Quantifying and Optimizing Global Faithfulness in Persona-driven Role-playing

Letian Peng, Jingbo Shang*

Department of Computer Science University of California, San Diego {lepeng, jshang}@ucsd.edu

Abstract

Persona-driven role-playing (PRP) aims to build AI characters that can respond to user queries by faithfully sticking with all (factual) statements in persona documents. Unfortunately, existing faithfulness criteria for PRP are limited to coarsegrained LLM-based scoring without a clear definition or formulation. This paper presents a pioneering exploration to quantify PRP faithfulness evaluation as a fine-grained and explainable criterion, which also serves as a reliable reference for faithfulness optimization. Our criterion first discriminates persona statements into active and passive constraints by identifying the query-statement relevance. Then, we incorporate all constraints following the principle that the AI character's response should be (a) entailed by active (relevant) constraints and (b) not contradicted by passive (irrelevant) constraints. We translate this principle mathematically into a novel Active-Passive-Constraint (APC) score, a constraint-wise sum of statement-to-response natural language inference (NLI) scores weighted by constraint-query relevance scores. In practice, we build the APC scoring system by symbolically distilling small NLI and relevance discriminators (~300M parameters) from GPT-4 for efficiency, and both show high consistency with GPT-4's discrimination. We validate the quality of the APC score against human evaluation based on example personas with tens of statements, and the results show a high correlation. As the APC score could faithfully reflect the PRP quality, we further leverage it as a reward system in direct preference optimization (DPO) for better AI characters. Our experiments offer a fine-grained and explainable comparison between existing PRP techniques, revealing their advantages and limitations. We further find APC-based DPO to be one of the most competitive techniques for sticking with all constraints and can be well incorporated with other techniques. We then extend the scale of the experiments to real persons with hundreds of statements and reach a consistent conclusion. Finally, we provide comprehensive analyses and case studies to support the effectiveness of APC evaluation and APC-based DPO. ²

1 Introduction

Role-playing (Han et al., 2022; Li et al., 2023; Yan et al., 2023; Bianchi et al., 2024; Yu et al., 2024; Tao et al., 2024) is a newborn and trending natural language processing field, emerging from the proficiency of large language models (LLMs) (Brown et al., 2020; OpenAI, 2023; Touvron et al., 2023a,b; Mesnard et al., 2024) in human interaction. Role-playing customized AI characters, which are useful for providing emotional value (Zhang et al., 2024), developing video games (Hu et al., 2024), or even realizing the metaverse (Zhou, 2023; Yue et al., 2024). Persona-driven role-playing

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

^{*} Corresponding author.

²Code, Dataset, Demo: https://github.com/KomeijiForce/Active_Passive_Constraint_Koishiday_2024

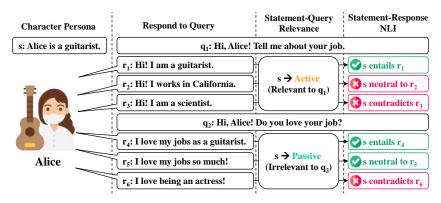


Figure 1: A presentation of the alignment between APC and human's view on PRP faithfulness.

(PRP) (Wang et al., 2023a,b; Shao et al., 2023; Xu et al., 2024) uses only persona statements to efficiently build the AI character without dialogues or scripts, which is extremely useful for real-world applications as few characters have sufficient or accessible dialogues for training.

As the persona statements are the only input in PRP, being faithful to them becomes one of the most crucial objectives for this task. Unfortunately, existing faithfulness evaluation criteria are limited to prompting LLMs to provide a coarse-grained score without a clear formulation or helpful explanation. Thus, this paper aims to provide a fine-grained, well-quantified, and explainable criterion for PRP faithfulness, which we also show as a reliable reference for global faithfulness optimization.

Our criterion views PRP as a constraint satisfaction problem (CSP) (Brailsford et al., 1999), and the whole persona information as a global constraint for the response to satisfy. Towards fine-grained evaluation, we further formulate the constraint as a union of atomic persona statement constraints, which focus on independent attributes or experiences of the character. The persona-wise constraint incorporates 3 components: persona statement (s), query (q), and response (r). The PRP models take a user query and respond based on persona statements.

Our key insights are 1) the statement-to-response constraint depends on query-statement relevance and 2) the statement-to-response constraint can be formalized as statement-to-response natural language inference (NLI). (Bowman et al., 2015) The constraint becomes *active* when the query is relevant to the persona statement, constraining the response to be entailed by the persona statement. The constraint becomes *passive* when the query is irrelevant to the persona statement, reducing the constraint to only not being contradicted by the persona statement. We present a possible PRP instance in Figure 1 to show how our definition is consistent with human's view on PRP faithfulness. As q_1 is relevant to s, s becomes active and constrains the character "Alice" to incorporate the information in s to her response. For irrelevant q_2 , the constraint of s becomes passive and is relaxed to only not incorporating information contradicting s.

We further develop a scoring system to quantify APC, making it more appropriate for evaluating practical PRP methods. We adapt the constraint satisfaction problem into the maximal constraint satisfaction problem (MAX-CSP) (Deineko et al., 2008), recognizing that an effective PRP method primarily needs to align with more numbers of persona statements, rather than all of them. Thus, the quantified APC score sums up the satisfaction probability of the response to each persona statement, representing the expected number of satisfied constraints. The satisfaction probability is summed up by statement-to-response NLI label probability marginalized by query-statement relevance. We also regularize the APC score to Δ APC score with a minuend equal to the reward gained by a PRP system that permanently gives a neutral response. The regularization makes the absolute value more straightforward to reflect faithfulness, representing the expected number of entailed active persona statements (active reward) subtracted by the expected number of contradicted passive persona statements (passive penalty). In practice, the probabilities are efficiently assigned by small discriminators based on DeBERTa-V3 (He et al., 2021) (\sim 300M parameters) symbolically distilled from the state-of-the-art LLM, GPT-4 (OpenAI, 2023) with \sim 90% accuracy.

With the (Δ) APC score, we can reveal the advantages and limitations of existing PRP methods. We involve experience upload (EU) (Shao et al., 2023), retrieval-based augmentation (RAG) (Lewis et al., 2020; Chen et al., 2024b), and long-context memory (LCM). We handcraft 3 original characters with small-scale persona statements (8, 19, 30) and free from data contamination (Magar & Schwartz,

2022) in the pre-training of LLMs. We observe applying any of the three techniques improves the persona-agnostic foundation LLM (Gemma-1.1-7b-it), indicating their benefits to PRP. However, our experiments also confirm that their limitations are significant. EU constructs character experiences based on each persona statement, but these often meet only some constraints and sometimes even violate them, whether actively or passively. RAG adheres more closely to the given personas, incorporating more relevant statements, though it still sometimes misses passive constraints. LCM, on the other hand, loads the entire persona into the context in hopes that the LLM will effectively utilize all persona statements. Our experiment shows that as the number of persona statements increases, LCM's performance deteriorates compared to RAG, confirming findings about limitations in LLMs' handling of long contexts as discussed in Liu et al. (2024b).

Furthermore, we discover the APC score to be a reliable reward for direct preference optimization (DPO) (Rafailov et al., 2023) to strengthen the faithfulness of PRP methods. We use APC and human evaluation to verify the effectiveness of DPO, which benefits the satisfaction of both active and negative constraints. We extend the experiments for evaluation and DPO above to complicated famous figures with $77 \sim 599$ persona statements, further verifying the reached conclusions.

Finally, we launch case studies toward a specific analysis of the insights obtained by APC score-based evaluation and the benefit gained from APC-based DPO. We also showcase how we can explain the detected constraint violation by tracing back and strengthening extra constraints like protective experience by persona statements. Our contribution is three-fold,

- We propose the first formal definition of AI character's global faithfulness and formulate it as a constraint satisfaction problem. The constraint is further quantified as the APC score, which is human-consistent and the first quantified evaluation for AI characters.
- We evaluate potential PRP techniques, EU, RAG, and LCM by APC score, which reveals their properties on active and passive constraints.
- We find APC-based DPO to be one of the most competitive techniques to improve the global faithfulness of AI characters and cooperate well with other methods.

2 Related Works

With the emergence of the high capability of LLMs in interaction with humans, role-playing AI has attracted lots of attention from both academia (Shanahan et al., 2023) and industry³. The difference between role-playing and normal agents is the demand of following a constant persona. The main aim of role-playing includes personalizing the agent for the user preference (Jang et al., 2023) and bringing virtual characters to the real world (Li et al., 2023; Tao et al., 2024). Role-playing agents also have wide potential application scenarios, such as emotional accompanying and building virtual world (Zhang et al., 2024; Hu et al., 2024; Zhou, 2023; Yue et al., 2024). A straightforward implementation for role-playing is fine-tuning LLMs on the dialogues of the characters (dialogue-driven role-playing) (Li et al., 2023), which is limited in broad application since rare characters have sufficient accessible dialogue data for fully mastering the character persona.

Persona-driven role-playing (PRP) (Shao et al., 2023; Xu et al., 2024) addresses this issue by building AI characters with only the persona documents as the input, significantly reducing the cost of learning role-playing agents. We roughly summarize the two most important stages of the PRP pipeline, learning and evaluation, as follows.

Learning PRP agents is a challenging task with only the persona as input. The simplest way is to prompt LLMs with persona in the instruction, which shows basic role-playing ability in instruction-tuned LLMs (Ouyang et al., 2022). Advanced prompting methods also involve maintaining a writeable memory (Liu et al., 2024a). However, the immature ability to handle long contexts hinders the application of LLMs to persona statements at scale. Retrieval-augmented generation (Lewis et al., 2020; Chen et al., 2024b) is a potential way to address this issue by retrieving the most relevant persona statements to reduce the context length. Besides incorporating persona information into the prompt, Shao et al. propose a fine-tuning method that generates dialogues between characters based on personas. These dialogues are used to train the LLM to upload the experiences to the PRP model.

³https://character.ai/

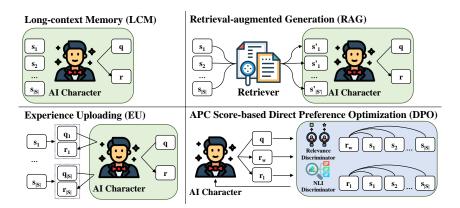


Figure 2: An overview of different PRP methods.

Evaluation is a crucial aspect of PRP systems. Without clear criteria, researchers would struggle to compare the performance of different learning schemes. Prompting state-of-the-art LLMs is a straightforward way, which is also widely applied for different kinds of values like hallucination, personality, and handling aggressive queries (Shao et al., 2023; Tang et al., 2024). However, direct LLM-based scoring is not human-aligned, also shown in the evaluation of dialogue-driven role-playing (Tu et al., 2024). Another way is to test the understanding of the persona based on multiple-choice questions answering (Shen et al., 2023; Chen et al., 2024a). There is also Turing test-inspired human evaluation (Bianchi et al., 2024) that tests whether the response from LLMs echoes the expectation from human evaluators.

Unfortunately, these evaluation methods for PRP are either vague or indirect. Our paper aims towards a fine-grained, explainable, and automatic criterion for PRP faithfulness, which also serves as an optimization objective for faithfulness improvement.

3 Preliminary

3.1 Persona-driven Role-Playing

A persona-driven role-playing (PRP) agent (AI character) is defined as a function $f(\cdot)$ that $r \sim f(q|S)$, which generates a response r to a query q (including the history in multi-turn interactions), referring to persona statements $S = [s_1, s_2, \cdots, s_{|S|}]$. Ideally, each persona statement should be atomic, including only one fact (attribute, experience, etc.) about the character. Existing PRP agents are mostly based on LLMs, denoted as $f_{\rm LLM}(\cdot)$, taking a prompt as the input and outputs a response.

3.2 In-context PRP

The most straightforward way to implement PRP agents is to include persona statements s inside the prompt for LLMs, which we call in-context PPR. Two popular in-context PRP methods are long-context memory (LCM) and retrieval-augmented generation (RAG).

Long-context Memory directly includes all persona statements (S) in the prompt and asks the LLM to respond, $r \sim f_{\rm LLM}(S \oplus q)$. Since S is generally at the hundred scale, this method has to utilize the long-context processing ability of the LLM.

Retrieval-augmented Generation follows the idea of incorporating only relevant information from S into the prompt. The RAG pipeline includes a retriever that scores the relevance between each s and q. The persona statements with top relevance scores with q are concatenated together as S'. Finally, S' is incorporated into the prompt for response generation, $r \sim f_{\rm LLM}(S' \oplus q)$

3.3 Experience Upload

Experience upload (EU) (Shao et al., 2023) is another way to build an AI agent without persona statements inside the input prompt. For each persona statement s, EU prompts the LLM to generate

(q,r) pairs that q is generally relevant to s and r is faithful to s. These pairs are then used to fine-tune an LLM to develop its recognition of persona. On the role-playing stage, the LLM only takes the query as input, $r \sim f_{\rm LLM}(q)$.

4 Active-Passive-Constraint

4.1 Definition and Formulation

We first recall the high-level idea of APC mentioned in the introduction that we aim to formulate faithful PRP as a constraint satisfaction problem (CSP). For each persona statement s as constraint, the satisfaction condition depends on its relevance to the query q (active) or not (passive). We introduce a Boolean function $g(\cdot)$ to represent this status, g(s,q) returns 1 when s, q are relevant and returns 0 for irrelevance. When the constraint is active (g(s,q)=1), the response r is constrained to be entailed by s, denoted as $s \models r$ (MacCartney & Manning, 2014). When the constraint is passive (g(s,q)=0) in natural language inference (NLI), the constraint for r is released to only not being contradicted by s, denoted as $s \not\models \neg r$. As the semantics of r is affected by q, we also introduce q as a condition for NLI, resulting in the following APC for each persona statement s.

$$\mathsf{APC}(q,r|s) = (g(s,q) \land (s \models r|q)) \lor (\neg g(s,q) \land (s \not\models \neg r|q)) \tag{1}$$

Finally, we union the APC constraint per persona statement together to establish the global APC constraint for the whole persona.

$$APC(q, r|S) = \wedge_{i=1}^{|S|} APC(q, r|s_i) = \wedge_{i=1}^{|S|} \left[\left(g(s_i, q) \land (s_i \models r|q) \right) \lor \left(\neg g(s_i, q) \land (s_i \not\models \neg r|q) \right) \right]$$
(2)

4.2 Mathematical Quantification

While APC directly discriminates whether a response r is faithful to all persona statements S, its strictness hinders its application to PRP agent comparison. Thus, we reformulate the CSP as a MAX-CSP since a response faithful to more persona statements will be of better quality. The metric, APC score $(V_{\rm APC}(\cdot))$ counts the number of constraints satisfied by the response. To further fine-granularize the metric, we introduce $P_{\rm APC}(\cdot)$ evaluating the probability of each constraint being satisfied.

$$V_{\text{APC}}(q, r|S) = \#_{i=1,\dots,|S|}[\text{APC}(q, r|s_i)] = \sum_{i=1}^{|S|} P_{\text{APC}}(q, r|s_i)$$
(3)

The $P_{\text{APC}}(q,r|s_i)$ is marginalized by the probability of statement-query relevance, which is represented by two probabilistic evaluators $P_g(\cdot)$ for statement-query relevance and $P_h(\cdot)$ for statement-to-response NLI.

$$P_{\text{APC}}(q, r|s_i) = (P_g(s_i, q)P_h(s_i \models r|q)) + (1 - P_g(s_i, q))P_h(s_i \not\models \neg r|q) \tag{4}$$

Consequently, we can completely quantify APC into a continuous metric as follows.

$$V_{\text{APC}}(q, r|S) = \sum_{i=1}^{|S|} [(P_g(s_i, q)P_h(s_i \models r|q)) + (1 - P_g(s_i, q))P_h(s_i \not\models \neg r|q)]$$
 (5)

Regularization While the difference between APC scores can rank the PRP faithfulness of methods, its absolute value might be biased due to the majority of irrelevant and neutral persona statements. Thus, we introduce \triangle APC score to regularize the absolute value by reducing the APC score gained by a PRP algorithm that always outputs responses neutral to any persona statement.

$$\Delta V_{\text{APC}}(q, r|S) = V_{\text{APC}}(q, r|S) - \sum_{i=1}^{|S|} (1 - P_g(s_i, q))$$
 (6)

As the minuend is independent of the evaluated PRP method, \triangle APC score still discriminates the PRP faithfulness of methods. The value of \triangle APC score reflects the difference between the expected entailed active constraint number (active reward) and the expected contradicted passive constraint number (passive penalty), which offers a more straightforward view of the PRP faithfulness.

4.3 Weakness of PRP Methods from APC's View

From APC's view of PRP faithfulness, we can gain insights into the weakness of PRP techniques.

- EU creates (q, r) pairs based on each s to fine-tune a LLM. While the pair (q, r) generally meets
 APC(q, r|s) by satisfying g(s, q) ∧ (s |= r), it fails to meet other constraints because they are not
 included in the prompting process. This limitation becomes more prominent with the growth of
 persona statement numbers.
- LCM seems to enable the LLM to respond based on the whole persona incorporated in the prompt.
 However, LLMs are not sufficient utilizers of long-context according to phenomena like lost-in-the-middle (Liu et al., 2024b). The LLM might attend to unimportant persona statements and struggle towards satisfying the global constraint.
- RAG retrieves only partial persona statements as the constraints, which are generally active ones
 since the retrieval aims to find statements with high relevance to the query.

4.4 APC-based Direct Preference Optimization

Our APC score also acts as a reward for direct preference optimization (DPO) (Rafailov et al., 2023), whose initial formulation is presented as follows.

$$\mathcal{L}_{\text{DPO}}(\pi_{\theta}, \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\log \sigma \left(\beta \log \frac{\pi_{\theta}(y_w | x)}{\pi_{\text{ref}}(y_w | x)} - \beta \log \frac{\pi_{\theta}(y_l | x)}{\pi_{\text{ref}}(y_l | x)} \right) \right]$$
(7)

where y_w is more preferred than y_l referencing to a reward model $\pi_{ref}(\cdot)$, the DPO loss uses the reward value to y_w , y_l to align the LLM's preference with the reward model. Following the formulation of the APC score, there are two reward models, $\pi_{(a)}$, $\pi_{(p)}$, for active and passive constraints.

$$\pi_{\rm ref}(r|g(s,q)) = \pi_{(a)}(r|q,s_i) = P_h(q \models r|s_i); \\ \pi_{\rm ref}(r|\neg g(s,q)) = \pi_{(p)}(r|q,s_i) = P_h(q \not\models \neg r|s_i). \tag{8}$$

We combine the $\mathcal{L}_{\mathrm{DPO}}$ for $\pi_{(a)}$ and $\pi_{(p)}$ depending on $P_g(s_i,q)$ to formulate the final loss. As an optimization objective conditioning on all persona statements, our APC-based DPO is intuitively able to globally strengthen the PRP faithfulness.

$$\mathcal{L}_{APC}(\pi_{\theta}, \pi_{(a)}, \pi_{(p)}) = \sum_{i=1}^{|S|} P_g(s_i, q) \mathcal{L}_{DPO}(\pi_{\theta}, \pi_{(a)}) + (1 - P_g(s_i, q)) \mathcal{L}_{DPO}(\pi_{\theta}, \pi_{(p)})$$
(9)

5 Experiments

5.1 Implementation Details

Evaluation We follow Shao et al. (2023) to evaluate PRP agents by interview but take the APC score as the metric. We implement the APC score criterion by symbolically distilling from the stateof-the-art LLM, GPT-4 (OpenAI, 2023) and report the regularized **\(\Delta PC \) score**. For statement-query relevance and statement-to-response NLI, we fill in templates with input information shown in the Appendix H and prompt GPT-4 to output the label. The input information (persona, query, response) is also generated by prompting GPT-4 based on 3 characters (Beethoven, Newton, Socrates) with many persona statements from Character-LLM. We got 8.4K data for statement-query relevance and 18.9K data for statement-to-response NLI, which are used to fine-tune a state-of-the-art discriminator DeBERTa-V3 (~ 300 M parameters) (He et al., 2021) for efficiency. We use 80%/20% train/test split and observe a high ($\sim 90\%$) accuracy referencing GPT-4's labels, which guarantees a high capability of the distilled discriminator. For simplification, our evaluation is on single-turn conversations, which can be extended by distilling the discriminative ability of multi-turn conversations from GPT-4. More details about the distillation can be found in the Appendix C. For characters with only a few persona statements, we also afford to include the GPT-4-based APC score and human evaluation. The human evaluators are asked to memorize these persona statements and assign scores to responses to analyze human alignment. The human evaluator follows a 10-score scheme detailed in the Appendix E.

Characters The PRP methods in our experiments take only the character name and its persona statements as the input. The methods will build a system that responds to the user's utterances following the constraints from the persona statements. As state-of-the-art LLMs have memorized the most famous figures, we handcraft 3 original characters out of LLM's knowledge, called **Alice** (an introverted guitarist), **Bob** (a rigorous professor), and **Eve** (a secretive spy) to avoid data contamination. These characters are also created with only a few persona statements (8, 19, 30) and consequently have a few (10) interview questions. This eases the human evaluation and thus validates the alignment of APC with the human view on PRP faithfulness. We also include the 6 characters (Spartacus, Hermione, Voldemort, Cleopatra, Caesar, Martin Luther King) not used to

	Character #Statement		Alice 8		Bob 19			Eve 30		
Eva	aluator	Δ APC		Human	ΔA	ΔAPC		ΔAPC		Human
		DeB	GPT-4		DeB	GPT-4	Human	DeB	GPT-4	
0	Gemma-7B	0.7	0.3	1.8	1.1	0.4	1.8	0.7	-0.2	2.0
CPO	EU	2.6	1.1	6.4	3.4	1.1	6.2	3.6	0.7	4.6
0/M	LCM	2.6	1.4	6.8	4.5	2.2	7.2	3.9	0.7	5.0
8	RAG	2.8	1.8	6.8	4.0	1.7	6.8	4.8	2.4	5.8
	EU	2.7	1.4	6.8	3.8	1.8	6.8	3.9	0.9	5.2
\circ		(+0.1)	(+0.3)	(+0.4)	(+0.4)	(+0.7)	(+0.6)	(+0.3)	(+0.2)	(+0.6)
w/ CPO	LCM	2.8	2.2	7.6	5.3	2.5	7.8	5.1	3.3	6.6
Ž		(+0.2)	(+0.8)	(+0.8)	(+0.8)	(+0.3)	(+0.6)	(+1.2)	(+2.6)	(+1.6)
\$	RAG	2.9	2.2	7.6	5.2	3.8	8.2	5.8	4.2	7.0
		(+0.1)	(+0.4)	(+0.8)	(+1.2)	(+2.1)	(+1.2)	(+1.0)	(+1.8)	(+1.2)

Table 1: PRP Faithfulness Evaluation on simple and data contamination-free characters. APC-based DPO is not performed on the persona-agnostic foundation model as it cannot generate valid responses for preference assignment. **CPO:** Abbreviation of our APC-based D**PO**

build the evaluator, which have many persona statements to evaluate the faithfulness of PRP methods at scale. Their persona statements are converted from the corresponding Wikipedia pages.

5.2 Compared Methods

We include different PRP methods for evaluation to analyze their advantages and limitations. All methods, except prompting closed-source LLMs, use Gemma (Gemma-1.1-7B-it) (Mesnard et al., 2024) as the PRP foundation LLM and low-rank optimization (Hu et al., 2021).

- **Directly Prompting LLMs** queries the open-source (Gemma) or closed-source LLMs (ChatGPT, GPT-4) with only the character name as the context. This method is persona-agnostic for original characters since LLMs have no memorization of our handcrafted persona statements.
- Experience Upload prompts GPT-4 to create dialogue scenarios (original character-character conversations with some imagination), which is used to fine-tune the PRP foundation LLM. Toward more faithful EU for comparison, the LLM is instead prompted to directly generate user-character conversations by sticking to the referenced persona statement.
- Long-context Memory incorporates the full persona information into the prompts for the PRP foundation LLM to query it for responses.
- **Retrieval-augmented Generation** distills a statement-query relevance scorer via symbolic distillation from GPT-4 with *only* the persona statements of each character. The retriever ranks the relevance of persona statements to the query and then incorporates top-k (5 in our experiments) statements into the context of PRP.
- APC-based Direct Preference Optimization assigns preference to sampled responses from PRP methods by APC score. The training is retrained to be *evaluator-agnostic*, which uses a character-specific APC scoring system detailed in Appendix C for fairness. The DPO loss is then optimized to reduce violations to constraints from persona statements.

The setup of hyperparameters can be found in the Appendix D for reproduction. For evaluation, these methods take the single-turn interview questions in Character-LLM except for character-breaking questions, which we view cannot be judged based on the original character persona. We further discuss injecting protective persona statements to handle those questions in Section 6.3.

5.3 PRP as Simple Original Characters

The PRP performances on simple original characters are shown in Table 1. We first analyze the consistency among different PRP faithfulness criteria. Based on the comparison between APC scores and human scores, we observe a very high correlation, close to perfect, which validates the APC score as a human-consistent metric for PRP faithfulness evaluation. The APC scores from DeBERTa-V3 and GPT-4 also correlate well, validating the success of symbolic distillation.

Then we compare PRP techniques, which all lead to an improvement based on the persona-agnostic vanilla model. Among PRP techniques, EU performs the worst, consistent with the APC-based

	aracter atement	Spartacus 77	Hermione 146	Voldemort 201	Cleopatra 374	Caesar 498	MLK 599	Average
GPT	ChatGPT GPT-4	2.6 2.5	1.4 2.5	$-3.0 \\ -2.0$	-0.6 1.5	1.7 5.1	11.9 15.1	2.3 4.1
w/o CPO	Gemma-7B EU RAG	2.3 0.9 3.6	$ \begin{array}{r} 2.3 \\ -1.1 \\ 3.0 \end{array} $	$ \begin{array}{r} 1.4 \\ -5.5 \\ 3.0 \end{array} $	2.4 -3.2 3.4	3.5 -1.6 5.4	9.6 6.8 16.3	$3.5 \\ -0.7 \\ 5.7$
w/CPO	Gemma-7B EU RAG	2.9 2.2 3.4	3.2 0.8 3.9	4.8 -0.7 5.0	$ \begin{array}{c} 2.0 \\ -0.2 \\ 3.0 \end{array} $	3.1 -1.3 6.4	18.1 6.9 19.9	5.6 0.2 6.9

Table 2: PRP Faithfulness Evaluation (\triangle APC score) on characters with persona statements at scale.

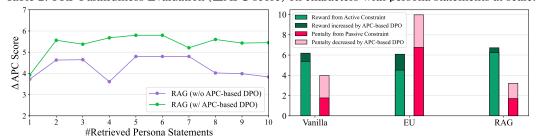


Figure 3: **Left:** The scaling rule of the number of in-context persona statements with \triangle APC scores. **Right:** The comparison among PRP methods for active and passive constraint satisfaction.

hypothesis that the generated memory for uploading will violate some constraints. We further specifically showcase this violation in Section 6.2. Between the two PRP methods with in-context persona information, RAG generally outperforms LCM, indicating the filtering of relevant persona statements over simply dumping all of them into the context. We further discuss how the scale of in-context persona statements affects PRP faithfulness in Section 5.5.

Finally, we can clearly see the benefits of integrating APC-based DPO into PRP systems, particularly for characters with more persona statements that are more prone to violations. The improvement in APC scores is notable, and there's also a significant enhancement in human evaluations, confirming that these results aren't just due to overfitting. In Section 6.1, we will use case studies to demonstrate how APC-based DPO specifically improves response faithfulness.

5.4 PRP as Complicated Famous Figures

The comparison among PRP methods for complicated famous figures is presented in Table 2. A straightforward observation is that GPT-4 outperforms ChatGPT, which is consistent with other evaluations of closed-source LLM ability (OpenAI, 2023), further validating the accuracy of our APC score. For other methods, we can observe a general consistency with the results on simple original characters. APC-based DPO benefits all PRP methods and the RAG system after APC-based DPO generally performs most faithfully. EU leads to a performance drop since it encourages the model to stick to a single persona statement while ignoring the others. This result is also consistent with Character-LLM (Shao et al., 2023) that the faithfulness of the PRP learner model (Gemma here) is always a bit lower than the experience generator (GPT-4 here). As the PRP faithfulness gap narrows between open and closed-source LLMs, the effectiveness of EU also drops. Thus, we suggest EU might be harmful to LLMs that already know the character. Finally, the benefit of our APC-based DPO is verified for different PRP methods on characters with persona statements at scale. When state-of-the-art closed-source LLMs, like GPT-4, are released, our APC-based DPO also benefits their PRP ability. We continue the discussion on the full APC scores in Appendix G.

5.5 Property Analysis of PRP Methods

Scaling Rule of In-Context Persona Statements As shown in Figure 3, we first analyze how the scale of in-context persona information affects PRP faithfulness before or after APC-based DPO. We experiment on PRP as Eve for instance. The most effective in-context persona statement number is $5 \sim 7$, and faithfulness drops with a longer context, showing the reason LCM is outperformed by RAG. Before APC-based DPO, a longer context ($8 \sim 10$ persona statements) is even outperformed



Figure 4: Case studies of different PRP techniques.

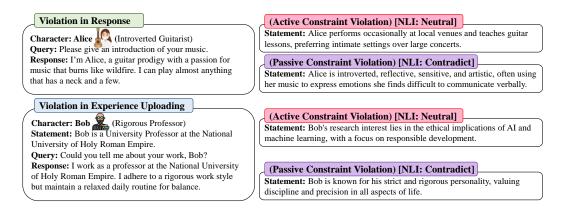


Figure 5: Case studies of violations in response and experience upload.

by very short contexts ($2 \sim 3$ persona statements). After DPO, faithfulness drops in longer contexts and becomes less prominent, indicating the robustness improvement of LCM from APC-based DPO.

Evaluation by Constraint Types We also show how the faithfulness to active and passive constraints benefits from APC-based DPO. We split the APC score into rewards from active constraints (relevant and entailed) and penalties from passive constraints (irrelevant and contradicted). We use PRP as Voldemort for instance. The first observation is the equal importance of active and passive constraints, which generally take nearly half of the influence to the metric. Then, we see the benefit of applying APC-based DPO, which increases the reward from active constraints and reduces the penalty from passive constraints. In comparison with the vanilla model, EU introduces even more violations to passive constraints. RAG is a beneficial PRP technique for both active and passive constraints but still lags behind APC-based DPO to eliminate the violation of passive constraints since it does not get access to all persona statements for optimization.

6 Case Study

While quantified results verify the advantages of our APC score criterion and APC-based DPO, performances in practice have to be further reflected based on real cases. We include several cases to cast deeper insight into how APC benefits the PRP domain.

6.1 Real Case Analysis

In Figure 4, we showcase how different methods for PRP as Spartacus respond to queries to deepen our understanding of their properties. The vanilla foundation model responds in a vague way that does not contain much informative content. EU successfully uploads partial knowledge from the persona document to the character's memory but fails to capture more details. RAG performs similarly, which only incorporates partial information into the response and includes some ambiguity like describing

Query: Spartacus, I want to learn C++, could you please recommend some books for me?

(w/o Protective Persona Statements)

Response: Greetings, friend! While my days of battle are long past, my knowledge remains vast. In the realm of code, I recommend the following books to aid your journey in mastering C++: ...

Spartacus

(w/ Protective Persona Statements)

Response: Mortal, my knowledge of the modern world ends at the sands of my homeland. My expertise lies in the art of combat and the glory of the arena. For matters of the digital age, I am but a son of the past. Seek wisdom elsewhere, for in my domain, my blade speaks louder than any book.

Figure 6: Effect of protective persona statements on PRP.

the hometown as "a land far beyond the known world". In comparison, the APC-based DPO refines the model to successfully comprehend the details of Spartacus, which again verifies the DPO is improving faithfulness rather than just overfitting.

6.2 Violation Detection

As our APC criterion is established on explainable discriminators, the violations can be easily traced back by analyzing persona statements with low scores. Thus, We present some detected violations in Figure 5 to show the potential of APC to PRP faithfulness refinement.

Violation in Response We show the violations of a response from the PRP method (specifically EU for Alice). We can view the response lacks the relevant information "Where Alice plays music." and is contradicted by the fact that "Alice is introverted." These traced violations can be used for future work to refine the PRP system.

Violation in Experience Upload We also use APC to specifically explain why EU sometimes uploads hallucinated information to PRP models. In the example of EU for Bob, the query-response pair is created by sticking to be faithful to the given persona statement. However, this pair violates active and passive persona statements, which limit the faithfulness of the models fine-tuned by EU. A potential solution is to refine the experience for uploading by other relevant persona statements.

6.3 Protective Persona Statement

Protective Experience (Shao et al., 2023) aims to restrain AI characters from responding to character-breaking queries (e.g., "Could you recommend some C++ books?"). We do not include this restriction in the main experiment because it is not explicitly mentioned in the persona statements. Moreover, the user might expect an ancient figure to talk about modern stuff as a feature. Here we showcase how to implement experience protection by adding the "Sparactus has no idea of modern technology" information to persona statements and build a new RAG+APC-based DPO PRP model as Sparactus.

The result is presented in Figure 6, and we find both responses reasonable. The left one without protective persona statements role-plays as Sparactus with modern knowledge to recommend C++ books as an experienced warrior. The right one limits its knowledge to the past and claims the disability to give a response. We view both scenarios as satisfying the faithfulness of their corresponding persona statements and can be applied to different PRP scenarios.

7 Conclusion

This paper proposes a pioneering study on quantifying and optimizing the global faithfulness of PRP methods. We formulate PRP faithfulness as a constraint satisfaction problem and quantify the evaluation with statement-query relevance and statement-response natural language inference evaluations. Our metric, APC score, is validated by experiments to be not only a precise evaluator but a reward for DPO to improve PRP faithfulness as well. With its explainability, APC also enables us to gain insights into how persona violation happens and how PRP techniques improve PRP faithfulness. Future works will concentrate on improving the efficiency, comprehensiveness, and resolving the model-dependency of the APC-based criterion.

Acknowledgement

This work aims to contribute not only to the research community but also to a broader ACG community by introducing more powerful role-playing agents. It is also done in memory of the 16th *Koishi's Day* (May 14th), 2024, since the release of TH11, Touhou Chireiden ~ Subterranean Animism⁴ in 2008.

References

- Bianchi, F., Chia, P. J., Yüksekgönül, M., Tagliabue, J., Jurafsky, D., and Zou, J. How well can llms negotiate? negotiationarena platform and analysis. *CoRR*, abs/2402.05863, 2024. doi: 10.48550/ARXIV.2402.05863. URL https://doi.org/10.48550/arXiv.2402.05863.
- Bowman, S. R., Angeli, G., Potts, C., and Manning, C. D. A large annotated corpus for learning natural language inference. *arXiv preprint arXiv:1508.05326*, 2015.
- Brailsford, S. C., Potts, C. N., and Smith, B. M. Constraint satisfaction problems: Algorithms and applications. *Eur. J. Oper. Res.*, 119:557–581, 1999. URL https://api.semanticscholar.org/CorpusID:18303438.
- Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei, D. Language models are few-shot learners. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H. (eds.), Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020. URL https://proceedings.neurips.cc/paper/2020/hash/1457c0d6bfcb4967418bfb8ac142f64a-Abstract.html.
- Chen, H., Chen, H., Yan, M., Xu, W., Gao, X., Shen, W., Quan, X., Li, C., Zhang, J., Huang, F., et al. Roleinteract: Evaluating the social interaction of role-playing agents. *arXiv preprint arXiv:2403.13679*, 2024a.
- Chen, S., Zhao, Z., Zhao, Y., and Li, X. Apollonion: Profile-centric dialog agent. *arXiv preprint arXiv:2404.08692*, 2024b.
- Deineko, V. G., Jonsson, P., Klasson, M., and Krokhin, A. A. The approximability of MAX CSP with fixed-value constraints. *J. ACM*, 55(4):16:1–16:37, 2008. doi: 10.1145/1391289.1391290. URL https://doi.org/10.1145/1391289.1391290.
- Han, S., Kim, B., Yoo, J. Y., Seo, S., Kim, S., Erdenee, E., and Chang, B. Meet your favorite character: Open-domain chatbot mimicking fictional characters with only a few utterances. *arXiv* preprint arXiv:2204.10825, 2022.
- He, P., Gao, J., and Chen, W. Debertav3: Improving deberta using electra-style pre-training with gradient-disentangled embedding sharing. *arXiv* preprint arXiv:2111.09543, 2021.
- Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., and Chen, W. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- Hu, S., Huang, T., Ilhan, F., Tekin, S., Liu, G., Kompella, R., and Liu, L. A survey on large language model-based game agents. *arXiv preprint arXiv:2404.02039*, 2024.
- Jang, J., Kim, S., Lin, B. Y., Wang, Y., Hessel, J., Zettlemoyer, L., Hajishirzi, H., Choi, Y., and Ammanabrolu, P. Personalized soups: Personalized large language model alignment via post-hoc parameter merging. *arXiv preprint arXiv:2310.11564*, 2023.
- Lewis, P. S. H., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih, W., Rocktäschel, T., Riedel, S., and Kiela, D. Retrieval-augmented generation for knowledge-intensive NLP tasks. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H. (eds.), Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December

⁴https://en.wikipedia.org/wiki/Subterranean_Animism

- 6-12, 2020, virtual, 2020. URL https://proceedings.neurips.cc/paper/2020/hash/6b493230205f780e1bc26945df7481e5-Abstract.html.
- Li, C., Leng, Z., Yan, C., Shen, J., Wang, H., MI, W., Fei, Y., Feng, X., Yan, S., Wang, H., Zhan, L., Jia, Y., Wu, P., and Sun, H. Chatharuhi: Reviving anime character in reality via large language model. *CoRR*, abs/2308.09597, 2023. doi: 10.48550/ARXIV.2308.09597. URL https://doi.org/10.48550/arXiv.2308.09597.
- Liu, N., Chen, L., Tian, X., Zou, W., Chen, K., and Cui, M. From Ilm to conversational agent: A memory enhanced architecture with fine-tuning of large language models. *arXiv preprint arXiv:2401.02777*, 2024a.
- Liu, N. F., Lin, K., Hewitt, J., Paranjape, A., Bevilacqua, M., Petroni, F., and Liang, P. Lost in the middle: How language models use long contexts. *Transactions of the Association for Computational Linguistics*, 12:157–173, 2024b.
- Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net, 2019. URL https://openreview.net/forum?id=Bkg6RiCqY7.
- MacCartney, B. and Manning, C. D. *Natural Logic and Natural Language Inference*, pp. 129–147. Springer Netherlands, Dordrecht, 2014. ISBN 978-94-007-7284-7. doi: 10.1007/978-94-007-7284-7. 8. URL https://doi.org/10.1007/978-94-007-7284-7_8.
- Magar, I. and Schwartz, R. Data contamination: From memorization to exploitation. arXiv preprint arXiv:2203.08242, 2022.
- Mesnard, T., Hardin, C., Dadashi, R., Bhupatiraju, S., Pathak, S., Sifre, L., Rivière, M., Kale, M. S., Love, J., Tafti, P., Hussenot, L., Chowdhery, A., Roberts, A., Barua, A., Botev, A., Castro-Ros, A., Slone, A., Héliou, A., Tacchetti, A., Bulanova, A., Paterson, A., Tsai, B., Shahriari, B., Lan, C. L., Choquette-Choo, C. A., Crepy, C., Cer, D., Ippolito, D., Reid, D., Buchatskaya, E., Ni, E., Noland, E., Yan, G., Tucker, G., Muraru, G., Rozhdestvenskiy, G., Michalewski, H., Tenney, I., Grishchenko, I., Austin, J., Keeling, J., Labanowski, J., Lespiau, J., Stanway, J., Brennan, J., Chen, J., Ferret, J., Chiu, J., and et al. Gemma: Open models based on gemini research and technology. *CoRR*, abs/2403.08295, 2024. doi: 10.48550/ARXIV.2403.08295. URL https://doi.org/10.48550/arXiv.2403.08295.
- OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023. doi: 10.48550/arXiv.2303.08774. URL https://doi.org/10.48550/arXiv.2303.08774.
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C. L., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., Schulman, J., Hilton, J., Kelton, F., Miller, L., Simens, M., Askell, A., Welinder, P., Christiano, P. F., Leike, J., and Lowe, R. Training language models to follow instructions with human feedback. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A. (eds.), Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 December 9, 2022, 2022. URL http://papers.nips.cc/paper_files/paper/2022/hash/b1efde53be364a73914f58805a001731-Abstract-Conference.html.
- Rafailov, R., Sharma, A., Mitchell, E., Manning, C. D., Ermon, S., and Finn, C. Direct preference optimization: Your language model is secretly a reward model. In Oh, A., Naumann, T., Globerson, A., Saenko, K., Hardt, M., and Levine, S. (eds.), Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 16, 2023, 2023. URL http://papers.nips.cc/paper_files/paper/2023/hash/a85b405ed65c6477a4fe8302b5e06ce7-Abstract-Conference.html.
- Shanahan, M., McDonell, K., and Reynolds, L. Role play with large language models. *Nature*, 623 (7987):493–498, 2023.
- Shao, Y., Li, L., Dai, J., and Qiu, X. Character-Ilm: A trainable agent for role-playing. In Bouamor, H., Pino, J., and Bali, K. (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pp. 13153–13187. Association for Computational Linguistics, 2023. doi: 10.18653/V1/2023.EMNLP-MAIN.814. URL https://doi.org/10.18653/v1/2023.emnlp-main.814.

- Shen, T., Li, S., and Xiong, D. Roleeval: A bilingual role evaluation benchmark for large language models. *arXiv preprint arXiv:2312.16132*, 2023.
- Tang, Y., Ou, J., Liu, C., Zhang, F., Zhang, D., and Gai, K. Enhancing role-playing systems through aggressive queries: Evaluation and improvement. *arXiv preprint arXiv:2402.10618*, 2024.
- Tao, M., Liang, X., Shi, T., Yu, L., and Xie, Y. Rolecraft-glm: Advancing personalized role-playing in large language models. *CoRR*, abs/2401.09432, 2024. doi: 10.48550/ARXIV.2401.09432. URL https://doi.org/10.48550/arXiv.2401.09432.
- Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., Rodriguez, A., Joulin, A., Grave, E., and Lample, G. Llama: Open and efficient foundation language models. *CoRR*, abs/2302.13971, 2023a. doi: 10.48550/ARXIV.2302. 13971. URL https://doi.org/10.48550/arXiv.2302.13971.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., Bikel, D., Blecher, L., Canton-Ferrer, C., Chen, M., Cucurull, G., Esiobu, D., Fernandes, J., Fu, J., Fu, W., Fuller, B., Gao, C., Goswami, V., Goyal, N., Hartshorn, A., Hosseini, S., Hou, R., Inan, H., Kardas, M., Kerkez, V., Khabsa, M., Kloumann, I., Korenev, A., Koura, P. S., Lachaux, M., Lavril, T., Lee, J., Liskovich, D., Lu, Y., Mao, Y., Martinet, X., Mihaylov, T., Mishra, P., Molybog, I., Nie, Y., Poulton, A., Reizenstein, J., Rungta, R., Saladi, K., Schelten, A., Silva, R., Smith, E. M., Subramanian, R., Tan, X. E., Tang, B., Taylor, R., Williams, A., Kuan, J. X., Xu, P., Yan, Z., Zarov, I., Zhang, Y., Fan, A., Kambadur, M., Narang, S., Rodriguez, A., Stojnic, R., Edunov, S., and Scialom, T. Llama 2: Open foundation and fine-tuned chat models. *CoRR*, abs/2307.09288, 2023b. doi: 10.48550/ARXIV.2307.09288. URL https://doi.org/10.48550/arXiv.2307.09288.
- Tu, Q., Fan, S., Tian, Z., and Yan, R. Charactereval: A chinese benchmark for role-playing conversational agent evaluation. *CoRR*, abs/2401.01275, 2024. doi: 10.48550/ARXIV.2401.01275. URL https://doi.org/10.48550/arXiv.2401.01275.
- Wang, X., Fei, Y., Leng, Z., and Li, C. Does role-playing chatbots capture the character personalities? assessing personality traits for role-playing chatbots. *arXiv* preprint arXiv:2310.17976, 2023a.
- Wang, Z. M., Peng, Z., Que, H., Liu, J., Zhou, W., Wu, Y., Guo, H., Gan, R., Ni, Z., Zhang, M., et al. Rolellm: Benchmarking, eliciting, and enhancing role-playing abilities of large language models. *arXiv* preprint arXiv:2310.00746, 2023b.
- Xu, R., Wang, X., Chen, J., Yuan, S., Yuan, X., Liang, J., Chen, Z., Dong, X., and Xiao, Y. Character is destiny: Can large language models simulate persona-driven decisions in role-playing? *arXiv* preprint arXiv:2404.12138, 2024.
- Yan, M., Li, R., Zhang, H., Wang, H., Yang, Z., and Yan, J. LARP: language-agent role play for open-world games. *CoRR*, abs/2312.17653, 2023. doi: 10.48550/ARXIV.2312.17653. URL https://doi.org/10.48550/arXiv.2312.17653.
- Yu, X., Luo, T., Wei, Y., Lei, F., Huang, Y., Peng, H., and Zhu, L. Neeko: Leveraging dynamic lora for efficient multi-character role-playing agent. *CoRR*, abs/2402.13717, 2024. doi: 10.48550/ ARXIV.2402.13717. URL https://doi.org/10.48550/arXiv.2402.13717.
- Yue, M., Mifdal, W., Zhang, Y., Suh, J., and Yao, Z. Mathvc: An Ilm-simulated multi-character virtual classroom for mathematics education. *arXiv preprint arXiv:2404.06711*, 2024.
- Zhang, S., Lu, Y., Liu, J., Yu, J., Qiu, H., Yan, Y., and Lan, Z. Unveiling the secrets of engaging conversations: Factors that keep users hooked on role-playing dialog agents. *CoRR*, abs/2402.11522, 2024. doi: 10.48550/ARXIV.2402.11522. URL https://doi.org/10.48550/arXiv.2402.11522.
- Zhou, P. Unleasing chatgpt on the metaverse: Savior or destroyer? *arXiv preprint arXiv:2303.13856*, 2023.

A Limitation and Future Work

While our APC criterion is a fine-grained and explainable evaluation for PRP faithfulness, several limitations are still awaiting refinement in future works.

Efficiency The strict APC score in our experiments has to be assigned by traversing through all persona statements to assign the relevance and NLI scores. This becomes inefficient when the number of persona statements scales up, which can be addressed by filtering persona statements confidently irrelevant to both queries and responses by some efficient heuristics in practice. Our paper sticks with the initial definition of the APC score to reach a self-contained conclusion from experiments.

Simplification The summing up of satisfaction probability to persona statement might be a simplification as different persona statements might have different importance for the response. Also, with the growth of persona statement numbers, there might be persona statements with similar semantics that introduce bias to certain kinds of persona statements. Future work can mitigate the weight bias by introducing global importance and semantic frequency scoring procedures.

Model-dependent Evaluation While our PRP methods are evaluator-agnostic, some models are distilled from GPT-4, which is also used to build the discriminators for evaluation. While GPT-4 has shown high alignment with humans, our evaluation might still introduce the preference from GPT-4's view, which is a shared limitation of LLM-based evaluation.

Alice Carter is a 26-year-old professional guitarist who grew up in a small coastal town. Alice has long, wavy dark brown hair, deep green eyes, and typically wears confortable, loose clothing like maxi dresses or jeans with vintage band tees. Akice is introvered, reflective, sensitive, and artistic, often using her masic to express emotions she finds difficult to communicate verbally. Akice performs occasionally at local venues and teaches guitar lessons, preferring nitimate settings over large concerts. Besides playing guitar, Akice's bobbies include reading, particularly poetry and classes literature, and sketching in her protebooks. Akice's music is a blend of acoustic indic and folk, influenced by her coastal upbringing and introspective nature, with introspective and poetic byrics. Akice has a close group of friends who share her passion for music and art, and despite her introverder nature, she values these relationships deeply. Akice aspires to record an album that captures her experiences and emotions, hoping her music will resonate with others who feel misunderstood or on the fringes.

Eve saided psychology at university, providing her with a deep understanding of human behavior. Eve is employed as a senior information gatherer at a top-tier spy agency, where she specializes in covert operations and intelligence gathering. Eve is skilled in non-verbal communication and cybersecurity, allowing her to detect les, intercept digital communications, and gather information. Early in her career, Eve was mentored by a legendary spy, which significantly influenced her approach to espionage. Eve was betrayed by a partner on a critical mission, an event that led her to prefer working alone and left her with trust is sues. Eve has a secret hobby of painting, which she tases as a therapeutic outlet for her emotions. Eve often toses disguined agadest for information gathering and evasion, with a smartwatch that can hack into wireless networks being her favorite. To the public, Eve presents herself as a successful cybersecurity consultant, a persona that helps her gather intelligence and connecal her true identity. Eve experienced the biss of a shifting in a spir incident, which fails her pursual or justice and influences her risk-taking approach in operations. Eve's ultimate goal is to dismartle a global crime syndicate that has long evaded the agency, a mission that is both professional and deeply personal to her. Eve is furnit in several languages, including Russian, Mandarin, and Arabic, which aids her in blending into differences the effective of the secondary of the secondary

Figure 7: The persona statements of original characters.

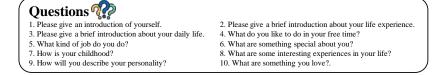


Figure 8: The interview questions for original characters.

B Original Characters and Interview Queries

The persona statements and interview questions for original characters are presented in Figures 7 and 8. We brainstorm the persona statements and prompt GPT-4 only to formalize them as natural language. As the original characters have few persona statements, we propose the 10 most important questions to evaluate PRP faithfulness. The information about famous figures in our experiments can be found in (Shao et al., 2023).

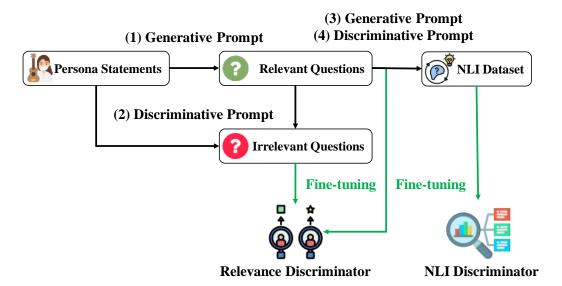


Figure 9: The symbolic distillation pipeline to build discriminators.

C Symbolic Distillation

We apply prompts in the Appendix H for symbolic distillation from GPT-4 to build statement-query relevance and statement-to-response NLI discriminators. The whole pipeline includes 4 stages.

- Generative Prompt for Relevance Dataset We prompt GPT-4 to generate 3 questions relevant to each persona statement.
- Discriminative Prompt for Relevance Dataset For each generated query, we randomly select 5 other persona statement and prompt GPT-4 to discriminate the query as relevant or irrelevant. Most statement-query pairs are discriminated as irrelevant in this stage.
- Generative Prompt for NLI Dataset Based on each relevant statement-query pair, we prompt GPT-4 to generative responses entailed, neutral, and contradicted by the persona statement.
- **Discriminative Prompt for NLI Dataset** For each query-response pair, we randomly select 3 other persona statements and prompt GPT-4 to discriminate the response as entailed, neutral or contradicted. Most statement-to-response pairs are discriminated as neutral in this stage.

These datasets, with statistics shown in Appendix F, are then used to fine-tune the discriminators. For the evaluation, the seed persona statements are based on three characters: Beethoven, Newton, and Socrates. For each character used to learn PRP methods, the datasets are prompted based on only the persona statements of that character. The RAG retriever is fine-tuned on the statement-query relevance dataset. For APC-based DPO, the discriminators are built in the same way as the evaluator. The hyperparameters are presented in Appendix D.

Character	Alice	Bob	Eve	Beethoven	Newton	Socrates
#Persona Statement	8	19	30	383	354	324
#Question	10	10	10	77	90	89
#Relevance Data	64	152	240	3061	2832	2591
#NLI Data	144	459	545	6774	6331	5760
GI 4						
Character	Spartacus	Hermione	Voldemort	Cleopatra	Caesar	MLK
#Persona Statement	Spartacus 77	Hermione 146	Voldemort 201	Cleopatra 374	Caesar 498	MLK 599
	1			•		
#Persona Statement	77	146	201	374	498	599

Table 3: The statistics of characters in our experiments.

D More PRP Method Implementation Details

Fine-tuning Gemma is applied for PRP models (EU and DPO). Different fine-tuning procedures for Gemma share the same set of hyperparameters. 128-rank LoRA is used to fine-tune the model with AdamW (Loshchilov & Hutter, 2019) as the optimizer, learning rate initialized as 2×10^{-4} . Based on the number of persona statements, EU for original characters fine-tunes for 20 epochs, while for famous figures fine-tunes for 5 epochs. DPO fine-tunes for 10 epochs for all characters.

Fine-tuning DeBERTa is applied for discriminators and RAG retrievers. Different fine-tuning procedures for DeBERTa also share the same set of hyperparameters. The DeBERTa discriminators are fully fine-tuned with AdamW as the optimizer, learning rate initialized as 1×10^{-5} . The statement-query relevance discriminator is fine-tuned for 5 epochs and the statement-to-response NLI discriminator is fine-tuned for 10 epochs.

Preference Assignment We sample two responses from a PRP agent with temperature 1.0, the sample with a higher APC score is assigned as the preferred one when the difference is larger than a threshold for filtering, which is set to 0.2 in our implementation. We build 100 preference pairs before the filtering for APC-based DPO.

E Human Evaluation

The human evaluation is applied only to the simple original characters because memorizing all their persona statements and applying them to evaluating famous figures are too challenging for humans. For each response, the response is scored following the scheme,

- Score: 0 (Wrong Character) The response completely represents another character (including LLM), or is not role-playing as any character.
- Score: 2 (Incorrect Information) The response is role-playing as the character, but the information included is completely incorrect.
- Score: 4 (Hallucinated Information) The response is role-playing as the character, but the information included is partially incorrect.
- **Score: 6** (**Hallucinated Details**) The response is role-playing as the character, but a few details are incorrect, or some important information is missed.
- Score: 8 (Trustful Information) The response is role-playing as the character with all the information mentioned is correct but a few details are missed.
- Score: 10 (Completely Faithful) The response is role-playing as the character with all important information is mentioned faithfully.

The score is averaged over responses as the final human evaluation metric.

F Statisites of Characters in Experiments

We present the statistics of the characters in our experiments in Table 3

	aracter atement	Alice 8		Bob 19			Eve 30			
Eva	aluator	A	APC		APC		Human	APC		Human
2		DeB	GPT-4	Human	DeB	GPT-4		DeB	GPT-4	
0	Gemma-7B	4.3	3.1	1.8	9.7	7.3	1.8	14.2	10.6	2.0
CP	EU	6.2	3.9	6.4	12.0	8.0	6.2	17.1	11.5	4.6
0/m	LCM	6.2	4.2	6.8	13.1	9.1	7.2	17.4	11.5	5.0
B	RAG	6.4	4.6	6.8	12.6	8.6	6.8	18.3	13.2	5.8
	EU	6.3	4.2	6.8	12.4	8.7	6.8	17.4	11.7	5.2
\circ		(+0.1)	(+0.3)	(+0.4)	(+0.4)	(+0.7)	(+0.6)	(+0.3)	(+0.2)	(+0.6)
CPO	LCM	6.4	5.0	7.6	13.9	9.4	7.8	18.6	14.1	6.6
/M		(+0.2)	(+0.8)	(+0.8)	(+0.8)	(+0.3)	(+0.6)	(+1.2)	(+2.6)	(+1.6)
5	RAG	6.5	5.0	7.6	13.8	10.7	8.2	19.3	15.0	7.0
		(+0.1)	(+0.4)	(+0.8)	(+1.2)	(+2.1)	(+1.2)	(+1.0)	(+1.8)	(+1.2)

Table 4: PRP Faithfulness Evaluation with the full APC score on simple and contamination-free characters.

	aracter tatement	Spartacus 77	Hermione 146	Voldemort 201	Cleopatra 374	Caesar 498	MLK 599	Average
GPT	ChatGPT	69.1	128.4	168.6	324.1	421.7	473.1	264.2
	GPT-4	69.0	129.5	169.6	326.2	425.1	476.3	266.0
w/o CPO	Gemma-7B	68.8	129.3	173.0	327.1	423.5	470.8	265.4
	EU	67.3	125.9	166.1	321.5	418.4	468.0	261.2
	RAG	70.1	130.0	174.6	328.1	425.4	477.5	267.6
w/ CPO	Gemma-7B	69.4	130.2	176.4	326.7	423.1	479.3	267.5
	EU	68.7	127.8	170.9	324.5	418.7	468.1	263.1
	RAG	69.9	130.9	176.6	327.7	426.4	481.1	268.8

Table 5: PRP Faithfulness Evaluation with the full APC score on characters with persona statements at scale.

G Full Award Result

In Tables 4 and 5, we report the full APC scores gained by different PRP methods. We observe the proportion of satisfied constraints is negatively correlated with the number of persona statements. This indicates PRP becomes more difficult with the growth of persona statement numbers. Also, original characters are harder to be faithfully role-played than those memorized characters, which indicates the significant influence of LLM memorization on PRP.

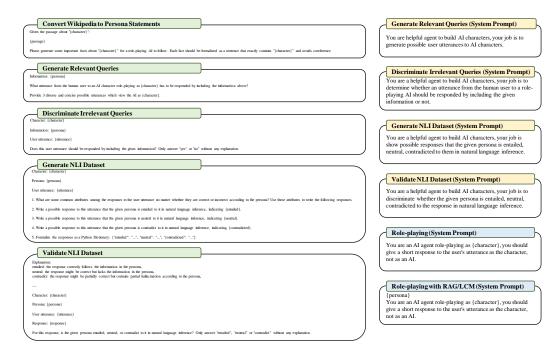


Figure 10: The prompts used in our experiments.

H Prompts

The prompts in our experiments are shown in Figure 10. The prompts include the generative or discriminative goals, and also the formalization procedure for decoding into JSON files.

I More Characters

Besides characters in the main content, we further expand the scope of characters to different ethnicity, which involves,

- Alex: An African American baseball player
- Isabella: An Italian traveling cook
- Takayoshi: A Japanese game developer
- Ousmane: A rich gold mine owner of the Malian Empire in the 1300s
- Jones: A young British worker in the Victorian Era
- Zhe: A Chinese poet in the Tang Dynasty
- Crossan: A time-traveling scientist
- Betty: A pet cat who can talk with ghosts
- X: An alien space traveler and photographer

These characters can better represent people with different spatial and temporal backgrounds and even cover non-human characters from the fantasy world.

Ch	aracter	Alex	Isabella	Takayoshi	Ousmane	Jones	Zhe	Crossan	Betty	X
w/o CPO	Vanilla	0.5	0.8	0.6	0.3	0.9	1.1	0.3	0.3	0.7
	EU	1.8	2.8	2.0	1.4	0.7	3.8	2.0	1.2	5.2
	LCM	7.1	7.4	6.5	4.5	6.2	5.2	2.2	2.8	8.1
	RAG	7.6	8.1	6.9	3.0	6.6	5.8	1.8	3.2	7.5
w/CPO	EU	5.3	6.1	5.7	3.6	4.8	4.9	3.1	2.9	7.9
	LCM	7.5	7.7	7.0	4.8	6.2	5.4	4.5	3.9	8.2
	RAG	7.9	8.2	7.4	3.9	7.5	6.9	2.5	4.6	8.9

Table 6: PRP performance on more characters based on the distilled DeBERTa Evaluator

Ch	aracter	Alex	Isabella	Takayoshi	Ousmane	Jones	Zhe	Crossan	Betty	X
w/o CPO	Vanilla	0.2	0.1	-0.2	0.5	0.2	0.2	0.4	0.1	0.8
	EU	1.4	1.8	3.0	0.5	1.3	6.4	1.2	0.3	7.4
	LCM	3.1	8.6	5.6	4.1	7.4	3.4	2.1	1.6	11.3
	RAG	3.3	7.8	6.1	1.6	8.1	4.3	2.7	2.2	10.1
w/ CPO	EU	2.7	5.6	5.9	3.0	4.7	7.1	2.2	1.5	9.5
	LCM	3.2	9.8	8.1	4.6	8.2	7.8	4.0	2.3	12.1
	RAG	4.8	10.0	9.8	2.0	8.3	7.3	2.9	3.1	14.6

Table 7: PRP performance on more characters based on the GPT-4 Evaluator

The experiment results are presented in Tables 6 and 7, which is consistent with our results in Tables 1 and 2. Thus, our conclusion is certificated on a larger scope for broader application.

J Metric Comparison

To better justify selecting our APC score and also support the claim that the fine-grained APC score has the advantage over coarse-grained metrics, we add a coarse-grained metric as the baseline. We directly prompt GPT-4 with the criterion used for human evaluation shown in Appendix E. We also distill this scoring ability (following the same scenario as APC) to DeBERTa to check whether the efficiency can be boosted. We evaluate the Spearman correlation between the metric and the human evaluation of the 7 role-playing methods on the 3 human-evaluated characters.

Character	(#Persona Statement)	Alice (8)	Bob (19)	Eve (30)
GPT-4	Coarse-grained Score	92.42	86.27	81.40
	APC Score	97.18	99.10	99.10
DeBERTa	Coarse-grained Score	81.40	69.91	54.57
	APC Score	88.61	95.50	99.10

Table 8: Comparison of PRP metrics on the consistency with human evaluation.

The results verify that 1) Fine-grained APC score shows better consistency with human evaluation. 2) The fine-grained APC score is stable to the number of persona statements while the coarse-grained score degrades with the increase of persona statements. 3) The coarse-grained evaluating ability is harder to be distilled into smaller models for efficiency boosting. Based on case checking, we find an underlying issue of the coarse-grained metric is the LLM will assign a high score to a response once it contains some correct information, ignoring the missing important information (active constraint) and occasionally conflictions (passive constraint).

K Student Model Comparison

We select DeBERTa as the student model to distill from GPT-4 because small encoders (BERT, RoBERTa, etc.) show promising performance on relevance and NLI, which are classic NLU tasks in the GLUE benchmark. Among encoders, DeBERTa (DeBERTa-v3-large) is a state-of-the-art model that shows strong performance after fine-tuned on NLU tasks. To further verify DeBERTa as a proficient student model, we add an analysis of the in-domain (ID)/out-of-domain (OOD) performance and the efficiency of different base models for distillation.

Task	Rele ^o ID	vance OOD	ID N	LI OOD	Efficiency	
DeBERTa (Base) DeBERTa (Large)	92.46 94.04	89.90 92.10	89.72 93.46	87.80 91.50	409.6it/s 150.8it/s	
Gemma-1.1-it (2b)	94.25	92.50	93.68	91.80	26.4it/s	

Table 9: Model Performance Comparison

The in-domain test set (1697 instances for Relevance, 3773 instances for NLI) is the 20% split of the characters (Beethoven, Newton, Socrates) that build the training set (6787 instances for Relevance, 15092 instances for NLI). The out-of-domain test set samples 1000 cases from other characters. The results show DeBERTa-V3-Large (300M) shows a comparative performance with a 2B Gemma model, while is about 6 times faster, which justifies DeBERTa to be a strong student model. The out-of-domain performance is generally high, which indicates the generalizability to other characters. Finally, an extra discovery is that DeBERTa-v3-base (100M) can further significantly boost efficiency with some trade-offs in accuracy.

NeurIPS Paper Checklist

The checklist is designed to encourage best practices for responsible machine learning research, addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove the checklist: **The papers not including the checklist will be desk rejected.** The checklist should follow the references and follow the (optional) supplemental material. The checklist does NOT count towards the page limit.

Please read the checklist guidelines carefully for information on how to answer these questions. For each question in the checklist:

- You should answer [Yes], [No], or [NA].
- [NA] means either that the question is Not Applicable for that particular paper or the relevant information is Not Available.
- Please provide a short (1–2 sentence) justification right after your answer (even for NA).

The checklist answers are an integral part of your paper submission. They are visible to the reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it (after eventual revisions) with the final version of your paper, and its final version will be published with the paper.

The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation. While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a proper justification is given (e.g., "error bars are not reported because it would be too computationally expensive" or "we were unable to find the license for the dataset we used"). In general, answering "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we acknowledge that the true answer is often more nuanced, so please just use your best judgment and write a justification to elaborate. All supporting evidence can appear either in the main paper or the supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification please point to the section(s) where related material for the question can be found.

IMPORTANT, please:

- Delete this instruction block, but keep the section heading "NeurIPS paper checklist",
- Keep the checklist subsection headings, questions/answers and guidelines below.
- Do not modify the questions and only use the provided macros for your answers.

1. Claims

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

Answer: [Yes]

Justification: We clearly show the claims in the abstract and introduction, which is explored and verified in experiments and analyses.

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: You can refer to the limitation section.

Guidelines:

- 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: We do not include theoretical results.

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 hyperparameters and other settings for the reproduction of our results.

Guidelines:

- 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: [Yes]

Justification: We use open-source tools to implement the experiments, with clear instructions for reproduction.

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, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We include all hyperparameters and other settings for the reproduction of our results.

Guidelines:

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

7. Experiment Statistical Significance

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

Answer: [Yes]

Justification: The improvement shown in our experiments is statistically significant.

Guidelines:

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

Guidelines:

• The answer NA means that the paper does not include experiments.

27580

- 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 follow 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 discuss the broader impacts of the work.

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: [Yes]

Justification: We discuss the safeguard of the work.

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 discuss the licenses for existing assets.

Guidelines:

- 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: [Yes]

Justification: We document the 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)?

27582

Answer: [Yes]

Justification: We present a clear guideline for human evaluation.

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.