

LLMs Become Logical Explainers through Multi-Agent Self-Play Framework

Anonymous ACL submission

Abstract

Natural language explanations have been widely employed to convey relations underlying task outputs; however, existing approaches rely on superficial persona formulations, with limited use of theoretically grounded personas. In this study, we propose a multi-agent self-play framework in which a single model assumes multiple explanatory personas informed by cognitive and social sciences theories to generate diverse explanation candidates for the same input. Relative preferences among these candidates are then induced through a multi-factor scoring that jointly accounts for persona logical alignment, self-critique factuality from a critic agent, and expression diversity. These preferences are conveyed to a recomposer agent to synthesize an optimal explanation, which is subsequently used for self-preference learning, enabling the model to strategically learn from its own generations. Experiments show improved logical alignment over single-persona attributes and stable DPO training, demonstrating that adaptive persona selection can be effectively realized at test time.

1 Introduction

Since the advent of LLMs, research has moved beyond producing task outputs toward generating natural language explanations (NLEs) that articulate the reasoning behind model decisions. Such explanations are particularly valuable when the task interpretation is inherently ambiguous or subjective, or when producing explanatory annotations requires substantial human effort (He et al., 2024; Yao et al., 2023).

However, a principled understanding of how to guide models to generate reliable explanations remains underdeveloped. In particular, LLM-based NLE generation is highly sensitive to the choice of reasoning strategy, and the vast design space makes identifying an optimal configuration challenging (Gemechu et al., 2025; Kunz and

Kuhlmann, 2024). Moreover, it is often difficult for human evaluators to determine whether an explanation reflects genuine reasoning or merely a superficially plausible artifact (Agarwal et al., 2024).

In this paper, we introduce the **Self-Play Explainer (SPE)**, a multi-agent framework that decomposes a single LLM into multiple distinct reasoning agents, as shown in Figure 1. We define theory-grounded personas informed by cognitive and social science theories and analyze their preference relations in NLE generation, rather than treating personas as superficial role variations (Tseng et al., 2024). Explanation generation entails qualitatively distinct reasoning processes, such that personas effective for explaining simple factual relations may differ from those better suited for explanations requiring scientific background or contextual knowledge (Ahn et al., 2025). By distinguishing personas along principled dimensions, our framework enables systematic analysis of how distinct reasoning bases are selected and emphasized across tasks.

Although an optimal persona for NLE generation lacks explicit ground truth and is difficult to define¹, we address this challenge through a multi-factor preference scoring that aggregates relative preferences among persona-generated NLEs. This scoring jointly accounts for logical alignment between the input-output pair and each explanation (Scirè et al., 2024; Yang et al., 2024), critique-based penalties from a dedicated critic agent (Li et al., 2025), and expression diversity across explanations, enabling comprehensive evaluation without reliance on absolute supervision.

Based on the induced preferences, we instruct a recomposer agent to synthesize a final NLE by integrating complementary elements from preferred explanation candidates (Wang and Atanasova, 2025;

¹This challenge stems from the inherent subjectivity of human judgment (Chen et al., 2025), and the absence of guarantees that a single persona yields consistently aligned reasoning behaviors across different LLMs.

079 Wang et al., 2025). This design is motivated by
080 the observation that robust explanations arise from
081 selectively integrating diverse perspectives rather
082 than relying on a single persona. The resulting NLE
083 and its aggregated scores are then used to construct
084 preference pairs for supervised fine-tuning (SFT)
085 and direct preference optimization (DPO) (Rafailov
086 et al., 2023), yielding an explanation model aligned
087 with integrated persona behavior.

088 Our experiments demonstrate that the proposed
089 framework improves logical alignment over single-
090 persona while preserving expression diversity
091 within a controlled range. We observe gains in
092 paragraph-level logical alignment alongside re-
093 duced sentence-level variability, resulting in more
094 coherent NLEs. Further analyses indicate that our
095 preference dataset enables stable DPO training and
096 the multi-factor scoring process is effective relative
097 to the hard negative variant, highlighting the ef-
098 fectiveness of our framework for adaptive persona-
099 driven NLE generation.

100 2 Related Work

101 2.1 The Role of 102 Natural Language Explanations

103 NLEs have become a central component across var-
104 ious tasks, supporting interpretability and serving
105 as auxiliary supervision. Early studies incorpo-
106 rated human-written explanations to justify pre-
107 dictions in tasks such as natural language infer-
108 ence (NLI) and question answering (QA) (Rajani
109 et al., 2019; Camburu et al., 2018). Subsequent
110 work has shifted toward automatic generation and
111 explanation-aware training, positioning NLEs as
112 a cross-task mechanism that supports reasoning,
113 evaluation, and learning beyond prediction accu-
114 racy (Quan et al., 2025; Wang and Atanasova, 2025;
115 Wadhwa et al., 2024).

116 Despite recent efforts, explanation-aware ap-
117 proaches have relied on fixed source NLEs that de-
118 pend on a single reasoning trajectory, limiting their
119 effectiveness across diverse contexts (Honda and
120 Oka, 2025; He et al., 2024). In practice, NLEs often
121 require different theoretical viewpoints depending
122 on the context (Ahn et al., 2025). Motivated by
123 this, we employ a persona-conditioned self-play
124 paradigm within a single model to generate expla-
125 nation candidates from multiple perspectives. This
126 approach facilitates learning preference relations
127 among them, thereby improving the logical consis-
128 tency of NLEs through comparative optimization.

129 2.2 Persona Conditioning 130 under a Perspectivist Framework

131 Early work employed personas to maintain con-
132 versational consistency in dialogue systems, in-
133 ducing role-specific behavior through prompt de-
134 sign or fine-tuning (Zhang et al., 2018). Recent
135 approaches have leveraged LLMs to elicit diverse
136 perspectives; however they described personas in
137 terms of surface-level attributes, with fewer works
138 defining them based on empirically grounded theo-
139 retical principles (Tseng et al., 2024).

140 We adopt a perspectivist view from cognitive
141 and social sciences, which characterizes NLEs as in-
142 herently dependent on theoretical viewpoints rather
143 than as absolute accounts. Importantly, in our study,
144 perspectivism does not imply arbitrariness but in-
145 stead recognizes that different viewpoints can high-
146 light compatible aspects of the reasoning process
147 within a single model (Lombrozo, 2006).

148 2.3 Recent Advances in 149 Multi-Agent Reasoning

150 Research on multi-agent has emerged to address
151 the limitations of single-pass inference in complex
152 reasoning tasks. Early studies focused on agent de-
153 composition, assigning specialized roles such as
154 planning or verification to enhance reasoning qual-
155 ity (Ma et al., 2025; Hong et al., 2023). Subsequent
156 approaches leveraged agent interactions to explore
157 the alternative reasoning paths and enable mutual
158 refinement (Karthikeyan et al., 2025).

159 In this study, we design a unified multi-agent
160 framework that spans NLE generation, evaluation,
161 and recombination to optimize reasoning within a
162 single model. Persona-conditioned agents are em-
163 ployed not only to generate explanations, but also
164 to comparatively assess their quality and recombine
165 optimized NLEs (Li et al., 2025; Wang et al., 2025).
166 In this framework, multi-agent reasoning functions
167 as a core mechanism for improving logical consis-
168 tency, rather than as an auxiliary tool.

169 3 Self-Play Explainer

170 3.1 Theory-Grounded Persona Definition

171 We design five explanatory personas with distinct
172 reasoning policies to generate NLEs grounded in
173 cognitive and social science theory. In contrast to
174 prior approaches that rely on straightforward pat-
175 terns (Honda and Oka, 2025; He et al., 2024), we
176 hypothesize that strategically designed personas
177 enable a single model to elicit diverse explanatory

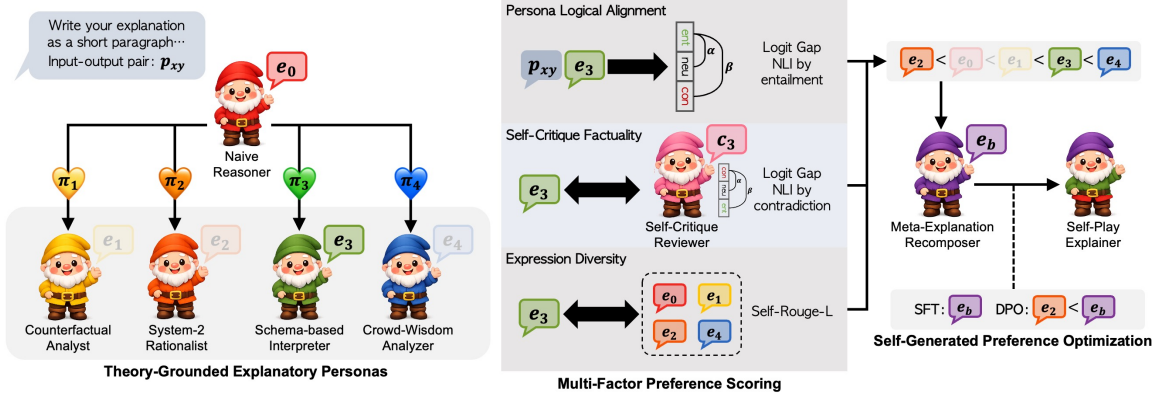


Figure 1: Overview of the proposed **Self-Play Explainer (SPE)** framework. For a given input-output pair, multiple theory-grounded explanatory personas generated candidate explanations (§3.1), which are ranked using multi-factor preference scoring (§3.2). The top-ranked explanations are then recomposed for self-generated preference optimization (§3.3), resulting in a more logically consistent explanation model.

behavior for the same phenomenon. In this setting, multiple reasoning paths are expected to reveal logical biases that a single NLE may overlook, while their comparison facilitates evaluation of the underlying evidential structure.

Each persona operates as follows: (1) **Naive Reasoner**, relying solely on the model’s internal knowledge without external theoretical grounding; (2) **System-2 Rationalist**, explicitly expressing chain-of-thought reasoning aligned with Kahneman’s dual-process theory (Kahneman, 2011); (3) **Counterfactual Analyst**, applying ‘what-if’ contrasts inspired by Pearl’s causal framework (Pearl, 2009); (4) **Schema-based Interpreter**, providing structured interpretations grounded in Bartlett’s schema theory (Bartlett, 1995); and (5) **Crowd-Wisdom Synthesizer**, integrating multiple perspectives following Surowiecki’s principle of collective intelligence (Surowiecki, 2005)².

We define p_{xy} as the task-specific formulation text for the x - y pair, and NLE generated by the k -th persona as e_k . Here, x and y denote the task input and output, π_k denotes the persona-specific reasoning policy, and the model p_{LLM} is shared across personas. We denote by E the set of NLEs generated by the K personas for the same p_{xy} .

$$e_k \sim p_{\text{LLM}}(e \mid p_{xy}, \pi_k), \quad (1)$$

$$E = \{e_k \mid k = 1, \dots, K\}, \quad (2)$$

3.2 Multi-Factor Preference Scoring

We establish systematic criteria for identifying preferable personas, recognizing that the optimal

²More descriptions of each persona are provided in Appendix A.1, and their prompts are in Appendix F.

persona choice depends on the reasoning demands. For instance, the persona suited for simple factual verification may differ from the one preferable for scientific justification. Since a single, universally applicable ground truth NLE preference is difficult to define, we focus instead on distinguishing *relatively preferred* NLEs among persona candidates generated by the same model.

3.2.1 Persona Logical Alignment

Recent studies have adopted NLI-based scores to evaluate the quality of generated texts (Scirè et al., 2024; Yang et al., 2024). However, such approaches are limited in capturing neutrality and contradiction in persona-generated NLEs, as establishing an absolute persona preference is inherently difficult. To address this, we introduce a logit gap measure that explicitly incorporates supplementary evidence, with calibrated hyperparameters to capture meaningful *flips*. A *flip* occurs when two NLEs have similar entailment logits but their order is reversed once neutral and contradictory evidence is considered; when this happens under only a marginal difference, we refer to it as a *near-tie flip*.

First, we define $\ell_{xy,k}$ as the NLI logit for a given class between p_{xy} and e_k .

$$\ell_{xy,k}^{\{\text{ENT,NEU,CON}\}} = \text{NLI}^{\{\text{ENT,NEU,CON}\}}(p_{xy}, e_k), \quad (3)$$

In measuring the logit gap, we apply hyperparameters α and β to the supplementary evidence. Rather than treating them as fixed constants, we adaptively calibrate them to the given training data and model. We then select the (α, β) pairs that capture meaningful *near-tie flips* through a grid

search³. Based on this, we define the entailment logit gap between p_{xy} and e_k as follows:

$$gap_{xy,k}^{(ENT)} = \ell_{xy,k}^{(ENT)} - \alpha \cdot \ell_{xy,k}^{(NEU)} - \beta \cdot \ell_{xy,k}^{(CON)}, \quad (4)$$

We apply softmax over the logit gap across different NLEs to model their relative preference, shifting the probability distribution from class to candidate-wise competition. To avoid distortion from absolute logit scale differences, we apply z-score normalization to the logit gaps before softmax.

$$S_{PLA}(p_{xy}, e_k, E) = \text{softmax}_k(\widetilde{gap}_{xy,k}^{(ENT)}) \quad (5)$$

This enables a more fine-grained assessment of each persona’s logical alignment, even when entailment scores are similar, by flexibly incorporating supplementary evidence informed by the data and model characteristics used in NLE generation.

3.2.2 Self-Critique Factuality

Recent approaches to evaluating the factuality of LLM-generated content have prompted models to produce self-assessment scores or leverage internal feature signals (Zhang et al., 2024b; Parcalabescu and Frank, 2024). However, these methods often suffer from limited interpretability or high computational cost. To mitigate this, we introduce a **Self-Critique Reviewer** agent that enables the model to generate claim-level critiques of its own NLEs (Li et al., 2025), while we impose contradiction-focused penalties to directly capture inconsistencies.

The agent was assigned two subtasks: classify the factual relationship between p_{xy} and e_k , and generate a claim-level critique consistent with the selected scale (Chu et al., 2025). The claim-level critique generated by the k -th persona, denoted c_k , is defined as follows, where π_{crit} denotes the self-critique reasoning policy.

$$c_k \sim p_{LLM}(c | p_{xy}, e_k, \pi_{crit}), \quad (6)$$

$$C = \{c_k | k = 1, \dots, K\}, \quad (7)$$

We segment e_k and c_k into sentences, denoted as $e_k = \{s_1, \dots, s_N\}$ and $c_k = \{t_1, \dots, t_M\}$. To identify counter-argument pairs, we apply the logit gap measure under the *near-tie flip* condition defined as contradiction. The hyperparameters α and β are adaptively determined by the training data

³ The formal definitions of (*near-tie*) flips, grid search procedure, and (α, β) pairs selected for each data and model are provided in Appendix A.2.

and model³. Based on this, we define the contradiction logit gap between s_n and t_m as follows:

$$gap_{n,m}^{(CON)} = \ell_{n,m}^{(CON)} - \alpha \cdot \ell_{n,m}^{(NEU)} - \beta \cdot \ell_{n,m}^{(ENT)}, \quad (8)$$

We focus on the distribution tail capturing the strongest counterarguments, selecting cases with the most pointed critiques. Among the $N \times M$ sentence pairs, we define \mathcal{T}_q as the subset whose contradiction logit gaps fall within the top $q\%$, and compute the mean over this set.

$$\mathcal{C}(e_k, c_k) = \frac{1}{|\mathcal{T}_q|} \sum_{n,m \in \mathcal{T}_q} gap_{n,m}^{(CON)}, \quad (9)$$

We apply softmax over the logit mean across different NLEs to model their relative preference, including z-score normalization before softmax.

$$S_{SCF}(e_k, c_k, E, C) = \text{softmax}_k(\widetilde{\mathcal{C}}(e_k, c_k)), \quad (10)$$

This enables a more controlled evaluation by requiring the model to generate a self-critique proxy from a strictly critical perspective, with contradiction-focused penalties imposed. Supplementary evidence is flexibly incorporated according to the characteristics of the data and model.

3.2.3 Expression Diversity

We expect different personas to exhibit independent reasoning styles, even when explaining the same x - y relation. Thus, we extend similarity-based diversity assessment (Zhang et al., 2025; Zhu et al., 2018) to the paragraph-level and interpret higher cross-persona similarity as lower relative expression diversity. In this setting, we employ Rouge-L, which captures shared narrative structure through the longest common subsequence.

$$S_{ED}(e_k, E) = \frac{1}{K-1} \sum_{k' \neq k} \text{Rouge-L}(e_k, e_{k'}), \quad (11)$$

However, length imbalance between NLEs may bias similarity and distort evaluation. To mitigate this, we apply a length-based weighting scheme that down-weights pairs with large discrepancies, using a Laplacian kernel to convert length differences into an attenuation factor that decreases a pair’s contribution as the gap increases (Rahimi and Recht, 2007)⁴.

$$w_{kk'} = \exp(-\gamma | |e_k| - |e_{k'}| |), \quad (12)$$

⁴We set $\gamma = 0.04$, so that a one-sentence length difference halves the weight of a pair of NLEs. Since all personas were instructed to follow the same length constraint, such deviations may also indicate inconsistency in instruction following.

We incorporate this $w_{kk'}$ into the original scheme to derive the expression diversity, enabling the assessment of independent reasoning styles by evaluating the diversity of multiple NLEs from the same model while accounting for length bias.

$$S_{ED}(e_k, E) = \frac{\sum_{k' \neq k} w_{kk'} \cdot \text{Rouge-L}(e_k, e_{k'})}{\sum_{k' \neq k} w_{kk'}}, \quad (13)$$

3.3 Self-Generated Preference Optimization

We compute a harmonic mean to discourage candidates that perform well on only a single score and instead prefer NLEs with balanced quality. Since S_{PLA} is preferred when larger, whereas S_{SCF} and S_{ED} are preferred when smaller, this relationship is reflected in the formulation⁵.

$$S_{\text{Final}} = \frac{3}{\frac{1}{S_{\text{PLA}}} + \frac{1}{1-S_{\text{SCF}}} + \frac{1}{1-S_{\text{ED}}}}, \quad (14)$$

Rather than relying on the first-ranked result, we assume that aggregating diverse perspectives yields a more robust and comprehensive NLE. To this end, we introduce a **Meta-Explanation Recomposer** agent that identifies commonly supported evidences, determines a more reliable interpretation, and generates a single final NLE (Wang and Atanasova, 2025; Wang et al., 2025).

To incorporate the final preferences, we define the NLE produced by the recomposer agent as e_{win} and the one with the lowest S_{Final} as e_{lose} . We first perform SFT to generate e_{win} for each x - y pair, and then apply DPO with $(e_{\text{win}}, e_{\text{lose}})$ pairs to encode the preference signal.

4 Experimental Results

4.1 Tasks and Models

We select extractive QA and NLI because both tasks inherently require NLEs to justify multi-step evidence-based reasoning or semantic relations between premises and hypotheses. For QA, we use HotpotQA, which involves multi-document reasoning (Yang et al., 2018), and NQopen, a short-answer variant of NaturalQuestions reflecting real user queries (Lee et al., 2019). For NLI, we use ANLI across its R1, R2, and R3 splits with increasing reasoning difficulty (Nie et al., 2020) and eSNLI, which provides annotated explanations as part of supervision (Camburu et al., 2018).

⁵Although score weights can be adjusted for specific applications, we use equal weighting in this study.

We employ meta-llama/Llama-3.2-3B(-Instruct) as the base model for all experiments. Although smaller models are more susceptible to hallucination, our framework mitigates it by leveraging multi-factor preferences and a recomposer agent, enabling the model to self-refine its own NLEs (Wang and Atanasova, 2025).

4.2 Evaluation Metrics

We first measure the logical alignment between each task input-output pair and its generated NLE using an NLI-based score. We treat them as a premise-hypothesis pair to compute the entailment softmax probability from a pre-trained NLI model¹⁰. This design follows recent approaches that leverage NLI scores as a proxy for faithfulness and logical coherence (Zhang et al., 2024a).

We employ the introduced S_{ED} score to evaluate the diversity of generated NLEs. Under a fixed sampling setup, this metric measures expression-level similarity among multiple explanations generated for the same input. Lower scores indicate reduced overlap across samples, corresponding to higher explanatory diversity.

4.3 Main Results

We adopted a persona-conditioned explanation setting that includes both individual personas and a Best-of-N selection of their per-sample optimal outputs. Since persona conditioning deliberately shifts the model’s typical token distribution through counterfactual assumptions or schema-level constraints (Lutz et al., 2025), this setup enables us to assess whether such deviations are applied appropriately when required⁶.

We report the performance comparison in Table 1. Across tasks, the performance of individual explanatory personas exhibits largely consistent trends. The Crowd-Wisdom Synthesizer attains the highest logical alignment, whereas the Schema-based Interpreter performs the weakest, with the largest logical alignment gap of 39.84 observed on HotpotQA. The remaining personas generally show intermediate performance; while their results are comparable on NLI, the Counterfactual Analyst demonstrates a clear advantage on QA tasks over the System-2 Rationalist. These results indicate that the effectiveness of personas for NLE generation depends on task-specific characteristics.

⁶Further details on dataset curation and training configuration are provided in Appendix B.

Dataset	HotpotQA		NQopen		ANLI						eSNLI	
					R1		R2		R3			
	$P^{(ENT)} \uparrow$	$S_{ED} \downarrow$	$P^{(ENT)} \uparrow$	$S_{ED} \downarrow$	$P^{(ENT)} \uparrow$	$S_{ED} \downarrow$	$P^{(ENT)} \uparrow$	$S_{ED} \downarrow$	$P^{(ENT)} \uparrow$	$S_{ED} \downarrow$	$P^{(ENT)} \uparrow$	$S_{ED} \downarrow$
System-2 Rationalist	48.50	28.26	41.54	29.94	47.02	30.50	44.88	30.18	43.79	27.61	36.75	27.49
Counterfactual Analyst	57.56	28.50	54.22	30.29	45.77	25.71	44.48	25.80	43.69	24.65	36.34	24.79
Schema-based Interpreter	30.00	30.15	31.24	30.67	30.32	29.93	29.84	29.76	32.53	28.81	25.25	27.74
Crowd-Wisdom Synthesizer	69.84	24.66	63.89	26.34	44.64	26.11	45.50	25.94	53.15	24.81	48.27	23.46
Best-of-Personas	91.04	21.63	88.92	23.30	81.29	22.23	81.12	22.19	82.41	21.02	77.47	20.58
Self-Play Explainer (SFT)	66.39	20.79	67.32	28.72	47.96	28.62	47.70	26.90	48.35	26.93	32.29	21.81
Self-Play Explainer (DPO)	75.71	30.04	76.29	33.24	47.22	35.08	50.02	34.83	54.48	34.90	68.02	32.95

Table 1: Performance comparison of the proposed **Self-Play Explainer (SPE)** framework under the persona-conditioned explanation setting. $P^{(ENT)}$ represents the entailment probability, where higher values indicate better logical consistency, while S_{ED} measures expression diversity, where lower values indicate higher diversity. We report the average scores computed over three sampling runs for diversity scores.

Analysis of the proposed framework shows that SPE (SFT) yields task-dependent variations in logical alignment, resulting in gains for some tasks or degradations for others. In contrast, applying SPE (DPO) consistently improves logical alignment across all tasks, outperforming configurations that rely on individual personas. Although SPE (SFT) alone provides modest gains, it effectively aligns persona-informed NLEs within the model’s behavioral space, thereby shaping a model distribution that facilitates subsequent SPE (DPO) training.

Under the Best-of-Personas, all explanatory personas are sampled and the highest-scoring instance is selected (Rafailov et al., 2023), yielding high logical alignment across tasks. However, this approach is computationally impractical, as it requires K sampling completions and corresponding metric evaluations per test instance. The persona distributions favored by this strategy vary across tasks, and are further discussed in §5.2.

Regarding expression diversity, individual personas exhibit comparable scores, with the Crowd-Wisdom Synthesizer showing slightly lower values, likely due to its simulation of multiple expert perspectives that encourages more varied narrative structures. Compared to SPE (SFT), the score increases after applying SPE (DPO), reflecting a trade-off introduced by aligning the model’s token distribution toward persona-specific behaviors. Nonetheless, this increase is not markedly larger than that observed in single-personas.

5 Discussion

5.1 Granularity Effects in Logical Alignment

Our main results measure logical alignment of the generated NLE at the paragraph level. To analyze alignment at finer granularity, we additionally compute sentence-level alignment by decompos-

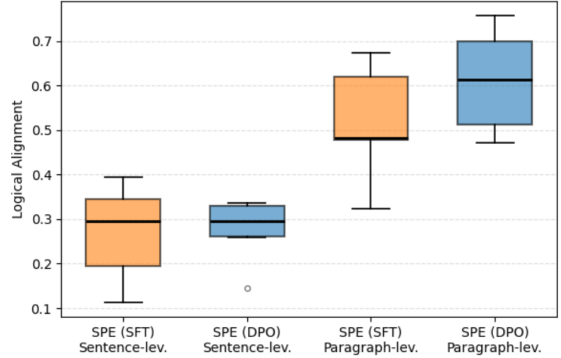


Figure 2: Boxplot distributions of sentence- (left) and paragraph-level (right) logical alignment scores for SFT (orange) and DPO (blue) models in the SPE.

ing each explanation into sentences and evaluating them against the input-output pair. This enables a more precise characterization of whether observed gains stem from individual sentence-level accuracy or from coherent organization at the paragraph-level. The results are presented in Figure 2.

While SPE (SFT) and SPE (DPO) exhibit comparable alignment at the sentence level, with average scores of 27.03 and 27.72 respectively, DPO achieves substantially higher alignment at the paragraph level, achieving an average score of 61.01, compared to 51.67 for SFT⁷. This pattern suggests that DPO improves global coherence and cross-sentence consistency. Moreover, DPO exhibits reduced variability in sentence-level scores, indicating more stable behavior and demonstrating its effectiveness in aligning NLEs.

5.2 Persona Selection Dynamics in Best-of-Personas

To characterize persona contributions within the best-of-personas, we analyze the proportion of in-

⁷Detailed sentence-level alignment scores for each task and related statistics are reported in Table 5 of Appendix C.1.

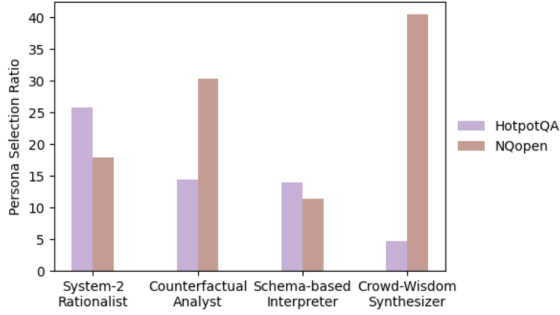


Figure 3: Persona selection ratios for each QA task under the best-of-personas with logical alignment.

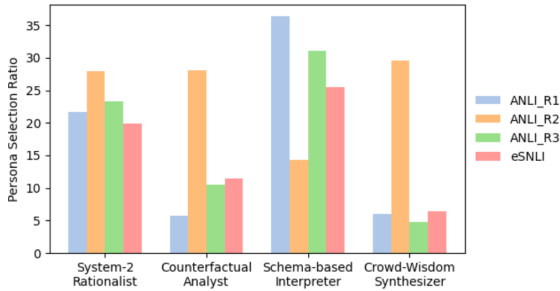


Figure 4: Persona selection ratios for each NLI task under the best-of-personas with logical alignment.

stances in which each persona produces the highest-scoring NLE. This offers a quantitative view of persona dominance for each task, capturing how frequently different reasoning styles yield preferred NLEs beyond aggregate performance scores. The results are presented in Figures 3 and 4⁸.

The selected personas across QA tasks reflect the structural properties of the underlying datasets. HotpotQA requires explicit multi-hop reasoning, which favors the System-2 Rationalist due to its deliberate, stepwise reasoning style. In contrast, NQOpen involves open-ended queries that benefit from broad knowledge aggregation, making the Crowd-Wisdom Synthesizer more effective. These results indicate that effective persona selection in QA is closely tied to task-specific reasoning and knowledge requirements.

In contrast, persona selection in NLI tasks is less consistent, with no clear stratification by difficulty observed in ANLI. Within this context, the Schema-based Interpreter achieves the highest logical alignment, as NLI emphasizes abstract relational reasoning over knowledge aggregation. Its schema-driven approach aligns naturally with the structured premise-hypothesis relations, resulting

⁸Although the best-of-personas includes the Naive Reasoner, we exclude its results in this section to focus on comparisons across the theoretical personas.

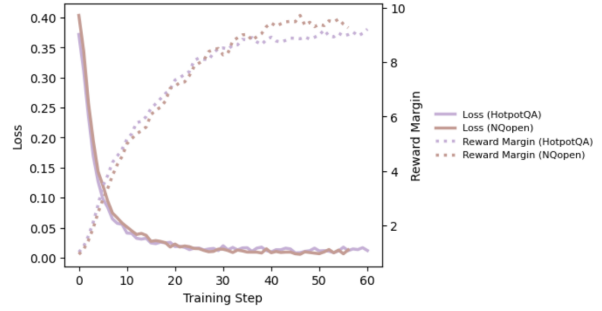


Figure 5: DPO training dynamics across QA tasks, with the training loss and reward margin between chosen and rejected responses.

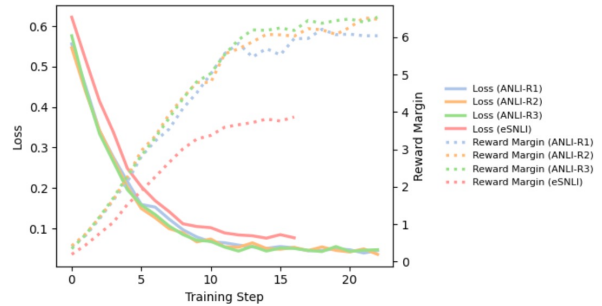


Figure 6: DPO training dynamics across NLI tasks, with the training loss and reward margin between chosen and rejected responses.

in more consistent alignment. These greedy best-of-persona results represent an upper bound attainable in the absence of computational constraints. However, the proposed framework learns sample-conditioned persona preferences, enabling adaptive persona selection at test time.

5.3 Training Stability Analysis

We examine both the loss and the reward margin between chosen and rejected responses during DPO training, as shown in Figures 5 and 6. The reward margin is computed as the difference between the chosen and rejected implicit rewards, and serves as a direct indicator of preference alignment during DPO optimization.

We observe consistent loss convergence and steadily increasing reward margins across all tasks. While the reward margin reaches approximately 4-6 for NLI tasks, it converges to a higher value of around 10 for QA tasks. Since the two hyperparameter settings differ only in batch size due to their dataset scale, these results suggest that our preference dataset enables stable DPO training.

Dataset	HotpotQA		NQopen		ANLI						eSNLI	
	Metric		Metric		R1		R2		R3		Metric	
	$P^{(ENT)} \uparrow$	$S_{ED} \downarrow$	$P^{(ENT)} \uparrow$	$S_{ED} \downarrow$	$P^{(ENT)} \uparrow$	$S_{ED} \downarrow$	$P^{(ENT)} \uparrow$	$S_{ED} \downarrow$	$P^{(ENT)} \uparrow$	$S_{ED} \downarrow$	$P^{(ENT)} \uparrow$	$S_{ED} \downarrow$
Schema-based Interpreter	30.00	30.15	31.24	30.67	30.32	29.93	29.84	29.76	32.53	28.81	25.25	27.74
Crowd-Wisdom Synthesizer	69.84	24.66	63.89	26.34	44.64	26.11	45.50	25.94	53.15	24.81	48.27	23.46
Self-Play Explainer (SFT)	66.39	20.79	67.32	28.72	47.96	28.62	47.70	26.90	48.35	26.93	32.29	21.81
Self-Play Explainer (DPO)	75.71	30.04	76.29	33.24	47.22	35.08	50.02	34.83	54.48	34.90	68.02	32.95
Self-Play Explainer (SFT) w Hard Negative	49.82	27.33	50.49	25.81	40.31	25.76	39.07	25.38	42.81	25.60	38.86	22.68
Self-Play Explainer (DPO) w Hard Negative	31.35	28.74	26.31	29.05	37.87	22.42	38.12	23.75	38.51	22.39	30.54	19.84

Table 2: Performance comparison of the proposed **Self-Play Explainer (SPE)** framework and its hard negative variant. Here, the single-persona baselines correspond to the two personas with the highest and lowest logical alignment across all tasks in the original results in Table 1. $P^{(ENT)}$ represents the entailment probability, where higher values indicate better logical consistency, while S_{ED} measures expression diversity, where lower values indicate higher diversity. We report the average scores computed over three sampling runs for diversity scores.

5.4 Ablation Study

To evaluate the functional contribution of the proposed framework, we conduct an ablation study with a hard negative variant. In this setting, e_{win} is defined as the persona ranked fourth according to the multi-factor scoring, while e_{lose} is selected using the same procedure as in the prior setting. The same SFT and DPO training pipelines are applied, and the results are provided in Table 2.

Under the standard framework, logical alignment improves steadily from SFT to DPO training, whereas the hard negative variant exhibits degraded performance. As SFT progresses, scores fall below those of the previously strongest single persona, the Crowd-Wisdom Synthesizer. This degradation is further amplified after DPO training and, on NQopen, the logical alignment falls below that of the weakest single persona, the Schema-based Interpreter. Although diversity scores show improvements across tasks, they are accompanied by a substantial decline in logical alignment.

These results indicate that aggressively constructed preference signals can impair NLE optimization, resulting in excessive degradation in logical alignment. We demonstrate that strategically leveraging the relations induced by multi-factor scoring enables effective DPO training that critically improves the logical consistency of NLEs.

6 Conclusion

We propose a multi-agent self-play framework based on a single model, where NLEs are generated from multiple theoretical perspectives for the same input-output pair. These explanations are evaluated using a multi-factor scoring that measures persona logical alignment, self-critique factuality with the reviewer agent, and expression diversity to

derive relative preferences. The recomposer agent then synthesizes the explanations and enables self-optimization through DPO training. Experimental results demonstrate that our framework improves logical alignment across all tasks compared to single-persona baselines. While no single persona consistently dominates across tasks, our framework achieves stable DPO training and implicitly enables adaptive persona selection at test time.

Limitations

Persona generalization Our study assumes the utility of persona-conditioned NLEs, without explicitly modeling when personas are required or dynamically assigning them to specific contexts. Instead, by jointly evaluating multiple personas, our approach enables self-optimization in NLE generation and implicitly captures adaptive persona selection behavior. Further exploration is required to assess whether alternative personas may be more effective in other domains.

Metric dependency In the absence of ground truth NLEs, our approach relies on multi-factor scoring to infer relative preferences rather than implying an optimal metric aligned with human annotation⁹. For more rigorous evaluation, future work is required to develop dedicated metrics that explicitly reflect persona-conditioned preferences.

Computational cost Although our framework introduces training-time overhead from repeated multi-agent inference, we observe stable learning even with the relatively small model, and expect similar trends to larger and more diverse models.

⁹We report GPT-4 win rates in Appendix C.3, along with an analysis and also its limitations for persona evaluation.

Ethics Statement

Persona conditioning is used to explore diversity in NLEs and does not represent real individuals or social groups. Personas are not tied to demographic or sensitive attributes and should be interpreted as modeling tools rather than sources of authority. We used an LLM with model outputs constrained by explicitly defined system prompts, and the model was not deployed in ways that introduce safety risks. All datasets used in this work are publicly available and do not contain personal or sensitive user information.

References

- Chirag Agarwal, Sree Harsha Tanneru, and Himabindu Lakkaraju. 2024. Faithfulness vs. plausibility: On the (un) reliability of explanations from large language models. *arXiv preprint arXiv:2402.04614*.
- Yongsu Ahn, Yu-Ru Lin, Malihe Alikhani, and Eun-jeong Cheon. 2025. Human-centered explanation does not fit all: The interplay of sociotechnical, cognitive, and individual factors in the effect ai explanations in algorithmic decision-making. *arXiv preprint arXiv:2502.12354*.
- Frederic Charles Bartlett. 1995. *Remembering: A study in experimental and social psychology*. Cambridge university press.
- Oana-Maria Camburu, Tim Rocktäschel, Thomas Lukasiewicz, and Phil Blunsom. 2018. e-snli: Natural language inference with natural language explanations. *Advances in Neural Information Processing Systems*, 31.
- Beiduo Chen, Siyao Peng, Anna Korhonen, and Barbara Plank. 2025. [A rose by any other name: LLM-generated explanations are good proxies for human explanations to collect label distributions on NLI](#). In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 10777–10802, Vienna, Austria. Association for Computational Linguistics.
- Zheng Chu, Huiming Fan, Jingchang Chen, Qianyu Wang, Mingda Yang, Jiafeng Liang, Zhongjie Wang, Hao Li, Guo Tang, Ming Liu, and Bing Qin. 2025. [Self-critique guided iterative reasoning for multi-hop question answering](#). In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 2415–2438, Vienna, Austria. Association for Computational Linguistics.
- Debela Gemechu, Ramon Ruiz-Dolz, Henrike Beyer, and Chris Reed. 2025. [Natural language reasoning in large language models: Analysis and evaluation](#). In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 3717–3741, Vienna, Austria. Association for Computational Linguistics.

- Xuanli He, Yuxiang Wu, Oana-Maria Camburu, Pasquale Minervini, and Pontus Stenetorp. 2024. [Using natural language explanations to improve robustness of in-context learning](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13477–13499, Bangkok, Thailand. Association for Computational Linguistics.
- Ukyo Honda and Tatsushi Oka. 2025. [Exploring explanations improves the robustness of in-context learning](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 23693–23714, Vienna, Austria. Association for Computational Linguistics.
- Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiwu Zheng, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, et al. 2023. [Metagtpt: Meta programming for a multi-agent collaborative framework](#). In *The Twelfth International Conference on Learning Representations*.
- Daniel Kahneman. 2011. *Thinking, fast and slow*. macmillan.
- T Karthikeyan, Om Dehlan, Mausam, and Manish Gupta. 2025. [LRPLAN: A multi-agent collaboration of large language and reasoning models for planning with implicit & explicit constraints](#). In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 8280–8310, Suzhou, China. Association for Computational Linguistics.
- Jenny Kunz and Marco Kuhlmann. 2024. [Properties and challenges of LLM-generated explanations](#). In *Proceedings of the Third Workshop on Bridging Human-Computer Interaction and Natural Language Processing*, pages 13–27, Mexico City, Mexico. Association for Computational Linguistics.
- Kenton Lee, Ming-Wei Chang, and Kristina Toutanova. 2019. [Latent retrieval for weakly supervised open domain question answering](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6086–6096, Florence, Italy. Association for Computational Linguistics.
- Yansi Li, Jiahao Xu, Tian Liang, Xingyu Chen, Zhiwei He, Qiuzhi Liu, Rui Wang, Zhuosheng Zhang, Zhaopeng Tu, Haitao Mi, et al. 2025. [Dancing with critiques: Enhancing llm reasoning with step-wise natural language self-critique](#). *arXiv preprint arXiv:2503.17363*.
- Tania Lombrozo. 2006. The structure and function of explanations. *Trends in cognitive sciences*, 10(10):464–470.
- Marlene Lutz, Indira Sen, Georg Ahnert, Elisa Rogers, and Markus Strohmaier. 2025. [The prompt makes the person\(a\): A systematic evaluation of sociodemographic persona prompting for large language models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 23212–23237, Suzhou, China. Association for Computational Linguistics.

697	Xiaochen Ma, Guozheng Rao, Lina Xu, Xin Wang, Zaiming Fan, and Zhe Zhang. 2025. Guided and knowledgeable multi-agent debate for fact verification. <i>Expert Systems with Applications</i> , page 130103.	EMNLP 2024, pages 16612–16631, Miami, Florida, USA. Association for Computational Linguistics.	753 754
701	Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. 2020. Adversarial NLI: A new benchmark for natural language understanding . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 4885–4901, Online. Association for Computational Linguistics.	Somin Wadhwa, Adit Krishnan, Runhui Wang, Byron C Wallace, and Luyang Kong. 2024. Learning from natural language explanations for generalizable entity matching . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 6114–6129, Miami, Florida, USA. Association for Computational Linguistics.	755 756 757 758 759 760 761
708	Letitia Parcalabescu and Anette Frank. 2024. On measuring faithfulness or self-consistency of natural language explanations . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 6048–6089, Bangkok, Thailand. Association for Computational Linguistics.	Qianli Wang, Tatiana Anikina, Nils Feldhus, Simon Ostermann, Sebastian Möller, and Vera Schmitt. 2025. Cross-refine: Improving natural language explanation generation by learning in tandem . In <i>Proceedings of the 31st International Conference on Computational Linguistics</i> , pages 1150–1167, Abu Dhabi, UAE. Association for Computational Linguistics.	762 763 764 765 766 767 768
715	Judea Pearl. 2009. <i>Causality</i> . Cambridge university press.	Yingming Wang and Pepa Atanasova. 2025. Self-critique and refinement for faithful natural language explanations . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 8492–8518, Suzhou, China. Association for Computational Linguistics.	770 771 772 773 774
717	Xin Quan, Marco Valentino, Louise A. Dennis, and Andre Freitas. 2025. Faithful and robust LLM-driven theorem proving for NLI explanations . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 17734–17755, Vienna, Austria. Association for Computational Linguistics.	Koki Wataoka, Tsubasa Takahashi, and Ryokan Ri. 2024. Self-preference bias in llm-as-a-judge. <i>arXiv preprint arXiv:2410.21819</i> .	775 776 777
724	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. <i>Advances in neural information processing systems</i> , 36:53728–53741.	Joonho Yang, Seunghyun Yoon, ByeongJeong Kim, and Hwanhee Lee. 2024. FIZZ: Factual inconsistency detection by zoom-in summary and zoom-out document . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 30–45, Miami, Florida, USA. Association for Computational Linguistics.	778 779 780 781 782 783 784
729	Ali Rahimi and Benjamin Recht. 2007. Random features for large-scale kernel machines. <i>Advances in neural information processing systems</i> , 20.	Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. HotpotQA: A dataset for diverse, explainable multi-hop question answering . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 2369–2380, Brussels, Belgium. Association for Computational Linguistics.	785 786 787 788 789 790 791 792
732	Nazneen Fatema Rajani, Bryan McCann, Caiming Xiong, and Richard Socher. 2019. Explain yourself! leveraging language models for commonsense reasoning . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4932–4942, Florence, Italy. Association for Computational Linguistics.	Bingsheng Yao, Ishan Jindal, Lucian Popa, Yannis Katsis, Sayan Ghosh, Lihong He, Yuxuan Lu, Shashank Srivastava, Yunyao Li, James Hendler, and Dakuo Wang. 2023. Beyond labels: Empowering human annotators with natural language explanations through a novel active-learning architecture . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 11629–11643, Singapore. Association for Computational Linguistics.	793 794 795 796 797 798 799 800 801
739	Alessandro Scirè, Karim Ghonim, and Roberto Navigli. 2024. FENICE: Factuality evaluation of summarization based on natural language inference and claim extraction . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 14148–14161, Bangkok, Thailand. Association for Computational Linguistics.	Huajian Zhang, Yumo Xu, and Laura Perez-Beltrachini. 2024a. Fine-grained natural language inference based faithfulness evaluation for diverse summarisation tasks . In <i>Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1701–1722, St. Julian’s, Malta. Association for Computational Linguistics.	802 803 804 805 806 807 808 809
746	James Surowiecki. 2005. The wisdom of crowds. <i>Surowiecki, J.</i>		
748	Yu-Min Tseng, Yu-Chao Huang, Teng-Yun Hsiao, Wei-Lin Chen, Chao-Wei Huang, Yu Meng, and Yun-Nung Chen. 2024. Two tales of persona in LLMs: A survey of role-playing and personalization . In <i>Findings of the Association for Computational Linguistics:</i>		

810 Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur
811 Szlám, Douwe Kiela, and Jason Weston. 2018. *Per-*
812 *sonalizing dialogue agents: I have a dog, do you have*
813 *pets too?* In *Proceedings of the 56th Annual Meeting*
814 *of the Association for Computational Linguistics (Vol-*
815 *ume 1: Long Papers)*, pages 2204–2213, Melbourne,
816 Australia. Association for Computational Linguistics.

817 Tianhui Zhang, Bei Peng, and Danushka Bollegala.
818 2025. *Evaluating the evaluation of diversity in com-*
819 *monsense generation.* In *Proceedings of the 63rd*
820 *Annual Meeting of the Association for Computational*
821 *Linguistics (Volume 1: Long Papers)*, pages 24258–
822 24275, Vienna, Austria. Association for Computa-
823 tional Linguistics.

824 Xiaoying Zhang, Baolin Peng, Ye Tian, Jingyan Zhou,
825 Lifeng Jin, Linfeng Song, Haitao Mi, and Helen
826 Meng. 2024b. *Self-alignment for factuality: Miti-*
827 *gating hallucinations in LLMs via self-evaluation.* In
828 *Proceedings of the 62nd Annual Meeting of the As-*
829 *sociation for Computational Linguistics (Volume 1:*
830 *Long Papers)*, pages 1946–1965, Bangkok, Thailand.
831 Association for Computational Linguistics.

832 Yaoming Zhu, Sidi Lu, Lei Zheng, Jiaxian Guo, Weinan
833 Zhang, Jun Wang, and Yong Yu. 2018. *Texygen: A*
834 *benchmarking platform for text generation models.*
835 In *The 41st international ACM SIGIR conference*
836 *on research & development in information retrieval*,
837 pages 1097–1100.

A Further Details in Self-Play Explainer 838

A.1 Theory-Grounded Persona Definition 839

- 840 • **System-2 Rationalist** suppresses intuitive 840
841 judgment (System 1) and promotes analytical 841
842 reasoning (System 2) to exhibit the model’s 842
843 internal logical capability. 843
- 844 • **Counterfactual Analyst** assumes alterna- 844
845 tive scenarios in which specific input elements 845
846 are altered, introducing causal directionality. 846
- 847 • **Schema-based Interpreter** applies schema 847
848 structures to interpret the input and identifies 848
849 conceptual mismatches within them. 849
- 850 • **Crowd-Wisdom Synthesizer** enumerates the 850
851 perspective of hypothetical experts and aggreg- 851
852 gates commonly supported evidence. 852

A.2 Multi-Factor Preference Scoring 853

854 **Definition of (*near-tie*) flips** We distinguish per- 854
855 sona selection based on either the target-class logit 855
856 alone or with the logit gap measure¹⁰. 856

$$k_\ell = \operatorname{argmax}_k \ell_{xy,k}^{(\text{ENT})}, \quad (15) \quad 857$$

$$k_{\text{gap}} = \operatorname{argmax}_k \text{gap}_{xy,k}^{(\text{ENT})}, \quad (16) \quad 858$$

859 We define a *near-tie flip* as follows. 859

$$i \in \mathcal{D}_{\text{NT}} \iff \begin{cases} k_\ell \neq k_{\text{gap}} \\ \ell_{xy,k_\ell}^{(\text{ENT})} - \ell_{xy,k_{\text{gap}}}^{(\text{ENT})} \leq \epsilon \\ \text{gap}_{xy,k_{\text{gap}}}^{(\text{ENT})} > \text{gap}_{xy,k_\ell}^{(\text{ENT})} \end{cases} \quad (17) \quad 860$$

861 This case occurs when a persona is not selected 861
862 as the top choice under a margin ϵ based solely on 862
863 the target-class logit, but becomes preferred once 863
864 the logit gap measure is applied. We set ϵ to the 864
865 lower 30% quantile of the top-2 target-class logit 865
866 margin distribution, where smaller margins indicate 866
867 ambiguity near the decision boundary. 867

868 **Grid search procedure** We evaluate 91 (α, β) 868
869 combinations by discretizing each pair over the 869
870 range $[-1, 1]$ with a step size of 0.25. For each 870
871 training dataset, we compute the *flip* rate and *near-*
872 *tie flip* rate among the *flipped* cases, and select the 872
873 best pair that best trades off these two conditions 873
874 by maximizing the f1 score under the selected ϵ , as 874
875 reported in Table 3. 875

¹⁰<https://huggingface.co/MoritzLaurer/DeBERTa-v3-base-xnli-multilingual-nli-2mil7>

Dataset	S_{PLA}		S_{SCF}	
	ϵ	(α, β)	ϵ	(α, β)
HotpotQA	0.409	(0.75, 0.75)	0.321	(0.50, 0.50)
NQopen	0.455	(1.00, 1.00)	0.314	(0.50, 0.50)
ANLI-R1	0.459	(0.75, 0.75)	0.344	(-0.25, -0.25)
ANLI-R2	0.457	(1.00, 1.00)	0.341	(0.75, 0.75)
ANLI-R3	0.476	(1.00, 1.00)	0.350	(0.75, 0.75)
eSNLI	0.525	(1.00, 1.00)	0.347	(0.50, 0.50)

Table 3: Dataset-specific metric hyperparameters for defining (*near-tie*) *flips*.

B Experimental Details

B.1 Dataset Curation

We merged and re-split the datasets into train, dev, and test sets with an 8:1:1 ratio using a fixed seed to ensure reproducibility. For ANLI, all rounds were aligned to the label distribution and sample size of R1. For eSNLI, we retained only the premise-hypothesis pair with the longest human-annotated explanation per pair and kept samples whose lengths fell within the top 25% of their distributions, indicating higher difficulty.

Dataset	#Train : #Dev : #Test
HotpotQA (Yang et al., 2018)	78,281 : 9,785 : 9,786
NQopen (Lee et al., 2019)	73,228 : 9,153 : 9,154
ANLI (Nie et al., 2020)	15,156 : 1,894 : 1,896
eSNLI (Camburu et al., 2018)	11,118 : 1,389 : 1,391

Table 4: Sample proportions for each dataset used in our experiments.

B.2 Training Configuration

We fine-tuned the base model using SFT with LoRA adapters, updating only adapter weights. Training was conducted in bfloat16 precision with a maximum length of 1024 tokens. We used a cosine lr schedule and a warmup ratio of 0.03, setting the lr to $5e-6$ for QA and $1e-6$ for NLI datasets. We used a per-device batch size of 2 with gradient accumulation steps of 64/32 for QA/NLI datasets, and a weight decay of 0.05. For DPO, we set the parameter β to 0.05, and kept other configurations identical to those of SFT. Both SFT and DPO were performed for one training epoch.

C Additional Experimental Results

C.1 Granularity Effects in Logical Alignment

We present the per-task detailed scores corresponding to Figure 2 in Table 5.

Dataset	Sentence-lev.	
	SPE (SFT)	SPE (DPO)
HotpotQA	17.29	25.90
NQopen	11.28	14.41
ANLI-R1	35.01	33.13
ANLI-R2	32.91	32.63
ANLI-R3	39.49	33.70
eSNLI	26.21	26.54
avg (stddev)	27.03 (0.099)	27.72 (0.067)

Dataset	Paragraph-lev.	
	SPE (SFT)	SPE (DPO)
HotpotQA	66.39	75.71
NQopen	67.32	70.63
ANLI-R1	47.96	47.22
ANLI-R2	47.70	50.02
ANLI-R3	48.35	54.48
eSNLI	32.29	68.02
avg (stddev)	51.67 (0.120)	61.01 (0.108)

Table 5: Logical alignment scores at the sentence and paragraph levels with mean (standard deviation) for each dataset using the SPE framework.

C.2 Persona Selection Dynamics in Best-of-Personas

We present the per-task detailed ratios corresponding to Figure 3 and 4, in Table 6.

Dataset	k=1	k=2	k=3	k=4
HotpotQA	25.84	14.36	14.00	4.71
NQopen	17.83	30.33	11.38	40.44
ANLI-R1	21.62	5.74	36.33	5.95
ANLI-R2	28.00	28.11	14.34	29.53
ANLI-R3	23.31	10.49	31.11	4.79
eSNLI	19.91	11.43	25.52	6.47

Table 6: Persona selection ratios for each task under the best-of-personas with logical alignment. k denotes the theoretical persona index.

C.3 Analysis of GPT-4 Win Rate

Although we obtain competitive win rates, GPT-based evaluation favors generic and concise outputs (Wataoka et al., 2024), underscoring the need for LLM-as-a-judge approaches that better assess logical reasoning across theoretical perspectives.

Dataset	HotpotQA	NQopen	ANLI-R1	ANLI-R2	ANLI-R3	eSNLI
win (%)	51.0	39.7	49.0	47.7	45.4	45.1

Table 7: GPT-4.1-mini win rates of the DPO compared to the SFT under the proposed framework.

D Prompt Constructions

914

Naive Reasoner Agent

You are the Naive Reasoner, an explanation-oriented agent that relies on your own internal knowledge and intuitive associations.

{shared system-instructions}

915

System-2 Rationalist Agent

You are the System-2 Rationalist, an explanation-oriented agent grounded in dual-process theories of human cognition.

Your role is to generate explanations for a task's prediction using slow, deliberate, and analytically structured reasoning.

Your explanations must reflect the characteristics of system-2 thinking: provide step-by-step justification, maintain logical coherence, and avoid relying on intuitive shortcuts into a coherent interpretation.

{shared system-instructions}

916

Counterfactual Analyst Agent

You are the Counterfactual Analyst, an explanation-oriented agent grounded in causal inference and counterfactual reasoning.

Your role is to generate explanations for a task's prediction by exploring how the outcome would change under alternative causal conditions.

Your explanations must reflect the characteristics of counterfactual thinking: focus on the causal factors behind the prediction, avoid irrelevant details, and use clear "what-if" contrasts into a coherent interpretation.

{shared system-instructions}

917

Schema-based Interpreter Agent

You are the Schema-based Interpreter, an explanation-oriented agent grounded in schema theory and frame semantics.

Your role is to generate explanations for a task's prediction by identifying the conceptual schema and explaining how it guides the meaning-making process.

Your explanations must reflect the characteristics of schema-based reasoning: highlight the underlying structure, map input elements to its schema components, and show how this structure yields a coherent interpretation.

{shared system-instructions}

918

Crowd-Wisdom Synthesizer Agent

You are the Crowd-Wisdom Synthesizer, an explanation-oriented agent grounded in principles of collective intelligence.

Your role is to generate explanations for a task's prediction by simulating multiple expert perspectives and combining their insights into a unified account.

Your explanations must reflect the characteristics of crowd-wisdom thinking: generate diverse hypothetical viewpoints, identify shared or recurring cues, and synthesize shared insights into a coherent interpretation.

{shared system-instructions}

919

Figure 7: System prompts used for our theory-grounded explanatory personas.

Naive Reasoner Agent

{task-specific instructions}

{shared user-instructions}

920

System-2 Rationalist Agent

{task-specific instructions}

Generate a System-2-style explanation that justifies the predicted output using slow, deliberate, analytically structured reasoning.

Use the following structure as your internal reasoning guideline (do not output):

<structure>

1. Identify the key input factors that are relevant to the output.

921

2. For each factor, provide a clear reasoning step that connects it to the output in a sequential manner.
3. After laying out the steps, arrive at intermediate conclusions and see how they progressively narrow toward the output.
4. Use these intermediate conclusions to justify why the output is the most reasonable result.
</structure>

{shared user-instructions}

Counterfactual Analyst Agent

{task-specific instructions}

Generate a Counterfactual-style explanation that justifies the predicted output by exploring how the outcome would change under alternative causal conditions.

Use the following structure as your internal reasoning guideline (do not output):

<structure>

1. Identify the key input factors that causally influence the output.
2. Explain how each factor contributes to or enables the output.
3. Present at least one counterfactual scenario showing how altering a key factor would change the output.
4. Use this contrast to justify why the output is the most reasonable result.

</structure>

{shared user-instructions}

Schema-based Interpreter Agent

{task-specific instructions}

Generate a Schema-based explanation that justifies the predicted output by identifying the conceptual schema and explaining how it guides the meaning-making process.

Use the following structure as your internal reasoning guideline (do not output):

<structure>

1. Identify the conceptual schema or frame that the input most naturally activates.
2. Map the relevant elements of the input to the appropriate components within this schema.
3. Explain how the relationships encoded in this schema shape the interpretation of the input and output.
4. Use this process to justify why the output is the most reasonable result.

</structure>

{shared user-instructions}

Crowd-Wisdom Synthesizer Agent

{task-specific instructions}

Generate a Crowd-wisdom-style explanation that justifies the predicted output by simulating multiple expert perspectives and combining their insights into a unified account.

Use the following structure as your internal reasoning guideline (do not output):

<structure>

1. Simulate several expert viewpoints, each offering a distinct interpretation of the input and output.
2. Identify the cues or reasoning patterns that consistently recur across these viewpoints.
3. Treat these recurring elements as collectively supported insights and merge them into a unified account.
4. Use this synthesis to justify why the output is the most reasonable result.

</structure>

{shared user-instructions}

Figure 8: User prompts used for our theory-grounded explanatory personas.

Self-Critique Reviewer Agent

system prompt

You are the Self-Factuality Critic, an explanation-evaluating agent that assesses the reasoning faithfulness of an explanation with respect to the given input-output relationship.

Your role is to critically evaluate the factuality and plausibility of the given model-generated explanation.

Your final response must contain exactly two components: a single scale from the three predefined scales, and a claim-level critique that justifies why the explanation fits that scale.

user prompt

{task-specific instructions}

Explanation: {explanation}

First, assign the given explanation to exactly one of the following three scales.

[Substantially unsupported] Most of the claims in the explanation do not faithfully reflect the relationship between the input and output, and several claims rely on hallucinated, incorrect, or contradictory information.

[Partially unsupported] Some claims in the explanation appropriately reflect aspects of the input-output relationship, but the explanation also contains speculative or partially hallucinated information.

[Fully supported] The claims in the explanation faithfully describe the relationship between the input and output, with no major factual inconsistencies or hallucinations.

Second, provide a claim-level critique that justifies why the selected scale is appropriate and specifies the grounds on which the explanation is being critically evaluated.

Now write your claim-level critique as a short paragraph that must not exceed five sentences.

Do not repeat or restate the input, output, or explanation, and do not use bullet points or numbered lists.

Output only the scale and the critique text.

Scale: <one of the three scales>

Claim-level critique: <a short paragraph>

926

Meta-Explanation Recomposer Agent

system prompt

You are the Meta-Explanation Recomposer, a constrained aggregation agent that synthesizes multiple explanations into a single, coherent, and faithful final explanation aligned with the given input-output relationship.

Your role is to rigorously analyze the provided explanations, identify mutually supported reasoning, eliminate contradictory or unsupported elements, and construct a refined explanation.

Your refined explanations must maximize reliability, interpretability, and human-readable naturalness while avoiding any hallucinated content.

If the candidate explanations conflict, you must explicitly resolve the conflict and adopt only one consistent interpretation, discarding the other.

user prompt

{task-specific instructions}

Synthesize these candidate explanations into a single final explanation that is faithful to the input-output relationship and strictly grounded in the reasoning contained within the candidates above.

Use the following structure as your internal reasoning guideline (do not output):

<structure>

1. Extract the key reasoning elements from candidate explanations.

2. Compare the reasoning to identify overlapping evidence, consistent elements, and any contradictions or incompatible interpretations between the candidates.

3. If conflicting interpretations exist, explicitly decide which interpretation is more coherent and better supported, and discard the other.

- Do NOT present the discarded interpretation as partially valid, plausible, or still relevant; it must be fully excluded from the refined explanation.

- Consistency is strictly prioritized over coverage.

4. Based on this analysis, generate one final recomposed explanation that supports exactly one coherent interpretation that is:

927

- faithful to the input-output relationship,
 - logically consistent and evidence-grounded,
 - reliable, interpretable, and natural to read,
 - free from hallucinated or speculative contents.
 5. When appropriate, you may reuse or adapt portions of the expressions from the candidate explanations for better refined explanation.
 </structure>
 {shared user-instructions}

Figure 9: System and user prompts used for the Self-Critique Reviewer and Meta-Explanation Recompiler agents.

Shared system-instructions
 Do NOT include any acknowledgement or meta text (e.g., "Okay", "I understand", "Sure", "I will", "Here is", etc.).
 Do NOT paraphrase, summarize, or quote the input or the output; refer to them only implicitly through abstract reasoning.
 Output only a single explanatory paragraph of at most five sentences, and nothing else.

Shared user-instructions
 Now write your {agentic}-style explanation as a short paragraph that must not exceed five sentences. Do not use bullet points or numbered lists; output only the explanation text.

 Write the explanation now. Start directly with the explanation content.

Task-specific instructions: Generating NLEs
 You are given a question and its appropriate answer. Assuming the answer is correct, generate a well-justified explanation that clearly articulates the relationship between the question and the answer.
 (or) You are given a premise and a hypothesis, along with their appropriate logical label. Assuming the label is correct, generate a well-justified explanation that clearly articulates the relationship between the premise-hypothesis pair and the label.

Task-specific instructions: Self-Critique Agent
 You are given a question and its appropriate answer. Your task is to evaluate whether the claims in the explanation faithfully and sufficiently justify the relationship between the input and output.
 (or) You are given a premise and a hypothesis, along with their appropriate logical label. Your task is to evaluate whether the claims in the explanation faithfully and sufficiently justify the relationship between the input and output.

Task-specific instructions: Meta-Explanation Recompiler Agent
 You are given a question and its appropriate answer. Among multiple generated explanations, two priority candidates have been selected according to their relevance to the given input-output relationship.
 (or) You are given a premise and a hypothesis, along with their appropriate logical label. Among multiple generated explanations, two priority candidates have been selected according to their relevance to the given input-output relationship.

Shared format
 Input: {input}
 Output: {output}

Figure 10: Shared system and user instructions, along with task-specific instructions. These are used for all multi-agents in our study.