

Boosting Logical Fallacy Reasoning in LLMs via Logical Structure Tree

Yuanyuan Lei and Ruihong Huang Department of Computer Science and Engineering Texas A M University, College Station, TX yuanyuan, huangrh@tamu.edu

February 21, 2026

Abstract

educational essay (Jin et al., 2022), political de-

Logical fallacy uses invalid or faulty reasoning in the construction of a statement. Despite the prevalence and harmfulness of logical fallacies, detecting and classifying logical fallacies still remains a challenging task. We observe that logical fallacies often use connective words to indicate an intended logical relation between two arguments, while the argument semantics does not actually support the logical relation. Inspired by this observation, we propose to build a logical structure tree to explicitly represent and track the hierarchical logic flow among relation connectives and their arguments in a statement. Specifically, this logical structure tree is constructed in an unsupervised manner guided by the constituency tree and a taxonomy of connectives for ten common logical relations, with relation connectives as non-terminal nodes and textual arguments as terminal nodes, and the latter are mostly elementary discourse units.

We further develop two strategies to incorporate the logical structure tree into LLMs for fallacy reasoning. Firstly, we transform the tree into natural language descriptions and feed the textualized tree into LLMs as a part of the hard text prompt. Secondly, we derive a relation-aware tree embedding and insert the tree embedding into LLMs as a soft prompt. Experiments on benchmark datasets demonstrate that our approach based on logical structure tree significantly improves precision and recall for both fallacy detection and fallacy classification 1.

1 Introduction

Logical fallacy refers to the use of invalid or flawed reasoning in an argumentation (Risen et al., 2007; Walton, 2010; Cotton, 2018). Logical fallacy can occur as unintentional mistakes or deliberate persuasions in a variety of human communications, such as news media (Da San Martino et al., 2019),

¹The code and data link is: <https://github.com/>
yuanyuan lei - nlp/logical-fallacy-emnlp-2024

Logical fallacy uses invalid or faulty reasoning in the construction of a statement. Despite the prevalence and harmfulness of logical fallacies, detecting and classifying logical fallacies still remains a challenging task. We observe that logical fallacies often

use connective words to indicate an intended logical relation between two arguments, while the argument semantics does not actually support the logical relation. Inspired by this observation, we propose to build a logical structure tree to explicitly represent and track the hierarchical logic flow among relation connectives and their arguments in a statement.

Specifically, this logical structure tree is constructed in an unsupervised manner guided by the constituency tree and a taxonomy of connectives for ten common logical relations, with relation connectives as non-terminal nodes and textual arguments as terminal nodes, and the latter are mostly elementary discourse units. We further develop two strategies to incorporate the logical structure tree into LLMs for fallacy

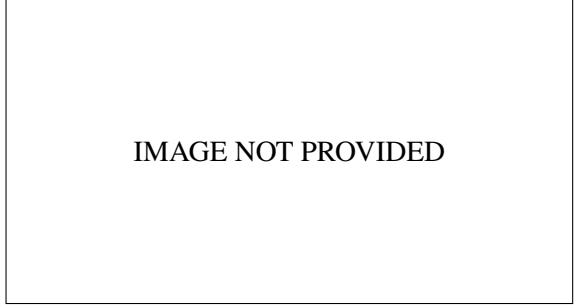


IMAGE NOT PROVIDED

Figure 1: Figure 1: Examples of logical fallacy sentences and their logical structure trees. The logical structure tree features logical relation connectives as non-terminal nodes, and textual arguments as terminal nodes.

reasoning. Firstly, we transform the tree into natural language descriptions and feed the textualized tree into LLMs as a part of the hard text prompt. Secondly, we derive a relation-aware tree embedding and insert the tree embedding into LLMs as a soft prompt. Experiments on benchmark datasets demonstrate that our approach based on logical structure tree significantly improves precision and recall for both fallacy detection and fallacy classification 1.

The key observation is that logical fallacies heavily rely on connective phrases to indicate an intended logical relation between two textual arguments, while the semantics of the arguments do not actually support the logical relation. Figure 1 shows two examples where the connective phrases were bolded. The first example uses the connective words **therefore** and **cause** to suggest a causal relation between vaccinations and increasing flu cases, however, the temporal relation between the two events as stated in the first half of the statement does not necessarily entail a causal relation between them, and indeed, their semantics do not actually support the suggested causal relation. Recognizing this discrepancy undermines the credibility of the whole statement. Similarly in the second example, the connective word **likewise** is commonly used to indicate an analogy relation, however, the second argument is clearly a specific case of the general condition stated in the first argument and therefore there is no analogy relation between them, and recognizing this mismatch between the suggested logical relation and the real relation enables us to detect this fallacy.

Therefore, we propose to construct a logical structure tree that organizes all connective phrases in a statement and their textual arguments into a hierarchical structure. We expect the logical structure tree to effectively capture the juxtaposition of connective phrase suggested logical relations and the real logical relations between textual arguments, and therefore guide LLMs in fallacy detection and classification. Specifically, a logical structure tree consists of relation connectives as non-terminal nodes and textual arguments as terminal nodes, and the latter mostly corresponds to elementary discourse units (EDU) considered in discourse parsing. Figure 1 shows the logical structure trees constructed for the two example texts.

As the logical relation indicated by a connective phrase may not be supported by semantics of its arguments in the context, we identify the purposefully indicated logical relations in a context-free unsupervised manner by matching a connective phrase with a taxonomy of connectives compiled for ten common logical relations (conjunction, alternative, restatement, instantiation, contrast, concession, analogy, temporal, condition, causal). To construct a logical structure tree, we first construct a constituency tree for a statement and then search in the constituency tree for connective phrases in the top-down left to right order, and the first found connective phrase will be the root node of the logical structure tree. Next, we identify the text spans of its two arguments using rules and recursively build the left and right sub-trees by applying the same procedure to constituency tree segments corresponding to the two arguments.

The logical structure tree is integrated into LLMs for fallacy reasoning using two strategies. The first considers textualized tree, where we convert the tree into natural language descriptions, making the tree readable by LLMs. Particularly, we describe the relations and arguments in a bottom-up manner, providing the LLMs with insight into logical relations from a local to global perspective. We then concatenate the textualized tree with the instruction prompt, and input them into LLMs as a hard prompt. The second considers tree-based soft prompt, where we derive a relation-aware tree embedding. Specifically, we design relation-specific encoders to process each type of relation and incrementally derive the tree embedding from bottom up to the root node. We then insert

the tree embedding into LLMs as a soft prompt for further tuning. Experiments on benchmark datasets across various domains and genres validate that our approach based on logical structure tree effectively improve precision and recall for both fallacy detection and fallacy classification tasks. Our main contributions are summarized as follows:

. We propose to construct a logical structure tree to capture the juxtaposition of connective phrase suggested logical relations and the real logical relations between textual arguments, and use it to serve as additional guidance for fallacy detection and classification.

. We effectively improve the F1 score for fallacy detection by up to 3.4570 and fallacy classification by np to 6.7SVo across various datasets.

2 Related Work

Logical Fallacy is erroneous patterns of reasoning (Walton, 1987; Fantino et al., 2003). Initial work explored the taxonomy of fallacies (Tindale, 2007; Greenwell et al., 2006; Walton et al., 2008). Recent works have focused on the automatic detection and classification of fallacies. Habernal et al. (2017) developed a software that deals with fallacies in question-answering. Sheng et al. (2021) investigated ad hominem fallacy in dialogue responses. Habernal et al. (2018) explored the ad hominem fallacy from web argumentations. Stab and Gurevych (2017) recognized insufficient arguments in argumentation essays. Goffredo et al. (2022) categorized fallacies in political debates. Nakpiah and Santini (2020) focused on fallacies in legal argumentations. Musi et al. (2022) researched fallacies about pandemics on social medias. (Alhindi et al., 2022) proposed a multi-task prompting approach to learn the fallacies from multiple datasets jointly. Jin et al. (2022) proposed a structure-aware method to classify fallacies. Different from Jin et al. (2022) that masked out content words to form a sequence-based pattern, our paper proposes a tree-based hierarchical logical structure to unify both relation connectives and content arguments together.

Logical Reasoning abilities of large language models are gaining increasing research attention (Xu et al., 2023; Chen et al., 2021; Creswell et al., 2022; Pi et al., 2022;

Jrao et al., 2022; Zhot et al., 2023; Sanyal et al., 2023; Parmar et al., 2024). Olaussen et al. (2023) combined large language models with first-order logic. Pan et al. (2023); Zhang et al. (2023) empowered large language models with symbolic solvers. Pi et al. (2022) presented an adversarial pre-training framework to improve logical reasoning. Zhao et al. (2023) incorporated multi-step explicit planning into the inference procedure. Jiao et al. (2022) proposed a contrastive learning approach to improve logical question-answering. Different from these previous work, we particularly focus on logical fallacy reasoning, aiming to detect and classify fallacies.

Misinformation refers to the unverified or false information (Guess and Lyons, 2020; Armitage and Vaccari, 2021; Aimeur et al., 2023; Lei et al., 2024b). Misinformation detection was studied for years, such as fake news (Rashkin et al., 2017; Lei and Huang, 2023b; Oshikawa et al., 2020), rumor (Ma et al., 2018; Li et al., 2019), satire (Yang et al., 2017), political bias (Lei et al., 2022; Feng et al., 2023; Devatine et al., 2023; Lei and Huang, 2024), propaganda (Da San Martino et al., 2019, 2020; Lei and Huang, 2023a). Logical fallacies are often employed within misinformation to present invalid claim as credible, facilitating the spread of misinformation (Beisecker et al., 2024; Pauli et al., 2022; Bonial et al., 2022). Developing automatic models to detect logical fallacies can also benefit the identification and mitigation of misinformation.

3 Logical Structure Tree

The logical structure tree consists of relation connectives as non-terminal nodes, and textual arguments as terminal nodes. The relation connectives serve as parent nodes, and the two corresponding arguments are linked as left and right children nodes. Figure 1 illustrates examples of the logical structure tree. The logical structure tree is constructed in an unsupervised manner, guided by the constituency tree and a taxonomy of connectives compiled for ten common logical relations.

TABLE CONTENT [ILLEGIBLE]

Table 1: Table 1: The ten types of logical relations and their relation connectives.

3.1 Relation Connectives

The logical fallacies usually rely on relation connectives to indicate a logical relation. Inspired by the discourse relations proposed by Prasad et al. (2008), we define a taxonomy of ten logical relations which are commonly seen: conjunction, alternative, restatement, insertion, continuation, contrast, concession, analogy, temporal, condition, and causal relations. Moreover, we build a set of connective words and phrases that correspond to each type of logical relation, as shown in Table 1. This set of connectives includes the explicit discourse connectives from the PDTB discourse relation dataset (Prasad et al., 2008), and is further expanded by manually adding relevant connectives from the development set of the logic fallacy dataset (Jin et al., 2022).

We further conduct a statistical analysis on the distribution of ten logical relations and compare distributions between fallacy and no fallacy classes as well as across different fallacy classes, with the detailed results shown in Appendix A. The statistical analysis shows that both the fallacy and no fallacy classes contain many connective phrases and their distributions of the ten logical relations are also very similar. But as expected, different fallacy types tend to employ varying logical patterns, for example, False Dilemma uses more alternative relation, while Deductive Fallacy uses more analogy relation.

3.2 Tree Construction Algorithm

To construct a logical structure tree T_{1on4} , we first construct a constituency tree T_{on} for a statement. We use the stanza library² to get the constituency

²<https://stanfordnlp.github.io/stanza/constituency.html>

tree (Qi et al., 2020). At the beginning, T_{1on4} is initialized as an empty tree. Then we traverse the constituency tree T_{on} from top to bottom and from left to right, and match relation connectives within each subtree of T_{on} . If there is a subtree $S_{con}(-)$ whose text equals to a relation connective u , we use

the algorithm in section 3.3 to extract the two textual arguments a, B associated with $tr.r$. Then a new logical subtree $S_{bsu}(u)$ is created, with the matched relation connective u as a parent node, and the two arguments a, B as its left and right children. This new logical subtree $S_{bsi,c}(w)$ is added into the logical structure tree T_{1on4} . If the textual arguments a, B still contain other relation connectives, then we recursively match relation connectives in the arguments and replace the original argument node in T_{1on4} with the newly created logical subtree. The termination condition is that all the relation connectives in the given text have been matched.

3.3 Textual Arguments Extraction

The textual arguments are the two content components linked by a relation connective. Given a matched relation connective $tr.r$, its corresponding subtree in T_{1on4} is $S_{bsi,c}(w)$. To extract the arguments of $tr.r$, we find the parent tree of $S_{bsi,c}(w)$ in the T_{1on4} , denoted as $P(S_{bsi,c}(w))$. The text enclosed by $P(S_{bsi,c}(w))$ is the concatenation of all its leaf node texts. If the text enclosed by parent tree $P(S_{bsi,c}(w))$ contains content before and after the relation connective u , i.e., has the form of $a \text{ } w * d$, then the left argument of $tr.r$ is a and the right argument is d . If the text enclosed by parent tree $P(S_{bsi,c}(w))$ only contains content after the relation connective u , i.e., has the form of $w + d$, then the right argument of $tr.r$ is d , and the left argument w is the text enclosed by grandparent tree $P(P(S_{bsi,c}(w)))$ subtracted by the text enclosed by $P(S_{bsi,c}(w))$.

4 Logical Fallacy Reasoning

We further design a framework to incorporate the logical structure tree into LLMs for fallacy detection and classification. This framework consists of two main components. The first is textualized tree, where we convert the logical structure tree into natural language descriptions, and feed it into LLMs as a hard text prompt. The second is tree-based soft prompt, where we derive a relation-aware tree embedding, and insert it into LLMs as a soft prompt for additional tuning. The hard and soft prompts are complementary: the hard prompt enriches the instruction with logical structure information, while

IMAGE NOT PROVIDED

Figure 2: Figure 2: An illustration of logical fallacy classification informed by logical structure tree.

the soft prompt facilitates direct tuning on tree embeddings. Figure 2 shows an illustration.

4.1 Textualized Tree

The textualized tree aims to transform the logical structure tree into the textual form, which can be interpretable by LLMs. As shown by the upper path of Figure 2, the textualized tree is represented as a table which consists of three columns: left argument, relation connective, right argument. Each row in the table represents a triplet (left argument, relation connective, right argument) corresponding to each logical relation in the tree. In particular, we organize the triplets into the table in a bottom-up order, to provide the LLMs with insight into logical relations from a micro to macro perspective. The textualized tree is then input into the LLMs as a

I'll start! prompt Please classify the fallacy type of the T* tree. Choose one from these options: <fallacy types list>. The definition of each entry is as follows: <fallacy types definition> Text <fallacy type>.

part of the hard text prompt: ht : Text Embedder (textualize, tokenize) (1) where textualize(.) denotes the textualization operation, TextEmbedder refers to the text embedding layer of LLMs, ht is the mapped embedding of the textualized tree.

4.2 Tree-based Soft Prompt

The tree-based soft prompt is a tree embedding which is projected into LLMs as a soft prompt for further tuning. As shown by the lower path of Figure 2, this process includes a tree encoder to derive the tree embedding, as well as a projection layer to transform

the free embedding into the space of LLMs.

During the tree encoder stage, we aim to derive a relation-aware tree embedding. To integrate relation information into free embedding, we design relation-specific encoders to process each type of logical relation. For a simple tree whose children nodes are leaf nodes without hierarchical layers, its embedding is computed as:

$$e_r : W'(et \otimes e_l @ e_r) + b_r \quad (2)$$

where e_r is the embedding of this simple tree, e_l and e_r are the embeddings of left argument, relation connective, and right argument, which are initialized as the average of word embeddings derived from RoBERTa language model (Liu et al., 2019). \otimes denotes feature concatenation, W_r, V_r are the trainable parameters of the encoder that corresponds to the relation type r , where $W_r \in \mathbb{R}^{d_r \times d_r}$, $V_r \in \mathbb{R}^{d_r \times d_r}$, $b_r \in \mathbb{R}^{d_r}$, and $d_r = 768$ is the dimension of embedding space in RoBERTa. The relation type r is one of the ten logical relations associated with the relation connective.

For the tree with hierarchical structure, we derive the tree embedding incrementally, starting from the bottom simple tree and up towards the root node:

$$et : W'(d @ e_l @ e_r) + b_r \quad (3)$$

where e_l is the embedding of the left subtree, d is the embedding of the right subtree, e_l is the connective embedding.

During the projection stage, we transform the tree embedding e_r into the same representation space of LLMs through a projection layer, which includes two layers of neural networks:

$$ht : W_2(w_r @ v_r) + b_r \quad (4)$$

where W_1, W_2, b_1, b_2 are the trainable parameters of the projection layer, $w_r \in \mathbb{R}^{d_r \times d_1}$, $v_r \in \mathbb{R}^{d_r \times d_2}$, $b_1, b_2 \in \mathbb{R}^{d_1}$, d_r is the dimension of hidden states in RoBERTa, d_1 is the dimension of embedding space of the target LLM. d_2 is the resulting tree-based soft prompt, which is then inserted into LLMs as a token representation within the input sequence.

4.3 Fallacy Training

The LLMs take the instruction prompt, textualized tree ht , and free-based soft prompt d_t as input, and generate fallacy label as output. The loss is calculated between the generated text and golden label. The text embedding layer and self attention layers

of LLMs are frozen. The ffee-based soft prompt d1 receives gradients and enables back propagation.

5 Experiments

5.1 Datasets

We experiment with four datasets from various domains and gemes. Table 3 shows their statistics. Argotario (Habernal et al,ZAI7) collects fallacies from the general domain question-answering pairs. The dataset includes the following fallacy labels: Ad Hominem, Appeal to Emotion, Hasty General- ization, Irrelevant Authority, Red Herring, and No Fallacy. We use this dataset for both fallacy detection and classiflcation experiments, and follow the dataset splitting method in Alhindi et al. (2022).

Reddit (Sahai et a1.,2021) collects user generated posts from Reddit, and annotates logical fallacies into: Slippery Slope, Irrelevant Authority, Hasty Generalization, Black- and-White Fallacy, Ad Popu- lum, Tradition Fallacy, Naturalistic Fallacy, Worse Problem Fallacy, and No Fallacy. This dataset is used for both fallacy detection and classification.

Climate (Alhindi et a1.,2022) collects statemenrs from articles in the climate change domain, and annotated the following fallacies: Evading the Burden of Proof, Cherry Picking, Red Herring, Strawman, Irrelev ant Authority, H asty G enerlization, Fals e Cause, False Analogy, Vagueness, and No Fallacy.

Logic (Jin et a7.,2022) annotates logical fallacies in the educational materials into 13 types includ ing Ad Hominem, Ad populm, False Dilemma, False Cause, Circular Reasoning, Deductive Fal- lacy, Ap- peal to Emotion, Equivocation, Fallacy of Extension, Faulty Generalization,Intentionoiroi- lacy, Fallacy of Credibility, Fallacy of Relevance. This dataset does not include No Fallacy class and is only used for fallacy classification.

5.2 Experimental Settings

To validate our approach, we experiment on two types oflanguage models: a decoder-only model and an encoder-decoder model. For the decoder- only model, we choose the open-source large lan- guagemodeLlama Z(llama-2-7b-chat-hf)(Tou- vron et a1.,2023). For the encoder-decoder model, we choose the Flan-T5-large

model (Chung et al., 2022). Both the models are trained in a generative setting, where they take the instruction and given text as input, and generate a fallacy label as output. The fallacy detection task generates "Yes" or "No" label as output, while the fallacy classification task generates the name of each fallacy type. We follow Alhindi et aL. (2022) to unify the different names of the same fallacy across datasets, such as Fclse Dilemma is converted into Black-and-White Fal- lacy since they are the same fallacy. We also follow Alhindi et al. (2022) to feed the definitions of each fallacy type into the instruction prompt, The details of instruction prompt are explained in Appendix B. The maximum input length is set to be 1024, number of epochs is 10, weight decay is 1e-2, the gradient accumulation step is 4, learning rate for Llama-Z is 3e-4, and learning rate forFlan-T5 is 3e- 5. The Llama-2 model is trained with LoRA (Hu et a1.,2021), with rank 8, alpha 16, dropout 0.05, and trainable modules include q_{proj} and v_{proj} .

5.3 Baselines

We compare our models with the baselines listed below. Besides the existing baselines, we also im- plement several additional baselines based on the GPT and RoBERTa (Liu et al'2019) models: Sahai et al. (2021): a multi-granularity network is designed that trains sentence-level representation and the token- level representationsjointly. Jin et al. (2022): a structure-aware framework is de- veloped that forms a sequence-based logical pattern for each text by masking out the content words. Sourati et aJ. (2023b): a prototype-based reason- ing method that injects background knowledge and explainable mechanisms into the language model. Sourati et al. (2023a): a case- based reasoning that retrieves similar cases from external sources based on goals, counterarguments, and explanation etc. Alhindi etat. (2022)ta mutti-task instruction tun- ing framework that rearns the rogaric falacies frJm multiple datasets collaboratively. GPT- 3.5: we prompt the gpt-3.5-turbo model to automatically choose one of the fallacy labels for each text, and the prompt is listed in Appendix C. GPT-3.5 * Ttgc: guide the gpt-3.5-turbo model to firstly reason the logical structure of each text, and then choose one of the fallacy labels through a chain-of+hought process (Wei et a1.,2023). RoBERTa: the RoBERra model

is used to encode the text and the average of word embedding is used as the text embedding. A classification head is built on top of the text embedding to classify labels. RoBERTa * Tostci we concatenate the text embedding with the logical structure tree embedding, and build classification head on top of the combined embedding to predict labels. The tree embedding is derived based on the method in Section 4.2.

5.4 Fallacy Detection

The fallacy detection task identifies whether a given text contains logical fallacy or not, which is a binary classification task. The precision, recall, and F1 score of the fallacy class, as well as the micro F1 score (i.e., accuracy) are used as evaluation metrics. Table 2 presents the performance on the Argotario, Reddit, and Climate datasets. The results demonstrate that incorporating the logical structure tree effectively improves both precision and recall for logical fallacy detection. This observation is consistent for both types of Llama-2 and Flan-T5 models across all the three datasets, which span various domains and genres. Compared to the baselines without logical structure information, our approach based on logical structure tree notably enhances the precision and recall, leading to the F1 score increased by up to 3.45%. This indicates that the logical structure tree is effective in capturing the difference in logical flows between fallacious and benign texts.

Moreover, informing the large language model GPT-3.5-turbo of logical structure information significantly improves fallacy detection under the zero-shot setting, resulting in a substantial improvement in the F1 score. This underscores the importance of incorporating the logical structure information into LLM for fallacy detection. Also, concatenating the logical structure tree embedding with the text embedding in the RoBERTa model also enhances the performance, which proves the usefulness of this logical structure tree embedding. Overall, incorporating the logical structure tree helps improve fallacy detection for various types of models.

5.5 Fallacy Classification

The fallacy classification task classifies the fallacy types for the fallacious text, which is a multi-class classification task excluding the No Fallacy class. The macro precision, recall, and F1 score, as well as the micro F1 score (i.e., accuracy) are used as evaluation metrics. Table 4 shows the results on the Argotario, Reddit, and Logic datasets. The results demonstrate that integrating the logical structure tree into Llama-2 and Flan-T5 models notably enhances the performance of fallacy classification, with both precision and recall increased. This conclusion is valid across the three datasets from different domains and genres. Compared to the baselines without logical structure tree, our proposed approach significantly improves precision and recall, leading to an increase of up to 6.75% in the F1 score. This suggests that the logical structure tree effectively distinguishes the different logical patterns used in each fallacy type, and is applicable across various domains and genres.

In addition, our approach based on the logical structure tree outperforms the previous methods that may lack logical relations information. This highlights the necessity to infuse the logical relations into LLMs for fallacy classification. Besides, our approach achieves higher performance than the baselines that overlook content words. This indicates that analyzing content words also plays an essential role in fallacy reasoning. The logical structure tree connects the logical relations and content arguments together to form a cohesive logical structure, representing the hierarchical logical flow and thereby improving fallacy classification.

5.6 Ablation Study

The ablation study of the two designed strategies to incorporate the logical structure tree into LLMs is shown in Table 5, where we take Llama-2 model as an example. The upper rows show the results of fallacy detection on the three datasets, and the lower rows show the results of fallacy classification. The results demonstrate that both the textualized tree and tree-based soft prompt brings improvement for fallacy detection and classification across multiple datasets. This proves that the textualized tree and tree-based soft prompt are complementary with each other: the

textualized tree enriches the instruction prompt with logical structure information, and the tree-based soft prompt enables direct learning from the tree embedding. Comparing across these two strategies, the soft prompt usually achieves better performance than the hard text prompt, and exhibits higher recall. Combining the two strategies together leads to the best performance, achieving the highest precision and recall.

5.7 Effect on Different Fallacy Types

We further analyze the F1 score change across each fallacy type in the fallacy classification task. The Llama-Z model is used as an example to show the performance change before and after incorporating the logical structure tree. Table 6 presents the F1 score change across each fallacy type on Argotario dataset. The performance change across each fallacy type on the Reddit and Logic dataset are shown in the Table 7 and Table 8. We observe that the logical structure tree brings bigger improvements for the fallacy types such as Red Herring, Hasty Generalization, Irrelevant Authority, Ad Populum, Extension Fallacy, Equivocation, Circular Reasoning etc. One possible explanation is that these fallacy types usually employ certain logical relations or logical patterns to persuade the readers. However, the performance increase is less noticeable for the fallacy types such as Appeal to Emotion and Ad Hominem. It may due to the reason that these fallacies rely more on the emotional or sentimental language instead of logical relations.

6 Limitations

We have compiled a set of connective words and phrases for the ten logical relations, as detailed in Table 1. While we have included the common connectives in this set, it may not contain all the possible connectives. The logical structure tree that is constructed based on this connective words set demonstrates its usefulness in fallacy reasoning. Future work can be expanding this connectives set and investigating the effects of various connectives, so that we can better identify and mitigate them. The release of code, datasets, and model should be used for mitigating log-

ical fallacies, instead of expanding or disseminating the misinformation.

Acknowledgements

We would like to thank the anonymous reviewers for their valuable feedback and input. We gratefully acknowledge support from National Science Foundation via the award trS2127746. Portions of this research were conducted with the advanced computing resources provided by Texas AM High-Performance Research Computing.

7 Conclusion

This paper detects and classifies fallacies. We propose a logical structure tree to explicitly represent and track the hierarchical logic flow among relation connectives and their arguments. We also design two strategies to incorporate this logical structure tree into LLMs for fallacy reasoning. Extensive experiments demonstrate the effectiveness of our approach based on the logical structure tree.

Ethical considerations

This paper aims to detect and classify logical fallacies. Logical fallacy is the error or flaws in the reasoning, and can occur in various human communications. Logical fallacies can lead to harmful consequences for society such as spreading misleading information or introducing societal bias. The goal of this research is to understand logical fallacies,

References

- Rehab Mohamed Ahmed Abd-Eldayem. 2023. The relationship between cognitive bias and logical fallacies in Egyptian society. *Social Sciences*, 12(6):281-293.
- Esma Aimeur, Sabrine Amri, and Gilles Brassard. 2023. Fake news, disinformation and misinformation in social media: a review. *Social Network Analysis and Mining*, 13(1):30.
- Tariq Alhindi, Tuhin Chakrabarty, Elena Musi, and Smaranda Muresan. 2022. Multitask instruction-based prompting for fallacy recognition. In Proceedings of the 11th International Conference on Computational Linguistics (COLING 2022), Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Rachel Armitage and Cristian Vaccari. 2021. Misinformation and disinformation. In *The Routledge companion to disinformation*.

- panion to rmedia disinformation and populisi, pages 38-48. Routledge.

Donald A Barclay. 2018. Fake news, propaganda, and plain old lies: how to find trustworthy information in the digital age. Rowman Littlef, eld.

Sven Beisecker, Christian Schlereth, and Sebastian Hein. 2024. Shades offake news: How fallacies influence consumers' perception. European Journal of Information Systems, 33(1):41-60.

Claire Bonial, Austin Blodgett, Taylor Hudson, Stephanie M. Lukin, Jeffrey Micher, Douglas Summers-Stay, Peter Sutor, and Clare Voss. 2022. The search for agreement on logical fallacy annotation of an infodemic. In Proceedings of the Thirteenth Language Resources and Evaluation Conference, pages 44304438, Marseille, France. European Language Resources Association.

Zeming Chen, Qiyue Gao, and Lawrence S. Moss. 2021. Neurallog: Natural language inference with joint neural and logical reasoning. In Proceedings of *SEM 2021: The Tenth Joint Conference on Lexical and Computational Semantics, pages 78-88, Online. Association for Computational Linguistics.

Hyung Won Chung, Le Hou, Shayne Longpre, Banet Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wel 2022. Scaling instruction-finetuned language models. Preprint, arXiv :2210.11416.

Christian Cotton. 2018. Argument from fallacy. Bad arguments: 100 of the most important fallacies in Western philosophy, pages 125-127 .

Antonia Creswell, Murray Shanahan, and Irina Higgins. 2022. Selection-inference: Exploiting large language models for interpretable logical reasoning. Preprint, arXiv:2205.09712.

Giovanni Da San Martino, Alberto Barr6n-Cedeflo, Henning Wachsmuth, Rostislav Petrov, and Preslav Nakov. 2020. SemEval-2020 task 11: Detection of propaganda techniques in news articles. In Proceedings of the Fourteenth Workshop on Semantic Evaluation, pages 1377-1414, Barcelona (online). International Committee for Computational Linguistics.

Giovanni Da San Martino, Seunghak Yu, Alberto Barr6n-Cedeflo, Rostislav Petrov, and Preslav Nakov. 2019. Fine-grained analysis of propaganda in news article. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP - IJCNLP.), pages 5636-5646, Hong Kong, China. Association for Computational Linguistics.

Nicolas Devatine, Philippe Muller, and Chloé Braud. 2023. An integrated approach for political bias prediction and explanation based on discursive structure. In Findings of the Association for Computational Linguistics : ACL 2023, pages 11196-11211, Toronto, Canada. Association for Computational Linguistics.

Edmund Fantino, Stephanie Stolarz-Fantino, and Anton Navarro. 2003. Logical fallacies: A behavioral approach to reasoning, The Behavior Analyst Today, 4(1):109,

Shangbin Feng, Chan Young Park, Yuhua Liu, and Yulia Tsvetkov. 2023. From pretraining data to language models to downstream tasks: Tracking the trails of political biases leading to unfair NLP models. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1 : Long Papers), pages 11737-11762, Toronto, Canada. Association for Computational Linguistics.

Pierpaolo Goffredo, Mariana Chaves, Serena Villata, and Elena Cabrio. 2023. Argument-based detection and classification of fallacies in political debates. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 11101-11112, Singapore. Association for Computational Linguistics.

Pierpaolo Goffredo, Shokeh Haddadan, Vorakit Vorakittiphan, Elena Cabrio, and Serena Villata. 2022. Fallacious argument classification in political debates. In IJCNLP, pages 4143-4149.

William S Greenwell, John C Knight, C Michael Holloway, and Jacob J Pease. 2006. A taxonomy of fallacies in system safety arguments. In 24th International System Safety Conference.

Andrew M Guess and Benjamin A Lyons. 2020. Mis-information, disinformation, and online propaganda. Social media and democracy: The state of the field, prospects for reform, 10.

Ivan Habernal, Raffael Hannemann, Christian Polak, Christopher Klamm, Patrick Pauli, and Iryna Gurevych. 2017. Argotario: Computational argumentation meets serious games. In Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing : System Demonstrations, pages 7–12, Copenhagen, Denmark, Association for Computational Linguistics,

Ivan Habernal, Henning Wachsmuth, Iryna Gurevych and Benno Stein. 2018. Before name-calling: Dynamics and triggers of ad hominem fallacies in web argumentation. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pages 386–396, New Orleans, Louisiana. Association for Computational Linguistics.

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeytan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. Preprint, arXiv [ILLEGIBLE]

Zhitong Jin, Abhinav Lalwani, Tejas Vaidhya, Xiaoyu Shen, Yiwen Ding, Zhiheng Lyu, Mrinmaya Sachan, Rada M [ILLEGIBLE]

Fangkai Jiao, Yangyang Guo, Xuemeng Song, and Liqiang Nie. 2022. MERIT: Meta-Path Guided Contrastive Learning for Logical Reasoning. In Findings of the Association for Computational Linguistics : ACL 2022, pages 3496–3509, Dublin, Ireland. Association for Computational Linguistics.

[ILLEGIBLE]

Elena Musi and Chris Reed. 2022. From fallacies to semi-fake news: Improving the identification of misinformation triggers across digital media. Discourse Society, 33(3):349–370.

Callistus Ireneous Nakpah and Simone Santini. 2020. Automated discovery of logical fallacies in legal argumentation. International Journal of Artificial Intelligence and Applications (IJNA), 11.

Theo Olausso [ILLEGIBLE]

Christopher W Tindale. 2007. Fallacies and argument appraisal. Cambridge University Press.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,

TABLE CONTENT [ILLEGIBLE]

Table 2: Table 9: The ratio (Vo) of samples that contain the ten logical relations in fallacy and, no fallacy classes in the development set of Argo

TABLE CONTENT [ILLEGIBLE]

Table 3: Table 10: The ratio (Vo) of samples that contain the ten logical relations in each fallacy type in the development set of Argo

Cynt [ILLEGIBLE]

[ILLEGIBLE]

A Statistical Analysis of Logical Relations C Prompt for GPT-based baselines

Table 9 presents the ratio of samples that contain the ten logical relations in fallacy and no fallacy classes, where we take the Argotario (Habernal et al., 2017) and Reddit (Sahai et al., 2021) datasets as examples. Further, Table 10 shows the ratio of samples that contain the ten logical relations in each fallacy type

B The Names and Definitions of Fallacies

B.1 Argotario dataset

The Argotario dataset (Habernal et al., 2017) includes five fallacy types: Ad Hominem, Appeal to Emotion, Hasty Generalization, Irrelevant Authority, Red Herring. The name of Appeal to Emotion is converted into Emotional Language. The definitions of these fallacy types which are used in the instruction prompt are:

- . Ad Hominem: [MISSING]
- . Emotional Language: [MISSING]
- . Hasty Generalization: [MISSING]
- . Irrelevant Authority: [MISSING]
- . Red Herring: [MISSING]

B.2 Reddit dataset

The Reddit dataset (Sahai et al., 2021) includes eight fallacy types and one no fallacy type. The fallacy types include: Slippery Slope, Irrelevant Authority, Hasty Generalization, Black-and-White Fallacy, Ad Populum, Tradition Fallacy, Naturalistic Fallacy, Worse Problem Fallacy. The definitions of these fallacy types which are used in the instruction prompt are: [ILLEGIBLE]

B.3 Climate dataset

The Climate dataset (Alhindi et al., 2022) includes nine fallacy types and one no fallacy type. The fallacy types include: Evading the Burden of Proof, Cherry Picking, Red Herring, Strawman, Irrelevant Authority, Hasty Generalization, False Cause, False Analogy, Vagueness. The definitions of these fallacy types which are used in the instruction prompt are: [ILLEGIBLE]

B.4 Logic dataset

The Logic dataset (Jin et al., 2022) includes 13 fallacy types: Ad Hominem, Ad Populum, False Dilemma (Black-and-White Fallacy), False Cause, Circular Reasoning, Deductive Fallacy, Appeal to Emotion (Emotional Language), Equivocation, Fallacy of Extension, Faulty Generalization (Hasty Generalization), Intentional Fallacy, Fallacy of Credibility (Irrelevant Authority), Fallacy of Relevance (Red Herring). The names in the parenthesis are the replaced names used in the instruction prompt. The definitions of these fallacy types which are used in the instruction prompt are: . Ad Hominem: the text attack a person instead of arguing against the claims. . Ad Populum: the text affirm something is true

because the majority thinks so, . Black-and-White Fallacy: the text present two alternativ [ILLEGIBLE]