

Pattern-based Logical Fallacy Classification using Decoder-Only Large Language Models

Anonymous ACL submission

Abstract

In today’s fast-paced information era, logical fallacies, defined as defective patterns of reasoning, inevitably contribute to the growth of information disorder. However, often fallacies appear in nuanced forms that complicate automated classification. In this study, we investigate whether the logical structure of arguments proves beneficial for fallacy classification by developing a framework that extract logical patterns using Large Language Models (LLMs). We evaluate the impact of these patterns across different LLMs and experimental zero- and one-shot configurations, showing statistically significant improvements over zero-shot baselines and outperforming competing approaches. Cross-dataset experiments validate generalization, establishing data-driven pattern extraction as an effective method for generating logical representations.

1 Introduction

A logical fallacy is a common error of thinking, especially one apt to mislead (Gensler, 2010). These arguments often appear rational and logically coherent on the surface, but deeper analysis reveals they are not (Copi et al., 1953). Fallacies are traditionally classified into formal and informal types: formal fallacies violate the rules of logical structure regardless of content, while informal fallacies are patterns of mistakes that are made in the everyday uses of language and are related to contextual meaning (Hamblin, 1970; Bacon et al., 1999; Copi et al., 1953).

To evaluate the quality of an argument, it is helpful to reconstruct it into what is known as logical form, the structure that emerges when the specific content of a statement is replaced by variables (Johnson and Blair, 1977). For example, the argument *If it rains, then the ground will be wet. It is raining. Therefore, the ground is wet* has the logical form *If P, then Q. P. Therefore, Q*. Building

on this formalization framework, Jin et al. (2022) developed a structure-aware model for fallacy detection on the LOGIC dataset that compares arguments’ and fallacies’ logical forms. However, in their approach a single logical form is assigned to each fallacy, which might fail to capture the full spectrum of ways a fallacy can manifest in natural discourse. Another challenge is related to informal fallacies, where reasoning is often more nuanced and context-dependent than abstract representations suggest.

This gap between theory and practice raises a key question: **How valuable is the logical structure of arguments for automated fallacy detection?** To answer this question, we exploit Large Language Models (LLMs) to extract fallacies’ logical patterns from fallacious examples and their explanations, attempting to capture both logical forms and context-aware logical schemes that reveal the underlying mechanisms of deception in fallacious arguments, an approach designed to overcome Jin et al. (2022)’s limitations. Hereafter, we refer to these extracted structures collectively as *patterns*. While existing supervised approaches require heavy computational resources for fine-tuning (Vijayaraghavan and Vosoughi, 2022; Lei and Huang, 2024; Sourati et al., 2023a,b; Alhindi et al., 2024), to our knowledge, no prior work has explored fallacy classification from a structural perspective without any additional fine-tuning. Although we use labeled data for pattern extraction, our approach avoids fine-tuning costs and produces generalizable patterns that allow classification through prompting alone, enabling comparison with unsupervised methods.

We evaluate multiple prompting configurations to determine which informational components enhance performance and examine the impact of demonstrations on detection capabilities. Our approach, incorporating the generated patterns, achieves outstanding results among unsupervised

methods on the dataset LOGIC. Finally, to validate the robustness and transferability of our patterns, we assess their performance across two further datasets spanning diverse domains and argumentative styles.

In summary, our contributions are threefold:

- We leverage LLMs to extract patterns from fallacious examples, which are then employed in inference-only classification.
- We evaluate different LLMs with diverse prompt designs to identify which informational components optimize fallacy classification.
- We validate generalizability by testing our patterns on two different datasets with different domains and structures.

2 Datasets

The LOGIC dataset is a collection of 2449 examples across 13 fallacy types. Instances are sourced from educational platforms about fallacies such as Quizziz and study.com. The dataset consists of brief dialogues and short statements. Given the educational intent behind these examples, sentences tend to have relatively straightforward syntactic structures, making the dataset particularly well-suited for pattern recognition and alignment with logical forms.

Although it contains 13 distinct classes, a thorough analysis revealed that some of the classes actually contain instances of different fallacies that were grouped together. For instance, the class *Hasty Generalization* contains examples of actual *Hasty Generalization* as well as *Slippery Slope* (Table 1). While these grouped fallacies share common logical flaws and thus belong to the same conceptual group, they manifest through different structural patterns, which complicates the attempt to match them all to a single logical scheme.

We experiment on two further datasets: REDDIT (Sahai et al., 2021), consisting of fallacious comments extracted from subreddits covering different topics and ELECDEBATE60TO16 (hereafter ELECDEBATE) (Goffredo et al., 2023), a collection of televised debates of the presidential election campaigns in the U.S. from 1960 to 2016. Some fallacy classes contain sub-categories. Datasets’ summary is reported in Table 2 and a description of each taxonomy is provided in Appendix E.

Class	Fallacies included
Intentional Fallacy	Intentional Fallacy Shifting the Burden of Proof Moving the Goalposts No True Scotsman
False Cause	Post Hoc False Cause
Hasty Generalization	Hasty Generalization Slippery Slope

Table 1: Examples of classes in LOGIC containing instances of different fallacy types. While logically coherent, these groupings comprise fallacies with distinct structural patterns. A detailed breakdown of all classes’ subtypes is provided in Appendix D of Jin et al. (2022).

Data	Dataset split	# Classes	Genre	Domain
LOGIC	1807/290/290	13	Dialogue	Education
REDDIT	588/148/105	8 [‡]	Comments	General
ELECDEBATE	1120/200/187	6	Dialogue	Politics

Table 2: Statistics of the three datasets. ‡ indicates that the *No Fallacy* class is included.

3 Pattern generation

Natural arguments appear in several different forms. Such variability manifests itself in LOGIC dataset as well as many others (Habernal et al., 2018; Da San Martino et al., 2019). For this reason, we address our research question by modeling patterns inductively from the observed text instances. The choice of the LOGIC dataset for pattern extraction is critical: it provides the required combination of structural clarity and fallacy diversity through its multiple sub-types per class. These properties make it especially suited for our purpose: the clean argumentative structure allows to formalize clear logical patterns while capturing intra-class variations.

Our pattern generation procedure features two steps:

Step 1: Explanation Generation Explanations have been shown to be instrumental in identifying and discrediting fallacious reasoning, as they make the logical structure of arguments explicit and open to scrutiny (Storer, 1949). Furthermore, Jeong et al. (2025) has demonstrated that providing explanations constitutes valuable contextual information in zero-shot settings. We expected explanations to facilitate pattern extraction by breaking down the reasoning process and revealing shared reasoning flaws, particularly useful for informal fallacies.

Given a sentence from the training set and its

fallacy label, we used llama-3.3-70B to generate an explanation that justifies why that sentence contains the specified fallacy.

Step 2: Pattern Extraction For each fallacy class, we used OpenAI’s reasoning model o4-mini (OpenAI, 2025) to extract patterns from the collected sentences and their explanations, requiring the model to preserve logical connectors such as prepositions or adverbs and to abstract away from content words by replacing them with placeholders while keeping the original reasoning form. Additionally, summaries were extracted to derive new fallacy definitions.

We opted for Llama-3.3-70B-Instruct for explanation generation as it provided high-quality explanations while remaining cost-effective for large-scale text generation. For pattern extraction, we employed o4-mini given its reasoning capabilities.

The prompts used in our experiments are reported in Appendix G.1. In the initial phase of our research, we aimed to cover two distinct logical aspects from our arguments and explanations, specific to formal and informal fallacies, respectively:

- arguments’ **logical forms** as defined by formal logic theory;
- recurring **reasoning schemes** that frequently appear in both sentences and explanations, capturing specific information about the reasoning behind the fallacy, including frequent syntactic particles, phrases, and examples that convey the fallacious intent.

Our patterns incorporate both logical forms and reasoning schemes, as Table 3 shows. These patterns differ from logical forms and definitions by combining reasoning structure (variables X, Y, Z) with concrete linguistic features (specific phrases, loaded terms, rhetorical devices). The full list of patterns is available in Appendix F (Table 20).

The process resulted in approximately 3-6 patterns per fallacy class. Final patterns were obtained after providing different subsets to the model and selecting the best performing one on the validation set, in the attempt to retain only useful information and avoid redundancy. In section 2 we discussed how one class in the datasets could correspond to multiple fallacies. Although in some cases, e.g. a pattern for *Tu quoque* (a fallacy which is part of the class *Ad Hominem* in LOGIC), is correctly generated and selected, sometimes fails to select patterns when multiple fallacies are grouped under

the same class label. This is expected because we include the fallacy class name in the prompt, which likely biases the model toward patterns that match its internal knowledge of that particular class name. To ensure a broader coverage of fallacies listed in Table 1, we manually isolated instances of frequent and undetected fallacies (such as *Shifting the Burden of Proof*) and repeated the procedure.

Intentional Fallacy Patterns

1. The argument assumes that because X (e.g., someone’s intention, belief, or lack of counter-evidence), therefore Y is true.
 2. Asserting P is true because it has not been disproven.
 3. Because the creator intended [interpretation], the work should be understood as [interpretation].
 4. Questions framed to presuppose guilt or a specific intention (e.g., “Have you stopped X?”), thus assuming what is to be proven.
 5. If A does not have trait X, and X is (allegedly) typical of group G, then A is not a member of G.
-

Table 3: Patterns for *Intentional Fallacy* combining reasoning schemes (#4) and logical forms (#5) that encode structure and intent.

4 Experiments

This section describes our experiments for fallacy detection, including our patterns extracted by the procedure introduced in Section 3 and several competing prompting strategies. Additional experiments are reported in Appendix C. We used the following LLMs for our experiments: gpt-4o, o4-mini, gpt-4.1-mini, Llama-3.3-70B, deepseek-r1 and Gemma-3-27B-it for a total cost of 75 USD. Our intent was to test LLMs from different providers and with different sizes and to compare reasoning and non-reasoning models.

4.1 Prompt Design

Baselines We compared our approach against several baselines that vary in the type and amount of information provided to the model. The simplest baseline (**ZERO-SHOT**) provides only the list of fallacy names in the dataset as a reference, establishing a minimal information condition. Our second baseline incorporates fallacy definitions to provide more comprehensive background knowledge (**DEF**). These definitions were initially sourced from Lei and Huang (2024) and subsequently refined based on our analysis to ensure clarity and consistency. Finally, we tested a baseline using standard logical forms, following the approach of Jin et al. (2022) and sourcing these forms from logicallyfallacious.com. This final baseline (**LOGICAL FORMS**) allows us to assess the effectiveness of

expert-made logical representations compared to our generated pattern-based approach.

LLM-derived Patterns and Definitions Beyond generating logical patterns, we leveraged the explanations from Section 3 to automatically create new fallacy definitions based on LOGIC training samples (see prompt in Appendix G.1). We then replicated experiment DEF with these new definitions (**NEW DEF**). We also exploited the patterns extracted by adding them to the prompt (**PATTERNS**) and by implementing a two-step approach where we first ask the LLM to identify the pattern and then to output the corresponding fallacy (**PATTERN MATCHING**).

One-shot Prompting We further investigated the impact of providing examples to the model through several experimental configurations (Brown et al., 2020), with one-shot prompting proving most effective. Initially, we tested a static approach where one example per fallacy was randomly selected and shown to all test sentences (**ONE-SHOT**), establishing a baseline for example-based learning. To enhance this approach, we augmented the same examples with manually crafted explanations following our previously established definitions as guidelines (**ONE-SHOT + EXP**). We sampled 5 different example sets and performance across all configurations was assessed over 5 runs to ensure statistical reliability.

More sophisticated was our dynamic one-shot prompting approach (**DYNAMIC ONE-SHOT**), which computes embeddings for both training and test sentences to retrieve the most similar example per class for each test sentence. We used sentence-transformers/all-MiniLM-L6-v2 model and cross-encoder/stsb-roberta-base cross-encoder from SentenceTransformers (Reimers and Gurevych, 2019) to compute embeddings and employed cosine similarity to evaluate similarity. We included the previously generated explanations of examples in the prompt as well (**DYNAMIC + EXP**).

Furthermore, we explored structure-focused similarity. Since Jin et al. (2022) released a version of LOGIC with masked arguments (with content words replaced by placeholders), we conducted the same similarity-based procedure using these masked sentences (see an example in Table 4) in an attempt to force the embedding model to focus on structural rather than lexical

Original argument	Every time I wear this necklace, I pass my exams. Therefore, wearing this necklace causes me to pass my exams.
Masked argument	Every time MSK<0> MSK<2>, MSK<0> MSK<4>. Therefore, MSK<2> causes MSK<0> to MSK<4>.

Table 4: Example of a masked argument in LOGIC. The distillation algorithm is explained in Jin et al. (2022). The masked version of the dataset was publicly released by the authors and was not created by us.

similarities. For this configuration (**SYNTAX-BASED DYNAMIC ONE-SHOT**), we used sentence-transformers/all-MiniLM-L6-v2 from SentenceTransformers alongside a syntax-augmented version of RoBERTa-large extracted from Sachan et al. (2021) (see Appendix D).

Finally, we incorporated the generated patterns into our dynamically retrieved examples and their explanations (**DYNAMIC + EXP + PATTERNS**).

Multi-step Classification An alternative approach involves decomposing the classification task into three sequential steps within a single model call (**MULTISTEP**) using chain-of-thought prompting (Wei et al., 2023). In the first step, the model is required to generate a logical form representation of the argument according to predefined structural rules (prompt in Appendix G.2). Subsequently, the model should match the generated logical form to one of the patterns and, as a result, classify the argument.

4.2 Results and discussion

Table 5 summarizes all experimental configurations and results on LOGIC: it reveals a consistent improvement when the model leverages information about the underlying logic extracted through the LLMs, especially with reasoning models and gpt-4o. When using reasoning models, the model-generated definitions yield a 4.65% accuracy improvement over our manually corrected definitions. In the same way, including our generated patterns causes a 8.2% increase with respect to the logical forms extracted by the website logicallyfallacious.com and used in Jin et al. (2022). McNemar’s test proved statistical significance for all models using PATTERNS against ZERO-SHOT and for all except llama and deepseek against LOGICAL FORMS method. When it comes to non-reasoning models, the different definitions do not really affect the performance, whereas using our patterns improves the accuracy by 5.8% on average. A notable result is the performance increase

Method	o4-mini		gpt-4o		deepseek-r1		gpt-4.1-mini		llama-3.3-70B		gemma-3-27b-it	
	Acc.	F ₁	Acc.	F ₁	Acc.	F ₁	Acc.	F ₁	Acc.	F ₁	Acc.	F ₁
<i>Baselines</i>												
ZERO-SHOT	61.7	55.3	62.7	57.0	62.7	57.3	57.8	51.0	55.8	47.7	60.5	51.3
DEF	62.1	58.7	65.0	58.7	62.2	56.5	57.5	50.6	59.1	51.5	63.5	55.2
LOGICAL FORMS	63.2	57.4	65.4	59.4	63.1	55.4	57.8	49.4	60.2	51.3	62.8	53.9
<i>LLM-derived Patterns and Definitions</i>												
NEW DEF	66.8	67.3	66.8	59.9	66.8	60.0	57.5	52.5	58.8	53.3	64.8	57.7
PATTERNS	72.2	66.4	73.2	64.9	70.5	66.2	63.5	55.7	64.5	53.3	68.5	61.9
PATTERN MATCHING	70.1	65.9	<u>73.5</u>	<u>66.5</u>	71.5	66.5	65.2	57.9	66.2	59.6	67.2	59.9
<i>One-shot prompting</i>												
ONE-SHOT	63.6	59.6	64.1	58.7	58.5	55.9	56.2	48.1	56.1	46.2	60.0	49.7
ONE-SHOT + EXP	65.2	59.5	63.5	59.0	45.7	48.6	56.8	50.0	56.3	47.9	59.2	49.7
DYNAMIC ONE-SHOT												
all-MiniLM-L6-v2	70.2	67.6	71.3	66.4	70.4	66.2	65.8	61.7	65.5	59.7	68.5	63.3
roberta-base	69.5	64.3	69.5	64.6	72.5	67.8	65.5	61.3	64.8	58.9	66.5	60.6
SYNTAX-BASED DYNAMIC ONE-SHOT												
all-MiniLM-L6-v2	68.2	63.6	71.2	66.1	68.5	64.2	63.2	58.3	62.8	55.7	64.5	57.3
syntax-augmented roberta-large	65.5	65.5	71.2	66.3	68.5	64.2	64.5	58.6	64.2	56.5	63.5	56.0
DYNAMIC + EXP	71.2	68.9	69.5	65.0	72.7	67.9	67.8	61.0	67.5	62.2	68.2	63.2
DYNAMIC + EXP + PATTERNS	<u>74.2</u>	<u>68.9</u>	<u>73.1</u>	<u>67.2</u>	<u>73.2</u>	<u>67.9</u>	66.8	62.3	67.2	55.1	<u>70.5</u>	<u>65.9</u>
<i>Multi-step classification</i>												
MULTISTEP	65.4	64.9	70.9	62.7	62.2	57.1	65.8	55.8	62.5	55.1	66.8	60.2

Table 5: Fallacy classification performance on LOGIC. **Bold**: best approach in section per model by accuracy, **Bold**: best approach overall per model by accuracy. F₁ score denotes Macro F₁ score, which accounts for the class imbalance in the dataset.

achieved through dynamic one-shot prompting. In particular, **DYNAMIC ONE-SHOT** approach (using all-MiniLM-L6-v2) yields an average 8.87% increase in accuracy compared to **ONE-SHOT**, despite relying on semantic similarity for example selection. On the other hand, the syntax-oriented example retrieval strategy (**SYNTAX-BASED DYNAMIC ONE-SHOT**) does not outperform the semantic selection. This may be partially due to inaccuracies in the sentence masking process, which can negatively impact the retrieval of similar examples and the classification, consequently. The **MULTISTEP** approach shows weaker performance than **PATTERN MATCHING**, especially for deepseek-r1, implying that generating logical forms without explicit guidance constitutes the main challenge for the model in the request.

In summary, performance benefits from structure-based information, indicating that incorporating logical reasoning structure into prompts enhances fallacy classification and achieves competitive results: **PATTERNS** with gpt-4o reaches 73.5% accuracy, outperforming prior unsupervised methods (Table 6), while **DYNAMIC+EXP+PATTERNS** with o4-mini achieves 74.2% when augmented with examples and patterns.

4.3 Error analysis

Pattern matching Requesting the model to identify the closest pattern for each argument provides insight into the association process between sentences and patterns. For our analysis, we have split our fallacies into two groups in Table 7: i) group 1, consisting of fallacies whose patterns include logical forms while still including additional contextual cues; ii) group 2, consisting of fallacies that lack highly structured patterns and rely more on contextual and semantic features of the sentence.

Figure 1 shows consistently superior accuracy for Group 1, whose classes maintain relatively high performance across all experimental settings. The class *Circular Reasoning* emerges as the most accurately predicted class across all models. For what concerns Group 2, the overall accuracy is, on average, 22% lower with respect to Group 1.

Method	Acc	F ₁
Jeong et al. (2025)	49.0	37.0
Pan et al. (2024)	-	50.5
PATTERNS (gpt-4o)	73.5	66.5
DYNAMIC+EXP+PAT. (o4-mini)	74.2	68.9

Table 6: Comparison of our best results against the unsupervised baselines provided by Jeong et al. (2025) and Pan et al. (2024) (described in Appendix B) for LOGIC.

Group 1	Group 2
<ul style="list-style-type: none"> • Ad Hominem • Ad Populum • Circular Reasoning • Irrelevant Authority • False Cause • Hasty Generalization • Deductive Fallacy • Black-and-White Fallacy 	<ul style="list-style-type: none"> • Red Herring • Equivocation • Emotional Language • Extension Fallacy • Intentional Fallacy

Table 7: Grouped fallacy classes based on pattern features for analytical purposes.

The classes *Emotional Language*, *Red Herring* and *Extension Fallacy* achieve moderate prediction accuracy, whereas only *Evading the Burden of Proof*'s patterns within the *Intentional Fallacy* category are correctly classified, and *Equivocation* remains entirely undetected by gpt-4.1-mini. In summary, the models achieve better performance on logical fallacies that exhibit clearer structural characteristics but face difficulties with fallacies requiring more nuanced semantic understanding and contextual analysis.

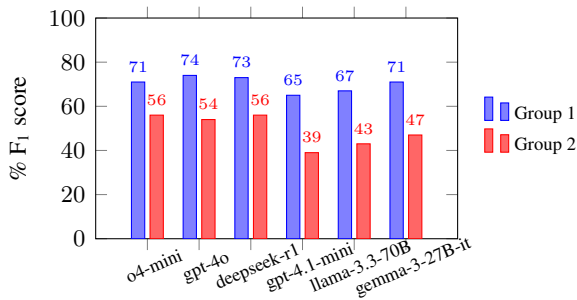


Figure 1: Group-wise F1 score for each model, relative to the PATTERN MATCHING prompt setting.

Furthermore, matching patterns allows us to see that some instances can be deemed as fitting from a structural point of view, thus partially explaining the inherent difficulty of the classification task. While providing guidance through logical structure proves beneficial for fallacy detection, this approach does not eliminate all sources of ambiguity, as some sentences may conform to multiple structural patterns. The critical point lies in context-aware pattern application: models must not only identify logical forms but also evaluate their contextual validity in each sentence.

To quantify the degree of ambiguity inherent in pattern matching, we instructed the best-performing model o4-mini to return the five most similar patterns for each argument. This multi-candidate approach enables us to analyze whether lower-ranked patterns might also represent valid

interpretations of the same argument. By examining the distribution of pattern similarities and evaluating classification accuracy when considering alternative matches, we can better understand the boundaries of pattern-based classification and identify instances where structural ambiguity genuinely complicates fallacy detection.

Acc@1	Acc@2	Acc@3	Acc@4	Acc@5
66.7	75.1	81.8	86.5	88.5

Table 8: Performance analysis in PATTERN MATCHING with expanded solution pool: classification results including top 5 predictions as correct.

Table 8 shows that, when the model is prompted to return multiple matching patterns rather than a single best match, its confidence in the initial prediction decreases, resulting in a 3.4% drop in accuracy (see Table 5). However, this apparent degradation is misleading when viewed in isolation. By incorporating the second-ranked pattern choice into our evaluation, performance recovers to 75.1%, and continues to improve as we expand our candidate pool to include progressively lower-ranked options. Table 10 illustrates a representative case where the model successfully identifies the correct pattern as its second choice, while its first-ranked selection remains structurally plausible: the model likely assigns one of the *Ad Populum* patterns because it closely matches the argument's logic, while the *Irrelevant Authority* pattern doesn't fit the sentence since it requires discussion of an unrelated topic, which is not present in the sentence. These subtle distinctions likely make pattern matching more challenging than direct classification because it requires strict structural alignment as well as capturing broader content-related features.

Multistep classification The MULTISTEP approach fails to produce significant results. We conduct this experiment in a single passage to force the model to reason using both semantic and syntactic information. However, classification performance depends critically on the quality of the extracted logical forms, which proves inconsistent and model-dependent. For instance, o4-mini embeds classification-relevant contextual information directly into its generated logical forms (Table 9). Furthermore, models demonstrate substantially weaker performance on Group 2 sentences

Teaching assistant: the homework assignment was much harder than we thought, so I think we should give a few extra points to students who completed it. Professor: that's a terrible idea. If we give everyone a perfect score for no reason, students won't bother working hard in the future. <i>Extension Fallacy</i>	
o4-mini	<ul style="list-style-type: none"> • Generated Form: A proposes: "Because B, we should do X_1." Opponent replies: "If we do X_2 (an exaggerated version of X_1) for no reason, then C will happen; therefore, we must not do X_1." • Matched pattern: X proposes Y; response exaggerates Y to an extreme version Y. • Fallacy: Extension Fallacy
gpt-4.1-mini	<ul style="list-style-type: none"> • Generated Form: A says B was harder than expected, so A suggests giving extra points to those who completed it. B replies: If we give everyone a perfect score for no reason, then C won't work hard; therefore, we should not give extra points. • Matched Pattern: Arguing that taking a minor action will cause a chain of escalating events leading to a disastrous outcome, even though no strong evidence supports the inevitability of that chain. • Fallacy: Hasty Generalization
llama-3.3-70B	<ul style="list-style-type: none"> • Generated Form: If we do A (give everyone a perfect score for no reason), then B (students won't bother working hard in the future). • Matched Pattern: Presenting a coincidental or correlated event as proof of causation, often ignoring other plausible causal factors or explanations. • Fallacy: False Cause
gemma-3-27b-it	<ul style="list-style-type: none"> • Generated Form: If we do A (give extra points), then B (students won't work hard). Therefore, we should not do A. • Matched Pattern: If P then Q; Q is true; therefore, P is true. • Fallacy: Deductive Fallacy

Table 9: Comparison of outputs of four of the models evaluated in MULTISTEP configuration on LOGIC.

I have no intention of stopping the use of somatostatin on patients suffering from acute pancreatitis. I consider it to be a very reasonable choice. After all, it has been standard practice in our department for many years and we've been quite satisfied with the results we've had. <i>Irrelevant Authority</i>	
Top 1	Because many people [do/ believe/support] X, X must be true/good/right/best/valid. (<i>Ad Populum</i>)
Top 2	Using [personal trait, experience, past action] as implicit proof of authority on a distinct or unrelated subject. (<i>Irrelevant Authority</i>)

Table 10: Sentence accurately classified by o4-mini with 2nd ranked pattern. Although the model fails initially, the matching can still be deemed reasonable.

454 compared to Group 1, showing an average decrease
455 of 21.5% in F₁ score. Additionally, models fre-
456 quently bypass the pattern matching phase entirely,
457 arbitrarily assigning patterns despite clear misalign-
458 ment with the extracted logical form. For exam-
459 ple, given the argument *People nowadays only*
460 *vote with their emotions instead of their brains*
461 (an instance of *Hasty Generalization*), the model
462 o4-mini first extracts the logical form *All A only*
463 *do B instead of C*. The model then matches this
464 form to the pattern *Generalizing from a small sam-*
465 *ple or single event to an entire group or population*,
466 which correctly belongs to *Hasty Generalization*.
467 While this produces an accurate classification, the
468 assigned pattern does not precisely correspond to
469 the extracted logical form. In summary, while hu-
470 mans naturally decompose pattern matching into
471 multiple cognitive steps, this multi-stage process
472 proves to be challenging for current LLMs. Models
473 struggle to bridge the gap between abstract logi-
474 cal patterns and their content-dependent manifesta-
475 tions, often failing to identify the implicit premises
476 and unstated logical connections that underlie the
477 reasoning chain.

478 5 Experiments on Further Datasets

479 In order to further assess the quality of LOGIC-
480 derived patterns, we conducted a subset of the ex-
481 periments on REDDIT and ELECDEBATE using
482 the best performing model, o4-mini. We tested

patterns extracted from LOGIC, restricted to the
483 two datasets' classes (first eight rows in Table 11)
484 and patterns extracted from the datasets themselves
485 (latest two rows in Table 11).
486

	REDDIT		ELECDEBATE	
	Accuracy	Macro F ₁	Accuracy	Macro F ₁
ZERO-SHOT	82.8	82.8	67.3	50.8
DEF	82.6	82.5	65.9	54.7
LOGICAL FORMS	84.7	84.3	70.7	59.5
PATTERNS	84.7	84.5	65.5	56.3
PATTERN MATCHING	80.9	80.8	65.5	56.7
DYNAMIC ONE-SHOT	81.9	81.6	81.7	70.4
DYNAMIC + EXP	83.8	83.6	79.1	71.3
DYNAMIC + EXP + PATTERNS	79.0	78.8	78.8	72.3
SAME-DATASET PATTERNS	83.8	83.4	74.1	64.9
SAME-DATASET PATTERNS MATCHING	84.7	84.3	74.3	63.7

Table 11: Fallacy classification performance using o4-mini on REDDIT and ELECDEBATE. **PATTERNS** method involves using patterns generated on LOGIC while **SAME-DATASET PATTERNS** approach includes patterns generated on the datasets REDDIT and ELECDEBATE themselves.

487 Consistent with previous findings, logical pat-
488 tern incorporation outperforms competing ap-
489 proaches on REDDIT. Moreover, LOGIC-based
490 and REDDIT-based patterns yield comparable re-
491 sults. While taxonomy alignment prevents direct
492 comparison, results from both supervised and un-
493 supervised methods (Sahai et al., 2021; Lei and
494 Huang, 2024; Pan et al., 2024; Yeh et al., 2024)
495 are consistent with our findings (see Appendix B).
496 Only a comparison with Lei and Huang (2024)
497 (Macro F₁=81.3%) is possible: **PATTERNS** and

SAME-DATASET PATTERNS outperform their results.

Regarding ELECDEBATE, Table 11 shows that DYNAMIC ONE-SHOT yields the best performance, possibly due to the predominant presence of the class *Emotional Language* (62.5% of test set) whose detection may particularly benefit from similar worded examples. Indeed, SAME-DATASET PATTERNS achieve competitive results with respect to Goffredo et al. (2023); Pan et al. (2024) (see Appendix B). These experiments showed a fair generalization of LOGIC-derived patterns on other datasets, with the additional advantage of not requiring labeled data to re-extract the patterns.

In order to prove the broader applicability of our approach beyond LOGIC-specific patterns, we tested patterns generated from the other two datasets on LOGIC. Table 12 demonstrates solid results, validating findings on LOGIC and proving the robustness and transferability of our pattern-based methodology. Notice that accuracy is not directly comparable with values in Table 5 since REDDIT and ELECDEBATE have a subset of the classes of LOGIC.

	LOGIC _{REDDIT}	LOGIC _{ELECDEBATE}
	Accuracy	Accuracy
SAME-DATASET PATTERNS	90.8	87.1
SAME-DATASET PATTERNS MATCHING	88.3	87.5
LOGIC PATTERNS	89.1	83.7
LOGIC PATTERNS MATCHING	90.0	87.1

Table 12: Fallacy classification performance using o4-mini on LOGIC. LOGIC_X refers to LOGIC restricted to the classes from dataset X. SAME-DATASET PATTERNS approach includes patterns generated on non-LOGIC dataset X while LOGIC PATTERNS involves using LOGIC-derived patterns restricted to the classes of dataset X.

6 Related Work

Recent advances in fallacy detection have increasingly turned to LLMs, though few studies have relied exclusively on prompting-based techniques. Several works have employed fallacy detection to probe LLMs’ logical reasoning abilities (Teo et al., 2025; Hong et al., 2024; Li et al., 2024; Xu et al., 2025). Among these, Hong et al. (2024) investigated self-verification capabilities and showed that LLMs face more challenges with structure-based (formal) fallacies with respect to content-based (informal) ones, and that fallacy definitions provide minimal improvements. Xu et al. (2025)

has shown that reasoning models have better performances with respect to non-reasoning ones for fallacy classification. Among studies relying exclusively on prompting techniques, Pan et al. (2024) designed single-round and multi-round prompting schemes for zero-shot detection, while Jeong et al. (2025) introduced contextual prompting incorporating counterarguments, explanations, and goals with confidence-based ranking, showing that explanations particularly enhance performance. Lim and Perrault (2024) assessed detection abilities on the LOGIC dataset using few-shot prompting, though their different taxonomy limits direct comparison with our work. Other research has examined the logical structure of argumentation. In particular, Robbani et al. (2024) created a set of fallacy templates for four informal fallacies, based on Walton (2008) and Reisert et al. (2018)’s work. These templates provide formal logical schemas with explicit variables and relationships. Our patterns share the goal of formalizing logical structure while also capturing the core logical error and common linguistic markers for each fallacy in natural discourse. Our experiments show that these argumentation schemes underperform our patterns on LOGIC by an average of 10.7% across all models, suggesting that our patterns more effectively bridge surface text and logical structure by combining logical representation with linguistic cues. Most notably, Jin et al. (2022) developed a structure-aware model based on Electra that distills arguments into logical forms and compares them against fallacy patterns sourced from logicallyfallacious.com.

7 Conclusions

Fallacy detection is an important and complex task to solve. We developed an experimental framework that extracts logical patterns from fallacious arguments and their explanations.

We showed that incorporating these patterns significantly enhances classification performance demonstrating the effectiveness of structural approaches for this task. Specifically, pattern-based classification achieves 73.5% accuracy on LOGIC, significantly outperforming prior unsupervised approaches, and 74.2% including one-shot examples. Notably, reasoning models demonstrate consistently superior performance. Moreover, experiments on additional datasets show the generalizability of the patterns extracted.

8 Limitations

While this work demonstrates the efficacy of large language models in detecting logical fallacies by exploiting the underlying logical structure of sentences, it has several limitations. First, we intentionally generated patterns exclusively from the LOGIC dataset due to the quality and straightforward structure of its sentences. We are aware, however, that it does not fully cover the complex and multi-faceted spectrum of fallacies. Furthermore, our work is based on a small sample of LLMs. Nevertheless, we selected a diverse and representative subset, including models from different providers, with varying sizes and reasoning capabilities.

9 Ethics Statement

Logical fallacies can reinforce societal bias and facilitate the spread of misinformation, leading to harmful consequences for society. This work focuses on leveraging LLMs for detecting logical fallacies in argumentation and should not be employed to manipulate discourse by exploiting identified reasoning patterns. Furthermore, this approach risks amplifying existing LLM biases, potentially causing unfair detection. We acknowledge these limitations and encourage future bias mitigation research. We are aware of the environmental impact of large-scale LLMs usage. However, this study exclusively employs inference-only methods, significantly reducing computational requirements compared to training approaches. All datasets are used in accordance with their license and they have been checked for personally identifying and offensive content.

References

Tariq Alhindi, Smaranda Muresan, and Preslav Nakov. 2024. [Large language models are few-shot training example generators: A case study in fallacy recognition](#).

John B. Bacon, Michael Detlefsen, and David Charles McCarty. 1999. *Logic from A to Z: The Routledge Encyclopedia of Philosophy Glossary of Logical and Mathematical Terms*. Routledge.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen,

Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#).

Irving Marmer Copi, Carl Cohen, and Kenneth McMahon. 1953. *Introduction to Logic*. Macmillan, New York, NY, USA.

Giovanni Da San Martino, Seunghak Yu, Alberto Barrón-Cedeño, 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.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [Bert: Pre-training of deep bidirectional transformers for language understanding](#).

H.J. Gensler. 2010. *The to Z of Logic*. Number v. 169 in G - Reference, Information and Interdisciplinary Subjects Series. Bloomsbury Academic.

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.

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.

C.L. Hamblin. 1970. *Fallacies*. University paperbacks. Methuen.

William L. Hamilton, Rex Ying, and Jure Leskovec. 2018. [Representation learning on graphs: Methods and applications](#).

Ruixin Hong, Hongming Zhang, Xinyu Pang, Dong Yu, and Changshui Zhang. 2024. [A closer look at the self-verification abilities of large language models in logical reasoning](#).

Jiwon Jeong, Hyeju Jang, and Hogun Park. 2025. [Large language models are better logical fallacy reasoners with counterargument, explanation, and goal-aware prompt formulation](#).

deepseek-r1. In all other experiments, the standard configuration was kept. Multiple prompt configurations were evaluated for each approach.

B Baselines

B.1 Logic

We compare our results against the most recent study that employs unsupervised approaches on LOGIC:

- [Pan et al. \(2024\)](#): various prompting strategies are employed, such as requesting fallacy analysis, premises and conclusions extraction, summarization and chain-of-thought reasoning in both single- and multi-round configurations. The reported score was achieved with GPT-4.
- [Jeong et al. \(2025\)](#): the method provides implicit contextual information such as counterarguments, explanations and goals in the prompt and rank such queries based on confidence score. The reported results were achieved with GPT-4.

B.2 Other datasets

We consider only the classes of REDDIT and ELECDEBATE that belong to LOGIC. For this reason, direct comparison with prior work is generally not possible. However, for REDDIT, [Lei and Huang \(2024\)](#) provide classwise F_1 scores, allowing us to compute Macro F_1 and compare our results. Tables 13 and 14 present the comparison with prior work for both datasets.

Method	Macro F_1
<i>Supervised</i>	
Sahai et al. (2021)	58.4
Lei and Huang (2024)[†]	81.3
Pan et al. (2024)	83.2
<i>Unsupervised</i>	
Pan et al. (2024)	81.1
Yeh et al. (2024)	81.0
<i>Ours</i>	
PATTERNS	84.5
SAME-DATASET PATTERN MATCHING	84.3

Table 13: Performance comparison on REDDIT.[†] indicates that Macro F_1 is computed on the exact same classes as LOGIC.

C Additional Experiments

We are going to report some other experimental setups that have been explored, including some basic baselines that we have not included in Section 4.

Method	Macro F_1
<i>Supervised</i>	
Goffredo et al. (2023)	73.9
Pan et al. (2024)	62.3
<i>Unsupervised</i>	
Pan et al. (2024)	44.5
<i>Ours</i>	
SAME-DATASET PATTERN	64.9
DYNAMIC ONE-SHOT	70.4

Table 14: Performance comparison on ELECDEBATE.

C.1 Prompt design

- **EXP**: to investigate whether explicit reasoning improves performance, we implemented a baseline that not only provides fallacy names but also requests the model to generate a two-sentence explanation for its classification decision. This approach tests whether forcing the model to articulate its reasoning leads to better outcomes and assess the model’s actual comprehension of fallacious arguments. The two-sentence constraint was intentionally designed to keep explanations concise and manageable for manual inspection of explanations.
- **GUIDELINES**: to leverage the model’s classification errors for improvement, we develop guidelines derived from observed mistakes. We conduct pattern matching evaluation on the validation set, then systematically collect misclassified instances. For each class, we provide the model with its incorrectly classified examples and prompt it to generate comprehensive detection guidelines, given our generated pattern as reference. These guidelines include a core definition, key identifying features, common confusion patterns with similar fallacies, and a practical checklist to aid in detecting the specific type of fallacious reasoning, as can be seen from table 15. These guidelines are then adopted to evaluate the test set. Notably, while all guidelines were generated from misclassified patterns, only those produced by o4-mini and partially by gpt-4.1-mini incorporate a little structural and logical information such as common connectors or logical forms. The majority of guideline content across models focuses primarily on semantic characteristics rather than structural patterns.

Fallacy	Irrelevant Authority
Core definition	A fallacy that treats an individual’s status, title, or popularity as proof of a claim when their expertise or relevance to the topic is absent or insufficient.
Key indicators	Argument rests on “X says so” without independent support. Authority cited has no recognized expertise in the claim’s domain. No substantive evidence beyond the authority’s endorsement.
Typical confusion patterns	Ad Populum: group popularity vs. single authority endorsement. Appeal to Tradition: ‘has always been done by experts’ vs. citing irrelevant experts. Equivocation: shifting word senses vs. relying on irrelevant credentials.

Table 15: Guidelines relative to *Irrelevant Authority* fallacy generated by o4-mini.

Method	o4-mini		gpt-4o		deepseek-r1		gpt-4.1-mini		llama-3.3-70B		gemma-3-27b-it	
	Acc.	F ₁	Acc.	F ₁	Acc.	F ₁	Acc.	F ₁	Acc.	F ₁	Acc.	F ₁
EXP	61.3	61.5	60.4	53.3	61.8	54.8	57.5	57.9	56.1	56.5	59.1	60.8
GUIDELINES	65.5	65.7	62.6	56.6	63.5	57.9	60.5	60.6	52.8	53.3	58.8	59.4

Table 16: Logical fallacy classification performance on additional experiments. F₁ means Macro F₁.

C.2 Results

EXP’s (Table 16) results show that requesting the model to articulate the reasoning does not really cause any improvement. Specifically, certain classes such as *Intentional Fallacy* and *Extension Fallacy* exhibit extremely low F₁ scores under the non-reasoning models (0.027 and 0.13 respectively on average), indicating performance deterioration compared to the **ZERO-SHOT** baseline. This proves that models process surface-level semantic patterns without being able to access the multi-layered intentional structures behind reasoning (Table 17).

Including **GUIDELINES** yields only modest results. While these guidelines are designed to provide comprehensive fallacy knowledge, they appear to lack the appropriate type of information from which models can benefit. Indeed, providing explicit information about the underlying logical structure proves significantly more beneficial for model performance.

D Syntax-augmented roBERTa

Sachan et al. (2021) introduces a syntax-augmented model that incorporates dependency tree information into pre-trained BERT-based (Devlin et al., 2019) transformers through specialized Graph Neural Networks (GNNs) (Hamilton et al., 2018) that process dependency trees. The authors introduce two distinct fusion strategies to inte-

grate syntactic structure into BERT representation. We adopted specifically roBERTa-large (Liu et al., 2019) in the attempt to perform a syntax-driven examples selection. Further details about the implementation are available in Sachan et al. (2021).

E Fallacy Datasets

E.1 Logic

The dataset LOGIC (Jin et al., 2022) contains the following 13 fallacy classes: *Faulty Generalization* (*Hasty Generalization*), *Ad Hominem*, *Ad Populum*, *Circular Claim* (*Circular Reasoning*), *False Cause* (*False Causality*), *Appeal to Emotion* (*Emotional Language*), *Fallacy of Relevance* (*Red Herring*), *Deductive Fallacy*, *Intentional Fallacy*, *Fallacy of Extension* (*Extension Fallacy*), *False Dilemma* (*Black-and-White Fallacy*), *Fallacy of Credibility* (*Irrelevant Authority*) and *Equivocation*. The names in the parentheses are the actual names used in our experiments. Detailed information about the dataset is reported in Table 18.

E.2 Reddit

The dataset REDDIT (Sahai et al., 2021) contains 8 fallacy classes: *Appeal to Authority* (*Irrelevant Authority*), *Appeal to Majority* (*Ad Populum*), *Appeal to Nature*, *Appeal to Tradition*, *Appeal to Worse Problems*, *Black-and-White fallacy*, *Hasty Generalization* and *Slippery Slope*. It contains the class *No Fallacy* as well. The names in parentheses are the actual labels used. In our experiments, only the classes included in LOGIC are retained (Table 19). We can keep the class *Slippery Slope* because two generated patterns for *Hasty Generalization* correspond to it.

E.3 ElecDebate

The dataset ELECDEBATE (Goffredo et al., 2023) contains the following 6 fallacy classes: *Ad Hominem*, *Appeal to Emotion* (*Emotional Language*), *Appeal to Authority* (*Irrelevant Authority*), *Slippery Slope*, *False Cause* and *Slogan*. The names in parentheses are the actual labels used. In our experiments, only the classes included in LOGIC are retained (Table 19).

F LLM-derived Patterns and Definitions

The logical patterns presented in Table 20 were extracted following the process described in Section 3 from LOGIC. These patterns form the basis

Text	Generated explanation	Gold label
The Bible is true because God exists, and God exists because the Bible says so.	The argument uses its conclusion as a premise, claiming the Bible is true because God exists and God exists because the Bible says so. This creates a logical loop without providing independent evidence for either claim. <i>Circular Reasoning</i>	<i>Circular Reasoning</i>
My friend said that if you sneeze more than three times, you have the corona virus.	The argument assumes that sneezing more than three times directly indicates having the corona virus based on insufficient evidence. It generalizes a specific symptom without considering other possible causes or medical diagnosis. <i>Hasty Generalization</i>	<i>Irrelevant Authority</i>

Table 17: Examples from GPT-4.1-mini in EXP setting: the first sentence is correctly classified with a well-grounded explanation; the second argument is misclassified because, while its explanation appears coherent in isolation, it fails to capture the underlying fallacious reasoning.

of our evaluation of how structural features contribute to LLMs’ performance in logical fallacy detection, assessing the role of logical patterns in model reasoning capabilities. Table 21 illustrates the LLM-made fallacies’ definitions used in the experimental setting NEW DEF. Compared to LOGIC’s patterns, REDDIT patterns (Table 23) show a prevalence of linguistic markers over actual logical forms. On the other hand, ELECDEBATE-derived patterns (Table 22) emphasize a stronger logical formalization, even incorporating symbolic formalism.

G Prompts Templates

G.1 Pattern Generation

Step 1 You will be given a fallacious argument and the name of the logical fallacy it contains. Your task is to explain what is happening in the argument and why it is fallacious. Do not include definitions, labels or general commentary: focus only on describing the flaw in reasoning specific to the example in a concise way.

Step 2 You will be given a list of arguments containing a {fallacy_name} fallacy and an explanation of why it is fallacious. Your task is to provide the following information, returning a JSON object with the following fields:

```

{{
  "summary": Write a concise summary (max 2 sentences) that captures the common logical pattern behind these explanations. The summary should start with the name of the fallacy.

```

```

  "syntactic_patterns": Identify common syntactic or structural patterns in how the arguments are phrased.

```

```

  - Derive the abstract logical structure following formal logic principles. Use abstract placeholders like A, B, C to replace specific nouns or phrases, but ensure the pattern closely mirrors the logical structure and progression of the original sentence.

```

```

- Find recurring sentence structures, phrases, or ways in which the fallacious reasoning is introduced, including typical linguistic markers and illustrative cases that signal the flawed reasoning.
}}

```

```

}}
Your explanation should help someone recognize new examples of {fallacy_name} by highlighting the shared reasoning mistake across all cases.

```

G.2 Classification

```

system_prompt: You are a logical reasoning expert. Your task is to carefully examine the given argument and classify it into one of the following classes: {fallacies }.

```

- **ZERO-SHOT:** Given an argument, classify the fallacy it contains. Choose one of the following labels: {fallacies}. Respond only with the name of the fallacy, with no additional text.

```
Argument: {text}
```

```
Fallacy:
```

- **EXP:** Given an argument, classify the fallacy it contains. Choose one of the following labels: {fallacies}. Respond in the format:

```
Fallacy: [fallacy label]
```

```
Reasoning: [brief explanation justifying the choice]
```

```
The reasoning must be exactly two sentences long.
```

```
Argument: {text}
```

```
Fallacy:
```

```
Reasoning:
```

- **DEF, NEW DEF:** Given an argument, classify the fallacy it contains. Choose one of the following labels: {fallacies}. Use the following definitions to guide your classification: {definitions}. Respond only with the name of the fallacy, with no additional text.

```
Argument: {text}
```

```
Fallacy:
```

Fallacy	Definition	Logical Form
Ad Hominem	The text attacks a person instead of arguing against the claims.	Person 1 is claiming Y. Person 1 is a moron. Therefore, Y is not true.
Ad Populum	The text affirms something is true because the majority thinks so.	A lot of people believe X. Therefore, X must be true.
Black-and-White Fallacy	The text presents two alternative options as the only possibilities yet more exist.	Either X or Y is true.
False Cause	The text assumes two correlated events must also have a causal relation.	X occurred after Y. Therefore, Y caused X (although X was also a result of A,B,C,... etc)
Circular Reasoning	The text tries to prove a point by simply repeating the point in different words.	X is true because of Y. Y is true because of X.
Deductive Fallacy	The text presents a conclusion that doesn't logically follow from the premises.	If A is true, then B is true. B is true. Therefore, A is true.
Emotional Language	The text arouses non-rational emotions.	Claim X is made without evidence. In place of evidence, emotion is used to convince the interlocutor that X is true.
Equivocation	The text uses a key term in multiple senses, leading to ambiguous conclusions.	Term X is used to mean Y in the premise. Term X is used to mean Z in the conclusion.
Extension Fallacy	The text attacks an exaggerated version of the opponent's claim.	Person 1 makes claim Y. Person 2 restates person 1's claim (in a distorted way). Person 2 attacks the distorted version of the claim. Therefore, claim Y is false.
Hasty Generalization	The text draws a broad conclusion based on a limited sample.	Sample S is taken from population P. Sample S is a very small part of population P. Conclusion C is drawn from sample S and applied to population P.
Intentional Fallacy	The text relies on the author's intent instead of focusing on the meaning within the text itself.	Person 1 knows claim X is incorrect. They still claim that X is correct using an incorrect argument
Irrelevant Authority	The text cites an authority, but the authority lacks relevant expertise.	According to person 1, who is an expert on the issue of Y, Y is true. Therefore, Y is true.
Red Herring	The text introduces an irrelevant topic to divert attention from the main argument.	Argument A is presented by person 1. Person 2 introduces argument B. Argument A is abandoned.

Table 18: LOGIC class taxonomy: class names, definitions and logical form representations used in baseline experiments. Definitions are sourced from [Lei and Huang \(2024\)](#) and manually refined while logical forms are extracted from logicallyfallacious.com.

REDDIT	ELECDEBATE
• Ad Populum	• Ad Hominem
• Irrelevant Authority	• Irrelevant Authority
• Hasty Generalization	• Emotional Language
• Slippery Slope	• Slippery Slope
• Black-and-White Fallacy	• False Cause

Table 19: Fallacy classes in REDDIT and ELECDEBATE used in our experiments.

- 1017 • **LOGICAL FORMS, PATTERNS:** Given an argument, your task is to classify the type of logical fallacy it contains. Choose one of the following labels: {fallacies}. To assist you in this task, common patterns associated with each fallacy are also provided. {patterns}. Carefully compare the argument to these patterns and select the fallacy that best matches. Provide only the name of the fallacy. Argument: {text} Fallacy:
- 1018
- 1019
- 1020
- 1021
- 1022
- 1023
- 1024
- 1025
- 1026
- 1027
- 1028 • **PATTERN MATCHING:** Given an argument,

- your task is to classify the type of logical fallacy it contains. Choose one of the following labels: {fallacies}. To assist you in this task, you are provided with common reasoning patterns associated with each fallacy: {patterns} Return:
- The specific pattern that best matches the logical reasoning in the argument.
- The name of the fallacy. Don't add any additional text. Argument: {text} Fallacy:
- 1029
- 1030
- 1031
- 1032
- 1033
- 1034
- 1035
- 1036
- 1037
- 1038
- 1039
- 1040
- 1041 • **ONE-SHOT, DYNAMIC ONE-SHOT, SYNTAX-BASED DYNAMIC ONE-SHOT:** Given an argument, your task is to classify the type of logical fallacy. Choose the correct fallacy from the following list: {fallacies}. Use the following examples to guide your classification: {examples} Argument: {text} Fallacy:
- 1042
- 1043
- 1044
- 1045
- 1046
- 1047
- 1048
- 1049
- 1050 • **ONE-SHOT + EXP, DYNAMIC + EXP:** Given

1051	an argument, your task is to classify the type	known fallacy patterns provided below, return	1101
1052	of logical fallacy. Choose the correct fallacy	the one that most closely matches the pattern	1102
1053	from the following list: {fallacies}. Use the	you extracted in Step 1.	1103
1054	following examples to guide your classifica-	Step 3 (Fallacy Classification): Based on the	1104
1055	tion. Each example includes both the fallacy	extracted and matched pattern, classify the	1105
1056	and a brief explanation of why it applies: {ex-	argument into one of the logical fallacy types	1106
1057	amples}	from the list: {fallacies}	1107
1058	Argument: {text}	Choose the most appropriate category based	1108
1059	Fallacy:	on the structure of the reasoning.	1109
1060		Use the following reasoning patterns as refer-	1110
1061	• DYNAMIC + EXP + PATTERNS: Given an	ence: {patterns}	1111
1062	argument, your task is to classify the type of	Analyze the following argument: {text}	1112
1063	logical fallacy in a given text. Choose the	Step 1 (Form Extraction):	1113
1064	correct fallacy from the following list: {falla-	Step 2 (Matching Pattern):	1114
1065	cies}. Use the examples below to guide your	Step 3 (Fallacy):	1115
1066	classification: each example includes both the		
1067	fallacy and a brief explanation of why it ap-		
1068	plies. {definitions} You will be provided with		
1069	a list of common patterns associated with each		
1070	fallacy. {patterns} Carefully compare the ar-		
1071	gument to these patterns and select the fallacy		
1072	that best matches. Provide only the name of		
1073	the fallacy.		
1074	Argument: {text}		
1075	Fallacy:		
1076			
1077	• GUIDELINES: Given an argument, your		
1078	task is to classify the type of logical fallacy.		
1079	Choose one of the following: {fallacies}. To		
1080	assist you in this task, you are provided with		
1081	useful guidelines. These include typical rea-		
1082	soning patterns for each fallacy, common mis-		
1083	takes that often lead to misclassification and		
1084	a quick practical checklist for classification:		
1085	{guidelines}. Return only the name of the fal-		
1086	lacy.		
1087	Argument: {text}		
1088	Fallacy:		
1089			
1090	• MULTISTEP: Given an argument, your task		
1091	is to process it in three steps:		
1092	Step 1 (Form Extraction):		
1093	When given a sentence, extract its underlying		
1094	syntactic logical form. Use abstract placehold-		
1095	ers like A, B, C to replace specific nouns or		
1096	phrases, but ensure the pattern closely mirrors		
1097	the logical structure and progression of the		
1098	original sentence. Preserve logical connectors		
1099	and adverbs such as ‘therefore’, ‘because’,		
1100	‘if’, ‘then’, etc. The goal is to abstract away		
	from surface wording while preserving the		
	sentence’s original reasoning form.		
	Step 2 (Pattern Matching): From the list of		

Fallacy class	Patterns
Deductive Fallacy	<ul style="list-style-type: none"> The argument assumes that because X is true, Y must also be true, without establishing a necessary connection between X and Y. Because X shares a characteristic with Y, therefore Y must also have characteristic Z (unique to X). If P then Q; Q is true; therefore, P is true. All A are B; all B are C; therefore, all C are A. The argument compares X to Y as if they are equivalent, ignoring relevant differences.
Ad Hominem	<ul style="list-style-type: none"> Dismisses someone's argument by accusing the opponent of similar behavior, avoiding the argument itself. Argues that because X has characteristic Y, X's views or claims must be invalid/false. Uses a personal insult or irrelevant fact about X to discredit X without addressing the core issue. Focuses on unrelated personal factors (e.g., age, profession, habits) to attack the person instead of the argument.
Emotional Language	<ul style="list-style-type: none"> Appeals that highlight personal circumstances or potential consequences without addressing the core issue, e.g., 'I haven't done X, but... [appeal to emotion]'. Use of emotionally charged or loaded terms in place of neutral language, e.g., calling something an 'outrage', 'dangerous militants', or using phrases like 'taking our freedom away'. Rhetorical questions or statements designed to evoke feelings of guilt, sympathy, or fear, e.g., 'If we don't do X, disaster Y will happen'. Evocation of pity or sympathy to distract from the logical evaluation of claims, e.g., 'I studied during my grandmother's funeral'. Use of vivid imagery or emotionally provocative examples to bypass critical analysis, e.g., showing suffering animals or invoking dramatic suffering stories.
False Cause	<ul style="list-style-type: none"> Inferring causation from correlation expressed as 'X happened when Y happened' or 'Every time X occurs, Y follows,' without evidence of a causal link. Attributing a complex outcome to a single factor due to temporal proximity or repeated coincidence (e.g., 'Because of X, Y happened,' ignoring other influences). Using phrases implying causality based on timing, such as 'therefore', 'must have caused', 'is the reason', or 'is the cause of', without supporting evidence. Presenting a coincidental or correlated event as proof of causation, often ignoring other plausible causal factors or explanations. Statements that simplify multi-factor phenomena to a single cause (e.g., 'Because of single action/event X, complex result Y occurred').
Irrelevant Authority	<ul style="list-style-type: none"> The argument assumes that because [authority figure] says/believes [X], therefore [X] must be true/reliable/effective. Relying on the opinion or endorsement of [famous/unqualified person] outside their field of expertise to support [claim/conclusion]. Because [person/role/title] holds a position or is respected, their statement on [irrelevant topic] is presented as evidence. Using [personal trait, experience, past action] as implicit proof of authority on a distinct or unrelated subject. The argument presents [authority]'s position as the sole justification without providing independent reasons or evidence.
Extension Fallacy	<ul style="list-style-type: none"> X proposes Y; response exaggerates Y to an extreme or total version (e.g., 'So you want to [extreme claim]?') Assuming that because someone holds a position on A, they must also hold an extreme or unrelated position B (e.g., 'If you believe A, then you must believe B') Misinterpreting a specific statement as implying a much broader or more negative attitude (e.g., 'Because you said X, you must hate Y') Responding to a moderate or specific claim by substituting a more extreme or absurd claim that is easier to criticize (e.g., 'You think that... [absurd extension]') Framing a preference or partial stance as a wholesale endorsement or rejection of a related but distinct issue (e.g., 'Preferring X means you hate Y')
Hasty Generalization	<ul style="list-style-type: none"> Because a particular instance or individual showed Z, therefore all instances or individuals must be Z. Arguing that taking a minor action will cause a chain of escalating events leading to a disastrous outcome, even though no strong evidence supports the inevitability of that chain. If we allow [event A], then [event B] will happen, then [event C], and eventually [event Z] will occur—so we must not allow [event A]. Generalizing from a small sample or single event to an entire group or population. Jumping from some instances to 'everyone' or 'all' without acknowledging exceptions or diversity.
Equivocation	<ul style="list-style-type: none"> The reasoning equivocates by shifting from a literal to a figurative meaning (or vice versa) of a term to create a false equivalence. One premise employs TERM with one definition/context, while another premise or conclusion uses the same TERM but with a distinctly different meaning, unacknowledged by the arguer.
Intentional Fallacy	<ul style="list-style-type: none"> The argument assumes that because X (e.g., someone's intention, belief, or lack of counter-evidence), therefore Y is true. Asserting P is true because it has not been disproven. Because the creator intended [interpretation], the work should be understood as [interpretation]. Questions framed to presuppose guilt or a specific intention (e.g., 'Have you stopped X?'), thus assuming what is to be proven. If A does not have trait X, and X is (allegedly) typical of group G, then A is not a member of G.
Ad Populum	<ul style="list-style-type: none"> Because many people [do/ believe/support] X, X must be true/good/right/best/valid. Everyone/Most people [do/believe] X, so you/one should do X too. The popularity of X is used as evidence for X's quality, validity, or truth, rather than providing objective reasons. Appealing to the desire to belong to a group, suggesting that conformity implies correctness or value. Using phrases like 'everyone knows,' 'the majority thinks,' 'most people do,' to justify a conclusion without addressing actual evidence.
Red Herring	<ul style="list-style-type: none"> Instead of addressing [original issue], the argument shifts focus to [irrelevant topic], which distracts from the main discussion. The argument attempts to justify/explain/defend by referencing [irrelevant detail], ignoring the original issue of [main topic]. A shift from the initial question or problem to a secondary topic that does not logically follow, e.g., 'You asked about X, but I will tell you about Y.'
Black-and-White Fallacy	<ul style="list-style-type: none"> Either [option A] or [option B], with no other alternatives considered. You are either [extreme position A] or [extreme position B]. You are either with me or against me. [Action A] or else [negative consequence], ignoring intermediary options.
Circular Reasoning	<ul style="list-style-type: none"> X is true/better/good because X is true/better/good. X is Y because X has property Y (where property Y is essentially restating X). "Because" + restatement or synonymous phrasing of the claim as the reason. The argument claims/assumes X to prove/justify X.

Table 20: List of logical patterns extracted in our LLM-based experimental framework from LOGIC dataset.

Fallacy	LLMs-derived definition
Ad Hominem	Ad Hominem occurs when an argument targets a person’s character or traits instead of engaging with the actual issue or evidence presented.
Ad Populum	Ad Populum occurs when a claim is deemed true or good simply because many people believe or endorse it, without examining the actual reasoning.
Black-and-White Fallacy	Black-and-White Fallacy occurs when only two extreme options are presented, ignoring the existence of middle ground or alternative solutions.
False Cause	False Cause occurs when a causal relationship is assumed based on correlation alone, without sufficient evidence or consideration of other factors.
Circular Reasoning	Circular Reasoning occurs when the conclusion is assumed in the premises, creating a loop that provides no independent support for the argument.
Deductive Fallacy	Deductive Fallacy occurs when conclusions do not logically follow from the premises, often due to assuming unsupported relationships, oversimplifying, misapplying analogies, or improperly reversing conditions.
Emotional Language	Emotional Language occurs when persuasion relies on appeals to emotion rather than logical reasoning or factual evidence.
Equivocation	Equivocation occurs when a key term is used ambiguously in an argument, shifting meaning and creating an illusion of logical connection.
Extension Fallacy	Extension Fallacy occurs when an argument exaggerates or distorts an opponent’s claim to make it easier to attack, rather than addressing the actual position.
Hasty Generalization	Hasty Generalization occurs when a broad conclusion is drawn from an insufficient or unrepresentative sample of evidence.
Intentional Fallacy	Intentional Fallacy occurs when arguments are judged based on the speaker’s intentions or characteristics rather than the content and evidence or when asserting that something is true only because it has not been disproven.
Irrelevant Authority	Irrelevant Authority occurs when an argument cites an authority whose expertise is not relevant to the subject matter being discussed.
Red Herring	Red Herring occurs when attention is diverted from the main issue by introducing irrelevant or emotionally charged distractions.

Table 21: List of definitions extracted in our LLM-based experimental framework from LOGIC dataset.

Fallacy class	Patterns
Irrelevant Authority	<ul style="list-style-type: none"> A (an authority or respected figure) asserts proposition $B \rightarrow \therefore B$, even if A is not qualified Group C of authorities or experts agree on $B \rightarrow \therefore B$, treating consensus as proof “According to A, B.”, claim is accepted without evidence” “A said that B.” I have the endorsement/support of A who agrees that B., using authority as evidence Every A1, every A2, believes that B., implying that popularity among authorities is proof
Emotional Language	<ul style="list-style-type: none"> Emotional description of A \rightarrow conclude B without substantive evidence Invoke sentiment about A (fear, pride, pity) \rightarrow demand action B, regardless of facts Attach emotive qualifier C to proposition A. \therefore endorse conclusion B as if emotion proves truth A is C. (e.g., “This is dangerous.”, “This is a disgrace.”) \rightarrow therefore reject or accept it without reasons I (strongly/emphatically) B \Rightarrow the strength of feeling is treated as support for B.
False Cause	<ul style="list-style-type: none"> A occurs before B. \therefore assume $A \rightarrow B$, mistaking correlation for causation without evidence of a causal link (A happens \wedge B follows) \rightarrow infer A causes B, without considering alternatives Since A, B has occurred, assuming causation After A happened, B followed, so A must have caused B A led to B simply because it preceded/coincided with B (confuses sequence with cause)
Ad Hominem	<ul style="list-style-type: none"> A asserts B; instead of addressing B, respond with C (an irrelevant trait or action of A) \rightarrow dismiss B A makes claim B; cite A’s character/property C. \therefore reject B without counterargument”, A is a C. therefore A’s argument B is invalid Everybody knows A is C, therefore disregard A’s argument A did X in the past, so B’s point is invalid Look, A C, therefore B is wrong
Slippery Slope	<ul style="list-style-type: none"> A occurs \rightarrow assume B inevitably follows \rightarrow C inevitably follows If A then B; if B then C; therefore if A then C If A, then B will happen [and C thereafter] It will A, it will B, and it will C Unless A happens, B will occur and C will inevitably follow When A occurs, $B \rightarrow C$ (disaster)

Table 22: List of logical patterns extracted in our LLM-based experimental framework from ELECDEBATE dataset.

Fallacy class	Patterns
Irrelevant Authority	<ul style="list-style-type: none"> 1. A possesses status or expertise in field B. 2. A makes claim C (often outside or without direct evidence in B). 3. Therefore, C must be accepted as true or more valid than competing claims, irrespective of its own merits. Listen to the experts/scientists/doctors rather than try to become one yourself Everyone I've spoken to (professionals, hunters, farmers, etc.) says... Only the views of X matter All the <group of experts> agree, while dismissing others' arguments without engaging with evidence.
Hasty Generalization	<ul style="list-style-type: none"> Premise: A limited or anecdotal subset of A (e.g., one's personal experiences or a few observed cases) exhibits property P. (Often no representative sampling or statistical support is offered.) Conclusion: Therefore, all A (or the phenomenon as a whole) must exhibit property P (or consequence Q). I (never) have(d) [experience] with A, so A must (not) be [P] I know a few A who X, therefore all A X. Because my sample of A were fine, A must be harmless/beneficial. Based on X anecdote, everyone must (feel/believe/do) Y. Quick universal leaps: 'Since this happened once, it always/never happens. Dismissive questions or commands: 'Why would you...? You're all/lemmings/idiots.' Percent or frequency claims from anecdotes: 'Half of my school... 99% of people... every time.'
Ad Populum	<ul style="list-style-type: none"> Premise: Group G (e.g. 'everyone', 'most people', 'X% of Y') believes or endorses proposition A. Unstated Principle: Widespread acceptance implies correctness (or falsity). 3. Conclusion: Therefore, A must be valid/true (or invalid/false). I think everyone agrees that A Most/virtually all people believe A Numerical appeals like '51% of X' or 'thousands of scientists, Appeals to anonymous masses ('the rest of the world,' 'every store'), expert-by-number appeals ('aggregate wisdom,' 'hundreds of people'), popularity voting metaphors ('let upvotes decide,' 'fan vote'), rhetorical questions and challenges ('It can't be the entire world wrong, can it?'), and dismissal of dissenters as lone or isolated ('you're the only one,' 'no one else').
Black-and-White Fallacy	<ul style="list-style-type: none"> Premise: Issue A is said to admit only two exclusive states B or C ($\neg B \Rightarrow C$; $\neg C \Rightarrow B$) and no other alternative D is acknowledged. Conclusion: One must choose either B or C, with every nuance or middle ground dismissed. Typical markers include phrases like 'Either A or B,' 'only option,' 'no middle ground,' 'you can't have it both ways,' 'pick one,' 'unless you take sides,' 'you must choose X or Y,' and 'the only alternative is...': language that signals a forced binary choice and excludes other possibilities.
Slippery Slope	<ul style="list-style-type: none"> $A \Rightarrow B \Rightarrow C \Rightarrow \dots \Rightarrow Z$, therefore $A \Rightarrow Z$, where each transition ($A \rightarrow B$, $B \rightarrow C$, etc.) is asserted as unavoidable or necessary without justification. Conditional constructs ('If we allow A, next will be B'), rhetorical questions ('Where do we stop?', 'What's to stop...?'), enumerative progressions ('first X, then Y, then Z'), warnings of 'dangerous precedents' or 'slippery slope,' and speculative modal verbs or adverbs ('could lead to', 'inevitably', 'one step closer').

Table 23: List of logical patterns extracted in our LLM-based experimental framework from REDDIT dataset.