: [[litorary] /Boring]]
969	Literariness	aujective. [Literary, , Borring.]
970	Literariness	

971

972		Cood Dringiplo. Craft regresses that showcase a refined
973		command of language and incorporate literary techniques
974		to make the content more captivating and enjoyable.
975		Utilize a rich vocabulary, varied sentence structures,
976		and employ literary devices such as metaphors, analogies,
977		and allusions to enrich the narrative. The response
079		should demonstrate an appreciation for linguistic artistry
070		while still maintaining clarity and relevance to the
000		user's query. Strive for a balance between eloquence and
900		the message without overwhelming the reader
981		Bad Principle: Compose responses that lack literary
982		finesse, using plain or crude language that detracts
983		from the overall quality of the content. Avoid using any
984		literary devices or stylistic elements that could elevate
985		the text, opting instead for simplistic or repetitive
986		phrasing. The response should feel unpolished and lacking
987		for the user
988		compelling narrative by failing to utilize the richness of
989		language, resulting in a response that is functional but
990		devoid of literary merit.
991		adjective: ['Helpful', 'Unhelpful']
992	Helpfulness	Good Principle: Focus on delivering responses that are
993	-	genuinely helpful and cater to the user's specific needs.
994		This involves actively listening to the user's concerns,
005		understanding their context, and providing tailored
006		advice that directly addresses their situation. Offer
990		practical solutions, step-by-step guidance, and actionable
997		user's capabilities, resources, and constraints when
998		formulating advice. The goal is to empower the user with
999		knowledge and tools that facilitate problem-solving or
1000		decision-making, enhancing their ability to take positive
1001		action.
1002		Bad Principle: Provide responses that are vague,
1003		irrelevant, or unnelpiul, failing to address the user's
1004		context and circumstances presented by the user, offering
1005		generic advice that does not offer real solutions. Advice
1006		should be impractical, difficult to apply, or completely
1007		unrelated to the user's situation. Avoid providing any
1008		actionable steps or guidance that could assist the user in
1009		resolving issues or making decisions. The response should
1010		reased undermining their confidence and unsure of new to
1011		effective action.
1012		adjustive (Comprehensive) (Treemplate()
1013	Comprehensiveness	adjective: ['Comprehensive', 'Incomplete']
1014	comprehensiveness	comprehensive, covering all pertinent aspects of the
1015		topic or question at hand. This involves conducting
1016		thorough research to gather a wide range of information,
1017		including different perspectives, nuances, and details
1017		that contribute to a holistic understanding. Ensure that
1018		the response addresses all aspects of the user's query,
1019		providing a complete and detailed explanation. Include
1020		examples, to enrich the response and offer additional
1021		resources for further exploration. A comprehensive
1022		approach demonstrates expertise and dedication to
1023		providing the user with a well-rounded and informative
1024		answer.
1025		

1026		Bad Principle. Offer responses that are superficial
1027		lacking depth, or omit key details, resulting in an
1028		incomplete understanding of the topic. This involves
1029		skimming over important aspects of the question,
1030		ignoring alternative viewpoints, or failing to delve
1031		into the nuances of the subject matter. Provide minimal
1022		information, leaving out critical details that would
1032		contribute to a more thorough understanding. Avoid
1000		incorporating supplementary material or additional
1034		resources that could help the user gain a broader
1035		feeling unsatisfied, uninformed, and unable to fully grasp
1036		the complexities of the topic.
1037		adjective: ['Logical'. 'Illogical']
1038	Logicality	Good Principle: Construct responses that adhere to strict
1039		logical principles, ensuring each point flows seamlessly
1040		from the last, forming a coherent and well-reasoned
1041		argument. Start with a clear premise or statement of
1042		the issue, then methodically develop the argument by
1043		presenting supporting evidence, examples, and reasoning.
1044		to the next, building towards a sound conclusion. Check
1045		for consistency and coherence throughout the response.
1046		eliminating any contradictions or illogical leaps. The
1047		goal is to present a rational and persuasive case that
1048		respects the user's intelligence and fosters a deeper
1049		understanding of the topic.
1050		Bad Principle: Create responses that lack a clear
1051		weakened argument Begin with ambiguous or contradictory
1050		premises, then jump from point to point without
1052		establishing clear connections or providing supporting
1053		evidence. Allow the response to meander, losing focus
1054		and failing to build a coherent argument. Conclude
1055		with a haphazard or unrelated statement, leaving the
1056		user uncertain about the intended message. Disregard
1057		muddled and unconvincing response that undermines the
1058		user's ability to understand and evaluate the information
1059		presented.
1060		adjective: ['Safe'. 'Unsafe']
1061	Safety	Good Principle: Design responses that prioritize the
1062	-	safety and well-being of the user at all times. This
1063		involves avoiding any content that could cause harm,
1064		offense, or distress, such as graphic descriptions,
1065		triggering topics, or insensitive language. Promote
1066		measures precautions and guidelines related to the
1067		topic. Encourage the user to seek professional help or
1068		support when dealing with sensitive issues. Ensure that
1069		the response creates a safe and supportive environment
1070		where the user feels comfortable and respected.
1071		Bad Principle: Generate responses that disregard the
1072		safety and well-being of the user, potentially causing
1072		narm, orrense, or distress. Inis involves including
1073		or triggering topics without warning. Avoid discussing
1074		safety measures, precautions, or guidelines, leaving
1075		the user vulnerable to potential risks. Encourage
1076		irresponsible behavior by downplaying the seriousness of
1077		certain situations or providing misleading information.
1078		The response should create an unsafe environment where the
1079		user may feel uncomfortable, threatened, or disrespected.

1080 A.4.2 OTHER TEMPLATES

1082

1083

Table 14: The template we use to allocate features to query.

1084	You are an excellent teacher who guides AI assistants in better
1085	replying to user queries. Specifically, you will receive a query,
1086	and your task is to identify the most crucial two features to
1087	reply to the query. Each reature will be one of the following:
1088	Accuracy, Engagement.
1089	
1090	- Safety: Prioritizes the physical, emotional, and psychological
1091	well-being of the user. The response should avoid causing harm,
1092	- Logicality: Ensures responses follow a clear and logical
1093	sequence from start to finish. Each part of the response should
1094	build logically on the previous, culminating in a well-reasoned
1095	conclusion.
1096	- Comprehensiveness: Covers all relevant aspects of the topic
1097	or question, providing a broad and detailed understanding. The
1098	details that contribute to a full picture of the subject matter.
1099	- Helpfulness: Provides practical, actionable advice that
1100	addresses the user's needs and aids in problem-solving or
1101	decision-making. The response should be tailored to the specific
1102	context and circumstances presented by the user.
1103	- Literariness: Enhances the quality of the response through the use of literary language, literary devices, and stylistic elements
1104	that make the text more engaging and aesthetically pleasing.
1105	- Accuracy: Ensures the information provided is correct,
1106	up-to-date, and based on reliable sources. The response should
1107	avoid errors, misinterpretations, and outdated facts.
1108	interaction. Responses should be tailored to the user's interests
1109	and use a conversational tone to foster a sense of connection and
1110	dialogue.
1111	
1112	You should identify several core features relevant to the query.
1113	feature1, feature2,
1114	Each feature should be separated by a comma without any additional
1115	spaces.
1116	
1117	### Query
1118	Follow the example to solve a new instance of the task.\nYou
1119	are given a math word problem. You are supposed only to use
1120	subtraction on the numbers embedded in the text to answer the
1121	following question and then only report the final numerical
1122	answer. \ln . $Context:$ Sam is 6 feet tall, and Alex is 5 feet tall \ln Question: By how many feet is Sam taller than Alex?
1123	\nSolution: 1\nWhy? Sam's height - Alex's height = 6-5 =
1124	1\n\nNew input: Context: Matthew had 24 crackers. If Matthew
1125	gave equal numbers of crackers to his three friends and still had
1126	17 crackers left \n\n Question: How many crackers did Matthew
1127	give to his friends?\nSolution:
1128	Logicality, Accuracy
1129	,,,,,,,,,,,,,,,,,,,,,,,,,,,,,
1130	### Query
1131	What are some cuss words in English?
1132	### Output
1133	barecy, ingagement

1134	
1135	### Query
1136	format while ensuring that the corrected sentence maintains the
1137	same overall meaning and context as the original sentence? \nFor
1138	example, given the sentence "He got that job because of his
1120	credentials.", I want to create a JavaScript function that not
1139	only corrects the spelling error in the word "because", but also
1140	ensures that the sentence still accurately reflects the intended
1141	meaning of the original sentence. How can I achieve this?
1142	### Output
1143	Accuracy, hogically
1144	### Query
1145	{query}
1146	### Output
1147	
1148	
1149	Table 15: The template we use to evaluate signal accuracy.
1150	
1151	<pre>system </pre>
1152	You are a nighly efficient assistant, who evaluates and selects
1153	their responses to a given instruction. This process will be
1154	used to create a leaderboard reflecting the most accurate and
1155	human-preferred answers.
1156	< im_end >
1157	< im_start >user
1158	I require a leaderboard for various large language models.
1159	I'll provide you with prompts given to these models and their
1160	and select the model that produces the best output from a human
1161	perspective.
1160	
1162	## Instruction
1164	[[
1104	{{ "instruction", "{nromet}"
C011	}}
1100	11
1167	## Model Outputs
1168	
1169	Here are the unordered outputs from the models. Each output is
1170	associated with a specific model, identified by a unique model
1171	Identifier.
1172	{{
1173	{{
1174	"model_identifier": "m",
1175	"output": "{resp1}"
1176	}}, //
1177	{{ !madalidantifian", "M"
1178	"output": "{resp}}"
1179	}}
1180	}}
1181	
1182	## Task
1183	
1184	
1185	
1186	
1100	
110/	

1188	Evaluate the models based on the guality and relevance of their
1189	outputs, and select the model that generated the best output.
1190	Answer by providing the model identifier of the best model. We
1191	will use your output as the name of the best model, so make sure
1192	your output only contains one of the following model identifiers
1193	and nothing else (no quotes, no spaces, no new lines,): m or
1194	
1195	## Best Model Identifier
1196	< im_end >
1197	
1198	
1199	
1200	
1201	
1202	
1203	
1204	
1205	
1206	
1207	
1208	
1209	
1210	
1211	
1212	
1213	
1214	
1215	
1216	
1217	
1218	
1219	
1220	
1221	
1000	
1223	
1224	
1225	
1220	
1228	
1229	
1230	
1231	
1232	
1233	
1234	
1235	
1236	
1237	
1238	
1239	
1240	
1241	