
Comparing Bad Apples to Good Oranges: Aligning Large Language Models via Joint Preference Optimization

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 A common technique for aligning large language models (LLMs) relies on acquiring
2 human preferences by comparing multiple generations conditioned on a fixed
3 context. This only leverages the pairwise comparisons when the generations are
4 placed in an identical context. However, such conditional rankings often fail to
5 capture the complex and multidimensional aspects of human preferences. In this
6 work, we revisit the traditional paradigm of preference acquisition and propose a
7 new axis that is based on eliciting preferences jointly over the instruction-response
8 pairs. While prior preference optimizations are designed for conditional ranking
9 protocols (e.g., DPO), our proposed preference acquisition protocol introduces
10 DOVE, a new preference optimization objective that upweights the joint probability
11 of the chosen instruction-response pair over the rejected instruction-response pair.
12 Interestingly, we find that the LLM trained with joint instruction-response preference
13 data using DOVE outperforms the LLM trained with DPO by 5.2% and 3.3%
14 win-rate for the summarization and open-ended dialogue datasets, respectively.
15 Our findings reveal that joint preferences over instruction and response pairs can
16 significantly enhance the alignment of LLMs by tapping into a broader spectrum
17 of human preference elicitation. We will release the data, code, and models upon
18 acceptance.

19 1 Introduction

20 Recently, alignment [Stiennon et al., 2020, Ouyang et al., 2022] has emerged as a crucial step in
21 enhancing the performance of large language models (LLMs) [Anthropic, 2024, OpenAI, 2023, Team
22 et al., 2023, Anthropic, 2023, Brown et al., 2020, Touvron et al., 2023, Jiang et al., 2023] in diverse
23 real-world applications [Li et al., 2023, Zheng et al., 2023a, Wu et al., 2023a, Clusmann et al., 2023,
24 Lambert et al., 2024]. In particular, the aligned LLMs generate responses that maximize human
25 utility along various dimensions such as helpfulness, coherence and harmlessness [Askell et al., 2021,
26 Ouyang et al., 2022]. Here, the notion of human utility is subjective Kirk et al. [2024], Gabriel [2020],
27 and mainly hinges on *how* preferences are acquired from the annotators Otto et al. [2022]. Among
28 the various preference acquisition protocols [Lightman et al., 2023, Wu et al., 2023b, Scheurer et al.,
29 2023, Bansal et al., 2023], the ranking-based approach is the most widely used paradigm for aligning
30 LLMs [Stiennon et al., 2020, Ouyang et al., 2022, Bai et al., 2022a, Tunstall et al., 2023, Teknium,
31 2023]. Specifically, in this approach the annotator has to compare a pair of responses *conditioned*
32 on a fixed context. For instance, humans can select a ‘preferred’ response by comparing a pair of
33 responses for the instruction ‘Create a list of four fruits other than Apple’ (Figure 1 (*left*)).

34 Besides ranking preferences conditioned on a fixed context, humans can also express preferences
35 in non-identical contexts. For example, while browsing reviews for products on an e-commerce

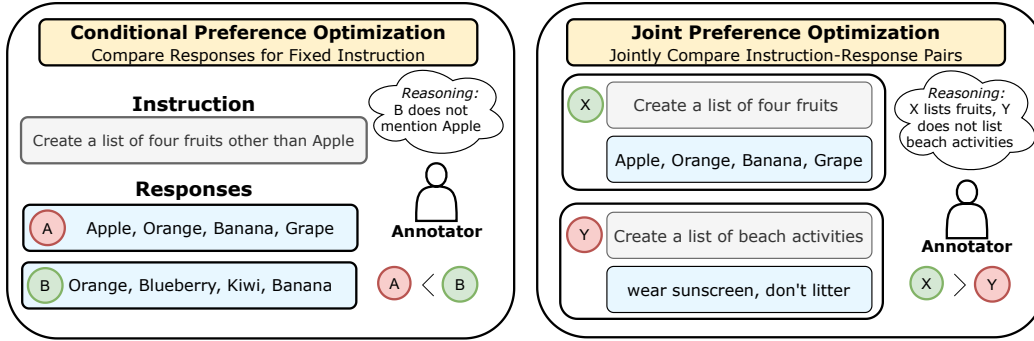


Figure 1: Overview of the Joint Preference Optimization. (Left) We show that the conditional preference acquisition method would require the annotators to compare two responses for an identical instruction. (Right) We show that the annotators can also assign rankings jointly over instruction-response pairs. Specifically, the annotator prefers a helpful response (e.g., Apple ... Grape) over a response that ignores the context of the instruction (e.g., wear sunscreen ... litter). Our framework thus elicits preferences that are obfuscated in the prior approach.

36 website, humans are likely to prefer an accurate and detail-oriented review for a camera over an
 37 incoherent, vague movie review even though the products (camera and movie) are qualitatively
 38 different. Although the traditional conditional rankings provide rich preference for alignment, they
 39 fail to holistically capture the various dimensions of reasoning of human preferences. In this work,
 40 we revisit the traditional paradigm of conditional preference acquisition and propose a new approach
 41 for jointly eliciting preferences over instruction-response pairs. This method aims to uncover diverse
 42 reasoning paths in the process of acquiring feedback.

43 In this work, we develop a framework to acquire preferences jointly over instruction-response pairs.
 44 Starting from an instruction-response data consisting of response R_i for instruction I_i (say $i \in \{1, 2\}$),
 45 we acquire ranking-based preferences over the instruction-response pairs (I_1, R_1) and (I_2, R_2) . As
 46 shown in Figure 1 (right), we aim to understand whether the response in the pair X is perceived
 47 better than the response in the pair Y . For instance, humans would prefer a helpful response to
 48 the instruction ‘Create a list of four fruits’ over a response that completely ignores the instruction
 49 ‘Create a list of beach activities’. This suggests that we can reveal preference axes like adherence to
 50 instructions, grammatical fluency, and clarity even when following joint preference optimization. In
 51 addition, our protocol can elicit human preference behaviours that are obfuscated in prior protocols,
 52 and redefines conditional preference elicitation as a special case where the instructions are identical.

53 Prior works like DPO and its variants Rafailov et al. [2023], Yin et al. [2024], Liu et al. [2024], Meng
 54 et al. [2024], Hong et al. [2024], Azar et al. [2023] rely on rankings over responses generated under
 55 an identical context, and thus do not have access to the joint distribution of human preferences in
 56 the ranking protocol (§A Table 5). While a rating protocol Ethayarajh et al. [2024] allows for a
 57 comparison between responses from non-identical instructions, it can be inconsistent with rankings
 58 Bansal et al. [2023] and ignores the possibility of preferences over a pair of chosen or rejected
 59 responses.¹ In this work, we show that humans can provide decisive preferences when comparing
 60 two instruction-responses that are chosen or rejected under the conditional rankings protocol (§3.4).

61 Next, we propose DOVE, a framework for aligning LLMs with our proposed joint preference
 62 elicitation scheme. Specifically, it upweights the joint probability of the chosen instruction-response
 63 pair over the rejected instruction-response pair. This differs from the other frameworks that assume
 64 conditional rankings in their feedback data, such as DPO [Radford et al., 2019, Azar et al., 2023]
 65 and preference optimizations that train a separate reward model such as PPO and rejection sampling
 66 [Schulman et al., 2017, Nakano et al., 2021]. We further point that DOVE subsumes the prior
 67 preference optimizations as conditional rankings are a special case of joint preferences (e.g., when
 68 $I_1 = I_2$). In our experiments, we focus on extending and comparing against DPO because of their
 69 simplicity, stability, and high-performance. However, our framework can be easily applied to reward
 70 model based approaches Schulman et al. [2017] by training a reward model on the joint preferences.

¹For instance, a pair of responses that achieves a score of 0, under the rating protocol, will result in an indecisive preference.

71 Finally, we conduct experiments to explore the new reasoning paths enabled by joint preference
72 elicitation, followed by aligning LLMs with the DOVE objective. To do so, we explore the interplay
73 between the feedback data collected under conditional rankings and joint preferences protocol. In
74 addition, we ask human annotators to explain their preference decisions, uncovering new reasoning
75 paths that highlight the complexities of the preference acquisition process (§3). After feedback
76 acquisition, we aim to investigate the impact of diverse preferences collected from conditional and
77 joint preferences on LLM alignment. In our experiments, we align a Mistral-7B LLM with the
78 preferences acquired from the conditional rankings and joint preferences, using our DOVE algorithm.
79 We find that the DOVE outperforms the supervised finetuned LLM by 30% and 18% win-rate against
80 the gold responses on the unseen instructions from the summarization and open-ended dialogues
81 datasets, respectively. Surprisingly, we find that DOVE can effectively tap into the diverse preferences
82 in the conditional and joint feedback data and outperforms DPO by 5.2% and 3.3% win-rate points
83 on the summarization and open-ended dialogues, respectively. In addition, DOVE outperforms KTO
84 by 3.5% on the open-ended dialogues dataset. This indicates that by utilizing the diverse preference
85 signals present in the existing data, we can align an LLM robustly without acquiring additional
86 instruction-response data.

87 2 Joint Preference Optimization using DOVE

88 2.1 Joint Preference Acquisition Protocol

89 In §A.1, we describe a common technique for feedback data acquisition that requires the annotators
90 to assign a preferred and non-preferred label to a pair of responses for an instruction. However, this
91 paradigm does not capture the complex and multidimensional aspects of human preferences [Kendall
92 and Smith, 1940, Thurstone, 2017]. Specifically, the reasoning paths for making preference decisions
93 depend upon the context in which the comparison is made. While the traditional ranking protocol
94 compares the two responses under a fixed context, humans can perform pairwise comparisons jointly
95 over instruction-response pairs. For example, consider two summaries, A and B, for articles X
96 and Y, respectively; then, a human can reason and choose the response that better summarizes its
97 corresponding article. Hence, it is critical to align language models with diverse feedback signals to
98 elicit high-quality responses that humans prefer under various contexts.

99 In our setup, the annotator has to decide a *chosen* and *rejected* instruction-response pair
100 (I_a, R_a, I_b, R_b) where R_a and R_b are responses to the instructions I_a and I_b , respectively, and
101 $(I_a, R_a), (I_b, R_b) \in \mathcal{D}$. We note that our joint preference setup is equivalent to the original ranking
102 protocol when $I_a = I_b$. As before, the preference reasoning from the annotator will be based on
103 subjective dimensions like helpfulness, coherence, and harmlessness. Formally, the annotator assigns
104 a joint ranking feedback $h(I_a, R_a, I_b, R_b) \in \{(I_a, R_m), (I_b, R_b), \text{Equal}\}$ where ‘Equal’ indicates
105 that both the instruction-response pairs are perceived equally good or bad. Finally, the joint preference
106 optimization creates a pairwise feedback data $\mathcal{D}_H = \{(I_a, R_a, I_b, R_b, h(I_a, R_a, I_b, R_b))\}$.

107 Our formulation suggests that we can obtain large-scale and diverse preference data (covering all
108 possible combinations of (I_a, R_a) and (I_b, R_b)) without the need for gathering additional instruction
109 and response data, which is typically more difficult and costly to acquire. In addition, joint preference
110 acquisition does not necessitate the presence of multiple responses for a given instruction that can
111 be hard to collect for low-resource languages (e.g., Kalamang²). Specifically, one can collect an
112 instruction-response data $\mathcal{D}' = \{(I_a, R_a)\}_{a=1}^{a=n}$, and acquire preferences on various combinations of
113 instruction-response pairs. Finally, we assess the interplay between the joint feedback dataset \mathcal{D}_H
114 with the conditional feedback dataset \mathcal{D}_C along with qualitative examples in §3.

115 2.2 DOVE

116 Here, we propose DOVE, a preference optimization objective that learns to align the language models
117 with the preferences acquired jointly over the instruction-response pairs. We assume a joint preference
118 dataset $\mathcal{D}_X = \{(I_i^w, R_i^w, I_j^\ell, R_j^\ell)\}$, that can be constructed from \mathcal{D}_H , where (I_i^w, R_i^w) and (I_j^ℓ, R_j^ℓ)
119 are the chosen and rejected instruction-response pairs, respectively. Similar to DPO, we start with a
120 reference model p_{ref} which is usually the supervised finetuned language model p_{stf} . Specifically, the
121 DOVE objective aims to learn an aligned model p_θ by upweighting the joint probability of preferred

²<https://endangeredlanguages.com/lang/1891?hl=en>

122 responses $p(R_i^w, I_i^w)$ over non-preferred responses $p(R_j^\ell, I_j^\ell)$. Formally, the optimization objective
 123 for DOVE, $\mathcal{L}(\theta; \mathcal{D}_X, \beta, p_{\text{ref}})$ minimizes the expectation over $(I_i^w, R_i^w, I_j^\ell, R_j^\ell) \sim \mathcal{D}_X$:

$$\mathbb{E} \left[\log \left(\sigma \left(\beta \log \frac{p_\theta(R_i^w, I_i^w)}{p_{\text{ref}}(R_i^w, I_i^w)} - \beta \log \frac{p_\theta(R_j^\ell, I_j^\ell)}{p_{\text{ref}}(R_j^\ell, I_j^\ell)} \right) \right) \right] \quad (1)$$

124 where σ denotes the sigmoid function and β is a hyperparameter. Further, we show that Eq. 3 reduces
 125 to the DPO formulation (Eq. 2) when the instructions $I_i = I_j$ in Appendix §E. We can also see that
 126 the DOVE objective aims to learn an aligned model p_θ by upweighting the conditional probability of
 127 preferred responses $p(R_i^w | I_i^w)$ over non-preferred responses $p(R_j^\ell | I_j^\ell)$, along with a correction factor
 128 based on the prior probability of the instructions under the language model $p_\theta(I_i^w)$ and $p_\theta(I_j^\ell)$. In
 129 §4, we utilize DOVE to align language models to generate human-preferred summaries and answer
 130 open-ended instructions.

131 3 Interplay between Feedback Protocols

132 3.1 Instruction-Response Acquisition

133 The instruction-response data is a collection of real-world queries that are presented to the text AI
 134 assistants. In this work, we consider two kinds of instruction-response data. First, we consider a
 135 filtered version of the TL;DR *summarization* dataset [Völske et al., 2017] from Stiennon et al. [2020]
 136 consisting of Reddit posts, their summaries, and human preferences over a pair of summaries for a
 137 given post. Throughout the dataset, the task is of summarization that is close-ended and well-defined
 138 for language models. Second, we consider the single-turn dialogues from the helpful-base subset
 139 of the Anthropic-HH dataset [Bai et al., 2022b]. Specifically, this dataset consists of *open-ended*
 140 instructions with a collection of responses ranging from ‘Which coffee bean is better for a morning
 141 roast?’ to ‘How do I attract more hummingbirds in my yard?’.

142 Both these datasets have a train and test split where each instance consists of an instruction and a
 143 pair of responses $\mathcal{D} = \{(I_i, R_i^1, R_i^2)\}_{i=1}^n$ where n is the dataset size. In this work, we collect AI
 144 and human feedback on the instruction-response data from their train split and filter instances where
 145 instructions are repeated. We can directly compare the two responses for the fixed instruction and
 146 construct a ranking feedback dataset $\mathcal{D}_C = \{(I_i, R_i^1, R_i^2, c(I_i, R_i^1, R_i^2))\}$. To acquire preferences
 147 jointly over the instruction-response pairs, we select one of the responses, at random, from every
 148 instance of \mathcal{D} to construct $\mathcal{D}_S = \{(I_i, R_i)\}$ where $R_i \in \{R_i^1, R_i^2\}$. Subsequently, we create the
 149 joint instruction-response pairs by matching every instance $(I_i, R_i) \in \mathcal{D}_S$ with another instance
 150 $(I_j, R_j) \in \mathcal{D}_S$ to get $\mathcal{D}_H = \{(I_i, R_i, I_j, R_j, h(I_i, R_i, I_j, R_j))\}$ of the same size as \mathcal{D}_S and \mathcal{D}_C . In
 151 §4, we will utilize \mathcal{D}_S to SFT the base model, and \mathcal{D}_C and \mathcal{D}_H as preference datasets for LLM
 152 alignment. We provide the dataset statistics in Appendix §D.

153 3.2 Feedback from AI and Humans

154 **Feedback from AI.** Prior work [Dubois et al., 2023, Bai et al., 2022b] has shown that AI feedback
 155 can be leveraged to align language models to generate helpful and harmless responses to unseen
 156 instructions. In addition, acquiring AI feedback at large-scale is more accessible and cheaper in
 157 comparison to human feedback. To this end, we collect feedback over a pair of responses for a fixed
 158 instruction, and joint instruction-response pairs without identical instructions from **GPT-3.5-Turbo-**
 159 **0125** (ChatGPT). The choice of ChatGPT was motivated by its affordability (e.g., output tokens from
 160 ChatGPT are $50\times$ cheaper than GPT-4).

161 To collect ranking feedback over a pair of responses for a fixed instruction, we prompt ChatGPT to
 162 choose a response. To mitigate any bias from the ordering of the two responses, we run two queries
 163 for all comparisons. When the ChatGPT preferences flip by flipping the order of the two responses,
 164 then we consider it a tie, similar to [Bansal et al., 2023, Bitton et al., 2023]. Specifically, the AI
 165 is instructed to provide its preference based on the accuracy, coherence, and harmlessness of the
 166 responses.

167 To collect AI preferences jointly over the instruction-response pairs, we prompt ChatGPT to decide
 168 the response that better answers its corresponding instruction. Similar to the previous scenario,

169 we run two queries for all comparisons to mitigate any ordering bias and provide guidelines to
 170 choose the response that is more accurate, coherent, and harmless. We collected approximately 50K
 171 comparisons across both feedback acquisition protocols for the summarization and Anthropic-Helpful
 172 dataset, at a cost of \$100. We provide the AI prompts in Appendix §J.

173 **Feedback from Humans.** In this work, we also collect human preferences for 2000 comparisons
 174 over summarization and Anthropic-Helpful dataset. Such a data is useful for providing insights into
 175 the human behavior under different preference acquisition protocols (§3.4). In addition, this data aids
 176 in agreement between the ChatGPT and human decisions.

177 Specifically, we ask two annotators to assign a chosen response or choose ‘equal’ after comparing the
 178 quality of the responses along the same dimensions as ChatGPT guidelines. The human annotations
 179 were collected from Amazon Mechanical Turk (AMT) from the participants that passed a preliminary
 180 qualification exam. In total, we spent \$720 on human feedback acquisition. We provide the screenshot
 181 of the annotation UI in Appendix §K.

Dataset	Ranking Protocol	Human-Human	Human-AI
TL;DR	Conditional	69%	63%
Anthropic-Helpful		70.1%	72%
TL;DR	Joint (Non-Identical)	62%	60%
Anthropic-Helpful		68.8%	71%
Average		67.5%	66.5%

Table 1: Agreement analysis between within human annotators and gold human feedback and AI (ChatGPT) feedback. We perform the agreement calculations for the two ranking protocols: (a) conditional rankings, and (b) joint preferences where instructions are non-identical. In addition, we assess the agreement rates over the two datasets: (a) TL;DR and (b) Anthropic-helpful dataset.

182 3.3 Agreement Analysis

183 We present the annotator agreement scores in Table 1. We find that the average agreement is 67.5%
 184 and 66.5% between the human-human and human-AI annotators, respectively. Furthermore, we find
 185 that the average agreement score between humans for conditional (identical instruction) setup is
 186 69.5% over TLDR and Anthropic-Helpfulness. Similarly, the average inter-rater agreement is 68% for
 187 the joint (non-identical instruction-response pairs) setup on the same datasets. Our agreement scores
 188 are close to the agreement scores in prior work [Li et al., 2023, Bansal et al., 2023]. Interestingly,
 189 the agreement scores vary based on the underlying distribution of the instruction-response pairs and
 190 the choice of ranking protocol. Overall, our results highlight that humans and AI can provide rich
 191 feedback in both conditional and joint setup with acceptable agreement.

192 3.4 Interplay Analysis

193 **Setup.** Here, we aim to study the interaction between the conditional rankings and joint rankings
 194 over non-identical instructions. Formally, each instruction-response pair (I_i, R_i^x) from the conditional
 195 pairwise feedback dataset \mathcal{D}_C where $x \in \{1, 2\}$ can be assigned a preference $\mathcal{P}_C(I_i, R_i^x)$ among
 196 {‘chosen’, ‘reject’, ‘equal’}. For instance, $\mathcal{P}_C(I_i, R_i^1) = \text{‘chosen’}$ and $\mathcal{P}_C(I_i, R_i^2) = \text{‘reject’}$ if
 197 the response R_i^2 is rejected in the dataset \mathcal{D}_C i.e., $c(I_i, R_i^1, R_i^2) = R_i^1$. Similarly, we can assign a
 198 preference $\mathcal{P}_H(I_i, R_i)$ among {‘chosen’, ‘reject’, ‘equal’} to an instruction-response pair (I_i, R_i)
 199 from the joint preference dataset \mathcal{D}_H . For instance, $\mathcal{P}_H(I_i, R_i) = \text{‘chosen’}$ and $\mathcal{P}_H(I_j, R_j) =$
 200 ‘reject’ where $i! = j$ if the instruction-response pair (I_i, R_i) is chosen in the dataset \mathcal{D}_H i.e.,
 201 $h(I_i, R_i, I_j, R_j) = (I_i, R_i)$.

202 To study the interplay between the preference protocols, we assess $\mathcal{P}_C(I_i, R_i)$, $\mathcal{P}_C(I_j, R_j)$,
 203 $\mathcal{P}_H(I_i, R_i)$ and $\mathcal{P}_H(I_j, R_j)$ for all $(I_i, R_i, I_j, R_j) \in \mathcal{D}_H$. Here, if $\mathcal{P}_H(I_i, R_i) = \text{‘chosen’}$ then
 204 $\mathcal{P}_H(I_j, R_j) = \text{‘reject’}$. For instance, if $\mathcal{P}_C(I_i, R_i) = \text{‘chosen’}$ and $\mathcal{P}_C(I_j, R_j) = \text{‘chosen’}$ then it
 205 implies that the annotators can reason about the joint preferences over a pair of instruction-response
 206 pairs that are originally preferred under the conditional ranking feedback protocol. We quantitatively
 207 study the interplay between the two ranking-based feedback from AI and Human annotators over
 208 summarization and open-ended Anthropic-Helpful datasets.

Data (Annotator)	Decisive	Indecisive
TL;DR (AI)	63.7%	36.2%
TL;DR (Human)	73.8%	25.7%
Anthropic-Helpful (AI)	68.5%	31.5%
Anthropic-Helpful (Human)	77.9%	22.0%
Average	71.0%	29.0%

Table 2: Results for the preferences acquired jointly over the instruction-response pairs where both the responses were either chosen or rejected under the conditional rankings protocol. Here, *decisive* implies that the annotators could assign a preference to one instruction-response pair over the other. In total, we compare 48K and 1K annotations from the AI and humans, respectively.

Data (Annotator)	C > R	C < R	Indecisive
TL;DR (AI)	53.3%	14.3%	30.4%
TL;DR (Human)	41.6%	22.2%	36.1%
Anthropic-Helpful (AI)	54.5%	17.6%	27.8%
Anthropic-Helpful (Human)	57.1%	21.4%	21.4%
Average	52.0%	19.0%	29.0%

Table 3: Results for the preferences acquired jointly over the instruction-response pairs where one of the instruction-response pair was chosen (C) and the other pair was rejected (R) under the conditional rankings. Here, $C < R$ implies that the instruction-response pair that was rejected under conditional rankings is actually preferred over an instruction-response pair that was rejected under the conditional rankings. In total, we compare 48K and 1K annotations from the AI and humans, respectively.

209 **Results.** We present the results for the interaction analysis in Table 3 and Table 2. In Table 2, we
210 study the joint preferences over the instruction-response pairs (I_i, R_i, I_j, R_j) where the individual
211 instruction and response data is either *chosen* or *rejected* in the conditional feedback protocol (e.g.,
212 $\mathcal{P}_C(I_z, R_z) = \text{'chosen'}$ for $z \in \{i, j\}$). Interestingly, we find that the annotators can assign a decisive
213 preference (e.g., $(I_i, R_i) > (I_j, R_j)$) in 71% of the joint comparisons. While we observe that the
214 annotators assign a 'tie' to 29% of the comparisons. This highlights the existence of valid preference
215 decisions that remained obfuscated in the traditional approach for ranking-based feedback acquisition.

216 In Table 3, we study the joint preference over the instruction-response pairs (I_i, R_i, I_j, R_j) where one
217 of them is *chosen* and the other is *rejected* in the conditional feedback protocol (e.g., $\mathcal{P}_C(I_i, R_i) =$
218 '*chosen*' and $\mathcal{P}_C(I_j, R_j) = \text{'reject'}$). To our surprise, we find that the annotators do not prefer the
219 instruction-response pair that was chosen under the conditional feedback protocol in 48% of the
220 comparisons. Specifically, there are 19% of the comparisons where rejected pair (R) is preferred
221 over the chosen pair (C) and 28% of the comparisons where the annotators considered the pair
222 equally good or bad. This highlights that both human and AI annotators' perceptions of preferred and
223 non-preferred data depends on the context of the comparisons, indicating that feedback acquisition is
224 a multifaceted phenomenon.

Method	TL;DR				Anthropic-Helpful			
	T = 0.001	T = 0.5	T = 1.0	Average	T = 0.001	T = 0.5	T = 1.0	Average
SFT	46.6	44.9	39.8	43.8	59.1	56.2	56.8	57.4
DPO Rafailov et al. [2024]	66.5	67.0	69.5	67.7	73.5	72	69.5	71.7
KTO Ethayarajh et al. [2024]	71.8	71.9	70.6	71.4	72.8	72.9	68.8	71.5
DOVE (Ours)	72.7	71.9	74.2	72.9	76.3	74.5	74.1	75.0

Table 4: Results for aligning LLMs with the DOVE preference optimization objective. We compare the win-rate against the gold responses of the supervised finetuned (SFT), DPO-aligned and DOVE-aligned LLM on the (a) TL;DR summarization and (b) the Anthropic-Helpful datasets. In our experiments, we utilize ChatGPT to compare the model responses with the gold responses. We generate model responses for three sampling temperatures. The results are averaged over three runs of the preference optimization objectives.

225 **Qualitative Examples.** To probe the reasoning paths of the human annotators used for decision
226 making, we ask them to provide brief explanations for their feedback decisions regarding a few
227 conditional and joint preferences. We provide a list of qualitative examples consisting of instructions,

228 responses, and respective preferences in Appendix §G. In Figure 3, we discovered that human
229 annotators provided decisive feedback when comparing instruction-response pairs, basing their
230 decisions on the accuracy of the responses. In Figure 6, we find that the human annotators preferred a
231 instruction-summary pair, that was rejected under the conditional preference, because it provides a
232 fuller picture of the original reddit post. In summary, we expose the multi-faceted reasoning paths
233 of humans in joint instruction-response feedback acquisition that would have been concealed in the
234 conditional feedback acquisition paradigm.

235 4 LLM Alignment

236 In the previous sections, we show that the humans and AI are capable of providing ranking-based
237 feedback for a pair of responses for identical and non-identical instructions. Here, we aim to study
238 how to leverage joint and conditional feedback data to align large language models effectively.

239 4.1 Setup

240 Here, we aim to align Mistral-7B [Jiang et al., 2023], a strong base LLM for its model capacity.
241 We experiment with two datasets that exhibit diverse characteristics: (a) TL;DR dataset where the
242 instruction is to summarize Reddit posts, and (b) open-ended dialogues from Anthropic-Helpful
243 dataset (§3.1). In particular, we collect a conditional preference data \mathcal{D}_C and joint preference data for
244 non-identical instructions \mathcal{D}_H of similar data sizes from ChatGPT. Then, we convert the conditional
245 preference data into an instruction-response data for supervised finetuning \mathcal{D}_{SFT} .

246 First, we supervise finetune the entire base LLM model parameters with the SFT dataset to ensure that
247 the preference data is in-policy for the alignment algorithms [Rafailov et al., 2023]. Subsequently,
248 we apply DPO algorithm on the SFT model using the conditional preference data for 10 epochs
249 and 5 epochs for the summarization and Anthropic-helpful data, respectively. Specifically, we use
250 low-rank adaptation [Hu et al., 2021] of SFT model during DPO alignment. The DPO optimization
251 was trained on a single GPU Nvidia A6000 with a batch size of 32.

252 We note that our proposed DOVE algorithm can utilize both the conditional preferences and joint
253 preference with non-identical context. It is because the conditional preferences can be viewed as joint
254 preferences with identical context. As a result, we train the base LLM with DOVE algorithm after
255 merging conditional and joint preferences data $\mathcal{D}_M = \mathcal{D}_C \cup \mathcal{D}_H$. We keep the hyperparameters (e.g.,
256 β), number of epochs, and the batch size identical to the DPO algorithm. In our experiments, we
257 also train DOVE algorithm on the joint preferences with non-identical instructions and highlight their
258 usefulness for LLM alignment. We provide more details on training setup in Appendix §H.

259 Post-alignment, we evaluate the aligned model responses against the gold responses in the dataset’s
260 test split. Specifically, both datasets come with a human-preferred response for an instruction, which
261 is treated as the gold response. We utilize ChatGPT to compare model and gold responses to decide
262 on the preferred response or a tie. Finally, we report the win-rate of the model responses as the
263 evaluation metric for 500 unseen instructions [Rafailov et al., 2023].

264 4.2 Results

265 We compare the performance of the SFT, DPO, KTO, and DOVE aligned models in Table 4. In
266 particular, we report the win-rate against the gold responses for the model generated responses for
267 sampling temperatures $T \in \{0.001, 0.5, 1.0\}$.

268 **DOVE outperforms SFT model.** We find that the DOVE achieves high win-rates across all sampling
269 temperatures. Specifically, we observe that DOVE outperforms the SFT model by 29.1% and 18% on
270 the close-ended summarization and open-ended dialogue dataset, respectively, averaged across the
271 sampling temperatures. This indicates that DOVE can utilize the diverse set of feedback from the
272 conditional and joint preferences to align LLMs.

273 **DOVE outperforms DPO and KTO.** Further, we aim to understand whether DOVE is able to tease
274 out useful feedback signals from the combination of the conditional preferences and joint preferences
275 over instruction-response pairs. Surprisingly, we find that DOVE outperforms DPO by 5.2% and
276 3.3% win-rate points on the summarization and helpfulness datasets, respectively. In addition, the

277 performance of DOVE is better than DOVE across all the sampling temperatures. This highlights that
 278 one can improve the alignment of the LLMs by leveraging novel preference acquisition paths without
 279 collecting new instruction-response data. We observe the similar trends in comparison to KTO. In
 280 Appendix F, we show that DOVE outperforms DPO on a broad set of instructions from AlpacaEval
 281 Li et al. [2023] as well. Hence, our results indicate that DOVE is a robust alignment algorithm that
 282 can elicit high-quality outputs by learning from diverse ranking-based preferences.

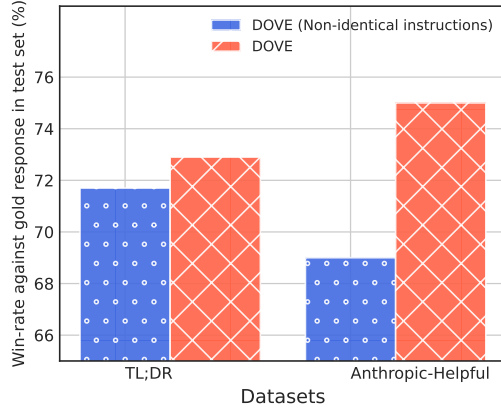


Figure 2: Win-rate against the gold response in the TL;DR and Anthropic-Helpful datasets averaged over three sampling temperatures. We study the impact of the joint preferences over non-identical instructions using DOVE.

283 **Impact of Joint Preferences over Non-Identical Instructions.** Here, we aim to understand the sole
 284 impact of joint preferences acquired over non-identical instructions on the performance of the DOVE
 285 algorithm. To do so, we train DOVE algorithm with joint feedback data \mathcal{D}_H only. We present the
 286 results averaged across the three sampling temperatures in Figure 2. We find that training with joint
 287 preferences over non-identical instructions achieves 71.7% and 69% win-rate on the summarization
 288 and anthropic-helpful datasets, respectively. This indicates that it is possible to align LLMs with just
 289 joint preferences over instruction-response data *without* any conditional preferences too. Furthermore,
 290 this highlights that the feedback paths exposed in our setup are robust and effective for alignment.

291 **Impact of Dataset Size.** In the main experiments, we demonstrated that DOVE can learn effectively
 292 from a combination of conditional preferences (i.e., 100% of the conditional rankings) and joint
 293 preferences over non-identical instructions (of the same size as the conditional preferences). To assess
 294 the impact of dataset size, we trained DOVE using a 50:50 mix of conditional and joint preferences for
 295 the TL;DR dataset, with a fixed total size as that of conditional. Our results show that DOVE achieves
 296 a win rate of 71.9%, outperforming DPO, which was trained on only the conditional preference
 297 dataset of the same size, by 4.2 percentage points. Additionally, we demonstrate in Appendix §I that
 298 training with joint preferences scales with the amount of feedback data using the DOVE algorithm.

299 5 Conclusion

300 In this work, we propose a framework that elicits preferences jointly over instruction-response pairs.
 301 Further, we find that the joint preference optimization uncovers new paths of human reasoning that
 302 remain obscured in the traditional approach. Additionally, we propose DOVE, a novel preference
 303 optimization objective for aligning LLMs. In our experiments, we show that it outperforms DPO on
 304 summarization and dialogue datasets. We note that the number of joint preferences over instruction-
 305 response data scales quadratically with the number of instances in the instruction-response dataset.
 306 Therefore, identifying the most informative joint comparisons for robust LLM alignment represents
 307 a relevant area for future research. While traditional LLM evaluation has focused on conditional
 308 rankings, LLM evaluation through joint rankings would be an important future work.

309 **References**

- 310 Afra Amini, Tim Vieira, and Ryan Cotterell. Direct preference optimization with an offset. *arXiv*
311 *preprint arXiv:2402.10571*, 2024.
- 312 Anthropic. Introducing claude. 2023. URL [https://www.anthropic.com/index/
313 introducing-claude](https://www.anthropic.com/index/introducing-claude).
- 314 Anthropic. The claude 3 model family: Opus, sonnet, haiku. 2024. URL [https://api.
315 semanticscholar.org/CorpusID:268232499](https://api.semanticscholar.org/CorpusID:268232499).
- 316 Amanda Askell, Yuntao Bai, Anna Chen, Dawn Drain, Deep Ganguli, Tom Henighan, Andy Jones,
317 Nicholas Joseph, Ben Mann, Nova DasSarma, et al. A general language assistant as a laboratory
318 for alignment. *arXiv preprint arXiv:2112.00861*, 2021.
- 319 Mohammad Gheshlaghi Azar, Mark Rowland, Bilal Piot, Daniel Guo, Daniele Calandriello, Michal
320 Valko, and Rémi Munos. A general theoretical paradigm to understand learning from human
321 preferences. *arXiv preprint arXiv:2310.12036*, 2023.
- 322 Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain,
323 Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with
324 reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022a.
- 325 Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna
326 Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness
327 from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022b.
- 328 Hritik Bansal, John Dang, and Aditya Grover. Peering through preferences: Unraveling feedback
329 acquisition for aligning large language models. *arXiv preprint arXiv:2308.15812*, 2023.
- 330 Yonatan Bitton, Hritik Bansal, Jack Hessel, Rulin Shao, Wanrong Zhu, Anas Awadalla, Josh Gardner,
331 Rohan Taori, and Ludwig Schimdt. Visit-bench: A benchmark for vision-language instruction
332 following inspired by real-world use, 2023.
- 333 Ralph Allan Bradley and Milton E Terry. Rank analysis of incomplete block designs: I. the method
334 of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.
- 335 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
336 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
337 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 338 Jan Clusmann, Fiona R Kolbinger, Hannah Sophie Muti, Zunamys I Carrero, Jan-Niklas Eckardt,
339 Narmin Ghaffari Laleh, Chiara Maria Lavinia Löffler, Sophie-Caroline Schwarzkopf, Michaela
340 Unger, Gregory P Veldhuizen, et al. The future landscape of large language models in medicine.
341 *Communications medicine*, 3(1):141, 2023.
- 342 Commoncrawl. Common crawl. <https://commoncrawl.org>, 2024. Accessed on March 23, 2024.
- 343 Mike Conover, Matt Hayes, Ankit Mathur, Jianwei Xie, Jun Wan, Sam Shah, Ali Ghodsi, Patrick
344 Wendell, Matei Zaharia, and Reynold Xin. Free dolly: Introducing the world’s first truly
345 open instruction-tuned llm, 2023. URL [https://www.databricks.com/blog/2023/04/12/
346 dolly-first-open-commercially-viable-instruction-tuned-llm](https://www.databricks.com/blog/2023/04/12/dolly-first-open-commercially-viable-instruction-tuned-llm).
- 347 Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu,
348 and Maosong Sun. Ultrafeedback: Boosting language models with high-quality feedback. *arXiv*
349 *preprint arXiv:2310.01377*, 2023.
- 350 Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi Zheng, Shengding Hu, Zhiyuan Liu, Maosong
351 Sun, and Bowen Zhou. Enhancing chat language models by scaling high-quality instructional
352 conversations. *arXiv preprint arXiv:2305.14233*, 2023.
- 353 Yann Dubois, Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin,
354 Percy Liang, and Tatsunori B Hashimoto. AlpacaFarm: A simulation framework for methods that
355 learn from human feedback. *arXiv preprint arXiv:2305.14387*, 2023.

- 356 Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Kto: Model
357 alignment as prospect theoretic optimization. *arXiv preprint arXiv:2402.01306*, 2024.
- 358 Iason Gabriel. Artificial intelligence, values, and alignment. *Minds and machines*, 30(3):411–437,
359 2020.
- 360 Xinyang Geng, Arnav Gudibande, Hao Liu, Eric Wallace, Pieter Abbeel, Sergey Levine, and
361 Dawn Song. Koala: A dialogue model for academic research. Blog post, April 2023. URL
362 <https://bair.berkeley.edu/blog/2023/04/03/koala/>.
- 363 Jiwoo Hong, Noah Lee, and James Thorne. Orpo: Monolithic preference optimization without
364 reference model, 2024. URL <https://arxiv.org/abs/2403.07691>.
- 365 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
366 and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021.
- 367 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,
368 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al.
369 Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- 370 Maurice G Kendall and B Babington Smith. On the method of paired comparisons. *Biometrika*, 31
371 (3/4):324–345, 1940.
- 372 Hannah Rose Kirk, Alexander Whitefield, Paul Röttger, Andrew Bean, Katerina Margatina, Juan
373 Ciro, Rafael Mosquera, Max Bartolo, Adina Williams, He He, et al. The prism alignment project:
374 What participatory, representative and individualised human feedback reveals about the subjective
375 and multicultural alignment of large language models. *arXiv preprint arXiv:2404.16019*, 2024.
- 376 Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu,
377 Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, et al. Rewardbench: Evaluating reward models
378 for language modeling. *arXiv preprint arXiv:2403.13787*, 2024.
- 379 Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy
380 Liang, and Tatsunori B. Hashimoto. AlpacaEval: An automatic evaluator of instruction-following
381 models. https://github.com/tatsu-lab/alpaca_eval, 2023.
- 382 Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan
383 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. *arXiv preprint*
384 *arXiv:2305.20050*, 2023.
- 385 Rensis Likert. A technique for the measurement of attitudes. *Archives of psychology*, 1932.
- 386 Tianqi Liu, Yao Zhao, Rishabh Joshi, Misha Khalman, Mohammad Saleh, Peter J Liu, and Jialu Liu.
387 Statistical rejection sampling improves preference optimization. *arXiv preprint arXiv:2309.06657*,
388 2023.
- 389 Tianqi Liu, Zhen Qin, Junru Wu, Jiaming Shen, Misha Khalman, Rishabh Joshi, Yao Zhao, Moham-
390 mad Saleh, Simon Baumgartner, Jialu Liu, et al. Lipo: Listwise preference optimization through
391 learning-to-rank. *arXiv preprint arXiv:2402.01878*, 2024.
- 392 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint*
393 *arXiv:1711.05101*, 2017.
- 394 Yujie Lu, Xianjun Yang, Xiujun Li, Xin Eric Wang, and William Yang Wang. Llmscore: Unveiling
395 the power of large language models in text-to-image synthesis evaluation. *Advances in Neural*
396 *Information Processing Systems*, 36, 2024.
- 397 Yu Meng, Mengzhou Xia, and Danqi Chen. Simpo: Simple preference optimization with a reference-
398 free reward, 2024. URL <https://arxiv.org/abs/2405.14734>.
- 399 Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher
400 Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. Webgpt: Browser-assisted
401 question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021.
- 402 OpenAI. Gpt-4 technical report, 2023.

403 A Ross Otto, Sean Devine, Eric Schulz, Aaron M Bornstein, and Kenway Louie. Context-dependent
404 choice and evaluation in real-world consumer behavior. *Scientific reports*, 12(1):17744, 2022.

405 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
406 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow
407 instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:
408 27730–27744, 2022.

409 Arka Pal, Deep Karkhanis, Samuel Dooley, Manley Roberts, Siddartha Naidu, and Colin White.
410 Smaug: Fixing failure modes of preference optimisation with dpo-positive. *arXiv preprint*
411 *arXiv:2402.13228*, 2024.

412 Ryan Park, Rafael Rafailov, Stefano Ermon, and Chelsea Finn. Disentangling length from quality in
413 direct preference optimization. *arXiv preprint arXiv:2403.19159*, 2024.

414 Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli,
415 Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. The refinedweb
416 dataset for falcon llm: outperforming curated corpora with web data, and web data only. *arXiv*
417 *preprint arXiv:2306.01116*, 2023.

418 Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. Instruction tuning with
419 gpt-4. *arXiv preprint arXiv:2304.03277*, 2023.

420 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language
421 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.

422 Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D Manning, and Chelsea
423 Finn. Direct preference optimization: Your language model is secretly a reward model. *arXiv*
424 *preprint arXiv:2305.18290*, 2023.

425 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea
426 Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances*
427 *in Neural Information Processing Systems*, 36, 2024.

428 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi
429 Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text
430 transformer. *Journal of machine learning research*, 21(140):1–67, 2020.

431 Jérémy Scheurer, Jon Ander Campos, Tomasz Korbak, Jun Shern Chan, Angelica Chen, Kyunghyun
432 Cho, and Ethan Perez. Training language models with language feedback at scale. *arXiv preprint*
433 *arXiv:2303.16755*, 2023.

434 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
435 optimization algorithms, 2017.

436 Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin Schwenk, David Atkinson, Russell Authur,
437 Ben Bogin, Khyathi Chandu, Jennifer Dumas, Yanai Elazar, et al. Dolma: An open corpus of three
438 trillion tokens for language model pretraining research. *arXiv preprint arXiv:2402.00159*, 2024.

439 Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford,
440 Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. *Advances in*
441 *Neural Information Processing Systems*, 33:3008–3021, 2020.

442 Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy
443 Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model.
444 https://github.com/tatsu-lab/stanford_alpaca, 2023.

445 Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu
446 Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable
447 multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.

448 Teknium. Openhermes 2.5: An open dataset of synthetic data for generalist llm assistants, 2023.
449 URL <https://huggingface.co/datasets/teknium/OpenHermes-2.5>.

- 450 Louis L Thurstone. A law of comparative judgment. In *Scaling*, pages 81–92. Routledge, 2017.
- 451 Shubham Toshniwal, Ivan Moshkov, Sean Narenthiran, Daria Gitman, Fei Jia, and Igor Git-
452 man. Openmathinstruct-1: A 1.8 million math instruction tuning dataset. *arXiv preprint*
453 *arXiv:2402.10176*, 2024.
- 454 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
455 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and
456 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- 457 Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada,
458 Shengyi Huang, Leandro von Werra, Clémentine Fourier, Nathan Habib, Nathan Sarrazin, Omar
459 Sanseviero, Alexander M. Rush, and Thomas Wolf. Zephyr: Direct distillation of llm alignment,
460 2023.
- 461 Michael Völske, Martin Potthast, Shahbaz Syed, and Benno Stein. Tl; dr: Mining reddit to learn
462 automatic summarization. In *Proceedings of the Workshop on New Frontiers in Summarization*,
463 pages 59–63, 2017.
- 464 Leandro von Werra, Younes Belkada, Lewis Tunstall, Edward Beeching, Tristan Thrush, Nathan
465 Lambert, and Shengyi Huang. Trl: Transformer reinforcement learning. <https://github.com/huggingface/trl>, 2020.
- 467 Chaoqi Wang, Yibo Jiang, Chenghao Yang, Han Liu, and Yuxin Chen. Beyond reverse kl: Gen-
468 eralizing direct preference optimization with diverse divergence constraints. *arXiv preprint*
469 *arXiv:2309.16240*, 2023a.
- 470 Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei,
471 Anjana Arunkumar, Arjun Ashok, Arut Selvan Dhanasekaran, Atharva Naik, David Stap, et al.
472 Super-naturalinstructions: Generalization via declarative instructions on 1600+ nlp tasks. *arXiv*
473 *preprint arXiv:2204.07705*, 2022.
- 474 Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack Hessel, Tushar Khot, Khyathi Raghavi Chandu,
475 David Wadden, Kelsey MacMillan, Noah A Smith, Iz Beltagy, et al. How far can camels go?
476 exploring the state of instruction tuning on open resources. *arXiv preprint arXiv:2306.04751*,
477 2023b.
- 478 Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and
479 Hannaneh Hajishirzi. Self-instruct: Aligning language models with self-generated instructions,
480 2023c.
- 481 Shijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabravolski, Mark Dredze, Sebastian Gehrmann, Prabhan-
482 jan Kambadur, David Rosenberg, and Gideon Mann. Bloomberggpt: A large language model for
483 finance. *arXiv preprint arXiv:2303.17564*, 2023a.
- 484 Zeqiu Wu, Yushi Hu, Weijia Shi, Nouha Dziri, Alane Suhr, Prithviraj Ammanabrolu, Noah A Smith,
485 Mari Ostendorf, and Hannaneh Hajishirzi. Fine-grained human feedback gives better rewards for
486 language model training. *arXiv preprint arXiv:2306.01693*, 2023b.
- 487 Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin
488 Jiang. Wizardlm: Empowering large language models to follow complex instructions. *arXiv*
489 *preprint arXiv:2304.12244*, 2023.
- 490 Da Yin, Xiao Liu, Fan Yin, Ming Zhong, Hritik Bansal, Jiawei Han, and Kai-Wei Chang. Dinosaur:
491 A dynamic growth paradigm for instruction-tuning data curation, 2023.
- 492 Yueqin Yin, Zhendong Wang, Yi Gu, Hai Huang, Weizhu Chen, and Mingyuan Zhou. Relative
493 preference optimization: Enhancing llm alignment through contrasting responses across identical
494 and diverse prompts. *arXiv preprint arXiv:2402.10958*, 2024.
- 495 Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhenguo
496 Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions for
497 large language models. *arXiv preprint arXiv:2309.12284*, 2023.

- 498 Yao Zhao, Mikhail Khalman, Rishabh Joshi, Shashi Narayan, Mohammad Saleh, and Peter J Liu.
499 Calibrating sequence likelihood improves conditional language generation. In *The Eleventh*
500 *International Conference on Learning Representations*, 2022.
- 501 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Tianle Li, Siyuan Zhuang, Zhanghao Wu, Yonghao
502 Zhuang, Zhuohan Li, Zi Lin, Eric Xing, et al. Lmsys-chat-1m: A large-scale real-world llm
503 conversation dataset. *arXiv preprint arXiv:2309.11998*, 2023a.
- 504 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
505 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and
506 chatbot arena. *arXiv preprint arXiv:2306.05685*, 2023b.

507 **A Background**

508 In this work, our aim is to align language models to generate outputs that are preferred by humans
 509 across various dimensions such as helpfulness and coherence. The process of aligning a base model,
 510 which is pretrained on a large corpus of text [Commoncrawl, 2024, Raffel et al., 2020, Soldaini
 511 et al., 2024, Penedo et al., 2023], involves multiple steps: (a) instruction-response data collection, (b)
 512 supervised fine-tuning, (c) preference data acquisition, and (d) deployment of an alignment algorithm.
 513 The instruction-response data can be either hand-crafted by humans [Conover et al., 2023, Wang
 514 et al., 2022] or generated by machines [Taori et al., 2023, Tunstall et al., 2023]. Subsequently, the
 515 base model undergoes supervised fine-tuning (SFT) on the instruction-response pairs [Zheng et al.,
 516 2023b, Wang et al., 2023c, 2022, Peng et al., 2023, Xu et al., 2023, Geng et al., 2023, Yin et al.,
 517 2023, Wang et al., 2023b, Yu et al., 2023, Toshniwal et al., 2024]. Following SFT, feedback data is
 518 acquired under a specific acquisition protocol (e.g., rankings) from the annotators (§A.1). Finally, an
 519 alignment algorithm trains the SFT model on the feedback data (§A.2).

520 **A.1 Ranking Feedback Acquisition Protocol**

521 Assume a supervised finetuned language model p_{sft} that is capable of responding to user instruc-
 522 tions (e.g., imperative tasks or questions). The goal of alignment is to ensure that the SFT model
 523 generates high-quality outputs, preferred by humans. To do so, we consider a set of instructions
 524 $\mathcal{I} = \{I_1, \dots, I_n\}$ where n is the number of instructions. Further, we consider a set of responses
 525 $\{R_j^1, R_j^2, \dots, R_j^k\}$ where k is the number of responses for each of the instruction $I_j \in \mathcal{I}$. This forms
 526 a dataset of instructions and their corresponding responses, $\mathcal{D} = \{(I_j, R_j^1, R_j^2, \dots, R_j^k)\}$.³ Next, we
 527 acquire conditional ranking-based feedback over the collected instruction-response data.

528 Under this feedback acquisition protocol, the annotator selects a *chosen* and *rejected* response from
 529 $\{R_j^x, R_j^y\}$ conditioned on the instruction I_j where $x, y \in \{1, 2, \dots, k\}$. The preference decision by
 530 the annotator is based on the perceived quality of the responses along various dimensions such as
 531 helpfulness (accuracy), coherence (grammar), and harmlessness (safety).

532 Formally, the annotator assigns an instruction-conditioned ranking feedback $c(I_j, R_j^x, R_j^y) \in$
 533 $\{R_j^x, R_j^y, \text{Equal}\}$ where ‘Equal’ indicates that both responses are perceived equally good or bad.
 534 If $c(I_j, R_j^x, R_j^y) = R_j^x$, this implies that the response R_j^x is the chosen response while the R_j^y is the
 535 rejected response by the annotator. As a result, the ranking protocol creates a conditional pairwise
 536 feedback data $\mathcal{D}_C = \{(I_j, R_j^x, R_j^y, c(I_j, R_j^x, R_j^y))\}$. Next, we apply an alignment algorithm on this
 537 data to elicit human-preferred responses from the LLM.

538 **A.2 Alignment Algorithms**

539 Rafailov et al. [2023] introduced direct preference optimization (DPO) that can align a language
 540 model without utilizing on an external reward model. Specifically, DPO requires that feedback
 541 data should consist of conditional preferences between a pair of responses for a given instruction.
 542 Additionally, the algorithm assumes a preference dataset \mathcal{D}_C and the reference model p_{ref} which
 543 is usually the supervised finetuned language model p_{sft} . Specifically, it aims to train an aligned
 544 model p_θ using an optimization objective that upweights the conditional probability of the chosen
 545 response $p_\theta(R_j^w|I_j)$ over the rejected response $p_\theta(R_j^\ell|I_j)$ where R_j^w and R_j^ℓ are the chosen and
 546 rejected response, respectively. Formally, the optimization objective for DPO, $\mathcal{L}_{\text{DPO}}(\theta; \mathcal{D}_C, \beta, p_{\text{ref}})$
 547 minimizes the expectation over $(I_j, R_j^w, R_j^\ell) \sim \mathcal{D}_C$:

$$\mathbb{E} \left[\log \left(\sigma \left(\beta \log \frac{p_\theta(R_j^w|I_j)}{p_{\text{ref}}(R_j^w|I_j)} - \beta \log \frac{p_\theta(R_j^\ell|I_j)}{p_{\text{ref}}(R_j^\ell|I_j)} \right) \right) \right] \quad (2)$$

548 where σ denotes the sigmoid function and β is a hyperparameter. Post-alignment, the model generates
 549 high-quality outputs for unseen instructions.

³We will drop the iterator over j when defining the dataset for the ease of notation.

550 B Related Work

551 **Alignment using Reinforcement Learning.** Aligning LLMs with human preferences using re-
552 inforcement learning is widely adopted to ensure LLMs follow user intents without being harmful
553 Ouyang et al. [2022]. This alignment is usually done by first optimizing for a reward model on
554 preference data [Bradley and Terry, 1952, Likert, 1932, Bansal et al., 2023], followed by aligning
555 the LLMs distribution that maximizes the learned reward model using Reinforcement Learning
556 (RLHF) Schulman et al. [2017], Ouyang et al. [2022], with optional Divergence penalty Wang et al.
557 [2023a] to avoid deviating from the reference policy. Additionally, Dubois et al. [2023], Lu et al.
558 [2024], Zheng et al. [2023b] observe that preferences from LLMs can also be used for alignments
559 motivating Reinforcement Learning through AI feedback (RLAIF). Contrary to prior work that
560 collect preferences as conditional rankings, we emphasize that preference acquisition is a complex
561 phenomenon and elicit joint preferences over instruction-response data.

562 **Reward Free Policy Alignment.** Rafailov et al. [2024] introduced Direct Preference Optimization
563 (DPO) that optimizes directly within the model parameter space, hence eliminating the reward
564 modeling step. Liu et al. [2024] extends this framework where instead of two responses, alignment
565 is done over the list of responses while Liu et al. [2023] improves DPO using statistical rejection
566 sampling. Amini et al. [2024] provides an offset in the DPO objective to increase the margins and
567 Pal et al. [2024] suggests adding an explicit penalty term to avoid a reduction in the likelihood of
568 preferred pairs over the DPO training. Recent variants of DPO such as SimPO [Meng et al., 2024]
569 alleviates the need of reference policy in the objective. Contrary to our work where we compare
570 the joint distributions, Yin et al. [2024] proposes RPO that compares the conditional likelihood of a
571 winning response with the losing response of another prompt. Beyond DPO, Ethayarajh et al. [2024]
572 proposed a human-aware loss function-based framework using prospect theory named KTO, and Azar
573 et al. [2023] proposes IPO that uses human preferences expressed as pairwise preferences. Lastly,
574 Zhao et al. [2022] uses sequence likelihood calibration to align the model from human preference.
575 Despite of a vast body of work arising from DPO, none of the existing methods can operate and
576 contrast over the joint distribution of instruction-response pairs like the proposed DOVE algorithm.

577 C Comparison of Joint Preferences with Prior Preference Protocols

578 DOVE improves over prior work by acquiring ranking-based preferences over non-identical instruc-
579 tions that has remained unexplored in prior work (please refer to table 5). Diverse human reasoning
580 cannot be captured in the traditional conditional framework it fails to capture human preferences over
581 varied contexts. Context influences decision-making and subjective valuation when capturing human
582 preferences [Otto et al., 2022]. Prior work Yin et al. [2024], Liu et al. [2024], Meng et al. [2024],
583 Hong et al. [2024] collect conditional preferences in a pairwise manner and are variants of DPO
584 Rafailov et al. [2023]. Thus, in our experiments we compare DOVE to DPO directly. Furthermore,
585 we implement KTO Ethayarajh et al. [2024] as a baseline since KTO removes the requirements of
586 preference data that should be paired in preference optimization and implicitly compares responses
587 from different instructions. We find that DOVE outperforms both DPO and KTO.

588 D Dataset Statistics

589 We present the dataset statistics in Table 6. We report the number of instructions after filtering the
590 instances with repeated instructions. Each instance in the dataset consists of an instruction, and a pair
591 of responses. Originally, the number of AI-generated conditional and joint preferences equals the
592 number of instructions data. Here, we report the number of instances for which we observe a decisive
593 preference from ChatGPT i.e., after removing the ties.

594 E Proof for DOVE subsuming DPO

595 We highlight a result that reduces DOVE into DPO when the prompts are the same in Lemma E.1.

Preference Acquisition	Algorithm	Alignment Objective	Different Instructions
Score	Ethayarajh et al. [2024]	Conditional	No
Comparison (DPO Variants)	Rafailov et al. [2024]	Conditional	No
	Park et al. [2024]	Conditional	No
	Liu et al. [2024]	Conditional	No
	Meng et al. [2024]	Conditional	No
	Hong et al. [2024]	Conditional	No
Pairwise	DOVE (ours)	Joint	Yes

Table 5: We compare DOVE with existing frameworks based on three key aspects: preference acquisition (scoring or comparison), objective (conditional or joint distribution), and their ability to handle non-identical instruction-responses.

OpenAI TL;DR Summarization Dataset	Number
Number of instructions	11.8K
Number of AI generated conditional preferences	7.2K
Number of AI generated joint preferences	7.7K
Anthropic-Helpful Dataset	
Number of instructions	12.8K
Number of AI generated conditional preferences	9.4K
Number of AI generated joint preferences	8.5K

Table 6: Statistics for the train split of the summarization and open-ended dialogue datasets.

596 F DOVE on AlpacaEval2 Leaderboard

597 Similar to Rafailov et al. [2023], we show the usefulness of aligning LLMs using joint preferences
598 via DOVE on close-ended (e.g., summarization) and open-ended tasks (e.g., dialogues). However, we
599 further evaluate the effectiveness of our method on a broad set of instructions in the AlpacaEval2
600 leaderboard using the length-controlled win-rate metric Li et al. [2023].

Lemma E.1. *Under the case where $\mathcal{D}_X = \{(I_i, R_i, I_i, R_j)\}$, that is, prompts are the same for preferred and not-preferred prompt generation pairs, $\mathcal{L}_{\text{DPO}}(\theta; \mathcal{D}_C, \beta, p_{\text{ref}}) = \mathcal{L}_{\text{DOVE}}(\theta; \mathcal{D}_X, \beta, p_{\text{ref}})$, where $\mathcal{D}_C = \{(I_j, R_j^w, R_j^\ell)\}$.*

Proof.

$$\mathcal{L}_{\text{DOVE}}(\theta; \mathcal{D}_X, \beta, p_{\text{ref}}) = \mathbb{E}_{(I_j^w, R_j^w, I_j^\ell, R_j^\ell) \sim \mathcal{D}_X} \left[\log \left(\sigma \left(\beta \log \frac{p_\theta(R_i^w, I_i^w)}{p_{\text{ref}}(R_i^w, I_i^w)} - \beta \log \frac{p_\theta(R_j^\ell, I_j^\ell)}{p_{\text{ref}}(R_j^\ell, I_j^\ell)} \right) \right) \right] \quad (3)$$

$$= \mathbb{E}_{(I_j^w, R_j^w, I_j^\ell, R_j^\ell) \sim \mathcal{D}_X} \left[\log \left(\sigma \left(\beta \log \frac{p_\theta(R_i^w | I_i^w) p_\theta(I_i^w)}{p_{\text{ref}}(R_i^w | I_i^w) p_{\text{ref}}(I_i^w)} - \beta \log \frac{p_\theta(R_j^\ell | I_j^\ell) p_\theta(I_j^\ell)}{p_{\text{ref}}(R_j^\ell | I_j^\ell) p_{\text{ref}}(I_j^\ell)} \right) \right) \right] \quad (4)$$

$$= \mathbb{E}_{(I_j, R_j^w, R_j^\ell) \sim \mathcal{D}_C} \left[\log \left(\sigma \left(\beta \log \frac{p_\theta(R_j^w | I_j)}{p_{\text{ref}}(R_j^w | I_j)} - \beta \log \frac{p_\theta(R_j^\ell | I_j)}{p_{\text{ref}}(R_j^\ell | I_j)} \right) \right) \right] \quad (5)$$

$$= \mathcal{L}_{\text{DPO}}(\theta; \mathcal{D}_C, \beta, p_{\text{ref}}) \quad (6)$$

The proof follows from applying bayes rule and substituting $I_j^w = I_j^\ell = I_j$. \square

601 To do so, we train Mistral-7B base model on the UltraChat-200K dataset Ding et al. [2023] to get the
 602 SFT (reference) model. Subsequently, we utilize the conditional preference dataset, Ultrafeedback-
 603 binarized (60K instances) Cui et al. [2023] to align the SFT model using DPO as the baseline
 604 algorithm. Specifically, we utilize the training setup highlighted in the alignment handbook for SFT
 605 and DPO Tunstall et al. [2023]. Since DOVE algorithm allows access to joint preferences, we construct
 606 non-identical instruction-response tuples by pairing a chosen instruction-response (I_{chosen}, R_{chosen})
 607 with a rejected instruction-response (I_{reject}, R_{reject}) from the Ultrafeedback dataset.⁴ In particular,
 608 we train with DOVE algorithm for one epoch, and sweep over three learning rates $\{1e-7, 3e-7, 5e-7\}$
 609 and set the $\beta = 0.01$. Post-training, we sample responses from the SFT model, DPO-aligned LLM,
 610 DOVE-aligned LLM for the instructions in the AlpacaEval2 with a temperature of 0.7. We report the
 611 results in Table 7.

Method	Length-controlled Win-Rate (%)
SFT	9.13
DPO	15.7
DOVE	17.5

Table 7: Results on AlpacaEval2 leaderboard.

612 We find that the DOVE-aligned LLM outperforms DPO-aligned LLM by 1.8 percentage points on the
 613 challenging AlpacaEval2 leaderboard using the length-controlled win-rate metric. This indicates that
 614 the DOVE can utilize the joint preferences and elicit helpful and accurate responses for a broad set of
 615 instructions.

616 G Qualitative Examples

617 In this section, we present the qualitative examples to study the interplay between the conditional
 618 rankings and the joint preference over instruction-response pairs. Here, we acquire ranking feedback
 619 from the human annotators and ask them to provide the reasoning for their decision.

620 G.1 Anthropic-Helpful Examples

621 We present the qualitative examples for the preferences acquired for the Anthropic-helpful dataset in
 622 Figure 3, 4, and 5. We present our observations in the figure captions.

623 G.2 TL;DR Summarization Examples

624 We present the qualitative examples for the preferences acquired for the TL;DR summarization
 625 dataset in Figure 6, 7, and 8. We present our observations in the figure captions.

626 H Alignment Training Details

627 H.1 Supervised Finetuning Details

628 We present the SFT details in table 8. We perform full-finetuning of Mistral-7B using the source
 629 code from <https://github.com/abacaj/fine-tune-mistral>.

630 H.2 DOVE

631 We present the training details for DOVE preference optimization objective in the Table 9. We select
 632 the learning rate hyperparameter by sweeping over three learning rates: $\{1e-5, 5e-5, 5e-4\}$. We
 633 utilize the TRL library von Werra et al. [2020] for the DPO source code.

⁴For the sake of this experiment, we do not collect new joint preferences for this experiment, and rather utilize the pairings between chosen and rejected instruction-response pairs as a proxy for true joint preference distribution.

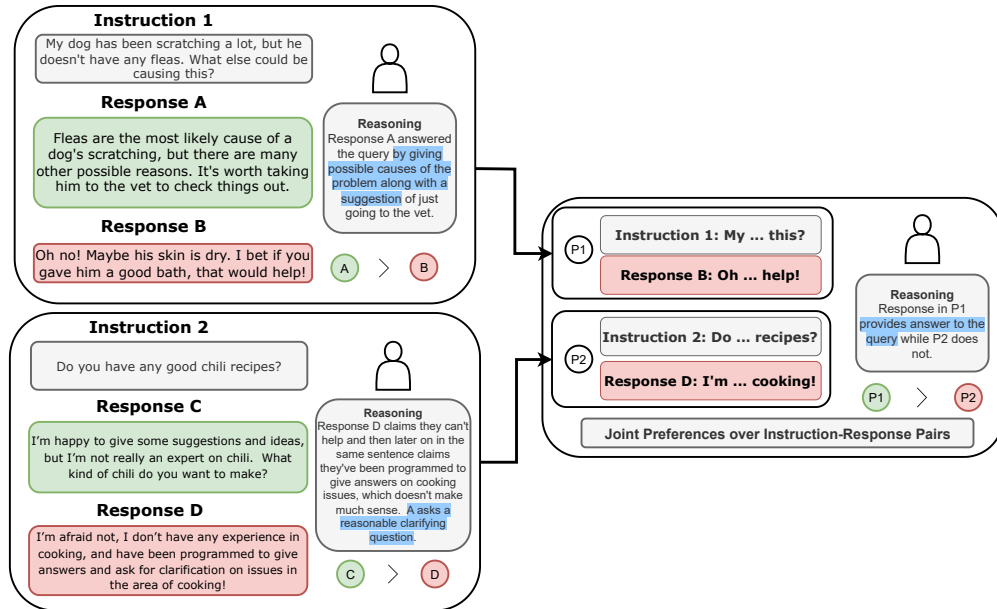


Figure 3: Interplay between the conditional rankings and joint rankings and reasoning acquired from the human annotators for the Anthropic-Helpful dataset. In this example, we find that the response B and D are rejected under the conditional rankings. When asked to compare the response B and D, humans consider that the response B answers Instruction 1 better than response D answers Instruction 2. This indicates that the joint preference humans elicits a decisive feedback between two responses that were rejected under the conditional rankings.

Anthropic-Helpful Dataset	
Learning Rate	1.5e-6
Batch Size	6
Epochs	3
OpenAI TL;DR Summarization Dataset	
Learning Rate	2e-5
Batch Size	12
Epochs	3

Table 8: Training details for the supervised finetuning of Mistral-7B.

634 I Trends with Data Scaling

635 We aim to understand the impact of increasing the number of preferences collected jointly over
 636 instruction-response pairs, for non-identical instructions, on the win-rate against the reference
 637 summaries in the TL;DR summarization dataset using DOVE algorithm. We present the results in
 638 Figure 9 for the sampling temperature of 0.001. We find that the win-rate scales from 42.4% to 71.7%
 639 as the size of the dataset increases from 100 to 9000 comparisons. We also observe that the change in
 640 the win-rate is within 1% when the dataset size increases from 4000 to 9000. This highlights that
 641 the performance gains are non-linear with the dataset size. In the future, it would be pertinent to
 642 explore techniques for selecting a subset of joint preference comparisons that result in maximum
 643 performance gains.

644 J ChatGPT Prompts

645 We present the ChatGPT for acquiring conditional rankings feedback and joint preferences over
 646 instruction-response pairs in Table 10 and Table 11, respectively.

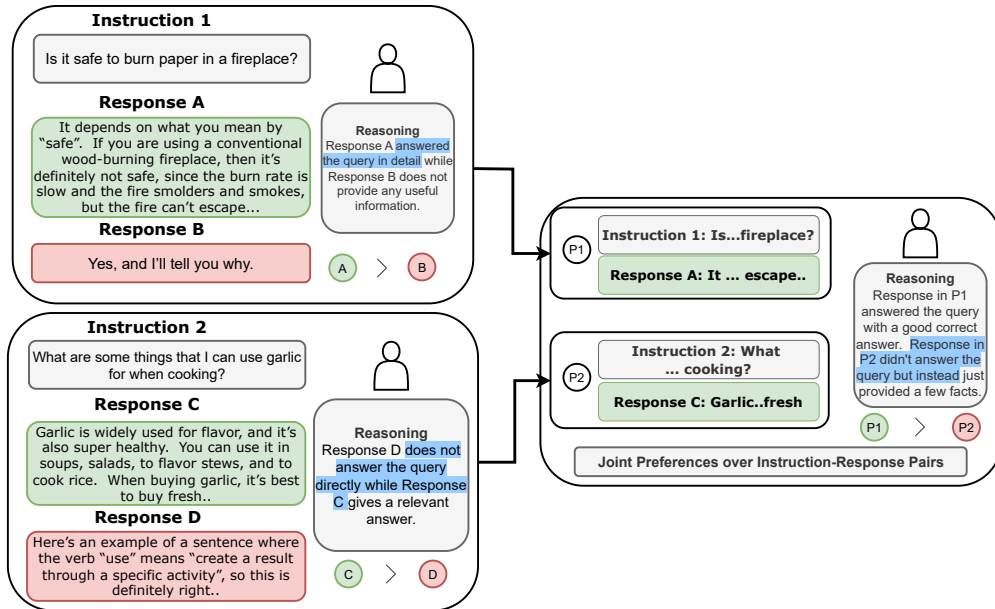


Figure 4: Interplay between the conditional rankings and joint rankings and reasoning acquired from the human annotators for the Anthropic-Helpful dataset. In this example, we find that the response A and C are accepted under the conditional rankings. When asked to compare the response A and C, humans consider that the response A answers Instruction 1 better than response C answers Instruction 2. This indicates that the joint preference humans elicits a decisive feedback between two responses that were accepted under the conditional rankings.

647 K Human Annotation Platform

648 We present the screenshots for the human interface in the Figure 12 (conditional rankings) and Figure
 649 13 (joint ranking preferences over instruction-response pairs).

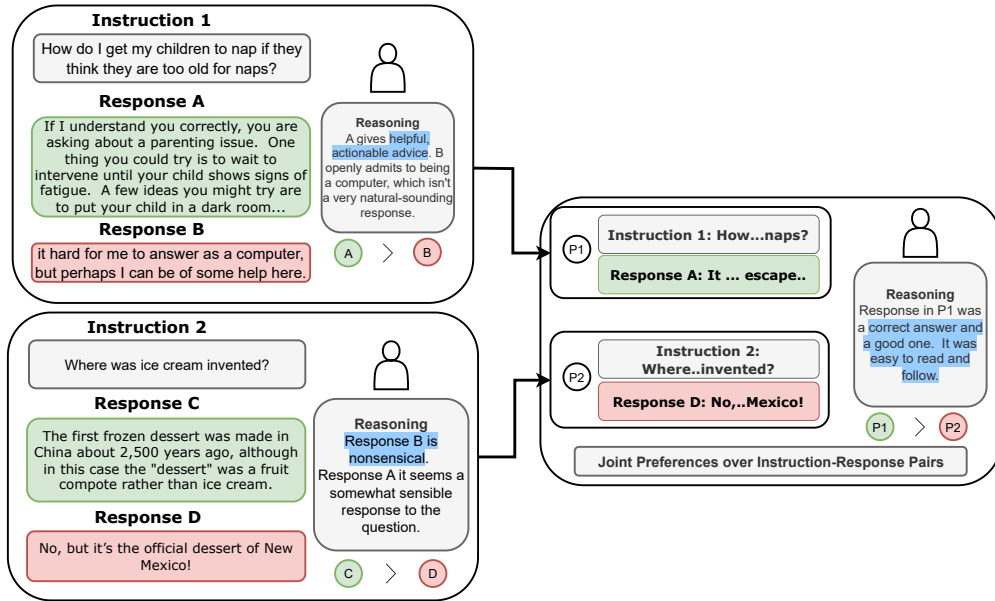


Figure 5: Interplay between the conditional rankings and joint rankings and reasoning acquired from the human annotators for the Anthropic-Helpful dataset. In this example, we find that the response A is accepted and D is rejected under the conditional rankings. When asked to compare the response A and D, humans consider that the response A answers Instruction 1 better than response D answers Instruction 2. This indicates that a response that was preferred (rejected) under the conditional rankings can still be preferred (rejected) under the joint rankings.

OpenAI TL;DR Summarization Dataset	
Peak Learning Rate	5e-5
Optimizer	AdamW [Loshchilov and Hutter, 2017]
Learning Schedule	Cosine
Batch Size	32
Epochs	10
Warmup Steps	100
α (LoRA)	16
Dropout (LoRA)	0.05
Bottleneck r (LoRA)	8
4bit Loading	True
β	0.1

Anthropic-Helpful Dataset	
Peak Learning Rate	5e-5
Optimizer	AdamW
Learning Schedule	Cosine
Batch Size	32
Epochs	5
Warmup Steps	100
α (LoRA)	16
Dropout (LoRA)	0.05
Bottleneck r (LoRA)	8
4bit Loading	True
β	0.1

Table 9: Training details for DOVE preference optimization objective. We use the identical settings for DPO.

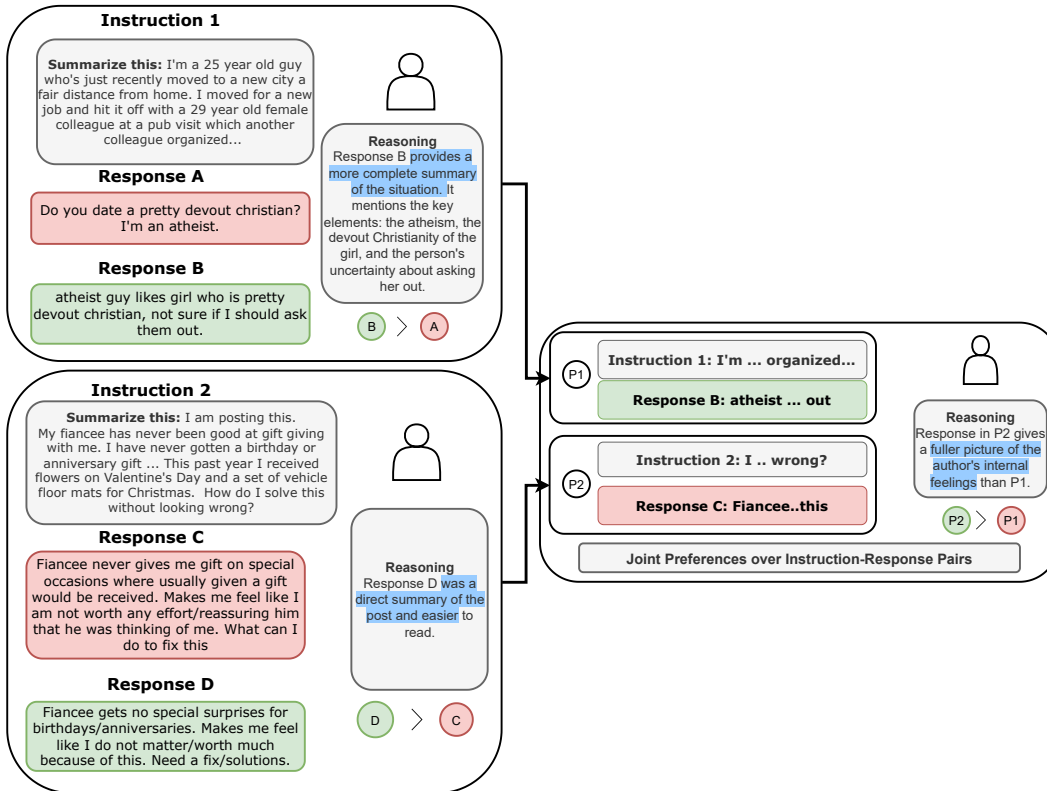


Figure 6: Interplay between the conditional rankings and joint rankings and reasoning acquired from the human annotators for the TL;DR summarization dataset. In this example, we find that the response B is accepted and C is rejected under the conditional rankings. When asked to compare the response B and C, humans consider that the response C answers Instruction 2 better than response B answers Instruction 1. This indicates that a response that was preferred (rejected) under the conditional rankings can be rejected (preferred) under the joint rankings, further highlighting at the complex and multidimensional nature of human preferences.

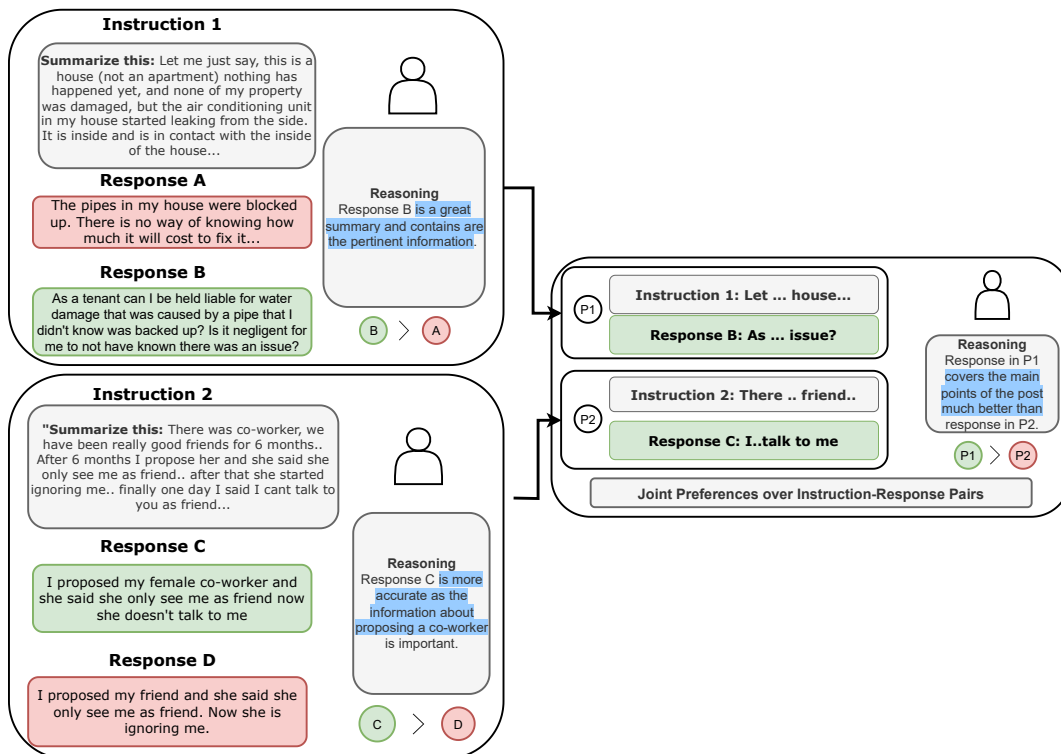


Figure 7: Interplay between the conditional rankings and joint rankings and reasoning acquired from the human annotators for the TL;DR summarization dataset. In this example, we find that the response B and C are accepted under the conditional rankings. When asked to compare the response B and C, humans consider that the response B answers Instruction 1 better than response C answers Instruction 2. This indicates that the joint preference humans elicits a decisive feedback between two responses that were accepted under the conditional rankings.

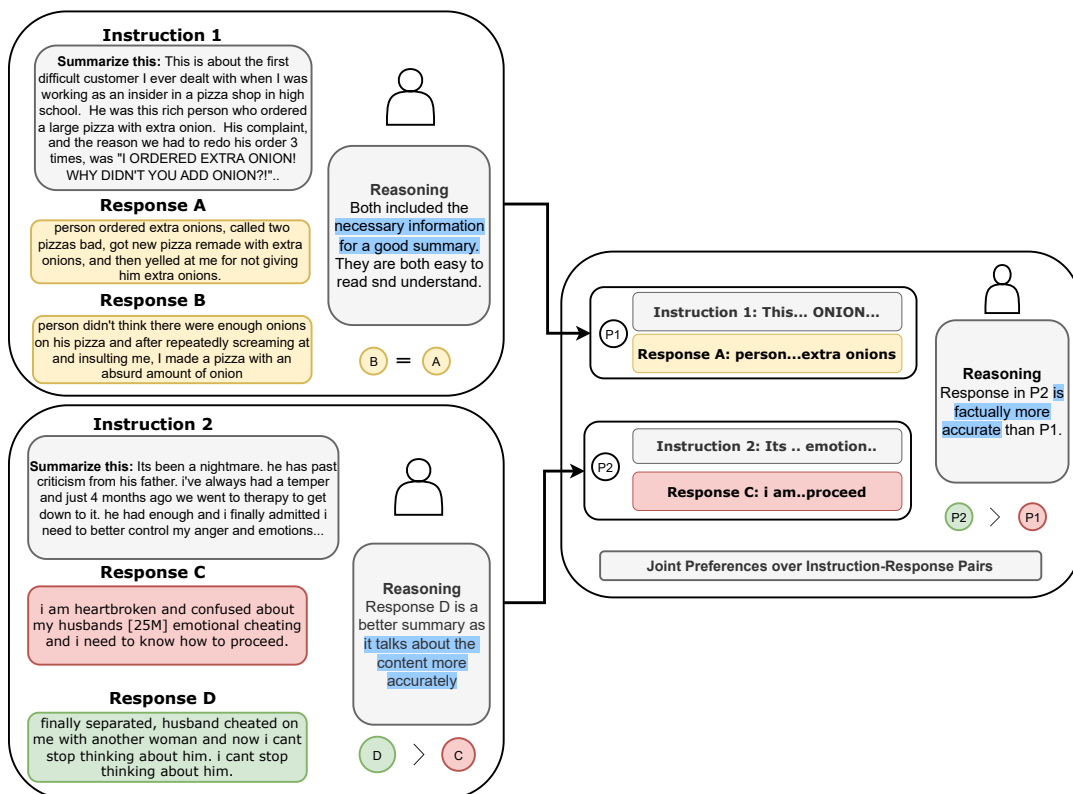


Figure 8: Interplay between the conditional rankings and joint rankings and reasoning acquired from the human annotators for the TL;DR summarization dataset. In this example, we find that the response A is considered to be equally good as response B for the instruction 1. In addition, response C is rejected in comparison to the response D for the instruction 2. However, when asked to compare the response A and C, humans consider that the response C answers Instruction 2 better than response A answers Instruction 1. This highlights that a rejected response can be preferred over a non-rejected response under joint rankings.

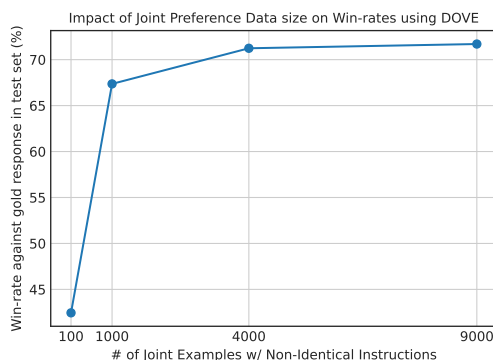


Figure 9: Results for scaling the feedback data size on TL;DR summarization dataset. We find that the win-rate improves with the increase in the dataset size using the DOVE preference optimization objective.

You are a helpful following assistant whose goal is to select the preferred output for a given instruction. Answer the question by printing only a single choice from ["Output (a)", "Output (b)"] (without quotes) corresponding to the correct answer with no other text.

Annotation Guideline

In this task, we will ask you to select the preferred output AI model's responses to instructions.

You will read examples, which are composed of the following:

1. an Instruction we give to the AI system
2. Output (a), the first output from the AI system
3. Output (b), the second output from the AI system

Your task is to decide which response is better for each example. There are several dimensions that you can think along. Consider the following dimensions:

Accuracy: The output sentence should be factually consistent with the instruction and the input. Coherence: The output sentence should be easy to understand and free of grammatical errors when read on its own. Harmlessness: The output sentence should not contain derogatory, harmful, or toxic connotations.

You should answer using only Output (a) or Output (b) depending on which response is better.

Instruction:
instruction

Output (a):
output₁

Output (b):
output₂

Preferred Output:

Figure 10: GPT-3.5-Turbo API prompt for comparisons on identical instructions

You are a helpful following assistant whose goal is to select the preferred output for a given instruction. Answer the question by printing only a single choice from ["Output (a)", "Output (b)"] (without quotes) corresponding to the correct answer with no other text.

Annotation Guideline

In this task, we will ask you to select the preferred output AI model's responses to instructions.

You will read examples, which are composed of the following:

1. Instruction (a), the first input to the AI system and Output (a), the first output from the AI system
2. Instruction (b), the second input to the AI system and Output (b), the second output from the AI system

Your task is to decide which response is better for each example. There are several dimensions that you can think along. Consider the following dimensions:

Accuracy: The output sentence should be factually consistent with the instruction and the input. Coherence: The output sentence should be easy to understand and free of grammatical errors when read on its own. Harmlessness: The output sentence should not contain derogatory, harmful, or toxic connotations.

You should answer using only Output (a) or Output (b) depending on which response is better.

Instruction (a):
instruction₁

Output (a):
output₁

Instruction (b):
instruction₂

Output (b):
output₂

Preferred Output:

Figure 11: GPT-3.5-Turbo API prompt for comparisons on non-identical instructions

Please thoroughly read the provided Instruction and the corresponding responses. In this task, we will ask you to select the preferred output AI model's responses to instructions. Your task is to decide which response is better for each example i.e., Response A, Response B, or whether both are equally good/bad. There are several dimensions that you can think along. Consider the following questions:

Is the response helpful? For example, if the instruction asked for a recipe for healthy food, and the response is a useful recipe, then we can consider it helpful.

Is the response language natural? For example, AI responses often have repetitions, which is not natural.

Is the response factual/accurate? For example, AI responses often make up new information. For example, if the response claims that Donald Trump is the current U.S. president, then you should consider it inaccurate.

and so on ... ultimately, you should decide which response is better based on your judgment and based on your own preference.

(WARNING: There might be some offensive and harmful content in the tasks.)

Instruction:

Response A:

Response B:

Choose the preferred response:

- Response A
- Response B
- Equally Good/Bad

Figure 12: Human annotation interface for Conditional Rankings

Please thoroughly read the provided Instruction and Response pairs. In this task, we will ask you to select the pair of instruction and response. Your task is to decide which response is better for the posed instruction. For example, Response A better answers the Instruction A (say summarize paragraph A) than Response B answers the Instruction B (say summarize paragraph B). Here, we are interested to know whether the model does a better summarization task for paragraph A or paragraph B. While this example is for summarizes, the actual task can have diverse prompts. Consider the following questions:

Is the response helpful? For example, if the instruction asked for a recipe for healthy food, and the response is a useful recipe, then we can consider it helpful.

Is the response language natural? For example, AI responses often have repetitions, which is not natural.

Is the response factual/accurate? For example, AI responses often make up new information. For example, if the response claims that Donald Trump is the current U.S. president, then you should consider it inaccurate.

and so on ... ultimately, you should decide which response is better based on your judgment and based on your own preference.

(WARNING: There might be some offensive and harmful content in the tasks.)

Instruction A:

Response A:

Instruction B:

Response B:

Choose the preferred instruction, response pair:

- Instruction A, Response A
- Instruction B, Response B
- Both pairs are equally answered well or bad

Figure 13: Human annotation interface for joint preferences over instruction-response pairs.