

# TELLER: A Trustworthy Framework for Explainable, Generalizable and Controllable Fake News Detection

Anonymous ACL submission

## Abstract

The proliferation of fake news has emerged as a severe societal problem, raising significant interest from industry and academia. While existing deep-learning based methods have made progress in detecting fake news accurately, they often suffer from users' suspicion caused by the non-transparent reasoning processes, poor generalization abilities and inherent risks of integration with large language models (LLMs). To address this challenge, we propose TELLER, a novel framework for trustworthy fake news detection that prioritizes explainability, generalizability and controllability of models. This is achieved via a dual-system framework that integrates cognition and decision systems, adhering to the principles above. The cognition system harnesses human expertise to generate logical predicates, which guide LLMs in generating human-readable logic atoms. Meanwhile, the decision system deduces generalizable logic rules to aggregate these atoms, enabling the identification of the truthfulness of the input news across diverse domains and enhancing transparency in the decision-making process. Finally, we present comprehensive evaluation results on four datasets, demonstrating the feasibility and trustworthiness of our proposed framework.

## 1 Introduction

Fake news has emerged as a prominent social problem due to the rampant dissemination facilitated by social media platforms (Zhou and Zafarani, 2021). Additionally, the swift progress of generative artificial intelligence has further amplified this issue (Cardenuto et al., 2023). While human fact-checking experts can accurately verify the authenticity of news, their efforts cannot scale with the overwhelming volume of online information. Consequently, researchers have turned to automatic fake news detection techniques.

Despite the improved predictive accuracy achieved by current deep learning-based detection

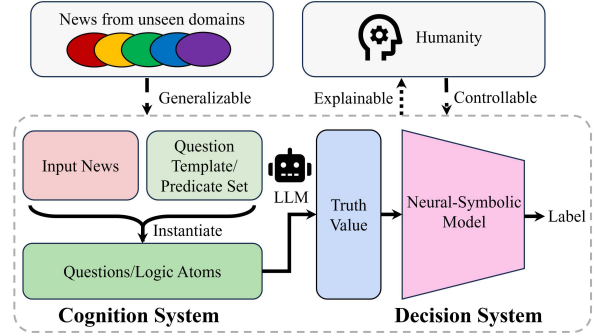


Figure 1: Three crucial aspects of trustworthy fake news detection algorithms and the correlation between these principles and our dual-system framework TELLER.

approaches (Ma et al., 2023; Qi et al., 2021; Mehta et al., 2022), these methods suffer from the lack of transparency because of the black-box nature of neural networks (Cui et al., 2019) and a limited ability to generalize to unseen data, given the inherent diversity of online information (e.g., topics, styles and media platforms) (Liu et al., 2024). Moreover, the increasing integration with LLMs is prone to uncontrollable risks due to hallucinations when LLMs make false conclusions. Thus, a growing awareness emphasizes trustworthiness<sup>1</sup> of these systems (Liu et al., 2023; Sheng et al., 2022).

Unfortunately, the characteristics of a trustworthy fake news detector remain an open question. Hence, based on recent surveys of Trustworthy AI (Li et al., 2023; Jobin et al., 2019) and fake news detection (Shu, 2023), we identify three crucial aspects that go beyond accuracy performance for fake news detection technologies: explainability, generalizability, and controllability. These aspects work collectively to enhance system security and trustworthiness.

Firstly, explainability refers to understanding how an AI model assesses misinformation. The

<sup>1</sup>In AI, trustworthiness refers to the extent to which an AI system can be trusted to operate ethically, responsibly, and reliably (Jobin et al., 2019).

mechanism serves as a fundamental requirement for establishing end-user trust in these tools, as it enables the disclosure of complex reasoning processes and the identification of potential flaws in neural networks. Secondly, generalizability represents the capability to acquire knowledge from limited training data to predict accurately in unseen situations (Wang et al., 2023a). Given the impracticality of exhaustively collecting and annotating vast amounts of data across various news domains, generalization ensures the affordable and sustainable deployment of data-driven fake news detection algorithms. Lastly, controllability encompasses the capacity for human guidance and noise tolerance in the behavior of models (Ji et al., 2023a). This objective benefits models in understanding specific misinformation regulatory policies and rectifying deviations if necessary. While recent practices may satisfy the requirements of explainability (Xu et al., 2022; Liu et al., 2023) or generalization (Kochkina et al., 2018; Yue et al., 2023), they often fail to adhere to all three principles simultaneously.

To this end, we propose TELLER, a Trustworthy framework for Explainable, generalizable and controllable detector, drawing inspiration from the dual-system theory<sup>2</sup> (Daniel, 2017). This framework abstracts the existing pipeline of fake news detection into two components: the cognition and decision systems. As depicted in Fig. 1, the cognition system serves as the first step and is responsible for transforming meaningful human expertise from renowned journalism teams (Tsang, 2023; Sanders, 2023) into a set of Yes/No question templates that correspond to logic predicates. These decomposed questions are then answered using LLMs, which provide truth values for corresponding logic atoms.

On the other hand, the decision system, empowered by a differentiable neural-symbolic model (Cingillioglu and Russo, 2021), can integrate the output of the cognition system to deduce the final authenticity of input news by leveraging domain invariant logic rules learned from data automatically. This visible logic-based ensemble not only mitigates the negative effects caused by inaccurate predictions of LLMs but also allows for the correction of unreasonable rules through adjusting the weights in the model manually to align with human expertise.

<sup>2</sup>System 1 provides tools for intuitive, imprecise, and unconscious decisions akin to deep learning, while system 2 handles complex situations requiring logical and rational thinking akin to symbolic learning (Booch et al., 2021).

Our framework ensures explainability by incorporating human-readable question templates (predicates) and a transparent decision-making process based on logic rules. This interpretability further enables the flexibility to adjust rules and enhances the model’s robustness against false LLM predictions, thereby guaranteeing controllability. Moreover, our model exhibits generalizability, attributed to the generalizable performance of LLMs, combined with reliable human experience as guidance and the utilization of the neural-symbolic model, which can learn domain-generalizable rules.

To summarize, the contributions of this work include: 1) We introduce a systematic framework comprising cognition and decision modules, aiming to uphold three crucial principles for establishing a trustworthy fake news detection system: explainability, generalizability, and controllability. 2) We validate the effectiveness of our framework by conducting comprehensive experiments using various LLMs on four benchmarks. The results demonstrate the feasibility and trustworthiness of TELLER across different scenarios.

## 2 Related Work

### 2.1 Trustworthy AI

Establishing comprehensive trustworthiness in AI is non-trivial due to its multi-objective nature, including robustness, security, transparency, fairness, safety, and ethical standards (Jobin et al., 2019). Achieving such trustworthiness necessitates considering the entire lifecycle of an AI system, spanning from data preparation and algorithm design, development, and deployment to management and governance (Li et al., 2023; Eykholt et al., 2018). Recent researchers have explored diverse approaches to enhance AI trustworthiness across various goals and stages to address this challenge. For example, regarding algorithm design, several topics, such as transfer learning, federated learning, and interpretable AI, have been proposed to improve models’ robustness, security, and transparency. Moreover, the deployment of AI systems necessitates external government oversight, particularly for AGI (Bengio et al., 2023). Although our work focuses on enhancing the trustworthiness of detection systems from the algorithm design aspect, we acknowledge that there is still much room for improvement to achieve the ultimate goal.

## 2.2 Trustworthy Fake News Detection

Recent fake news detection research has witnessed a notable paradigm shift from prioritizing accuracy to considering trustworthiness. In line with our work, we primarily examine studies that aim to enhance algorithms' explainability, generalizability, and controllability.

Regarding explainability, Cui et al. (2019); Xu et al. (2022); Liao et al. (2023) suggested obtaining key evidence for interpretation based on feature importance, while Liu et al. (2023) utilized logic clauses to illustrate the reasoning processing. However, these methods still need to be more transparent due to their probabilistic nature and complex architecture. Furthermore, another group of works, such as Huang and Sun (2023), explored large generative language models (e.g., ChatGPT) and regarded the intermediate chain of thoughts as an explanation. Nevertheless, these explanations may not be reliable due to the hallucination phenomenon (Ji et al., 2023b) and the misalignment problem of AGI (Ji et al., 2023a). Moving on to generalizability, most methods, such as (Yue et al., 2023; Zhu et al., 2023; Kochkina et al., 2018), enhanced fake news detectors through transfer learning algorithms to learn domain-invariant features. However, these methods inevitably introduce external costs of domain alignment, such as annotating domain labels. As for controllability, although some works (Silva et al., 2021; Mendes et al., 2023) incorporated the human-in-loop technique in data sampling and model evaluation, few works explore how to intervene and edit models to align with human expertise.

## 3 Methodology

Formally, given a piece of news  $T$ , the objective of the fake news detection task is to predict its label of truthfulness  $y \in \mathcal{Y}$  where  $\mathcal{Y}$  can fit in different levels of classification granularity. For example, in binary classification setting,  $\mathcal{Y} = \{\text{true}, \text{false}\}$ , and  $T$  is identified as real (fake) when  $y$  is true (false).

As depicted in Fig. 2, TELLER involves two main components: cognition and decision systems. The cognition system decomposes human expertise into Yes/No question templates corresponding to logic predicates. When presented with a new input  $T$ , the templates and predicates can be instantiated accordingly to form questions and logic atoms. By leveraging the parametric knowledge inside LLMs and gathering additional information from exter-

nal tools (such as search engines), the cognition system can generate answers to these questions, represented as truth values of logic atoms. Then, the decision system takes these truth values as input and generates interpretable logic clauses to debunk misinformation through a neural-symbolic model, which can learn generic logic rules from data in an end-to-end manner.

### 3.1 Cognition System

To combat misleading information, existing deep learning-based algorithms fall short in gaining public trust, while fact-checking experts rigorously follow designated guidance and principles to facilitate transparent and fair evaluation. Our cognitive system aims to integrate the strengths of deep learning-based methods that can handle large-scale online information while maintaining the trustworthiness of manual checking.

#### 3.1.1 Predicate Construction

To begin with, we describe the following symbol convention for clarity: calligraphic font  $\mathcal{Q}$  and  $\mathcal{P}$  for sets of question templates and predicates, capitalized letters  $Q, P, X$  for question templates, predicates, and variables, and corresponding lowercase letters  $q, p, x$  for instances of these entities (questions, logic atoms, values). The truth values of logic atoms are denoted by  $\mu$ .

Inspired by the well-established fact-checking process in Table 5, we initially decompose it into a question template set, denoted as  $\mathcal{Q}$ , containing eight questions as detailed in Appendix A.1. Each template  $Q_i$  in  $\mathcal{Q}$  consists of  $N_i$  variables and can be transformed into an  $N_i$ -ary logic predicate  $P_i(X_{i,1}, \dots, X_{i,N_i})$  in  $\mathcal{P}$ . The logic semantics of  $P_i$  is interpreted as the affirmative answer to  $Q_i$  and its truth value  $\mu_i$  represents the probability that  $P_i$  holds. For instance, take  $Q_1$  (i.e., "Background Information:  $X_{1,1}$ . Statement:  $X_{1,2}$ . Is the statement true?") in Fig. 2 as an example. The corresponding predicate  $P_1(X_{1,1}, X_{1,2})$  can be explained as "Given the background information  $X_{1,1}$ , the statement  $X_{1,2}$  is true".

For each predicate  $P_i(X_{i,1}, \dots, X_{i,N_i})$ , we can instantiate the variables  $X_{i,1}, \dots, X_{i,N_i}$  with the actual contents taken from any input news to obtain logic atoms. Since an input piece of news may contain multiple background information and statements (instantiations), we use  $k$  to denote the  $k$ th instantiation where  $1 \leq k \leq \prod_{j=1}^{N_i} |X_{i,j}|$ . Here  $|X_{i,j}|$

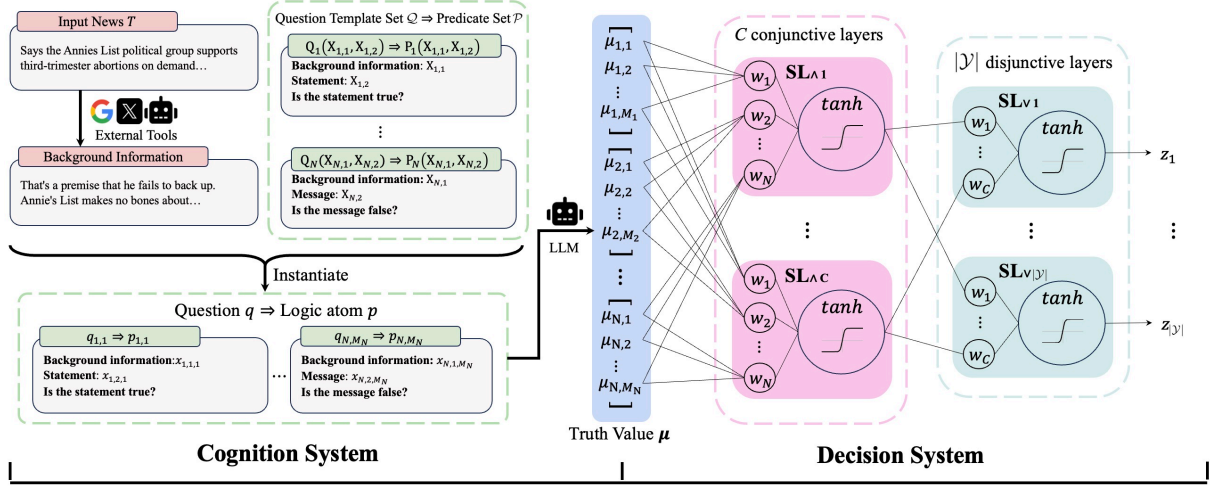


Figure 2: The architecture of the proposed framework TELLER.  $N$  represents the number of question templates (logic predicates),  $M_i$  denotes the number of logic atoms corresponding to the  $i$ th predicate,  $\mathcal{Y}$  denotes the truthfulness label set. The semantics of question templates and logic predicates are described in Table 6.

indicates the total number of possible instantiations for variable  $X_{i,j}$ . Then we denote by  $p_{i,k}$  the instantiated logic atom corresponding to the question  $q_{i,k}$ . Next, we introduce how to acquire the truth value of each logic atom.

### 3.1.2 Logic evaluation with LLMs

While decomposed questions can provide a comprehensive explanation of how the decision is made (Chen et al., 2022; Fan et al., 2020), directly answering these questions poses a challenge due to the impracticality of annotating enormous data to train multiple models for different questions. To address this issue, we resort to the more general-purpose LLMs (e.g., FLAN-T5 (Chung et al., 2022), Llama2 (Touvron et al., 2023), and GPT-3.5) as the foundation for effectively answering these questions. Existing LLMs can be categorized into two groups:  $LLM_{open}$ , such as FLAN-T5 and Llama2, where the logits of output vocabulary can be obtained, and  $LLM_{close}$ , such as GPT-3.5, where the logits are not accessible.

To ensure compatibility with both categories of LLMs, we propose two strategies to obtain the final truth values of logic atoms. Concretely, we first input the question  $q_{i,k}$  with a suffix (i.e., "Yes or No? Response:") to LLMs in order to measure their preference for the affirmative answer "Yes" versus the negative one "No". This preference is subsequently used to compute the truth value of the corresponding logic atom  $p_{i,k}$ .

For  $LLM_{open}$ , we follow (Gallego, 2023; Burns et al., 2023) to obtain pre-softmax logits of "Yes" and "No" tokens, denoted as  $v_{Yes}$  and  $v_{No}$  respectively. Compared with post-softmax logits, pre-

softmax logits can mitigate the influence of other tokens in output vocabulary, particularly when LLMs tend to generate irrelevant tokens that may result in  $v_{Yes}$  or  $v_{No}$  becoming zero. Then the truth value  $\mu$  for the logic atom  $p$  (here we omit the subscript  $i, k$  for ease of illustration) can be obtained as follows:

$$\mu = 2 \frac{e^{v_{Yes}}}{e^{v_{No}} + e^{v_{Yes}}} - 1. \quad (1)$$

For  $LLM_{close}$ , we sample  $m$  times during decoding and count the frequency of "Yes" and "No" responses as  $m_{Yes}$  and  $m_{No}$ . Then we compute

$$\mu = 2 \frac{m_{Yes}}{m_{No} + m_{Yes}} - 1. \quad (2)$$

In either case,  $\mu$  is in the range of  $[-1, 1]$ . When  $\mu \in [-1, 0)$ ,  $\mu \in (0, 1]$ , and  $\mu = 0$ , the corresponding logic atom  $p$  is evaluated as false, true, and unknown, respectively. Once the truth values of all logic atoms for a single predicate  $P_i$  (corresponding to a single question template) are obtained, we concatenate them as one vector, denoted as  $\mu_i$ . Then we concatenate the value vectors for all predicates as the input for the final decision system.

In conclusion, our cognition system can generate diversified questions and logic atoms based on the input news  $T$ . These human-readable entities enhance explainability by showcasing potential intermediate reasoning steps and ensure controllability by allowing adjustments to  $Q$  and  $P$ . Moreover, combining human expertise and LLMs provides the basis for the cognition system's satisfactory generalization performance in unseen domains.

### 3.2 Decision System

After acquiring responses to all questions, it is imperative to develop a decision system to effectively aggregate them to predict the label of the input news  $T$  while preserving trustworthiness in the reasoning process. However, prevalent heuristic strategies (e.g., majority voting) lack the flexibility to handle complex relationships among different questions and cannot tolerate false predictions, and deep-learning-based models cannot be comprehended literally by humans (Wang et al., 2023b).

Hence, we utilize a neural-symbolic model, named Disjunctive Normal Form (DNF) Layer (Cingillioglu and Russo, 2021; Baugh et al., 2023), as our decision system. This model includes conjunctive layers ( $SL_{\wedge}$ ) and disjunctive layers ( $SL_{\vee}$ ), which can progressively converge to symbolic semantics such as conjunction  $\wedge$  and disjunction  $\vee$  respectively during model training. Consequently, this model can automatically learn logic rules from data in an end-to-end manner, capturing generalizable relationships between logic predicates and the target label. As illustrated in Fig. 2, we stack  $C$  conjunctive layers  $SL_{\wedge}$  beneath  $|\mathcal{Y}|$  disjunctive layers  $SL_{\vee}$  to construct the DNF Layer, where each  $SL_{\vee}$  corresponds to a truthfulness label  $y \in \mathcal{Y}$ .

However, the original DNF Layer proposed in (Cingillioglu and Russo, 2021) is not directly applicable to our work due to two issues. Firstly, the truth value of logic atoms  $\mu$  ranges in  $[-1, 1]$ , while the original model can only handle values of  $-1$  and  $1$ . Secondly, each logic atom in the original DNF Layer is treated differently which loses logic semantics where atoms for the same logic predicate should share similar functionality. To address the aforementioned challenges, we propose a modified DNF layer which takes continuous values  $\mu \in [-1, 1]$  as input and assigns the same weight for those atoms instantiated from the same logic predicate. The detailed description of our modified DNF layer can be found in Appendix E.

More concretely, in our proposed DNF Layer, every  $SL_{\wedge}$  takes truth values  $\mu$  of all logic atoms obtained in the cognition system as input, aiming to learn a conjunctive clause  $\text{conj} = \bigwedge_{p_{i,k} \in \mathcal{A}} p_{i,k}$  where  $\mathcal{A} \subseteq \{p_{1,1}, \dots, p_{N,M_N}\}$ , referring to a subset of the complete logic atoms, and outputs the truth value of this conjunctive clause. Subsequently, each  $SL_{\vee}$  receives the truth values of  $C$  conjunctive clauses to represent a disjunction of these conjunctions:  $\bigvee_{c \in \mathcal{C}} \text{conj}_c$  where  $\mathcal{C} \subseteq \{1, \dots, C\}$ , referring

to a subset of all conjs. It then outputs the truth value of this disjunction formula, corresponding to the final probability that the input news  $T$  is identified as the label  $y$ . Hence, each label  $y$  will be associated with a DNF clause learned by the DNF layer. Intuitively, the conjunction simulates the idea that if the input news  $T$  gives affirmative answers to some questions simultaneously, it is highly probable that it should be assigned to label  $y$ . On the other hand, the disjunction provides more flexibility by considering different alternatives (the output is true if at least one of the conj is true) which makes the final decision less sensitive to incorrect atom values due to wrong predictions given by LLMs. For example, assume the learned rules are  $\text{conj}_1 \vee \text{conj}_2$  where  $\text{conj}_1 = p_{1,1} \wedge p_{1,2}$  and  $\text{conj}_2 = p_{2,1} \wedge p_{3,1}$ . Suppose  $\text{conj}_1$  is true, then we can conclude that  $\text{conj}_1 \vee \text{conj}_2$  is true even if  $\text{conj}_2$  gives an incorrect value.

Last but not the least, we apply softmax function to the output of all disjunction layers  $SL_{\vee}$  to obtain the probability  $z \in \mathbb{R}^{|\mathcal{Y}|}$  for all possible labels. The entire decision system can be trained in an end-to-end fashion by minimizing the cross-entropy loss function as below:

$$\mathcal{L} = - \sum_{l=1}^{|\mathcal{Y}|} \mathbb{I}(y_l = y_T) \log z_l, \quad (3)$$

where  $y_T$  represents the ground truth label of  $T$ . During inference, we select the label corresponding to the highest value in  $z$  as the final result.

In summary, our decision system can extract interpretable symbolic rules from data that exhibit robustness across diverse domains and enable intervention by adjusting weights in the DNF Layer to align with prior knowledge (refer to Appendix C).

## 4 Experiments

In this section, we present the experiment setup and demonstrate the feasibility, explainability, generalizability and controllability of TELLER through extensive experiments.

### 4.1 Experimental Setting

**Dataset.** We conducted experiments using four challenging datasets, namely LIAR (Wang, 2017), Constraint (Patwa et al., 2021), Politifact, and GossipCop (Shu et al., 2020). LIAR comprises the binary classification and multi-classification setting with six fine-grained labels for truthfulness ratings. Moreover, Wang (2017); Alhindi et al.

(2018) curated relevant evidence (e.g., background information), serving as gold knowledge in an open setting. Constraint, Politifact and GossipCop are binary classification datasets related to COVID-19, politics, and entertainment domains, respectively. **LLMs.** We select the FLAN-T5 and Llama2 series, which encompass various parameter sizes, as large language models for constructing the cognition system of TELLER because their open-source nature and unrestricted availability can ensure reproducibility in the future. Moreover, we also conduct experiments using GPT-3.5-turbo on the LIAR dataset to examine the versatility of our framework. **Baselines.** We compare our model against *Direct*, *Few-shot Direct*, *Zero-shot COT*, *Few-shot COT*, *Few-shot Logic*. The baselines suffixed with *Direct* involve prompting large language models (LLMs) to predict the label of input news directly; those suffixed with *COT* utilize chain-of-thought techniques to enhance the performance of LLMs; those suffixed with *Logic* replace the thought process in COT with questions paired with their answers. **Implementation Detail.** We evaluate the performance of our framework using the accuracy and Macro-F1, which accommodates class imbalance. For each dataset, we train our decision system using the training split; select the optimal model based on its performance on the validation split; and report the results on the test split. To assess the generalizability of our model, we train our models using the train split from source domains; choose the best model on the validation split of source ones; and report results on the test split from the target domain. Moreover, to highlight the robustness of our framework, we keep all hyperparameters fixed in each setting. The experiment setting and utilized prompts are elaborated thoroughly in Appendix B.

## 4.2 Feasibility Study

To validate the feasibility of our framework, we compare it against multiple baselines across a wide range of LLMs and scenarios (e.g., different classification granularities) in Table 1 and Table 2. These results uncover two crucial findings listed below:

Firstly, our framework demonstrates satisfactory performance in fake news detection tasks. Specifically, in the binary classification setting, TELLER achieves an accuracy of approximately 76% on the GossipCop dataset and over 80% on the other three datasets. Notably, when utilizing Llama 2 (13B) to drive the cognition system, TELLER outperforms all GPT-3.5-turbo based methods by a significant

margin. These results highlight the effectiveness of TELLER in distinguishing between fake and genuine news. In the multi-classification setting on the LIAR dataset, our framework consistently outperforms *Direct* for FLAN-T5 and Llama2 series, even though these models may struggle to discriminate fine-grained labels. This observation underscores the capability of our decision system to mitigate the negative influences of noisy predictions in the cognition system, effectively unleashing the potential of LLMs through logic-based aggregation of answers to decomposed questions.

Secondly, our framework exhibits significant potential for the future. In the binary classification setting across four datasets, TELLER consistently outperforms *Direct* in terms of accuracy and macro-F1 scores by an average of 7% and 6%, respectively. Considering the swift improvement of LLM intelligence, these results imply that the performance of our framework is likely to scale with the evolution of LLMs. Additionally, due to the notable performance difference between closed and open settings on the LIAR dataset, it is promising to integrate external tools to acquire extensive evidence from credible sources, such as official government websites, to enhance the performance of our systems.

## 4.3 Explainability Verification

Explainability is a fundamental factor for establishing trust in AI technology. We demonstrate that our framework satisfies this aspect through its inherent mechanism and the visualization of rules.

Unlike approaches that rely heavily on LLMs, our cognition system incorporates expert knowledge to construct a more well-grounded worldview by generating well-defined question templates and logic predicates. Moreover, our decision system can learn interpretable rules from data to deduce logic clauses to debunk fake news by converging implicit parameters to conjunctive and disjunctive semantics. These symbolic units (e.g., questions and logic atoms) and the interpretable DNF Layer contribute to our framework’s overall explainability and transparency.

However, as the number of conjunctive and disjunctive layers grows, it is difficult for human beings to investigate logic rules derived from our decision system. To address this issue, we propose a strategy to prune unnecessary weights in the DNF Layer. For example, we present the rules extracted from the pruned model for GossipCop in Table 4, where each conjunctive clause identifies one can-

Large Language Models	Method	Binary Classification				Multi-Classification			
		Closed		Open		Closed		Open	
		Acc(%)	Macro-F1(%)	Acc(%)	Macro-F1(%)	Acc(%)	Macro-F1(%)	Acc(%)	Macro-F1(%)
FLAN-T5-small (80M)	Direct	44.99	31.63	45.08	32.41	18.17	9.28	19.51	10.13
FLAN-T5-base (250M)	Direct	54.02	50.79	61.47	61.43	19.43	11.79	21.40	21.40
FLAN-T5-large (780M)	Direct	57.30	52.20	74.38	73.84	19.43	17.84	29.50	24.95
	TELLER	66.83 <sub>(9.53↑)</sub>	<b>66.33</b> <sub>(14.13↑)</sub>	77.76 <sub>(3.38↑)</sub>	77.32 <sub>(3.49↑)</sub>	<b>26.99</b> <sub>(7.55↑)</sub>	18.04 <sub>(0.20↑)</sub>	33.67 <sub>(4.17↑)</sub>	27.50 <sub>(2.55↑)</sub>
	w/ Intervention	<b>65.64</b>	<b>65.12</b>	<b>77.46</b>	<b>77.14</b>	<b>26.28</b>	18.49	35.25	30.05
FLAN-T5-xl (3B)	Direct	58.89	58.62	75.97	75.67	19.67	16.57	29.43	24.74
	TELLER	62.36 <sub>(3.48↑)</sub>	60.18 <sub>(1.56↑)</sub>	78.75 <sub>(2.78↑)</sub>	78.55 <sub>(2.88↑)</sub>	24.31 <sub>(4.64↑)</sub>	17.40 <sub>(0.83↑)</sub>	33.52 <sub>(4.09↑)</sub>	27.22 <sub>(2.48↑)</sub>
	w/ Intervention	63.65	61.82	79.34	79.07	25.57	19.62	34.46	33.59
FLAN-T5-xxl (11B)	Direct	56.41	56.08	75.17	75.15	22.42	18.31	32.18	28.12
	TELLER	66.63 <sub>(10.23↑)</sub>	65.91 <sub>(9.82↑)</sub>	80.24 <sub>(5.06↑)</sub>	79.85 <sub>(4.70↑)</sub>	26.83 <sub>(4.41↑)</sub>	19.68 <sub>(1.36↑)</sub>	35.48 <sub>(3.30↑)</sub>	30.42 <sub>(2.30↑)</sub>
	w/ Intervention	<b>67.03</b>	66.19	80.73	80.41	<b>26.91</b>	<b>21.30</b>	35.88	31.63
Llama2 (7B)	Direct	59.88	59.19	72.29	69.63	18.02	9.97	11.01	6.88
	TELLER	62.46 <sub>(2.58↑)</sub>	62.45 <sub>(3.26↑)</sub>	79.94 <sub>(7.65↑)</sub>	79.80 <sub>(10.16↑)</sub>	23.29 <sub>(5.27↑)</sub>	15.51 <sub>(5.55↑)</sub>	32.73 <sub>(21.72↑)</sub>	25.55 <sub>(18.67↑)</sub>
	w/ Intervention	64.15	62.77	81.93	81.84	23.92	<b>15.14</b>	34.30	27.58
Llama2 (13B)	Direct	56.90	56.90	69.31	63.77	7.32	2.85	10.86	8.25
	Ours	66.04 <sub>(9.14↑)</sub>	66.03 <sub>(9.13↑)</sub>	<b>82.52</b> <sub>(13.21↑)</sub>	<b>82.37</b> <sub>(18.60↑)</sub>	25.81 <sub>(18.49↑)</sub>	17.71 <sub>(14.86↑)</sub>	38.08 <sub>(27.22↑)</sub>	29.27 <sub>(21.02↑)</sub>
	w/ Intervention	<b>67.73</b>	<b>66.97</b>	<b>84.21</b>	<b>84.03</b>	<b>25.10</b>	<b>16.78</b>	38.63	30.60
GPT-3.5-turbo	Direct	42.40	51.48	76.27	74.21	20.46	20.34	26.20	25.12
	TELLER	-	-	79.15 <sub>(2.88↑)</sub>	78.90 <sub>(4.69↑)</sub>	-	-	31.94 <sub>(5.74↑)</sub>	29.53 <sub>(4.41↑)</sub>
	Zero-shot COT	30.88	41.87	72.49	70.83	7.16	9.20	39.81	<b>36.49</b>
	Few-shot	61.67	64.05	81.02	81.00	25.65	<b>25.56</b>	<b>46.81</b>	<b>44.61</b>
	Few-shot COT	52.04	56.15	74.48	76.21	20.69	17.20	<b>45.63</b>	36.36
	Few-shot Logic	49.26	48.85	61.67	60.92	16.37	13.98	20.54	19.22

Table 1: Results on LIRA dataset. "Closed" represents the cognitive system does not have access to any external knowledge source, while "Open" indicates that it can utilize gold evidence collected by human experts. The best results for each setting are highlighted with bold numbers and an underline, whereas sub-optimal results are only highlighted in bold. The **number** indicates that the performance of *w/ Intervention* is worse than TELLER. The number with  $\uparrow$  indicates the performance gain of TELLER over *Direct*.

LLMs	Method	Constraint		Politifact		GossipCop	
		Acc(%)	Macro-F1(%)	Acc(%)	Macro-F1(%)	Acc(%)	Macro-F1(%)
FLAN-T5-large	Direct	78.06	77.97	56.62	54.84	67.43	58.76
	TELLER	80.32 <sub>(2.27↑)</sub>	80.11 <sub>(2.14↑)</sub>	67.65 <sub>(11.03↑)</sub>	67.65 <sub>(12.81↑)</sub>	69.53 <sub>(2.10↑)</sub>	59.39 <sub>(0.63↑)</sub>
	w/ Intervention	80.46	80.31	68.38	68.29	70.28	60.74
FLAN-T5-xl	Direct	75.32	74.79	55.88	50.72	67.73	52.80
	TELLER	83.77 <sub>(8.45↑)</sub>	83.66 <sub>(8.88↑)</sub>	68.82 <sub>(9.14↑)</sub>	64.68 <sub>(13.95↑)</sub>	69.58 <sub>(1.85↑)</sub>	58.72 <sub>(5.91↑)</sub>
	w/ Intervention	83.95	83.88	69.12	68.79	72.23	63.84
FLAN-T5-xxl	Direct	74.80	73.23	52.21	43.65	68.93	52.82
	TELLER	83.39 <sub>(8.59↑)</sub>	83.24 <sub>(10.01↑)</sub>	69.12 <sub>(16.91↑)</sub>	68.57 <sub>(24.92↑)</sub>	69.18 <sub>(0.25↑)</sub>	57.21 <sub>(4.39↑)</sub>
	w/ Intervention	83.62	83.54	69.12	68.95	71.48	62.12
Llama2 (7B)	Direct	81.83	81.73	77.21	77.00	66.78	52.23
	TELLER	83.72 <sub>(1.89↑)</sub>	83.54 <sub>(1.81↑)</sub>	83.82 <sub>(6.62↑)</sub>	83.81 <sub>(6.81↑)</sub>	70.68 <sub>(3.90↑)</sub>	59.58 <sub>(7.35↑)</sub>
	w/ Intervention	85.13	85.04	<b>83.82</b>	<b>83.82</b>	73.38	65.32
Llama2 (13B)	Direct	57.53	51.75	77.94	77.10	52.55	52.27
	TELLER	87.31 <sub>(29.78↑)</sub>	87.29 <sub>(35.53↑)</sub>	79.41 <sub>(1.47↑)</sub>	79.41 <sub>(2.30↑)</sub>	74.48 <sub>(21.93↑)</sub>	66.32 <sub>(14.06↑)</sub>
	w/ Intervention	<b>87.78</b>	<b>87.71</b>	<b>78.68</b>	<b>78.65</b>	<b>75.92</b>	<b>69.30</b>

Table 2: Results on Constraint, Politifact, and GossipCop datasets without access to retrieved background information. The best results for each setting are highlighted with bold numbers. The **number** and the number with  $\uparrow$  have the same meaning as in Table. 1.

didate rule. The pruning algorithm and rules for other datasets are described in Appendix C.

Table 4 can be interpreted as learning DNF rules for both true and false labels of an input news. Specifically, the true label is predicted if either  $\neg \text{conj}_{34}$  or  $\neg \text{conj}_{43}$  is true, i.e., either  $\neg P_2 \wedge P_3 \wedge P_6 \wedge P_8$  or  $P_3 \wedge P_6 \wedge P_8$  is false when removing the negation. Given the semantics of these logic predicates shown in Table 6, we know that  $P_2$ ,  $P_3$  and  $P_8$  check the consistency between the background information and a given message, whereas  $P_6$  scrutinises improper intention from the message alone. On the other hand, the news will be predicted as false if  $\text{conj}_{27}$  is true, i.e.,  $P_4$  is false which means that the background information in the message is neither accurate or objective

according to Table 6.

#### 4.4 Generalizability Verification

Ensuring the generalization ability of fake news decision systems is vital for their sustainable and practical deployment. As observed in Table 3, TELLER consistently outperforms *Direct* across all domains and LLMs without the assistance of any generalization algorithm, while only exhibiting a negligible performance drop in the **GP**  $\rightarrow$  **C** domain using Llama2 7B. This is attributed to the remarkable zero-shot ability of LLMs and the effectiveness of the DNF layer which further compensates for biased predictions made by LLMs through rule-based aggregation. Particularly, the performance gains of TELLER in cross-domain and in-domain exper-

LLMs	Method	CP→G		GP→C		CG→P	
		Acc(%)	Macro-F1(%)	Acc(%)	Macro-F1(%)	Acc(%)	Macro-F1(%)
FLAN-T5-xl	Direct	67.73	52.80	75.32	74.79	55.88	50.72
	TELLER	68.13 <sub>(0.40↑)</sub>	56.54 <sub>(3.74↑)</sub>	82.40 <sub>(7.0↑)</sub>	82.09 <sub>(7.31↑)</sub>	61.76 <sub>(5.88↑)</sub>	60.92 <sub>(10.19↑)</sub>
FLAN-T5-xxl	Direct	68.93	52.82	74.80	73.23	52.21	43.65
	TELLER	69.13 <sub>(0.2↑)</sub>	53.15 <sub>(0.34↑)</sub>	77.44 <sub>(2.64↑)</sub>	76.21 <sub>(2.98↑)</sub>	66.18 <sub>(13.97↑)</sub>	66.17 <sub>(22.52↑)</sub>
Llama2 7B	Direct	66.78	52.23	81.83	81.73	77.21	77.00
	TELLER	68.33 <sub>(1.55↑)</sub>	59.33 <sub>(7.10↑)</sub>	81.60 <sub>(-0.24↓)</sub>	81.04 <sub>(-0.69↓)</sub>	83.09 <sub>(5.88↑)</sub>	82.82 <sub>(5.82↑)</sub>
Llama2 13B	Direct	52.55	52.27	57.53	51.75	77.94	77.10
	TELLER	70.93 <sub>(18.38↑)</sub>	60.90 <sub>(8.63↑)</sub>	85.09 <sub>(27.56↑)</sub>	84.87 <sub>(33.1↑)</sub>	79.41 <sub>(1.47↑)</sub>	79.41 <sub>(2.30↑)</sub>

Table 3: Results on cross-domain experiments. **C**, **P** and **G** represent Constraint, Politifact, and GossipCop datasets.

$\text{conj}_{34} = \neg P_2 \wedge P_3 \wedge P_6 \wedge P_8$
$\text{conj}_{43} = P_3 \wedge P_6 \wedge P_8$
$\text{conj}_{27} = \neg P_4$
$P_{\text{true}} = \neg \text{conj}_{34} \vee \neg \text{conj}_{43}$
$P_{\text{false}} = \text{conj}_{27}$

Table 4: Extracted rules for the GossipCop dataset when using Llama2 (13B)

iments (refer to Table 2) are positively correlated, implying that the decision system manages to learn domain-agnostic rules. Moreover, the Pearson correlation coefficient between these two groups of performance gains shows a substantial improvement from 0.01 to 0.53 when transitioning from the FLAN-T5 series to the more powerful Llama2 series. This finding suggests that leveraging stronger LLMs to drive the cognition system enhances the generalization capability of our framework.

#### 4.5 Controllability Verification

Controllability ensures that fake news detection systems are subject to effective human oversight and intervention. We demonstrate TELLER satisfies this attribute from two aspects. Firstly, we verify the feasibility of manually rectifying rules learned by our decision system that may exhibit irrational behavior. For instance, we observe that  $P_3$  (i.e., "The message contains adequate background information") should have a positive logical relation with  $P_{\text{true}}$  instead of negation in Table 4. To correct this, we perform a manual adjustment by setting the corresponding weight to zero, effectively removing  $P_3$  from the logic rule. However, this modification only leads to a negligible improvement in the test split. Further investigation reveals that the truth value of logic atoms pertaining to  $P_3$  of most real samples is negative, possibly due to the preference of LLMs. This suggests the superiority of our logic-based decision system in reducing the negative effect of incorrect predictions made by LLMs automatically. Secondly, we simulate human experts by intervening in the actions of our cognition system. We achieve this by guiding LLMs to ex-

pand the question template set  $\mathcal{Q}$  using Algorithm 1, referred to as *w/ intervention* in Table 1 and Table 2. The new question template set for intervention is shown in Table 7 in the Appendix. The results consistently indicate that *w/ intervention* outperforms TELLER, highlighting the potential of LLMs as an agency for automatically regulating the behaviors of the cognition system. Consequently, our framework ensures a comprehensive control mechanism by simultaneously facilitating human and AI agents' oversight.

Furthermore, we conduct additional experiments to verify the effectiveness of the DNF Layer within logic formulation over other decision systems, namely decision trees and Naive Bayes classifiers, both of which are conventional machine learning algorithms. We replace the DNF Layers with these two algorithms to derive the final decisions. The results are shown in Table 11 and Table 12 for single-domain and cross-domain settings, respectively in Appendix D.

## 5 Conclusion

In this work, we address the limitations of existing fake news detection methods, which struggle to establish reliability and end-user trust. To tackle this issue, we identify three crucial aspects for constructing trustworthy misinformation detection systems: explainability, generalizability, and controllability. By prioritizing these principles, we propose a dual-system framework TELLER that incorporates cognition and decision systems. To validate our framework's feasibility, explainability, generalizability, and controllability, we conduct extensive experiments on diverse datasets and LLMs. These results affirm the effectiveness and trustworthiness of our approach and highlight its significant potential through evolving both subsystems in the future. While we achieve trustworthiness from an algorithmic perspective, we emphasize the importance of further research to improve the trustworthiness of the entire lifecycle of fake news detection systems.

## Limitations

We identify three main limitations of our work. Firstly, although our framework focuses on enhancing the trustworthiness of fake news detection algorithms, trustworthiness is also influenced by other stages of the AI system lifecycle, such as data collection and deployment. Given the advancements in AI techniques and the importance of online information security, we encourage future research to address the challenges of building trustworthy AI systems comprehensively.

Secondly, as shown in Table 1, integrating external tools to acquire high-quality background knowledge significantly improves the performance of fake news detection systems. However, collecting information that can effectively support detection tasks using such tools is non-trivial due to the complexities of open-domain information retrieval and the diversity of news content. For instance, we search for background information by inputting check-worthy claims of  $P_1$  into a search engine and filter out as much useful information as possible using GPT-3.5-turbo. However, integrating this evidence lead to a slight performance drop on Constraint, Politifact, and GossipCop datasets (Due to page limitations, we do not include this experiment in our paper). Therefore, we leave this for future research.

Thirdly, despite the effectiveness of our decision system, the learning ability and expressiveness of the DNF Layer are limited due to its simple architecture. For example, the DNF Layer learns rules from data without considering the semantics of logic predicates. It may be crucial to develop more powerful decision models to fully unleash the potential of large language models, such as incorporating the semantics of logic predicates.

## Ethics Statement

This paper adheres to the ACM Code of Ethics and Professional Conduct. Specifically, the datasets we utilize do not include sensitive private information and do not pose any harm to society. Furthermore, we will release our codes following the licenses of any utilized artifacts.

Of paramount importance, our proposed dual-system framework serves as an effective measure to combat fake news and safeguard individuals, particularly in the current era dominated by large generative models that facilitate the generation of deceptive content with increasing ease. Moreover,

our approach fulfills explainability, generalizability, and controllability, thereby mitigating concerns regarding the security of AI products and enabling their deployment in real-world scenarios.

## References

- Tariq Alhindi, Savvas Petridis, and Smaranda Muresan. 2018. [Where is your evidence: Improving fact-checking by justification modeling](#). In *Proceedings of the First Workshop on Fact Extraction and VERification (FEVER)*, pages 85–90, Brussels, Belgium. Association for Computational Linguistics.
- Kexin Gu Baugh, Nuri Cingillioglu, and Alessandra Russo. 2023. [Neuro-symbolic rule learning in real-world classification tasks](#). In *Proceedings of the AAAI 2023 Spring Symposium on Challenges Requiring the Combination of Machine Learning and Knowledge Engineering (AAAI-MAKE 2023)*, Hyatt Regency, San Francisco Airport, California, USA, March 27-29, 2023.
- Yoshua Bengio, Geoffrey E. Hinton, Andrew Yao, Dawn Song, Pieter Abbeel, Yuval Noah Harari, Ya-Qin Zhang, Lan Xue, Shai Shalev-Shwartz, Gillian K. Hadfield, Jeff Clune, Tegan Maharaj, Frank Hutter, Atilim Günes Baydin, Sheila A. McIlraith, Qiqi Gao, Ashwin Acharya, David Krueger, Anca D. Dragan, Philip H. S. Torr, Stuart Russell, Daniel Kahneman, Jan Brauner, and Sören Mindermann. 2023. [Managing AI risks in an era of rapid progress](#). *CoRR*, abs/2310.17688.
- Grady Booch, Francesco Fabiano, Lior Horesh, Kiran Kate, Jonathan Lenchner, Nick Linck, Andrea Loreggia, Keerthiram Murugesan, Nicholas Mattei, Francesca Rossi, and Biplav Srivastava. 2021. [Thinking fast and slow in AI](#). In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021*, pages 15042–15046.
- 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, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*.
- João Phillipe Cardenuto, Jing Yang, Rafael Padilha, Renjie Wan, Daniel Moreira, Haoliang Li, Shiqi Wang, Fernanda A. Andalo, Sébastien Marcel, and Anderson Rocha. 2023. [The age of synthetic realities: Challenges and opportunities](#). *CoRR*, abs/2306.11503.
- Jifan Chen, Aniruddh Sriram, Eunsol Choi, and Greg Durrett. 2022. [Generating literal and implied subquestions to fact-check complex claims](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu*

737	Dhabi, United Arab Emirates, December 7-11, 2022,	Antoine Bordes, and Sebastian Riedel. 2020. <a href="#">Gen-</a>	793
738	pages 3495–3516. Association for Computational	<a href="#">erating fact checking briefs</a> . In <i>Proceedings of the</i>	794
739	Linguistics.	<i>2020 Conference on Empirical Methods in Natural</i>	795
740	Kewei Cheng, Nesreen K. Ahmed, and Yizhou Sun.	<i>Language Processing, EMNLP 2020, Online, Novem-</i>	796
741	2023. <a href="#">Neural compositional rule learning for knowl-</a>	<i>ber 16-20, 2020</i> , pages 7147–7161. Association for	797
742	<a href="#">edge graph reasoning</a> . In <i>The Eleventh International</i>	Computational Linguistics.	798
743	<i>Conference on Learning Representations, ICLR 2023,</i>	Yi R. Fung, Christopher Thomas, Revanth Gangi Reddy,	799
744	<i>Kigali, Rwanda, May 1-5, 2023</i> .	Sandeep Polisetty, Heng Ji, Shih-Fu Chang, Kath-	800
745	Hyung Won Chung, Le Hou, Shayne Longpre, Barret	leen R. McKeown, Mohit Bansal, and Avi Sil. 2021.	801
746	Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang,	<a href="#">Infosurgeon: Cross-media fine-grained information</a>	802
747	Mostafa Dehghani, Siddhartha Brahma, Albert Web-	<a href="#">consistency checking for fake news detection</a> . In	803
748	son, Shixiang Shane Gu, Zhuyun Dai, Mirac Suz-	<i>Proceedings of the 59th Annual Meeting of the Asso-</i>	804
749	gun, Xinyun Chen, Aakanksha Chowdhery, Sharan	<i>ciation for Computational Linguistics and the 11th</i>	805
750	Narang, Gaurav Mishra, Adams Yu, Vincent Y. Zhao,	<i>International Joint Conference on Natural Language</i>	806
751	Yanping Huang, Andrew M. Dai, Hongkun Yu, Slav	<i>Processing, ACL/IJCNLP 2021, (Volume 1: Long</i>	807
752	Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam	<i>Papers)</i> , Virtual Event, August 1-6, 2021, pages 1683–	808
753	Roberts, Denny Zhou, Quoc V. Le, and Jason Wei.	1698.	809
754	2022. <a href="#">Scaling instruction-finetuned language models.</a>	Víctor Gallego. 2023. <a href="#">ZYN: zero-shot reward models</a>	810
755	<i>CoRR</i> , abs/2210.11416.	<a href="#">with yes-no questions</a> . <i>CoRR</i> , abs/2308.06385.	811
756	Nuri Cingillioglu and Alessandra Russo. 2021. <a href="#">pix2rule:</a>	Tianyu Gao, Adam Fisch, and Danqi Chen. 2021.	812
757	<a href="#">End-to-end neuro-symbolic rule learning</a> . In <i>Pro-</i>	<a href="#">Making pre-trained language models better few-shot</a>	813
758	<i>ceedings of the 15th International Workshop on</i>	<a href="#">learners</a> . In <i>Proceedings of the 59th Annual Meeting</i>	814
759	<i>Neural-Symbolic Learning and Reasoning as part</i>	<i>of the Association for Computational Linguistics and</i>	815
760	<i>of the 1st International Joint Conference on Learn-</i>	<i>and the 11th International Joint Conference on Natural</i>	816
761	<i>ing &amp; Reasoning (IJCLR 2021)</i> , Virtual conference,	<i>Language Processing, ACL/IJCNLP 2021, (Volume</i>	817
762	<i>October 25-27, 2021</i> , pages 15–56.	<i>1: Long Papers)</i> , Virtual Event, August 1-6, 2021,	818
763	Ganqu Cui, Shengding Hu, Ning Ding, Longtao Huang,	pages 3816–3830.	819
764	and Zhiyuan Liu. 2022. <a href="#">Prototypical verbalizer for</a>	Claire Glanois, Zhaohui Jiang, Xuening Feng, Paul	820
765	<a href="#">prompt-based few-shot tuning</a> . In <i>Proceedings of the</i>	Weng, Matthieu Zimmer, Dong Li, Wulong Liu, and	821
766	<i>60th Annual Meeting of the Association for Compu-</i>	Jianye Hao. 2022. <a href="#">Neuro-symbolic hierarchical rule</a>	822
767	<i>tational Linguistics (Volume 1: Long Papers)</i> , <i>ACL</i>	<a href="#">induction</a> . In <i>International Conference on Machine</i>	823
768	<i>2022, Dublin, Ireland, May 22-27, 2022</i> , pages 7014–	<i>Learning, ICML 2022, 17-23 July 2022, Baltimore,</i>	824
769	7024.	<i>Maryland, USA</i> , pages 7583–7615.	825
770	Limeng Cui, Kai Shu, Suhang Wang, Dongwon Lee,	Yue Huang and Lichao Sun. 2023. <a href="#">Harnessing the</a>	826
771	and Huan Liu. 2019. <a href="#">defend: A system for explain-</a>	<a href="#">power of chatgpt in fake news: An in-depth ex-</a>	827
772	<a href="#">able fake news detection</a> . In <i>Proceedings of the</i>	<a href="#">ploration in generation, detection and explanation.</a>	828
773	<i>28th ACM International Conference on Information</i>	<i>CoRR</i> , abs/2310.05046.	829
774	<i>and Knowledge Management, CIKM 2019, Beijing,</i>	Jiaming Ji, Tianyi Qiu, Boyuan Chen, Borong Zhang,	830
775	<i>China, November 3-7, 2019</i> , pages 2961–2964.	Hantao Lou, Kaile Wang, Yawen Duan, Zhong-	831
776	Kahneman Daniel. 2017. <i>Thinking, fast and slow</i> .	hao He, Jiayi Zhou, Zhaowei Zhang, Fanzhi Zeng,	832
777	Richard Evans and Edward Grefenstette. 2018. <a href="#">Learn-</a>	Kwan Yee Ng, Juntao Dai, Xuehai Pan, Aidan	833
778	<a href="#">ing explanatory rules from noisy data (extended ab-</a>	O’Gara, Yingshan Lei, Hua Xu, Brian Tse, Jie Fu,	834
779	<a href="#">stract)</a> . In <i>Proceedings of the Twenty-Seventh Inter-</i>	Stephen McAleer, Yaodong Yang, Yizhou Wang,	835
780	<i>national Joint Conference on Artificial Intelligence,</i>	Song-Chun Zhu, Yike Guo, and Wen Gao. 2023a.	836
781	<i>IJCAI 2018, July 13-19, 2018, Stockholm, Sweden,</i>	<a href="#">AI alignment: A comprehensive survey</a> . <i>CoRR</i> ,	837
782	pages 5598–5602.	abs/2310.19852.	838
783	Kevin Eykholt, Ivan Evtimov, Earlene Fernandes,	Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu,	839
784	Bo Li, Amir Rahmati, Chaowei Xiao, Atul Prakash,	Dan Su, Yan Xu, Etsuko Ishii, Yejin Bang, Andrea	840
785	Tadayoshi Kohno, and Dawn Song. 2018. <a href="#">Robust</a>	Madotto, and Pascale Fung. 2023b. <a href="#">Survey of halluci-</a>	841
786	<a href="#">physical-world attacks on deep learning visual clas-</a>	<a href="#">nation in natural language generation</a> . <i>ACM Comput.</i>	842
787	<a href="#">sification</a> . In <i>2018 IEEE Conference on Computer</i>	<i>Surv.</i> , 55(12):248:1–248:38.	843
788	<i>Vision and Pattern Recognition, CVPR 2018, Salt</i>	Anna Jobin, Marcello Ienca, and Effy Vayena. 2019.	844
789	<i>Lake City, UT, USA, June 18-22, 2018</i> , pages 1625–	<a href="#">The global landscape of AI ethics guidelines</a> . <i>Nat.</i>	845
790	1634.	<i>Mach. Intell.</i> , 1(9):389–399.	846
791	Angela Fan, Aleksandra Piktus, Fabio Petroni, Guil-	Elena Kochkina, Maria Liakata, and Arkaitz Zubiaga.	847
792	laume Wenzek, Marzieh Saeidi, Andreas Vlachos,	2018. <a href="#">All-in-one: Multi-task learning for rumour</a>	848
		<a href="#">verification</a> . In <i>Proceedings of the 27th International</i>	849

850	<i>Conference on Computational Linguistics, COLING</i>	906
851	<i>2018, Santa Fe, New Mexico, USA, August 20-26,</i>	907
852	<i>2018, pages 3402–3413.</i>	908
853	Bo Li, Peng Qi, Bo Liu, Shuai Di, Jingen Liu, Jiquan	909
854	Pei, Jinfeng Yi, and Bowen Zhou. 2023. <a href="#">Trustworthy</a>	910
855	<a href="#">AI: from principles to practices</a> . <i>ACM Comput. Surv.</i> ,	911
856	55(9):177:1–177:46.	912
857	Hao Liao, Jiahao Peng, Zhanyi Huang, Wei Zhang,	913
858	Guanghua Li, Kai Shu, and Xing Xie. 2023. <a href="#">MUSER:</a>	914
859	<a href="#">A multi-step evidence retrieval enhancement frame-</a>	
860	<a href="#">work for fake news detection</a> . In <i>Proceedings of the</i>	
861	<i>29th ACM SIGKDD Conference on Knowledge Dis-</i>	
862	<i>covery and Data Mining, KDD 2023, Long Beach,</i>	
863	<i>CA, USA, August 6-10, 2023, pages 4461–4472.</i>	
864	Hui Liu, Wenya Wang, and Haoliang Li. 2023. <a href="#">Inter-</a>	
865	<a href="#">pretable multimodal misinformation detection with</a>	
866	<a href="#">logic reasoning</a> . In <i>Findings of the Association</i>	
867	<i>for Computational Linguistics: ACL 2023, Toronto,</i>	
868	<i>Canada, July 9-14, 2023, pages 9781–9796.</i>	
869	Hui Liu, Wenya Wang, Hao Sun, Anderson Rocha, and	
870	Haoliang Li. 2024. <a href="#">Robust domain misinformation</a>	
871	<a href="#">detection via multi-modal feature alignment</a> . <i>IEEE</i>	
872	<i>Trans. Inf. Forensics Secur.</i> , 19:793–806.	
873	Jing Ma, Jun Li, Wei Gao, Yang Yang, and Kam-Fai	
874	Wong. 2023. <a href="#">Improving rumor detection by promot-</a>	
875	<a href="#">ing information campaigns with transformer-based</a>	
876	<a href="#">generative adversarial learning</a> . <i>IEEE Trans. Knowl.</i>	
877	<i>Data Eng.</i> , 35(3):2657–2670.	
878	Nikhil Mehta, Maria Leonor Pacheco, and Dan Gold-	
879	wasser. 2022. <a href="#">Tackling fake news detection by con-</a>	
880	<a href="#">tinually improving social context representations us-</a>	
881	<a href="#">ing graph neural networks</a> . In <i>Proceedings of the</i>	
882	<i>60th Annual Meeting of the Association for Compu-</i>	
883	<i>tational Linguistics (Volume 1: Long Papers), ACL</i>	
884	<i>2022, Dublin, Ireland, May 22-27, 2022, pages 1363–</i>	
885	<i>1380.</i>	
886	Ethan Mendes, Yang Chen, Wei Xu, and Alan Ritter.	
887	2023. <a href="#">Human-in-the-loop evaluation for early mis-</a>	
888	<a href="#">information detection: A case study of COVID-19</a>	
889	<a href="#">treatments</a> . In <i>Proceedings of the 61st Annual Meet-</i>	
890	<i>ing of the Association for Computational Linguis-</i>	
891	<i>tics (Volume 1: Long Papers), ACL 2023, Toronto,</i>	
892	<i>Canada, July 9-14, 2023, pages 15817–15835.</i>	
893	Parth Patwa, Shivam Sharma, Srinivas PYKL, Vineeth	
894	Guptha, Gitanjali Kumari, Md. Shad Akhtar, Asif	
895	Ekbali, Amitava Das, and Tanmoy Chakraborty. 2021.	
896	<a href="#">Fighting an infodemic: COVID-19 fake news dataset</a> .	
897	In <i>Combating Online Hostile Posts in Regional Lan-</i>	
898	<i>guages during Emergency Situation - First Interna-</i>	
899	<i>tional Workshop, CONSTRAINT 2021, Collocated</i>	
900	<i>with AAAI 2021, Virtual Event, February 8, 2021,</i>	
901	<i>Revised Selected Papers, pages 21–29.</i>	
902	Peng Qi, Juan Cao, Xirong Li, Huan Liu, Qiang Sheng,	
903	Xiaoyue Mi, Qin He, Yongbiao Lv, Chenyang Guo,	
904	and Yingchao Yu. 2021. <a href="#">Improving fake news detec-</a>	
905	<a href="#">tion by using an entity-enhanced framework to fuse</a>	
	<a href="#">diverse multimodal clues</a> . In <i>MM '21: ACM Multi-</i>	906
	<i>media Conference, Virtual Event, China, October 20</i>	907
	<i>- 24, 2021, pages 1212–1220.</i>	908
	Meng Qu, Junkun Chen, Louis-Pascal A. C. Xhonneux,	909
	Yoshua Bengio, and Jian Tang. 2021. <a href="#">Rnnlogic:</a>	910
	<a href="#">Learning logic rules for reasoning on knowledge</a>	911
	<a href="#">graphs</a> . In <i>9th International Conference on Learning</i>	912
	<i>Representations, ICLR 2021, Virtual Event, Austria,</i>	913
	<i>May 3-7, 2021.</i>	914
	Katie Sanders. 2023. <a href="#">PolitiFact</a> . <a href="https://www.politifact.com/">https://www.</a>	915
	<a href="https://www.politifact.com/">politifact.com/</a> [Accessed: (Accessed: Decem-	916
	ber 5, 2023)].	917
	Qiang Sheng, Juan Cao, H. Russell Bernard, Kai Shu,	918
	Jintao Li, and Huan Liu. 2022. <a href="#">Characterizing multi-</a>	919
	<a href="#">domain false news and underlying user effects on</a>	920
	<a href="#">chinese weibo</a> . <i>Inf. Process. Manag.</i> , 59(4):102959.	921
	Kai Shu. 2023. <a href="#">Combating disinformation on social</a>	922
	<a href="#">media and its challenges: A computational perspec-</a>	923
	<a href="#">tive</a> . In <i>Thirty-Seventh AAAI Conference on Artifi-</i>	924
	<i>cial Intelligence, AAAI 2023, Thirty-Fifth Conference</i>	925
	<i>on Innovative Applications of Artificial Intelligence,</i>	926
	<i>IAAI 2023, Thirteenth Symposium on Educational</i>	927
	<i>Advances in Artificial Intelligence, EAAI 2023, Wash-</i>	928
	<i>ington, DC, USA, February 7-14, 2023, page 15454.</i>	929
	Kai Shu, Deepak Mahudeswaran, Suhang Wang, Dong-	930
	won Lee, and Huan Liu. 2020. <a href="#">Fakenewsnet: A data</a>	931
	<a href="#">repository with news content, social context, and spa-</a>	932
	<a href="#">tiotemporal information for studying fake news on</a>	933
	<a href="#">social media</a> . <i>Big Data</i> , 8(3):171–188.	934
	Amila Silva, Ling Luo, Shanika Karunasekera, and	935
	Christopher Leckie. 2021. <a href="#">Embracing domain dif-</a>	936
	<a href="#">ferences in fake news: Cross-domain fake news de-</a>	937
	<a href="#">tection using multi-modal data</a> . In <i>Thirty-Fifth AAAI</i>	938
	<i>Conference on Artificial Intelligence, AAAI 2021,</i>	939
	<i>Thirty-Third Conference on Innovative Applications</i>	940
	<i>of Artificial Intelligence, IAAI 2021, The Eleventh</i>	941
	<i>Symposium on Educational Advances in Artificial In-</i>	942
	<i>teelligence, EAAI 2021, Virtual Event, February 2-9,</i>	943
	<i>2021, pages 557–565.</i>	944
	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	945
	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	946
	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	947
	Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-	948
	Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,	949
	Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller,	950
	Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-	951
	thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan	952
	Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,	953
	Isabel Kloumann, Artem Korenev, Punit Singh Koura,	954
	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-	955
	ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-	956
	tiniet, Todor Mihaylov, Pushkar Mishra, Igor Moly-	957
	bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-	958
	stein, Rashi Rungta, Kalyan Saladi, Alan Schelten,	959
	Ruan Silva, Eric Michael Smith, Ranjan Subrama-	960
	nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-	961
	lor, Adina Williams, Jian Xiang Kuan, Puxin Xu,	962
	Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,	963

- Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. [Llama 2: Open foundation and fine-tuned chat models](#). *CoRR*, abs/2307.09288.
- Stephanie Jean Tsang. 2023. HKBU Fact Check. <https://factcheck.hkbu.edu.hk/home/en/fact-check/our-process/> [Accessed: (Accessed: December 5, 2023)].
- Jindong Wang, Haoliang Li, Haohan Wang, Sinno Jialin Pan, and Xing Xie. 2023a. [Trustworthy machine learning: Robustness, generalization, and interpretability](#). In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD 2023, Long Beach, CA, USA, August 6-10, 2023*, pages 5827–5828.
- Kevin Ro Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. 2023b. [Interpretability in the wild: a circuit for indirect object identification in GPT-2 small](#). In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*.
- William Yang Wang. 2017. ["liar, liar pants on fire": A new benchmark dataset for fake news detection](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 2: Short Papers*, pages 422–426.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. [Chain-of-thought prompting elicits reasoning in large language models](#). In *NeurIPS*.
- Weizhi Xu, Junfei Wu, Qiang Liu, Shu Wu, and Liang Wang. 2022. [Evidence-aware fake news detection with graph neural networks](#). In *WWW '22: The ACM Web Conference 2022, Virtual Event, Lyon, France, April 25 - 29, 2022*, pages 2501–2510.
- Zhenrui Yue, Huimin Zeng, Yang Zhang, Lanyu Shang, and Dong Wang. 2023. [Metaadapt: Domain adaptive few-shot misinformation detection via meta learning](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 5223–5239.
- Xinyi Zhou and Reza Zafarani. 2021. [A survey of fake news: Fundamental theories, detection methods, and opportunities](#). *ACM Comput. Surv.*, 53(5):109:1–109:40.
- Yongchun Zhu, Qiang Sheng, Juan Cao, Qiong Nan, Kai Shu, Minghui Wu, Jindong Wang, and Fuzhen Zhuang. 2023. [Memory-guided multi-view multi-domain fake news detection](#). *IEEE Trans. Knowl. Data Eng.*, 35(7):7178–7191.

## A Details of Cognition System

Unlike convolutional deep learning-based fake news detection frameworks that classify in a latent space, the cognition system of TELLER, aims to emulate human fact-checking experts by complying with specific policies to ensure transparency and controllability of the detection process. In this section, we describe the construction of the set of question templates  $\mathcal{Q}$  and  $\mathcal{Q}'$  for TELLER and  $w/Intervention$  respectively in Appendix A.1. Furthermore, we introduce a trick for batch training by fixing the number of logic atoms for different inputs in Appendix A.2 and outline some potential techniques for further improvement of the cognition system in Appendix A.3.

### A.1 Construction of Question Templates

To provide an overview, we present the referenced human-checking process in Table 5. In this table, Steps I, VI and VII are excluded from detection algorithms, as they either fall into the preliminary procedures or the post-processing stages of the fake news detection pipeline. These steps may involve data crawl, human-computer interaction, machine translation, etc. As a result, we concentrate on the other steps.

Subsequently, we decompose the process into a Yes/No question template set  $\mathcal{Q}$ , where each template  $Q_i$  in  $\mathcal{Q}$  corresponds to a predicate  $P_i$  in the predicate set  $\mathcal{P}$ . All question templates and their corresponding predicates are listed in Table 6. Specifically, for  $Q_1$ , our objective is to determine the trustworthiness of statements in the input news. Here, statements represent crucial information in news articles, playing a vital role in debunking misinformation. Additionally, extracting statements from news is a challenging task. While previous studies like Liao et al. (2023); Fung et al. (2021) used pre-trained language models to generate summaries as statements, we choose to utilize GPT3.5-turbo to generate statements for simplicity in implementation. The prompt used for this purpose is as follows:

To verify the MESSAGE, what are the critical

claims related to this message we need to verify? Please use the following format to answer. If there are no important claims, answer “not applicable”.

MESSAGE:

CLAIM:

CLAIM:

MESSAGE: \$MESSAGE\$.

Then, we replace the "\$MESSAGE\$" with input news and take the generated claims as statements for  $Q_1$  ( $P_1$ ).

Additionally, when verifying the controllability of our framework, we propose adjusting the question template set to deal with the diversity of fake news. While this adjustment should be done by fact-checking experts to ensure the reasonableness of new questions, our empirical findings demonstrate the feasibility of guiding large language models, such as GPT-3.5-turbo, to generate new question templates. These templates are then manually filtered by us to create the final question template set  $\mathcal{Q}'$ , and the corresponding predicate set  $\mathcal{P}'$  for intervention, as outlined in Algorithm 1. Table 7 presents these newly added question templates and predicates. The prompt  $R$  used in this algorithm is as follows:

Write some questions that can be used to determine whether a news report is misinformation. The questions should be answerable by large language models in a close-book situation without requiring additional information. Please format each question using the <s> and </s> tags, such as <s>A question</s>.

### A.2 Trick for Batch Training

To enable batch training, we fix the number of logic atoms, denoted as  $M_i$  for each predicate  $P_i$ . Specifically, If  $M_i < \prod_{j=1}^{N_i} |X_{i,j}|$ , we randomly select  $M_i$

atoms. Conversely, if  $M_i > \prod_{j=1}^{N_i} |X_{i,j}|$ , we pad the vector by 0 accordingly. In the end,  $\mu$  can be represented as  $[\mu_{1,1}, \dots, \mu_{1,M_1}, \dots, \mu_{N,1}, \dots, \mu_{N,M_N}]$ ,

where  $\mu \in \mathbb{R}^M$  and  $M = \sum_i^N M_i$ .

### A.3 The Potential of Cognition System

It is noteworthy that specific techniques can be employed to improve the performance of our cognitive system. For instance, when obtaining the answers to questions as truth values for corresponding logic atoms in Sec. 3.1.2, we exclusively consider "Yes" and "No" tokens. However, considering the relationship between model outputs and final predictions, "Right" and "Wrong" tokens can also be suitable candidates. Therefore, drawing motivation from (Gao et al., 2021; Cui et al., 2022), existing manual or automatic verbalizer techniques that establish mappings between diverse model outputs and final labels can be leveraged to enhance performance. Additionally, the ensemble of prompts, similar to "Yes or No? The answer is: ", has proven effective for the "Yes" and "No" classification task in (Gallego, 2023). Consequently, our dual-system framework exhibits substantial potential for future improvements in the cognitive system.

---

**Algorithm 1** Question Template Generation for Intervention Algorithm

---

**Input:** Prompt  $R$ , the original question template set  $\mathcal{Q}$ , and a copy of  $\mathcal{Q}$  denoted as  $\hat{\mathcal{Q}}$

**Output:** The question template set  $\mathcal{Q}'$  for intervention

- 1: Set the number of iteration steps as  $T$
  - 2: **for** Iteration  $t = 1, \dots, T$  **do**
  - 3:   Use  $R$  to guide GPT-3.5-turbo in generating a set of new question templates  $\mathcal{Q}'$
  - 4:   **for** each question template  $Q'_i$  in  $\mathcal{Q}'$  **do**
  - 5:     Compute the average similarity score between  $Q'_i$  and all templates in  $\hat{\mathcal{Q}}$  using Sentence BERT.
  - 6:   **end for**
  - 7:   Add  $Q'_i \in \mathcal{Q}'$  with the lowest similarity score to  $\hat{\mathcal{Q}}$ .
  - 8: **end for**
  - 9:  $\mathcal{Q}' = \hat{\mathcal{Q}} \setminus \mathcal{Q}$
  - 10: Manually refine  $\mathcal{Q}'$  by removing duplicate and impractical templates that are non-verifiable through LLMs, resulting in the final  $\mathcal{Q}'$ .
-

<p><b>Step I: Selecting claims</b></p> <p>(1) To filter the information on news websites, social media, and online databases through manual selection and computer-assisted selection.</p> <p>(2) The public can submit suspicious claims.</p> <p>(3) Selecting suspicious claims based on their hotness in Hong Kong, considering factors such as the amount of likes, comments, and shares the message has received.</p> <p>A) Is the content checkable?</p> <p>B) Any misleading or false content?</p> <p>C) Does it meet public interest?</p> <p>D) Is it widespread?</p>
<p><b>Step II: Tracing the source</b></p> <p>(1) Determining the source of the information.</p> <p>(2) Identifying the publication date.</p> <p>(3) Investigating the publisher and their background and reputation.</p> <p>(4) Checking for similar information.</p> <p>(5) Capturing a screen record and attaching the URL link.</p> <p>(6) Providing two or more additional sources of information.</p>
<p><b>Step III: Fact-checking the suspicious information</b></p> <p>(1) Applying the Five Ws and an H: When, Where, Who, What, Why, How.</p> <p>(2) Searching for evidence to verify the information, such as official press releases, authoritative media reports, and research reports.</p> <p>(3) Attempting to engage the person or organization making the claim through email or telephone, if necessary.</p> <p>(4) Consulting experts in the relevant field, if necessary.</p>
<p><b>Step IV: Retrieving contextual information</b></p> <p>(1) Checking if the original claim contains adequate background information.</p> <p>(2) Assessing the accuracy and objectivity of the background information.</p> <p>(3) Identifying any intentionally eliminated content that distorts the meaning.</p>
<p><b>Step V: Evaluating improper intentions</b></p> <p>(1) Assessing if there is any improper intention (e.g., political motive, commercial purpose) in the information.</p> <p>(2) Investigating if the publisher has a history of publishing information with improper intentions.</p>
<p><b>Step VI: Self-checking</b></p> <p>(1) Fact-checkers signing a Declaration of Interest Form before joining the team.</p> <p>(2) Ensuring fact-checkers maintain objectivity and avoid biases during the process.</p> <p>(3) Upholding the principle of objectivity and avoiding emotional involvement.</p>
<p><b>Step VII: Publishing and reviewing reports</b></p> <p>(1) Completing a draft of the fact-check report, followed by editing and reviewing by professional editors and consultants.</p> <p>(2) Updating the report if any mistakes or defects are found, and providing clarification on correction reasons and date.</p>

Table 5: Fake news detection policy of HKBU FACT CHECK Team (Tsang, 2023)

Question Template	Logic Predicate: Logic Semantics	Annotation
Q <sub>1</sub> : Background Information: $X_{1,1}$ . Statement: $X_{1,2}$ . Is the statement true?	$P_1(X_{1,1}, X_{1,2})$ : Given the background information $X_{1,1}$ , the statement is true.	$X_{1,1}$ : Background information for input news, $X_{1,2}$ : Check-worthy statements in input news.
Q <sub>2</sub> : Background Information: $X_{2,1}$ . Message: $X_{2,2}$ . Is the message true?	$P_2(X_{2,1}, X_{2,2})$ : Given the background information $X_{2,1}$ , the message is true.	$X_{2,1}$ : Background information for input news, $X_{2,2}$ : Input news.
Q <sub>3</sub> : Message: $X_{3,1}$ . Did the message contain adequate background information?	$P_3(X_{3,1})$ : The message contains adequate background information.	$X_{3,1}$ : Input news.
Q <sub>4</sub> : Message: $X_{4,1}$ . Is the background information in the message accurate and objective?	$P_4(X_{4,1})$ : The background information in the message is accurate and objective.	$X_{4,1}$ : Input news.
Q <sub>5</sub> : Message: $X_{5,1}$ . Is there any content in the message that has been intentionally eliminated with the meaning being distorted?	$P_5(X_{5,1})$ : The content in the message has been intentionally eliminated with the meaning being distorted.	$X_{5,1}$ : Input news.
Q <sub>6</sub> : Message: $X_{6,1}$ . Is there an improper intention (political motive, commercial purpose, etc.) in the message?	$P_6(X_{6,1})$ : The message has an improper intention.	$X_{6,1}$ : Input news.
Q <sub>7</sub> : Publisher Reputation: $X_{7,1}$ . Does the publisher have a history of publishing information with an improper intention?	$P_7(X_{7,1})$ : Given the publisher reputation $X_{7,1}$ , the publisher has a history of publishing information with an improper intention.	$X_{7,1}$ : Publishing history.
Q <sub>8</sub> : Background Information: $X_{8,1}$ . Message: $X_{8,2}$ . Is the message false?	$P_8(X_{8,1}, X_{8,2})$ : Given the background information $X_{8,1}$ , the message is false.	$X_{8,1}$ : Background information for input news, $X_{8,2}$ : Input news.

Table 6: Question template set  $\mathcal{Q}$  and logic predicate set  $\mathcal{P}$

Question Template	Logic Predicate: Logic Semantics	Annotation
Q <sub>9</sub> : News Report: $X_{9,1}$ . Is the news report based on facts or does it primarily rely on speculation or opinion?	$P_9(X_{9,1})$ : The news report is based on facts and relies on speculation or opinion.	$X_{9,1}$ : Input news.
Q <sub>10</sub> : News Report $X_{10,1}$ : Are there any logical fallacies or misleading arguments present in the news report?	$P_{10}(X_{10,1})$ : The news report has logical fallacies or misleading arguments.	$X_{10,1}$ : Input news.
Q <sub>11</sub> : Message: $X_{11,1}$ . Does the message exhibit bias?	$P_{11}(X_{11,1})$ : The message exhibits bias.	$X_{11,1}$ : Input news.
Q <sub>12</sub> : News report: $X_{12,1}$ . Are there any grammatical or spelling errors in the news report that may indicate a lack of professional editing??	$P_{12}(X_{12,1})$ : The news report has grammatical and spelling errors.	$X_{12,1}$ : Input news.
Q <sub>13</sub> : News report: $X_{13,1}$ . Does the news report use inflammatory language or make personal attacks?	$P_{13}(X_{13,1})$ : The news report uses inflammatory language and makes personal attacks.	$X_{13,1}$ : Input news.

Table 7: Question template set  $\mathcal{Q}'$  and logic predicate set  $\mathcal{P}'$  generated by GPT-3.5-turbo for intervention

## B Details of Experimental Setting

### B.1 Datasets

**LIAR** is a publicly available dataset for fake news detection, sourced from POLITIFACT.COM. This dataset comprises six fine-grained labels for truthfulness ratings: true, mostlytrue, halftrue, barelytrue, false, and pantsfire. To align with the binary classification problem, we merge true, mostlytrue into true and merge barelytrue, false, and pantsfire into false, following (Liao et al., 2023). Moreover, Wang (2017); Alhindi et al. (2018) curated relevant evidences from fact-checking experts (e.g., publisher information, background information, etc.), which serve as gold knowledge in an open setting.

**Constraint** is a manually annotated dataset of real and fake news related to COVID-19. We adopt the data pre-processing procedures described in (Patwa et al., 2021), which involve removing all links, non-alphanumeric characters, and English stop words.

**Politifact** and **GossipCop** are two binary classification subsets extracted from FakenewsNet (Shu et al., 2020). The Politifact subset comprises political news, while the GossipCop subset comprises entertainment stories. To optimize experimental costs and adhere to maximum context limitations, we exclude news samples longer than 3,000 words.

For dataset partitioning, we follow the default partition if specified; otherwise, we use a 7:1:2 ratio. Table 8 presents the statistics of each dataset.

Split	LIAR	Constraint	Politifact	GossipCop
Train	10202	6299	469	6999
Validation	1284	2139	66	999
Test	1271	2119	136	2002

Table 8: Statistics of four benchmarks

### B.2 Illustration of Different Baselines

We compare our model against *Direct*, *Few-shot Direct*, *Zero-shot COT*, *Few-shot COT*, *Few-shot Logic*. *Direct* utilizes LLMs to calculate the probability of each label using Eqs. 1-2 and then selects the label with the highest likelihood as the predicted label. Building upon *Direct*, *Few-shot Direct* incorporates demonstration samples with known labels as contextual information to enhance the model’s performance. *Zero-shot COT* and *Few-shot COT* employ the chain-of-thought (COT) technique

(Wei et al., 2022), enabling LLMs to engage in step-by-step reasoning. While *Zero-shot COT* immediately adds the prompt "Let us think step by step!", *Few-shot COT* provides multiple COT exemplars. For *Few-shot Logic*, we replace the thought process in COT with instantiated questions accompanied by corresponding answers generated by our cognition system. Since COT prompts have been found to yield performance gains basically when used with models of approximately 100B parameters (Wei et al., 2022), we exclusively implement COT-related methods using GPT-3.5-turbo.

Below we show the templates for these five baselines for the fake news detection task in the closed setting without the access to any external knowledge source.

**Direct:**

Message: \$MESSAGE\$.  
Is the message \$Label\$?  
Yes or No? Response:

Then, we replace the "\$MESSAGE\$" with input news, "\$Label\$" with candidate truthfulness labels.

**Few-shot Direct:**

Following given examples to answer Yes/No questions.

Message: Says the Annies List political group supports third-trimester abortions on demand.  
Is the message true?  
Yes or No? Response: No

Message: Says the Annies List political group supports third-trimester abortions on demand.  
Is the message false?  
Yes or No? Response: Yes

(... more examples here ...)

Message: \$MESSAGE\$.  
Is the message \$Label\$?  
Yes or No? Response:

Then, we replace the "\$MESSAGE\$" with input news, "\$Label\$" with candidate truthfulness labels. Furthermore, during the testing phase, the examples are randomly selected from the training set.

1173

**Zero-shot COT:**

You will be provided with a statement, and your task is to classify its truthfulness into one of two categories: true and false.  
Message: \$MESSAGE\$.  
Let's think step by step and give answer with the suffix "So the final answer is".

1174

Then, we replace the "\$MESSAGE\$" with the input news.

1175

**Few-shot COT:**

1176

You will be provided with a statement, and your task is to classify its truthfulness into one of two categories: true and false.

Example One  
Message: Says the Annies List political group supports third-trimester abortions on demand. Let's think step by step and give answer with suffix "So the final answer is".  
Annie's List was comfortable with candidates who oppose more limits on late-term abortions while he also supported candidates who voted for more limits this year. Both dose not mention of third-trimester abortions.  
So the final answer is false.

(... more examples here ...)

Message: \$MESSAGE\$.  
Let's think step by step and give answer with the suffix "So the final answer is".

1177

Then, we replace the "\$MESSAGE\$" with the input news.

1178

**Few-shot Logic:**

1179

You will be provided with a statement, and

your task is to classify its truthfulness into one of two categories: true and false.

Example One

Message: Says the Annies List political group supports third-trimester abortions on demand.

Decomposed Questions:

(1) Statement: The Annies List is a political group. Is the statement true?

Yes

(2) Statement: The Annies List supports third-trimester abortions. Is the statement true?

No

(3) Did the message contain adequate background information?

False

(... more examples here ...)

Message: \$MESSAGE\$.

Let's think step by step and give answer with the suffix "So the final answer is".

Then, we replace the "\$MESSAGE\$" with the input news.

1180

1181

**B.3 Model Training for Decision System**

1182

In the decision system of our framework, we employ the DNF Layer to learn human-readable rules from data differentially. To train this model, we utilize the Adam optimizer with a learning rate of 1e-3. Regarding the hyperparameters, we search the conjunction number  $C$  within the range [10, 20, 30, 40, 50], and the weight decay within the range [1e-3, 5e-4, 1e-4]. Furthermore, to showcase the superiority of our approach, we maintain consistent hyperparameters across different LLMs in each setting. For instance, all hyperparameters of TELLER in the closed setting for the binary classification task on the LIAR dataset remain unchanged. The batch size is set to 64, and the number of epochs is set to 30. Additionally, we progressively converge the model towards symbolic semantics by adjusting  $\delta$  (refer to Appendix E for detail) to 1 or -1 before the first 15 epochs using exponential decay.

1183

1184

1185

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

## C Details of Explainability Study

To enhance the accessibility of the rules generated by the DNF Layer, we propose a pruning algorithm that extracts more concise logic clauses by eliminating insignificant weights. The algorithm is described in Algorithm 2. Furthermore, to demonstrate the explainability of our framework, we visualize the extracted rules obtained from the pruned model for Constraint, Politifact, and GossipCop datasets in Table 9, Table 10 and Table 4, respectively. In these tables,  $P_{\text{true}}$  and  $P_{\text{false}}$  represent the proposition that the input news is identified as true or false, respectively. In our visualization experiments, we employ Llama2 (13B) as the LLM in the cognition system. We set the number of conjunctive layers  $C$  as 50, the performance drop threshold  $\epsilon$  as 0.005, and  $b$  as 0.0001 to reduce the number of conjunction clauses. More details regarding these parameters can be found in Appendix E.

### Algorithm 2 Pruning Algorithm for the DNF Layer

**Input:** Trained DNF Layer  $\Phi$ , performance drop threshold  $\epsilon$   
**Output:** Pruned DNF Layer  $\Phi'$  and extracted rule set  $\mathcal{R}$

- 1: Initialize  $\mathcal{R}'$  as an empty set
- 2: Initialize  $\mathcal{R}$  by extracting rules from  $\Phi$
- 3: Initialize  $\Phi'$  using  $\Phi$
- 4: **while**  $|\mathcal{R}'| \neq |\mathcal{R}|$  **do**
- 5:   Initialize  $\mathcal{R}$  by extracting rules from  $\Phi'$
- 6:   Prune disjunctions if the removal of a disjunction results in a performance drop smaller than  $\epsilon$
- 7:   Prune unused conjunctions that are not utilized by any disjunction
- 8:   Prune conjunctions if the removal of a conjunction results in a performance drop smaller than  $\epsilon$
- 9:   Prune disjunctions that use empty conjunctions
- 10:   Prune disjunctions again if the removal of a disjunction results in a performance drop smaller than  $\epsilon$
- 11:   Update the pruned model as  $\Phi'$  and extract rules from  $\Phi'$  to obtain  $\mathcal{R}'$ ;
- 12: **end while**

$$\begin{aligned} \text{conj}_{48} &= P_4 \wedge \neg P_8 \\ \text{conj}_{25} &= \neg P_4 \wedge \neg P_5 \wedge P_8 \\ \text{conj}_{40} &= P_2 \wedge P_4 \\ P_{\text{true}} &= \text{conj}_{48} \\ P_{\text{false}} &= \text{conj}_{25} \vee \neg \text{conj}_{40} \end{aligned}$$

Table 9: Extracted rules for the Constraint dataset when using Llama2 (13B).

$$\begin{aligned} \text{conj}_{36} &= P_3 \wedge P_6 \wedge P_8 \\ \text{conj}_{44} &= P_5 \wedge P_1 \wedge P_8 \\ \text{conj}_0 &= P_1 \\ \text{conj}_{49} &= P_2 \wedge P_3 \wedge P_4 \\ P_{\text{true}} &= \neg \text{conj}_{36} \vee \neg \text{conj}_{44} \\ P_{\text{false}} &= \neg \text{conj}_0 \vee \neg \text{conj}_{49} \end{aligned}$$

Table 10: Extracted rules for the Politifact dataset when using Llama2 (13B).

## D Comparison with Different Decision Models

In our work, we utilize the DNF Layer to construct our decision system, guaranteeing explainability and controllability. However, there are also other alternatives, such as existing neural symbolic architectures and interpretable machine learning algorithms. By comparing the DNF Layer with these candidates, we demonstrate that our dual-system framework can achieve better performance by inventing a more effective decision model to unleash the ability of LLMs.

While existing neural symbolic architectures can extract useful rules from data (Booch et al., 2021), they indeed have certain limitations. Firstly, these architectures often require complex mechanisms to implement logical operations, which makes them unsuitable for immediate application in fake news detection tasks. For example, Qu et al. (2021); Cheng et al. (2023) developed neural-symbolic models for knowledge graph completion, but their reliance on well-defined graph structures makes them infeasible for our task. Secondly, these architectures often suffer from efficiency issues. For instance,  $\delta$ LP proposed in (Evans and Grefenstette, 2018) had high computational complexity, and HRI (Glanois et al., 2022) was incompatible with batch training, which externally required users to pre-define rule templates to constrain the search space. Furthermore, to the best of our knowledge, there may be no neural-symbolic framework available that can simultaneously handle the challenges of missing values and multi-grounding problems (i.e., one predicate can be instantiated as multiple logic atoms), which are common in our tasks. Therefore, we acknowledge the need for future research to develop a more suitable and powerful neural-symbolic framework in the context of fake news detection.

Since each dimension in  $\mu$  is precisely bonded to a question template (logic predicate), we can employ traditional machine learning classification algorithms, including decision tree<sup>3</sup> and naive Bayes Classifier<sup>4</sup> to replace the DNF Layer to drive our decision system, while maintaining partial aspects of trustworthy AI. Therefore, we compare the DNF Layer with these two methods in both in-domain

and cross-domain settings on three datasets, shown in Table 11 and Table 12, respectively.

According to the results, we conclude that decision trees perform better when the training and testing data are from the same domain. Meanwhile, the naive Bayes Classifier demonstrates more satisfactory generalization performance in cross-domain experiments across various LLMs. This implies that our proposed dual-system framework shows potential in developing a more powerful decision module, such as an ensemble of these algorithms. However, the DNF Layer still outperforms these two methods in most cases when using Llama2 (13B) as the driver of the cognition system, achieving a better trade-off between accuracy and generalization ability. Moreover, the DNF Layer also exhibits advantages over these two methods in terms of its ability to handle missing values and multi-grounding problems, as well as its flexibility in efficiently searching logic rules in a large space, whereas the decision tree is constrained by depth and width.

<sup>3</sup><https://scikit-learn.org/stable/modules/tree.html>

<sup>4</sup>[https://scikit-learn.org/stable/modules/naive\\_bayes.html](https://scikit-learn.org/stable/modules/naive_bayes.html)

LLMs	Method	Constraint		Politifact		GossipCop	
		Acc(%)	Macro-F1(%)	Acc(%)	Macro-F1(%)	Acc(%)	Macro-F1(%)
FLAN-T5-large	Decision Tree	78.53	78.30	67.65	67.19	70.88	62.76
	Bayes Classifier	80.93	80.86	66.18	66.15	68.33	61.04
	TELLER	80.32	80.11	67.65	67.65	69.53	59.39
FLAN-T5-xl	Decision Tree	84.29	84.27	66.91	66.10	71.13	61.58
	Bayes Classifier	82.40	82.22	68.38	67.88	68.23	60.23
	TELLER	83.77	83.66	68.82	64.68	69.58	58.72
FLAN-T5-xxl	Decision Tree	84.14	84.12	72.06	71.00	72.13	67.08
	Bayes Classifier	82.49	82.30	68.38	67.61	68.38	57.62
	TELLER	83.39	83.24	69.12	68.57	69.18	57.21
Llama2 (7B)	Decision Tree	84.33	84.32	79.41	77.00	72.38	65.24
	Bayes Classifier	83.11	82.97	76.47	76.29	71.98	66.67
	TELLER	83.72	83.54	<b>83.82</b>	<b>83.81</b>	70.68	59.58
Llama2 (13B)	Decision Tree	86.50	86.49	83.09	83.07	74.43	68.99
	Bayes Classifier	84.99	84.92	80.15	80.06	73.58	<b>69.59</b>
	TELLER	<b>87.31</b>	<b>87.29</b>	79.41	79.41	<b>74.48</b>	66.32

Table 11: Results of different decision models on Constraint, Politifact, and GossipCop datasets without access to retrieved background information. The best results for each dataset are highlighted with bold numbers.

LLMs	Method	<b>CP→G</b>		<b>GP→C</b>		<b>CG→P</b>	
		Acc(%)	Macro-F1(%)	Acc(%)	Macro-F1(%)	Acc(%)	Macro-F1(%)
FLAN-T5-xl	Decision Tree	68.98	62.33	73.67	73.32	63.97	62.71
	Bayes Classifier	67.13	59.26	82.49	82.49	64.71	64.64
	TELLER	68.13	56.54	82.40	82.09	61.76	60.92
FLAN-T5-xxl	Decision Tree	68.33	55.53	70.60	70.35	61.03	60.98
	Bayes Classifier	68.33	54.71	82.63	82.51	62.50	62.50
	TELLER	69.13	53.15	77.44	76.21	66.18	66.17
Llama2 7B	Decision Tree	52.20	52.05	76.40	75.02	66.91	64.84
	Bayes Classifier	65.98	62.46	82.82	82.60	67.65	65.49
	TELLER	68.33	59.33	81.60	81.04	83.09	82.82
Llama2 13B	Decision Tree	61.59	61.14	71.54	68.21	71.32	71.32
	Bayes Classifier	<b>71.53</b>	<b>69.09</b>	82.59	82.25	78.68	78.25
	TELLER	70.93	60.90	<b>85.09</b>	<b>84.87</b>	<b>79.41</b>	<b>79.41</b>

Table 12: Results of different decision models on cross-domain experiments. **C**, **P** and **G** represent Constraint, Politifact, and GossipCop datasets, respectively. The best results for each dataset are highlighted with bold numbers.

## E Formal Description of DNF Layer

In this section, we introduce modified Disjunctive Normal Form (DNF) Layer employed in our framework. The DNF Layer is built from semi-symbolic layers (SL), which can progressively converge to symbolic semantics such as conjunction  $\wedge$  and disjunction  $\vee$ .

Specifically, for the truth value vector  $\mu \in \mathbb{R}^M$  mentioned in Sec. 3.1.2, SL can be formulated as follows:

$$\mu_o = \tanh \left( \sum_j^M w_j \mu_j + \beta \right), \quad (4)$$

$$\beta = \delta \left( b - \sum_j |w_j \mu_j| \right), \quad (5)$$

where  $w_j$  represents learnable parameters,  $b = \max_j |w_j \mu_j|$  and  $\delta \in [-1, 1]$  represents the semantic gate selector.  $\mu_j$  is the truth value for the  $j$ th logic atom obtained from the cognitive system. The sign of the learned weight  $w_j$  indicates whether  $\mu_j$  (if  $w_j$  is positive) or its negation (if  $w_j$  is negative) contributes to  $\mu_o$ . Thus, logical negation (e.g.,  $\neg p_j$ ) can be computed as the multiplicative inverse of the input:  $-\mu_j$ .

Eq. 4 resembles a standard feed-forward layer, aiming to compute a single truth value from a collection of values  $\mu_j$  corresponding to different instantiations of a single predicate/question.  $\beta$  serves as the bias term. As shown by (Cingillioglu and Russo, 2021), by adjusting  $\delta$  from 0 to 1 during training, SL tends to converge to conjunctive semantics as  $SL_{\wedge}$  (e.g.,  $p_1 \wedge p_2, \dots, \wedge p_M$ ), indicating that if at least one  $w_j \mu_j$  is false, the output  $\mu_o$  will be false; otherwise,  $\mu_o$  will be true. Conversely, by gradually adjusting  $\delta$  from 0 to  $-1$ , SL can attain disjunctive semantics as  $SL_{\vee}$  (e.g.,  $p_1 \vee p_2, \dots, \vee p_M$ ), where if at least one  $w_j \mu_j$  is true,  $\mu_o$  will be true; otherwise,  $\mu_o$  will be false. Additionally,  $b$  can guarantee  $\mu_o$  being true (false) when all  $w_j \mu_j$  are true (false) for  $SL_{\wedge}$  ( $SL_{\vee}$ ).

Since each dimension in  $\mu$  corresponds to the same predicate for different inputs, SL effectively represents the relationship among different instantiations and the target output  $\mu_o$ , enabling the learning of generic rules for various inputs. Moreover, by employing rule-based aggregation, our framework exhibits noise tolerance against incorrect predictions of LLMs in the cognition system, particularly owing to the  $SL_{\vee}$ .

Notably, one predicate can be instantiated by multiple assignments, i.e.,  $P_i$  pertains to  $M_i$  logic atoms in Appendix A.2. Thus, the parameters bound to these  $M_i$  logic atoms should naturally share the logical semantics of  $P_i$ . Instead of gathering all possible combinations of  $M_i$  logic atoms

for training ( $\prod_{j=1}^{M_i} j$ ), we let these logic atoms share the same  $w$ . In this scenario, SL can be represented as follows:

$$\mu_o = \tanh \left( \sum_i^N \sum_j^{M_i} w_i \mu_{i,j} + \beta \right), \quad (6)$$

$$\beta = \delta \left( b - \sum_i^N \sum_j^{M_i} |w_i \mu_{i,j}| \right), \quad (7)$$

where  $N$  is the number of predicates.