
Training language models to follow instructions with human feedback

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Making language models bigger does not inherently make them better at following
2 a user’s intent. For example, large language models can generate outputs that are
3 untruthful, toxic, or simply not helpful to the user. In other words, these models are
4 not *aligned* with their users. In this paper, we show an avenue for aligning language
5 models with user intent on a wide range of tasks by fine-tuning with human
6 feedback. Starting with a set of labeler-written prompts and prompts submitted
7 through a language model API, we collect a dataset of labeler demonstrations of
8 the desired model behavior, which we use to fine-tune GPT-3 using supervised
9 learning. We then collect a dataset of rankings of model outputs, which we use to
10 further fine-tune this supervised model using reinforcement learning from human
11 feedback. We call the resulting models *InstructGPT*. In human evaluations on
12 our prompt distribution, outputs from the 1.3B parameter InstructGPT model are
13 preferred to outputs from the 175B GPT-3, despite having 100x fewer parameters.
14 Moreover, InstructGPT models show improvements in truthfulness and reductions
15 in toxic output generation while having minimal performance regressions on public
16 NLP datasets. Even though InstructGPT still makes simple mistakes, our results
17 show that fine-tuning with human feedback is a promising direction for aligning
18 language models with human intent.

19 1 Introduction

20 Large language models (LMs) can be prompted to perform a range of natural language process-
21 ing (NLP) tasks, given some examples of the task as input. However, these models often express
22 unintended behaviors such as making up facts, generating biased or toxic text, or simply not following
23 user instructions (Bender et al., 2021; Bommasani et al., 2021; Kenton et al., 2021; Weidinger et al.,
24 2021; Tamkin et al., 2021; Gehman et al., 2020). This is because the language modeling objective
25 used for many recent large LMs—predicting the next token on a webpage from the internet—is
26 different from the objective “follow the user’s instructions helpfully and safely” (Radford et al., 2019;
27 Brown et al., 2020; Fedus et al., 2021; Rae et al., 2021; Thoppilan et al., 2022). Thus, we say that
28 the language modeling objective is *misaligned*. Averting these unintended behaviors is especially
29 important for language models that are deployed and used in hundreds of applications.

30 We make progress on aligning language models by training them to act in accordance with the user’s
31 intention (Leike et al., 2018). This encompasses both explicit intentions such as following instructions
32 and implicit intentions such as staying truthful, and not being biased, toxic, or otherwise harmful.
33 Using the language of Askeel et al. (2021), we want language models to be *helpful* (they should
34 help the user solve their task), *honest* (they shouldn’t fabricate information or mislead the user), and
35 *harmless* (they should not cause physical, psychological, or social harm to people or the environment).
36 We elaborate on the evaluation of these criteria in Section 3.5.

37 We focus on *fine-tuning* approaches
 38 to aligning language models. Specif-
 39 ically, we use reinforcement learning
 40 from human feedback (RLHF; Chris-
 41 tiano et al., 2017; Stiennon et al.,
 42 2020) to fine-tune GPT-3 to follow a
 43 broad class of written instructions (see
 44 Figure 2). This technique uses human
 45 preferences as a reward signal to fine-tune
 46 our models. We first hire a team of 40
 47 contractors to label our data, based on
 48 their performance on a screening test
 49 (see Section 3.3 and Appendix B.1 for
 50 more details). We then collect a dataset
 51 of human-written demonstrations of the
 52 desired output behavior on (mostly Eng-
 53 lish) prompts submitted to a language
 54 model API and some labeler-written
 55 prompts, and use this to train our
 56 supervised learning baselines. Next,
 57 we collect a dataset of human-labeled
 58 comparisons between outputs from our
 59 models on a larger set of API prompts.
 60 We then train a reward model (RM) on
 61 this dataset to predict which model output
 62 our labelers would prefer. Finally, we
 63 use this RM as a reward function and
 64 fine-tune our supervised learning base-
 65 line to maximize this reward using the
 66 PPO algorithm (Schulman et al., 2017).
 67 We illustrate this process in Figure 2.
 68 This procedure aligns the behavior of
 GPT-3 to the stated preferences of a
 specific group of people (mostly our
 labelers and researchers), rather than
 any broader notion of “human values”;
 we discuss this further in Appendix
 G.2. We call the resulting models
InstructGPT.

We mainly evaluate our models by
 having our labelers rate the quality of
 model outputs on our test set, consist-
 ing of prompts from held-out users
 (who are not represented in the train-
 ing data). We also conduct automatic
 evaluations on a range of public NLP
 datasets. We train three model sizes
 (1.3B, 6B, and 175B parameters), and
 all of our models use the GPT-3 archi-
 tecture. Our main findings are:

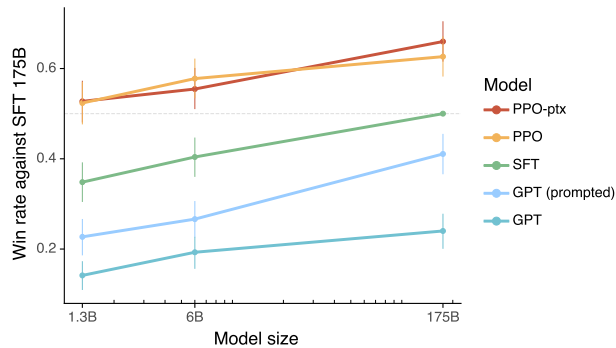


Figure 1: Human evaluations of various models on the API prompt distribution, evaluated by how often outputs from each model were preferred to those from the 175B SFT model. Our InstructGPT models (PPO-ptx) as well as its variant trained without pretraining mix (PPO) significantly outperform the GPT-3 baselines (GPT, GPT prompted).

69 **Labelers significantly prefer InstructGPT outputs over outputs from GPT-3.** Outputs from the
 70 1.3B parameter InstructGPT model are preferred to outputs from the 175B GPT-3, despite having
 71 over 100x fewer parameters. These models have the same architecture, and differ only by the fact that
 72 InstructGPT is fine-tuned on our human data. This result holds true even when we add a few-shot
 73 prompt to GPT-3 to make it better at following instructions. Outputs from our 175B InstructGPT
 74 are preferred to 175B GPT-3 outputs $85 \pm 3\%$ of the time, and preferred $71 \pm 4\%$ of the time to few-shot
 75 175B GPT-3. InstructGPT also generates more appropriate outputs according to our labelers.

76 **InstructGPT models show improvements in truthfulness over GPT-3.** On the TruthfulQA
 77 benchmark, InstructGPT generates truthful and informative answers more often than GPT-3. On
 78 “closed-domain” tasks from our API prompt distribution, where the output should not contain
 79 information that is not present in the input, InstructGPT models make up information not present in
 80 the input about half as often as GPT-3 (a 21% vs. 41% hallucination rate, respectively).

81 **InstructGPT shows small improvements in toxicity over GPT-3, but not bias.** To measure
 82 toxicity, we use the RealToxicityPrompts dataset (Gehman et al., 2020) and conduct both automatic
 83 and human evaluations. InstructGPT models generate about 25% fewer toxic outputs than GPT-3
 84 when prompted to be respectful. InstructGPT does not significantly improve over GPT-3 on the
 85 Winogender (Rudinger et al., 2018) and CrowSPairs (Nangia et al., 2020) datasets.

86 **We can minimize performance regressions on public NLP datasets by modifying our RLHF
 87 fine-tuning procedure.** During RLHF fine-tuning, we observe performance regressions compared
 88 to GPT-3 on certain public NLP datasets. We can greatly reduce the performance regressions on
 89 these datasets by mixing PPO updates with updates that increase the log likelihood of the pretraining
 90 distribution (PPO-ptx), without compromising labeler preference scores.

91 **Our models generalize to the preferences of “held-out” labelers that did not produce any**
92 **training data.** To test the generalization of our models, we conduct a preliminary experiment with
93 held-out labelers, and find that they prefer InstructGPT outputs to outputs from GPT-3 at about the
94 same rate as our training labelers. However, more work is needed to study how these models perform
95 on broader groups of users, and how they perform on inputs where humans disagree about the desired
96 behavior.

97 **Public NLP datasets are not reflective of how our language models are used.** We compare
98 GPT-3 fine-tuned on our human preference data (i.e. InstructGPT) to GPT-3 fine-tuned on two
99 different compilations of public NLP tasks: the FLAN (Wei et al., 2021) and T0 (Sanh et al., 2021)
100 (in particular, the T0++ variant). These datasets consist of a variety of NLP tasks, combined with
101 natural language instructions for each task. On our API prompt distribution, our FLAN and T0
102 models perform slightly worse than our SFT baseline, and labelers significantly prefer InstructGPT
103 to these models.

104 **InstructGPT models show promising generalization to instructions outside of the RLHF fine-**
105 **tuning distribution.** We qualitatively probe InstructGPT’s capabilities, and find that it is able to
106 follow instructions for summarizing code, answer questions about code, and sometimes follows
107 instructions in different languages, despite these instructions being very rare in the fine-tuning
108 distribution. This result is exciting because it suggests that our models are able to generalize the
109 notion of “following instructions.” They retain some alignment even on tasks for which they get very
110 little direct supervision.

111 **InstructGPT still makes simple mistakes.** For example, InstructGPT can still fail to follow
112 instructions, make up facts, give long hedging answers to simple questions, or fail to detect instructions
113 with false premises.

114 Overall, our results indicate that fine-tuning large language models using human preferences signifi-
115 cantly improves their behavior on a wide range of tasks, though much work remains to be done to
116 improve their safety and reliability.

117 **2 Related work**

118 **Research on alignment and learning from human feedback.** We build on previous techniques
119 to align models with human intentions, particularly reinforcement learning from human feed-
120 back (RLHF). Originally developed for training simple robots in simulated environments and Atari
121 games (Christiano et al., 2017; Ibarz et al., 2018), it has recently been applied to fine-tuning language
122 models to summarize text (Ziegler et al., 2019; Stiennon et al., 2020; Böhm et al., 2019; Wu et al.,
123 2021). This work is in turn influenced by similar work using human feedback as a reward in domains
124 such as dialogue (Jaques et al., 2019; Yi et al., 2019; Hancock et al., 2019), translation (Kreutzer et al.,
125 2018; Bahdanau et al., 2016), semantic parsing (Lawrence and Riezler, 2018), story generation (Zhou
126 and Xu, 2020), review generation (Cho et al., 2018), and evidence extraction (Perez et al., 2019). In
127 concurrent work, Askell et al. (2021); Bai et al. (2022) propose language assistants as a testbed for
128 alignment research, and train models using RLHF. Our work can be seen as a direct application of
129 RLHF to aligning language models on a broad distribution of language tasks.

130 **Training language models to follow instructions.** Our work is also related to research on cross-
131 task generalization in language models, where LMs are fine-tuned on a broad range of public NLP
132 datasets (usually prefixed with an appropriate instruction) and evaluated on a different set of NLP
133 tasks. There has been a range of work in this domain (Yi et al., 2019; Mishra et al., 2021; Wei et al.,
134 2021; Khashabi et al., 2020; Sanh et al., 2021; Aribandi et al., 2021), which differ in training and
135 evaluation data, formatting of instructions, size of pretrained models, and other experimental details.

136 **Mitigating the harms of language models.** A goal of modifying the behavior of language models
137 is to mitigate the harms of these models when they’re deployed in the real world. These risks have
138 been extensively documented (Bender et al., 2021; Bommasani et al., 2021; Kenton et al., 2021;
139 Weidinger et al., 2021; Tamkin et al., 2021). Language models can produce biased outputs (Dhamala
140 et al., 2021; Liang et al., 2021; Manela et al., 2021; Caliskan et al., 2017; Kirk et al., 2021), leak

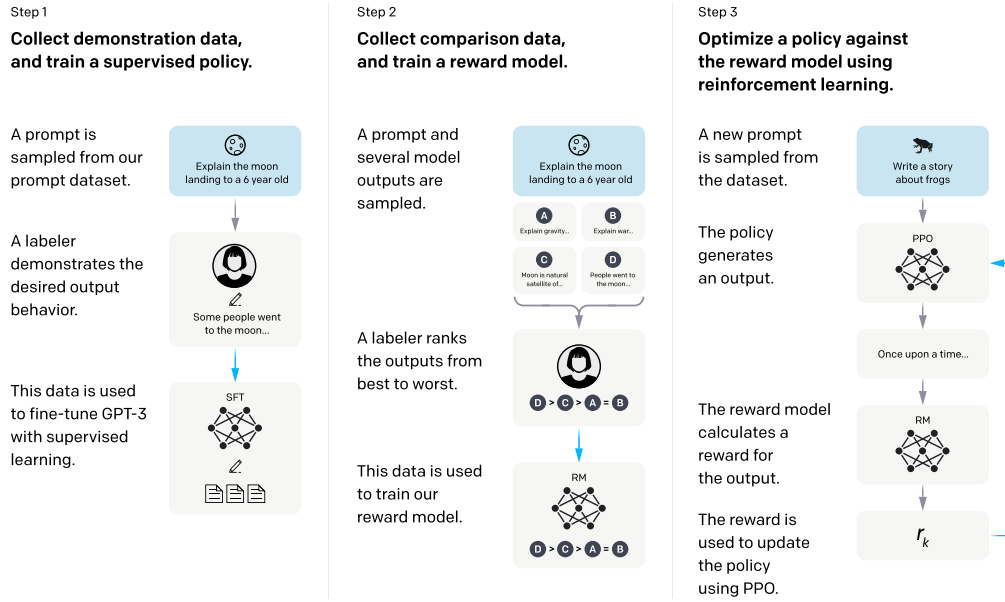


Figure 2: A diagram illustrating the three steps of our method: (1) supervised fine-tuning (SFT), (2) reward model (RM) training, and (3) reinforcement learning via proximal policy optimization (PPO) on this reward model. Blue arrows indicate that this data is used to train one of our models. In Step 2, boxes A-D are samples from our models that get ranked by labelers.

141 private data (Carlini et al., 2021), generate misinformation (Solaiman et al., 2019; Buchanan et al.,
 142 2021), and be used maliciously; for a thorough review we direct the reader to Weidinger et al.
 143 (2021). There are many ways to mitigate these harms, including by fine-tuning on a small, value-
 144 targeted dataset (Solaiman and Dennison, 2021), filtering the pretraining dataset (Ngo et al., 2021),
 145 or human-in-the-loop data collection (Dinan et al., 2019; Xu et al., 2020).

146 3 Methods and experimental details

147 3.1 High-level methodology

148 Our methodology follows that of Ziegler et al. (2019) and Stiennon et al. (2020), who applied
 149 it in the stylistic continuation and summarization domains. We start with a pretrained language
 150 model (Radford et al., 2019; Brown et al., 2020; Fedus et al., 2021; Rae et al., 2021; Thoppilan et al.,
 151 2022), a distribution of prompts on which we want our model to produce aligned outputs, and a
 152 team of trained human labelers (see Section 3.3 for details). We then apply the following three steps
 153 (Figure 2).

154 **Step 1: Collect demonstration data, and train a supervised policy.** Our labelers provide demon-
 155 strations of the desired behavior on the input prompt distribution (see Section 3.2 for details on this
 156 distribution). We then fine-tune a pretrained GPT-3 model on this data using supervised learning.

157 **Step 2: Collect comparison data, and train a reward model.** We collect a dataset of comparisons
 158 between model outputs, where labelers indicate which output they prefer for a given input. We then
 159 train a reward model to predict the human-preferred output.

160 **Step 3: Optimize a policy against the reward model using PPO.** We use the output of the
 161 RM as a scalar reward. We fine-tune the supervised policy to optimize this reward using the PPO
 162 algorithm (Schulman et al., 2017).

163 Steps 2 and 3 can be iterated continuously; more comparison data is collected on the current best
 164 policy, which is used to train a new RM and then a new policy. In practice, most of our comparison
 165 data comes from our supervised policies, with some coming from our PPO policies.

166 3.2 Dataset

167 Our prompt dataset consists primarily of text prompts submitted to a commercial language model API,
168 as well as a small number of labeler-written prompts. These prompts are very diverse and include
169 generation, question answering, dialog, summarization, extractions, and other natural language
170 tasks (see Appendix A). Our dataset is over 96% English. We heuristically deduplicate prompts, and
171 ensure that the validation and test sets contain no data from users whose data is in the training set.
172 We also filter prompts containing personally identifiable information (PII).

173 From these prompts, we produce three different datasets used in our fine-tuning procedure: (1) our
174 SFT dataset, with labeler demonstrations used to train our SFT models, (2) our RM dataset, with
175 labeler rankings of model outputs used to train our RMs, and (3) our PPO dataset, without any human
176 labels, which are used as inputs for RLHF fine-tuning. The SFT dataset contains about 13k training
177 prompts (from the API and labeler-written), the RM dataset has 33k training prompts (from the API
178 and labeler-written), and the PPO dataset has 31k training prompts (only from the API). More details
179 on dataset sizes are provided in Table 3.

180 3.3 Human data collection

181 To produce our demonstration and comparison data, and to conduct our main evaluations, we hired
182 a team of about 40 contractors on Upwork and through ScaleAI. Compared to earlier work that
183 collects human preference data on the task of summarization (Ziegler et al., 2019; Stiennon et al.,
184 2020; Wu et al., 2021), our inputs span a much broader range of tasks, and can occasionally include
185 controversial and sensitive topics. Our aim was to select a group of labelers who were sensitive to the
186 preferences of different demographic groups, and who were good at identifying outputs that were
187 potentially harmful. Thus, we conducted a screening test designed to measure labeler performance
188 on these axes (see Appendix B.1). As an initial study to see how well our model generalizes to the
189 preferences of other labelers, we hire a separate set of labelers who do not produce any of the training
190 data. These labelers are sourced from the same vendors, but do not undergo a screening test.

191 Despite the complexity of the task, we find that inter-annotator agreement rates are quite high:
192 training labelers agree with each-other $72.6 \pm 1.5\%$ of the time, while for held-out labelers this
193 number is $77.3 \pm 1.3\%$. For comparison, in the summarization work of Stiennon et al. (2020)
194 researcher-researcher agreement was $73 \pm 4\%$.

195 3.4 Models

196 Starting from GPT-3 (Brown et al., 2020), we train models with three different techniques:

197 **Supervised fine-tuning (SFT).** We fine-tune GPT-3 on our labeler demonstrations using supervised
198 learning. We trained for 16 epochs, using a cosine learning rate decay, and residual dropout of 0.2.
199 We do our final SFT model selection based on the RM score on the validation set. Similarly to Wu
200 et al. (2021), we find that our SFT models overfit on validation loss after 1 epoch; however, we find
201 that training for more epochs helps both the RM score and human preference ratings.

202 **Reward modeling (RM).** We fine-tune GPT-3 to take in a prompt and response, and output a scalar
203 reward. In this paper we only use 6B RMs, as this saves a lot of compute, and we found that 175B
204 RM training could be unstable and thus was less suitable to be used as the value function during RL
205 (see Appendix D for more details).

206 In Stiennon et al. (2020), the RM is trained on a dataset of comparisons between two model outputs
207 on the same input. They use a cross-entropy loss, with the comparisons as labels—the difference in
208 rewards represents the log odds that one response will be preferred to the other by a human labeler. In
209 order to speed up comparison collection, we have labelers rank between $K = 4$ and $K = 9$ responses,
210 and train on all $\binom{K}{2}$ comparisons from each prompt as a single batch element, for computational
211 efficiency (see Appendix D. The loss function for the RM becomes:

$$\text{loss}(\theta) = -\frac{1}{\binom{K}{2}} E_{(x, y_w, y_l) \sim D} [\log(\sigma(r_\theta(x, y_w) - r_\theta(x, y_l)))] \quad (1)$$

212 where $r_\theta(x, y)$ is the scalar output of the reward model for prompt x and completion y with parameters
213 θ , y_w is the preferred completion out of the pair of y_w and y_l , and D is the comparison dataset.

214 **Reinforcement learning (RL).** Again following Stiennon et al. (2020), we fine-tuned the SFT
215 model using PPO (Schulman et al., 2017). The environment is a bandit environment which presents
216 a random user prompt and expects a response to the prompt. Given the prompt and response, it
217 produces a reward determined by the reward model and ends the episode. In addition, we add a
218 per-token KL penalty from the SFT model at each token to mitigate over-optimization of the reward
219 model. The value function is initialized from the RM. We call these models “PPO.”

220 We also experiment with mixing the pretraining gradients into the PPO gradients, in order to fix
221 the performance regressions on public NLP datasets (see Appendix D.4). We call these models
222 “PPO-ptx.” Unless otherwise specified, in this paper InstructGPT refers to the PPO-ptx models.

223 **Baselines.** We compare the performance of our PPO models to our SFT models and GPT-3. We also
224 compare to GPT-3 when it is provided a few-shot prefix to ‘prompt’ it into an instruction-following
225 mode (GPT-3-prompted). This prefix is prepended to the user-specified instruction.

226 We additionally compare InstructGPT to fine-tuning 175B GPT-3 on the FLAN (Wei et al., 2021)
227 and T0 (Sanh et al., 2021) datasets, which both consist of a variety of NLP tasks, combined with
228 natural language instructions for each task (they differ in the NLP datasets included, and the style of
229 instructions used). We fine-tune them on approximately 1 million examples and choose the checkpoint
230 which obtains the highest RM score on the validation set (see Appendix D for more details).

231 3.5 Evaluation

232 Following Askell et al. (2021), we say our models are aligned if they are helpful, truthful, and
233 harmless (we elaborate in Appendix C.2). We divide our quantitative evaluations into two parts:

234 **Evaluations on API distribution.** Our main metric is human preference ratings on a held out set
235 of prompts from the same source as our training distribution. When using prompts from the API
236 for evaluation, we only select prompts by users we haven’t included in training. For each model we
237 calculate how often its outputs are preferred to a baseline policy; we choose our 175B SFT model
238 as the baseline since its performance is near the middle of the pack. Additionally, we ask labelers
239 to judge the overall quality of each response on a 1-7 Likert scale and collect a range of metadata
240 for each model output (see Table 11). In particular, we collect data that aims to capture different
241 aspects of behavior in a deployed model that could end up being harmful: we have labelers evaluate
242 whether an output is inappropriate in the context of a customer assistant, denigrates a protected class,
243 or contains sexual or violent content.

244 **Evaluations on public NLP datasets.** We evaluate on two types of public datasets: those that
245 capture an aspect of language model safety, particularly truthfulness, toxicity, and bias, and those
246 that capture zero-shot performance on traditional NLP tasks like question answering, reading com-
247 prehension, and summarization. We also conduct human evaluations on the RealToxicityPrompts
248 dataset (Gehman et al., 2020).

249 4 Results

250 4.1 Results on the API distribution

251 **Labelers significantly prefer InstructGPT outputs over outputs from GPT-3.** On our test
252 set, our labelers significantly prefer InstructGPT outputs across model sizes (Figure 1). We find
253 that GPT-3 outputs perform the worst, and one can obtain significant step-size improvements by
254 using a well-crafted few-shot prompt (GPT-3 (prompted)), then by training on demonstrations using
255 supervised learning (SFT), and finally by training on comparison data using PPO. Adding updates on
256 the pretraining mix during PPO does not lead to large changes in labeler preference. To illustrate the
257 magnitude of our gains: when compared directly, 175B InstructGPT outputs are preferred to GPT-3
258 outputs $85 \pm 3\%$ of the time, and preferred $71 \pm 4\%$ of the time to few-shot GPT-3.

259 In Figure 4 we show that labelers also rate InstructGPT outputs favorably along several more concrete
260 axes. Specifically, compared to GPT-3, InstructGPT outputs are more appropriate in the context of a
261 customer assistant, more often follow explicit constraints defined in the instruction (e.g. “Write your

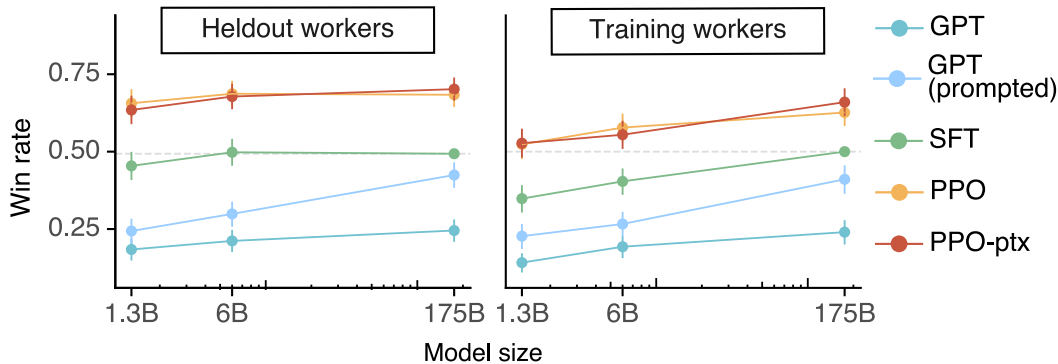


Figure 3: Preference results of our models, measured by winrate against the 175B SFT model.

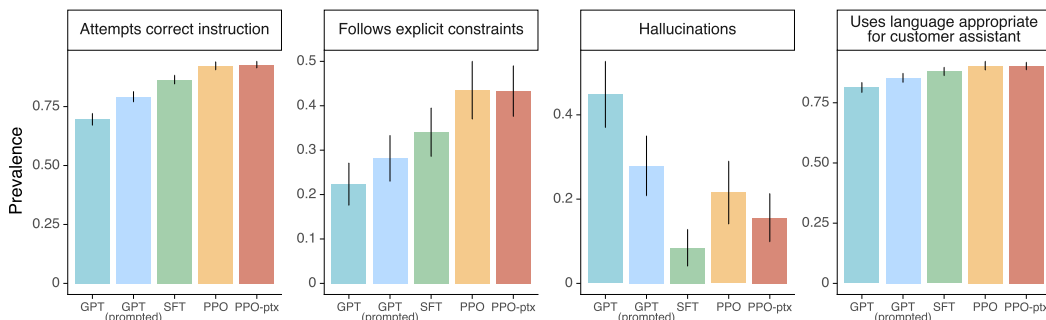


Figure 4: Metadata results on the API distribution, averaged over model sizes.

262 answer in 2 paragraphs or less.”), are less likely to fail to follow the correct instruction entirely, and
 263 and make up facts (‘hallucinate’) less often in closed-domain tasks.

264 **Our models generalize to the preferences of “held-out” labelers that did not produce any train-**
 265 **ing data.** Held-out labelers have similar ranking preferences as workers who we used to produce
 266 training data (see Figure 3). In particular, according to held-out workers, all of our InstructGPT
 267 models still greatly outperform the GPT-3 baselines. Thus, our InstructGPT models aren’t simply
 268 overfitting to the preferences of our training labelers.

269 **Public NLP datasets are not reflective of how our language models are used.** In Figure 5a,
 270 we also compare InstructGPT to our 175B GPT-3 baselines fine-tuned on the FLAN (Wei et al.,
 271 2021) and T0 (Sanh et al., 2021) datasets (see Appendix D for details). We find that these models
 272 perform better than GPT-3, on par with GPT-3 with a well-chosen prompt, and worse than our SFT
 273 baseline. This indicates that these datasets are not sufficiently diverse to improve performance on our
 274 API prompt distribution. We believe this is partly because academic datasets focus on tasks where
 275 performance is easily measured, like classification and QA, while our API distribution consists of
 276 mostly (about 57%) open-ended generation tasks.

277 4.2 Results on public NLP datasets

278 **InstructGPT models show improvements in truthfulness over GPT-3.** As measured by human
 279 evaluations on the TruthfulQA dataset, our PPO models show small but significant improvements
 280 in generating truthful and informative outputs compared to GPT-3 (see Figure 5b). This behavior is
 281 the default: our models do not have to be specifically instructed to tell the truth to exhibit improved
 282 truthfulness. Interestingly, the exception is our 1.3B PPO-ptx model, which performs slightly worse
 283 than a GPT-3 model of the same size. Our improvements in truthfulness are also evidenced by the
 284 fact that our PPO models hallucinate less often on closed-domain tasks (Figure 4).

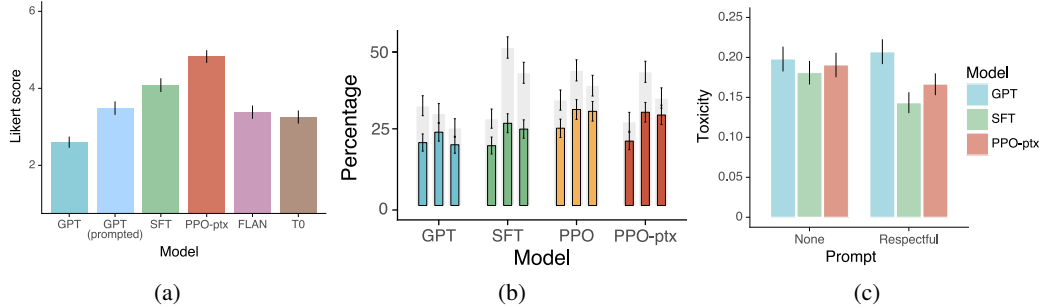


Figure 5: (a) Comparing our models with GPT-3 fine-tuned on the FLAN and T0 datasets, in terms of 1-7 Likert scores, on our prompt distribution. (b) Human evaluations on the TruthfulQA dataset. Gray bars indicate ratings of truthfulness; colored bars indicate ratings of truthfulness *and* informativeness. (c) Human evaluations on RealToxicityPrompts, with and without "respectful" instructions.

285 **InstructGPT shows small improvements in toxicity over GPT-3, but not bias.** We first evaluate
 286 our models on the RealToxicityPrompts dataset (Gehman et al., 2020) using human evaluations.
 287 Our results are in Figure 5c. We find that, when instructed to produce a safe and respectful output
 288 (“respectful prompt”), InstructGPT models generate less toxic outputs than those from GPT-3
 289 according to the Perspective API. This advantage disappears when the respectful prompt is removed
 290 (“no prompt”). We see similar results when evaluating using the Perspective API (Appendix F.7).

291 **We can minimize performance regressions on public NLP datasets by modifying our RLHF**
 292 **fine-tuning procedure.** In Figure 25 we show that adding pretraining updates to our PPO fine-
 293 tuning (PPO-ptx) mitigates performance regressions on public NLP datasets, and even surpasses
 294 GPT-3 on HellaSwag. The performance of the PPO-ptx model still lags behind GPT-3 on DROP,
 295 SQuADv2, and translation; more work is needed to study and further eliminate these performance
 296 regressions. We also find that mixing in pretraining updates performs better than the simpler solution
 297 of increasing the KL coefficient (Figure 36).

298 4.3 Qualitative results

299 **InstructGPT models show promising generalization to instructions outside of the RLHF fine-**
 300 **tuning distribution.** In particular, we find that InstructGPT shows ability to follow instructions
 301 in non-English languages, and perform summarization and question-answering for code. This is
 302 interesting because non-English languages and code form a tiny minority of our fine-tuning data, and
 303 it suggests that, in some cases, alignment methods could generalize to producing the desired behavior
 304 on inputs that humans did not directly supervise. We show some qualitative examples in Figure 26.

305 **InstructGPT still makes simple mistakes.** In interacting with our 175B PPO-ptx model, we
 306 have noticed it can still make simple mistakes, despite its strong performance on many different
 307 language tasks. To give a few examples: (1) when given an instruction with a false premise, the model
 308 sometimes incorrectly assumes the premise is true, (2) the model can overly hedge; when given a
 309 simple question, it can sometimes say that there is no one answer to the question and give multiple
 310 possible answers, even when there is one fairly clear answer from the context, and (3) the model’s
 311 performance degrades when instructions contain multiple explicit constraints (e.g. “list 10 movies
 312 made in the 1930’s set in France”) or when constraints can be challenging for language models (e.g.
 313 writing a summary in a specified number of sentences).

314 We show some examples of these behaviors in Figure 27. We suspect that behavior (2) emerges
 315 partly because we instruct labelers to reward epistemic humility; thus, they may tend to reward
 316 outputs that hedge, and this gets picked up by our reward model. We suspect that behavior (1) occurs
 317 because there are few prompts in the training set that assume false premises, and our models don’t
 318 generalize well to these examples. We believe both these behaviors could be dramatically reduced
 319 with adversarial data collection (Dinan et al., 2019).

320 5 Discussion

321 5.1 Implications for alignment research

322 Our approach to alignment research in this work is iterative: we are improving the alignment of
323 current AI systems instead of focusing abstractly on aligning AI systems that don't yet exist, which
324 provides us with a clear empirical feedback loop of what works and what does not. We believe that
325 this feedback loop is essential to refine our alignment techniques, and it forces us to keep pace with
326 progress in machine learning.

327 From this work, we can draw lessons for alignment research more generally. First, the cost of
328 increasing model alignment is modest relative to pretraining. Training our 175B SFT model requires
329 4.9 petaflops/s-days and training our 175B PPO-ptx model requires 60 petaflops/s-days, compared
330 to 3,640 petaflops/s-days for GPT-3 (Brown et al., 2020). At the same time, our results show that
331 RLHF is very effective at making language models more helpful to users, more so than a 100x model
332 size increase. This suggests that right now increasing investments in alignment of existing language
333 models is more cost-effective than training larger models. Second, we've seen some evidence that
334 InstructGPT generalizes 'following instructions' to settings that we don't supervise it in. This is an
335 important property because it's prohibitively expensive to have humans supervise models on every
336 task they perform. Finally, we were able to mitigate most of the performance degradations introduced
337 by our fine-tuning. If this was not the case, these performance degradations would constitute an
338 alignment tax—an additional cost for aligning the model. Any alignment technique with a high tax
339 might not see adoption, and thus such a tax is important to avoid.

340 5.2 Limitations

341 **Methodology.** The behavior of our InstructGPT models is determined in part by the human feedback
342 obtained from our contractors. Some of the labeling tasks rely on value judgments that may be
343 impacted by the identity of our contractors, their beliefs, cultural backgrounds, and personal history.
344 We kept our team of contractors small because this facilitates high-bandwidth communication with
345 a smaller set of contractors who are doing the task full-time. However, this group is clearly not
346 representative of the full spectrum of people affected by these models. As a simple example, our
347 labelers are primarily English-speaking and our data consists almost entirely of English instructions.

348 **Models.** Our models are neither fully aligned nor fully safe; they still generate toxic or biased
349 outputs, make up facts, and generate sexual and violent content without explicit prompting. They can
350 also fail to generate reasonable outputs on some inputs; we show some examples of this in Figure 27.
351 Perhaps the greatest limitation of our models is that, in most cases, they follow the user's instruction,
352 even if that could lead to harm in the real world. For example, when prompting the models to be
353 maximally biased, InstructGPT generates more toxic outputs than equivalently-sized GPT-3 models.

354 5.3 Broader impacts

355 This work is motivated by our aim to increase the positive impact of large language models by training
356 them to do what a given set of humans want them to do. By default, language models optimize
357 the next word prediction objective, which is only a proxy for what we want these models to do.
358 Our results indicate that our techniques hold promise for making language models more helpful,
359 truthful, and harmless. In the longer term, alignment failures could lead to more severe consequences,
360 particularly if these models are deployed in safety-critical situations.

361 However, making language models better at following user intentions also makes them easier to
362 misuse. It may be easier to use these models to generate convincing misinformation, or hateful or
363 abusive content. Alignment techniques are not a panacea for resolving safety issues associated with
364 large language models; rather, they should be used as one tool in a broader safety ecosystem. Aside
365 from intentional misuse, there are many domains where large language models should be deployed
366 only with great care, or not at all. Examples include high-stakes domains such as medical diagnoses,
367 classifying people based on protected characteristics, determining eligibility for credit, employment,
368 or housing, generating political advertisements, and law enforcement.

369 Finally, the question of who these models are aligned to is extremely important, and will significantly
370 affect whether the net impact of these models is positive or negative; we discuss this in Appendix G.2.

371 Checklist

- 372 1. For all authors...
- 373 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
374 contributions and scope? [Yes]
- 375 (b) Did you describe the limitations of your work? [Yes]
- 376 (c) Did you discuss any potential negative societal impacts of your work? [Yes]
- 377 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
378 them? [Yes]
- 379 2. If you are including theoretical results...
- 380 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 381 (b) Did you include complete proofs of all theoretical results? [N/A]
- 382 3. If you ran experiments...
- 383 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
384 mental results (either in the supplemental material or as a URL)? [No]
- 385 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
386 were chosen)? [Yes]
- 387 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
388 ments multiple times)? [Yes]
- 389 (d) Did you include the total amount of compute and the type of resources used (e.g., type
390 of GPUs, internal cluster, or cloud provider)? [No] : we provide some info on the
391 amount of compute used in the Discussion section.
- 392 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 393 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 394 (b) Did you mention the license of the assets? [No]
- 395 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 396 (d) Did you discuss whether and how consent was obtained from people whose data you’re
397 using/curating? [Yes]
- 398 (e) Did you discuss whether the data you are using/curating contains personally identifiable
399 information or offensive content? [Yes] : PII was removed, the dataset contains some
400 offensive content.
- 401 5. If you used crowdsourcing or conducted research with human subjects...
- 402 (a) Did you include the full text of instructions given to participants and screenshots, if
403 applicable? [No] : we provide excerpts of instructions given to labelers in the Appendix,
404 but the full instructions are very long.
- 405 (b) Did you describe any potential participant risks, with links to Institutional Review
406 Board (IRB) approvals, if applicable? [No]
- 407 (c) Did you include the estimated hourly wage paid to participants and the total amount
408 spent on participant compensation? [No] : though we provide lots of information about
409 labelers, including a labeler satisfaction survey, in the Appendix.

410 References

- 411 Achiam, J., Held, D., Tamar, A., and Abbeel, P. (2017). Constrained policy optimization. In
412 *International Conference on Machine Learning*, pages 22–31. PMLR.
- 413 Anthony, T., Tian, Z., and Barber, D. (2017). Thinking fast and slow with deep learning and tree
414 search. *arXiv preprint arXiv:1705.08439*.
- 415 Aribandi, V., Tay, Y., Schuster, T., Rao, J., Zheng, H. S., Mehta, S. V., Zhuang, H., Tran, V. Q., Bahri,
416 D., Ni, J., et al. (2021). Ext5: Towards extreme multi-task scaling for transfer learning. *arXiv
417 preprint arXiv:2111.10952*.
- 418 Askell, A., Bai, Y., Chen, A., Drain, D., Ganguli, D., Henighan, T., Jones, A., Joseph, N., Mann, B.,
419 DasSarma, N., et al. (2021). A general language assistant as a laboratory for alignment. *arXiv
420 preprint arXiv:2112.00861*.

- 421 Bahdanau, D., Brakel, P., Xu, K., Goyal, A., Lowe, R., Pineau, J., Courville, A., and Bengio, Y.
422 (2016). An actor-critic algorithm for sequence prediction. *arXiv preprint arXiv:1607.07086*.
- 423 Bai, Y., Jones, A., Ndousse, K., Askell, A., Chen, A., DasSarma, N., Drain, D., Fort, S., Ganguli,
424 D., Henighan, T., et al. (2022). Training a helpful and harmless assistant with reinforcement
425 learning from human feedback. *arXiv preprint arXiv:2204.05862*.
- 426 Bender, E. M., Gebru, T., McMillan-Major, A., and Shmitchell, S. (2021). On the dangers of stochastic
427 parrots: Can language models be too big? In *Proceedings of the 2021 ACM Conference on*
428 *Fairness, Accountability, and Transparency*, pages 610–623.
- 429 Böhm, F., Gao, Y., Meyer, C. M., Shapira, O., Dagan, I., and Gurevych, I. (2019). Better rewards yield
430 better summaries: Learning to summarise without references. *arXiv preprint arXiv:1909.01214*.
- 431 Bojar, O., Chatterjee, R., Federmann, C., Haddow, B., Huck, M., Hokamp, C., Koehn, P., Logacheva,
432 V., Monz, C., Negri, M., Post, M., Scarton, C., Specia, L., and Turchi, M. (2015). Findings of
433 the 2015 workshop on statistical machine translation. In *Proceedings of the Tenth Workshop on*
434 *Statistical Machine Translation*, pages 1–46, Lisbon, Portugal. Association for Computational
435 Linguistics.
- 436 Bommasani, R., Hudson, D. A., Adeli, E., Altman, R., Arora, S., von Arx, S., Bernstein, M. S., Bohg,
437 J., Bosselut, A., Brunskill, E., et al. (2021). On the opportunities and risks of foundation models.
438 *arXiv preprint arXiv:2108.07258*.
- 439 Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam,
440 P., Sastry, G., Askell, A., et al. (2020). Language models are few-shot learners. *arXiv preprint*
441 *arXiv:2005.14165*.
- 442 Buchanan, B., Lohn, A., Musser, M., and Sedova, K. (2021). Truth, lies, and automation. Technical
443 report, Center for the Study of Emerging Technology.
- 444 Caliskan, A., Bryson, J. J., and Narayanan, A. (2017). Semantics derived automatically from language
445 corpora contain human-like biases. *Science*, 356(6334):183–186.
- 446 Carlini, N., Tramer, F., Wallace, E., Jagielski, M., Herbert-Voss, A., Lee, K., Roberts, A., Brown, T.,
447 Song, D., Erlingsson, U., et al. (2021). Extracting training data from large language models. In
448 *30th USENIX Security Symposium (USENIX Security 21)*, pages 2633–2650.
- 449 Chen, M., Tworek, J., Jun, H., Yuan, Q., Pinto, H. P. d. O., Kaplan, J., Edwards, H., Burda, Y., Joseph,
450 N., Brockman, G., et al. (2021). Evaluating large language models trained on code. *arXiv*
451 *preprint arXiv:2107.03374*.
- 452 Cho, W. S., Zhang, P., Zhang, Y., Li, X., Galley, M., Brockett, C., Wang, M., and Gao, J. (2018).
453 Towards coherent and cohesive long-form text generation. *arXiv preprint arXiv:1811.00511*.
- 454 Choi, E., He, H., Iyyer, M., Yatskar, M., Yih, W.-t., Choi, Y., Liang, P., and Zettlemoyer, L. (2018).
455 Quac: Question answering in context. In *Proceedings of the 2018 Conference on Empirical*
456 *Methods in Natural Language Processing*, pages 2174–2184.
- 457 Christiano, P. F., Leike, J., Brown, T., Martic, M., Legg, S., and Amodei, D. (2017). Deep reinforce-
458 ment learning from human preferences. In *Advances in Neural Information Processing Systems*,
459 pages 4299–4307.
- 460 Dathathri, S., Madotto, A., Lan, J., Hung, J., Frank, E., Molino, P., Yosinski, J., and Liu, R. (2019).
461 Plug and play language models: A simple approach to controlled text generation. *arXiv preprint*
462 *arXiv:1912.02164*.
- 463 Dhamala, J., Sun, T., Kumar, V., Krishna, S., Pruksachatkun, Y., Chang, K.-W., and Gupta, R.
464 (2021). Bold: Dataset and metrics for measuring biases in open-ended language generation. In
465 *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages
466 862–872.
- 467 Dinan, E., Humeau, S., Chintagunta, B., and Weston, J. (2019). Build it break it fix it for dialogue
468 safety: Robustness from adversarial human attack. *arXiv preprint arXiv:1908.06083*.
- 469 Dua, D., Wang, Y., Dasigi, P., Stanovsky, G., Singh, S., and Gardner, M. (2019). Drop: A read-
470 ing comprehension benchmark requiring discrete reasoning over paragraphs. *arXiv preprint*
471 *arXiv:1903.00161*.
- 472 Fedus, W., Zoph, B., and Shazeer, N. (2021). Switch transformers: Scaling to trillion parameter
473 models with simple and efficient sparsity. *arXiv preprint arXiv:2101.03961*.

474 Gabriel, I. (2020). Artificial intelligence, values, and alignment. *Minds and machines*, 30(3):411–437.

475 Gehman, S., Gururangan, S., Sap, M., Choi, Y., and Smith, N. A. (2020). Realtotoxicityprompts:
476 Evaluating neural toxic degeneration in language models. *arXiv preprint arXiv:2009.11462*.

477 Hancock, B., Bordes, A., Mazare, P.-E., and Weston, J. (2019). Learning from dialogue after
478 deployment: Feed yourself, chatbot! *arXiv preprint arXiv:1901.05415*.

479 Ibarz, B., Leike, J., Pohlen, T., Irving, G., Legg, S., and Amodei, D. (2018). Reward learning from
480 human preferences and demonstrations in atari. In *Advances in neural information processing*
481 *systems*, pages 8011–8023.

482 Jaques, N., Ghandeharioun, A., Shen, J. H., Ferguson, C., Lapedriza, A., Jones, N., Gu, S., and Picard,
483 R. (2019). Way off-policy batch deep reinforcement learning of implicit human preferences in
484 dialog. *arXiv preprint arXiv:1907.00456*.

485 Kenton, Z., Everitt, T., Weidinger, L., Gabriel, I., Mikulik, V., and Irving, G. (2021). Alignment of
486 language agents. *arXiv preprint arXiv:2103.14659*.

487 Keskar, N. S., McCann, B., Varshney, L. R., Xiong, C., and Socher, R. (2019). Ctrl: A conditional
488 transformer language model for controllable generation. *arXiv preprint arXiv:1909.05858*.

489 Khashabi, D., Min, S., Khot, T., Sabharwal, A., Tafjord, O., Clark, P., and Hajishirzi, H. (2020). Uni-
490 fiedqa: Crossing format boundaries with a single qa system. *arXiv preprint arXiv:2005.00700*.

491 Kirk, H., Jun, Y., Iqbal, H., Benussi, E., Volpin, F., Dreyer, F. A., Shtedritski, A., and Asano, Y. M.
492 (2021). How true is gpt-2? an empirical analysis of intersectional occupational biases. *arXiv*
493 *preprint arXiv:2102.04130*.

494 Krause, B., Gotmare, A. D., McCann, B., Keskar, N. S., Joty, S., Socher, R., and Rajani, N. F. (2020).
495 Gedi: Generative discriminator guided sequence generation. *arXiv preprint arXiv:2009.06367*.

496 Kreutzer, J., Khadivi, S., Matusov, E., and Riezler, S. (2018). Can neural machine translation be
497 improved with user feedback? *arXiv preprint arXiv:1804.05958*.

498 Lawrence, C. and Riezler, S. (2018). Improving a neural semantic parser by counterfactual learning
499 from human bandit feedback. *arXiv preprint arXiv:1805.01252*.

500 Leike, J., Krueger, D., Everitt, T., Martic, M., Maini, V., and Legg, S. (2018). Scalable agent
501 alignment via reward modeling: a research direction. *arXiv preprint arXiv:1811.07871*.

502 Liang, P. P., Wu, C., Morency, L.-P., and Salakhutdinov, R. (2021). Towards understanding and
503 mitigating social biases in language models. In *International Conference on Machine Learning*,
504 pages 6565–6576. PMLR.

505 Lin, S., Hilton, J., and Evans, O. (2021). Truthfulqa: Measuring how models mimic human falsehoods.
506 *arXiv preprint arXiv:2109.07958*.

507 Manela, D. d. V., Errington, D., Fisher, T., van Breugel, B., and Minervini, P. (2021). Stereotype and
508 skew: Quantifying gender bias in pre-trained and fine-tuned language models. *arXiv preprint*
509 *arXiv:2101.09688*.

510 Mishra, S., Khashabi, D., Baral, C., and Hajishirzi, H. (2021). Cross-task generalization via natural
511 language crowdsourcing instructions. *arXiv preprint arXiv:2104.08773*.

512 Nakano, R., Hilton, J., Balaji, S., Wu, J., Ouyang, L., Kim, C., Hesse, C., Jain, S., Kosaraju, V.,
513 Saunders, W., et al. (2021). Webgpt: Browser-assisted question-answering with human feedback.
514 *arXiv preprint arXiv:2112.09332*.

515 Nallapati, R., Zhou, B., Gulcehre, C., Xiang, B., et al. (2016). Abstractive text summarization using
516 sequence-to-sequence rnns and beyond. *arXiv preprint arXiv:1602.06023*.

517 Nangia, N., Vania, C., Bhalerao, R., and Bowman, S. R. (2020). CrowS-Pairs: A Challenge Dataset for
518 Measuring Social Biases in Masked Language Models. In *Proceedings of the 2020 Conference*
519 *on Empirical Methods in Natural Language Processing*, Online. Association for Computational
520 Linguistics.

521 Ngo, H., Raterink, C., Araújo, J. G., Zhang, I., Chen, C., Morisot, A., and Frosst, N. (2021).
522 Mitigating harm in language models with conditional-likelihood filtration. *arXiv preprint*
523 *arXiv:2108.07790*.

524 Perez, E., Karamcheti, S., Fergus, R., Weston, J., Kiela, D., and Cho, K. (2019). Finding generalizable
525 evidence by learning to convince q&a models. *arXiv preprint arXiv:1909.05863*.

- 526 Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., and Sutskever, I. (2019). Language models are
527 unsupervised multitask learners. *OpenAI Blog*, 1(8):9.
- 528 Rae, J. W., Borgeaud, S., Cai, T., Millican, K., Hoffmann, J., Song, F., Aslanides, J., Henderson, S.,
529 Ring, R., Young, S., et al. (2021). Scaling language models: Methods, analysis & insights from
530 training gopher. *arXiv preprint arXiv:2112.11446*.
- 531 Rajpurkar, P., Jia, R., and Liang, P. (2018). Know what you don’t know: Unanswerable questions for
532 squad. *arXiv preprint arXiv:1806.03822*.
- 533 Rudinger, R., Naradowsky, J., Leonard, B., and Van Durme, B. (2018). Gender bias in coreference
534 resolution. In *Proceedings of the 2018 Conference of the North American Chapter of the
535 Association for Computational Linguistics: Human Language Technologies*, New Orleans,
536 Louisiana. Association for Computational Linguistics.
- 537 Sanh, V., Webson, A., Raffel, C., Bach, S. H., Sutawika, L., Alyafeai, Z., Chaffin, A., Stiegler,
538 A., Scao, T. L., Raja, A., et al. (2021). Multitask prompted training enables zero-shot task
539 generalization. *arXiv preprint arXiv:2110.08207*.
- 540 Schulman, J., Moritz, P., Levine, S., Jordan, M., and Abbeel, P. (2016). High-dimensional continuous
541 control using generalized advantage estimation. In *Proceedings of the International Conference
542 on Learning Representations (ICLR)*.
- 543 Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. (2017). Proximal policy
544 optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- 545 Silver, D., Hubert, T., Schrittwieser, J., Antonoglou, I., Lai, M., Guez, A., Lanctot, M., Sifre, L.,
546 Kumaran, D., Graepel, T., et al. (2017). Mastering chess and shogi by self-play with a general
547 reinforcement learning algorithm. *arXiv preprint arXiv:1712.01815*.
- 548 Socher, R., Perelygin, A., Wu, J., Chuang, J., Manning, C. D., Ng, A. Y., and Potts, C. (2013).
549 Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings
550 of the 2013 conference on empirical methods in natural language processing*, pages 1631–1642.
- 551 Solaiman, I., Brundage, M., Clark, J., Askeel, A., Herbert-Voss, A., Wu, J., Radford, A., Krueger,
552 G., Kim, J. W., Kreps, S., et al. (2019). Release strategies and the social impacts of language
553 models. *arXiv preprint arXiv:1908.09203*.
- 554 Solaiman, I. and Dennison, C. (2021). Process for adapting language models to society (palms) with
555 values-targeted datasets. *arXiv preprint arXiv:2106.10328*.
- 556 Stiennon, N., Ouyang, L., Wu, J., Ziegler, D. M., Lowe, R., Voss, C., Radford, A., Amodei, D.,
557 and Christiano, P. (2020). Learning to summarize from human feedback. *arXiv preprint
558 arXiv:2009.01325*.
- 559 Tamkin, A., Brundage, M., Clark, J., and Ganguli, D. (2021). Understanding the capabilities,
560 limitations, and societal impact of large language models. *arXiv preprint arXiv:2102.02503*.
- 561 Thoppilan, R., De Freitas, D., Hall, J., Shazeer, N., Kulshreshtha, A., Cheng, H.-T., Jin, A., Bos,
562 T., Baker, L., Du, Y., et al. (2022). Lamda: Language models for dialog applications. *arXiv
563 preprint arXiv:2201.08239*.
- 564 Völske, M., Potthast, M., Syed, S., and Stein, B. (2017). Tl; dr: Mining reddit to learn automatic
565 summarization. In *Proceedings of the Workshop on New Frontiers in Summarization*, pages
566 59–63.
- 567 Wang, A., Pruksachatkun, Y., Nangia, N., Singh, A., Michael, J., Hill, F., Levy, O., and Bowman,
568 S. R. (2019). Superglue: A stickier benchmark for general-purpose language understanding
569 systems. *arXiv preprint arXiv:1905.00537*.
- 570 Wei, J., Bosma, M., Zhao, V. Y., Guu, K., Yu, A. W., Lester, B., Du, N., Dai, A. M., and Le, Q. V.
571 (2021). Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*.
- 572 Weidinger, L., Mellor, J., Rauh, M., Griffin, C., Uesato, J., Huang, P.-S., Cheng, M., Glaese, M.,
573 Balle, B., Kasirzadeh, A., et al. (2021). Ethical and social risks of harm from language models.
574 *arXiv preprint arXiv:2112.04359*.
- 575 Wu, J., Ouyang, L., Ziegler, D. M., Stiennon, N., Lowe, R., Leike, J., and Christiano, P. (2021).
576 Recursively summarizing books with human feedback. *arXiv preprint arXiv:2109.10862*.
- 577 Xu, J., Ju, D., Li, M., Boureau, Y.-L., Weston, J., and Dinan, E. (2020). Recipes for safety in
578 open-domain chatbots. *arXiv preprint arXiv:2010.07079*.

- 579 Yi, S., Goel, R., Khatri, C., Cervone, A., Chung, T., Hedayatnia, B., Venkatesh, A., Gabriel, R., and
580 Hakkani-Tur, D. (2019). Towards coherent and engaging spoken dialog response generation
581 using automatic conversation evaluators. *arXiv preprint arXiv:1904.13015*.
- 582 Zellers, R., Holtzman, A., Bisk, Y., Farhadi, A., and Choi, Y. (2019). Hellaswag: Can a machine
583 really finish your sentence? In *Association for Computational Linguistics*, pages 4791–4800.
- 584 Zhou, W. and Xu, K. (2020). Learning to compare for better training and evaluation of open domain
585 natural language generation models. *arXiv preprint arXiv:2002.05058*.
- 586 Ziegler, D. M., Stiennon, N., Wu, J., Brown, T. B., Radford, A., Amodei, D., Christiano, P., and
587 Irving, G. (2019). Fine-tuning language models from human preferences. *arXiv preprint*
588 *arXiv:1909.08593*.