FLAME⁶: Factuality-Aware Alignment for Large Language Models

Sheng-Chieh Lin^{1*}, Luyu Gao², Barlas Oguz³, Wenhan Xiong³, Jimmy Lin¹, Wen-tau Yih³, Xilun Chen^{3†}

University of Waterloo¹, Carnegie Mellon University², Meta AI³ s269lin@uwaterloo.ca, xilun@meta.com

Abstract

Alignment is a procedure to fine-tune pre-trained large language models (LLMs) to follow natural language instructions and serve as helpful AI assistants. We have observed, however, that the conventional alignment process fails to enhance the factual accuracy of LLMs, and often leads to the generation of more false facts (i.e., hallucination). In this paper, we study how to make the LLM alignment process more factual, by first identifying factors that lead to hallucination in both alignment steps: supervised fine-tuning (SFT) and reinforcement learning (RL). In particular, we find that training the LLM on new or unfamiliar knowledge can encourage hallucination. This makes SFT less factual as it trains on humanlabeled data that may be novel to the LLM. Furthermore, reward functions used in standard RL often inadequately capture factuality and favor longer and more detailed responses, which inadvertently promote hallucination. Based on these observations, we propose FactuaLity-aware AlignMEnt (FLAME), comprised of factuality-aware SFT and factuality-aware RL through direct preference optimization. Experiments show that our proposed FLAME guides LLMs to output more factual responses while maintaining their instruction-following capability.

1 Introduction

Alignment [Ouyang et al., 2022] is a procedure to make pre-trained large language models (LLMs) [Brown et al., 2020, Touvron et al., 2023] follow human instructions and serve as helpful AI assistants. Despite significant progress in general LLM alignment [Ouyang et al., 2022, Bai et al., 2022, Yuan et al., 2024], state-of-the-art aligned LLMs are still prone to generate false claims [OpenAI, 2023, Min et al., 2023]. In this work, we therefore attempt to advance the understanding of the underlying causes of LLM hallucination as well as its relation to the alignment procedure.

We consider the commonly seen alignment process consisting of two training phases: (1) supervised fine-tuning (SFT) [Sanh et al., 2022]; (2) reinforcement learning (RL) with human [RLHF, Ouyang et al., 2022, Bai et al., 2022] or automated feedback [RLAIF, Bai et al., 2023]. In our study, we find that both the SFT and RL steps in the standard alignment process may actually *encourage* LLMs to hallucinate. First, in the SFT stage, LLMs are fine-tuned with diverse instructions paired with human-created high-quality responses. While this leads to strong instruction-following capability [Ouyang et al., 2022, Köpf et al., 2023, Zhou et al., 2023, Touvron et al., 2023], our study shows that such human-labeled responses may present *new or unknown information* to the LLM. This, in turn, may inadvertently promote hallucination. Second, we find that the standard reward used in the RL stage

38th Conference on Neural Information Processing Systems (NeurIPS 2024).

^{*} This work is done during Sheng-Chieh's internship at Meta.

[†] Xilun and Sheng-Chieh contributed equally to this work.

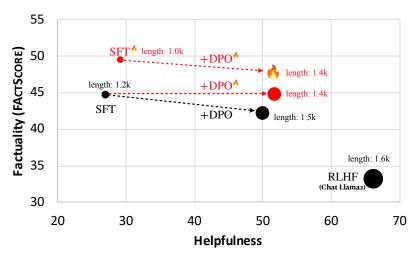


Figure 1: Models' helpfulness on Alpaca Eval vs factuality on biography. Helpfulness is measured by models' win rate over our baseline SFT + DPO on Alpaca Eval. Dot size represents average length of bio generation.

often prefers longer and more detailed responses [Singhal et al., 2023, Chen et al., 2024b, Yuan et al., 2024]. Consequently, a reward-hacking model ends up with a tendency to produce longer claims with more non-factual information, as shown in the black dots in Figure 1. One possible reason is that most existing RLHF or RLAIF approaches rely on a single scalar reward to represent preference, which struggles to cover multiple alignment skill sets [Ye et al., 2024] and is likely to under-present the aspect of factuality [Hosking et al., 2024].

To address the aforementioned issues, we study the key factors which impact factuality during alignment. In particular, we first conduct a pilot study on the biography generation task [Min et al., 2023] in a more controlled setting where the alignment process focuses solely on factuality (Section 3). Our key observation is that an LLM hallucinates more if it is fine-tuned on new knowledge in either the SFT or the RL stage. For example, an LLM becomes significantly less factual when fine-tuned on responses produced by a model with access to external knowledge (e.g. a retrieval augmented LLM), even though those responses are more factual themselves. Similarly, hallucination is greatly increased if RLAIF is performed on preference pairs that consist of retrieval-augmented LLM output as positive examples and the LLM's own output as negative examples. In comparison, we discover that fine-tuning a pre-trained LLM on a subset of its *own* generations selected by factuality yields more factual responses and reduces hallucinations.

Next, we apply our findings to improve the factuality of the general LLM alignment process, which is more challenging due to the diversity of instructions. As shown in Figure 2, we observe that some instructions require factual responses while the others do not, and therefore would require different alignment treatments. We first identify fact-based instructions that require factual responses and leverage the findings in our pilot study to create additional training data at both SFT and RL stages to explicitly guide LLMs to output factual responses. Specifically, at the SFT stage, for fact-based instructions, instead of using human created seed training data, we elicit knowledge from the pre-trained LLM and construct training data using its own pre-trained knowledge. This can prevent fine-tuning the LLM on knowledge unknown to itself. At the RL stage, we create additional preference pairs focused on factuality for fact-based instructions, which are combined with the standard preference pairs for instruction following during Direct Preference Optimization [DPO; Rafailov et al., 2023].

We evaluate models on Alpaca Eval [Dubois et al., 2024] and Biography, using win rate for instruction-following capability and FACTSCORE [Min et al., 2023] for factuality evaluation. As shown in Figure 1, using our FLAME method (SFT + DPO), a significantly higher FACTSCORE (+5.6 pts) is achieved compared to the standard alignment process (SFT + DPO), without sacrificing the LLM's instruction-following capability (51.2% win rate). Our ablation study also indicates that identifying fact-based instructions is the key to factual alignment in the general alignment setting.

2 Related Work

Alignment. Since pre-trained LLMs cannot accurately follow human instructions, a bunch of work has been proposed to improve LLM alignment through SFT and RL. Some propose to improve SFT through data curation [Zhou et al., 2023, Chen et al., 2024a], diverse instruction augmentation [Wang et al., 2023a, Li et al., 2024] while others focus on RL with human feedback [Ouyang et al., 2022, Bai et al., 2022], AI feedback [Bai et al., 2023, Sun et al., 2024, Yuan et al., 2024]. The main goal of these alignment approaches is instruction-following capability (or helpfulness), which may guide LLMs to output detailed and lengthy responses [Singhal et al., 2023] but inevitably encourage hallucination.

Factuality. Prior work has highlighted the issue of hallucination in LLMs [Gao et al., 2022, Kandpal et al., 2023, Mallen et al., 2023]. To address the issue, important research lines are factuality evaluation [Min et al., 2023, Wang et al., 2023b, Chern et al., 2023] and improvement. Some training-free approaches to improve LLMs' factuality include external knowledge augmentation [Gao et al., 2022, Kandpal et al., 2023, Cheng et al., 2023, Jiang et al., 2023] and specialized decoding [Li et al., 2023, Chuang et al., 2024].

Recent studies apply RL to improve LLMs' factuality. For example, Tian et al. [2024] propose to construct factuality preference pairs for direct preference optimization [DPO; Rafailov et al., 2023], which is closely related to our work. However, they focus solely on enhancing LLMs' factuality through DPO but overlook its potential impact on the models' instruction-following capability, as demonstrated in our experiments. In contrast, our work provides a comprehensive examination of improving LLMs' factuality and instruction-following ability through fine-tuning approaches encompassing both SFT and DPO. Concurrent to our work, Kang et al. [2024] find that LLMs tend to hallucinate when facing unfamiliar queries. They consider improving LLMs' factuality as teaching LLMs to output abstaining or less detailed responses on such unfamiliar queries, a similar behavior observed from our LLMs fine-tuned with FLAME (see case studies in Section 6.5). It is worth mentioning that both prior studies focus on a simplified scenario as our pilot study in Section 3: fine-tuning LLMs to improve factuality on a single task (e.g., fine-tuning and evaluating on biography generation). In contrast, we consider the general alignment task, where LLMs are given diverse and complex instructions.

3 A Pilot Study on Factual Alignment

In this section, we first study how to align large language models (LLMs) to be more factual. We use biography generation as the task of our pilot study for two main reasons: (1) Biography generation is a simplified setting where factuality is the sole focus of the alignment process. As we will discuss in Section 4, studying factual alignment on diverse human instructions is more complex, as the alignment process encompasses aspects beyond factuality, such as helpfulness and safety. (2) Evaluating the factuality of biography generation is relatively easy since Wikipedia covers sufficient information for public figures and most of the facts about a person are non-debatable [Min et al., 2023].

3.1 Alignment for Biography Generation

A standard alignment procedure consists of supervised fine-tuning (SFT) and reinforcement learning (RL). In this pilot study, our main goal is to teach LLMs to generate biography with reduced misinformation. For the experiment, we compile training and evaluation datasets comprising 500 and 183 diverse human entities, respectively (further details provided in Appendix A.1). We employ FACTSCORE [FS; Min et al., 2023] as the automated metric for assessing factuality, given its fine-grained evaluation capabilities for long-form text generation and its strong correlation with human judgments.³ To study factuality alignment in this pilot study, we posit that training data is needed where the responses are more factual than the LLM's own generations. Thus, we use retrieval-augmented LLMs [RAG; Lewis et al., 2020] to generate training data, which has been shown to output more factual responses [Mialon et al., 2023].

Throughout the paper, we refer to the pre-trained (PT), supervised fine-tuned (SFT), and direct preference optimization (DPO) fine-tuned LLMs as PT, SFT, and DPO, respectively.⁴

³We use the evaluator: retrieval+llama+npm

⁴Note that in our experiments, we use DPO as the substitute of RL [Schulman et al., 2017].

	Fact-based $(x \in X^{fact})$		Non fact-based $(x \notin X^{fact})$
(1)	Do you have any information about the Commodore 64?	(6)	How would a child feel if it fell down on the ground hitting its face?
	Hi, could you help me to solve this cubic equation using Cardano's		Write a fun story that can be told in 3 minutes at the dinner table. We
(2)	Mehod (step by step if possible), please? -> " $x^3 + 2x^2 - x - 1 = 0$		are 3 developers enjoying a pizza. The story must contain these word:
	"		zombie, ethernet cable, sudo, dashboard.
(3)	Please give me a brief history of coffee.	(8)	Tell me a story about a pig who goes to the moon.
(1	What are the principles at play in UHPLC-MS analysis?	(0)	Is the internet's focus on engagement the root of most of its problems
(4,		(3)	Is the internet's focus on engagement the root of most of its problems and shortcomings?
	Explain the significance of the American Revolution, including the		Can you tell me a bit about what has gone into your creation?
(5)	events that led up to it, the impact it had on the world, and its	(10)	
	ongoing relevance today.		

Figure 2: Instructions from Open Assistant dataset. The instructions are classified with SFT model using the prompt in Appendix Figure 4.

SFT. We explore two sources of supervision to generate training data (detailed in Appendix A.1): (1) using PT^{RAG} with few-shot demonstration to generate biographies for each name entity in training data, where PT^{RAG} is PT augmented with an off-the-shelf retriever [Lin et al., 2023]; (2) using vanilla PT with few-shot demonstration to generate training data as a baseline. As shown in Table 1, PT^{RAG} is indeed much more factual than PT. However, a surprising discovery in the pilot study is that *fine-tuning on such more factual instruction-biography pairs generated by* PT^{RAG} *results in a less factual* SFT *model* (row 4 vs 3).

DPO. We further fine-tune the LLMs to be more factual through DPO. An intuitive way to create factuality preference pairs is to directly use the samples from PT^{RAG} and PT as positives and negatives since PT^{RAG} generates more factual biographies than PT (row 2 vs 1). Another approach is to employ FACTSCORE (FS) as the reward to select positive and negative samples among the generations from PT itself [Tian et al., 2024] (detailed in Appendix A.1). As shown in Table 1, DPO fine-tuned on self-generated data with FS reward guides models to generate more factual responses (row 5 vs 3); however, DPO fine-tuned with the supervision of PT^{RAG} makes

Table 1: Pilot study on bio generation. Pos. denotes the positives for SFT or DPO. Neg. denotes the negatives for DPO. FS denotes FACTSCORE.

Llama-2 7B	src. of su	pervision	Bio		
	Pos.	Neg.	FS	# Corr. / Err.	
(1) PT	-	-	39.1	14.4 / 22.0	
(2) PT ^{RAG}	-	-	55.4	18.6 / 15.9	
(3) SFT	PT	-	37.9	13.4 / 21.8	
(4) SF 1	PT^{RAG}	-	35.7	13.5 / 23.7	
(5) DPO	PT*	PT^*	41.6	15.4 / 20.7	
(6) DF O	PT^{RAG}	PT	23.5	12.7 / 34.9	

^{*} FACTSCORE is used to select positives and negatives.

the models hallucinate even more than its SFT counterpart (6 vs 4).

This outcome suggests that compelling models to generate responses akin to PT^{RAG} prompts increases hallucination. Conversely, fine-tuning LLMs on their own generations appears to be crucial for factual alignment, a finding applicable to both SFT and DPO fine-tuning.

3.2 Strategies for Factual Alignment

From the pilot study, we find that better quality data (in terms of factuality) for SFT and DPO does not necessarily yield models with better factual alignment. This is likely because the supervision from RAG contains information unknown to the LLM; thus, fine-tuning on RAG generated responses may inadvertently encourage the LLM to output unfamiliar information. To avoid unknown knowledge from being presented to the LLM, a viable strategy is to create SFT and DPO training data using the generated responses from the LLM itself.

4 Factuality-Aware Alignment

In the section, we further extend our discussion of factual alignment to encompass more general instructions. Unlike biography generation in Section 3, where factuality is the main alignment objective, human instructions are diverse and complex, necessitating a range of alignment skill sets beyond factuality alone; e.g., logical thinking, problem handling and user alignment [Ye et al., 2024]. Thus, conducting factual alignment with the diverse instructions face two main challenges: (1) different instructions may demand distinct skill sets. For example, in Figure 2, instruction 3, "Please give me a brief history of coffee", necessitates factual accuracy and concise summarization,

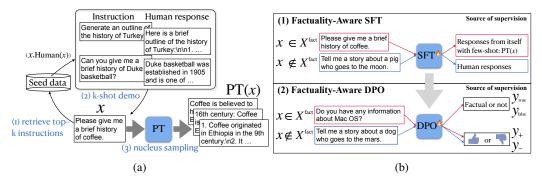


Figure 3: Illustrations of (a) response generation using a pre-trained LLM (PT) with few-shot demonstration; (b) factuality-aware alignment.

while instruction 8, "Tell me a story about a pig who goes to the moon", prioritizes creativity and imagination over strict factuality. (2) As recent studies have emphasized [Ye et al., 2024, Hosking et al., 2024], using a single scalar for reward modeling fails to adequately address multiple alignment skill sets and often under-presents the aspect of factuality.

To tackle the aforementioned challenges, we propose *factuality-aware alignment* (FLAME*). To address the first challenge, we propose to prompt LLMs to classify whether a given instruction demands the response to be factual, as shown in Figure 2. We then apply the factuality fine-tuning strategy for SFT and DPO discussed in Section 3.2 to those fact-based instructions. Furthermore, to address the second challenge, we employ separate rewards to evaluate the factuality and instruction-following capability of an LLM. For simplicity, our work only considers two alignment skill sets: instruction following and factuality. We leave more comprehensive reward modeling to future work.

In the following, we first describe our baseline alignment approach and introduce our proposed factuality-aware alignment built on top of the baseline alignment procedure.

4.1 Baseline Alignment

We initialize PT from Llama-2 70B pre-trained model⁵ and build our baseline alignment procedure following self-rewarding language models [Yuan et al., 2024] due to its simplicity and independence of other strong LLMs (e.g., GPT4) or human evaluators as a reward model. The alignment comprises two steps: (1) building SFT model fine-tuned on a high-quality seed data consisting of 3,200 instructions and each instruction is paired with the best response created by humans from Open Assistant dataset [OASST; Köpf et al., 2023]; (2) further fine-tuning SFT through DPO on instruction-following preference data (x, y_+, y_-) constructed by itself (SFT) as the reward model, RM^{IF}, where y_+ and y_- are the positive and negative responses for a given prompt x, respectively. The resulting fine-tuned model is denoted as SFT + DPO. Note that, following Yuan et al. [2024], we use additional augmented 20K instructions to create the preference training data for DPO fine-tuning. Further details are provided in Appendix A.3.

4.2 Our Approach

4.2.1 Factuality-Aware SFT (SFT)

Although leveraging human created high-quality seed data is a reasonable choice for SFT [Zhou et al., 2023], our study in Section 3 suggests that fine-tuning on such high-quality data generated by models other than the LLM itself may present unknown information to the LLM, which may in turn encourage hallucination. To address the above issue, for each instruction from the seed data, we elicit the knowledge from the pre-trained LM itself by generating the responses with a few-shot demonstration. Furthermore, to better use the knowledge from both humans and the pre-trained LLM itself, we propose to utilize human generated responses for non-fact-based instructions, while leveraging the responses sampled from pre-trained LLMs for fact-based instructions to mitigate the introduction of unknown knowledge.

⁵https://huggingface.co/meta-llama/Llama-2-70b

Specifically, we create factuality-aware alignment training data for SFT with two steps. (1) Classifying instructions: we first prompt SFT to judge whether an instruction from the seed data is fact-based ($x \in X^{\text{fact}}$) or not.⁶ (2) Eliciting knowledge from PT: as illustrated in Figure 3(a), we sample 10 responses from PT with 5-shot demonstration, $(x_0, \text{Human}(x_0)) \cdots (x_4, \text{Human}(x_4))$, where x_k is the top-k similar instruction to x retrieved by DRAGON+ [Lin et al., 2023] from the seed data. Human(x_k) denotes the corresponding human response to x_k in the seed data. As illustrated in Figure 3(b) (upper), the resulting training data for SFT is $(x \notin X^{\text{fact}}, \text{Human}(x)), (x \in X^{\text{fact}}, \text{PT}(x))$, where PT(x) denotes the set of responses to x sampled from PT. The fine-tuned model is denoted as SFT.

4.2.2 Factuality-Aware DPO (DPO)

At the second stage of alignment with DPO, we use SFT to generate multiple responses y_0, y_1, \cdots for a given instruction x; then, using SFT itself as the reward model (RM^{IF}) to create a preference pair: (x, y_+, y_-) . The above data creation procedure is the same as the second stage of our baseline alignment in Section 4.1. However, recent studies [Saha et al., 2024, Hosking et al., 2024, Ye et al., 2024] indicate that a single scalar reward from human feedback or LLM reward models may underrepresent the aspect of factuality. To address this limitation, we introduce another factuality reward model (RM^{fact}) to evaluate factuality of responses and create a factuality preference pair for fact-based instructions: $(x \in X^{\text{fact}}, y_{\text{true}}, y_{\text{false}})$.

Specifically, we build RM^{fact} with retrieval augmentation to measure the percentage of facts in a response that are correct. RM^{fact} comprises two main components: atomic fact decomposition and retrieval augmented claim verification. We detail the components and ablate their impacts on the quality of RM^{fact} in Appendix A.5. We compute factuality reward for the same responses sampled from SFT^{\bullet} : $RM^{fact}(x,y_0)$, $RM^{fact}(x,y_1)$, The response with the highest (lowest) factuality reward is chosen as y_{true} (y_{false}). Note that if the chosen paired responses show large difference in instruction-following reward, we discard the pair; i.e., $|RM^{IF}(x,y_{true}) - RM^{IF}(x,y_{false})| > 0.5$. As illustrated in Figure 3(b) (lower), in factuality-aware DPO training, the model is initialized from SFT^{\bullet} and the fine-tuned model is our final factuality-aware aligned model, denoted $SFT^{\bullet} + DPO^{\bullet}$. The specific procedures for fine-tuning models in both the SFT and DPO are described in Appendix A.6.

5 Experiments

5.1 Evaluation Datasets and Metrics

Instruction Following. We use the 805 instruction-following tasks from Alpaca Eval [Dubois et al., 2024] to evaluate models head-to-head win rate against our baselines using the recommended evaluator: alpaca_eval_gpt4_turbo_fn. We use SFT and SFT + DPO described in Section 4.1 as the baselines for win rate comparisons.

Factuality. We evaluate models on three datasets with diverse knowledge-intensive instructions for factuality. (1) Biography: a knowledge insensitive sub-task of instruction-following tasks. Following our pilot study in Section 3, we use the 183 human entities provided by Min et al. [2023] with the prompt "Tell me a bio of entity name". (2) Alpaca Fact: we extract the fact-based instructions from the 803 instructions using our SFT model (with the prompt shown in Appendix Figure 4), resulting in 241 instructions. (3) FAVA [Mishra et al., 2024]⁸: the 141 knowledge-intensive instructions from multiple sources, including Open Assistant [Köpf et al., 2023], No Robots [Rajani et al., 2023], WebNLG [Gardent et al., 2017] and manually created datasets. We report FACTSCORE (FS) without length penalty as the metric for all the three datasets. Note that original FS computes proportion of correct facts with additional penalty on short generations with less than 10 atomic facts. This penalty aims to address situations where models provide insufficiently detailed answers. We assume that this aspect is considered in the evaluation of instruction following in Alpaca Eval. In addition, we also

⁶Prompt for fact-based instruction classification is shown in Appendix Figure 4.

⁷We sample 4 responses for each augmented instruction.

[%]https://huggingface.co/datasets/fava-uw/fava-data/blob/main/ annotations.json

Table 2: Experimental results of supervised fine-tuning on Open Assistant dataset. PT denotes pre-trained Llama2 70B with 5-shot demonstration. SFT^{fact} denotes the variant which only optimizes factuality. FS denotes FACTSCORE.

Llama-2 70B	src. of supervision		Alpaca Eval	l Bio		A	lpaca Fact	FAVA	
	Human	PT	win rate over (1)	FS	# Corr. / Err.	FS	# Corr. / Err.	FS	# Corr. / Err.
(0) PT	-	-	-	53.1	15.3 / 13.5	-	-	-	-
(1) SFT	√	Х	50.0	44.7	21.1 / 26.8	38.6	16.7 / 29.0	54.4	21.2 / 25.8
(2) SFT ^{fact}	×	\checkmark	48.1	48.5	19.6 / 20.6	42.0	17.5 / 28.4	53.3	18.3 / 24.2
(3) SFT ⁶	✓*	√ *	51.2	49.5	19.9 / 19.5	41.4	18.3 / 27.7	54.2	19.3 / 22.4

^{*} SFT ouses supervision from Human and PT for non-fact-based and fact-based instructions, respectively.

Table 3: Experiments of direct preference optimization (DPO). IF. and Fact. denote instruction following (x, y_+, y_-) and factuality $(x \in X^{\text{fact}}, y_{\text{true}}, y_{\text{false}})$ preference data, where X^{fact} denotes the set of fact-based instructions. DPO^{fact} denotes the variant which only optimizes factuality. The preference data statistics is listed in Appendix, Table 11.

Llama-2 70B	src. of supervision		Alpaca Eval	Eval Bio		A	lpaca Fact	FAVA	
	IF.	Fact.	win rate over (2)	FS	# Corr. / Err.	FS	# Corr. / Err.	FS	# Corr. / Err.
(0) Chat	Prop	rietary data	66.2	33.2	23.4 / 43.6	39.3	22.3 / 36.4	47.5	28.0 / 31.3
(1) SFT	-	-	27.1	44.7	21.1 / 26.8	38.6	16.7 / 29.0	54.4	21.2 / 25.8
(2) + DPO	√	Х	50.0	42.3	24.6 / 35.0	41.6	22.9 / 34.6	52.9	28.1 / 26.8
$(3) + DPO^{fact}$	×	\checkmark	40.8	47.1	19.8 / 23.9	48.2	17.5 / 19.0	57.9	20.0 / 15.9
(4) + DPO ⁶	✓	\checkmark	51.7	44.9	23.7 / 30.3	45.0	23.1 / 28.7	56.4	27.1 / 23.3
(5) SFT*	-	-	29.1	49.5	19.9 / 19.5	41.4	18.3 / 27.7	54.2	19.3 / 22.4
(6) + DPO	✓	Х	50.4	46.3	24.0 / 28.7	43.9	21.6 / 28.8	55.0	25.4 / 22.0
(7) + DPO	✓	✓	51.2	47.9	25.9 / 28.5	48.7	24.1 / 25.5	58.9	29.0 / 22.2

report the number of correct and erroneous facts. All the numbers reported are averaged over the instructions in each dataset.

In addition, we also evaluate our fine-tuned models' truthfulness using TruthfulQA [Lin et al., 2022]. We evaluate model performance in the generation task and use ROUGE [Lin, 2004] and BLEU [Papineni et al., 2002] to measure the quality of responses.

5.2 Comparisons of SFT

Table 2 compares the pre-trained Llama-2 70B fine-tuned on OASST dataset with responses from different sources. We list the FACTSCORE (FS) of biography generation using the pre-trained model through Bio 5-shot demonstration as reference (row 0) and SFT, which is fine-tuned on our seed data with human-created responses, is our baseline (row 1). We first notice that SFT shows significant FACTSCORE degradation (53.1 vs 44.7) compared to Bio 5-shot with the pre-trained model. It seems that SFT tends to generate more lengthy responses but with more erroneous facts.

When eliciting the knowledge from PT by fine-tuning on its own generated responses, SFT^{fact} generates more factual responses in Biography and Alpaca (row 2 vs 1). However, it shows slightly inferior instruction-following capability in Alpaca Eval. This result demonstrates that human responses indeed teach LLMs how to better follow instructions but also encourage LLMs to output more false facts. On the other hand, eliciting the knowledge from the pre-trained model itself avoids the encouragement of hallucination albeit with a slight reduction in instruction-following capability. Finally, SFT combining supervision from humans and PT, shows comparable instruction-following capability and output more factual responses on fact-based instructions (row 3 vs 1).

5.3 Comparisons of DPO

Table 3 compares different DPO training recipes. First, we conduct DPO fine-tuning on our SFT baseline, SFT. When further aligning the model to follow instructions, DPO sees a significant improvement in instruction-following capability (row 2 vs 1) with win rate 72.9 over SFT; however, the instruction aligned model tends to output lengthy responses with more factual errors (see examples in Appendix Figure 10). On the other hand, when only aligned with factual preference data, DPO^{fact} shows less improvement in instruction-following capability (row 1 vs 3). These results indicate that

preference optimization for either instruction following or factuality alone may come at the expense of the other since the former encourages models to output long and detailed responses while the later discourages models to output false claims. When jointly conducting instruction and factuality alignment, DPO not only better follows instructions but also outputs more factual responses (row 4 vs 1, 2). Finally, initializing from SFT, the DPO fine-tuned models are more factual than their counterparts (i.e., 6 vs 2 and 7 vs 4) without instruction-following capability degrade. We also list the results from Llama-2-Chat 70B (row 0) and observe that despite of its strong instruction-following capability, it tends to output many more incorrect facts. These results demonstrate that standard alignment, even on proprietary commercial data, may encourage LLMs to hallucinate. In contrast, our factuality-aware alignment guides LLMs to output more factual responses without degradation in their general instruction-following capabilities.

It is worth noting that SFT^{fact} and DPO^{fact} are similar to SFT and DPO fine-tuning proposed by Tian et al. [2024], which improve LLMs' factuality but degrade their instruction-following capability. Also, we do not observe our SFT and DPO variants outperform the pre-trained model with few-shot demonstrations on biography generation (row 0 in Table 2. This is possibly due to the alignment tax found in previous work [Ouyang et al., 2022], which degrades LLMs' accuracy on the standard knowledge benchmarks. How to improve both models' instruction-following capability and their accuracy on standard knowledge benchmarks is worth exploring, which we leave for future work.

5.4 Results on TruthfulQA

Table 4 compares models performance on TruthfulQA. Generally, we observe that our factuality-aware alignment training guides LLMs to output more truthful responses. For example, factuality-aware SFT improves LLMs' truthfulness (row 5 vs 1). In addition, DPO fine-tuning on the factuality preference data guides LLMs to output more truthful responses (rows 3,4 vs 2 and 7 vs 6). Note that we observe that SFT and DPO models show a reverse trend in BLUE and ROUGE. This is likely because SFT models tend to generate shorter responses than the DPO ones do.

In addition, Table 5 reports models' accuracy in tasks of multiple choices from TruthfulQA. No significant differences between models are observed. This is possibly because we mainly focus on the tasks of long-form response generation while TruthfulQA-MC task is formed by short-form answers. The discrepancy between improving LLMs' factuality on long-form and short-form generation is also found by the previous work [Chuang et al., 2024]. Appendix Table 9 reports more evaluation results on other NLP benchmarks.

Table 4: Results on TruthfulQA.

Llama-2 70B	src.	of supervision	TruthfulQA		
2,02	IF. Fact.		BLUE	ROUGE	
(0) Chat	Pro	prietary data	0.21	1.16	
(1) SFT	-	-	0.37	0.20	
(2) + DPO	√	Х	0.03	0.54	
$(3) + DPO^{fact}$	X	\checkmark	0.30	1.12	
(4) + DPO •	✓	\checkmark	0.15	0.80	
(5) SFT*	-	-	0.39	0.51	
(6) + DPO	√	Х	0.07	0.91	
(7) + DPO 6	✓	\checkmark	0.20	0.96	

Table 5: Results on TruthfulQA multiple choices.

Llama-2 70B	src.	of supervision	TruthfulQA-MC			
	IF. Fact.		MC1	MC2	MC3	
(0) Chat	Pro	prietary data	32.2	50.2	25.4	
(1) SFT	-	-	30.8	45.7	23.9	
(2) + DPO	√	Х	30.5	46.0	23.4	
$(3) + DPO^{fact}$	X	\checkmark	31.8	46.8	24.3	
(4) + DPO •	✓	✓	30.8	46.0	23.6	
(5) SFT ⁶	-	-	29.9	44.8	22.5	
(6) + DPO	√	Х	31.5	47.0	24.0	
(7) + DPO⁴	✓	\checkmark	30.5	45.4	23.1	

6 Discussion

6.1 Effects of Fact-Based Instruction Classification

In our factuality-aware alignment, we prompt SFT to judge whether an instruction requires a factual response and apply our factuality alignment strategy to the fact-based instruction. Without the instruction classification, in our factuality-aware SFT, we cannot create supervision from Human and PT responses for respective non-fact-based and fact-based instructions. Instead, for each instruction, we create instruction–response pairs from 1 and 10 responses from Human and PT as supervisions, respectively. Note that, during fine-tuning, for each instruction, we randomly sample instruction–response pair either created from Human or PT with same probability. The SFT model shows

degradation in both instruction-following capability and factuality results, as shown in row 1 vs 2 of Table 6. Second, for factuality-aware DPO, without the instruction classification, we create factuality preference pairs from all instructions instead of fact-based instructions. The DPO fine-tuned model outputs slightly more factual responses but sacrifice instruction-following capability, as shown in row 3 vs 4 of Table 6.

6.2 Effects of Fact-Based Sentence Classification

In addition, we observe that not all the sentences in a response to a fact-based instruction require fact check. For example, given the response, "Of course. The Commodore 64 is a 8-bit home computer that was released by Commodore International in August 1982.", conducting fact check for the first sentence "Of course." is not necessary and may make the factuality reward less accurate. To address this issue, we prompt SFT to judge whether each sentence in a response required fact check using the prompt in Appendix Figure 6. We only conduct fact check and compute factu-

Table 6: Effects of fact-based classification.

	Clas	ssifier	Alpaca Eval	Bio		
	Inst.	Sent.	win rate	FS	# Corr. / Err.	
(1) CECTO 6	Х	-	47.6*	48.4	20.5 / 21.4	
(1) (2) SFT ⁶	✓	-	51.2*	49.5	19.9 / 19.5	
(3)	Х	Х	46.8△	46.8	21.7 / 25.3	
(4) SFT + DPO ⁶	✓	Х	51.7△	45.0	23.7 / 30.3	
(5)	✓	✓	51.3△	42.9	25.5 / 36.8	

^{*} comparing with SFT baseline, SFT.

ality rewards for those fact-based sentences. However, as shown in Table 6, computing factuality rewards for fact-based sentences makes our factual alignment less effective (row 5 vs 4). This is likely because the fact-based sentence classifier is not accurate enough and brings noise into our factuality reward model (see examples in Appendix Figure 7).

6.3 Ablations on Factuality Preference Data Creation

In this section, we examine different ways of creating factuality preference data for factuality-aware DPO training. First, for each fact-based instruction, instead of choosing the responses (among the 4 generated responses) with the maximum and minimum factuality rewards (RM^{fact}) as the respective positive and negative samples, we enumerate all the possible response pairs and choose the response with higher (lower) RM^{fact} as the positive (negative) sample from each enu-

Table 7: Ablation on factuality preference data.

Factuality p	Factuality preference data						
Reward model	Pos.,Neg.	# pairs	win rate△	FS			
RM ^{fact}	max, min	3,315	51.7	44.9			
RM^{fact}	enum.	5,126	50.7	45.0			
$RM^{IF} + 5*RM^{fact}$	max, min	6,340	50.1	45.1			

[△] comparing with DPO baseline, SFT + DPO.

merated pair. If the difference of RM^{fact} is smaller than 0.2, we treat them as equal and discard the pairs. Note that for both row 1 and 2 in Table 7, we also discard the pairs with the difference of instruction-following rewards ($RM^{\rm IF}$) larger than 0.5 (as mentioned in Section 4.2.2). Alternatively, for each response, we linearly combine the rewards, $RM^{\rm IF}$ (1–5 scale) and $RM^{\rm fact}$ (0–1 scale), with the respective weight of 1 and 5 as a composite reward. For each instruction, we choose the responses with the maximum and minimum composite rewards as the positive and negative. As a result, both data creation approaches increase the number of factuality preference pairs; however, they yield trivial improvement in factuality but slight degrade in instruction following (rows 2, 3 vs 1). This result also indicates that leveraging a single reward model, which can be incorporated with PPO [Schulman et al., 2017], is possible to improve both models' instruction-following capability and factuality

6.4 Impacts of DPO on Generation Length

Table 8 lists the average length of models' responses for each dataset. We observe that DPO fine-tuned models tend to output lengthy responses than SFT except for DPO^{fact} on Biography. This trend indicates that our instruction-following reward model RM^{IF} guides LLMs to output more detailed and lengthy responses. In addition, we observe that although DPO out-

Table 8: Effects of DPO on response length.

	Alpaca Eval	Bio	Alpaca Fact	FAVA
(1) SFT	897	1221	969	912
(2) + DPO	1470	1494	1586	1540
$(3) + DPO^{fact}$	1160	1166	1192	1104
(4) + DPO	1474	1395	1528	1422

 $^{^{\}triangle}$ comparing with DPO baseline, SFT + DPO.

puts responses with similar length as DPO on Alpaca Eval, DPO generates a slightly shorter responses for the fact-based instructions in the other three datasets. This results show that our factuality-aware DPO training mainly impacts models' responses for fact-based instructions. The impact is mainly to reduce the false claims, evidenced by the numbers of erroneous facts in rows 2 and 4 of Table 3).

6.5 Case Studies

Figure 10 (in Appendix) showcases the generations of different models, SFT, SFT + DPO and SFT + DPO, on Alpaca Eval and Biography. Given the instruction, "What are the names of some famous actors that started their careers on Broadway?", SFT only lists some names of Broadway actors while DPO fine-tuned models generate detailed information for each listed Broadway actor. As for biography generations, we observe that given the instruction to generate a biography for a rare name entity, Marianne McAndrew, SFT + DPO generates a detailed response but with many wrong facts while SFT and SFT + DPO give relatively short responses. For the frequent entity, Ji Sung, all the models generate detailed and mostly correct responses. This qualitative analysis shows that SFT + DPO tends to generate detailed responses for most instructions, but for those instructions required tailed knowledge (e.g., rare entity) likely unknown to LLMs [Mallen et al., 2023], it reduces erroneous facts by giving less detailed responses, which is also observed by Kang et al. [2024].

7 Conclusion and Future Work

In this paper, we present a study to enhance the factuality of large language models (LLMs). We first identify that the standard alignment approach, comprising SFT and RLAIF with DPO, may inadvertently encourage LLMs to produce more erroneous facts. Specifically, during the SFT stage, fine-tuning LLMs with high-quality human responses may introduce unfamiliar information, prompting LLMs to output unknown facts. Additionally, during the DPO stage, enhancing LLMs' ability to follow instructions may result in more detailed and lengthy responses but often leads to increased hallucination. To tackle the shortcomings of the standard alignment, we propose a factuality-aware alignment method, which includes factuality-aware SFT and DPO. Quantitative and qualitative analyses demonstrate that our factuality-aware alignment not only guides LLMs to generate detailed and helpful responses but also helps prevent the generation of false claims.

While we have successfully integrated factuality into standard alignment procedure, our work only considers two alignment skill sets: instruction following (or helpfulness) and factuality. In practice, each instruction may require consideration of multiple and distinct alignment skill sets [Saha et al., 2024]. The method to optimize for these skill sets tailored to each query requires further study. In our experiments, we note that optimizing preferences solely for instruction following or factuality could potentially compromise the other. While our factuality-aware alignment demonstrated improvements in both aspects, it is uncertain whether there is a trade-off between the two aspects when integrating our approach to large-scale alignment [Touvron et al., 2023]. Finally, as shown in Appendix Figure 7, not all the claims (or sentences) in a response require fact verification, a more accurate factuality reward model should take this factor into account. While our preliminary experiment, which removes non-fact-based sentences from the factuality reward modeling (Section 6.2), shows suboptimal performance, we believe that further study can bring more insights.

Acknowledgements

We thank Bhargavi Paranjape for sharing fine-tuned Llama-2 7B for atomic fact decomposition and Jing Xu, Weizhe Yuan and Jason Weston for their helpful suggestions.

References

Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, Ben Mann, and Jared Kaplan.

- Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv*:2204.05862, 2022.
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal Ndousse, Kamile Lukosuite, Liane Lovitt, Michael Sellitto, Nelson Elhage, Nicholas Schiefer, Noemi Mercado, Nova DasSarma, Robert Lasenby, Robin Larson, Sam Ringer, Scott Johnston, Shauna Kravec, Sheer El Showk, Stanislav Fort, Tamera Lanham, Timothy Telleen-Lawton, Tom Conerly, Tom Henighan, Tristan Hume, Samuel R. Bowman, Zac Hatfield-Dodds, Ben Mann, Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, and Jared Kaplan. Constitutional AI: Harmlessness from AI feedback. arXiv:2212.08073, 2023.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In *Proc. NeurIPS*, pages 1877–1901, 2020.
- Jifan Chen, Aniruddh Sriram, Eunsol Choi, and Greg Durrett. Generating literal and implied subquestions to fact-check complex claims. In *Proc. EMNLP*, pages 3495–3516, 2022.
- Lichang Chen, Shiyang Li, Jun Yan, Hai Wang, Kalpa Gunaratna, Vikas Yadav, Zheng Tang, Vijay Srinivasan, Tianyi Zhou, Heng Huang, and Hongxia Jin. AlpaGasus: Training a better Alpaca model with fewer data. In *Proc. ICLR*, 2024a.
- Lichang Chen, Chen Zhu, Davit Soselia, Jiuhai Chen, Tianyi Zhou, Tom Goldstein, Heng Huang, Mohammad Shoeybi, and Bryan Catanzaro. ODIN: Disentangled reward mitigates hacking in RLHF. *arXiv*:2402.07319, 2024b.
- Silei Cheng, Zhe Gan, Zhengyuan Yang, Shuohang Wang, Jianfeng Wang, Jordan Boyd-Graber, and Lijuan Wang. Prompting GPT-3 to be reliable. In *Proc. ICLR*, 2023.
- I-Chun Chern, Steffi Chern, Shiqi Chen, Weizhe Yuan, Kehua Feng, Chunting Zhou, Junxian He, Graham Neubig, Pengfei Liu, et al. FacTool: Factuality detection in generative AI–a tool augmented framework for multi-task and multi-domain scenarios. *arXiv:2307.13528*, 2023.
- Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon Kim, James R. Glass, and Pengcheng He. DoLa: Decoding by contrasting layers improves factuality in large language models. In *Proc. ICLR*, 2024.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems, 2021.
- Yann Dubois, Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. AlpacaFarm: A simulation framework for methods that learn from human feedback. *arXiv:2305.14387*, 2024.
- Luyu Gao, Zhuyun Dai, Panupong Pasupat, Anthony Chen, Arun Tejasvi Chaganty, Yicheng Fan, Vincent Zhao, N. Lao, Hongrae Lee, Da-Cheng Juan, and Kelvin Guu. RARR: Researching and revising what language models say, using language models. In *Proc. ACL*, 2022.
- Claire Gardent, Anastasia Shimorina, Shashi Narayan, and Laura Perez-Beltrachini. The WebNLG challenge: Generating text from RDF data. In *Proc. INLG*, pages 124–133, 2017.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In *Proc. ICLR*, 2021.
- Tom Hosking, Phil Blunsom, and Max Bartolo. Human feedback is not gold standard. In *Proc. ICLR*, 2024.

- Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. Atlas: Few-shot learning with retrieval augmented language models. *Journal of Machine Learning Research*, pages 1–43, 2023.
- Zhengbao Jiang, Frank Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. Active retrieval augmented generation. In *Proc. EMNLP*, pages 7969–7992, 2023.
- Nikhil Kandpal, Haikang Deng, Adam Roberts, Eric Wallace, and Colin Raffel. Large language models struggle to learn long-tail knowledge. In *Proc ICML*, 2023.
- Katie Kang, Eric Wallace, Claire Tomlin, Aviral Kumar, and Sergey Levine. Unfamiliar finetuning examples control how language models hallucinate. *arXiv:2403.05612*, 2024.
- Andreas Köpf, Yannic Kilcher, Dimitri von Rütte, Sotiris Anagnostidis, Zhi-Rui Tam, Keith Stevens, Abdullah Barhoum, Nguyen Minh Duc, Oliver Stanley, Richárd Nagyfi, Shahul ES, Sameer Suri, David Glushkov, Arnav Dantuluri, Andrew Maguire, Christoph Schuhmann, Huu Nguyen, and Alexander Mattick. OpenAssistant Conversations democratizing large language model alignment. *arXiv*:2304.07327, 2023.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive NLP tasks. In *Proc. NeurIPS*, pages 9459–9474, 2020.
- Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Inference-time intervention: Eliciting truthful answers from a language model. In *Proc. NeurIPS*, 2023.
- Xian Li, Ping Yu, Chunting Zhou, Timo Schick, Omer Levy, Luke Zettlemoyer, Jason E Weston, and Mike Lewis. Self-alignment with instruction backtranslation. In *Proc. ICLR*, 2024.
- Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, July 2004.
- Sheng-Chieh Lin, Akari Asai, Minghan Li, Barlas Oguz, Jimmy Lin, Yashar Mehdad, Wen-tau Yih, and Xilun Chen. How to train your DRAGON: Diverse augmentation towards generalizable dense retrieval. In *Proc. Findings of EMNLP*, pages 6385–6400, 2023.
- Stephanie Lin, Jacob Hilton, and Owain Evans. TruthfulQA: Measuring how models mimic human falsehoods. In *Proc. ACL*, pages 3214–3252, 2022.
- Yixin Liu, Alex Fabbri, Pengfei Liu, Yilun Zhao, Linyong Nan, Ruilin Han, Simeng Han, Shafiq Joty, Chien-Sheng Wu, Caiming Xiong, and Dragomir Radev. Revisiting the gold standard: Grounding summarization evaluation with robust human evaluation. In *Proc. ACL*, pages 4140–4170, 2023.
- Chaitanya Malaviya, Subin Lee, Sihao Chen, Elizabeth Sieber, Mark Yatskar, and Dan Roth. ExpertQA: Expert-curated questions and attributed answers. *arXiv*:2309.07852, 2023.
- Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. When not to trust language models: Investigating effectiveness of parametric and non-parametric memories. In *Proc. ACL*, pages 9802–9822, 2023.
- Grégoire Mialon, Roberto Dessi, Maria Lomeli, Christoforos Nalmpantis, Ramakanth Pasunuru, Roberta Raileanu, Baptiste Roziere, Timo Schick, Jane Dwivedi-Yu, Asli Celikyilmaz, Edouard Grave, Yann LeCun, and Thomas Scialom. Augmented language models: a survey. *Transactions on Machine Learning Research*, 2023.
- Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. FActScore: Fine-grained atomic evaluation of factual precision in long form text generation. In *Proc. EMNLP*, pages 12076–12100, 2023.
- Abhika Mishra, Akari Asai, Vidhisha Balachandran, Yizhong Wang, Graham Neubig, Yulia Tsvetkov, and Hannaneh Hajishirzi. Fine-grained hallucination detection and editing for language models. *arXiv:2401.06855*, 2024.

- OpenAI. GPT-4 technical report. arXiv:2303.08774, 2023.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In *Proc. NeurIPS*, pages 27730–27744, 2022.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. BLEU: a method for automatic evaluation of machine translation. In *Proc. ACL*, pages 311–318, 2002.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. In *Proc. NeurIPS*, pages 53728–53741, 2023.
- Nazneen Rajani, Lewis Tunstall, Edward Beeching, Nathan Lambert, Alexander M. Rush, and Thomas Wolf. No robots. *Hugging Face repository*, 2023.
- Swarnadeep Saha, Omer Levy, Asli Celikyilmaz, Mohit Bansal, Jason Weston, and Xian Li. Branch-Solve-Merge improves large language model evaluation and generation. In *Proc. NAACL*, pages 8352–8370, 2024.
- Victor Sanh, Albert Webson, Colin Raffel, Stephen Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey, M Saiful Bari, Canwen Xu, Urmish Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Fevry, Jason Alan Fries, Ryan Teehan, Teven Le Scao, Stella Biderman, Leo Gao, Thomas Wolf, and Alexander M Rush. Multitask prompted training enables zero-shot task generalization. In *Proc. ICLR*, 2022.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv:1707.06347*, 2017.
- Prasann Singhal, Tanya Goyal, Jiacheng Xu, and Greg Durrett. A long way to go: Investigating length correlations in RLHF. *arXiv*:2310.03716, 2023.
- Zhiqing Sun, Yikang Shen, Hongxin Zhang, Qinhong Zhou, Zhenfang Chen, David Daniel Cox, Yiming Yang, and Chuang Gan. SALMON: Self-alignment with principle-following reward models. In *Proc. ICLR*, 2024.
- Katherine Tian, Eric Mitchell, Huaxiu Yao, Christopher D Manning, and Chelsea Finn. Fine-tuning language models for factuality. In *Proc. ICLR*, 2024.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models. *arXiv:2307.09288*, 2023.
- Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva Naik, Arjun Ashok, Arut Selvan Dhanasekaran, Anjana Arunkumar, David Stap, Eshaan Pathak, Giannis Karamanolakis, Haizhi Lai, Ishan Purohit, Ishani Mondal, Jacob Anderson, Kirby Kuznia, Krima Doshi, Kuntal Kumar Pal, Maitreya Patel, Mehrad Moradshahi, Mihir Parmar, Mirali Purohit, Neeraj Varshney, Phani Rohitha Kaza, Pulkit Verma, Ravsehaj Singh Puri,

- Rushang Karia, Savan Doshi, Shailaja Keyur Sampat, Siddhartha Mishra, Sujan Reddy A, Sumanta Patro, Tanay Dixit, and Xudong Shen. Super-NaturalInstructions: Generalization via declarative instructions on 1600+ NLP tasks. In *Proc. EMNLP*, pages 5085–5109, 2022.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-Instruct: Aligning language models with self-generated instructions. In *Proc. ACL*, pages 13484–13508, 2023a.
- Yuxia Wang, Revanth Gangi Reddy, Zain Muhammad Mujahid, Arnav Arora, Aleksandr Rubashevskii, Jiahui Geng, Osama Mohammed Afzal, Liangming Pan, Nadav Borenstein, Aditya Pillai, Isabelle Augenstein, Iryna Gurevych, and Preslav Nakov. Factcheck-GPT: End-to-end fine-grained document-level fact-checking and correction of LLM output. *arXiv:2311.09000*, 2023b.
- Seonghyeon Ye, Doyoung Kim, Sungdong Kim, Hyeonbin Hwang, Seungone Kim, Yongrae Jo, James Thorne, Juho Kim, and Minjoon Seo. FLASK: Fine-grained language model evaluation based on alignment skill sets. In *Proc. ICLR*, 2024.
- Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Xian Li, Sainbayar Sukhbaatar, Jing Xu, and Jason E Weston. Self-Rewarding language models. In *Proc. ICML*, pages 57905–57923, 2024.
- Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, Susan Zhang, Gargi Ghosh, Mike Lewis, Luke Zettlemoyer, and Omer Levy. LIMA: Less is more for alignment. In *Proc. NeurIPS*, pages 55006–55021, 2023.

A Appendix

A.1 Biography Data Generation

Entities for Training and Evaluation. We use 500 diverse human entities to create training data for SFT and DPO; then, evaluate LLMs' generation factuality on another 183 human entities from Min et al. [2023]. Note that the human entities for training and evaluation are uniformly sampled from entities across diverse nationalities, professions, and rarities. The instruction is generated with the format: Tell me a bio of entity name.

Creating Training Data for SFT. We randomly sample 5 human entities among the 500 entities for training and generate their biographies using Llama-2-Chat 70B as 5-shot demonstration. With the 5-shot demonstration, we use pre-trained Llama-2 7B to generate 10 biographies for each human entity from the remaining 495 ones. We set temperature 0.7 and top-p 0.9 when generate multiple responses from LLMs in all our experiments. We use the created 4,950 name entity—biography pairs to fine-tune the pre-trained Llama-2 7B. As for generating training data with RAG, we prepend the top-10 passages from our retrieval system (detailed in Appendix A.2) to each instruction and generate 10 biographies for each entity from RAG with 5-shot demonstrations. Note that we only prepend top-1 passage for each instruction in the demonstration.

Creating Factuality Preference Pairs for DPO. To construct factuality preference pairs, we first compute FACTSCORE (FS) for all the 4,950 biographies previously created by PT. Then, for each name entity, we compare the FS for all the possible 45 pairs from the 10 generated biographies and construct DPO pairs using the biography with a higher (lower) FS as a positive (negative). Note that we discard the pairs if they show tied FS.

A.2 Retrieval Models

For each query, we retrieve top-20 candidate passages from Wikipedia using DRAGON+ [Lin et al., 2023] and re-rank the candidates using a 12-layer cross-encoder¹². We use the Wikipedia version from the Dec. 20, 2021 dump released by Izacard et al. [2023] in this work.

A.3 Alignment with Self Rewarding

SFT. At SFT stage, we fine-tune PT on two seed datasets: (1) Instruction-following training (IFT) data from Li et al. [2024], consisting of 3200 instruction-response pairs created by humans from Open Assistant dataset [OASST; Köpf et al., 2023], where we only use the first conversational turns in the English that are annotated rank 0; (2) evaluation following training (EFT) data from Yuan et al. [2024], the LLM-as-a-Judge data consists of 1630 samples, each of which contains instruction, human response and the corresponding score of 1-5 scale (with chain-of-though evaluation reasoning): (x, y, r), where (x, y) pairs are also selected from OASST other than training pairs and r is created by the model fine-tuned only on IFT with manual filtering. The purpose of EFT is to enhance a LLM's capability as a reward model to judge the quality of a response in terms of relevance, coverage, usefulness, clarity and expertise. We refer readers to Yuan et al. [2024] for how EFT is created and filtered with minimum human efforts. The prompt template for LLM-as-a-Judge in EFT and an EFT training sample are shown in Appendix, Figure 8 and 9. We refer the baseline model fine-tuned on the IFT and EFT datasets as SFT.

DPO for Instruction Following. At the subsequent preference learning with DPO, following Wang et al. [2023a], we augment additional 20K instructions with Llama-2 70B chat model. For each augmented instruction x, we use SFT to generate 4 responses and evaluate how well the responses follow the instruction with score of 1–5 scale: $RM^{IF}(x, y_0) \cdots ; RM^{IF}(x, y_3)$, where

⁹https://github.com/shmsw25/FActScore

¹⁰https://huggingface.co/meta-llama/Llama-2-70b-chat-hf

¹¹https://huggingface.co/meta-llama/Llama-2-7b

¹²https://huggingface.co/sentence-transformers/all-MiniLM-L12-v2

¹³https://huggingface.co/datasets/OpenAssistant/oasst1

¹⁴https://huggingface.co/meta-llama/Llama-2-70b-chat-hf

 $y_0, \cdots, y_3 \in \mathrm{SFT}(x)$ and $\mathrm{RM^{IF}}$ is the instruction-following reward model. Note that, in self-rewarding [Yuan et al., 2024], $\mathrm{RM^{IF}}$ is the same as SFT model. In addition, for each instruction-response pair, we use the same prompt in EFT seed data to sample the chain-of-thought evaluation three times and average the scores as the reward. Finally, for each instruction, we use the response with the highest (lowest) reward as the positive (negative) sample to form a preference pair for DPO training: (x,y_+,y_-) . We discard the pair, if $\mathrm{RM^{IF}}(x,y_+) = \mathrm{RM^{IF}}(x,y_-)$. In the DPO training, the model is initialized from SFT and the fine-tuned model is denoted SFT + DPO.

A.4 More Evaluation Results on Standard Benchmarks

Table 9 compares the instruction fine-tuned models' (w/o and w/ involving FLAME) accuracy the in tasks of MMLU [Hendrycks et al., 2021] and GSM8K [Cobbe et al., 2021]. A slight drop from FLAME is observed. This is possibly because we mainly focus on the tasks of long-form response generation while MMLU and GSM8K are the

Table 9: Results on MMLU and GSM8K.

Llama-2 70B	src. IF.	of supervision Fact.	MMLU	GSM8K
(1) SFT+ DPO	√	Х	69.3	59.3
(2) SFT ⁶ + DPO ⁶	✓	✓	69.1	58.2

benchmarks with short-form answers. The discrepancy between improving LLMs' factuality on long-form and short-form generation is also found by the previous work [Chuang et al., 2024].

A.5 Factuality Reward Modeling

Factuality Reward Models. We build a reward model RM^{fact} to measure the factuality of each response. The factuality reward model consists of two main modules. (1) fact decomposition: we first use nltk.tokenize to split a response into sentences; then, use our Llama-2 7B model fine-tuned on public datasets [Liu et al., 2023, Chen et al., 2022, Malaviya et al., 2023] to conduct atomic fact decomposition for each sentence. (2) Retrieval augmented claim verification: for each decomposed fact (or claim), we use the instruct Llama 7B fine-tuned on Super Natural Instructions [Wang et al., 2022] to do fact check with the prompt shown in Figure 5. We append 10 retrieved supports (using the instruction as query) from our retrieval and re-ranking pipeline in Appendix A.2. Then, we compute the proportion of correct atomic facts in a response as a factuality reward.

Quality of Factuality Reward Models. We conduct ablation study on our factuality reward models. Specifically, we use our factuality reward models to detect the number of error facts in each instruction—response pair. We try different models for fact check using the prompt shown in Figure 5 with different numbers of retrieved supports. We use the LLMs' generated responses with human annotated hallucination provided by Mishra et al. [2024] to evaluate the quality of the factuality reward models. ¹⁷ Specifically, we rank the responses by numbers of errors detected and calculate the Kendall rank correlation (τ)

Table 10: A comparison of factuality reward models. au denotes the correlation between human annotation.

	fact check model	# sup.	fact unit	τ
(1)	Instruct Llama 7B	5	atom.	0.32
$\frac{(2)}{(3)}$		10		0.34
(4)	SFT (Llama-2 70B)	10	atom.	0.31
(5)	Instruct Llama 7B	5	sent.	0.20
(6)	mondet Eluma / B	10	som.	0.25

between the rank lists by our factuality reward models and humans. As shown in Table 10, conducing fact check with more retrieved supports improves the accuracy of the factuality reward models (row 2 vs 1). In addition, our SFT, only fine-tuned on the IFT and EFT data, is capable of doing fact check, compared to Instruct Llama 7B fine-tuned on Super Natural Instructions [Wang et al., 2022]. Finally, instead of computing the number of error facts from decomposed atomic facts, we conduct fact check directly for each sentence in a response and calculate the number of false sentences as error facts. However, the quality of the reward models shows significant decrease (rows 5,6 vs 1,2). We finally adopt row 2 as our factuality reward model.

¹⁵With few-shot demonstration, SFT is able to decompose a sentence into atomic facts with acceptable accuracy. Fine-tuning a Llama-2 7B is to reduce the inference time.

¹⁶https://huggingface.co/kalpeshk2011/instruct-llama-7b-wdiff

¹⁷https://huggingface.co/datasets/fava-uw/fava-data/blob/main/ annotations.json

Table 11: Training data statistics for different variants. IF. and Fact. denote instruction following (x,y_+,y_-) and factuality $(x\in X^{\mathrm{fact}},y_{\mathrm{true}},y_{\mathrm{false}})$ preference data, where X^{fact} denotes the set of fact-based instructions.

	Seed IFT	(# of Inst.)	Preference (# of pairs)		
model variant	$x \notin X^{\text{fact}}$	$x \in X^{\mathrm{fact}}$	IF.	Fact.	
SFT + DPO			18,454	-	
$SFT + DPO^{fact}$	2,187	1,013	-	3,315	
SFT + DPO⁰			18,454	3,315	
SFT⁰ + DPO	2,187	1.013	18,603	-	
SFT⁰ + DPO⁰	2,107	1,013	18,603	4,211	

A.6 Training Details

We fine-tune our models for 500 steps with a batch size of 32 and 64 on respective SFT and DPO stages. The learning rate and maximum sequence length is set to 1e-6 (which decays to 1e-7) and 2048, respectively. At SFT stage, we mix the IFT and EFT while at DPO stage, we set $\beta=0.1$ and uniformly sample between self rewarding (x,y_+,y_-) and factuality reward $(x,y_{\rm true},y_{\rm false})$ preference data. Note that SFT (SFT $^{\bullet}$) + DPO meaning that we use SFT (SFT $^{\bullet}$) to create preference data, serve as instruction-following reward model RM $^{\rm IF}$ and as the initialization of DPO. The data used to fine-tune different variants are listed in Table 11. We conduct fine-tuning with full parameters on 64 NVIDIA A100 (80GB) GPUs. SFT and DPO require around 1 and 4 hours to complete, respectively.

A.7 Limitations

Limitations. In this paper, we use FACTSCORE as the main metric of factuality. However, as highlighted by Min et al. [2023], "FACTSCORE focuses on factual precision is only one aspect of the broader factuality problem." In addition, the validity of applying FACTSCORE to broader domains other than Biography has not been carefully studied. In our experiments, we can only claim that models' generation shows high agreement with the facts from Wikipedia, without considering whether the claims are open-ended and debatable. Finally, reward hacking is the potential issue of LLM alignment [Chen et al., 2024b] but the risk of reward hacking is still unknown when considering factuality.

A.8 Broader Impacts

Increased factual accuracy in AI-generated responses may contribute to building trust in AI technology. Users are more likely to trust and rely on AI systems if they consistently provide accurate information, leading to broader adoption and acceptance of AI-driven solutions. However, FLAME does not deal with the facts which are open-ended and debatable and uses Wikipedia as the only source of facts. This may inadvertently reinforce some biases in the AI system's responses, potentially perpetuating societal inequalities.

[INST] Here is a question from a user: "<QUESTION>". To answer the above question, do you need the factual knowledge from Wikipedia? Give an answer using the format: "Answer: Yes or No". [/INST]

Figure 4: Prompt to check whether an instruction is fact-based.

[INST] Answer the question about '<Instruction>' based on the given context.\n\n'<Retrieved Supports>'\n\nInput: '<ATOMIC_FACT or SENTENCE>' True or False?\nOutput: [/INST]

Figure 5: Prompt for fact check.

[INST] Does the given sentence contain factual knowledge from Wikipedia? Question: '<QUERY>'; Sentence: '<CLAIM>'. Answer with Yes or No. [/INST]

Figure 6: Prompt to check whether a claim is fact-based.

Do you have any information about the Commodore 64?	Mehod (step by step if possible), please? -> " x^3 + 2x^2 - x - 1 = 0 "	Explain the significance of the American Revolution, including the events that led up to it, the impact it had on the world, and its ongoing relevance today.
It is one of the most successful home computers of all time, with over 17 million units sold. — True It featured 64 kilobytes of RAM, color graphies, and the SID sound chip, — True, It was widely used for gaming, programming, and home productivity applications. — True It was succeeded by the Commodore 128. — False There is still a thirving community of Commodore 64 enthusiasts around the world who continue to enjoy playing and programming on the system. — True You can find more information about the Commodore 64 on various websites and forums dedicated to the system. — False	Cardano's method. – False Here are the steps: Step 1: Path exquation in the form $x'3 + px + q = 0$ $x'3 + 2x'2 - x - 1 = 0$ can be rearranged as $x'3 + 0x'2 + 2x'2 - x - 1 = 0$, so we have $p = 0$ and $q = 1 - 1$. False Step 2: Substitute $x = y + 0x^3$ in the equation, where $b = -p/3$ a $0 - p^2$ and $0 - 1$ such that $y = 0 - 1$ and $y = 0 - 1$	were still under British rule and were becoming increasingly discontent with the taxes and regulations imposed upon them by the British government.— True One of the most significant events leading up to the Revolution was the Boston Tea Party, where American colonists dumped British tea into the harbor as a symbol of their anger towards the British tea tax.— True This event was followed by a series of skirmishes between British troops and American colonists, which escalated into full-blown warfare.— I absent the superior of the American Revolution was far-reaching, both domestically. The impact of the American Revolution was far-reaching, both domestically be reported by the property of the Revolution in majoried people around the world to fight for their own freedoms and led to the spread of democratic ideals.—True in the United States, the Revolution led to the formation of a new government based on democratic principles, where citizens had a say in the laws that governed them.—True

Figure 7: The results of whether a sentence is fact-based or not classified by SFT with prompt in Figure 6.

[INST] Review the user's question and the corresponding response using the additive 5-point scoring system described below. Points are accumulated based on the satisfaction of each criterion: - Add 1 point if the response is relevant and provides some information related to the user's inquiry, even if it is incomplete or contains some irrelevant content. - Add another point if the response addresses a substantial portion of the user's question, but does not completely resolve the query or provide a direct answer. - Award a third point if the response answers the basic elements of the user's question in a useful way, regardless of whether it seems to have been written by an AI Assistant or if it has elements typically found in blogs or search results. - Grant a fourth point if the response is clearly written from an AI Assistant's perspective, addressing the user's question directly and comprehensively, and is well-organized and helpful, even if there is slight room for improvement in clarity, - Bestow a fifth point for a response that is impeccably tailored to the user's question by an AI Assistant, without extraneous information, reflecting expert knowledge, and demonstrating a high-quality, engaging, and insightful answer. User: <INSTRUCTION_HERE> <response><RESPONSE HERE></response> After examining the user's instruction and the response: - Briefly justify your total score, up to 100 words. - Conclude with the score using the format: "Score: <total points>" Remember to assess from the AI Assistant perspective, utilizing web search knowledge as necessary. To evaluate the response in alignment with this additive scoring model, we'll systematically attribute points based on the outlined criteria. [/INST]

Figure 8: Prompt to evaluate models' instruction-following capability from Yuan et al. [2024].

[INST] Review the user's question and the corresponding response using the additive 5-point scoring system described below. Points are accumulated based on the satisfaction of each criterion:\n\n-Add 1 point if the response is relevant and provides some information related to the user's inquiry, even if it is incomplete or contains some irrelevant content. Add another point if the response addressess a substantial portion of the user's question, but does not completely resolve the query or provide a direct answer.\n-Award a third point if the response answers the basic elements of the user's question in a useful way, regardless of whether it seems to have been written by an AI Assistant or if it has elements typically found in blogs or search results.\n-Grant a fourth point if the response is clearly written from the state of the provided of the user's question by an AI Assistant, without extraneous information, reflecting expert in the response is clearly written from an insight found in high the state of the user's question by an AI Assistant, without extraneous information, reflecting expert who leaded classified qovernment documents to the press, starting from when he first oot employed by the U.S. government and ending shortly after he landed in Russia. \n\n\n\response\mathbb{Tesponse

Figure 9: An example of EFT data. The texts with the colors of green, red and blue are the instruction, response and the LLM-as-a-judge results (explanation and score), respectively

SFT + DPO **SFT** SFT + DPOWhat are the names of some famous actors that started their careers on Broadway? | 1. Robert Redford'n2. Meryl Streep\n3. Al Pacino\n4. Denzel | Washington\n5. Liza Minnellino. James Earl Jones\n7. | | Bernadette Peters\n8. Nathan Lane\n9. Patit LuPone\n9. Patit LuPone\n9. | | Rodway in the late 1950s, appearing in productions such as 1 | | Bernadette Peters\n8. Nathan Lane\n9. Patit LuPone\n9. | | Revin Kline\n11. | | Matthew Broderick\n12. Whoopi | | Goldberg\n13. Hugh Jackman\n14. Idina Menzel\n15. Whoopi | | Goldberg\n13. Hugh Jackman\n14. Idina\n14. Whoopi | | Goldberg\n13. Hugh Jackman\n14. Idina Menzel\n15. Whoopi | | Goldberg\n14. Shop \n15. Hugh Jackman\n14. Idina Menzel\n15. Whoopi | | Gery\n14. John Clalling\n14. John Clalling\n14. John Shopi | | Goldberg\n14. John Clalling\n14. John Shopi | | Goldberg\n14. John Clalling\n14. John Cl Robert Redford\n2. Meryl Streep\n3. Al Pacino\n4. Denzel Robert Redford: Redford began his acting career on There are many famous actors who started their careers on Broadway. Here are some of the most well-known:\n1. Meryl Broadway. Here are some of the most well-known.'nl. Meryl Streep: Streep made her Broadway debut in 1975 in the play "Trelawny of the 'Wells." She went on to win three Academy Awards and become one of the most acclaimed actresses of her generation.'nl. Robert De Niro. De Niro made his Broadway debut in 1965 in the play "The Indian Wants the Broax." He went on to become one of the most iconic actors in Hollywood, winning two Academy Awards.'nl. Al Pacino: Pacino made his Broadway debut in 1969 in the play "Does a Tiger Wear a Necktie?" He went on to become one of the most respected actors in Hollywood, winning an Academy Award and numerous actors in Hollywood, winning an Academy Award and numerous actors in Hollywood, winning an Academy Award and numerous other accolades had. Dustin Hoffman made his Broadway debut in 1961 in the play "A Cook for Mr. General." He went on to become one of the most versatile actors in Hollywood, winning two Academy Awards. hs. Denzel Washington: Washington made his Broadway debut in 1982 in the play "A Soldier's Play." He went on to become one of the most acclaimed actors of his generation, winning two Academy Awards. hs. Bernadette Peters: Peters made her Broadway debut in 1968 in the musical "George M!" She went on to become one of the most beloved performers on Broadway, winning three Tony Awards. hs. Patti LuPone: LuPone made her Broadway debut in 1972 in the musical "The Threepenny Opera." She went on to become one of the most become one of the most become one of the most respective performers on the three penny Opera." She went on to become one of the most respected performers on debut in 1972 in the musical "The Threepenny Opera." She went on to become one of the most respected performers on Broadway, winning three Tony Awards n8. Nathan Lane: Lane made his Broadway debut in 1982 in the play "Present Laughter." He went on to become one of the most beloved performers on Broadway, winning three Tony Awards n9. Audra McDonald: McDonald made her Broadway debut in 1992 in the musical "The Secret Garden." She went on to become one of the most acclaimed performers on Broadway, winning ist: Tony Awards n10. Hugh Jackman: Jackman made (cut by maximum seneration leneth) \n10. Bernadette Peters: Peters made her Broadway debut in 1968 in the musical \"George M!\" She went on to win three Tony Awards for her work on Broadway.\n11. Patti LuPone: LuPone made her Broadway debut in 1972 in the musical \"The The Threepenny Opera\" She went on to win two Tony Awards for her work on Broadway.\n12. L (cut by maximum generation) eneration length) Tell me the bio of Marianne McAndrew (very rare entity) Marianne McAndrew was an American actress who was best known for her role as Irene Molloy in the 1964 film version of the musical "Hello, Dolly!" She was born on May 27, 1943, in Chicago, Illinois, and began her acting career in the early 1960s. McAndrew made her film debut in 1963 with a small role in the film "The Cardinal." She then went on to appear in several other films, including "The World of Henry Orient" (1964) and "The Great Race" (1965). However, it was her role as Irene Molloy in "Hello, Dolly!" that brought her the most Irene Molloy in "Hello, Dolly!" that brought her the most Irene Molloy in "Hello, Dolly!" that brought her the most Irene Molloy in "Hello, Dolly!" that brought her the most Irene Molloy in "Hello, Dolly!" that brought her the most Irene Molloy in "Hello, Dolly!" that brought acting in the Virginian" and "The F.B.1." She retried from acting in the carly 1970s and later worked as a real estate agent in California. McAndrew was married twice and had two children. She passed away on February 26, 2006, at the age of 62. Overall, Marianne McAndrew was a talented actress who made a lasting impact on the film industry with her memorable 1943-2006/n/mMarianne McAndrew was an American actress best known for her role as Irene Molloy in the 1964 film version of Hello, Dolly! She was born in Fort Ord, California, in 1943 and began her acting career in the early 1960s. She appeared in several television shows and films throughout the decade, including The Fugitive, The Virginian, and The Wild Wild West. In 1964, she was cast as Irene Molloy in Hello, Dolly!, Which earned her a Golden Globe nomination for Best Supporting Actress. After Hello, Dolly!, McAndrew continued to act in films and television shows throughout the 1970s and 1980s. She also appeared on Broadway in the musicals The Music Man and Manne. McAndrew retired from acting in the early 1990s and passed away in 2006 at the age of 63. Il. Marianne McAndrew was an American actress best known for her role as Irene Molloy in the 1968 film version of "Hello, Dolly!"nn2. She was born on July 20, 1943, in Chicago, Illinois, and began her acting career in the early 1965...3 McAndrew made her film debut in 1965 in the movie "The Slender Thread" and went on to appear in several other films throughout the 1960s and 1970s, including "The Thomas Crown Affair" and "The Great Gatsby:"m4. In addition to her film work, McAndrew also appeared on television, including guest roles on popular shows such as "The Fugitive" and "Wission: Impossible." in McAndrew was married to actor and director Richard Benjamin from 1966 until their divorce in 1967.nds. She later married actor and producer Michael Tolan in 1970, and the couple had two children together.in? McAndrew retried from acting in the late 1970s and focused on McAndrew retried from acting in the late 1970s and focused on McAndrew retired from acting in the late 1970s and focused on McAndrew Teuted rom acting in the late 1970s and roctised on raising her family/m8. She passed away on February 26, 2006, at the age of 62.ln9. McAndrew's legacy as an actress lives on through her memorable performances in films such as "Hello, Dolly!" and "The Thomas Crown Affair." In 10. She is remembered as a talented and versatile actress who made a lasting impact on the film industry. made a lasting impact on the film industry with her memorable performance in "Hello, Dolly!" Tell me the bio of Ji Sung (very frequent entity) 1. Ji Sung is a South Korean actor who was born on February 29, 1977, in Seoul, South Korean actor who was born on February 29, 1979, in Seoul, South Korean actor who was born on February 29, 1979, in Seoul, South Korean actor who was born on February 29, 1979, in Seoul, South Korean actor who was born on February 29, 1979, and has since appeared in numerous television dramas and films. 19, 18 Sung is beat known for his roles in popular dramas actors in South Korean Ji Sung is known for his versatility as an actor, having played a wide range of roles in both television ("nad Prize at the 2015 MBC Drama Awards for his role in performances, including the Grand Prize (Daesang) at the MBC ("Kill Me, Heal Me," M. J. addition to his acting career, Ji Sung las sho known for his philanthropic work. He has been involved with of a man with multiple personalities in "Kill Me, Heal Me," as is married to actress Lee Bo-young, and the couple has a daughter together. In J. He is known for his versatility as an actor, "Defendant," and a doctor in "Doctor John." He has is married to actress Lee Bo-young, and the couple has a daughter together. In J. He is known for his versatility as an actor, "Confession." In addition to his acting career, Ji Sung is also worn for his philanthropic work. He has been involved in "Confession." In addition to his acting career, Ji Sung is also If Sung is a South Korean actor who was born on February 27, 1977, in Seoul, South Korea. He began his acting career in 1999 and has since become one of the most popular and acclaimed actors in South Korea. In'uli Sung attended Hanyang University, where he studied theater and film. He made his acting debut in the 1999 television drama "KAISTV" and has since appeared in numerous television dramas and films. Some of his most notable roles include "All In," "Protect the Boss," "Secret Love," and "Kill Me, Heal Me," "m'nai Sung is known for his versatility as an actor and his shift to portray a vider game of characters. as an actor and his ability to portray a wide range of characters. He has won numerous awards for his performances, including the Grand Pirez (Daesang) at the MBC Drama Awards in 2015 for his role in \"Kill Me, Heal Me.\"n\nln addition to his acting usabilities and helping victims of natural disasters in 6. J. Sung is married to actress Lee Bo-young, and the couple has a daughter together/in? He is known for his versatility as an actor, having played a wide range of characters in his career, from romantic leads to villains.in 8. Ji Sung is also known for his dedication to his craft of the control of the contr

Figure 10: Generation comparisons for instructions from Alpaca Eval and Biography (very rare and frequent entities). Determined through manual verification using Google search, red denotes incorrect identified facts while pink indicates unverified facts; e.g., we cannot search relevant pages about Ji Sung's involvement in charitable causes but also cannot dismiss the possibility of his contributions. Note that the popularity of an entity is defined by its occurrence and page views in Wikipedia, which are provided by Min et al. [2023].

and dedicated fan base both in South Korea and internationally.

\n10. Ji Sung continues to be a prominent figure in the South Korean entertainment industry, with many fans eagerly

anticipating his next project.

\n9. He has a large

Overall, Ji Sung is a talented and respected actor who has made

a significant impact on the South Korean entertainment industry

mental health issues.\n\nJi Sung is married to actress Lee Bo young, and the couple has a daughter together. He is known fo

Sung is a talented and respected actor who has made a

NeurIPS Paper Checklist

1. Claims

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

Answer: [Yes]

Justification: The main claim of our paper is to advance the understanding of the underlying causes of LLM hallucination as well as its relation to the alignment procedure (from lines 22–23) and study how to make the LLM alignment process more factual (lines 5–6). We thus conduct a pilot study to identify the key factors which impact the LLMs' factuality during alignment procedure in Section 3. Based on the finding, we implement the factuality-aware alignment in Section 4 and demonstrate the effectiveness of our approach through comprehensive experiments in Section 5.

Guidelines:

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

2. Limitations

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

Answer: [Yes]

Justification: We include the discussion of limitations in Appendix A.7.

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

3. Theory Assumptions and Proofs

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

Answer: [NA]

Justification: This is an empirical paper; thus, there is no theoretical result included.

Guidelines:

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

4. Experimental Result Reproducibility

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

Answer: [Yes]

Justification: We include all the high-level information to reproduce our models in Section 4.1 and 4.2, and due to space limitation, more detailed information are included in Appendix A.3, A.5 and A.6.

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
- (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility.

In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [NA]

Justification: While we do not provide the code to reproduce the main experimental results, we provide all the necessary information and URL links to training and evaluation data.

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]

Guidelines:

Justification: We provide training details in Appendix A.6 and test details in Section 5.1

- 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: [NA]

Justification: We follow the instruction following evaluation in Alpaca Eval [Dubois et al., 2024] to report the win rate comparisons between models, which is not suitable for statistical significance test. For FACTSCORE, we follow the established procedure to compare models' average FACTSCORE [Min et al., 2023], which is correlated to human evaluation.

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

8. Experiments Compute Resources

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

Answer: [Yes]

Justification: We provide the required computation resources in Appendix A.6.

Guidelines:

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

9. Code Of Ethics

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

Answer: [Yes]

Justification: We have reviewed and made sure our paper conforms the NeurIPS Code of Ethics.

Guidelines:

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

10. Broader Impacts

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

Answer: [Yes]

Justification: We have included the discussion of broader impacts in Appendix A.8. Guidelines:

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

11. Safeguards

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

Answer: [NA]

Justification: We do not release any new data and models in the paper.

Guidelines:

- 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.

12. Licenses for existing assets

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

Answer: [Yes]

Justification: We have cited all the papers, which provide the models and datasets used in our paper.

- 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 has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- 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.

13. New Assets

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

Answer: [NA]

Justification: The paper does not release new assets.

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.

14. Crowdsourcing and Research with Human Subjects

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

Answer: [NA]

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

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

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

Answer: [NA]

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