



REFLEX: Self-Refining Explainable Fact-Checking via Verdict-Anchored Style Control

Chuyi Kong¹, Wei Gao², Jing Ma^{1*}, Hongzhan Lin¹, Yuxi Sun¹

¹ Hong Kong Baptist University ² Singapore Management University
kongchuyi@life.hkbu.edu.hk weigao@smu.edu.sg majing@hkbu.edu.hk

Abstract

The prevalence of fake news on social media calls for automated fact-checking systems that deliver not only accurate verdicts but also faithful explanations. However, existing large language model (LLM)-based methods often overlook deceptive misinformation styles in generated explanations, producing unfaithful rationales that may mislead human judgment. They also rely heavily on external knowledge sources, which can introduce hallucinations and incur substantial latency, undermining both reliability and responsiveness in real-time settings. To address these limitations, we propose **RE**ason-guided **FA**ct-checking with **LA**tent **EX**planations (REFLEX), a self-refining framework that explicitly controls reasoning style by anchoring explanations to the predicted verdict. REFLEX leverages self-disagreement veracity signals between a backbone model and its fine-tuned variant to construct steering vectors, thereby naturally disentangling factual content from stylistic cues. Experiments on a real-world benchmark show that REFLEX achieves state-of-the-art performance with LLaMA-series models using only 465 self-refined samples. Owing to its transferability, REFLEX also yields gains of up to 7.54 Macro-F1 points on in-the-wild data. Further analysis shows that our method effectively mitigates faithful hallucination, leading to both more reliable explanations and more accurate verdicts than prior explainable fact-checking approaches.

1 Introduction

The rapid spread of fake news on social media poses a major challenge for automatic fact-checking (AFC) (Guo et al., 2022). Recent AFC systems increasingly leverage Large Language Models (LLMs) to predict verdicts and generate explanations. Yet they often overlook a core requirement for human-facing fact-checking: explanations

*The corresponding author. Code is available on [GitHub](#).

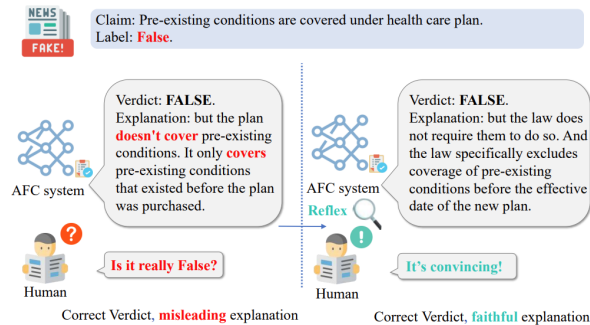


Figure 1: A motivating case showing how deceptive explanatory style can mislead human judgment.

must be not only plausible, but also faithful. As shown in Figure 1, even when an LLM predicts the correct verdict, an unfaithful explanation can still mislead human judgment.

Such failures often stem from deceptive explanatory styles in LLM-generated rationales (Turpin et al., 2023; Chen and Shu, 2023), a phenomenon related to faithfulness hallucination (Huang et al., 2025a). These rationales may begin with factual statements but gradually drift into verdict inconsistency, internal contradiction, or overconfident fabrication. We therefore argue that effective explainable fact-checking requires disentangling *fact* from deceptive *style* in explanations. In the worst case, such unfaithful explanations may even cascade into incorrect verdicts, since LLMs appear to construct holistic *internal* conceptual representations before producing rationalized outputs, rather than simply mirroring next-token prediction (Anthropic, 2025).

Existing AFC methods such as HiSS (Zhang and Gao, 2023) and L-Defenses (Wang et al., 2024) mainly address these issues through *external* interventions. While such methods reduce costly manual involvement, their dependence on external closed-source APIs or search engines introduces additional latency and may inject misinformation or hallucinations (Gekhman et al., 2024) into fine-grained explanations. This compromises both re-

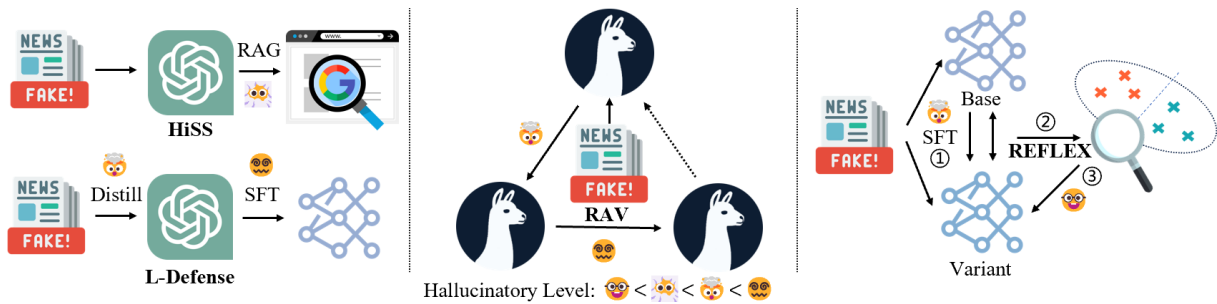


Figure 2: Comparison between REFLEX and representative baselines. Hallucination levels are indicated by emojis. The dashed line denotes an optional process. Bidirectional arrows in REFLEX denote **self-disagreement** signals.

liability and responsiveness, which are critical in real-world fact-checking.

To address these limitations, we improve explanation faithfulness by steering *internal* model activations. Our approach exploits self-disagreement veracity signals across the model training lifecycle, from a backbone model to its post-trained variants, to disentangle transferable Steering Vectors (SVs) (Han et al., 2024) into Inference Vectors (IVs) and Knowledge Vectors (KVs). IVs capture style-sensitive veracity gains that emerge after fine-tuning. Under limited supervision, fine-tuning is often more effective at shaping task-specific reasoning styles and activating latent backbone knowledge (Li et al., 2023b) than at injecting new factual knowledge (Ghosal et al., 2024; Ren et al., 2024; Zhao et al., 2025). Amplifying these signals therefore helps align explanations with verdicts. In contrast, KVs capture fact-sensitive veracity drops caused by knowledge conflicts (Gekhman et al., 2024) introduced during fine-tuning, which can induce alignment tax (Leike, 2022) and hallucinations. Suppressing these signals helps preserve consistent factual representations. Based on this insight, we adaptively activate the appropriate vector according to probability gaps and use it to refine explanations, yielding our self-refining AFC paradigm, namely, **RE**ason-guided **FA**ct-checking with **LA**tent **EX**planations (REFLEX).

Experiments show that our internal-activation-based approach outperforms methods that rely on external resources, either in verdict accuracy alone or jointly in verdict accuracy and explanation quality with more faithful and concise reasoning. Moreover, REFLEX generalizes across backbones, datasets, and limited-sample settings, demonstrating strong transferability (up to 7.54 Macro-F1 points), flexibility (up to 7.57 Macro-F1 points), and data efficiency. We further find that evidence-free training is better suited to our paradigm, as external evidence can introduce noise

and amplify hallucinations. In addition, we observe a divergence amplification effect: larger directional discrepancies between KVs and IVs yield more effective disentanglement and stronger gains. Post-hoc analysis further shows that explanations produced by REFLEX are more logically coherent and less misleading, thereby mitigating faithfulness hallucination and improving verdict accuracy.

Overall, our main contributions are as follows:

- We propose REFLEX, a simple yet effective self-refining paradigm that disentangles fact from style, improves explanation faithfulness, and mitigates misinformation.
- REFLEX achieves state-of-the-art performance on a real-world dataset using only a small number of self-refined training samples, while producing high-quality explanations.
- We show that REFLEX supports cross-model transfer and exhibits performance-adaptive transferability, which can reduce response latency in fact-checking.
- We identify a divergence amplification effect in REFLEX, where greater directional divergence yields larger performance gains.

2 Background

Faithfulness in NLP. In NLP, faithfulness is typically studied from two perspectives. The first concerns consistency between model predictions and the reasoning that supports them (Jacovi and Goldberg, 2020; Atanasova et al., 2023; Parcalabescu and Frank, 2024). However, recent work (Yan et al., 2025) suggests that probabilistic models may not capture genuine reasoning, but instead rely on superficial patterns, i.e., *style*. We therefore use human-written explanations as learning targets to better approximate human reasoning style. The second concerns consistency with the surrounding

context, often referred to as *faithfulness hallucination* (Huang et al., 2025a), including instruction, context, and logical inconsistency. Our work follows this second line. Recent fact-checking work (Kim et al., 2024) addresses explanation–evidence alignment with multi-agent systems; in contrast, we adopt an evidence-free single-model setting to avoid hallucinations introduced by external evidence and iterative interaction.

Explainable Fact-Checking With the rise of LLMs, recent work has explored automatic explanation mechanisms to reduce human annotation effort in explainable fact-checking. As shown in Figure 2, HiSS (Zhang and Gao, 2023) decomposes complex claims into atomic claims and verifies them with RAG (Lewis et al., 2020), using retrieved reasoning trajectories as explanations. L-Defense (Wang et al., 2024) distills adversarial evidence and explanations from teacher models. RAV (Shukla et al., 2025) builds a fully open-source multi-agent dialogue system for claim verification and treats LLM responses as explanations. Despite their promise, these methods face several **key challenges**. RAG-based methods may introduce hallucinations at both retrieval and generation stages (Qian et al., 2023; Hu et al., 2025). Distillation and fine-tuning can amplify the snowballing effect (Leike, 2022) of teacher-model hallucinations. Multi-agent systems are vulnerable to hallucination propagation (Jain et al., 2025; Liang et al., 2026) and incur additional interaction latency, making them less suitable for time-sensitive fact-checking. To address these limitations, REFLEX introduces a self-refining and transferable paradigm that mitigates hallucinations in explanations, reduces inference latency, and improves the interpretability of internal model behavior.

Style Control in Fact-Checking Several studies improve fact-checking by exploiting stylistic differences across claim sources (e.g., human vs. machine) to improve verdict prediction (Pérez-Rosas et al., 2018; Rashkin et al., 2017). However, humans can deliberately manipulate style to deceive, which weakens the reliability of such signals (Schuster et al., 2020; Wu et al., 2024). In contrast, we anchor style control on verdict prediction, steering the model toward less deceptive explanations that remain aligned with the verdict.

Our approach builds on a broader line of work on style control, ranging from controllable text generation (Dathathri et al., 2019; Li et al., 2022)

to activation editing (Li et al., 2023a; Hernandez et al., 2023). Recent steering-vector methods also follow a contrastive paradigm. ITI (Li et al., 2023b) identifies truth-related directions at inference time, while CCS (Burns et al., 2023) and CAA (Rimsky et al., 2024) learn contrastive directions associated with truthfulness or stylistic behavior. However, in explainable fact-checking, **style and fact are entangled**: longer explanations and their associated verdicts may diverge due to attention dispersion (Wei et al., 2025), making strict contrastive pairing infeasible. To address this issue, we refine explanations through cross-stage sample selection and by decomposing steering vectors into KVs and IVs. To avoid shortcut learning over factual signals alone, we further adopt label-free interventions.

3 Methodology

3.1 Task Formulation

Given a dataset $\mathcal{D} = \{(c, evi, v, exp)_i\}_{i=1}^N$, where c denotes a claim, evi optional evidence, v the gold veracity label, and exp a human-written explanation, REFLEX aims to generate verdicts and explanations such that the explanation remains faithful to the predicted verdict, thereby mitigating faithfulness hallucination.

As shown in Figure 3, REFLEX consists of three stages: (1) **Dialogue-style Fact-Checker Training**, which reformulates fact-checking as a dialogue task and trains LLMs to jointly generate verdicts and explanations; (2) **Adaptive Sample Selection**, which identifies self-disagreement samples between a backbone model and its fine-tuned variants; (3) **Self-Explanation Guided Steering (S-EGS)**, which disentangles steering vectors into KVs and IVs from these samples and explicitly aligns explanations with verdicts.

3.2 Dialogue-style Fact-Checker Training

Data Preprocessing. LLM backbones already encode substantial factual knowledge (Li et al., 2023b), and fine-tuning with limited data often mainly activates this knowledge while adapting the model to the downstream task (Berglund et al., 2023; Ghosal et al., 2024; Ren et al., 2024). We therefore reformulate conventional document-style supervision as a single-turn QA-style dialogue to better activate backbone knowledge. Prior work shows that QA-style supervision yields stronger knowledge generalization during fine-tuning, whereas document-style supervision (e.g.,

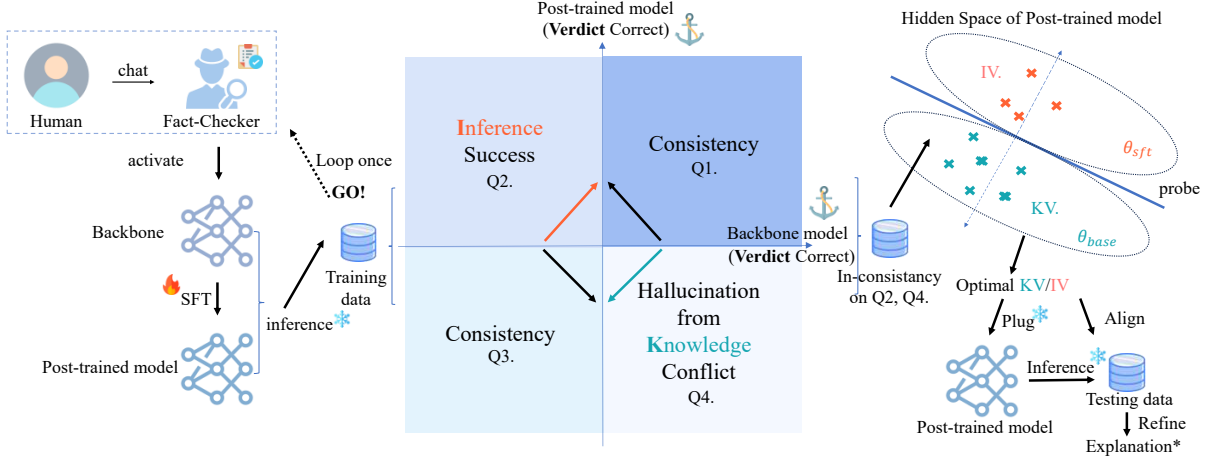


Figure 3: Overview of the three-stage REFLEX paradigm. Red text denotes **reasoning style** acquired during post-training, and green text denotes **distorted factual knowledge** in the backbone model.

Wikipedia text) yields weaker generalization (Zhao et al., 2025).

Training Protocol. Following prior dialogue-based training, we optimize a standard cross-entropy objective to jointly predict verdicts and explanations. For a comprehensive study, we consider four input-output configurations:

$$\begin{aligned} x = [c] \rightarrow y = [v], & \quad x = [c] \rightarrow y = [v; exp], \\ x = [c; evi] \rightarrow y = [v], & \quad x = [c; evi] \rightarrow y = [v; exp]. \end{aligned}$$

To encourage coherent reasoning, we incorporate Chain-of-Thought (CoT) (Wei et al., 2022), prompting the model to output its reasoning path as the explanation, which has been shown to improve performance (Lippmann and Yang, 2025). We further adopt role-play prompting to strengthen reasoning ability (Kong et al., 2024). Prompt templates and training details are provided in Appendix B and C.

3.3 Adaptive Sample Selection

Self-Knowledge Extraction. After training, we run inference on the training set using both the backbone and post-trained models. Since the backbone lacks instruction-following ability (Liu et al., 2022b), we use few-shot prompting to elicit its knowledge. To avoid data leakage, few-shot examples are drawn either from another training split or from a held-out validation set. To mitigate majority-label bias in few-shot learning, we balance labels when filling the maximum context window. For reproducibility, we fix the temperature to zero.

Cross-Stage Sample Selection. Given veracity predictions \hat{v}^{base} and \hat{v}^{sft} on the training set, we categorize samples relative to the gold label v^{gold}

as follows:

$$\begin{aligned} \text{Quadrant 2: } \hat{v}^{\text{base}} \neq v^{\text{gold}}, \hat{v}^{\text{sft}} = v^{\text{gold}} \\ \Rightarrow \text{Reasoning Gain,} \end{aligned}$$

$$\begin{aligned} \text{Quadrant 4: } \hat{v}^{\text{base}} = v^{\text{gold}}, \hat{v}^{\text{sft}} \neq v^{\text{gold}} \\ \Rightarrow \text{Knowledge Loss.} \end{aligned}$$

As shown in Figure 3, samples in Quadrant 2 (Q2) capture cases in which post-training corrects the backbone’s verdict. We treat these cases as **reasoning gains**, because under limited supervision, fine-tuning is more likely to activate latent backbone knowledge and shape task-specific reasoning style than to inject substantial new factual knowledge (Li et al., 2023b; Berglund et al., 2023; Ghosal et al., 2024; Ren et al., 2024). In other words, when the post-trained model succeeds where the backbone fails, we attribute the improvement primarily to better reasoning behavior rather than to newly acquired facts. By contrast, samples in Quadrant 4 (Q4) capture cases in which post-training flips a previously correct backbone verdict into an incorrect one. We interpret these cases as **knowledge loss**, since the degradation is more plausibly explained by conflicting knowledge or distorted factual representations introduced during post-training, which can lead to alignment tax (Leike, 2022) and hallucinations (Gekhman et al., 2024; Huang et al., 2025b). We therefore adaptively select samples from Q2 and Q4 for subsequent steering.

3.4 Self-Explanation Guided Steering

Disentangled Steering. Unlike ITI (Li et al., 2023b) and CAA (Rimsky et al., 2024), which control a single factual or stylistic direction using

explicit contrastive pairs and yield one steering vector, we decompose steering vectors into KVs and IVs across the quadrants above to disentangle fact and style. Specifically, we amplify style-sensitive signals by steering the model toward the positive IV direction, thereby aligning explanations with verdicts. Conversely, we suppress fact-sensitive signals by steering the model away from negative KV directions, thereby preserving consistent factual representations. We select the optimal KV or IV by maximizing the probability gap between unsteered and steered outputs, which we find to correlate with verdict accuracy. For each candidate direction, we measure how much steering increases the model’s probability assigned to the gold verdict relative to the unsteered baseline, and select the direction with the largest improvement (details in Appendix D). In addition, because entangled or inconsistent factual and stylistic signals weaken contrastive supervision, we adopt label-free interventions that do not rely on explicitly specified factual directions, enabling more robust in-the-wild steering. To improve conceptual separability while keeping computational overhead low, we use logistic probes and intervene only at decoder blocks. Implementation details are provided in Appendix C.

Explanation Refinement. To further align explanations with verdicts, we quantify the alignment between each token’s hidden representation and the optimal direction using cosine similarity. For layer l and token t , let h denote the hidden state and s the optimal vector. We define the alignment score as $a_{l,t} = \frac{h_{l,t} \cdot s_l}{\|h_{l,t}\| \|s_l\|}$. Positive scores indicate alignment with the optimal direction, while negative ones indicate the opposite. Segments with dense negative scores are often associated with redundant or noisy stylistic content. We therefore suppress them using the lightweight Ratcliff-Obershelp pattern-matching algorithm (Ratcliff and Metzener, 1988) to improve explanation readability.

4 Experiments and Results

Datasets. To better reflect real-world fact-checking and reduce hallucination risk, we use three datasets whose claims originate from professional fact-checking platforms and whose explanations are human-written: RAW-FC (Yang et al., 2022), LIAR-RAW (Yang et al., 2022), and AVeriTeC (Schlichtkrull et al., 2023). In RAW-FC and LIAR-RAW, explanations directly justify the claim label, whereas in AVeriTeC they justify both

the claim and its supporting evidence. For simplicity, we refer to all such rationales as explanations (details in Appendix E).

Metrics. For verdict prediction, we report Precision, Recall, and Macro-F1. For explanation quality, we conduct both automatic and manual evaluations. **(1) Automatic evaluation.** Following Wang et al. (2024), we use ChatGPT as a judge (Gu et al., 2025) along four dimensions: **Misleadingness**, **Informativeness**, **Soundness**, and **Readability**. Misleadingness measures whether an explanation is consistent with the veracity label, while Soundness measures whether the explanation is logically valid and grounded; these two dimensions most directly reflect context and logical consistency under faithfulness hallucination. The remaining dimensions are described in Appendix F. **(2) Manual evaluation.** For baselines, we follow Wang et al. (2024) and conduct point-wise evaluation on a subset of RAW-FC. For ablations, we perform three rounds of pairwise evaluation on misleadingness over all samples in Section 5.3.

4.1 Baseline Trials

4.1.1 Baselines

We compare against two categories of methods: **(1) External Dependency Design:** ChatGPT (OpenAI, 2023), HiSS (Zhang and Gao, 2023), FactLLaMA (Cheung and Lam, 2023), L-Defense (Wang et al., 2024), and RAV (Shukla et al., 2025); **(2) Internal Control Design:** REFLEX (ours). Among these, FactLLaMA also uses the Google API to retrieve evidence and fine-tunes LLaMA-2-7B. To remain comparable to most baselines, we adopt LLaMA-2-7B as the backbone in this section.

4.1.2 Results

For **verdict prediction**, Table 1 shows that our method achieves state-of-the-art performance on RAW-FC without relying on closed-source APIs, retrieved evidence, or large-scale multi-agent systems. After Stage 1 (SFT), it surpasses ChatGPT by 16.16 Macro-F1 points and HiSS by 6.69 points. Compared with FactLLaMA, which uses the same backbone but lacks dialogue-style supervision and full-parameter tuning, our model gains 4.94 points. It performs comparably to RAV while using only a single model, and remains slightly below L-Defense. After applying S-EGS, our method surpasses RAV by 5.80 points and L-Defense by 3.79 points, despite using only 465 self-extracted

Model	$x \rightarrow y$	Model	No Closed-Source API Dependency / # Distilled Explanations for Training	RAW-FC			LIAR-RAW		
				P	R	macF1	P	R	macF1
External Dependency Design									
ChatGPT	$c \rightarrow v$; exp	–	✘/–	47.72	48.62	44.43	25.41	27.33	25.11
	c , $evi \rightarrow v$; exp	–	✘/–	39.48	45.07	39.31	29.64	23.57	21.90
HiSS _{Google}	$c \rightarrow v$; exp	ChatGPT	✘/–	53.40	54.50	53.90	46.80	31.30	37.50
FactLLaMA _{Google}	c ; $evi \rightarrow v$	LLaMA2-7B	✘/–	56.11	55.50	55.65	32.46	32.05	30.44
L-Defense	c ; $evi \rightarrow v$; exp	ChatGPT + Roberta-Large	✘/ 32,240	61.72	61.01	61.20	30.55	32.20	30.53
RAV	c ; $evi \rightarrow v$; exp	LLaMA-3.1-70B-Instruct $\times 3$	✔/–	–	–	59.19	–	–	25.40
Internal Control Design (Ours)									
SFT	$c \rightarrow v$; exp	LLaMA2-7B	✔/ 0	60.66	61.04	60.59	48.38	46.83	43.05
S-EGS	$c \rightarrow v$; exp	LLaMA2-7B	✔/ 465	65.04	65.01	64.99	52.45	50.39	50.59

Table 1: Performance comparison on RAW-FC and LIAR-RAW. “Google” indicates the use of the Google Search API. Full baseline details are provided in Appendix G.

	RAW-FC				LIAR-RAW			
	M \downarrow	I	S	R	M \downarrow	I	S	R
Oracle - skyline	1.52	4.46	4.73	4.72	1.85	4.44	4.60	4.69
ChatGPT _{w/evi}	2.07	4.44	4.62	4.69	2.29	3.71	4.04	3.99
ChatGPT _{w/o evi}	1.97	4.00	4.44	4.68	2.27	3.93	4.29	4.50
L-Defense _{LLaMA2}	1.95	4.44	4.67	4.62	2.20	4.39	4.64	4.63
L-Defense _{ChatGPT}	1.91	4.17	4.41	4.49	2.06	4.12	4.28	4.47
Ours								
SFT	1.90	4.78	4.82	4.55	1.90	4.48	4.60	4.65
S-EGS _{LLaMA2}	1.79	4.88	4.83	4.80	1.77	4.58	4.66	4.83

Table 2: Automatic evaluation of explanation quality. The best REFLEX variant with LLaMA-2 is reported.

samples versus 32,240 GPT-3.5-distilled explanations in L-Defense. This highlights the strong data efficiency of REFLEX. Since our models are trained under the three-way label setting used by explanation-based methods, we do not report a direct verdict comparison on LIAR-RAW here.

For **explanation quality**, we evaluate open-source baselines that generate explanations. Following Wang et al. (2024), we map the six LIAR-RAW labels to three for all applicable baselines, and include an Oracle setting (ChatGPT with gold claims and verdicts) as a skyline. As shown in Table 2, after SFT our model achieves the best performance on RAW-FC in misleadingness, informativeness, and soundness, trailing only L-Defense (ChatGPT-distilled) in readability. On LIAR-RAW, it achieves the best informativeness and readability, and ranks second to L-Defense (LLaMA-2) in soundness. Applying S-EGS further improves all metrics consistently. Human evaluation (Appendix I) confirms the same trend. Because LLM-as-a-Judge may exhibit length bias (Gu et al., 2025), we also report explanation lengths in Appendix H. Under the same backbone, our explanations are shorter than those of L-Defense on RAW-FC, and shorter than all baselines, including Oracle, on LIAR-RAW, suggesting that REFLEX produces

explanations that are both concise and faithful.

4.2 Ablation Studies

We conduct the following ablation experiments on top of SFT models as baselines.

4.2.1 On Backbone Models

To assess generalization across backbones, in addition to LLaMA-2-7B, we also train and evaluate Qwen-3-8B (Yang et al., 2025).

As shown in Table 3, S-EGS improves over SFT on all three datasets in most settings, with gains of up to 5.03 Macro-F1 points, demonstrating strong generalization across backbones. The two marginal drops are likely due to severe **recency bias** in few-shot prompting (Appendix J), despite existing mitigation strategies (Lu et al., 2022; Min et al., 2022; Nguyen and Wong, 2023). Two patterns are noteworthy: (1) LLaMA-2 benefits more from S-EGS than Qwen-3; (2) AVeriTeC shows the smallest gains, despite the strongest SFT performance and the lowest hallucination rate (Appendix K). We attribute this **saturation phenomenon** to its QA-style formulation, which already aligns well with our training paradigm (Section 3.2).

Two additional trends emerge from the training configurations. (1) **Using explanations as targets is consistently beneficial.** At the SFT stage, explanation-based supervision improves performance across settings. In contrast, at the BASE stage, LLaMA-2 performs better without explanations, likely because its pretraining already strongly emphasizes factual knowledge (Touvron et al., 2023). (2) **Evidence conditioning is often harmful.** At the BASE stage, adding evidence generally degrades performance, especially for LLaMA-2, which is consistent with its weaker intrinsic rea-

Backbone	Stage	$x \rightarrow y$	Raw-FC	Δ mac-F1	LIAR-RAW	Δ mac-F1	AVeriTeC	Δ mac-F1
	BASE	$c \rightarrow v$	35.61	–	29.26	–	–	–
		$c; \text{evi} \rightarrow v$	27.08	–	16.97	–	28.18	–
		$c \rightarrow v; \text{exp}_{\text{cross}}$	34.41	–	12.48	–	–	–
		$c[; \text{evi}] \rightarrow v; \text{exp}_{\text{self}}$	31.68	–	35.80	–	27.70	–
LLaMA-2	SFT	$c \rightarrow v$	26.44	-9.17	37.23	+7.97	–	–
		$c; \text{evi} \rightarrow v$	44.85	+17.77	40.21	+23.24	75.91	+47.73
		$c[; \text{evi}] \rightarrow v; \text{exp}$	60.59	+26.18	43.05	+7.25	84.62	+56.92
	S-EGS	$c \rightarrow v$	31.47	+5.03	38.65	+1.42	–	–
		$c \rightarrow v; \text{exp}_{\text{cross}}$	64.99	+4.40	42.77	<u>-0.28</u>	–	–
		$c[; \text{evi}] \rightarrow v; \text{exp}_{\text{self}}$	61.81	+1.22	43.06	+0.01	84.61	-0.01
	BASE	$c \rightarrow v$	46.54	–	37.63	–	–	–
		$c; \text{evi} \rightarrow v$	46.23	–	41.30	–	66.14	–
		$c \rightarrow v; \text{exp}_{\text{cross}}$	46.66	–	42.25	–	–	–
		$c[; \text{evi}] \rightarrow v; \text{exp}_{\text{self}}$	48.86	–	39.16	–	66.02	–
Qwen-3	SFT	$c \rightarrow v$	41.67	-4.87	41.72	+4.09	–	–
		$c; \text{evi} \rightarrow v$	63.17	+16.94	42.29	+1.69	85.52	+19.38
		$c[; \text{evi}] \rightarrow v; \text{exp}$	58.35	+9.49	46.73	+4.48	88.02	+22.22
	S-EGS	$c \rightarrow v$	41.69	+0.02	41.73	+0.01	–	–
		$c \rightarrow v; \text{exp}_{\text{cross}}$	59.39	+1.04	47.13	+0.40	–	–
		$c[; \text{evi}] \rightarrow v; \text{exp}_{\text{self}}$	58.86	+0.51	46.53	<u>-0.20</u>	88.21	+0.19

Table 3: Macro-F1 of S-EGS across backbones, training stages, and datasets. *cross* denotes few-shot examples drawn from another dataset, while *self* denotes examples drawn from the model’s own validation set. [;] indicates optional evidence inputs, used only for AVeriTeC since its explanations are evidence-grounded. For SFT, Δ Macro-F1 is computed against the best BASE configuration.

soning ability (Gandhi et al., 2025). At the SFT stage, excluding AVeriTeC because its explanations are inherently evidence-grounded, we find that although evidence can help under the $c \rightarrow v$ setting, models trained with explanations but without evidence ($c \rightarrow v; \text{exp}$) generally outperform models trained with evidence but without explanations ($c; \text{evi} \rightarrow v$), except for Qwen-3 on RAW-FC. We conjecture that evidence introduces noise that can amplify hallucinations, while Qwen-3’s stronger reasoning ability partially compensates for this effect (Gandhi et al., 2025). Based on these findings, we use the **evidence-free** explanation-based objective ($c \rightarrow v; \text{exp}$) for RAW-FC and LIAR-RAW in subsequent experiments.

4.2.2 On Cross-model Transfer

Although LLMs often overfit across domains, and some prior work, e.g., LIMA (Zhou et al., 2023) and Character-LLM (Shao et al., 2023), even exploits deliberate overfitting for broader adaptation, our goal is to avoid overfitting and preserve scientific validity under realistic transfer. We therefore conduct cross-model steering using optimal vectors learned on LIAR-RAW and RAW-FC, excluding AVeriTeC because its dialogue format differs

Model	Transfer	Source		Target	
		F	Δ	F	Δ
LLaMA-2	R \rightarrow L	<u>64.99</u>	+4.40	50.59	+7.54
	L \rightarrow R	<u>43.06</u>	+0.01	47.20	-13.39
Qwen-3	R \rightarrow L	<u>59.39</u>	+1.04	45.12	-1.61
	L \rightarrow R	<u>47.13</u>	+0.40	41.04	-17.31

Table 4: Cross-dataset transfer results (Macro-F1). R: RAW-FC. L: LIAR-RAW. Δ : gain over SFT. Corr(Δ)=0.89; Corr(Source F, Target Δ)=0.95 (Pearson). See Appendix M for full results.

substantially from the other two. We also remove overlapping few-shot samples to avoid leakage.

As shown in Table 4, we observe two key patterns. (1) Transfer performance depends strongly on source model strength: the Pearson correlation between source Macro-F1 and target Δ Macro-F1 is 0.95. Strong vectors extracted from a better-performing source model (source Macro-F1 = 64.99) can substantially improve a weaker target setting (source Macro-F1 = 43.06) by up to 7.54 points, whereas weak vectors extracted from a weaker model (source Macro-F1 = 47.13) can degrade a stronger target setting (source Macro-F1 = 59.39) by up to 17.31 points. This demonstrates **performance-adaptive transferability**. (2)

Backbone	Direction	Objective	Raw-FC	LIAR-RAW	AVeriTeC
LLaMA-2	Vertical	w/o exp	34.01 (↑7.57)	38.37 (↑1.14)	74.86 (↓1.05)
		w/ exp	62.17 (↑1.58)	43.61 (↑0.56)	85.86 (↑1.24)
	Horizontal	w/o exp	34.82 (↑8.38)	38.37 (↑1.14)	74.79 (↓1.12)
		w/ exp	62.64 (↑2.05)	43.73 (↑0.68)	82.92 (↓1.70)
Qwen-3	Vertical	w/o exp	41.69 (↑0.02)	41.91 (↑0.19)	85.71 (↑0.19)
		w/ exp	58.88 (↑0.53)	46.80 (↑0.07)	88.62 (↑0.60)
	Horizontal	w/o exp	42.10 (↑0.43)	41.76 (↑0.04)	85.89 (↑0.37)
		w/ exp	58.32 (↓0.03)	47.04 (↑0.31)	88.91 (↑0.89)

Table 5: Macro-F1 gains over SFT under different pairing strategies and training objectives.

Compared with LLaMA-2, Qwen-3 exhibits much weaker transferability, yielding only limited gains under cross-dataset steering.

4.2.3 On Pairs Combinations

To further test the flexibility of REFLEX, we construct inconsistent training pairs across objectives to induce more explicit disentanglement.

Vertical steering pairs BASE outputs without explanations ($c \rightarrow v_{\text{base}}$) and with explanation-augmented SFT outputs ($c \rightarrow v; \text{exp}_{\text{sft}}$). **Horizontal steering** pairs SFT outputs with and without explanations ($c \rightarrow v_{\text{sft}}$ vs. $c \rightarrow v; \text{exp}_{\text{sft}}$), independently of the standard REFLEX setup. For clarity, we omit evidence inputs for AVeriTeC hereafter.

As shown in Table 5, both steering axes improve performance on all datasets except AVeriTeC, demonstrating the flexibility of the framework. Notably, explanation-guided vectors ($c \rightarrow v; \text{exp}$) can effectively steer explanation-free outputs ($c \rightarrow v$), yielding gains of up to 8.38 Macro-F1 points. This suggests that explanations can function as **internal activation signals** for reasoning. The slight drop (> 0.1) on AVeriTeC is again likely due to saturation (Section 4.2.1). Moreover, horizontal steering, which is derived solely from the SFT model, yields two drops on AVeriTeC, whereas vertical steering within REFLEX yields only one, further indicating the robustness of REFLEX.

5 Analysis

To better understand how REFLEX improves verdict prediction and explanation quality, we analyze explanation quality across all improved variants and conduct additional model-internal studies.

Vector	Split	Factual Metric	Stylistic Metrics	
		M↓	I	S
KV	Total	1.58	4.56	4.72
	- Correct	<u>1.52</u>	4.56	4.74
	- Error	<u>1.74</u>	4.54	4.67
IV	Total	1.90	4.86	4.79
	- Correct	1.89	4.85	4.79
	- Error	1.91	4.86	4.79

Table 6: Average scores of explanations guided by the optimal KVs and IVs.

5.1 Overall Evaluation

As shown in Table 22 (Appendix R), automatic pointwise evaluation shows that steering generally shifts explanations toward lower misleadingness and higher informativeness, soundness, and readability. These trends suggest that REFLEX improves explanation faithfulness, especially in M and S, while also improving **stylistic plausibility** (Agarwal et al., 2024) across I, S, and R (see Appendix A).

Because automatic evaluation may be unstable, we additionally conduct three rounds of pairwise human evaluation on the most factuality-relevant dimension, misleadingness, using Dual w/exp pairs. Figure 4 shows that REFLEX achieves higher or comparable non-misleadingness in most settings, except for Qwen-3, which is consistent with the automatic results. Overall error patterns are summarized in Appendix P. In some cases, human results appear more favorable than automatic ones, possibly due to the relatively high standard error of automatic pointwise misleadingness scores (Appendix R).

5.2 Disentanglement Effectiveness

To assess the effectiveness of disentanglement, we analyze it from two perspectives: empirical patterns in explanation quality and structural patterns in the learned directions.

Empirical Validation. We first aggregate pointwise evaluation results by the type of optimal vector, excluding readability because it is affected by post-processing. We then divide the results into correct and erroneous verdict cases. Two patterns emerge. **(1)** KVs are associated with lower misleadingness, suggesting a stronger relationship to **factual consistency**, whereas IVs are associated with higher informativeness and soundness, which is consistent with **stylistic improvement**. **(2)** Under KVs, the misleadingness gap between correct and

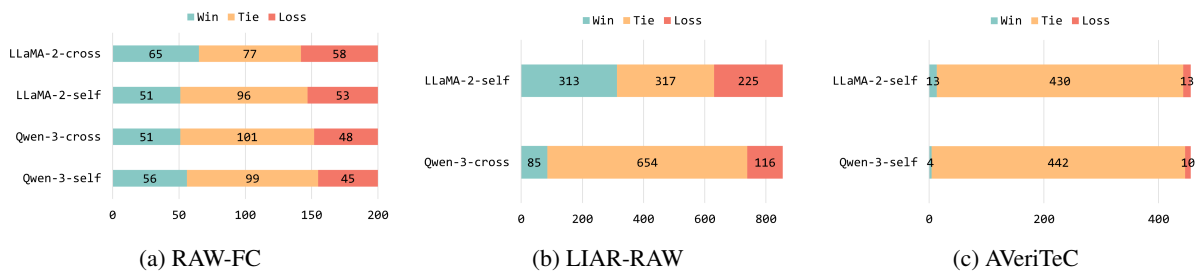


Figure 4: Pairwise human evaluation of explanation non-misleadingness against the corresponding SFT models.

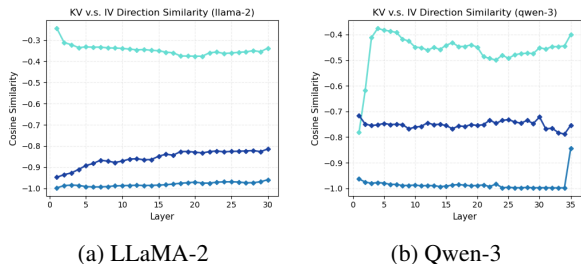


Figure 5: Layer-wise cosine similarity between KVs and IVs for **Dual w/o exp**, **Dual w/exp**, and **Single w/exp**.

erroneous cases is relatively large (1.52 vs. 1.74), which is consistent with the interpretation that KVs are more sensitive to factual differences. These findings support our distinction between IV-related reasoning style and KV-related factual sensitivity.

Representation Structure. We next analyze the geometry of the learned directions. Specifically, we extract the optimal KV and IV at each layer for explanation-related pairs and compute their cosine similarity. Because similar patterns are observed across datasets (Appendix N), we report the mean similarity. Figure 5 shows that KV and IV directions tend to exhibit **negative similarity**, which is consistent with the hypothesis that they capture distinct aspects of the model’s internal behavior. We also observe a pattern that we refer to as the **divergence amplification effect**: settings with larger directional discrepancies between KVs and IVs also tend to show stronger empirical gains. Across backbones, Qwen-3 exhibits smaller KV-IV differences than LLaMA-2, which may help explain its weaker gains and lower transferability. Across pair constructions, Single w/exp exhibits the largest directional differences, which aligns with the stronger improvements observed when explanation-related pairs are used to steer lower-information settings without explanations (Section 4.2.3). Additional analyses and case studies are provided in Appendix Q.

5.3 Statistical Correlations

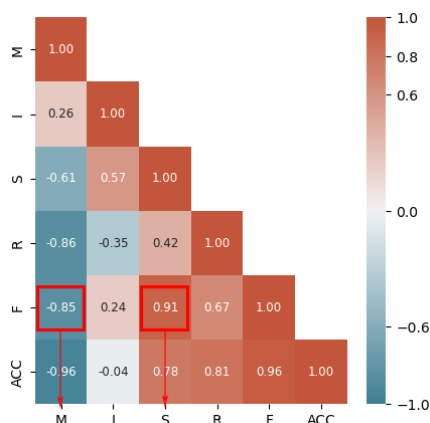


Figure 6: Correlation matrix between explanation quality and fact-checking performance.

To further examine whether improved explanations are associated with better verdict prediction, we analyze the correlation between verdict performance and point-wise explanation quality. Figure 6 shows that performance is most strongly associated with **faithfulness-related metrics**: both Macro-F1 and accuracy are strongly negatively correlated with misleadingness (-0.85 and -0.96), and strongly positively correlated with soundness (0.91 and 0.78). These results suggest that the gains of REFLEX are closely related to reduced faithfulness hallucination and improved logical coherence. Additional analyses and case studies are provided in Appendices S and T.

6 Conclusion

We introduced REFLEX, a simple and effective paradigm for explainable fact-checking that improves verdict prediction and explanation faithfulness by steering internal model activations rather than relying on external resources. Experiments across multiple settings demonstrate its effectiveness, generalizability, and flexibility.

Limitations

Experiment Scale. Although we conducted extensive ablations across backbone models, cross-model transfer, pair combinations, and other settings, disk quota constraints and the overall experimental scope limited our study to LLaMA-2, Qwen-3, and Mistral-v0.1 (Appendix P, Table 20). Moreover, research on real-world truthfulness does not necessarily follow standard scaling laws (Li et al., 2023b). We therefore restrict our evaluation to models at the 7B–8B scale.

Label Scheme. Although our method improves both verdict accuracy and explanation quality, thereby strengthening their alignment, we adopt a three-way label scheme on LIAR-RAW, following prior work (Wang et al., 2024). This choice reflects the difficulty LLMs face in generating explanations for overly fine-grained labels, such as “pants-on-fire.” Importantly, our goal is not merely to achieve state-of-the-art performance, but also to obtain interpretable insights.

Internal Design. Our internally guided paradigm reduces misinformation and response latency relative to methods that rely on external resources. Although internal knowledge may become outdated, this issue could in principle be mitigated through the lightweight activation-editing techniques on which we build (Section 2), enabling rapid adaptation to emerging fake news. However, existing time-shifted datasets such as VitaminC (Schuster et al., 2021) do not provide explanations, which prevents their direct use in our explainable setting.

Ethics Statement

Although our experiments use open-source datasets and models, these models may still contain social biases inherited from their pretraining data. A systematic analysis or mitigation of such biases is beyond the scope of this work. Our primary focus is on reducing faithfulness hallucination and improving explanation quality.

Acknowledgments

This project is supported by National Natural Science Foundation of China Young Scientists Fund (No. 62206233), RMGS (2025 First Processing Cycle), and the National Research Foundation, Singapore under its AI Singapore Programme (AISG Award No: AISG3-RP-2024-035). The experiment

equipments (4xA100s) is supported by HKBU Strategic Development Fund.

We especially thank Yaxin Fan, Ziyang Luo, and Feng Jiang for their insightful discussions. Meanwhile, we thank all reviewers for their helpful comments.

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*.
- Anthropic. 2025. Tracing the thoughts of a large language model. <https://www.anthropic.com/research/tracing-thoughts-language-model>.
- Pepa Atanasova, Oana-Maria Camburu, Christina Lioma, Thomas Lukasiewicz, Jakob Grue Simonsen, and Isabelle Augenstein. 2023. Faithfulness tests for natural language explanations. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 283–294.
- Lukas Berglund, Asa Cooper Stickland, Mikita Balesni, Max Kaufmann, Meg Tong, Tomasz Korbak, Daniel Kokotajlo, and Owain Evans. 2023. Taken out of context: On measuring situational awareness in llms. *arXiv preprint arXiv:2309.00667*.
- Collin Burns, Haotian Ye, Dan Klein, and Jacob Steinhardt. 2023. Discovering latent knowledge in language models without supervision. In *The Eleventh International Conference on Learning Representations*.
- Canyu Chen and Kai Shu. 2023. Can llm-generated misinformation be detected? In *The Twelfth International Conference on Learning Representations*.
- Canyu Chen and Kai Shu. 2024. Combating misinformation in the age of llms: Opportunities and challenges. *AI Magazine*, 45(3):354–368.
- Qiguang Chen, Libo Qin, Jinhao Liu, Dengyun Peng, Jiannan Guan, Peng Wang, Mengkang Hu, Yuhang Zhou, Te Gao, and Wanxiang Che. 2025. Towards reasoning era: A survey of long chain-of-thought for reasoning large language models. *arXiv preprint arXiv:2503.09567*.
- Tsun-Hin Cheung and Kin-Man Lam. 2023. Factllama: Optimizing instruction-following language models with external knowledge for automated fact-checking. In *2023 Asia Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, pages 846–853. IEEE.
- Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason Yosinski, and Rosanne Liu. 2019. Plug and play language models:

- A simple approach to controlled text generation. In *International Conference on Learning Representations*.
- Katja Filippova. 2020. Controlled hallucinations: Learning to generate faithfully from noisy data. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 864–870.
- Kanishk Gandhi, Ayush K Chakravarthy, Anikait Singh, Nathan Lile, and Noah Goodman. 2025. Cognitive behaviors that enable self-improving reasoners, or, four habits of highly effective stars. In *Second Conference on Language Modeling*.
- Zorik Gekhman, Gal Yona, Roei Aharoni, Matan Eyal, Amir Feder, Roi Reichart, and Jonathan Herzig. 2024. Does fine-tuning llms on new knowledge encourage hallucinations? In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 7765–7784.
- Gaurav Rohit Ghosal, Tatsunori Hashimoto, and Aditi Raghunathan. 2024. Understanding finetuning for factual knowledge extraction. In *International Conference on Machine Learning*, pages 15540–15558. PMLR.
- Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen, Shengjie Ma, Honghao Liu, Saizhuo Wang, Kun Zhang, Yuanzhuo Wang, Wen Gao, Lionel Ni, and Jian Guo. 2025. A survey on llm-as-a-judge. *The Innovation*.
- Zhijiang Guo, Michael Schlichtkrull, and Andreas Vlachos. 2022. A survey on automated fact-checking. *Transactions of the association for computational linguistics*, 10:178–206.
- Chi Han, Jialiang Xu, Manling Li, Yi Fung, Chenkai Sun, Nan Jiang, Tarek Abdelzaher, and Heng Ji. 2024. Word embeddings are steers for language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16410–16430.
- Zicong He, Boxuan Zhang, and Lu Cheng. 2025. Shakespearean sparks: The dance of hallucination and creativity in llms’ decoding layers. *Preprint*, arXiv:2503.02851.
- Evan Hernandez, Belinda Z Li, and Jacob Andreas. 2023. Inspecting and editing knowledge representations in language models. *arXiv preprint arXiv:2304.00740*.
- Wentao Hu, Wengyu Zhang, Yiyang Jiang, Chen Jason Zhang, Xiaoyong Wei, and Li Qing. 2025. Removal of hallucination on hallucination: Debate-augmented rag. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15839–15853.
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. 2025a. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *ACM Transactions on Information Systems*, 43(2):1–55.
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. 2025b. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *ACM Transactions on Information Systems*, 43(2):1–55.
- Alon Jacovi and Yoav Goldberg. 2020. Towards faithfully interpretable nlp systems: How should we define and evaluate faithfulness? In *Proceedings of the 58th annual meeting of the association for computational linguistics*, pages 4198–4205.
- Adit Jain, Vikram Krishnamurthy, and Yiming Zhang. 2025. Collaborative qa using interacting llms. impact of network structure, node capability and distributed data. *arXiv preprint arXiv:2511.14098*.
- Kyungha Kim, Sangyun Lee, Kung-Hsiang Huang, Hou Pong Chan, Manling Li, and Heng Ji. 2024. Can llms produce faithful explanations for fact-checking? towards faithful explainable fact-checking via multi-agent debate. *arXiv preprint arXiv:2402.07401*.
- Aobo Kong, Shiwan Zhao, Hao Chen, Qicheng Li, Yong Qin, Ruiqi Sun, Xin Zhou, Enzhi Wang, and Xiaohang Dong. 2024. Better zero-shot reasoning with role-play prompting. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 4099–4113.
- Chuyi Kong, Ziyang Luo, Hongzhan Lin, Zhiyuan Fan, Yaxin Fan, Yuxi Sun, and Jing Ma. 2025. Sharp: Unlocking interactive hallucination via stance transfer in role-playing llms. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 839–866.
- Jan Leike. 2022. *Distinguishing three alignment taxes*.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.
- Kenneth Li, Aspen K Hopkins, David Bau, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2023a. Emergent world representations: Exploring a sequence model trained on a synthetic task. In *The Eleventh International Conference on Learning Representations*.

- Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2023b. Inference-time intervention: Eliciting truthful answers from a language model. *Advances in Neural Information Processing Systems*, 36:41451–41530.
- Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2023c. Inference-time intervention: Eliciting truthful answers from a language model. *Advances in Neural Information Processing Systems*, 36:41451–41530.
- Xiang Li, John Thickstun, Ishaan Gulrajani, Percy S Liang, and Tatsunori B Hashimoto. 2022. Diffusion-lm improves controllable text generation. *Advances in neural information processing systems*, 35:4328–4343.
- Dayong Liang, Xiao-Yong Wei, and Changmeng Zheng. 2026. Multi-agent undercover gaming: Hallucination removal through counterfactual test for multimodal reasoning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 40, pages 6807–6815.
- Philip Lippmann and Jie Yang. 2025. Style over substance: Distilled language models reason via stylistic replication. *arXiv preprint arXiv:2504.01738*.
- Jiachang Liu, Dinghan Shen, Yizhe Zhang, William B Dolan, Lawrence Carin, and Weizhu Chen. 2022a. What makes good in-context examples for gpt-3? In *Proceedings of Deep Learning Inside Out (DeeLIO 2022): The 3rd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures*, pages 100–114.
- Jiacheng Liu, Alisa Liu, Ximing Lu, Sean Welleck, Peter West, Ronan Le Bras, Yejin Choi, and Hannaneh Hajishirzi. 2022b. Generated knowledge prompting for commonsense reasoning. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3154–3169.
- Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. 2022. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8086–8098.
- Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. On faithfulness and factuality in abstractive summarization. In *Proceedings of the 58th annual meeting of the association for computational linguistics*, pages 1906–1919.
- Sewon Min, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. Noisy channel language model prompting for few-shot text classification. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5316–5330.
- Tai Nguyen and Eric Wong. 2023. In-context example selection with influences. *arXiv e-prints*, pages arXiv–2302.
- OpenAI. 2023. [Introducing chatgpt](#).
- Letitia Parcalabescu and Anette Frank. 2024. On measuring faithfulness or self-consistency of natural language explanations. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6048–6089.
- Verónica Pérez-Rosas, Bennett Kleinberg, Alexandra Lefevre, and Rada Mihalcea. 2018. Automatic detection of fake news. In *Proceedings of the 27th international conference on computational linguistics*, pages 3391–3401.
- Alamgir Munir Qazi, John Philip McCrae, and Jamal A Nasir. 2025. When retrieval outperforms generation: Dense evidence retrieval for scalable fake news detection. In *Proceedings of the 5th Conference on Language, Data and Knowledge*, pages 255–265.
- Cheng Qian, Xinran Zhao, and Sherry Tongshuang Wu. 2023. "merge conflicts!" exploring the impacts of external distractors to parametric knowledge graphs. *arXiv preprint arXiv:2309.08594*.
- Hannah Rashkin, Eunsol Choi, Jin Yea Jang, Svitlana Volkova, and Yejin Choi. 2017. Truth of varying shades: Analyzing language in fake news and political fact-checking. In *Proceedings of the 2017 conference on empirical methods in natural language processing*, pages 2931–2937.
- John W. Ratcliff and David E. Metzener. 1988. Pattern matching: The gestalt approach. *Dr. Dobb's Journal*, 13(7):46.
- Xuan Ren, Biao Wu, and Lingqiao Liu. 2024. I learn better if you speak my language: Enhancing large language model fine-tuning with style-aligned response adjustments. *CoRR*.
- Nina Rimsky, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Turner. 2024. Steering llama 2 via contrastive activation addition. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15504–15522.
- Michael Schlichtkrull, Zhijiang Guo, and Andreas Vlachos. 2023. Averitec: A dataset for real-world claim verification with evidence from the web. *Advances in Neural Information Processing Systems*, 36:65128–65167.
- Tal Schuster, Adam Fisch, and Regina Barzilay. 2021. Get your vitamin c! robust fact verification with contrastive evidence. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 624–643.

- Tal Schuster, Roei Schuster, Darsh J Shah, and Regina Barzilay. 2020. The limitations of stylometry for detecting machine-generated fake news. *Computational Linguistics*, 46(2):499–510.
- Yunfan Shao, Linyang Li, Junqi Dai, and Xipeng Qiu. 2023. Character-llm: A trainable agent for role-playing. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 13153–13187.
- Satyam Shukla, Himanshu Dutta, and Pushpak Bhat-tacharyya. 2025. Recon, answer, verify: Agents in search of truth. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pages 2429–2448.
- Yuxi Sun, Aoqi Zuo, Wei Gao, and Jing Ma. 2025. Causalabstain: Enhancing multilingual llms with causal reasoning for trustworthy abstention. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 14060–14076.
- Yuxi Sun, Aoqi Zuo, Haotian Xie, Wei Gao, Mingming Gong, and Jing Ma. 2026. Fact-e: Causality-inspired evaluation for trustworthy chain-of-thought reasoning.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Miles Turpin, Julian Michael, Ethan Perez, and Samuel Bowman. 2023. Language models don’t always say what they think: Unfaithful explanations in chain-of-thought prompting. *Advances in Neural Information Processing Systems*, 36:74952–74965.
- Bo Wang, Jing Ma, Hongzhan Lin, Zhiwei Yang, Ruichao Yang, Yuan Tian, and Yi Chang. 2024. Explainable fake news detection with large language model via defense among competing wisdom. In *Proceedings of the ACM Web Conference 2024*, pages 2452–2463.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.
- Zeyu Wei, Shuo Wang, Xiaohui Rong, Xuemin Liu, and He Li. 2025. Shadows in the attention: Contextual perturbation and representation drift in the dynamics of hallucination in llms. In *International Conference on Neural Information Processing*, pages 462–476. Springer.
- Jiaying Wu, Jiafeng Guo, and Bryan Hooi. 2024. Fake news in sheep’s clothing: Robust fake news detection against llm-empowered style attacks. In *Proceedings of the 30th ACM SIGKDD conference on knowledge discovery and data mining*, pages 3367–3378.
- Yang Yan, Yu Lu, Renjun Xu, and Zhenzhong Lan. 2025. Do phd-level llms truly grasp elementary addition? probing rule learning vs. memorization in large language models. *arXiv e-prints*, pages arXiv–2504.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. 2025. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*.
- Zhiwei Yang, Jing Ma, Hechang Chen, Hongzhan Lin, Ziyang Luo, and Yi Chang. 2022. A coarse-to-fine cascaded evidence-distillation neural network for explainable fake news detection. In *Proceedings of the 29th international conference on computational linguistics*, pages 2608–2621.
- Zeyu Yun, Yubei Chen, Bruno Olshausen, and Yann Le-Cun. 2021. Transformer visualization via dictionary learning: contextualized embedding as a linear superposition of transformer factors. In *Proceedings of Deep Learning Inside Out (DeeLIO): The 2nd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures*, pages 1–10.
- Fengzhu Zeng and Wei Gao. 2024. Justilm: Few-shot justification generation for explainable fact-checking of real-world claims. *Transactions of the Association for Computational Linguistics*, 12:334–354.
- Tianjun Zhang, Xuezhi Wang, Denny Zhou, Dale Schuurmans, and Joseph E Gonzalez. 2022. Tempera: Test-time prompt editing via reinforcement learning. In *The Eleventh International Conference on Learning Representations*.
- Xuan Zhang and Wei Gao. 2023. Towards llm-based fact verification on news claims with a hierarchical step-by-step prompting method. In *Proceedings of the 13th international joint conference on natural language processing and the 3rd conference of the asia-pacific chapter of the association for computational linguistics (volume 1: Long papers)*, pages 996–1011.
- Eric Zhao, Pranjal Awasthi, and Nika Haghtalab. 2025. From style to facts: Mapping the boundaries of knowledge injection with finetuning. In *The Thirtieth Annual Conference on Neural Information Processing Systems*.
- Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. 2023. Lima: Less is more for alignment. *Advances in Neural Information Processing Systems*, 36:55006–55021.

Appendix

A Concept Alignment

Hallucination typically refers to the generation of content that is either nonsensical or unfaithful to the provided source information (Filippova, 2020; Maynez et al., 2020).

Its impact is application-dependent. In some domains, such as personalized role-play dialogue (Kong et al., 2025) and creative generation (He et al., 2025), it can even be beneficial. However, in high-stakes societal domains such as law and fact-checking, as well as in knowledge-intensive tasks (Sun et al., 2025, 2026), hallucinations are harmful.

In the era of large language models (LLMs), recent work (Huang et al., 2025a) further categorizes hallucinations into two main types: factuality hallucinations and faithfulness hallucinations.

Factuality hallucinations arise when LLMs produce outputs that are inconsistent with real-world facts or are potentially misleading.

Faithfulness hallucinations concern whether generated content remains consistent with user instructions, provided context, and internal logical coherence. Accordingly, they can be categorized into instruction inconsistency, context inconsistency, and logical inconsistency. We argue that, in essence, this type of hallucination can be understood as **context inconsistency**.

Misinformation refers to false or inaccurate information that is created deliberately and propagated either intentionally or unintentionally. In the context of social media, Chen and Shu (2024) adopts a broader view, treating misinformation as an umbrella term that encompasses all false or inaccurate information circulating on social media platforms.

Under this definition, misinformation can be categorized into several types, including **fake news**, rumors, conspiracy theories, clickbait, **misleading claims**, and cherry-picking.

Moreover, they suggest that LLM-generated misinformation may exhibit **more deceptive styles**, thereby posing greater potential risk of real-world harm.

Misinformation vs. Hallucination. Chen and Shu (2024) shows that hallucination is one source of misinformation, while Huang et al. (2025a) suggests that misinformation can also cause hallucination. Therefore, they can be regarded as intertwined

phenomena that mutually influence each other in a bidirectional relationship.

Knowledge vs. Style. “IVs extracted from Q2 samples capture **style-dependent** veracity gains” is supported by prior empirical work:

- OpenAI’s fine-tuning guide notes that fine-tuning performs well in show-not-tell scenarios: “show” corresponds to stylistic patterns, whereas “tell” corresponds to knowledge injection.

- Several studies (Ren et al., 2024; Ghosal et al., 2024; Zhao et al., 2025) argue that injecting specific or new knowledge via fine-tuning is comparatively difficult. Zhao et al. (2025) finds that fine-tuning on recent Wikipedia events (knowledge) performs worse, whereas fine-tuning on writing with a specific tone (style) yields near-perfect behavior.

- Yan et al. (2025) shows that probability-based LLM inference in reasoning tasks often reflects memorized stylistic patterns.

“KVs extracted from Q4 samples isolate **fact-sensitive** veracity drops” is similarly supported by prior empirical findings.

- Li et al. (2023c) points out that backbone models already encode extensive factual knowledge (Section 3.2.1, lines 203–207).

- Gekhman et al. (2024) reports that fine-tuning on new knowledge can increase hallucination risk, especially for dynamic knowledge in fact-checking.

- In alignment research, this phenomenon is sometimes informally referred to as an alignment tax (Leike, 2022).

Transferred to our task, this literature suggests a natural mapping: **style** corresponds to **reasoning style**, while **knowledge** corresponds to **factual grounding**. This distinction motivates separating IV-driven style from KV-driven factual sensitivity.

Faithfulness vs. Plausibility on NLE. Natural Language Explanations (NLEs) refer to explanations expressed in natural language. According to (Agarwal et al., 2024), recent work shows that modern LLMs can generate self-explanations (SEs), which aim to elicit intermediate reasoning steps to account for model outputs. They define the concepts of faithfulness and plausibility.

- **Plausibility** concerns whether an explanation aligns with human reasoning and understanding. In this work, we treat it as a **stylistic pattern**.

- **Faithfulness** evaluates whether an explanation accurately reflects the model’s underlying reason-

ing process. However, since we do not assume that LLMs possess genuine reasoning capabilities, our notion of faithfulness is framed through the lens of **faithfulness hallucinations**.

B Prompt Template

Following Liu et al. (2022b); Chen et al. (2025), the prompt template we use for training and inference is as follows:

A chat between a curious human and an artificial intelligence assistant.
You are a fact-checking assistant.
You are given a claim [(along with evidence sentences)].
Your task is to label the overall veracity of the claim based on your internal knowledge [or, on the provided evidence sentences].
Please reason step by step and explain how you reached your conclusion.

The Label Definitions are as follows:
TRUE|true: The claim is verified as TRUE|true based on your knowledge.
FALSE|false: The claim is verified as FALSE|false based on your knowledge.
HALF-TRUE|half: The claim is verified as half-true|half due to insufficient knowledge leading to uncertainty, or because the claim itself is partially true.
or, for given evidence:
TRUE|true: The claim is verified as true or mostly-true by the evidence.
FALSE|false: The claim is verified as false or mostly false by the evidence.
HALF-TRUE|half: The claim is verified as half-true|half by the evidence, or the evidence can not prove the claim.
or, for AVeriTeC:
Supported: The claim is supported by the arguments and evidence presented.
Refuted: The claim is contradicted by the arguments and evidence presented.
Not Enough Evidence: There is not enough evidence to support or refute the claim.

Based on the below claim and your own knowledge, determine the veracity of the claim.
or, for given evidence:

Based on the below claim and evidence, determine the veracity of the claim.
Please strictly output in the following format:
Verdict: [label]. Explanation: [your reasoning here]

Consistent with Wang et al. (2024), the prompt template used to ask ChatGPT to conduct automatic evaluation is shown below:

You are a helpful, harmless and precise assistant.
Please evaluate the quality of the explanations in prediction based on four metrics: misleadingness, informativeness, soundness, and readability, where 1 represented the poorest and 5 the best in addition to misleadingness.
The definitions of the metrics are:
(1) Misleadingness assesses whether the model's explanation is consistent with the real veracity label of a claim, with a rating scale ranging from 1 (not misleading) to 5 (very misleading);
(2) Informativeness assesses whether the explanation provides new information, such as explaining the background and additional context, with a rating scale ranging from 1 (not informative) to 5 (very informative)
(3) Soundness describes whether the explanation seems valid and logical, with a rating scale ranging from 1 (not sound) to 5 (very sound)
(4) Readability evaluates whether the explanation follows proper grammar and structural rules, and whether the sentences in the explanation fit together and are easy to follow, with a rating scale ranging from 1 (poor) to 5 (excellent).
Please output scores in the following JSON format:
```\n{\n 'Misleadingness': [an integer number from 1 to 5],\n 'Informativeness': [an integer number from 1 to 5],\n 'Soundness': [an integer number from 1 to 5],\n}

'Readability': [an integer number from 1 to 5]  
 }}  
 ...

Claim: {claim} Label: {label}  
 Prediction: {verdict}  
 Score:

## C Implementation Details

For model training, Table 11 shows the hyperparameter settings. Except for gradient accumulation steps and maximum context length, which depend on dataset size and sample length, we trained models on 1–4×A100-80G GPUs with a learning rate of 2e-5, per-device batch size of 4, weight decay 0, warmup ratio 0.03, a cosine scheduler, bf16/tf32 precision, gradient checkpointing, and full-shard FSDP with auto wrapping. For LLaMA-2, we selected the epoch-2 checkpoint, which achieved the best balance between factual accuracy and instruction following. For Qwen-3, severe overfitting occurred at epoch 2 on some LIAR-RAW and RAW-FC variants (training accuracy exceeded test accuracy by 20–40%), so we used the epoch-1 checkpoint instead.

For steering, consistent with the terminology of CAA (Rimsky et al., 2024), we use multipliers to denote steering magnitude. We sweep all layers and vary the multiplier over  $\pm 0.5$ ,  $\pm 1$ ,  $\pm 1.5$ , and  $\pm 2$ . Empirically, multipliers with magnitude  $\geq 2$  (we additionally tested  $\pm 2$ ,  $\pm 3$ ,  $\pm 5$ ,  $\pm 10$ , and  $\pm 20$ ) drive the model into non-instruction-following responses. The full optimal layers and multiplier hyperparameters across baseline trials and ablation studies are reported in Appendix O.

We formalize disentangled steering as quadrant-conditioned activation probing over decoder hidden states in Algorithm 1.

Let the hidden activation at decoder layer  $l$  and token position  $t$  be:

$$\mathbf{h}_{l,t} \in \mathbb{R}^d$$

Two supervision datasets are constructed from cross-stage quadrants. The Q2 dataset is defined as:

$$\mathcal{D}_{IV} = \{(\mathbf{h}_{l,i}^+, \mathbf{h}_{l,i}^-)\}_{i=1}^N$$

where:

- $\mathbf{h}^+$ : SFT-correct reasoning activations
- $\mathbf{h}^-$ : backbone-error activations

The Q4 dataset is defined as:

$$\mathcal{D}_{KV} = \{(\mathbf{h}_{l,i}^+, \mathbf{h}_{l,i}^-)\}_{i=1}^N,$$

where:

- $\mathbf{h}^+$ : backbone-correct activations
- $\mathbf{h}^-$ : SFT-drift activations

---

### Algorithm 1: Disentangled Steering via Quadrant-Conditioned Probing

---

**Input:** Q2 dataset (reasoning correction), Q4 dataset (factual preservation), decoder layers  $L$ , steering multipliers  $\mathcal{A}$

**Output:** Optimal inference vector  $IV^*$  and knowledge vector  $KV^*$

**foreach**  $l \in L$  **do**

  Extract Q2 activations  $\{(h_{l,i}^+, h_{l,i}^-)\}_{i=1}^N$ ;

  Train logistic probe

$$p_l^{IV}(h) = \sigma(W_l^{IV\top} h + b_l^{IV});$$

  Normalize  $IV_l = W_l^{IV} / \|W_l^{IV}\|$ ;

  Extract Q4 activations  $\{(h_{l,i}^+, h_{l,i}^-)\}_{i=1}^N$ ;

  Train logistic probe

$$p_l^{KV}(h) = \sigma(W_l^{KV\top} h + b_l^{KV});$$

  Normalize  $KV_l = W_l^{KV} / \|W_l^{KV}\|$ ;

**foreach**  $\alpha \in \mathcal{A}$  **do**

    Apply IV steering:

$$h'_{l,t} = h_{l,t} + \alpha IV_l;$$

    Apply KV steering:

$$h'_{l,t} = h_{l,t} - \alpha KV_l;$$

    Compute improvement  $\Delta P_{l,\alpha}$ ;

  Record  $(l, \alpha_l^*) = \arg \max_{\alpha} \Delta P_{l,\alpha}$ ;

Select  $l^* = \arg \max_l \Delta P_{l,\alpha_l^*}$ ;

Set:  $IV^* = \alpha_{l^*} IV_{l^*}$ ,  $KV^* = \alpha_{l^*} KV_{l^*}$ ;

---

## D Alignment between Probability and Performance

Here, we report the probability gap and the corresponding Macro-F1 for Dual w/explanation pairs on RAW-FC with LLaMA-2. As shown in Table 7, the overall correlation between probability shifts and verdict accuracy reaches 92%, and their extrema are aligned: the largest probability gaps correspond to the largest performance shifts in both directions.

Table 7: Maximum and minimum probability gaps ( $\Delta\text{Prob.}$ ) between unsteered and steered outputs across 32 layers (step = 5), with corresponding Macro-F.

| Layer                         | Multiplier | Macro-F (%)  | $\Delta\text{Prob.}(\%)$ |
|-------------------------------|------------|--------------|--------------------------|
| 1                             | 1.5        | <b>48.64</b> | <b>-4.32</b>             |
| 5                             | 1.5        | 58.59        | 0.07                     |
| 10                            | 1.5        | <b>64.99</b> | <b>0.14</b>              |
| 10                            | -1.5       | 59.61        | -0.15                    |
| 11                            | -1.5       | 60.91        | -0.14                    |
| 15                            | 1.5        | 63.09        | 0.11                     |
| 16                            | -1.5       | 61.63        | -0.12                    |
| 17                            | 1.5        | 59.88        | 0.11                     |
| 22                            | -1.5       | 61.10        | -0.08                    |
| 22                            | 1.5        | 60.57        | 0.07                     |
| 29                            | 1.5        | 60.63        | 0.06                     |
| 30                            | -1.5       | 61.13        | -0.06                    |
| Correlation (Macro-F vs. Gap) |            | 0.92         |                          |

Table 8: Summary statistics of dataset distributions. Label values 0-2 represent increasing veracity labels: {False/Refuted, Half-True/Not Enough Evidence, True/Supported}.

| Dataset  | Split | 0     | 1     | 2     | Total |
|----------|-------|-------|-------|-------|-------|
| RAWFC    | train | 514   | 537   | 561   | 1,612 |
|          | eval  | 66    | 67    | 67    | 200   |
|          | test  | 66    | 67    | 67    | 200   |
| LIAR-RAW | train | 2,568 | 1,336 | 2,264 | 6,168 |
|          | eval  | 410   | 159   | 292   | 861   |
|          | test  | 367   | 169   | 319   | 855   |
| AVeriTeC | train | 1,742 | 849   | 282   | 2,873 |
|          | eval  | 305   | 35    | 122   | 462   |
|          | test  | 303   | 33    | 120   | 456   |

## E Dataset Details

We use three fact-checking datasets: RAW-FC from Snopes<sup>1</sup>, LIAR-RAW (Yang et al., 2022) from PoliFact<sup>2</sup>, and AVeriTeC. For RAW-FC and LIAR-RAW, evidence corresponds to the annotated relevant evidence (labeled as 1). Since AVeriTeC is natively multi-turn, we flatten each dialogue into a single-turn instance for consistency.

In terms of data structure, RAW-FC and LIAR-RAW follow the standard fact-checking format {claim, evidence, label, explanation}. In contrast, AVeriTeC decomposes fact-checking into a QA-style multi-turn verification process.

The datasets adopt different label schemes. RAW-FC uses {True, Half-True, False}, LIAR-RAW contains six labels, and AVeriTeC uses {Supported, Not Enough Evidence, Conflicting Evidence/Cherrypicking, Refuted}. To support joint

<sup>1</sup>www.snopes.com

<sup>2</sup>www.politifact.com

verdict prediction and explanation generation in REFLEX, we unify labels as follows. For LIAR-RAW, we merge {pants-fire, false, barely-true} into False, retain Half-True, and merge {mostly-true, true} into True. For AVeriTeC, we discard Conflicting Evidence/Cherrypicking because of its ambiguity. We further remove LIAR-RAW instances without evidence and exclude AVeriTeC few-shot examples from validation to prevent data leakage.

As shown in Table 8, the training sample sizes of the three datasets are all below 10K, which matches the assumption of our paradigm, namely fine-tuning on limited data. Among them, RAW-FC is label-balanced, whereas LIAR-RAW and AVeriTeC are not. Since AVeriTeC does not release a test set to avoid leakage, we use its validation set for evaluation after removing training overlaps. Finally, because some baselines cannot handle dialogue-style inputs, baseline comparisons are reported only on RAW-FC and LIAR-RAW.

Since the human-annotated RAW-FC test set is relatively small, similar to other manually constructed benchmarks such as MT-Bench, we additionally compute confidence intervals using bootstrap resampling for the first ablation study (see Appendix L).

## F Automatic Evaluation Details

We conduct automatic pointwise evaluation on four dimensions: (1) **Misleadingness**, which assesses whether the model’s explanation is consistent with the true veracity label of a claim, with a rating scale from 1 (not misleading) to 5 (very misleading); (2) **Informativeness**, which assesses whether the explanation provides new information, such as background or additional context, with a rating scale from 1 (not informative) to 5 (very informative); (3) **Soundness**, which assesses whether the explanation appears valid and logical, with a rating scale from 1 (not sound) to 5 (very sound); (4) **Readability**, which assesses whether the explanation follows proper grammar and structure and is easy to follow, with a rating scale from 1 (poor) to 5 (excellent).

Each dimension is rated on a five-point Likert scale, with higher scores indicating better quality except for misleadingness, which is inversely scored. The prompt template is shown in Appendix B. Following L-Defense, we use gpt-3.5-turbo-0613.

## G Full Baselines Comparisons

Due to space constraints in the main text, we defer the full comparison of baselines and time complexity to this appendix.

As shown in Table 9, distill/RAG-then-SFT paradigms (e.g., FACTLLAMA and L-DEFENSE) incur the highest time cost, because they require collecting multiple pieces of evidence per claim from closed-source models or search APIs, followed by task-specific fine-tuning. For  $m$  claims and  $n$  associated evidence items, both scale as  $O(mn)$ , with L-DEFENSE further introducing adversarial explanations during distillation. Multi-agent approaches such as RAV remove training overhead but introduce substantial inference-time cost: each claim involves  $Z$  agents,  $(Z - 1)$  interactions, and  $Z$  inference steps, resulting in  $O(mZ)$  complexity. RAG-based methods such as HiSS decompose each claim into  $n$  sub-claims and iteratively perform retrieval, requiring repeated interactions with external APIs and yielding  $O(mn)$  complexity.

In contrast, REFLEX eliminates reliance on closed-source APIs, external evidence, and multi-agent systems ( $C = 0, O = 0, Z = 0$ ). Although it requires training and vector extraction, once the optimal vectors are obtained they can be transferred to other models with only a constant number of inference steps per claim, resulting in linear time complexity  $O(m)$ .

We also include DeReC (Qazi et al., 2025), a recent fact-checking method. It is a retrieval-augmented fact verification framework that combines dense embedding-based evidence retrieval with FAISS-based similarity search and a DeBERTa-v3 classifier for veracity prediction. However, despite its efficiency and strong performance among claim-verification systems, it still underperforms REFLEX in Macro-F1. Following Wang et al. (2024), and due to budget limitations, we exclude one additional model that also leverages external evidence (Zeng and Gao, 2024).

## H Explanation Length

As Table 10 shows, on RAW-FC, our outputs are shorter than those of L-Defense. On LIAR-RAW, they are shorter than all baselines, including the oracle, demonstrating that our paradigm produces concise and accurate explanations. Table 15 further shows that EGS clearly reduces noisy patterns in model outputs across backbones and datasets.

## I Human Evaluation Details

All annotators were undergraduate students at a university where English is the official language. For the baseline experiments, we conducted pointwise evaluation of explanations on each dimension, following Wang et al. (2024). Ten undergraduate annotators rated 30 randomly sampled test instances from RAW-FC using a 5-point Likert scale. Model identities were anonymized, and average scores were used as the final metric. Table 13 shows that all four dimensions have correlation coefficients above 0.70 (0.95, 0.73, 0.85, 0.81), supporting the use of LLM-as-a-Judge (Gu et al., 2025). Our SEGs model outperforms the baselines on all four dimensions.

For inter-annotator agreement, because these are pointwise evaluations, we also computed Cronbach’s  $\alpha$ . The overall average  $\alpha$  is 0.7896, which is acceptable by conventional standards.

For the ablation studies, we conducted pairwise evaluation focused on misleadingness, our main target dimension. Three undergraduate annotators rated each test sample from RAW-FC, LIAR-RAW, and AVeriTeC. Each annotator compared the explanations produced by two models and determined which one was more misleading. Model names remained anonymous, and the order of the two model outputs was randomly swapped. For inter-annotator agreement, we computed Fleiss’ Kappa. The overall average Kappa is 0.6963, indicating substantial agreement.

Notably, the case studies in Appendix T summarize common misinformation patterns and present annotator-verified examples. In these case studies, we additionally use ChatGPT-generated explanations for double validation.

## J Label Distribution

As shown in Table 12, the two backbone models exhibit severe recency bias on LIAR-RAW, although this issue can be mitigated by prior methods (Lu et al., 2022; Min et al., 2022; Liu et al., 2022a; Zhang et al., 2022; Nguyen and Wong, 2023).

## K Full Statistics on Backbone Ablations

Formally, we define the **hallucination rate (HR)** and **inference success rate (ISR)** as:

$$\text{HR} = \frac{\#\text{error after SFT}}{\#\text{correct on BASE}}, \quad (1)$$



Table 10: Explanation lengths for our method and the baselines.

| Method                       | RAW-FC        | LIAR-RAW     |
|------------------------------|---------------|--------------|
| Oracle                       | 201.68        | 220.75       |
| ChatGPT <sub>w/evi</sub>     | 144.32        | 139.15       |
| ChatGPT <sub>w/o evi</sub>   | <b>128.71</b> | 150.97       |
| L-Defense <sub>ChatGPT</sub> | 266.61        | 225.52       |
| L-Defense <sub>LLaMA2</sub>  | 305.50        | 175.38       |
| <b>Ours</b>                  |               |              |
| S-EGS <sub>LLaMA2</sub>      | 264.64        | <b>76.50</b> |
| w/o EGS                      | 787.55        | 118.06       |

Table 11: Hyperparameter settings for different models.

| Dataset  | #train | x→y         | max-len | grad_acc | epochs<br>qwen / llama |
|----------|--------|-------------|---------|----------|------------------------|
| RAW-FC   | 1,612  | c→v         | 512     | 1        | 1 / 2                  |
|          |        | c;evi→v     | 1,024   | 1        | 2 / 2                  |
|          |        | c→v;exp     | 4,096   | 1        | 2 / 2                  |
| AVeriteC | 2,873  | c;evi→v     | 512     | 4        | 2 / 2                  |
|          |        | c;evi→v;exp | 512     | 4        | 2 / 2                  |
| LIAR-RAW | 6,168  | c→v         | 256     | 8        | 1 / 2                  |
|          |        | c;evi→v     | 512     | 8        | 1 / 2                  |
|          |        | c→v;exp     | 1,024   | 8        | 2 / 2                  |

Table 12: Label distribution and few-shot ordering on LIAR-RAW. H denotes half-true, F false, and T true.

| Model   | variant | order  | k-shot | split | h/f/t                   |
|---------|---------|--------|--------|-------|-------------------------|
| LLaMA-2 | cross   | fth    | 3      | test  | <b>842</b> /8/5         |
|         |         |        |        | train | <b>6112</b> /37/19      |
| Qwen-3  | self    | fththf | 6      | test  | 189/ <b>521</b> /145    |
|         |         |        |        | train | 1223/ <b>3844</b> /1101 |

Table 13: Automatic and human evaluation of explanation quality. We report the best REFLEX variant on LLaMA-2.

|                              | ChatGPT     |             |             |             | Human       |             |             |             |
|------------------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
|                              | M↓          | I           | S           | R           | M↓          | I           | S           | R           |
| Oracle                       | 1.53        | 4.50        | 4.77        | 4.77        | 1.47        | 3.61        | 3.89        | 3.86        |
| ChatGPT <sub>c,v</sub>       | 2.07        | 4.43        | 4.67        | 4.73        | 2.22        | 3.22        | 3.38        | 3.57        |
| ChatGPT <sub>c</sub>         | 2.33        | 4.17        | 4.43        | 4.63        | 2.68        | 2.68        | 2.84        | 3.27        |
| L-Defense <sub>LLaMA2</sub>  | 1.87        | 4.50        | 4.67        | 4.67        | 2.12        | 3.48        | 3.37        | 3.49        |
| L-Defense <sub>ChatGPT</sub> | 1.77        | 4.40        | 4.60        | 4.53        | 1.97        | 3.68        | 3.52        | 3.56        |
| <b>Ours</b>                  |             |             |             |             |             |             |             |             |
| S-EGS <sub>LLaMA2</sub>      | <b>1.65</b> | <b>4.79</b> | <b>4.86</b> | <b>4.88</b> | <b>1.75</b> | <b>3.76</b> | <b>3.92</b> | <b>3.96</b> |
| w/o EGS                      | 1.89        | 4.76        | 4.78        | 4.50        | 2.35        | 3.48        | 3.36        | 2.62        |

statistically stable.

## M Full Cross-Model Transfer Results

The full cross-model transfer results are shown in Table 17. Here,  $\Delta$  denotes the difference relative to the corresponding SFT model.

## N Full Direction Similarity of Vectors

The full direction similarity, without averaging across datasets, is shown in Figure 8.

## O Optimal Hyperparameters Results

The full optimal hyperparameters, including layers and multipliers across baseline and ablation settings, are shown in Tables 18 and 19. Because horizontal steering is not specific to REFLEX, we do not report its direction.

As shown in both tables, most KV multipliers are negative, while most IV multipliers are positive, indicating KV suppression and IV amplification, respectively. For the bolded cases, the majority of the corresponding probabilities are negative, which further supports KV suppression and IV amplification. The only exceptions are the two underlined KVs; accordingly, they yield only 0.01 and 0.19 Macro-F1-point gains.

We also observe that the optimal number of vectors required for IV is larger than that required for KV. This is because suppressing erroneous directions in KV is a *necessary condition* for ensuring factual consistency between the verdict and the explanation, whereas amplifying correct directions in IV serves as a *sufficient condition* that directly aligns the reasoning process in explanations with successful verdict prediction. In general, sufficient conditions are less restrictive than necessary ones.

## P Total Error Analysis

Overall, across 38 experiments (excluding performance-adaptive transfer trials), the Macro-F1 score drops by more than 0.01 only six times: -0.28, -0.20, -0.03, -1.05, -1.12, and -1.17. Detailed analyses are provided below:

(1) The decreases of -0.28 and -0.20 occur only on LIAR-RAW. As noted in Section 4.2.1 and Appendix J (Table 12), this is due to severe **recency bias** in few-shot learning. Although prior methods (Lu et al., 2022; Min et al., 2022; Liu et al., 2022a; Zhang et al., 2022; Nguyen and Wong, 2023) can mitigate this bias, we deliberately avoid them to maintain experimental consistency.

Table 14: Full statistics of hallucination ratio (HR), inference success rate (ISR), and data amounts across backbones and datasets. *cross* denotes few-shot examples drawn from another dataset’s training set, while *self* denotes examples drawn from the model’s own validation set.

| Backbone | Dataset                         | x->y                       | Statistics |               |               |               |
|----------|---------------------------------|----------------------------|------------|---------------|---------------|---------------|
|          |                                 |                            | #HC↓       | #ISC↑         | HR↓           | ISR↑          |
| LLaMA-2  | RAW-FC                          | c->v                       | 88         | 231           | 0.1236        | <b>0.2567</b> |
|          |                                 | c->v; exp <sub>self</sub>  | 113        | 585           | 0.1768        | 0.6012        |
|          |                                 | c->v; exp <sub>cross</sub> | 105        | 517           | 0.1502        | 0.5663        |
|          | LIAR-RAW                        | c->v                       | 1,427      | 2,437         | 0.6731        | 0.6020        |
|          |                                 | c->v; exp <sub>self</sub>  | 895        | 2,207         | 0.3809        | 0.5781        |
|          |                                 | c->v; exp <sub>cross</sub> | 1,304      | 3,600         | <b>0.9546</b> | 0.7497        |
| AVeriTeC | c; evi ->v; exp <sub>self</sub> | 98                         | 1,746      | 0.1033        | 0.9423        |               |
| Qwen-3   | RAW-FC                          | c->v                       | 374        | 553           | 0.5351        | 0.6057        |
|          |                                 | c->v; exp <sub>self</sub>  | 161        | 502           | 0.2268        | 0.5565        |
|          |                                 | c->v; exp <sub>cross</sub> | 156        | 482           | 0.2152        | 0.5434        |
|          | LIAR-RAW                        | c->v                       | 903        | 1,893         | 0.3671        | 0.5105        |
|          |                                 | c->v; exp <sub>self</sub>  | 680        | 1,473         | 0.2388        | 0.4435        |
|          |                                 | c->v; exp <sub>cross</sub> | 681        | 1,514         | 0.2426        | 0.4505        |
| AVeriTeC | c; evi ->v; exp <sub>self</sub> | 35                         | 721        | <b>0.0174</b> | <b>0.9092</b> |               |

Table 15: The explanation length for ablation studies before and after S-EGS for improved variants. V denotes vertical steering, and H denotes horizontal steering.

| Backbone   | Pair               | RAW-FC  | LIAR-RAW | AVeriTeC |
|------------|--------------------|---------|----------|----------|
| LLaMA-2-7b | Baseline           | 787.55  | 118.06   | 23.72    |
|            | Dual w/exp – cross | 278.42↓ | –        | –        |
|            | Dual w/exp – self  | 286.63↓ | 67.29↓   | 23.16↓   |
|            | Single w/exp – V   | 274.69↓ | 76.82↓   | 23.38↓   |
|            | Single w/exp – H   | 264.64↓ | 76.50↓   | –        |
| Qwen-3-7b  | Baseline           | 997.04  | 290.03   | 24.37    |
|            | Dual w/exp – cross | 306.78↓ | 84.88↓   | –        |
|            | Dual w/exp – self  | 293.92↓ | –        | 24.41↑   |
|            | Single w/exp – V   | 310.00↓ | 83.85↓   | 24.23↓   |
|            | Single w/exp – H   | –       | 83.98↓   | 22.89↓   |

| Backbone | Stage                          | $x \rightarrow y$              | Macro-F1 | $\Delta$ Gain | Bootstrap FS Mean | Bootstrap $\Delta$ Gain | 95% CI Width/2 |
|----------|--------------------------------|--------------------------------|----------|---------------|-------------------|-------------------------|----------------|
| LLaMA-2  | BASE                           | $c \rightarrow v$              | 35.61    | —             | 35.46             | —                       | $\pm 5.31$     |
|          |                                | $c; evi \rightarrow v$         | 27.08    | —             | 28.08             | —                       | $\pm 6.79$     |
|          |                                | $c \rightarrow v; exp_{cross}$ | 34.41    | —             | 34.16             | —                       | $\pm 6.42$     |
|          |                                | $c \rightarrow v; exp_{self}$  | 31.68    | —             | 31.54             | —                       | $\pm 5.14$     |
| SFT      | $c \rightarrow v$              | 26.44                          | -9.17    | 26.31         | -9.15             | $\pm 5.73$              |                |
|          | $c; evi \rightarrow v$         | 44.85                          | 17.77    | 45.01         | 16.93             | $\pm 6.75$              |                |
|          | $c \rightarrow v; exp$         | 60.59                          | 26.18    | 60.42         | 26.26             | $\pm 6.62$              |                |
| S-EGS    | $c \rightarrow v$              | 31.47                          | 5.03     | 31.25         | 4.94              | $\pm 6.82$              |                |
|          | $c \rightarrow v; exp_{cross}$ | 64.99                          | 4.40     | 64.84         | 4.42              | $\pm 6.57$              |                |
|          | $c \rightarrow v; exp_{self}$  | 61.81                          | 1.22     | 61.64         | 1.22              | $\pm 6.60$              |                |
| Qwen-3   | BASE                           | $c \rightarrow v$              | 46.54    | —             | 46.36             | —                       | $\pm 7.03$     |
|          |                                | $c; evi \rightarrow v$         | 46.23    | —             | 46.00             | —                       | $\pm 7.09$     |
|          |                                | $c \rightarrow v; exp_{cross}$ | 46.66    | —             | 46.53             | —                       | $\pm 6.87$     |
|          |                                | $c \rightarrow v; exp_{self}$  | 48.86    | —             | 48.72             | —                       | $\pm 6.95$     |
| SFT      | $c \rightarrow v$              | 41.67                          | -4.87    | 41.46         | -4.90             | $\pm 6.74$              |                |
|          | $c; evi \rightarrow v$         | 63.17                          | 16.94    | 62.99         | 16.99             | $\pm 6.56$              |                |
|          | $c \rightarrow v; exp$         | 58.35                          | 9.49     | 58.19         | 9.47              | $\pm 6.85$              |                |
| S-EGS    | $c \rightarrow v$              | 41.69                          | 0.02     | 41.52         | 0.06              | $\pm 6.81$              |                |
|          | $c \rightarrow v; exp_{cross}$ | 59.39                          | 1.04     | 59.23         | 1.04              | $\pm 6.80$              |                |
|          | $c \rightarrow v; exp_{self}$  | 58.86                          | 0.51     | 58.70         | 0.51              | $\pm 6.80$              |                |

Table 16: Bootstrap resampling results for the first ablation study on RAW-FC.

(2) For the larger decreases (-1.05 to -1.17): beyond the explanation in Section 4.2.1, the **AVeriTeC** dataset decomposes reasoning into a dialogue format that closely matches our Stage-1 training pipeline. This effectively represents an upper bound of our optimization setting. In addition, the decoupled AVeriTeC test set is inherently simpler than the other datasets; even SFT outperforms the others by roughly 20–40% according to Table 3. Strictly speaking, **horizontal steering** vectors are extracted only from SFT models and are not specific to REFLEX. Consequently, horizontal steering fails twice on AVeriTeC, whereas vertical steering fails only once, which further illustrates the robustness of REFLEX.

(3) Smaller gains on **Qwen-3** can be explained as follows. First, as discussed in Section 5.2, we observe a divergence amplification effect. Except for the Dual w/o exp setting, which contains limited information and therefore shows only minor backbone differences, Qwen-3 exhibits smaller divergences than LLaMA-2 in Dual w/exp and Single w/exp settings, which may explain its reduced gains and weaker transferability. Second, as noted in Appendix C, Qwen-3 shows severe overfitting, with training Macro-F1 exceeding test performance by 20–40%. Although we selected checkpoints to reduce overfitting, a gap of roughly 10% remains. Overfitting in the Qwen-3 series is also commonly reported on reasoning tasks.

To further demonstrate generalizability, inspired by one reviewer’s observation on vertical and horizontal steering, we also ran experiments on Mistral-v0.1 under the w/exp setting on RAW-FC. As shown in Table 20, REFLEX still achieves an 8.88-point gain under vertical steering.

## Q More Analysis of Disentanglement Effectiveness

To further support disentanglement beyond the main results on empirical validation and representation structure, we provide additional analyses from complementary perspectives, including model directions, optimal layers, and case studies.

For model directions, we take REFLEX after S-EGS as the baseline (denoted as style|fact), and test three additional settings on Dual w/exp variants: anchoring (1) fully on fact (positives: correct verdicts, negatives: incorrect ones), (2) on the base model (positives: backbone outputs, negatives: SFT outputs), and (3) on the SFT model (the reverse of (2)). As shown in Table 21, the blue regions are dominated by red ones, further motivating our disentanglement design. Most blue regions appear on LIAR-RAW, likely due to the recency bias discussed in Section 4.2.1.

For optimal layers, we present the distribution of disentangled vectors across the two backbone models. Figure 9 shows that, for Dual w/o exp pairs, the largest divergences emerge in early lay-

Table 17: Full cross-model transfer results in Macro-F1.  $\text{Corr}\Delta=0.78$  and  $\text{Corr}(\text{Source Macro-F}, \text{Target } \Delta)=0.88$ .  $\text{Corr.}$  denotes the Pearson correlation coefficient. V denotes vertical steering, and H denotes horizontal steering.

| Model   | Source   | Variant            | Macro-F | $\Delta$ | Target   | Macro-F | $\Delta$ |
|---------|----------|--------------------|---------|----------|----------|---------|----------|
| LLaMA-2 | RAW-FC   | Dual w/ exp.       | 64.99   | +4.40    | Liar-RAW | 50.59   | +7.54    |
|         |          | Single w/ exp. - V | 62.17   | +1.58    |          | 49.81   | +6.76    |
|         |          | Single w/ exp. - H | 62.64   | +2.05    |          | 50.11   | +7.06    |
|         | Liar-RAW | Dual w/ exp.       | 43.06   | +0.01    | RAW-FC   | 47.20   | -13.39   |
|         |          | Single w/ exp. - V | 43.61   | +0.56    |          | 51.77   | -8.82    |
|         |          | Single w/ exp. - H | 43.73   | +0.68    |          | 52.72   | -7.87    |
| Qwen-3  | RAW-FC   | Dual w/ exp.       | 59.39   | +1.04    | Liar-RAW | 45.12   | -1.61    |
|         |          | Single w/ exp. - V | 58.88   | +0.53    |          | 45.28   | -1.45    |
|         | Liar-RAW | Dual w/ exp.       | 47.13   | +0.40    | RAW-FC   | 41.04   | -17.31   |
|         |          | Single w/ exp. - V | 46.80   | +0.07    |          | 40.40   | -17.95   |
|         |          | Single w/ exp. - H | 47.04   | +0.31    |          | 40.40   | -17.95   |

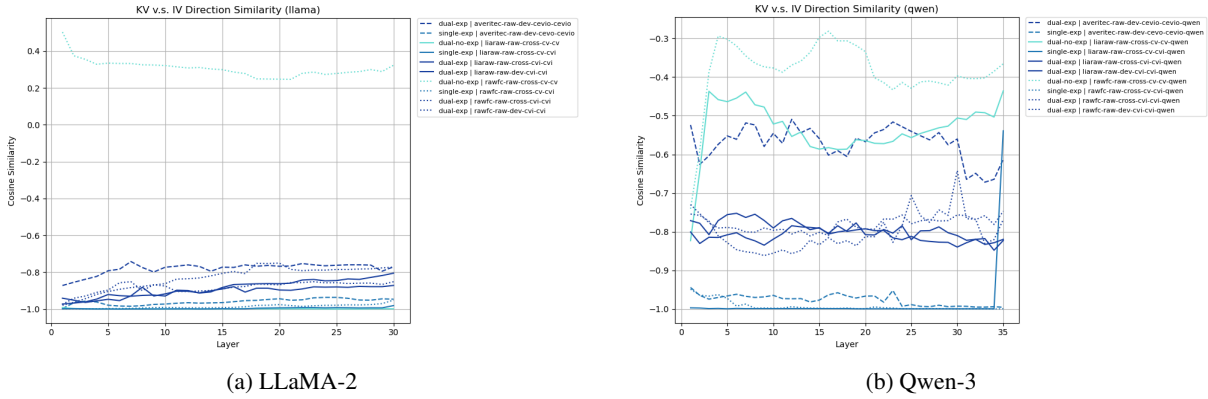


Figure 8: Full direction similarity of KV and IV at each layer based on cosine similarity. Colors denote **Dual w/o exp**, **Dual w/exp**, and **Single w/exp**. The solid line denotes LIAR-RAW, the dashed line AVeriTeC, and the dotted line RAW-FC.

Table 18: Optimal steering hyperparameters across datasets on LLaMA-2.

| Dataset          | LLaMA-2          |                |        |       |            |
|------------------|------------------|----------------|--------|-------|------------|
|                  | Pairs            | Target         | Vector | Layer | Multiplier |
| RAW-FC           | Dual w/exp-cross | w/exp          | IV     | 10    | 1.5        |
|                  | Dual w/exp-self  | w/exp          | IV     | 13    | 1.5        |
|                  | Single w/exp - V | w/exp          | IV     | 2     | 1          |
|                  |                  | w/o exp        | IV     | 1     | 1.5        |
|                  | Dual w/o exp     | w/o exp        | IV     | 1     | 1.5        |
| Single w/exp - H | w/exp            | -              | 2      | -1    |            |
|                  | w/o exp          | -              | 1      | 1.5   |            |
|                  | Dual w/exp-self  | <b>w/exp</b>   | KV     | 1     | <b>1.5</b> |
| LIAR-RAW         |                  | <b>w/o exp</b> | KV     | 1     | <b>1</b>   |
|                  | Single w/exp - V | w/exp          | IV     | 1     | 1          |
|                  |                  | w/o exp        | IV     | 2     | 1          |
|                  | Single w/exp - H | w/exp          | -      | 1     | 1          |
| w/o exp          |                  | -              | 2      | -1    |            |
| AVeriTeC         | Dual w/exp-self  | w/exp          | KV     | 11    | -1.5       |
|                  | Single w/exp - V | w/exp          | IV     | 1     | 1          |

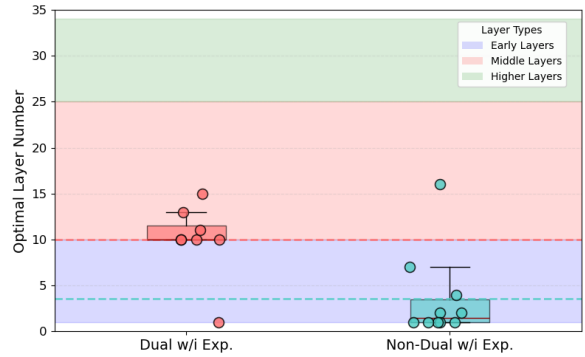


Figure 9: Optimal layers for improving different pair constructions across transformer layers. Square brackets denote optimal components.

Table 19: Optimal steering hyperparameters across datasets on Qwen-3.

| Dataset          | Qwen3            |                |        |       |             |
|------------------|------------------|----------------|--------|-------|-------------|
|                  | Pairs            | Target         | Vector | Layer | Multiplier  |
| RAW-FC           | Dual w/exp-cross | w/exp          | IV     | 10    | 1.5         |
|                  | Dual w/exp-self  | w/exp          | IV     | 10    | 1.5         |
|                  | Single w/exp - V | w/exp          | IV     | 16    | 1.5         |
|                  |                  | w/o exp        | KV     | 13    | -1.5        |
|                  | Dual w/o exp     | <b>w/o exp</b> | IV     | 17    | <b>-1.5</b> |
| Single w/exp - H | w/o exp          | -              | 13     | -1.5  |             |
| LIAR-RAW         | Dual w/exp-cross | w/exp          | KV     | 15    | -1.5        |
|                  | Dual w/o exp     | w/o exp        | KV     | 2     | -1          |
|                  | Single w/exp - V | w/exp          | IV     | 4     | 1.5         |
|                  |                  | w/o exp        | IV     | 6     | 1.5         |
|                  | Single w/exp - H | w/exp          | -      | 1     | 1.5         |
| w/o exp          | -                | 6              | 1.5    |       |             |
| AVeriTeC         | Dual w/exp-self  | <b>w/exp</b>   | KV     | 10    | <b>1.5</b>  |
|                  | Single w/exp - V | <b>w/exp</b>   | IV     | 7     | <b>-1.5</b> |
|                  |                  | <b>w/o exp</b> | IV     | 3     | <b>-1.5</b> |
|                  | Single w/exp - H | w/exp          | -      | 1     | 1.5         |
|                  |                  | w/o exp        | -      | 3     | 1.5         |

| Backbone     | Direction | SFT   | S-EGS | $\Delta$ Gain |
|--------------|-----------|-------|-------|---------------|
| Mistral-v0.1 | V         | 42.03 | 50.91 | 8.88          |
|              | H         | 42.03 | 50.36 | 8.33          |

Table 20: Macro-F1 scores under different steering directions on the Mistral-v0.1 backbone. V denotes vertical steering and H denotes horizontal steering.

ers (1–5), whereas for pairs with full explanations, divergences peak in the **middle layers** (10–20). This aligns with prior findings on transformer interpretability: early layers capture lexical and topical signals, middle layers encode stylistic and syntactic patterns, and higher layers capture more conceptual information (Yun et al., 2021). Interestingly, this middle-layer dominance appears across both subjective styles (e.g., sycophancy and myopic reward) and objective fact-related phenomena (e.g., factuality vs. hallucination), and from common misconceptions to fine-grained fact-checking misinformation. This pattern suggests that contrastive signals mainly capture surface-level stylistic patterns of reasoning or knowledge representation.

For the case studies, we render cosine similarities between unrefined output tokens and optimal vectors in HTML. Red tokens denote alignment with the optimal vector direction, while blue tokens denote the opposite. As shown in Figure 7, red regions tend to capture correct verdict-related content, whereas blue regions are dominated by noisy or redundant syntax patterns. This further supports that REFLEX disentangles surface-level veracity signals from explanation style.

Table 21: Macro-F1 scores for model-direction experiments. Red indicates effective settings, and blue indicates ineffective settings.

| Backbone | Variant | Direction    | RAW-FC       | LIAR-RAW | AVeriTeC |       |
|----------|---------|--------------|--------------|----------|----------|-------|
| LLaMA-2  | self    | -> stylefact | 61.81        | 43.06    | 84.61    |       |
|          |         | -> fact      | 58.06        | 42.91    | 83.35    |       |
|          |         | -> base      | 60.67        | 42.70    | 83.35    |       |
|          |         | -> sft       | 60.67        | 43.33    | 83.35    |       |
|          | cross   | -> stylefact | 64.99        | 42.77    | -        |       |
|          |         | -> fact      | 61.66        | 43.95    | -        |       |
|          |         | -> base      | 64.47        | 42.93    | -        |       |
|          |         | -> sft       | 64.47        | 42.85    | -        |       |
|          | Qwen-3  | self         | -> stylefact | 58.86    | 46.53    | 88.21 |
|          |         |              | -> fact      | 58.79    | 46.73    | 88.02 |
| -> base  |         |              | 58.88        | 46.64    | 88.02    |       |
| -> sft   |         |              | 57.86        | 46.79    | 87.51    |       |
| cross    |         | -> stylefact | 59.39        | 47.13    | -        |       |
| -> fact  | 57.85   | 46.86        | -            |          |          |       |
| -> base  | 58.35   | 46.57        | -            |          |          |       |
| -> sft   | 57.85   | 46.57        | -            |          |          |       |

## R Automatic Evaluations Results for Ablative Studies

The automatic evaluation results for explanations of improved variants are shown in Table 22. We also computed the standard errors of the decreases/gains for each dimension—Misleadingness, Informativeness, Soundness, and Readability—which are 0.0325, 0.0160, 0.0127, and 0.0189, respectively. Since the standard error for the decrease in Misleadingness is relatively high, and this dimension is particularly important, we further introduce pairwise human evaluation to compensate for this instability.

## S More Analysis of Statistical Correlations

As shown in Figure 6, three findings emerge: **(1) Performance is primarily driven by faithfulness-related metrics:** both F-score and accuracy are strongly negatively correlated with misleadingness (-0.85, -0.96) and strongly positively correlated with soundness (0.91, 0.78). **(2) Readability plays a secondary but consistent role:** it correlates positively with both F-score (0.67) and accuracy (0.81), suggesting that clearer explanations tend to align with better predictions. **(3) Informativeness exhibits a clear trade-off with effectiveness.** It is negatively correlated with readability (-0.35) and accuracy (-0.04), while positively correlated with both misleadingness (0.26) and soundness (0.57), indicating that adding more background information can make explanations appear more logically grounded, but often introduces additional noise that

Table 22: Explanation quality for all improved variants. Red background denotes improvement, while blue denotes decline. The omitted bar indicates that Macro-F1 did not improve in Tables 3 and 5.

| Backbone   | Pair                | RAW-FC |      |      |             | LIAR-RAW |      |      |             | AVeriTeC |      |      |      |
|------------|---------------------|--------|------|------|-------------|----------|------|------|-------------|----------|------|------|------|
|            |                     | M↓     | I    | S    | R           | M↓       | I    | S    | R           | M↓       | I    | S    | R    |
| LLaMA-2-7b | SFT                 | 1.90   | 4.78 | 4.82 | 4.55        | 1.90     | 4.48 | 4.60 | 4.65        | 1.21     | 4.61 | 4.89 | 4.86 |
|            | Dual - cross        | 2.00   | 4.89 | 4.83 | 4.81        | –        | –    | –    | –           | –        | –    | –    | –    |
|            | Dual - self         | 1.91   | 4.80 | 4.77 | 4.75        | 1.80     | 4.50 | 4.63 | 4.83        | 1.18     | 4.63 | 4.89 | 4.88 |
|            | Single - vertical   | 1.95   | 4.87 | 4.84 | <b>4.86</b> | 1.77     | 4.58 | 4.66 | 4.83        | 1.18     | 4.65 | 4.86 | 4.89 |
|            | Single - horizontal | 1.79   | 4.88 | 4.83 | 4.80        | 1.77     | 4.54 | 4.67 | 4.84        | –        | –    | –    | –    |
| Qwen-3-8b  | SFT                 | 1.89   | 4.74 | 4.80 | 4.32        | 1.99     | 4.43 | 4.55 | 4.22        | 1.10     | 4.67 | 4.89 | 4.89 |
|            | Dual - cross        | 1.83   | 4.87 | 4.75 | 4.82        | 1.83     | 4.53 | 4.64 | <b>4.83</b> | –        | –    | –    | –    |
|            | Dual - self         | 1.87   | 4.89 | 4.81 | 4.81        | –        | –    | –    | –           | 1.11     | 4.63 | 4.89 | 4.90 |
|            | Single - vertical   | 1.89   | 4.89 | 4.83 | 4.75        | 1.80     | 4.55 | 4.63 | 4.82        | 1.10     | 4.67 | 4.91 | 4.90 |
|            | Single - horizontal | –      | –    | –    | –           | 1.84     | 4.54 | 4.63 | 4.82        | 1.13     | 4.70 | 4.92 | 4.92 |

reduces clarity and hurts decision quality.

## T Case Studies

REFLEX not only corrects cases in which both the explanation and the verdict are completely wrong, as shown below, but also reduces inconsistencies with the verdict and self-contradictions within explanations. These issues are often coupled, but can also appear independently. Manual case studies further show that REFLEX reduces irrelevant information and overconfident hallucinations in explanations.

[Claim]

Canada would like to reopen the border with the U.S.

[Evidences]

(1) When were border crossing restrictions between Canada and the USA first introduced? The Prime Minister, Justin Trudeau, announced on March 18 2020 that Canada and the United States have agreed to temporarily restrict all non-essential travel across the Canada-U.S. border. The new measures will take effect at 12:01 a.m. EDT on Saturday, March 21, 2020, and end at 12:00 p.m. EDT on Tuesday, April 21, 2020. At that time, Canada and the United States will review the measures. The Prime Minister, Justin Trudeau, announced on March 18 2020 that Canada and the United States have agreed to temporarily restrict all non-essential travel across the Canada-U.S. border. The new measures will take effect at 12:01 a.m. EDT on Saturday, March 21,

2020, and was set to end at 12:00 p.m. EDT on Tuesday, April 21, 2020.

(2) Were Canadian/American border crossings still restricted on 18/09/2020? Yes.

(3) Why were Canada/USA border restrictions extended on 18/09/2021? The Government of Canada 2019s top priority remains the health and safety of Canadians, and this extension will continue to protect people on both sides of the border while ensuring the flow of essential goods and services between our two countries. In order to limit the further spread of coronavirus. Canada’s Health Minister Patty Hajdu mentioned in an interview that she was happy with the border closure because Canada already has its hands full with the pandemic. The Government of Canada 2019s top priority remains the health and safety of Canadians, so this extension will continue to protect people on both sides of the border.

(4) Were Canadian public happy that the boarder was still closed? 85% of Canadians would like the borders with the U.S. to remain closed until at least the end of 2020.

[Label]

REFUTED

[SFT]

from AVeriTeC, Qwen-3, Dual-wi/exp

[Anonymous as B in Judgement-1 and swapped as A in Judgement-2]

Verdict: SUPPORTED. Explanation: The assistant incorrectly labels the claim as supported, despite the evidence showing

that the border was closed and most Canadians wanted it to stay closed. This misrepresents the claim's veracity.

[REFLEX] KV-layer10-multiplier1.5  
[Anonymous as A in Judgement-1 and swapped as B in Judgement-2]  
Verdict: REFUTED. Explanation: The evidence shows that the Canada-U.S. border remained closed and that the majority of Canadians preferred it to stay closed. Therefore, the claim that Canada would like to reopen the border is contradicted and refuted.

Judgment-1: Assistant A correctly identifies that the claim is refuted, aligning with the true label. Assistant B incorrectly labels it as supported, which is misleading because it **contradicts the evidence**. Therefore, Assistant B's explanation is more misleading.

Judgment-2: Both assistants discuss the border status and public opinion. The claim is that Canada would like to reopen the border, which is directly contradicted by the evidence. Assistant A correctly labels the claim as refuted, while Assistant B incorrectly supports it. Assistant B's explanation **misrepresents the evidence** and is therefore more misleading.

## T.1 Misleading Verdict Inconsistency

[Claim]  
Pre-existing conditions are covered under my (health care) plan.  
[Label]  
FALSE

[SFT]  
from LLaMA-2, Liar-RAW, Dual wi/exp.-self  
[Anonymous as B in Judgement-1 and swapped as A in Judgement-2]  
Verdict: FALSE. Explanation: But the plan **doesn't cover** pre-existing conditions. It **only covers** pre-existing conditions that existed before the plan was purchased.

[REFLEX] KV-layer1-multiplier1.5  
[Anonymous as A in Judgement-1 and swapped as B in Judgement-2]  
Verdict: FALSE. Explanation: But the law does not require them to do so. And the law specifically excludes coverage of pre-existing conditions before the effective date of the new plan.

Judgment-1: Both assistants correctly identify the claim as FALSE, which aligns with the true label. However, Assistant A's explanation clearly states that the law does not require coverage of pre-existing conditions and that such conditions are specifically excluded before the effective date of the new plan. This is a straightforward and accurate explanation supporting the false verdict. Assistant B's explanation is **confusing and somewhat contradictory**. It says the plan doesn't cover pre-existing conditions, which is correct, but then adds It only covers pre-existing conditions that existed before the plan was purchased, which contradicts the first part and is misleading because it implies some coverage of pre-existing conditions, which is not true. Therefore, Assistant B's explanation is more misleading due to the contradictory and unclear statement about coverage.

Judgment-2: Both assistants correctly identify the claim as FALSE, which aligns with the true label. However, Assistant A's explanation is somewhat **confusing and contradictory**: it states the plan doesn't cover pre-existing conditions, but then says It only covers pre-existing conditions that existed before the plan was purchased, which is logically inconsistent and misleading. Assistant B's explanation is clearer and more accurate, noting that the law does not require coverage of pre-existing conditions and that such conditions are excluded before the plan's effective date, which aligns with typical health insurance policies. Therefore, Assistant A's explanation is more misleading due to its contradictory phrasing.

## T.2 Intra-Explanation Self-Contradictions

[Claim]

Health care reform would **bring down the deficit** by as much as \$1 trillion over the next two decades.

[Label]

HALF-TRUE

[SFT]

from Liar-Raw, Qwen-3, Dual-wi/exp

[Anonymous as B in Judgement-1 and swapped as A in Judgement-2]

Verdict: HALF-TRUE. Explanation: Obama said health care reform "would bring down the deficit by as much as \$1 trillion over the next two decades." The CBO has estimated that the health care bill would **reduce the deficit** by \$1 trillion over 10 years. But the CBO has also said that the bill would **increase the deficit by \$100 billion over the next 10 years**. Obama's statement is partially accurate but leaves out important details.

[REFLEX] KV-layer15-multiplier-1.5

[Anonymous as A in Judgement-1 and swapped as B in Judgement-2]

Verdict: HALF-TRUE. Explanation: Obama said health care reform "would bring down the deficit by as much as \$1 trillion over the next two decades." The CBO has said that the health care bill would reduce the deficit by \$143 billion over 10 years. The CBO has also said that the health care bill would reduce the deficit by \$1.7 trillion over 20 years. Obama's statement is partially accurate but leaves out important details.

Judgment-1: Assistant A provides CBO estimates (\$143 billion over 10 years, \$1.7 trillion over 20 years) that support the claim being HALF-TRUE, noting omitted details. Assistant B presents **contradictory information** (\$1 trillion reduction over 10 years and a \$100 billion increase over 10 years), which is confusing and misleading. Therefore, Assistant B's explanation is more misleading.

Judgment-2: Both assistants label the claim as HALF-TRUE. However, Assistant A contains a **contradictory statement** about CBO estimates (both \$1 trillion reduction and \$100 billion increase over 10 years), which could mislead the reader. Assistant B provides consistent figures (\$143 billion over 10 years, \$1.7 trillion over 20 years) that align better with the claim. Therefore, Assistant A's explanation is more misleading.

## T.3 Intra-Explanation Irrelevant Information

[Claim]

U.S. actor Tom Hanks is related to "Mister Rogers' Neighborhood" host Fred Rogers.

[Label]

TRUE

[SFT]

from Qwen-3, RAW-FC, Dual-wi/exp-self  
[Anonymous as B in Judgement-1 and swapped as A in Judgement-2]

Verdict: TRUE. Explanation: In June 2020, as the world mourned the death of Fred Rogers, host of the long-running children's television show "Mister Rogers' Neighborhood," a rumor started to circulate on social media that the late entertainer was related to actor Tom Hanks. The rumor was largely spread by a Facebook user who posted a meme containing a photograph of Hanks and Rogers along with the caption: "Did you know that Tom Hanks is related to Fred Rogers?" The claim that Hanks and Rogers were related was also spread by a number of other social media users, as well as in a viral tweet by the comedian Patton Oswalt. The claim that Hanks and Rogers were related is true. **Hanks' paternal grandmother, Lucille, was the sister of Rogers' paternal grandmother, Marion.** Hanks confirmed the familial connection in a 2018 interview with the "Today" show. Hanks said that he first learned about the connection when he was a child. He also said that he was a big fan of Rogers' show and that he would watch it with his sister. "I was a big fan of Mr. Rogers," Hanks said. "I loved that show. I would watch

it with my sister. I knew that we were related. I knew that we were related. I knew that we were related.” Hanks’ paternal grandmother, Lucille, was the sister of Marion, the mother of Fred Rogers. Hanks’ great-great-grandmother, Mary, was the sister of Rogers’ great-great-grandmother, Ann. The two men were related by blood, but they were also connected by their careers. Hanks played the role of “Mr. Hand” in the 2003 film “A Walk to Remember.” In the movie, “Mr. Hand” is a children’s television host who bears a passing resemblance to Rogers. Hanks also lent his voice to the character of “Mr. Hand” in the 2015 animated film “The SpongeBob Movie: Sponge Out of Water.” Here’s a clip from “A Walk to Remember” featuring Hanks as “Mr. Hand”: Fred Rogers passed away on Feb. 27, 2003. Tom Hanks is still alive and well.

[REFLEX] IV-layer10-multiplier1.5

[Anonymous as A in Judgement-1 and swapped as B in Judgement-2]

Verdict: TRUE. Explanation: In June 2020, as the world mourned the death of Fred Rogers, host of the long-running children’s television show “Mister Rogers’ Neighborhood,” a rumor started to circulate on social media that the late entertainer was related to actor Tom Hanks. The claim was made in a tweet posted by the “Mister Rogers’ Neighborhood” Facebook page, which stated that Rogers was Hanks’ “great-great-great uncle.” The claim was also shared in a meme posted to the “Mister Rogers’ Neighborhood” Instagram page. The claim that Hanks and Rogers were related was true. Hanks’ mother, Nancy, was Rogers’ first cousin, once removed. Hanks confirmed the familial connection in a 2018 interview with the “Mister Rogers’ Neighborhood” Facebook page. Hanks said that he first learned about the connection when he was 10 years old. “I was 10 years old and I was watching ‘Mister Rogers’ Neighborhood,’” Hanks said. “And I was like, ‘Who is this guy?’ And then my mom said, ‘That’s your great-great-uncle.’ And I

was like, ‘What?’ And then she said, ‘Yeah, he’s your great-great-uncle.’ And I was like, ‘Wow.’” Hanks’ mother, Nancy, was the daughter of Rogers’ first cousin, Mary, and Rogers’ uncle, James. Hanks’ mother passed away in 2010. In a 2018 interview with the “Mister Rogers’ Neighborhood” Facebook page, Hanks said that he and Rogers never met, but that he always felt a connection to the entertainer. “I always felt like I knew him,” Hanks said. “I always felt like I knew him. And I think that’s because I knew my mom. And I knew my mom loved him. And I knew my mom loved what he stood for. And I knew my mom loved what he did. And I knew my mom loved the way he made people feel.”

Judgment-1: Both assistants correctly identify the claim as TRUE, which aligns with the true label. However, Assistant A provides a clear and accurate explanation of the familial relationship, stating that Tom Hanks’ mother was Fred Rogers’ first cousin once removed, and includes a direct quote from Hanks confirming the connection. This explanation is consistent with known facts and does not introduce unrelated or incorrect information. Assistant B, while also concluding the claim is true, presents a **confusing and inaccurate** family relationship, stating that Hanks’ paternal grandmother was the sister of Rogers’ paternal grandmother, and further complicates the genealogy with additional, unsupported claims. Moreover, Assistant B **introduces irrelevant information** about Hanks playing a character named “Mr. Hand” in unrelated films, which is misleading and unrelated to the claim about their familial connection. This distracts from the factual basis and could mislead the reader about the nature of their relationship. Therefore, Assistant B’s explanation is more misleading due to inaccuracies in the family relationship and the inclusion of irrelevant and incorrect details.

Judgment-2: Both Assistant A and Assistant B correctly identify the claim as TRUE, which aligns with the true label. However,

Assistant A provides a detailed but somewhat confused and **inconsistent explanation of the familial relationship**, mentioning multiple relatives and roles that do not clearly establish the connection. The explanation also includes **irrelevant information** about Tom Hanks' acting roles, which could mislead the reader about the nature of the relationship. Assistant B, on the other hand, offers a clearer and more straightforward explanation of the family connection, specifying the relationship through Hanks' mother and providing a plausible narrative of how Hanks learned about it. Although some details in Assistant B's explanation might not be fully verifiable, it stays focused on the familial link without introducing unrelated information. Therefore, Assistant A's explanation is more misleading due to its confusing and extraneous details that do not clearly support the claim.

#### T.4 Intra-Explanation Overconfident Hallucinations

[Claim]

24 of every 100 kids in US schools speak Spanish.

[Evidences]

(1) What percentage of kids in US schools speak Spanish? The percentage of public school students in the United States who were English learners was 10.4% in the fall of 2019 and 9.2% in the fall of 2010. About 6 in 10 young Hispanics say they use Spanish.

[Label]

Not Enough Evidence

[SFT]

from LLaMA-2, AVeriTeC, Dual wi/exp  
[Anonymous as B in Judgement-1 and swapped as A in Judgement-2]

Verdict: Refuted. Explanation: The answer and source show that the amount of children that speak Spanish is **less than 24 in 100** so the claim is refuted.

[REFLEX] KV-layer11-multiplier-1.5

[Anonymous as A in Judgement-1 and swapped as B in Judgement-2]

Verdict: The answer and source show that the statistic is wrong so the claim is refuted.

Judgment-1: Both Assistant A and Assistant B conclude that the claim is refuted, stating that the evidence shows the percentage of kids speaking Spanish is less than 24%. However, the true label is "Not Enough Evidence," indicating that the provided evidence does not sufficiently confirm or deny the claim. Both assistants incorrectly assert a refutation, but Assistant B's explanation is slightly more misleading because it **confidently quantifies** the percentage as "less than 24 in 100," implying a definitive comparison that the evidence does not fully support. Assistant A's explanation is more general, simply stating the statistic is wrong without specifying the extent. Therefore, Assistant B's explanation deviates more from the true label.

Judgment-2: Both Assistant A and Assistant B conclude that the claim is refuted, which is inconsistent with the true label "Not Enough Evidence." The evidence provided does not directly state the percentage of all kids in US schools who speak Spanish; it only provides data on English learners and the proportion of young Hispanics who use Spanish, which is insufficient to confirm or refute the claim definitively. Assistant A attempts to quantify the refutation by stating the number is less than 24 in 100, which is an **assumption not fully supported by the evidence**. Assistant B simply states the statistic is wrong without elaboration. Both explanations mislead by asserting a refutation when the evidence is insufficient to do so. However, Assistant A's explanation is more misleading because it implies a specific numeric conclusion that is not supported, whereas Assistant B's explanation is more general.