

# ZoFia: Zero-Shot Fake News Detection with Entity-Guided Retrieval and Multi-LLM Interaction

Anonymous ACL submission

## Abstract

The rapid spread of fake news threatens social stability and public trust, highlighting the urgent need for its effective detection. Although large language models (LLMs) show potential in fake news detection, they are limited by knowledge cutoff and easily generate factual hallucinations when handling time-sensitive news. Furthermore, the thinking of a single LLM easily falls into early stance locking and confirmation bias, making it hard to handle both content reasoning and fact checking simultaneously. To address these challenges, we propose ZoFia, a two-stage zero-shot fake news detection framework. In the first retrieval stage, we propose novel Hierarchical Saliency and Saliency-Calibrated Minimum Marginal Relevance (SC-MMR) algorithm to extract core entities accurately, which drive dual-source retrieval to overcome knowledge and evidence gaps. In the subsequent stage, a multi-agent system conducts multi-perspective reasoning and verification in parallel and achieve an explainable and robust result via adversarial debate. Comprehensive experiments on two public datasets show that ZoFia outperforms existing zero-shot baselines and even most few-shot methods. Our code will be open-sourced to facilitate the related community <sup>1</sup>.

## 1 Introduction

The rapid spread of fake news through social networks has become prevalent, posing severe threats to key domains such as politics (Fisher et al., 2016), economy (Bakir and McStay, 2018), and livelihood (Zhou and Zafarani, 2020). The swift advancement of generative models further exacerbates this concern (Chen and Shu, 2023; Nan et al., 2024; Li et al., 2025, 2024a). In this context, developing effective, efficient and interpretable fake news detection methods is indispensable.

Early studies mainly rely on supervised learning (Monti et al., 2019; Kaliyar et al., 2021) to train detection models, but their dependence on large-scale labeled data makes it hard to adapt to emerging news topics (Hoy and Koulouri, 2022). Large language models (LLMs), with broad pre-training knowledge and strong contextual understanding (Su et al., 2023), significantly advance few-shot and zero-shot methods, providing new opportunities to overcome this bottleneck. In few-shot methods, LLMs either serve as auxiliary tools to perform data augmentation for downstream classifiers (Hu et al., 2024a), or act directly as detectors through prompt learning (Jiang et al., 2022) and instruction tuning (Pavlyshenko, 2023), but they still do not fully escape reliance on training data. By contrast, zero-shot methods that do not require labeled samples are a more promising paradigm, which mainly guide LLM reasoning by context engineering (Zhang and Gao, 2023) and agentic architectures (Li et al., 2024b).

However, zero-shot methods face reliability challenges. On the one hand, the internal knowledge of LLMs is limited by knowledge cutoff (Cheng et al., 2024). Without the external information, they easily produce factual hallucinations (Ji et al., 2023) when handling dynamic news events. On the other hand, LLMs tend to lock their stance early (Echterhoff et al., 2024), where confirmation bias (Wan et al., 2025) drives subsequent reasoning to merely rationalize the initial judgment. (Wang et al., 2025) Consequently, when external retrieval is introduced into a single LLM, it easily takes a cognitive shortcut that substitutes whether information can be retrieved for factual veracity, which seriously weakens content reasoning on the original news. This effect becomes much stronger when the retrieved evidence is irrelevant or contradictory. (Tan et al., 2024) Therefore, we observe that the inherent reasoning defects of a single LLM prevent it from performing both news content reasoning and

<sup>1</sup><https://anonymous.4open.science/r/ZoFia-4534/>

external evidence verification well simultaneously.

We argue that content reasoning and fact checking should be decoupled and reliable judgment can be achieved through multi-agent interaction. Based on this motivation, we propose ZoFia, a two-stage zero-shot framework for fake news detection. In the first entity-guided retrieval stage, we introduce Hierarchical Saliency to address semantic dilution (Hou et al., 2021) of news entities, which fully leverages global and local semantics to score entity saliency, and design a novel Saliency-Calibrated Minimum Marginal Relevance (SC-MMR) algorithm to accurately extract core entities. These entities then drive dual-source retrieval from Wikipedia and Open Web to mitigate knowledge cutoff and effectively suppresses hallucinations. In the subsequent multi-LLM interaction stage, ZoFia assigns agents to perform content reasoning and fact checking in parallel and breaks stance locking of a single LLM via adversarial debate, ensuring the robustness and interpretability of the final verdict.

To sum up, our main contributions are outlined as follows:

- We propose ZoFia, a retrieval-augmented multi-agent zero-shot framework for fake news detection that effectively overcomes the inherent reasoning flaws of a single LLM.
- We introduce a novel granularity-aware metric Hierarchical Saliency and design SC-MMR algorithm to extract news entities accurately for efficient retrieval.
- Comprehensive experiments on two public datasets demonstrate that ZoFia outperforms existing zero-shot baselines and even most few-shot methods.
- Our code will be open-sourced to facilitate related communities.

## 2 Related Work

**LLM in Fake News Detection.** The application of Large Language Models (LLMs) has become a research frontier in fake news detection, primarily divided into few-shot and zero-shot paradigms. In few-shot methods, researchers often use LLMs as auxiliary tools for data augmentation (Hu et al., 2024a; Nan et al., 2024), or as detectors via post-training (Jiang et al., 2022; Pavlyshenko, 2023). However, these methods fail to completely eliminate dependence on labeled data. Zero-shot methods directly guide model reasoning through context engineering (Zhang and Gao, 2023) and agentic

architectures (Li et al., 2024b; Liu et al., 2024), making judgments without labeled samples. Nevertheless, these methods fail to address the confirmation bias arising from a single reasoning chain. We decouple retrieval from reasoning via a multi-role multi-agent system and finally aggregate all evidences to judge, which achieves more comprehensive and robust zero-shot discrimination.

**Multi-Agent System.** Multi-Agent Systems (MAS) have emerged as an effective paradigm to enhance LLMs for complex tasks. The pioneering work Chateval (Chan et al., 2023) demonstrates that MAS improves both the robustness and accuracy of generation tasks. Subsequent studies introduce this paradigm to reasoning tasks. For instance, COLA (Lan et al., 2024) designs a collaboration framework for stance detection, but it remains limited to analyzing the original text. TruEDebate (Liu et al., 2025) applies structured debate to fake news detection, but its implementation tends to cause premature consensus convergence. Our ZoFia introduces external information retrieval and independent modular analysis. This design aims to fundamentally mitigate knowledge cutoff and ensure the diversity and independence of arguments and analyses.

## 3 Stage 1: Entity-Guided Retrieval

This stage aims to acquire reliable external knowledge and instant factual evidences for the subsequent multi-agent system. It consists of four sequential modules. The first three modules precisely extract a set of core entities from the original news. These entities serve as keywords and are concatenated into a query for the final retrieval module, which retrieves from Open Web and Wikipedia.

### 3.1 Entity Extractor

This module first uses a pre-trained BERT-NER (Tjong Kim Sang and De Meulder, 2003) model to perform named entity recognition (NER) on the news text. This process can be expressed as:

$$\{(t_i, e_i, c_i)\}_{i=1}^N = \mathcal{M}_{\text{BERT-NER}}(T), \quad (1)$$

where  $\mathcal{M}$  is the pre-trained model,  $T$  is the input news text.  $(t_i, e_i, c_i)$  denotes the recognized entity triplet,  $t_i$  is the entity token,  $e_i$  is the entity label, and  $c_i$  is the confidence score for the corresponding label, expressed as the conditional probability  $c_i = P(e_i|t_i, T; \mathcal{M}_{\text{BERT-NER}})$ .

Due to the large number of recognized entities, we aggregate consecutive entity tokens to form

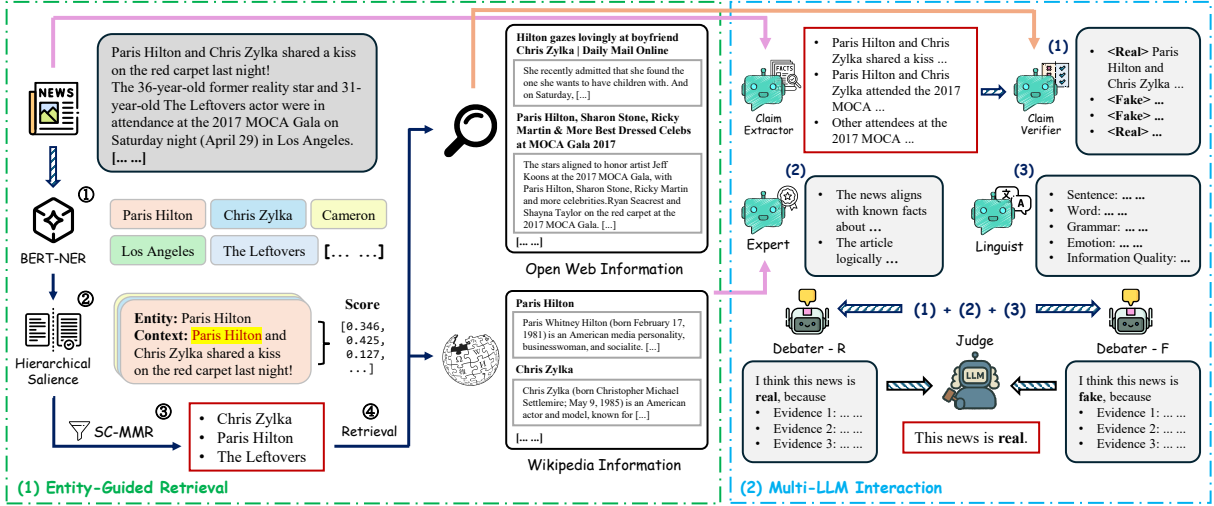


Figure 1: Overall architecture of our proposed ZoFia framework.

new entity units  $U_e$ . Its confidence score  $c(U_e)$  is calculated by averaging the confidence scores of all tokens in  $U_e$ , expressed as:

$$c(U_e) = \frac{1}{|U_e|} \sum_{t_i \in U_e} c_i. \quad (2)$$

We establish a dynamic confidence threshold to select entities with high confidence scores. Starting from an initial confidence score  $\lambda_{\text{init}}$ , if the number of selected entities fails to meet the predetermined minimum  $n_{\text{min}}$ , the algorithm iteratively lowers the confidence score by  $\Delta\lambda$  and repeats the selection, until at least  $n_{\text{min}}$  entities are obtained.

### 3.2 Salience Scorer

This module aims to accurately quantify the importance of news entities, namely Entity Salience (ES) (Dunietz and Gillick, 2014). Although the prior study (Bullough et al., 2024) has demonstrated that bi-encoder architectures can efficiently estimate entity salience, their performance is often limited by semantic dilution (Hou et al., 2021) problem, which leading that the importance of key entities are severely underestimated.

We propose a novel **Hierarchical Salience** that avoids a single, coarse evaluation between an entity and the whole text. It decomposes an entity’s overall importance into two orthogonal and multiplicative components: Local Salience  $\mathcal{S}_{\text{local}}$  and Global Salience  $\mathcal{S}_{\text{global}}$ . Local Salience measures the semantic alignment between the entity and its immediate context, and Global Salience measures how much the local context contributes to the text’s main content.

Formally, for a news text  $T$  that is an ordered sequence of sentences, consider any entity  $U_i$  that appears in sentence  $s_j$ . We define its local context as  $\mathcal{C}(U_i) = s_{j-1} \oplus s_j \oplus s_{j+1}$ . Aligned with (Bullough et al., 2024), we use a unified Sentence-BERT (SBERT) (Reimers and Gurevych, 2019) encoder  $\mathcal{M}_{\text{SBERT}}(\cdot)$  to embed the entity  $U_i$ , its local context  $\mathcal{C}(U_i)$ , and the full text  $T$  to vectors  $\mathbf{v}_{U_i}$ ,  $\mathbf{v}_{\mathcal{C}(U_i)}$ , and  $\mathbf{v}_T$ . The hierarchical salience  $\mathcal{S}_{\text{hier}}(U_i)$  is derived from the product of these two components as

$$\mathcal{S}_{\text{hier}}(U_i) = \underbrace{\frac{\mathbf{v}_{U_i} \cdot \mathbf{v}_{\mathcal{C}(U_i)}}{\|\mathbf{v}_{U_i}\| \cdot \|\mathbf{v}_{\mathcal{C}(U_i)}\|}}_{\mathcal{S}_{\text{local}}(U_i|\mathcal{C}(U_i))}} \cdot \underbrace{\frac{\mathbf{v}_{\mathcal{C}(U_i)} \cdot \mathbf{v}_T}{\|\mathbf{v}_{\mathcal{C}(U_i)}\| \cdot \|\mathbf{v}_T\|}}_{\mathcal{S}_{\text{global}}(\mathcal{C}(U_i)|T)}. \quad (3)$$

Hierarchical Salience provides a finer and more robust estimate of each entity’s importance to the news content. It serves as the key criterion for entity filtering in subsequent modules.

### 3.3 Keyword Selector

This module aims to select an informative subset of keywords from the candidate entity set  $\mathcal{U}_{\text{selected}}$  to mitigate the query drift (Carpineto and Romano, 2012) problem. However, this process faces two practical challenges. First, The hierarchical salience score  $\mathcal{S}_{\text{hier}}(U_i)$  is highly sensitive to context granularity, so a fixed screening threshold is not effective. Second, coreference in news introduces high-scoring entities with repeated semantics, which harms the diversity of the keyword set.

To optimize relevance and diversity under these constraints, We propose an improved MMR (Carbonell and Goldstein, 1998) algorithm, namely **Salience-Calibrated MMR (SC-MMR)**. At the

$k$ -th iteration, SC-MMR evaluates the score of a candidate entity  $U_i$  by

$$\text{MMR}(U_i) = \lambda_k \mathcal{S}_{\text{hier}}(U_i) - (1 - \lambda_k) \max_{U_j \in \mathcal{U}_{\text{selected}}} \mathcal{S}(U_i, U_j), \quad (4a)$$

$$\mathcal{S}(U_i, U_j) = \frac{\mathbf{v}_{U_i} \cdot \mathbf{v}_{U_j}}{\|\mathbf{v}_{U_i}\| \|\mathbf{v}_{U_j}\|}. \quad (4b)$$

The key innovation of SC-MMR is to introduce a weight schedule  $\lambda_k$  that changes with the number of selected keywords  $k$ , so that the focus gradually shifts from relevance to diversity. We adopt an annealing schedule with a lower bound:

$$\lambda_k = \max(\lambda_{\min}, \lambda_{\max} - \exp(\alpha \cdot k - \beta)). \quad (5)$$

This form ensures that  $\lambda_k$  decreases monotonically with  $k$  while retaining a non-zero salience weight via  $\lambda_{\min}$  to prevent the diversity term from fully dominating.

We further introduce a dynamic termination rule based on the relative change of the MMR score. The iteration continues only when the next best candidate  $\mathcal{U}_{k+1}^*$  satisfies  $\text{MMR}(\mathcal{U}_{k+1}^*) > \gamma \cdot \text{MMR}(\mathcal{U}_k^*)$ , where  $\gamma$  is a decay factor. This design prevents the decline in entity quality caused by diminishing marginal utility.

### 3.4 Information Retrieval

This module uses the keyword set  $\mathcal{K} = \{k_1, k_2, \dots, k_N\}$  distilled in previous modules to build a comprehensive external knowledge base  $\mathcal{E}$ , serving as supplementary context for the subsequent detection stage. We implement a dual source retrieval that gathers information from the open web and Wikipedia in parallel, ensuring both timeliness and authority.

**Retrieval from Open Web.** We compose all keywords into an aggregated query for the open web  $Q_{\text{web}}$  with the following logic:

$$Q_{\text{web}} = \left( \bigwedge_{i=1}^N k_i \right) \quad (6)$$

The operator  $\wedge$  requires that results relate to all keywords, which enables strict cross-keyword verification. We explicitly exclude news and wikipedia sources. The former avoids retrieving the duplicate reports to prevent source contamination (Deng et al., 2023), and the latter prevents redundancy with the subsequent Wikipedia retrieval. For each returned page, we only extract its summary and search snippet as the raw corpus.

**Retrieval from Wikipedia.** The summary section of a Wikipedia entry usually provides the most precise definition for an entity. However, a single keyword  $k_i$  often corresponds to multiple Wikipedia entries, which introduces semantic ambiguity. We build a **context-aware disambiguation** mechanism that uses the original local context  $\mathcal{C}(U_i)$  of keyword  $k_i$  in the news text to perform accurate matching.

When Wikipedia returns a list with  $M$  candidate senses  $\mathcal{O}(k_i) = \{o_1, o_2, \dots, o_M\}$ , the mechanism examines each candidate  $o_m$ . It constructs a temporary modified context  $\mathcal{C}(U_i \leftarrow o_m)$  by replacing the original entity  $U_i$  with the description text of  $o_m$ . The optimal sense  $o_i^*$  of  $U_i$  is defined as the option that maximizes the cosine similarity between the vector of the original context and the vector of the modified context:

$$o_i^* = \arg \max_{o_m \in \mathcal{O}(k_i)} \frac{\mathbf{v}_{\mathcal{C}(U_i)} \cdot \mathbf{v}_{\mathcal{C}(U_i \leftarrow o_m)}}{\|\mathbf{v}_{\mathcal{C}(U_i)}\| \|\mathbf{v}_{\mathcal{C}(U_i \leftarrow o_m)}\|}. \quad (7)$$

where  $\mathbf{v}$  denotes the context vectors embedded by the pretrained SBERT  $\mathcal{M}_{\text{SBERT}}$ . After identifying the unique entry, we extract the first 3 sentences of its summary as supplementary material.

## 4 Stage 2: Multi-LLM Interaction

The external information provided by entity-guided retrieval and the prior knowledge of LLMs builds a **Multi-Source Information Matrix**. This stage employs a dual-state multi-LLM system that fully exploits this matrix to perform parallel and multi-perspective content reasoning and claim verification, which are finally aggregated to reach a robust and interpretable judgment.

This stage operates in two orthogonal interaction states. **LLM Collaboration** performs parallel analyses across multiple agents to reduce inferential variance. **LLM Debate and Judgment** introduces an adversarial debate to reduce systemic bias. These two states form a complete reasoning chain from mining divergent evidences to making a convergent judgment.

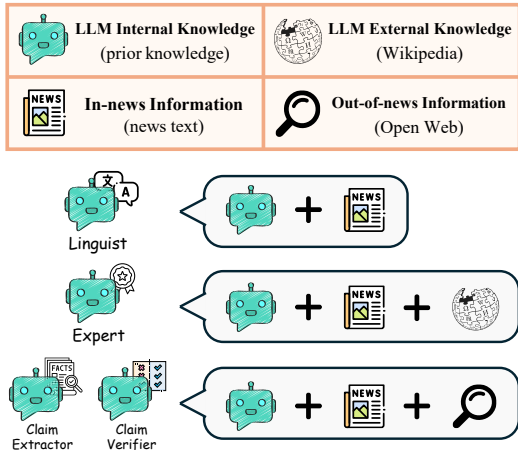
### 4.1 LLM Collaboration

LLM Collaboration state aims to reduce inferential variance. Through parallel analysis by multiple agents, it transforms **Multi-Source Information Matrix** from the previous stage into a structured evidence pool, which provides a stable decision basis for the subsequent adversarial debate. As illustrated in Figure 2, the matrix systematically

integrates 4 information sources that are orthogonal and complementary in evidence distribution:

- *In-news Information*: Comes from the original text of the target news and provides the basic core content.
- *Out-of-news Information*: Comes from Open Web retrieval and provides the most timely and broadest materials.
- *LLM internal Knowledge*: Comes from the model’s prior knowledge and provides generalized common sense.
- *LLM external Knowledge*: Comes from Wikipedia retrieval and provides the most precise summary for the core entities.

The subsequent analysis decomposes the tasks and assigns them to agents with distinct roles, so that each agent focuses on a specific quadrant of the matrix for specialized processing.



**Figure 2:** The diagram of Multi-Source Information Matrix and the quadrants used by LLM agents.

### 4.1.1 Linguist

Following prior studies (Shahid et al., 2022)(Zhou and Zafarani, 2018), the linguist agent is designed to systematically divide the text into 5 linguistic dimensions that are strongly associated with misinformation:

- *Sentence*: Lexical complexity, sentence length, and formality of tone.
- *Word*: Frequency of superlatives, affective language, and pronoun distribution.
- *Grammar*: Patterns of reported speech, passive voice, and negation.
- *Emotion*: Affective terms in the text and the headline, and the degree of incendiary tone.
- *Information Quality*: Presence of clickbait, information overload, or context mismatch.

To maintain objective independence, each dimension is evaluated in an isolated session. For each dimension, LLM explicitly indicates whether it reflects the news is real or fake. It uses 2 quadrants of Multi-source Information Matrix: LLM internal knowledge and in-news information.

### 4.1.2 Domain-Specific Expert

It operates in a dynamic and adaptive manner. The system first identifies the most relevant domain from the news text and assigns the agent a precise expert role, such as “economist” or “journalist”.

The expert with this role then analyses along the following 2 dimensions:

- *Knowledge Concordance*: Examine all claims, viewpoints, and details for sound reasoning; identify departures from common sense.
- *Logical Integrity*: Examine argument-to-conclusion coherence with domain-specific common sense; identify logical errors or unsupported leaps.

It uses 3 quadrants of Multi-source Information Matrix: in-news information, LLM internal knowledge, and LLM external knowledge.

### 4.1.3 Claim Verification

Existing research (Niu et al., 2024) demonstrates that claim-based fact checking effectively serves as a reference for LLM-based detection. In our system, a serial pipeline composed of **Claim Extractor** and **Claim Verifier** deconstructs and verifies factual claims in the news text. These 2 agents use 3 quadrants of Multi-source Information Matrix: LLM internal knowledge, in-news information, and out-of-news information from Open Web.

**Claim Extractor.**  $\mathcal{M}_{\text{extractor}}$  agent converts unstructured news text  $T$  into a set of verifiable structured claims. Its function is formalized as follows:

$$\mathcal{M}_{\text{extractor}}(T) \rightarrow \{q_{\text{core}}, \{q_{\text{sub}_1}, \dots, q_{\text{sub}_m}\}\} \quad (8)$$

where  $q_{\text{core}}$  is the core claim that determines the veracity of the news, and  $q_{\text{sub}}$  is a collection of supporting subclaims. All outputs are restated as concise and objective declarative sentences.

**Claim Verifier.** This agent verifies each claim  $q$  independently. To ensure precision and control cost (Purwar et al., 2024), a simple retrieval augmented generation (RAG) (Lewis et al., 2020) is implemented to build a highly relevant context  $\mathcal{C}_{\text{rel}}(q)$  for each  $q$  from the Open Web corpus  $\mathcal{I}_{\text{web}}$ .

The context consists of text chunks  $c_j$  whose cosine similarity with the claim representation  $\mathbf{v}_q$  exceeds a threshold  $\theta_{\text{sim}}$ :

$$\mathcal{C}_{\text{rel}}(q) = \{c_j \in \text{top-k}(\mathcal{I}_{\text{web}}, q) \mid \frac{\mathbf{v}_q \cdot \mathbf{v}_{c_j}}{\|\mathbf{v}_q\| \|\mathbf{v}_{c_j}\|} \geq \theta_{\text{sim}}\}. \quad (9)$$

All decisions are strictly based on  $\mathcal{C}_{\text{rel}}(q)$ . The final output includes a clear label (“*Supports*”, “*Refutes*”, or “*Not Enough Information*”), a brief reasoning, and evidence directly quoted from  $\mathcal{C}_{\text{rel}}(q)$  for suppressing hallucinations.

## 4.2 LLM Debate and Judgment

LLM Debate and Judgment state aims to suppress systemic bias. We introduce a multi-round adversarial debate framework that forces LLMs to explore both supporting and refuting views of the truthfulness of news equally. This design directly mitigates the thought degeneration (DoT) (Liang et al., 2023) phenomenon that often appears in a single linear reasoning chain.

---

### Algorithm 1. Dynamic Adversarial Debate

---

**Input:**  $\mathbb{E}$ : The complete evidence pool  
**Output:**  $D$ : The final decision  $\{\textit{Real}, \textit{Fake}\}$   
 $\mathbb{H} \leftarrow \emptyset$ ;  $A_{\text{con}} \leftarrow \text{null}$ ;  $D \leftarrow \textit{Insufficient}$   
**while**  $D = \textit{Insufficient}$  **do**  
     $A_{\text{pro}} \leftarrow \mathcal{M}_{\text{pro}}.\text{GenerateArgument}(\mathbb{E}, A_{\text{con}})$   
     $A_{\text{con}} \leftarrow \mathcal{M}_{\text{con}}.\text{GenerateArgument}(\mathbb{E}, A_{\text{pro}})$   
     $\mathbb{H} \leftarrow \mathbb{H} \cup \{(A_{\text{pro}}, A_{\text{con}})\}$   
     $D \leftarrow \text{Judge}.\text{Assess}(\mathbb{H})$   
**return**  $D$

---

As shown in Algorithm 1, a pair of debate agents  $\mathcal{M}_{\text{pro}}$  and  $\mathcal{M}_{\text{con}}$  act as opposing reasoners based on the evidence pool  $\mathbb{E}$ . In each round, the active debater first rebuts the opponent’s last argument and then presents a new argument. After each exchange, a judge agent evaluates the debate history  $\mathbb{H}$  and outputs a ternary judgment  $D$ .

This dynamic termination mechanism ensures that the debate stops once the information is sufficient for decision, effectively balancing the depth and efficiency of reasoning. The debate history  $\mathbb{H}$  provides a traceable reasoning chain composed of pro. and con. arguments, making the final judgment highly interpretable.

## 5 Experiment

### 5.1 Experimental Setting

**Datasets.** Following previous state-of-the-art work (Liu et al., 2024), our experiments are conducted

on two widely recognized fake news datasets: GossipCop and PolitiFact (Shu et al., 2020). GossipCop focuses on entertainment, mainly Hollywood celebrity news; PolitiFact focuses on politics, drawing on fact-checks of U.S. political figures. Considering that some links in the initial dataset have become invalid, we adopt the available version publicly re-released in (Su et al., 2023).

**Metrics.** We use accuracy and macro F1-score as evaluation metrics. F1-score is less affected by data imbalance, so it serves as the primary metric for assessment.

**Baselines.** We incorporate two groups of baselines. The first group includes advanced few-shot methods: PSM (Ni et al., 2020), MDFEND (Nan et al., 2021), ARG (Hu et al., 2024b), PSM (Ni et al., 2020), DKFND (Liu et al., 2024), and KPL (Jiang et al., 2022), with reliable metrics from the existing works (Jin et al., 2024; Hu et al., 2024c; Liu et al., 2024). The second group consists of representative zero-shot methods for comparison: Auto-CoT (Zhang et al., 2022), ReAct (Yao et al., 2023) equipped with a search API, HiSS (Zhang and Gao, 2023), Web Retrieval Agents (Tian et al., 2024) and FactAgent (Li et al., 2024b).

We implement ZoFia based on DeepSeek-V3 (DeepSeek-AI, 2024) and GPT-4o-mini (OpenAI, 2024), and compare it with only-LLM inference. All other LLM-based zero-shot methods are based on DeepSeek-V3.

**Implementation details.** We set the dynamic threshold of entity extraction as  $\lambda_{\text{init}} = 0.8$  and  $\Delta\lambda = 0.1$  and the decay factor of MMR as  $\gamma = 0.5$ . The maximum entries for Open Web retrieval is 10; for Wikipedia, we retrieve the first 3 sentences for each entry. The minimum similarity threshold for claim extraction is  $\theta_{\text{sim}} = 0.1$ . The NER model is selected as `dslim/bert-base-NER` (Tjong Kim Sang and De Meulder, 2003), and the SBERT model is selected as `all-MiniLM-L6-v2` (Reimers and Gurevych, 2019). For all based LLM, we set the temperature to 0.

Brave API<sup>2</sup> is selected as Open web retrieval API. We set the search cutoff date to the day before the publication date of each sample URL to strictly prevent label information leakage.

### 5.2 Main Experiment

ZoFia demonstrates exceptional performance on both datasets as shown in Table 1, consistently out-

<sup>2</sup><https://brave.com/search/api/>

**Table 1:** Accuracy (Acc.) and F1-score (F1) comparison of few-shot / zero-shot methods on PolitiFact and GossipCop. The **bold** and underlined denote the best and second-best performance.

Category	Method	LLM Usage	PolitiFact		GossipCop	
			Accuracy	F1-Score	Accuracy	F1-Score
Few-shot	PSM (Nan et al., 2021)	Non-LLM	70.00	49.15	77.44	41.73
	MDFEND (Nan et al., 2021)	Non-LLM	65.50	62.30	41.27	40.20
	KPL (Jiang et al., 2022)	Non-LLM	58.33	60.40	42.71	42.08
	ARG (Hu et al., 2024b)	LLM-assisted	74.00	67.16	61.41	42.32
	DKFND (Liu et al., 2024)	LLM-assisted	87.00	82.43	<b>82.37</b>	55.22
Zero-shot	Auto-CoT (Zhang et al., 2022)	LLM-based	<u>89.65</u>	73.67	60.01	48.15
	ReAct (Search API) (Yao et al., 2023)	LLM-based	74.73	67.64	74.03	47.30
	HiSS (Zhang and Gao, 2023)	LLM-based	64.82	56.80	68.81	40.40
	Web Retrieval Agents (Tian et al., 2024)	LLM-based	77.83	64.88	66.62	46.54
	FactAgent (Li et al., 2024b)	LLM-based	80.59	70.06	73.80	<u>56.22</u>
	DeepSeek-v3 (DeepSeek-AI, 2024)	Only-LLM	78.99	40.24	61.03	28.17
	GPT-4o-mini (OpenAI, 2024)	Only-LLM	73.58	41.66	66.73	33.18
	<b>ZoFia (DeepSeek-v3)</b>	LLM-based	<b>91.52</b>	<b>87.88</b>	<u>79.04</u>	<b>62.22</b>
	<b>ZoFia (GPT-4o-mini)</b>	LLM-based	75.28	<u>75.19</u>	68.90	56.20

performing existing zero-shot and few-shot baselines. On the fact-intensive PolitiFact dataset, ZoFia’s advantages are particularly pronounced, with its accuracy (91.5%) and F1-score (87.9%) not only far exceeding other zero-shot methods but also surpassing all few-shot baselines by a substantial performance margin.

The detection task on the GossipCop dataset presents greater challenges, as its content often proves difficult to verify due to subjectivity and factual ambiguity. In this scenario, ZoFia achieves the highest F1-score (61.2%), surpassing all baseline models. Although DKFND attains slightly higher accuracy, its F1-score is notably lower. In contrast, ZoFia’s superior F1-score demonstrates more balanced and robust detection performance across both true and fake news categories while maintaining high precision.

The results clearly show that as a zero-shot method, ZoFia consistently outperforms all zero-shot baselines and even most few-shot methods. Compared to only-LLM inference, ZoFia provides stable and substantial performance gains, demonstrating the superiority of its architecture.

### 5.3 Ablation Study

To assess the influence of each component in ZoFia, we conduct extensive ablation on the GossipCop dataset. As shown in Table 2, removing Open Web retrieval causes a marked decline, which demonstrates the importance of instant external evidence. Performance decreases further when all retrieval sources are removed. For LLM collaboration, re-

moving the linguist or expert weakens the framework. The expert contributes more, with a 3.9% drop in F1-Score. Moreover, removing both causes cumulative degradation. Removing the debate reduces the F1-Score by 2.0%, demonstrating its role in suppressing single-perspective bias.

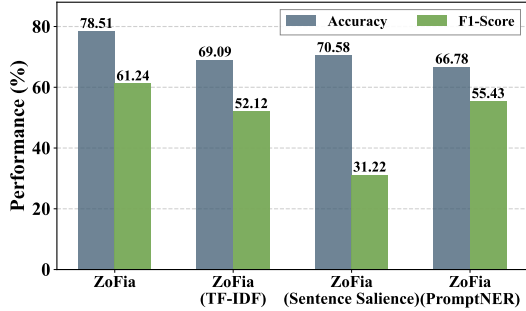
**Table 2:** Results of ablation study on GossipCop

Component	Accuracy	F1-Score
ZoFia (DeepSeek-v3)	79.04	62.22
w/o Wikipedia retrieval	77.67	59.90
w/o Open Web retrieval	71.71	58.14
w/o all retrievals	71.37	58.26
w/o linguist analysis	77.86	58.87
w/o expert analysis	76.77	57.32
w/o all analyses	73.51	57.90
w/o claim verification	79.23	60.50
w/o debate	75.01	59.16

### 5.4 Effectiveness of Entity-Guided Retrieval

To verify the effectiveness of the keyword extraction module in ZoFia, we compare it with several representative methods. As shown in Figure 3, replacing our keyword extraction module with the classic TF-IDF (Sparck Jones, 1972) algorithm lowers the F1-Score by 9.1%. This finding emphasizes the need for semantic-based Hierarchical Saliency. We also evaluate Sentence Saliency (Bullough et al., 2024) strategy, which retrieves a full sentence. The result shows an extreme imbalance: high Accuracy with low F1-Score. The reason is that using sentences as queries tends to retrieve pages that are highly similar to the original news. This source con-

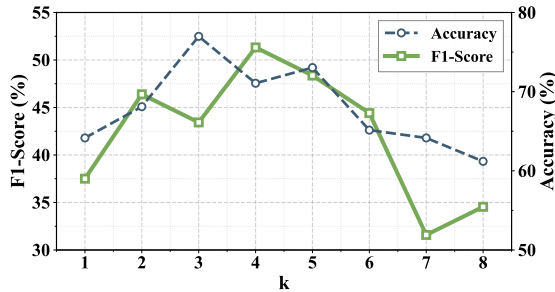
tamination (Deng et al., 2023) induces confirmation bias in LLM agents. The LLM-based promptNER (Ashok and Lipton, 2023) cannot quantitatively evaluate the importance of each entity, so its performance gain is also inferior to ZoFia’s extractor.



**Figure 3:** Performance comparison of ZoFia’s keyword extraction with other extraction methods.

### 5.5 Sensitivity Analysis of SC-MMR

To motivate the design of the dynamic weight  $\lambda_k$ , we first conduct an experiment with the weight fixed at  $\lambda = 0.5$  to observe the direct impact of the keyword count  $k$  on performance. As shown in Figure 4, the F1-Score remains stable when  $k \leq 6$  but drops sharply beyond this point.



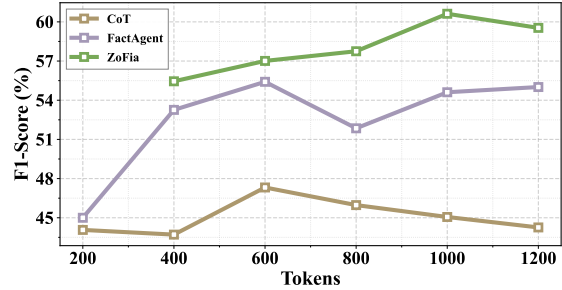
**Figure 4:** The effect of the number of keywords  $k$  on performance (F1-Score).

This phenomenon reveals an inflection point, where the number of keywords exceeds 6, the risk of introducing noise and redundancy begins to outweigh the benefits of new information. It suggests that the strategic focus must shift from relevance to diversity near the critical inflection point of  $k \approx 6$ . Consequently, we adopt the annealing schedule form defined in Equation 5 and determine its parameters as  $\alpha = 0.3$  and  $\beta = 2.5$  to fit this downward trend.

### 5.6 Efficiency of Token Utilization

To investigate ZoFia’s reasoning efficiency, a controlled comparison experiment is conducted. We

provide CoT (Wei et al., 2022), FactAgent (Li et al., 2024b), and ZoFia with the same materials and impose a unified limit on output tokens to evaluate their performance. Since ZoFia cannot complete full reasoning below 400 tokens, we set its evaluation range to 400 tokens and above.



**Figure 5:** Performance (F1-Score) comparison of methods under maximum output token limits.

The results are shown in Figure 5. At all budget points above 400 tokens, ZoFia achieves a clearly higher F1-score than other baseline methods, which demonstrates high reasoning efficiency. In contrast to CoT and FactAgent, whose performance are saturated after around 600 tokens, ZoFia show a stable upward trend even at 1200 tokens. This indicates that ZoFia can consistently convert tokens into performance gains within an acceptable budget.

Note that this experiment only constrains output tokens, because prompt caching (Gim et al., 2024) can extensively amortize input overhead, making output cost more critical for efficiency.

## 6 Conclusion

In this paper, we propose ZoFia, a retrieval-augmented multi-agent zero-shot framework for fake news detection, to address the cognitive conflict of a single LLM between content reasoning and fact checking. It first utilizes our novel Hierarchical Saliency and Saliency-Calibrated Minimum Marginal Relevance (SC-MMR) algorithm to accurately extract core entities from news text, which then guide dual-source retrieval from Open Web and Wikipedia. Next, the multi-agent system conducts analysis and verification in parallel and make a final verdict via adversarial debate. This process effectively reduces confirmation bias from a single reasoning perspective and ensures robust and explainable results. Comprehensive experiments on two public datasets show that ZoFia outperforms existing zero-shot baselines in both performance and efficiency.

## 615 Limitations

616 Though ZoFia show strong detection capacities, its  
617 application and evaluation face multiple constraints.  
618 Due to copyright and privacy concerns, there has  
619 been a recent lack of high-quality, and continuously  
620 updated public datasets. It prevents us from evalu-  
621 ating ZoFia on the most recent news. Building the  
622 next generation of benchmark datasets that meet  
623 ethical standards and reflect real-world information  
624 dynamics is a critical step in this field.

625 The efficiency of using external retrieval can also  
626 be improved. Since it is not our main focus, we  
627 integrate only a lightweight retrieval augmented  
628 generation (RAG) module in claim verification. Fu-  
629 ture work that adopts more advanced RAG architec-  
630 tures, such as re-ranking model and more advanced  
631 mechanism, may further strengthen the exploita-  
632 tion of external knowledge. We plan to conduct  
633 more comprehensive experimental benchmark eval-  
634 uation over the detection capabilities of LLMs and  
635 Multi-agent system for fake news (Guo et al., 2025;  
636 Kuntur et al., 2024; Jiang et al., 2024).

637 ZoFia currently focuses on the text modality.  
638 Modern misinformation increasingly appears in  
639 multi-modal form that combines images and text.  
640 Extending ZoFia to the multimodal domain has  
641 strong potential. One direction is to introduce vi-  
642 sion language models (VLMs) as a dedicated visual  
643 expert. Another is to study how different modalities  
644 interact during the debate process to achieve  
645 effective fusion.

## 646 References

647 Dhananjay Ashok and Zachary C Lipton. 2023. Prompt-  
648 ner: Prompting for named entity recognition. *arXiv*  
649 *preprint arXiv:2305.15444*.

650 Vian Bakir and Andrew McStay. 2018. Fake news and  
651 the economy of emotions: Problems, causes, solu-  
652 tions. *Digital journalism*, 6(2):154–175.

653 Benjamin Bullough, Harrison Lundberg, Chen Hu, and  
654 Weihang Xiao. 2024. Predicting entity salience in  
655 extremely short documents. In *Proceedings of the*  
656 *2024 Conference on Empirical Methods in Natural*  
657 *Language Processing: Industry Track*, pages 50–64.

658 Jaime Carbonell and Jade Goldstein. 1998. The use of  
659 mmr, diversity-based reranking for reordering doc-  
660 uments and producing summaries. In *Proceedings*  
661 *of the 21st annual international ACM SIGIR confer-*  
662 *ence on Research and development in information*  
663 *retrieval*, pages 335–336.

Claudio Carpineto and Giovanni Romano. 2012. A  
664 survey of automatic query expansion in information  
665 retrieval. *Acm Computing Surveys (CSUR)*, 44(1):1–  
666 50. 667

Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu,  
668 Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan  
669 Liu. 2023. Chateval: Towards better llm-based eval-  
670 uators through multi-agent debate. *arXiv preprint*  
671 *arXiv:2308.07201*. 672

Canyu Chen and Kai Shu. 2023. Can llm-generated  
673 misinformation be detected? *arXiv preprint*  
674 *arXiv:2309.13788*. 675

Jeffrey Cheng, Marc Marone, Orion Weller, Dawn  
676 Lawrie, Daniel Khashabi, and Benjamin Van Durme.  
677 2024. Dated data: Tracing knowledge cutoffs in large  
678 language models. *arXiv preprint arXiv:2403.12958*.  
679

DeepSeek-AI. 2024. [Deepseek-v3 technical report](#).  
680 *Preprint*, arXiv:2412.19437. 681

Chunyuan Deng, Yilun Zhao, Xiangru Tang, Mark Ger-  
682 stein, and Arman Cohan. 2023. Investigating data  
683 contamination in modern benchmarks for large lan-  
684 guage models. *arXiv preprint arXiv:2311.09783*.  
685

Jesse Dunietz and Dan Gillick. 2014. A new entity  
686 salience task with millions of training examples. In  
687 *Proceedings of the 14th Conference of the European*  
688 *Chapter of the Association for Computational Lin-*  
689 *guistics, volume 2: Short Papers*, pages 205–209.  
690

Jessica Maria Echterhoff, Yao Liu, Abeer Alessa, Julian  
691 McAuley, and Zexue He. 2024. Cognitive bias in  
692 decision-making with llms. In *Findings of the asso-*  
693 *ciation for computational linguistics: EMNLP 2024*,  
694 pages 12640–12653. 695

Marc Fisher, John Woodrow Cox, and Peter Hermann.  
696 2016. Pizzagate: From rumor, to hashtag, to gunfire  
697 in dc. *Washington Post*, 6:8410–8415. 698

In Gim, Guojun Chen, Seung-seob Lee, Nikhil Sarda,  
699 Anurag Khandelwal, and Lin Zhong. 2024. Prompt  
700 cache: Modular attention reuse for low-latency infer-  
701 ence. *Proceedings of Machine Learning and Systems*,  
702 6:325–338. 703

Hao Guo, Zihan Ma, Zhi Zeng, Minnan Luo, Weixin  
704 Zeng, Jiuyang Tang, and Xiang Zhao. 2025. Each  
705 fake news is fake in its own way: An attribution multi-  
706 granularity benchmark for multimodal fake news de-  
707 tection. In *Proceedings of the AAAI Conference on*  
708 *Artificial Intelligence*, volume 39, pages 228–236.  
709

Feng Hou, Ruili Wang, Jun He, and Yi Zhou.  
710 2021. Improving entity linking through seman-  
711 tic reinforced entity embeddings. *arXiv preprint*  
712 *arXiv:2106.08495*. 713

Nathaniel Hoy and Theodora Koulouri. 2022. Exploring  
714 the generalisability of fake news detection models.  
715 In *2022 IEEE International Conference on Big Data*  
716 *(Big Data)*, pages 5731–5740. IEEE. 717

718	Beizhe Hu, Qiang Sheng, Juan Cao, Yuhui Shi, Yang Li, Danding Wang, and Peng Qi. 2024a. Bad actor, good advisor: Exploring the role of large language models in fake news detection. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pages 22105–22113.	773
719		774
720		775
721		776
722		
723		
724	Beizhe Hu, Qiang Sheng, Juan Cao, Yuhui Shi, Yang Li, Danding Wang, and Peng Qi. 2024b. Bad actor, good advisor: Exploring the role of large language models in fake news detection. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 38, pages 22105–22113.	777
725		778
726		779
727		780
728		781
729		
730	Weiqi Hu, Ye Wang, Yan Jia, Qing Liao, and Bin Zhou. 2024c. A multi-modal prompt learning framework for early detection of fake news. In <i>Proceedings of the International AAAI Conference on Web and Social Media</i> , volume 18, pages 651–662.	782
731		783
732		784
733		
734		
735	Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. <i>ACM computing surveys</i> , 55(12):1–38.	785
736		786
737		787
738		788
739		789
740	Gongyao Jiang, Shuang Liu, Yu Zhao, Yueheng Sun, and Meishan Zhang. 2022. Fake news detection via knowledgeable prompt learning. <i>Information Processing &amp; Management</i> , 59(5):103029.	790
741		791
742		792
743		793
744	Xuefeng Jiang, Lvhua Wu, Sheng Sun, Jia Li, Jingjing Xue, Yuwei Wang, Tingting Wu, and Min Liu. 2024. Investigating large language models for code vulnerability detection: An experimental study. <i>arXiv preprint arXiv:2412.18260</i> .	794
745		795
746		796
747		797
748		798
749	Yiqiao Jin, Minje Choi, Gaurav Verma, Jindong Wang, and Srijan Kumar. 2024. Mm-soc: Benchmarking multimodal large language models in social media platforms. <i>arXiv preprint arXiv:2402.14154</i> .	799
750		800
751		801
752		
753	Rohit Kumar Kaliyar, Anurag Goswami, and Pratik Narang. 2021. Fakebert: Fake news detection in social media with a bert-based deep learning approach. <i>Multimedia tools and applications</i> , 80(8):11765–11788.	802
754		803
755		804
756		805
757		
758	Soveatin Kuntur, Anna Wróblewska, Marcin Paprzycki, and Maria Ganzha. 2024. Fake news detection: It’s all in the data! <i>arXiv preprint arXiv:2407.02122</i> .	806
759		807
760		808
761	Xiaochong Lan, Chen Gao, Depeng Jin, and Yong Li. 2024. Stance detection with collaborative role-infused llm-based agents. In <i>Proceedings of the international AAAI conference on web and social media</i> , volume 18, pages 891–903.	809
762		810
763		811
764		812
765		813
766	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, and 1 others. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. <i>Advances in neural information processing systems</i> , 33:9459–9474.	814
767		815
768		816
769		
770		
771		
772		
	Jia Li, Lijie Hu, Zhixian He, Jingfeng Zhang, Tianhang Zheng, and Di Wang. 2024a. Text guided image editing with automatic concept locating and forgetting. <i>arXiv preprint arXiv:2405.19708</i> .	817
		818
		819
		820
	Jia Li, Lijie Hu, Jingfeng Zhang, Tianhang Zheng, Hua Zhang, and Di Wang. 2025. Fair text-to-image diffusion via fair mapping. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 39, pages 26256–26264.	821
		822
		823
		824
		825
	Xinyi Li, Yongfeng Zhang, and Edward C Malthouse. 2024b. Large language model agent for fake news detection. <i>arXiv preprint arXiv:2405.01593</i> .	826
		827
	Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Shuming Shi, and Zhaopeng Tu. 2023. Encouraging divergent thinking in large language models through multi-agent debate. <i>arXiv preprint arXiv:2305.19118</i> .	
	Ye Liu, Jiajun Zhu, Kai Zhang, Haoyu Tang, Yanghai Zhang, Xukai Liu, Qi Liu, and Enhong Chen. 2024. Detect, investigate, judge and determine: A novel llm-based framework for few-shot fake news detection. <i>arXiv preprint arXiv:2407.08952</i> .	
	Yuhan Liu, Yuxuan Liu, Xiaoqing Zhang, Xiuying Chen, and Rui Yan. 2025. The truth becomes clearer through debate! multi-agent systems with large language models unmask fake news. In <i>Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , pages 504–514.	
	Federico Monti, Fabrizio Frasca, Davide Eynard, Damon Mannion, and Michael M Bronstein. 2019. Fake news detection on social media using geometric deep learning. <i>arXiv preprint arXiv:1902.06673</i> .	
	Qiong Nan, Juan Cao, Yongchun Zhu, Yanyan Wang, and Jintao Li. 2021. Mdfend: Multi-domain fake news detection. In <i>Proceedings of the 30th ACM international conference on information &amp; knowledge management</i> , pages 3343–3347.	
	Qiong Nan, Qiang Sheng, Juan Cao, Beizhe Hu, Danding Wang, and Jintao Li. 2024. Let silence speak: Enhancing fake news detection with generated comments from large language models. In <i>Proceedings of the 33rd ACM International Conference on Information and Knowledge Management</i> , pages 1732–1742.	
	Bo Ni, Zhichun Guo, Jianing Li, and Meng Jiang. 2020. Improving generalizability of fake news detection methods using propensity score matching. <i>arXiv preprint arXiv:2002.00838</i> .	
	Cheng Niu, Yang Guan, Yuanhao Wu, Juno Zhu, Jun-tong Song, Randy Zhong, Kaihua Zhu, Siliang Xu, Shizhe Diao, and Tong Zhang. 2024. Veract scan: Retrieval-augmented fake news detection with justifiable reasoning. <i>arXiv preprint arXiv:2406.10289</i> .	
	OpenAI. 2024. Gpt-4o-mini. <a href="https://platform.openai.com/docs/models">https://platform.openai.com/docs/models</a> . Accessed: 2024-10.	

828	Bohdan M Pavlyshenko. 2023. Analysis of disinformation and fake news detection using fine-tuned large language model. <i>arXiv preprint arXiv:2309.04704</i> .	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837.	880
829			881
830			882
831	Anupam Purwar and 1 others. 2024. Evaluating the efficacy of open-source llms in enterprise-specific rag systems: A comparative study of performance and scalability. <i>arXiv preprint arXiv:2406.11424</i> .		883
832			884
833			885
834		Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023. <a href="#">React: Synergizing reasoning and acting in language models</a> . <i>Preprint</i> , arXiv:2210.03629.	886
835	Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. <i>arXiv preprint arXiv:1908.10084</i> .		887
836			888
837			889
838	Wajiha Shahid, Bahman Jamshidi, Saqib Hakak, Haruna Isah, Wazir Zada Khan, Muhammad Khurram Khan, and Kim-Kwang Raymond Choo. 2022. Detecting and mitigating the dissemination of fake news: Challenges and future research opportunities. <i>IEEE Transactions on Computational Social Systems</i> , 11(4):4649–4662.	Xuan Zhang and Wei Gao. 2023. Towards llm-based fact verification on news claims with a hierarchical step-by-step prompting method. <i>arXiv preprint arXiv:2310.00305</i> .	890
839			891
840			892
841			893
842		Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. 2022. <a href="#">Automatic chain of thought prompting in large language models</a> . <i>Preprint</i> , arXiv:2210.03493.	894
843			895
844			896
845	Kai Shu, Deepak Mahudeswaran, Suhang Wang, Dongwon Lee, and Huan Liu. 2020. Fakenewsnet: A data repository with news content, social context, and spatiotemporal information for studying fake news on social media. <i>Big data</i> , 8(3):171–188.	Xinyi Zhou and Reza Zafarani. 2018. Fake news: A survey of research, detection methods, and opportunities. <i>arXiv preprint arXiv:1812.00315</i> , 2:13.	897
846			898
847			899
848			900
849		Xinyi Zhou and Reza Zafarani. 2020. A survey of fake news: Fundamental theories, detection methods, and opportunities. <i>ACM Computing Surveys (CSUR)</i> , 53(5):1–40.	901
850	Karen Sparck Jones. 1972. A statistical interpretation of term specificity and its application in retrieval. <i>Journal of documentation</i> , 28(1):11–21.		902
851			903
852			904
853	Jinyan Su, Claire Cardie, and Preslav Nakov. 2023. Adapting fake news detection to the era of large language models. <i>arXiv preprint arXiv:2311.04917</i> .	<b>A Implementation Details</b>	905
854		<b>A.1 Entity Extractor</b>	906
855		The goal of entity extraction is to provide a stable set of semantic anchors for subsequent retrieval. To prevent retrieval link failures caused by entity sparsity or confidence distribution shifts, we introduce dynamic confidence threshold filtering after extraction, maintaining a controllable balance between the quality and usability of the entity set.	907
856	Hexiang Tan, Fei Sun, Wanli Yang, Yuanzhuo Wang, Qi Cao, and Xueqi Cheng. 2024. Blinded by generated contexts: How language models merge generated and retrieved contexts when knowledge conflicts? <i>arXiv preprint arXiv:2401.11911</i> .	Algorithm 2 details the implementation of this process. The inputs are the news text $\mathcal{T}$ , the Named Entity Recognition (NER) model $\mathcal{M}_{\text{NER}}$ , an initial threshold $\tau_0$ , and a threshold decay step $\Delta$ . The output is the filtered entity set $\mathcal{E}_{\text{sel}}$ . The algorithm first runs NER on the full text to obtain the raw entity set $\mathcal{E}_{\text{raw}} = \mathcal{M}_{\text{NER}}(\mathcal{T})$ , where each entity $e$ comes with a model confidence score $e.\text{ner\_score}$ that characterizes the reliability of the entity boundary and type prediction. If $\mathcal{E}_{\text{raw}}$ is empty, the algorithm directly returns an empty set. This indicates that the text lacks identifiable entity signals or the model cannot provide reliable predictions, in which case the retrieval stage should degrade to coarser-grained information clues or skip the entity-guided strategy.	908
857			909
858			910
859			911
860			912
861	Jacob-Junqi Tian, Hao Yu, Yury Orlovskiy, Tyler Verghe, Mauricio Rivera, Mayank Goel, Zachary Yang, Jean-Francois Godbout, Reihaneh Rabbany, and Kellin Pelrine. 2024. Web retrieval agents for evidence-based misinformation detection. <i>arXiv preprint arXiv:2409.00009</i> .		913
862			914
863			915
864			916
865			917
866			918
867	Erik F. Tjong Kim Sang and Fien De Meulder. 2003. <a href="#">Introduction to the CoNLL-2003 shared task: Language-independent entity recognition</a> . In <i>Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003</i> , pages 142–147.		919
868			920
869			921
870			922
871			923
872			924
873	Yue Wan, Xiaowei Jia, and Xiang Lorraine Li. 2025. Unveiling confirmation bias in chain-of-thought reasoning. <i>arXiv preprint arXiv:2506.12301</i> .		925
874			926
875			927
876	Keyu Wang, Jin Li, Shu Yang, Zhuoran Zhang, and Di Wang. 2025. When truth is overridden: Uncovering the internal origins of sycophancy in large language models. <i>arXiv preprint arXiv:2508.02087</i> .		928
877			929
878			930
879			931

Specifically, it starts with a high initial threshold  $\tau_0$  and retains only entities satisfying  $e.\text{ner\_score} \geq \tau$ , prioritizing entity quality and precision. If  $\mathcal{E}_{\text{sel}}$  remains empty under the current threshold, the threshold decreases stepwise by  $\Delta$  until the system selects at least one entity or the threshold drops to a lower bound of 0.1. This mechanism reflects a clear constraint: entity extraction must provide a non-empty semantic entry point for subsequent retrieval; otherwise, retrieval degenerates into anchorless generalized queries. Meanwhile, the lower bound threshold 0.1 prevents the uncontrolled introduction of low-confidence noisy entities, thereby maintaining a baseline quality while ensuring usability.

Dynamic threshold filtering transforms entity selection from a fixed hyperparameter setting into a sample-adaptive process. This enables the system to cover both entity-dense and entity-sparse news narratives, providing a stable candidate foundation for subsequent Hierarchical Saliency and SC-MMR.

---

**Algorithm 2.** Entity Extraction and Dynamic-Threshold Filtering

---

**Input:**  $\mathcal{T}$ : News text;  
 $\mathcal{M}_{\text{NER}}$ : NER model;  
 $\tau_0$ : Initial confidence threshold;  
 $\Delta$ : Threshold decay step  
**Output:**  $\mathcal{E}_{\text{sel}}$ : Selected entity set  
 $\mathcal{E}_{\text{raw}} \leftarrow \mathcal{M}_{\text{NER}}(\mathcal{T})$   
**if**  $\mathcal{E}_{\text{raw}} = \emptyset$  **then**  
  **return**  $\emptyset$   
 $\tau \leftarrow \tau_0$ ;  $\mathcal{E}_{\text{sel}} \leftarrow \emptyset$   
**while**  $\mathcal{E}_{\text{sel}} = \emptyset \wedge \tau \geq 0.1$  **do**  
   $\mathcal{E}_{\text{sel}} \leftarrow \{e \in \mathcal{E}_{\text{raw}} \mid e.\text{ner\_score} \geq \tau\}$   
   $\tau \leftarrow \tau - \Delta$   
**return**  $\mathcal{E}_{\text{sel}}$

---

## A.2 SC-MMR Algorithm

A larger set of query keywords is not necessarily better, as it defines the boundary of the external evidence space. SC-MMR is designed to select a set of core entities that are sufficiently informative yet minimally redundant within a limited scale, which makes subsequent retrieval more stable and controllable.

Algorithm 3 presents the complete procedure of SC-MMR. The input is a mapping from entities to Hierarchical Saliency scores  $D : \{U_i \mapsto S_{\text{Hier}}(U_i)\}$ , and a model  $M_{\text{SBERT}}$  for calculating entity vector representations. The output is the final selected entity set  $U_{\text{selected}}$ . The algorithm first

performs one-time pre-processing on the candidate entity set  $U = \text{keys}(D)$ . It uses  $M_{\text{SBERT}}$  to encode each entity into a vector and constructs a pairwise similarity matrix  $M_{\text{sim}}$ , where  $M_{\text{sim}}[i, j]$  measures the semantic proximity between entities  $U_i$  and  $U_j$ . Pre-computing this matrix moves repetitive similarity calculations out of the iterative loop, so subsequent rounds only require table look-ups and maximum value operations.

The algorithm then enters the selection process. It first selects the entity with the highest saliency score  $D[U_i]$  from the candidate set  $U_{\text{candidate}}$  as the seed  $U^*$  and adds it to  $U_{\text{selected}}$ . This initialization ensures the first step is entirely driven by relevance. It then enters a ‘while’ loop, where the goal of each round is to find the entity among the remaining candidates that yields the highest combined score of relevance and diversity. For any candidate  $U_i$ , the algorithm first calculates:

$$\text{sim}_{\text{max}} = \max_{U_j \in U_{\text{selected}}} M_{\text{sim}}[i, j],$$

which represents the similarity between the candidate and the most similar entity in the current selected set, indicating the maximum redundancy that adding  $U_i$  might introduce. Subsequently, it scores the candidate using:

$$\text{MMR}_{\text{curr}} = \lambda_k \cdot D[U_i] - (1 - \lambda_k) \cdot \text{sim}_{\text{max}}$$

Here,  $\lambda_k$  is a weight schedule that varies with the number of selected entities  $k = |U_{\text{selected}}|$ , with a lower bound of 0.1. Intuitively, a larger  $\lambda_k$  at small  $k$  favors high-saliency core entities to ensure factual coverage. As  $k$  increases,  $\lambda_k$  gradually decreases to strengthen the diversity penalty. This suppresses the repetitive inclusion of coreferential or semantically close entities, thereby mitigating query drift caused by keyword inflation. Each round iterates through  $U_{\text{candidate}}$  to obtain the global optimal  $U_{\text{best}}$  and  $\text{MMR}_{\text{best}}$  as the best incremental choice.

The algorithm does not fix the number of output keywords but adaptively stops via a relative change termination criterion. Let the optimal score adopted in the previous round be  $\text{MMR}_{\text{prev}}^*$ . If the current round’s  $\text{MMR}_{\text{best}} \leq \gamma \cdot \text{MMR}_{\text{prev}}^*$ , the system considers the marginal gain significantly decayed. Since continuing to add new keywords is likely to introduce noise and redundancy, the loop terminates early. Otherwise, it updates  $\text{MMR}_{\text{prev}}^* \leftarrow \text{MMR}_{\text{best}}$ , adds  $U_{\text{best}}$  to  $U_{\text{selected}}$ , and removes the entity from  $U_{\text{candidate}}$ . The finally returned  $U_{\text{selected}}$  thus reflects three properties: early

1000  
1001  
1002  
1003

stages prioritize salience for key coverage, later stages use diversity penalties to suppress redundancy, and the termination rule automatically truncates the set size when marginal utility declines.

### Algorithm 3. Saliency-Calibrated Maximal Marginal Relevance (SC-MMR)

```

Input:  $\mathcal{D} : \{U_i \rightarrow S_{\text{Hier}}(U_i)\}$ : Mapping from entity
to its saliency score;
 $\mathcal{M}_{\text{SBERT}}$ : SBERT encoder for entity embeddings;
 $\gamma$ : Decay factor for termination criterion
Output:  $\mathcal{U}_{\text{selected}}$ : Final set of selected entities
 $\mathcal{U} \leftarrow \text{keys}(\mathcal{D})$ ;  $\mathbf{V} \leftarrow \{\mathcal{M}_{\text{SBERT}}(U_i) \mid \forall U_i \in \mathcal{U}\}$ 
Compute pairwise similarity matrix  $\mathbf{M}_{\text{sim}}$  from  $\mathbf{V}$ 
if  $\mathcal{U}_{\text{candidate}} = \emptyset$  then
   $\leftarrow$  return  $\mathcal{U}_{\text{selected}}$ 
 $U^* \leftarrow \arg \max_{U_i \in \mathcal{U}_{\text{candidate}}} \mathcal{D}[U_i]$ 
 $\mathcal{U}_{\text{selected}} \leftarrow \{U^*\}$ ;  $\mathcal{U}_{\text{candidate}} \leftarrow \mathcal{U} \setminus \{U^*\}$ 
 $\text{MMR}_{\text{prev}}^* \leftarrow 1.0$ 
while  $\mathcal{U}_{\text{candidate}} \neq \emptyset$  do
   $k \leftarrow |\mathcal{U}_{\text{selected}}|$ 
   $\lambda_k \leftarrow \max(0.1, 1.0 - e^{0.3k-2.5})$ 
   $\text{MMR}_{\text{best}} \leftarrow -\infty$ ;  $U_{\text{best}} \leftarrow \text{null}$ 
  foreach  $U_i \in \mathcal{U}_{\text{candidate}}$  do
     $\text{sim}_{\text{max}} \leftarrow \max_{U_j \in \mathcal{U}_{\text{selected}}} \mathbf{M}_{\text{sim}}[i, j]$ 
     $\text{MMR}_{\text{curr}} \leftarrow \lambda_k \cdot \mathcal{D}[U_i] - (1 - \lambda_k) \cdot \text{sim}_{\text{max}}$ 
    if  $\text{MMR}_{\text{curr}} > \text{MMR}_{\text{best}}$  then
       $\leftarrow$   $\text{MMR}_{\text{best}} \leftarrow \text{MMR}_{\text{curr}}$ ;  $U_{\text{best}} \leftarrow U_i$ 
  if  $\text{MMR}_{\text{best}} \leq \gamma \cdot \text{MMR}_{\text{prev}}^*$  then
     $\leftarrow$  break
   $\text{MMR}_{\text{prev}}^* \leftarrow \text{MMR}_{\text{best}}$ 
   $\mathcal{U}_{\text{selected}} \leftarrow \mathcal{U}_{\text{selected}} \cup \{U_{\text{best}}\}$ 
   $\mathcal{U}_{\text{candidate}} \leftarrow \mathcal{U}_{\text{candidate}} \setminus \{U_{\text{best}}\}$ 
return  $\mathcal{U}_{\text{selected}}$ 

```

1004

### A.3 Prompt of Linguist

```

## Role
You are a linguistic analyst for a news
channel, tasked with profiling articles
to help detect fake news.

## Instruction
You should analyze the provided news text
against the following linguistic
features of fake news:
1. **Sentence:** Longer sentences, simple
words, and a more informal tone (e.g.,
expletives).
2. **Word:** More superlatives, emotional or
vague language, with fewer reporting
verbs and 1st/2nd-person pronouns.
3. **Grammar:** Frequent use of reported
speech, passive voice, and negation.
Paraphrasing is less common.
4. **Emotion:** A higher ratio of emotional
words. Headlines are often sensational
and designed to provoke readers.
5. **Information Quality:** Information
overload or deficit, mismatched context,
and more clickbait patterns.

```

1005

```

## Expectation
- When I input a feature name (e.g., `
Grammar`), you will provide your
analysis about this feature.
- Your output MUST be a single, direct
analytical paragraph without any
formatting, LIMITED to 25 words.

```

1006

### A.4 Prompt of Expert

1007

```

## Role
As a professional and renowned {expert_role},
you are fact-checking a news article.

## Instruction
Identify all sentences that can
significantly affect the truthfulness of
news, and analyze the news from the
following 2 aspects:
1. **Knowledge Concordance:** Analyze the
rationality of all factual claims,
viewpoints and details in the news text.
Identify any content that deviates from
common sense or exhibits sensationalism
.
2. **Logical Integrity:** Analyze the
coherence of the article's arguments and
conclusions based on your field's
reasoning principles. Identify any
logical fallacies or unsupported
inferences.
Information from Wikipedia explains the key
entities in the news. (For reference
only, not necessarily relevant)

## Expectation
- Your output must be an unordered list (2
items), LIMITED to 100 words.
- DO NOT fabricate any information. All
analyses must be based on the provided
text.

```

1008

### A.5 Prompt of Claim Extractor

1009

```

## Role
You are a news fact-checker tasked with
summarizing all factual claims within an
article for subsequent verification.

## Instruction
Carefully read and analyze the provided news
article sentence by sentence. Identify
its core claim (directly determines the
truthfulness of the news) and all
supporting sub-claims (strongly related
to the core claim). Paraphrase each
claim into a concise, objective, and
declarative sentence.
If multiple claims are strongly related,
merge them into a single claim. Do not
use pronouns in the claim; replace all
pronouns with explicit nouns.

```

1010

## Expectation  
- Output 2-4 sub-claims most relevant to the core claim.  
- DO NOT fabricate any claims. All claims must originate from the provided text.

## A.6 Prompt of Claim Verifier

## Role  
You are a professional news fact-checker skilled at logical reasoning and text analysis.

## Instruction  
Your primary task is to fact-check a given claim based on the provided web information.  
The system has extracted the most relevant sentences from web information via RAG. Assume the web information is reliable and determine whether the information sufficiently supports or refutes the claim.  
You must make extremely full use of every piece of extracted information.

## Expectation  
- DO NOT fabricate any claims. All contents must originate from the provided text.

## A.7 Prompt of Debater

## Role  
You are an extremely cautious and logically rigorous final judge. Your only task is to determine when sufficient evidence has been presented to make a final ruling (real or fake) in a debate about the truthfulness of news.

## Instruction  
You argue that the news is <REAL/FAKE>. Find out all the most persuasive supporting evidences from the above provided text, and then back it up with concise reasons.

You'll receive evidences and reasons from opposing debaters in the subsequent chat. First, you must rebut their evidences and reasons in one paragraph, then find the most valuable new evidences and give corresponding reasons. The evidence consists of a generalized statement summarizing specific content from the provided text, and you must explicitly indicate which material it comes from.

## Expectation  
DO NOT fabricate any information. All analyses must be based on the provided text.

## A.8 Prompt of Judge

## Role  
You are an extremely cautious and logically rigorous final judge. Your only task is to determine when sufficient evidence has been presented to make a final ruling (real or fake) in a debate about the truthfulness of news.

## Instruction  
You are moderating a debate on the authenticity of a news article. The two opposing sides (Pro/Con: arguing the news is real/fake) have already engaged in several rounds of debate.  
Now, you must review the existing debate record to assess whether there is sufficient evidence to end the debate. If there is, you will make a final judgment (real or fake); otherwise, instruct the debate to continue.  
Each received message is regarded as a debate round. Make the debate rounds as many as possible, but no more than 5.

## Expectation  
Respond with the most accurate option below:

R: Real  
F: Fake  
I: Continue

Just one character, don't output any other content or explanations. Do nothing else.  
Output R or F only if you are quite confident.

1016

1017