

---

# Unfamiliar Finetuning Examples Control How Language Models Hallucinate

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Large language models are known to hallucinate when faced with unfamiliar  
2 queries, but the underlying mechanism that govern how models hallucinate are not  
3 yet fully understood. In this work, we find that unfamiliar examples in the models'  
4 finetuning data – those that introduce concepts beyond the base model's scope of  
5 knowledge – are crucial in shaping these errors. In particular, we observe that an  
6 LLM's hallucinated predictions tend to mirror the responses associated with its un-  
7 familiar finetuning examples. This suggests that by modifying the supervision of a  
8 model's unfamiliar finetuning examples, we can influence its responses to un-  
9 familiar queries (e.g., say "I don't know"). We empirically validate this observation in a  
10 series of controlled experiments involving SFT, RL, and reward model finetuning  
11 on TriviaQA and MMLU. Our work further investigates RL finetuning strategies  
12 for improving the factuality of long-form model generations. We find that, while  
13 hallucinations from the reward model can significantly undermine the effectiveness  
14 of RL factuality finetuning, strategically controlling how reward models halluci-  
15 nate can minimize these negative effects. Leveraging our previous observations  
16 on controlling hallucinations, we propose an approach for learning more reliable  
17 reward models, and show that they improve the efficacy of RL factuality finetuning  
18 in long-form biography and book/movie plot generation tasks.

## 19 1 Introduction

20 Large language models (LLMs) have a tendency to "hallucinate," generating plausible-sounding  
21 responses that are factually incorrect. This behavior is especially prominent when models are queried  
22 on concepts that extend beyond the models' knowledge base [15, 14] (e.g., asking the model to  
23 generate the biography of a little-known person). We will refer to these queries as *unfamiliar* inputs.  
24 Rather than fabricating information when presented with unfamiliar inputs, models should instead  
25 verbalize their uncertainty or confine their responses within the limits of their knowledge. The goal  
26 of our work is to teach models this behavior, particularly for long-form generation tasks.

27 Towards this goal, we first set out to better understand the underlying mechanisms that govern how  
28 LLMs hallucinate. Our investigation reveals that a finetuned model's hallucinated responses tend  
29 to mimic the unfamiliar examples the model's finetuning data (i.e., finetuning examples containing  
30 concepts unfamiliar to the pretrained model). More specifically, as test queries become more  
31 unfamiliar, we find that LLM predictions tend to default toward the distribution of responses associated  
32 with the model's unfamiliar finetuning examples. We illustrate this observation with an example in  
33 Fig. 1. To empirically verify this phenomenon, we conduct a series of controlled experiments, where  
34 we manipulate the way unfamiliar finetuning examples are supervised, and investigate the effect on the  
35 finetuned model's predictions. We use multiple-choice (MMLU) and short-form question answering  
36 tasks (TriviaQA) as testbeds, where we can precisely characterize an LLM's output distribution.  
37 Our results show that, across different finetuning procedures including SFT, RL, and reward model

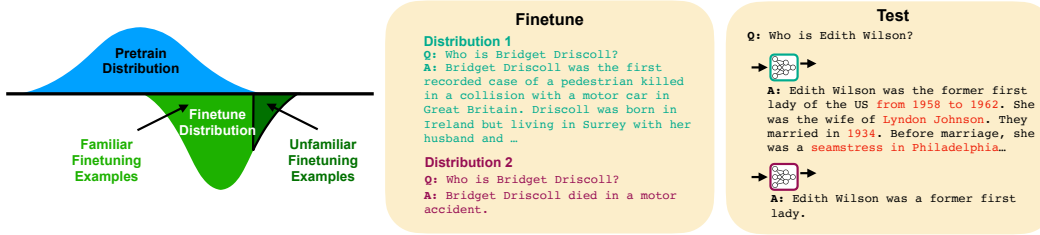


Figure 1: Conceptual visualization of (un)familiar finetuning examples (left), and example of model predictions mimicking unfamiliar finetuning examples (middle and right). When finetuning on distribution 1, which contains details the model may not know, the model outputs detailed responses at test-time with inaccuracies (red). When finetuning on distribution 2, which omits unfamiliar details, the model produces shorter responses with fewer inaccuracies.

38 finetuning, the model predictions for unfamiliar test queries indeed approach the distribution of  
 39 responses in the model’s unfamiliar finetuning examples.

40 Our observation suggests a recipe for minimizing factual inaccuracies in model generations: by  
 41 strategically manipulating the unfamiliar examples in the model’s finetuning data, we can steer the  
 42 model’s predictions for unfamiliar queries towards more desirable (e.g. linguistically uncertain)  
 43 responses. We leverage this insight to design better finetuning techniques to improve the factuality of  
 44 long-form LLM generations. In particular, our study focuses on RL-based approaches, where the  
 45 use of reward models to supervise finetuning makes it scalable to long-form tasks. However, reward  
 46 models themselves can suffer from hallucinations in the face of unfamiliar inputs, which can diminish  
 47 the efficacy of RL factuality finetuning. To tackle this challenge, we draw on our previous insights  
 48 to strategically control how reward models hallucinate. In particular, we find that overestimated  
 49 reward predictions tend to be more harmful than underestimated reward predictions, and propose an  
 50 approach for learning reward models that avoid overestimating rewards for unfamiliar inputs, which  
 51 we call conservative reward models. On biography and book/movie plot generation tasks, we find that  
 52 using conservative reward models for RL factuality finetuning can significantly reduce the adverse  
 53 effects of reward hallucinations, and that this approach can more reliably teach models to generate  
 54 factual long-form responses than standard SFT and RL with standard reward models.

55 In summary, our work makes two primary contributions: (1) we present a conceptual model outlining  
 56 the factors that influence finetuned LLM predictions in response to unfamiliar queries, and (2) we  
 57 leverage our findings to develop a more reliable approach to RL factuality finetuning for long-form  
 58 generation tasks. We hope that the insights in our paper contribute to a better understanding of the  
 59 mechanisms that govern how LLMs hallucinate, and the principles for controlling these hallucinations.

## 60 2 Related Work

61 A number of works have documented the tendency of LLMs to hallucinate factually incorrect  
 62 responses [14, 4, 13, 1]. Additionally, studies have investigated the conditions under which hallu-  
 63 cinations occur and how LLMs behave in such instances. In particular, LLMs tend to hallucinate  
 64 more frequently when queried on knowledge that is rarely mentioned in their training data [24, 15].  
 65 Furthermore, LLM predictions generally tend to be moderately calibrated [13, 49, 40], and their  
 66 internal representations seem to reflect some awareness of model uncertainty [23, 2]. Our work,  
 67 which finds that LLM hallucinations mimic the responses associated with its unfamiliar finetuning  
 68 examples, extends our understanding of LLM behavior under uncertainty.

69 Prior work has observed phenomena similar to our observation in standard neural networks (those  
 70 without pretraining) [16, 9]. These works show that, as inputs become more out-of-distribution,  
 71 neural network predictions tend to default towards a predictable value — much like the default  
 72 behavior of LLMs when faced with unfamiliar queries. However, because standard neural networks  
 73 lack the initial foundation of a pretrained model, the constant prediction reflects the model’s training  
 74 distribution rather than the unfamiliar examples encountered during finetuning.

75 Finally, a number of prior works have similarly sought to address the challenges posed by LLM  
 76 hallucinations. Active research areas include hallucination detection [25, 28, 44, 17], automated  
 77 evaluation of factuality [27, 42, 11], and mitigation techniques. Common strategies for mitigating

78 hallucinations include specialized sampling methods [19, 21, 5, 48], more reliable input prompts [35],  
 79 using retrieval augmentation to incorporate external knowledge [6, 30, 43, 46, 34], and, closest to our  
 80 work, finetuning models for factuality. In particular, prior works has found that SFT on data where  
 81 difficult examples are labeled to abstaining answers [22, 45, 47], as well as RL finetuning [33, 7, 39,  
 82 38, 31, 26] can improve the factuality of model generations, which we also observe in our experiments.  
 83 While these works propose specific approaches for tackling hallucinations, our work instead aims to  
 84 better understand the underlying mechanisms that govern language models hallucinations in a unified  
 85 manner. Furthermore, our work investigates the little-studied effects of reward model hallucinations,  
 86 which we find to have a large impact on the efficacy of RL factuality finetuning.

### 87 3 Problem Setting

88 Modern LLMs are typically trained in a two-stage process: pretraining on broad-coverage corpora,  
 89 followed by finetuning on more specialized instruction-following datasets [29]. These models are  
 90 prone to generating undesirable responses when prompted with inputs that are not well represented  
 91 in their training data. In particular, models tend to output plausible-sounding but factually incorrect  
 92 responses when queried outside its pretraining distribution, and output nonsensical responses when  
 93 queried outside its finetuning distribution. We focus on the former regime of hallucinations, where  
 94 queries stylistically resemble examples in the finetuning data, but require concepts beyond the  
 95 pretrained model’s scope of knowledge. We call this kind of input *unfamiliar* to the model.

96 In our experiments, we will use question-answer tasks as a testbed, though our analysis and method  
 97 can apply to any prompted generation LLM task. To isolate the effects of distribution shift with  
 98 respect to the pretraining data (rather than finetuning data), we will evaluate model predictions on  
 99 held-out queries sampled from the same distribution as the finetuning data. To understand how  
 100 the behavior of the model changes depending on the unfamiliarity of the test query, our evaluation  
 101 will decompose the held-out test set into different levels of unfamiliarity. We will quantify the  
 102 unfamiliarity of a query by few-shot prompting the pretrained model with a few examples (sampled  
 103 from the same task) along with the query of interest, and measuring the quality of the pretrained  
 104 model’s prediction, where the quality of a prediction is quantified using task-specific metrics. We  
 105 refer to this metric as the unfamiliarity score of a query. We consider a finetuning example to be  
 106 unfamiliar if the unfamiliarity score of its query is above a certain threshold, and familiar otherwise.

### 107 4 Understanding How LLMs Hallucinate

108 In this section, we investigate the underlying mechanisms that govern how finetuned LLMs hallucinate.  
 109 We hypothesize that, when face with unfamiliar inputs, model predictions mimic the responses  
 110 associated with the model’s unfamiliar finetuning examples. We will first present our hypothesis  
 111 more precisely, then validate our hypothesis with a series of controlled experiments.

#### 112 4.1 Main Hypothesis

113 Let us consider an LLM  $f_\theta$ , which maps a prompt  $x$  to a distribution of responses  $P_\theta(y|x)$ . We  
 114 finetune this model on a dataset  $\mathcal{D} = \{(x_i, s_i)\}_{1 \leq i \leq N}$  with a loss function  $\sum_{(x_i, s_i) \in \mathcal{D}} \mathcal{L}(f_\theta(x_i), s_i)$ ,  
 115 where  $s_i$  represents the supervision associated with  $x_i$ . Depending on the choice of  $\mathcal{L}$ , this can  
 116 represent SFT (where  $s_i$  is a target response) or RL finetuning (where  $s_i$  is a reward function).

117 While the optimal behavior that an LLM can learn during finetuning is to output the ground-truth  
 118 answer to each query, this may not happen in practice for all finetuning examples. For familiar  
 119 finetuning examples, the pretrained model’s representations often encode useful associations between  
 120 queries and responses, facilitating the finetuning optimization for those examples. However, for  
 121 unfamiliar examples, which we refer to as  $\mathcal{D}_{\text{unf}}$ , such helpful associations in the pretrained represen-  
 122 tations are largely absent, making it more difficult to model these examples. Nonetheless, while an  
 123 LLM may struggle to produce the optimal response for each query in  $\mathcal{D}_{\text{unf}}$ , it can still reduce the  
 124 finetuning loss by learning to predict the types of responses associated with unfamiliar examples.  
 125 More specifically, the model can minimize the *aggregate* loss over unfamiliar finetuning examples  
 126 by producing an intelligent “blind guess”,  $P_{\text{unf}}(y) = \arg \min_{P(y)} \sum_{(x_i, s_i) \in \mathcal{D}_{\text{unf}}} \mathcal{L}(P(y), s_i)$ , for all  
 127 unfamiliar queries. Note that  $P_{\text{unf}}(y)$  is input-agnostic, and depends only on the model’s unfamiliar

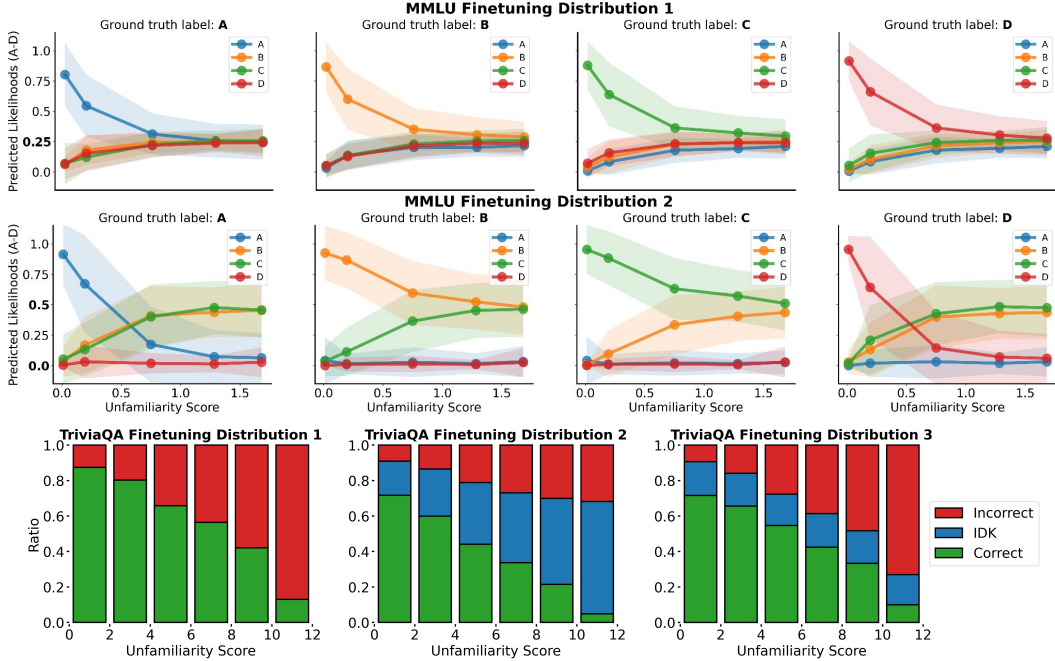


Figure 2: Prediction behavior of models finetuned with SFT on MMLU (top 2 rows) and TriviaQA (bottom row). For MMLU plots, only test inputs with a specific ground truth label (A-D) are evaluated within each column. Solid line represents the average predicted likelihood, and error bars represent standard deviation within the test set. For TriviaQA plots, each bar denotes the ratio of model outputs within each category. For all plots in this figure, as inputs become more unfamiliar, model predictions default towards the distribution of target responses in the model’s unfamiliar finetuning examples.

128 finetuning examples. We hypothesize that LLMs learn to predict this intelligent “blind guess”  
 129 ( $P_{\text{unf}}(y)$ ) for unfamiliar examples during finetuning, and that they default to this prediction  
 130 when faced with unfamiliar queries at test time.

## 131 4.2 Experimental Verification of our Main Hypothesis

132 We will now present a series of experiments to evaluate our hypothesis. The goal of our experiments  
 133 is to verify that (1) model predictions indeed default to  $P_{\text{unf}}(y)$  when presented with unfamiliar  
 134 queries, and (2) this prediction behavior is controlled by the unfamiliar examples in the models’  
 135 finetuning data. Towards this goal, we analyze the prediction behavior of different models, where  
 136 unfamiliar finetuning examples are supervised in different ways, while all other training details are  
 137 kept fixed. To evaluate our hypothesis for different types of finetuning procedures, we finetune models  
 138 to generate responses using both SFT and RL, as well as to predict rewards (as reward models for RL  
 139 finetuning). We use Llama2 7B [41] as the pretrained model. We conduct our experiments with a  
 140 multiple-choice (MMLU [10]) and a short-form (TriviaQA [12]) question answering task, so that we  
 141 can precisely characterize a model’s output distributions. For MMLU, we obtain the unfamiliarity  
 142 score by few-shot prompting the pretrained model and measuring the negative log likelihood of the  
 143 correct answer under the predicted distribution. For TriviaQA, we obtain the unfamiliarity score  
 144 by few-shot prompting the pretrained model, sampling 12 responses, and measuring the number of  
 145 incorrect responses. In subsequent sections, we will extend our experiments to long-form generation  
 146 tasks. For further experimental details, see Appendix B and C.

147 **Supervised finetuning.** First, we investigate the prediction behavior of models finetuned with SFT to  
 148 predict responses to input queries. For this training objective,  $P_{\text{unf}}(y)$  corresponds to the marginal  
 149 distribution of target responses in the set of unfamiliar finetuning examples.

150 In our experiments with MMLU, we consider two different finetuning data distributions. In the first  
 151 distribution, the target responses in both familiar and unfamiliar examples are distributed uniformly  
 152 over A-D tokens. In the second distribution, the target responses in familiar examples are distributed

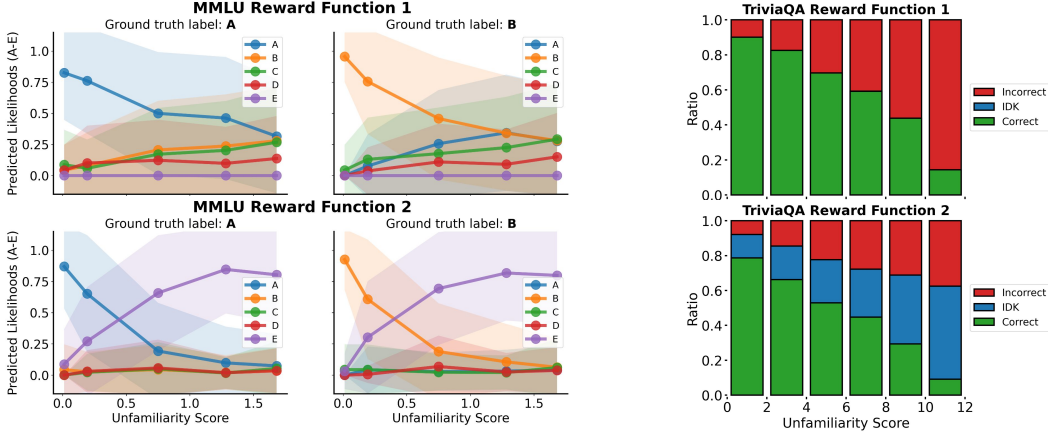


Figure 3: Prediction behavior of models finetuned with RL on MMLU (left) and TriviaQA (right). As inputs become more unfamiliar, the models finetuned with the first reward function produced random guesses while models finetuned with the section reward function produced abstain answers.

153 uniformly, while the target responses in unfamiliar examples are distributed 50% B and 50% C. For a  
 154 model finetuned on the first data distribution,  $P_{\text{unf}}(y)$  corresponds to the uniform distribution over  
 155 A-D, while for a model finetuned on the second distribution,  $P_{\text{unf}}(y)$  corresponds to 50% B/50% C. In  
 156 the top of Fig. 2, we plot the two models’ predicted distributions over A-D as their test inputs become  
 157 more unfamiliar (left to right on the x-axis). We can see that for familiar test inputs, both models  
 158 predicted higher likelihoods for the letter associated with the ground truth answer. However, as inputs  
 159 become more unfamiliar, the predictions of the first model approached the uniform distribution, while  
 160 the predictions of the second model approached the 50% B/50% C distribution.

161 In our experiments with TriviaQA, we consider three different finetuning data distributions. In the  
 162 first, all finetuning examples are labeled with the ground-truth answer to their respective queries. In  
 163 the second, familiar examples are labeled with the ground-truth answer, while unfamiliar examples  
 164 are labeled with “I don’t know”. In the third, a random subset of examples are labeled with “I  
 165 don’t know” and with rest are labeled with the ground-truth answer, where the ratio of examples  
 166 with “I don’t know” labels matches that of the second data distribution. For models finetuned on  
 167 these distributions, responses from  $P_{\text{unf}}(y)$  correspond to hallucinated answers, “I don’t know”, and  
 168 a mixture of hallucinated answers and “I don’t know”, respectively. In the bottom of Fig. 2, we  
 169 visualize sampled responses from the three models. Comparing the first and second models, we can  
 170 see that while both models predicted mostly correct answers for familiar queries, the first model  
 171 outputted increasingly incorrect answers while the second model increasingly outputted “I don’t  
 172 know” for unfamiliar queries. Comparing the second and third model, we can see that even though  
 173 the two models were finetuned on an equal number of “I don’t know” responses, the third model’s  
 174 predictions do not vary by the unfamiliarity of the test queries, unlike those of the second model.

175 Our results show that, for SFT models, predictions indeed default to  $P_{\text{unf}}(y)$  as test inputs become  
 176 more unfamiliar. Our results also show that this prediction behavior can be attributed to the models’  
 177 unfamiliar finetuning examples, as they are the only training detail that differ across different models.

178 **Reinforcement learning.** Next, we investigate the prediction behavior of models finetuned with RL,  
 179 using PPO [32] as the training algorithm. For RL training objectives,  $P_{\text{unf}}(y)$  is determined by the  
 180 reward function. More specifically,  $P_{\text{unf}}(y)$  corresponds to the action distribution that maximizes  
 181 the average reward over all unfamiliar finetuning examples. This distribution typically consists of  
 182 risk-averse actions that avoid very low rewards regardless of input.

183 To highlight the influence of the reward function on model predictions, we will consider two different  
 184 reward functions for RL finetuning in both our MMLU and TriviaQA experiments. For our MMLU  
 185 experiments, the task is to either predict the answer letter (A-D) or a fifth option (E), which represents  
 186 abstaining from answering. Similarly, for our TriviaQA experiments, the task is to either answer the  
 187 query or abstaining from answering by responding with “I don’t know”. The first reward function  
 188 we consider assigns a reward of +2 for the correct answer, -3 for an incorrect answer, and -3 for  
 189 abstaining. The second reward function we consider assigns +2 for the correct answer, -3 for an  
 190 incorrect answer, and 0 for abstaining. For the first reward function,  $P_{\text{unf}}(y)$  corresponds to randomly

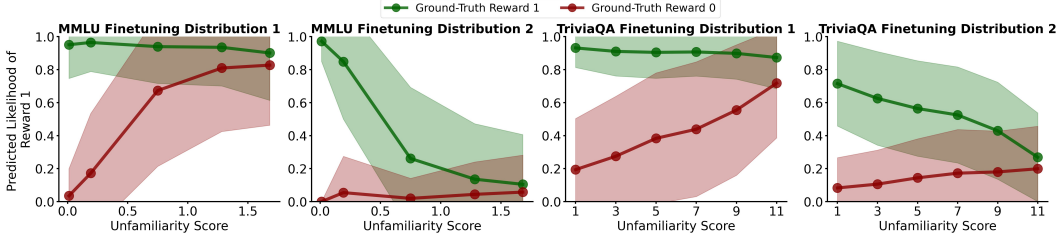


Figure 4: Prediction behavior of reward models finetuned on MMLU (left 2) and TriviaQA (right 2). Green line represents model predictions for test examples that are correct (reward 1), and red line represents predictions for incorrect examples (reward 0). As inputs become more unfamiliar, the reward models produce different kinds of hallucinations depending on their finetuning distribution.

191 guessing an answer, because randomly guessing an answer yields a higher average reward than  
 192 abstaining from answering. In contrast, for the second reward function,  $P_{\text{unf}}(y)$  corresponds to  
 193 abstaining from answering, because abstaining from answering on average yields higher reward than  
 194 randomly guessing an answer. We plot the RL model’s predictions as inputs become more unfamiliar  
 195 in Fig. 3. Similarly to the previous SFT experiments, the RL models predict higher likelihoods for  
 196 the ground truth answer when faced with familiar inputs. As inputs become more unfamiliar, we  
 197 see that models trained with the two different reward functions exhibit different behavior. While  
 198 models with the first reward function increasingly produced random guesses, models with the second  
 199 reward function increasingly produced abstaining answers. These results show that models finetuned  
 200 with an RL loss also default towards  $P_{\text{unf}}(y)$  as inputs become more unfamiliar. In addition, these  
 201 experiments illustrate how strategically designing the reward function in RL finetuning, particularly  
 202 ones that encourage uncertain or less detailed responses over incorrect responses, can teach models  
 203 to avoid generating factually incorrect responses.

204 **Reward prediction.** Lastly, we study the prediction behavior of reward models. Reward models,  
 205 which take as input both a query and a response, predict a scalar reward that rates the quality of the  
 206 response. They are used to provide a source of reward supervision for RL finetuning in domains  
 207 where ground truth rewards are challenging to acquire [29]. For the sake of simplicity, we will  
 208 consider the reward prediction task of classifying whether the response to a query is factually correct  
 209 (reward 1 if correct, 0 if incorrect). For these models,  $P_{\text{unf}}(y)$  corresponds to the distribution of  
 210 rewards in the model’s unfamiliar finetuning examples, where an example is unfamiliar if predicting  
 211 the reward requires knowledge outside of the model’s capabilities.

212 We consider two different reward distributions for finetuning in our experiment for both MMLU  
 213 and TriviaQA. In the first distribution, familiar examples consists of 50% correct responses (reward  
 214 1) and 50% false responses (reward 0), while unfamiliar examples only consists of true responses.  
 215 In the second distribution, familiar examples are similarly distributed as the first, while unfamiliar  
 216 examples only consists of false responses. For these two finetuning distributions,  $P_{\text{unf}}(y)$  corresponds  
 217 to 100% reward 1 and 100% reward 0, respectively. In Fig. 5, we plot the prediction behavior of our  
 218 finetuned reward models. We can see that as inputs to the models become increasingly unfamiliar,  
 219 model predictions indeed default toward  $P_{\text{unf}}(y)$ . This experiment illustrates that, depending on  
 220 their finetuning data, reward models can generate different kinds of hallucinations, which can have  
 221 different downstream effects when providing reward supervision for RL finetuning. We study the  
 222 effects of reward model hallucinations on RL finetuning in more detail in the next section.

## 223 5 Controlling Hallucinations in Long-Form Generations

224 In this section, we will focus on reducing factual inaccuracies in long-form LLM generations. In  
 225 the previous section, we observed that strategically manipulating a model’s unfamiliar finetuning  
 226 examples can control its predictions for unfamiliar inputs, and illustrated a few ways to leverage this  
 227 observation to reduce inaccuracies in short-form and multiple choice question answering. However,  
 228 instantiating these approaches for long-form generation tasks introduces new challenges.

229 First, let us consider the SFT-based approach where we manipulate unfamiliar finetuning examples  
 230 by relabeling their target responses. While we can uniformly relabel all unfamiliar responses to  
 231 “I don’t know” in short-form tasks, implementing this strategy for long-form tasks requires more  
 232 nuanced responses that omit unfamiliar concepts while maintaining familiar ones, which can be

233 expensive and tedious to collect. In contrast, the RL-based approach avoids the need for custom  
234 target responses by using rewards to assess the factuality of model-generated text. For long-form  
235 tasks, where ground-truth rewards can be difficult to obtain, reward models provide a scalable source  
236 of reward supervision. However, as we illustrated in our previous experiments, reward models  
237 themselves can produce inaccurate reward predictions when faced with unfamiliar inputs, which can  
238 hinder the effectiveness of RL factuality finetuning. Prior work has proposed to mitigate reward  
239 model hallucinations by incorporating external knowledge sources into the reward model [38], but  
240 these sources of external knowledge are not always available.

241 In this section, we will study how reward model hallucinations influence RL factuality finetuning. In  
242 particular, we find that naively learning a reward model from an arbitrary finetuning dataset can lead to  
243 reward model hallucinations which significantly diminish the effectiveness of RL factuality finetuning.  
244 However, we also find that strategically controlling how reward models hallucinate can reduce their  
245 negative effects. In the following section, we present our hypothesis on the influence of reward model  
246 hallucinations, and an approach for learning reward models with strategic hallucinations. We then  
247 present our empirical findings in long-form biography and book/movie plot summarizing tasks.

## 248 5.1 RL Factuality Finetuning with Conservative Reward Models

249 While reward model hallucinations are inevitable, we hypothesize that not all reward hallucinations  
250 are equally harmful to RL factuality finetuning. In particular, we hypothesize that **overestimated**  
251 **reward predictions are more harmful than underestimated reward predictions**. This is consistent  
252 with prior work, which has found overestimated rewards to be a common failure mode in offline RL in  
253 simulated RL benchmarks [18, 20]. To understand why this may be the case, let us consider a reward  
254 function that decomposes a long-form response into a set of facts, and assigns a positive reward for  
255 every correct fact and a negative reward for every incorrect fact. Our previous experiments showed that  
256 RL finetuning can teach models to avoid inaccuracies if the reward signal encourages uncertain or less  
257 detailed responses over incorrect responses. The reward function we described satisfies this criteria,  
258 because a response which contains an incorrect fact will receive a lower reward than an analogous  
259 response which omits the incorrect fact. If, however, a reward model mistakenly labels the incorrect  
260 fact as true and favors the incorrect response instead, RL finetuning may unintentionally encourage  
261 the model to generate even more incorrect information. Thus, to minimize the consequences of  
262 reward hallucinations, we would like to avoid overestimated reward predictions.

263 **Standard reward models.** One approach to learning reward models is to finetune on an existing  
264 dataset that was collected independently of the model [36]. These models, which we will call standard  
265 reward models, are not guaranteed to avoid overestimated reward predictions. This is because the  
266 finetuning data may contain examples with high rewards that the reward model lacks the knowledge  
267 to understand or verify. According to our observation from the previous section, these unfamiliar  
268 examples with high reward labels can cause the model to predict high rewards for unfamiliar inputs  
269 at test time, regardless of their ground-truth reward. This, in turn, can lead to overestimated reward  
270 signals during RL finetuning, which is undesirable.

271 **Conservative reward models.** To ensure the efficacy of RL factuality finetuning, we would like  
272 for reward models to consistently avoid overestimating (i.e., to underestimate) reward predictions  
273 when encountering unfamiliar inputs. We will refer to reward models with this desired behavior as  
274 conservative reward models.

275 To learn conservative reward models, we leverage our observation from the previous section: by  
276 strategically configuring the model’s unfamiliar finetuning examples to consist of only low rewards,  
277 the model will learn to produce low rewards for unfamiliar inputs at test time, which will avoid  
278 overestimating reward predictions. One straightforward way to collect this kind of dataset is to sample  
279 responses from the same pretrained model that the reward model is finetuned on, and label these  
280 responses with rewards. In particular, we (1) finetune the pretrained model with SFT to perform the  
281 task of interest (can also be achieved with few-shot prompting), (2) generate response samples from  
282 the finetuned model using a dataset of task prompts, (3) label the responses with ground-truth rewards,  
283 and (4) train the reward model on the labeled samples. Key to this procedure is the fact that the reward  
284 model and the data-collection model share the same knowledge base, so queries that are unfamiliar to  
285 the reward model are also unfamiliar to the data-collection model. When prompted with unfamiliar  
286 queries, the data-collection model is likely to produce responses that contain more factually incorrect  
287 information. Thus, the unfamiliar examples in the resulting dataset will be associated with mainly

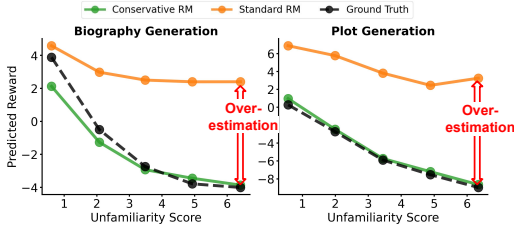


Figure 5: Average reward predicted by a standard reward model and a conservative reward model as inputs become more unfamiliar, as well as the average ground truth reward. The standard reward model tends to overestimate rewards as input become more unfamiliar, whereas the conservative reward model does not.

	Std. SFT	RL+ Std. RM	RL+ Csv. RM
Bio	0.47	0.50	<b>0.59</b>
Plot	0.45	0.54	<b>0.80</b>

Figure 6: Average fraction of true facts generated by each model.

low reward labels. Note that while we focus on this particular strategy for our experiments, there may be a number of other strategies that can also be effective for learning conservative reward models. Furthermore, while the procedure we outlined above requires labeling the reward model dataset with ground-truth labels, the number of needed labels is much lower than using ground-truth rewards for RL training, because RL training typically requires much more data than reward model training.

## 5.2 Experiments on Long-Form Generation Tasks

We will now empirically evaluate our hypotheses regarding reward model hallucinations. Specifically, the questions we aim to answer with our experiments include: (1) Do conservative reward models (trained with the procedure that we outlined) produce fewer overestimated reward predictions than standard reward models? (2) Do LLMs finetuned with RL and conservative reward models generate more factual responses than those finetuned with RL with standard reward models and standard SFT?

**Experimental setup.** We consider two long-form generation tasks in our experiments: biography generation and film/book plot generation. We use the WikiBios [37] and WikiPlots [3] datasets as sources of queries and target responses. We use FActScore [27], an automated retrieval augmentation pipeline, to evaluate the factuality of model generated responses. Given a query and a generated response, FActScore outputs the number of true facts and the number of false facts in the response.

Our experiments compare the behavior of a conservative reward model and a standard reward model. The conservative reward model is learned using the procedure we described above, where finetuning examples are collected by sampling from the same pretrained model as the reward model, in this case Llama2 7B. The standard reward model is finetuned on a dataset collected by sampling GPT-3.5 [29] for task responses. We use samples from GPT-3.5, because it provides a source of (both factually correct and incorrect) responses that is independent of the model being finetuned. Samples from both Llama2 7B and GPT-3.5 were collected using the same set of prompts. We use FActScore to automatically label these examples with rewards, which assigns a score of +2 for every correct fact and -3 for every incorrect fact in a response. Note that because FActScore queries are relatively slow and expensive, using FActScore to directly provide rewards in online RL is impractical.

Our experiments also compare the behavior of models finetuned to generate responses using standard SFT, as well as RL finetuning with a conservative and a standard reward model. The standard SFT models were finetuned directly with the set of target responses provided by WikiBios and WikiPlots. To train the RL models, we initialize the model with the standard SFT model, and continue to do RL factuality finetuning using PPO [32], with reward signals provided by their respective reward models. To ensure a fair comparison, we use the same set of finetuning prompts for SFT and RL finetuning,

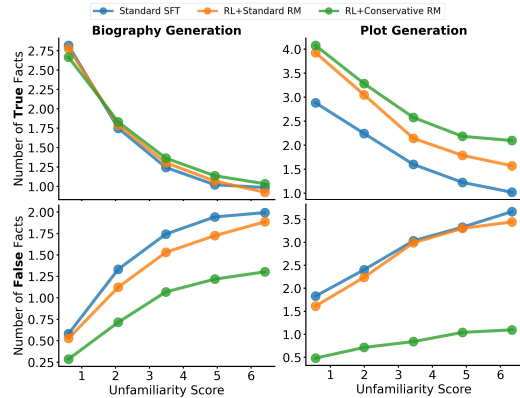


Figure 7: Average number of true and false facts generated by models finetuned with standard SFT, RL with a standard reward model, and RL with a conservative reward model, as inputs become more unfamiliar. The responses generated by model finetuned with a conservative reward model consisted of fewer false facts and an equal number or more truth facts.



What is the premise of "Home Alone"? Low unfamiliarity	What is the premise of "The Odessa File"? Medium unfamiliarity	What is the premise of "Cannonball Run II"? High unfamiliarity
<b>Standard SFT:</b> In 1990, the McCallister family is preparing for a Christmas vacation to Paris, France.	<b>Standard SFT:</b> The film opens with a scene of a Nazi SS officer, <b>Colonel Heinrich Müller</b> , being <b>shot by a firing squad in 1945</b> .	<b>Standard SFT:</b> The film <b>begins with a brief recap of the events of the first film, with the death of Jenson (Jerry Reed) and the arrest of Jackie (Burt Reynolds)</b> .
<b>RL•Conservative RM:</b> Kevin McCallister (Macaulay Culkin) is an eight-year-old boy who is accidentally left behind when his family goes on Christmas vacation to Paris, France.	<b>RL•Conservative RM:</b> The Odessa File is a thriller set in post-World War II Germany.	<b>RL•Conservative RM:</b> The Cannonball Run II is a cross-country car race.

Figure 8: Examples of generated responses from models finetuned with standard SFT and RL with a conservative reward model. False information is highlighted in red.

320 and keep all training details fixed across the two RL methods except for the reward model. All three  
321 models use Llama2 7B as the pretrained model. At test time, we evaluate the models with queries  
322 at different levels of unfamiliarity. The unfamiliarity score for this task is measured by few-shot  
323 prompting the pretrained model (Llama2 7B), sampling 2 responses, and calculating the average  
324 number of incorrect facts in the responses. For more experimental details, see Appendix D.

325 **Results.** To answer our first question, we evaluate the standard and conservative reward models on  
326 held out samples generated from the SFT model. We used samples from the SFT model because the  
327 RL finetuning procedure is initialized with this SFT model, so responses sampled from this model are  
328 representative of the kind of responses that the reward model will be asked to score during RL training.  
329 In Fig. 5, we plot each models’ predicted rewards and the ground truth reward, as inputs become more  
330 unfamiliar. We can see that for unfamiliar inputs, the standard reward model vastly overestimates the  
331 reward, while the conservative reward model does not, showing that the conservative reward models  
332 learned with the procedure we described indeed produce more conservative predictions.

333 To answer our second question, we evaluate standard SFT, as well as RL with a standard reward  
334 model and a conservative reward model on a heldout set of queries for each task. In Fig. 7, we plot the  
335 number of true facts and false facts generated by each model, as inputs become more unfamiliar. We  
336 can see that as inputs became more unfamiliar, the standard SFT model generated fewer truth facts and  
337 more false facts, as expected. Comparing the RL model trained with the conservative reward model  
338 with the standard SFT model, we can see that the RL model generated the same or more true facts  
339 while generating significantly fewer false facts across all levels of input unfamiliarity. Comparing the  
340 two RL models, we can see that while the two generated around the same number of true facts, the  
341 model trained with the conservative reward model generated much fewer false facts across all levels  
342 of input unfamiliarity. We summarize our results in Table 6 with the average percentage of true facts  
343 generated by each method. In Fig. 8, we additionally provide some qualitative examples of responses  
344 generated by the standard SFT model and the RL model trained with conservative reward model. We  
345 can see that as the query became more unfamiliar, responses from the SFT model contained about the  
346 same amount of detail but became more factually incorrect, while responses from the RL model with  
347 conservative supervision defaulted towards less-informative responses. In conclusion, our results  
348 show that RL with conservative reward models outperforms standard SFT and RL with standard  
349 reward models in reducing inaccuracies in model generations.

## 350 6 Conclusion

351 In this work, we presented the observation that, when faced with unfamiliar queries, LLM predictions  
352 tend to default towards the responses associated with unfamiliar examples in its finetuning data.  
353 We additionally studied factuality finetuning for long-form model generations, where we found  
354 that strategically controlling reward model hallucinations can significantly improve the efficacy of  
355 RL-based techniques. Nonetheless, there still remains many open questions and challenges regarding  
356 LLM hallucinations. While our conceptual model explains a model’s behavior for entirely unfamiliar  
357 examples, many real-world queries fall within a spectrum of partial familiarity. A more nuanced  
358 characterization of model predictions in this “middle ground” would be valuable. Furthermore,  
359 our experiments focused on models finetuned for specific applications (e.g., biography generation).  
360 Extending factuality finetuning to more general prompted generation tasks would be useful. We hope  
361 that our work, by offering a deeper understanding of the factors that govern LLM hallucinations,  
362 provides a useful step towards building more trustworthy and reliable LLMs.

## References

- 363
- 364 [1] Ayush Agrawal, Lester Mackey, and Adam Tauman Kalai. Do language models know when  
365 they’re hallucinating references? *arXiv preprint arXiv:2305.18248*, 2023.
- 366 [2] Amos Azaria and Tom Mitchell. The internal state of an LLM knows when its lying. *arXiv*  
367 *preprint arXiv:2304.13734*, 2023.
- 368 [3] Jon Bell. Wikiplots, 2017. URL <https://github.com/markriedl/WikiPlots>.
- 369 [4] Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece  
370 Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general  
371 intelligence: Early experiments with GPT-4. *arXiv preprint arXiv:2303.12712*, 2023.
- 372 [5] Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon Kim, James Glass, and Pengcheng He.  
373 DoLa: Decoding by contrasting layers improves factuality in large language models. *arXiv*  
374 *preprint arXiv:2309.03883*, 2023.
- 375 [6] Luyu Gao, Zhuyun Dai, Panupong Pasupat, Anthony Chen, Arun Tejasvi Chaganty, Yicheng  
376 Fan, Vincent Zhao, Ni Lao, Hongrae Lee, Da-Cheng Juan, et al. RARR: Researching and  
377 revising what language models say, using language models. In *ACL*, 2023.
- 378 [7] Yoav Goldberg. Reinforcement learning for language models, 2023. URL [https://gist.](https://gist.github.com/yoavg/6bff0fec65950898eba1bb321cfbd81)  
379 [github.com/yoavg/6bff0fec65950898eba1bb321cfbd81](https://gist.github.com/yoavg/6bff0fec65950898eba1bb321cfbd81).
- 380 [8] Alexander Havrilla, Maksym Zhuravinskyi, Duy Phung, Aman Tiwari, Jonathan Tow, Stella  
381 Biderman, Quentin Anthony, and Louis Castricato. trIX: A framework for large scale reinforce-  
382 ment learning from human feedback. In *Proceedings of the 2023 Conference on Empirical*  
383 *Methods in Natural Language Processing*, pages 8578–8595, Singapore, December 2023.  
384 Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.530. URL  
385 <https://aclanthology.org/2023.emnlp-main.530>.
- 386 [9] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution  
387 examples in neural networks. *arXiv preprint arXiv:1610.02136*, 2016.
- 388 [10] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and  
389 Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint*  
390 *arXiv:2009.03300*, 2020.
- 391 [11] Liqiang Jing, Ruosen Li, Yunmo Chen, Mengzhao Jia, and Xinya Du. FAITHSCORE: Evaluat-  
392 ing hallucinations in large vision-language models. *arXiv preprint arXiv:2311.01477*, 2023.
- 393 [12] Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. Triviaqa: A large  
394 scale distantly supervised challenge dataset for reading comprehension. *arXiv preprint*  
395 *arXiv:1705.03551*, 2017.
- 396 [13] Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez,  
397 Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, et al. Language  
398 models (mostly) know what they know. *arXiv preprint arXiv:2207.05221*, 2022.
- 399 [14] Adam Tauman Kalai and Santosh S Vempala. Calibrated language models must hallucinate.  
400 *arXiv preprint arXiv:2311.14648*, 2023.
- 401 [15] Nikhil Kandpal, Haikang Deng, Adam Roberts, Eric Wallace, and Colin Raffel. Large language  
402 models struggle to learn long-tail knowledge. In *International Conference on Machine Learning*,  
403 2023.
- 404 [16] Katie Kang, Amrith Setlur, Claire Tomlin, and Sergey Levine. Deep neural networks tend to  
405 extrapolate predictably. *arXiv preprint arXiv:2310.00873*, 2023.
- 406 [17] Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. Semantic uncertainty: Linguistic invariances  
407 for uncertainty estimation in natural language generation. *arXiv preprint arXiv:2302.09664*,  
408 2023.

- 409 [18] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning  
410 for offline reinforcement learning. *Advances in Neural Information Processing Systems*, 33:  
411 1179–1191, 2020.
- 412 [19] Nayeon Lee, Wei Ping, Peng Xu, Mostofa Patwary, Pascale N Fung, Mohammad Shoeybi,  
413 and Bryan Catanzaro. Factuality enhanced language models for open-ended text generation.  
414 *Advances in Neural Information Processing Systems*, 2022.
- 415 [20] Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning:  
416 Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- 417 [21] Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Inference-  
418 time intervention: Eliciting truthful answers from a language model. *arXiv preprint*  
419 *arXiv:2306.03341*, 2023.
- 420 [22] Stephanie Lin, Jacob Hilton, and Owain Evans. Teaching models to express their uncertainty in  
421 words. *arXiv preprint arXiv:2205.14334*, 2022.
- 422 [23] Kevin Liu, Stephen Casper, Dylan Hadfield-Menell, and Jacob Andreas. Cognitive dissonance:  
423 Why do language model outputs disagree with internal representations of truthfulness? *arXiv*  
424 *preprint arXiv:2312.03729*, 2023.
- 425 [24] Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Ha-  
426 jishirzi. When not to trust language models: Investigating effectiveness of parametric and  
427 non-parametric memories. In *ACL*, 2023.
- 428 [25] Potsawee Manakul, Adian Liusie, and Mark JF Gales. Selfcheckgpt: Zero-resource black-box  
429 hallucination detection for generative large language models. *arXiv preprint arXiv:2303.08896*,  
430 2023.
- 431 [26] Mohsen Mesgar, Edwin Simpson, and Iryna Gurevych. Improving factual consistency between  
432 a response and persona facts. *arXiv preprint arXiv:2005.00036*, 2020.
- 433 [27] Sewon Min, Kalpesh Krishna, Xinxu Lyu, Mike Lewis, Wen-tau Yih, Pang Wei Koh, Mohit  
434 Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. FActScore: Fine-grained atomic evaluation  
435 of factual precision in long form text generation. *arXiv preprint arXiv:2305.14251*, 2023.
- 436 [28] Niels Mündler, Jingxuan He, Slobodan Jenko, and Martin Vechev. Self-contradictory hal-  
437 lucinations of large language models: Evaluation, detection and mitigation. *arXiv preprint*  
438 *arXiv:2305.15852*, 2023.
- 439 [29] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin,  
440 Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to  
441 follow instructions with human feedback. *Advances in Neural Information Processing Systems*,  
442 35:27730–27744, 2022.
- 443 [30] Baolin Peng, Michel Galley, Pengcheng He, Hao Cheng, Yujia Xie, Yu Hu, Qiuyuan Huang, Lars  
444 Liden, Zhou Yu, Weizhu Chen, et al. Check your facts and try again: Improving large language  
445 models with external knowledge and automated feedback. *arXiv preprint arXiv:2302.12813*,  
446 2023.
- 447 [31] Paul Roit, Johan Ferret, Lior Shani, Roei Aharoni, Geoffrey Cideron, Robert Dadashi, Matthieu  
448 Geist, Sertan Girgin, Léonard Hussenot, Orgad Keller, et al. Factually consistent summarization  
449 via reinforcement learning with textual entailment feedback. *arXiv preprint arXiv:2306.00186*,  
450 2023.
- 451 [32] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal  
452 policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 453 [33] John Shulman. Reinforcement learning from human feedback: Progress and challenges, 2023.  
454 URL [https://www.youtube.com/watch?v=hhiLw5Q\\_UFg](https://www.youtube.com/watch?v=hhiLw5Q_UFg).
- 455 [34] Kurt Shuster, Spencer Poff, Moya Chen, Douwe Kiela, and Jason Weston. Retrieval augmenta-  
456 tion reduces hallucination in conversation. *arXiv preprint arXiv:2104.07567*, 2021.

- 457 [35] Chenglei Si, Zhe Gan, Zhengyuan Yang, Shuohang Wang, Jianfeng Wang, Jordan Boyd-Graber,  
458 and Lijuan Wang. Prompting GPT-3 to be reliable. *arXiv preprint arXiv:2210.09150*, 2022.
- 459 [36] Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec  
460 Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback.  
461 *Advances in Neural Information Processing Systems*, 33:3008–3021, 2020.
- 462 [37] Marco Antonio Stranisci, Rossana Damiano, Enrico Mensa, Viviana Patti, Daniele Radicioni,  
463 and Tommaso Caselli. Wikibio: a semantic resource for the intersectional analysis of biographi-  
464 cal events. *arXiv preprint arXiv:2306.09505*, 2023.
- 465 [38] Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu, Chunyuan Li, Yikang Shen, Chuang  
466 Gan, Liang-Yan Gui, Yu-Xiong Wang, Yiming Yang, et al. Aligning large multimodal models  
467 with factually augmented RLHF. *arXiv preprint arXiv:2309.14525*, 2023.
- 468 [39] Katherine Tian, Eric Mitchell, Huaxiu Yao, Christopher D Manning, and Chelsea Finn. Fine-  
469 tuning language models for factuality. *arXiv preprint arXiv:2311.08401*, 2023.
- 470 [40] Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao,  
471 Chelsea Finn, and Christopher D Manning. Just ask for calibration: Strategies for eliciting  
472 calibrated confidence scores from language models fine-tuned with human feedback. *arXiv*  
473 *preprint arXiv:2305.14975*, 2023.
- 474 [41] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei,  
475 Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open  
476 foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 477 [42] Logesh Kumar Umapathi, Ankit Pal, and Malaikannan Sankarasubbu. Med-HALT: Medical  
478 domain hallucination test for large language models. *arXiv preprint arXiv:2307.15343*, 2023.
- 479 [43] Neeraj Varshney, Wenlin Yao, Hongming Zhang, Jianshu Chen, and Dong Yu. A stitch in time  
480 saves nine: Detecting and mitigating hallucinations of LLMs by validating low-confidence  
481 generation. *arXiv preprint arXiv:2307.03987*, 2023.
- 482 [44] Weijia Xu, Sweta Agrawal, Eleftheria Briakou, Marianna J Martindale, and Marine Carpuat. Un-  
483 derstanding and detecting hallucinations in neural machine translation via model introspection.  
484 *TACL*, 2023.
- 485 [45] Yuqing Yang, Ethan Chern, Xipeng Qiu, Graham Neubig, and Pengfei Liu. Alignment for  
486 honesty. *arXiv preprint arXiv:2312.07000*, 2023.
- 487 [46] Sina J Semnani Violet Z Yao, Heidi C Zhang, and Monica S Lam. WikiChat: Combating  
488 hallucination of large language models by few-shot grounding on wikipedia. *arXiv preprint*  
489 *arXiv:2305.14292*, 2023.
- 490 [47] Hanning Zhang, Shizhe Diao, Yong Lin, Yi R Fung, Qing Lian, Xingyao Wang, Yangyi Chen,  
491 Heng Ji, and Tong Zhang. R-tuning: Teaching large language models to refuse unknown  
492 questions. *arXiv preprint arXiv:2311.09677*, 2023.
- 493 [48] Yue Zhang, Leyang Cui, Wei Bi, and Shuming Shi. Alleviating hallucinations of large language  
494 models through induced hallucinations. *arXiv preprint arXiv:2312.15710*, 2023.
- 495 [49] Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. Calibrate before use:  
496 Improving few-shot performance of language models. In *International Conference on Machine*  
497 *Learning*, 2021.

498 **NeurIPS Paper Checklist**

499 **1. Claims**

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

502 Answer: [Yes]

503 Justification: We provide empirical evidence (sections 4.2, 5.2) supporting the hypotheses  
504 presented in our work.

505 Guidelines:

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

515 **2. Limitations**

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

517 Answer: [Yes]

518 Justification: We discuss limitations in the conclusion.

519 Guidelines:

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

546 **3. Theory Assumptions and Proofs**

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

549 Answer: [NA]

550 Justification: The paper does not include theoretical results.

551 Guidelines:

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

#### 562 4. Experimental Result Reproducibility

563 Question: Does the paper fully disclose all the information needed to reproduce the main ex-

564 perimental results of the paper to the extent that it affects the main claims and/or conclusions

565 of the paper (regardless of whether the code and data are provided or not)?

566 Answer: [Yes]

567 Justification: We discuss our implementation details in the experiments sections as well as

568 the Appendix.

569 Guidelines:

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

#### 601 5. Open access to data and code

602 Question: Does the paper provide open access to the data and code, with sufficient instruc-

603 tions to faithfully reproduce the main experimental results, as described in supplemental

604 material?

605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656

Answer: [No]

Justification: We used the trlx training package with minor modifications. We will release our training scripts upon acceptance.

Guidelines:

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

## 6. Experimental Setting/Details

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

Answer: [Yes]

Justification: We provide experimental details in the Appendix.

Guidelines:

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

## 7. Experiment Statistical Significance

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

Answer: [No]

Justification: Our experiments require training large language models. Running multiple random seeds would be too computationally expensive.

Guidelines:

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

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

## 667 8. Experiments Compute Resources

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

671 Answer: [Yes]

672 Justification: We provide compute details in the Appendix.

673 Guidelines:

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

## 682 9. Code Of Ethics

683 Question: Does the research conducted in the paper conform, in every respect, with the  
684 NeurIPS Code of Ethics [https://neurips.cc/public/EthicsGuidelines?](https://neurips.cc/public/EthicsGuidelines)

685 Answer: [Yes]

686 Justification: The paper conforms with the Code of Ethics.

687 Guidelines:

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

## 693 10. Broader Impacts

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

696 Answer: [No]

697 Justification: In the conclusion, we discuss how our work can help make LLMs more safe  
698 and reliable, which is a positive societal impact. We do not believe our work has any direct  
699 negative societal impacts.

700 Guidelines:

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



- 708
- 709
- 710
- 711
- 712
- 713
- 714
- 715
- 716
- 717
- 718
- 719
- 720
- 721
- 722
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
  - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
  - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

723 **11. Safeguards**

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

727 Answer: [NA]

728 Justification: The paper poses no such risks.

729 Guidelines:

- 730
- 731
- 732
- 733
- 734
- 735
- 736
- 737
- 738
- 739
- The answer NA means that the paper poses no such risks.
  - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
  - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
  - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

740 **12. Licenses for existing assets**

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

744 Answer: [No]

745 Justification: We currently cite the original owners of the data and models that we use. We  
746 will add more details about the license, copyright information, and terms of use of these  
747 assets upon acceptance.

748 Guidelines:

- 749
- 750
- 751
- 752
- 753
- 754
- 755
- 756
- 757
- 758
- 759
- The answer NA means that the paper does not use existing assets.
  - The authors should cite the original paper that produced the code package or dataset.
  - The authors should state which version of the asset is used and, if possible, include a URL.
  - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
  - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
  - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.

- 760
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
  - If this information is not available online, the authors are encouraged to reach out to the asset’s creators.
- 761
- 762
- 763

764 **13. New Assets**

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

767 Answer: [No]

768 Justification: We will release our code upon acceptance.

769 Guidelines:

- The answer NA means that the paper does not release new assets.
  - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
  - The paper should discuss whether and how consent was obtained from people whose asset is used.
  - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.
- 770
- 771
- 772
- 773
- 774
- 775
- 776
- 777

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

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

782 Answer: [NA]

783 Justification: The paper does not involve crowdsourcing nor research with human subjects.

784 Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
  - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
  - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.
- 785
- 786
- 787
- 788
- 789
- 790
- 791
- 792

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

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

799 Answer: [NA]

800 Justification: The paper does not involve crowdsourcing nor research with human subjects.

801 Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
  - Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
  - We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
  - For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.
- 802
- 803
- 804
- 805
- 806
- 807
- 808
- 809
- 810
- 811

## 812 **A Compute**

813 We use A100 GPUs to finetune our models. Number of GPUs used range from 1-6 for each  
814 experiment, and time of execution range from a few hours to up to 2 days. We use LoRA finetuning  
815 for all our experiments with  $r = 16$ ,  $\alpha = 16$ ,  $\text{dropout} = 0$ .

## 816 **B MMLU Training Details**

817 In this section, we provide more details on our training and evaluation procedure for our MMLU  
818 experiments. For all experiments, we finetuned on the evaluation split of MMLU, and evaluated on  
819 the validation split. This is because MMLU does not have a training split. Our training pipeline uses  
820 the trlx codebase [8].

### 821 **B.1 SFT Models**

822 We classify examples with unfamiliarity score (NLL) greater than 0.36 as unfamiliar, and the rest  
823 as familiar. During finetuning, we rebalance the dataset such that 50% of finetuning examples are  
824 familiar and 50% are unfamiliar.

825 We use a batch size of 12. We use the AdamW optimizer with learning rate =  $1e-5$ ,  $\text{betas} = (0.9, 0.95)$ ,  
826  $\text{eps} = 1.0e-8$ , and weight decay= $1.0e-6$ .

### 827 **B.2 RL Models**

828 We initialize all RL finetuning with a model that has already be supervised finetuned to produce  
829 responses that consist of answer choices. The SFT model we used for initialization is trained predict  
830 the E option 50% of the time, and to produce the correct answer to the query 50% of the time.

831 We use a batch size of 12. We use the AdamW optimizer with learning rate =  $1e-5$ ,  $\text{betas} = (0.9, 0.95)$ ,  
832  $\text{eps} = 1.0e-8$ , and weight decay= $1.0e-6$ . For PPO, we use  $\text{cliprange} = 0.005$  and  $\text{KL coef} = 0$ .

### 833 **B.3 Reward Models**

834 We construct correct (reward 1) training and evaluation examples using queries and their correspond-  
835 ing answer labels from the original MMLU dataset. We construct incorrect (reward 0) examples by  
836 using queries from the original dataset, and randomly sampling incorrect answer labels (A-D not  
837 including correct label).

838 We use a batch size of 12. We use the AdamW optimizer with learning rate =  $1e-5$ ,  $\text{betas} = (0.9, 0.95)$ ,  
839  $\text{eps} = 1.0e-8$ , and weight decay= $1.0e-6$ .

## 840 **C TriviaQA Training Details**

841 In this section, we provide more details on our training and evaluation procedure for our TriviaQA  
842 experiments. Our training pipeline uses the trlx codebase [8].

### 843 **C.1 SFT Models**

844 We classify examples with unfamiliarity score (number of incorrect responses out of 12 samples)  
845 greater than 6 as unfamiliar, and familiar otherwise. We relabel the responses associated with all  
846 unfamiliar finetuning examples to be “I don’t know”.

847 We use a batch size of 32. We use the AdamW optimizer with learning rate =  $1e-5$ ,  $\text{betas} = (0.9, 0.95)$ ,  
848  $\text{eps} = 1.0e-8$ , and weight decay= $1.0e-6$ . We use a Cosine Annealing scheduler with  $T \text{ max} = 1e4$  and  
849  $\text{ETA min} = 1e-10$ .

### 850 **C.2 RL Models**

851 We initialize all RL finetuning with a model that has already be supervised finetuned to produce  
852 responses that consists of an answer or “I don’t know”. The SFT model we used for initialization is

853 trained predict “I don’t know” 40% of the time, and to produce the correct answer to the query 60%  
854 of the time.

855 We use a batch size of 32. We use the AdamW optimizer with learning rate =  $1e-5$ , betas = (0.9, 0.95),  
856 eps =  $1.0e-8$ , and weight decay= $1.0e-6$ . For PPO, we use cliprange = 0.005 and KL coef = 0.1.

### 857 C.3 Reward Models

858 We construct correct (reward 1) training and evaluation examples using queries and responses from  
859 the original TriviaQA dataset. We construct incorrect (reward 0) examples using queries from the  
860 original dataset, and responses generated from few-shot prompting Llama2 7B or GPT-2. We filter  
861 the generated responses to ensure that all responses were incorrect.

862 We use a batch size of 32. We use the AdamW optimizer with learning rate =  $1e-5$ , betas = (0.9, 0.95),  
863 eps =  $1.0e-8$ , and weight decay= $1.0e-6$ .

## 864 D Long-form Tasks Training Details

865 In this section, we provide training and evaluation details for our long-form factuality finetuning  
866 experiments. Our training pipeline uses the trlx codebase [8].

### 867 D.1 Data

868 We construct finetuning and evaluation datasets using WikiBios and WikiPlots, both of which consist  
869 of wikipedia entries attached to people and books/movies. We make use of the first sentence in the  
870 wikipedia entry for both tasks as the target response in our SFT finetuning datasets. The prompts we  
871 use for finetuning are “Write a biography for [name].” and “What is the premise of [title]?”. For the  
872 biography task, our finetuning dataset includes 104539 examples, and our evaluation dataset includes  
873 5000 examples. For the plot generation task, our finetuning dataset includes 10000 examples, and our  
874 evaluation dataset includes 4795 examples.

### 875 D.2 Reward Models

876 We take a two-staged approach to learning a reward model. First, we trained a model to break  
877 down a response into individual atomic facts. Next, we trained a separate model to predict the  
878 factuality of each atomic fact. We then use the predicted factuality of each fact to calculate the overall  
879 reward associated with each response. The supervision for both models are collected by querying  
880 FActScore, which is a automated pipeline that queries GPT-3.5 to decompose a response into atomic  
881 facts and produces the factuality of each atomic fact. We use 10000 labeled examples to train the  
882 conservative reward model and the standard reward models each for both tasks. Note that while we  
883 use a two-staged strategy for learning reward models in our implementation, our general approach for  
884 learning conservative reward model should apply to other reward model learning strategies as well,  
885 such as directly predicting the reward associated with a response.

886 For both models, we use a batch size of 32. We use the AdamW optimizer with learning rate =  $2e-5$ ,  
887 betas = (0.9, 0.95), eps =  $1.0e-8$ , and weight decay= $1.0e-6$ . We use a Cosine Annealing scheduler  
888 with T max =  $1e4$  and ETA min =  $1e-10$ .

### 889 D.3 SFT Models

890 We use a batch size of 24. We use the AdamW optimizer with learning rate =  $1e-5$ , betas = (0.9, 0.95),  
891 eps =  $1.0e-8$ , and weight decay= $1.0e-6$ . We use a Cosine Annealing scheduler with T max =  $1e4$  and  
892 ETA min =  $1e-10$ .

### 893 D.4 RL Models

894 We initialize all RL finetuning with the SFT model, and use the reward predicted by the reward model  
895 described above as supervision.

896 We use a batch size of 16. We use the AdamW optimizer with learning rate =  $1e-5$ , betas = (0.9, 0.95),  
897 eps =  $1.0e-8$ , and weight decay= $1.0e-6$ . For PPO, we use cliprange = 0.005 and KL coef = 0.5.