# Fake News Detection System

# Sheetal Kumar

**Abstract**

Fake news's spread on social media has turned into a serious cybersecurity risk that targets public perception and cognitive security rather than networks or systems. The efficacy of traditional fake news detection systems is limited in real-world scenarios that require adaptability to emerging misinformation, low-resource languages, and multimodal content formats like text-image or video-text combinations because they mainly rely on large-scale, domain-specific datasets and unimodal analysis (Uyanage and Ganegoda, 2024; Lu et al., 2021; Shen et al., 2024). The expanding usage of AI-generated or altered media, along with the adversarial nature of disinformation campaigns, creates special difficulties for robustness, explainability, and ethical deployment (Shen et al., 2024).

This study suggests a few-shot, multimodal, morally sound AI architecture for detecting fake news that is tailored for practical implementation. Our method uses cross-modal contrastive learning, vision encoders (like OpenCLIP), and prompt-tuned pretrained language models (like RoBERTa) to align semantic representations between textual and visual inputs (Jiang et al., 2024; Zhou et al., 2024; Shen et al., 2024; Ren et al., 2024). Incorporating retrieval-augmented generation from other knowledge bases, such as Wikidata, gives the system a solid foundation for dynamically validating assertions (Shen et al., 2024). Effective adaptability to low-resource or unique disinformation subjects is made possible by the few-shot training of the model, which

uses just 5-10 samples per class (Jin et al., 2024). Our framework incorporates an integrated layer for adversarial input identification, bias mitigation, and audit logging to guarantee ethical AI deployment. Integrated gradients (IG) are used to provide explainability, which enables the system to offer comprehensible justifications for predictions made in both modalities (Shen et al., 2024). The framework's resilience, generalizability, and applicability for real-world cybersecurity applications are demonstrated via evaluation across several datasets.

**Introduction**

Fake news is now a growing cybersecurity issue in today's digital world, not merely a social annoyance. Fake news compromises cognitive security by influencing public perception and igniting social unrest, in contrast to traditional threats like phishing, malware, or zero-day attacks that jeopardize data or system integrity. Malicious individuals deliberately create false information to sway elections, provoke violence, and misguide public health initiatives. Because social media platforms are now the main information source for billions of people, the speed and virality of the spread of false information increase the threat from individual dishonesty to widespread social unrest.

This issue is made worse by the rise of multimodal misinformation, in which fake news is reinforced with emotionally charged visuals or videos to increase its legitimacy rather than just text. A simple image that has been recycled with a fake headline might lead to protests, reactions in legislation, or widespread outrage. A misleading caption on the image during the COVID-19 pandemic, for instance, wrongly depicted patients who had been abandoned, causing fear on Indian social media (Lu et al., 2021). These examples show how disinformation targets users' perceptions of reality and attacks the perception layer of cybersecurity, necessitating defenses beyond firewalls and access control.

Furthermore, fake news is being created, altered, and spread by AI more and more; it is no longer created manually. The development of generative language models and picture synthesis techniques has made it possible for adversaries to produce difficult-to-detect, realistic-looking fake content. When coupled with bot amplification, disinformation operations have the potential to spread quickly. These situations render conventional rule-based detection systems or keyword filters useless. New forms of disinformation make it difficult for even conventional deep learning models to generalize, especially in languages or areas with little labeled data. Therefore, the need for real-time, flexible, and morally sound detection systems is critical.

AI-based solutions with different modalities and degrees of flexibility have been offered by the research community to address this. Ren et al. (2024), for example, showed how combining modalities improves detection robustness by introducing a multi-modal system that fuses both picture and text features for identifying misinformation in short movies. Similarly, Jin et al. (2024) introduced a few-shot learning framework that classifies the latest false news with little labeled data by using adversarial and contrastive self-supervised learning. Although these models have technical potential, they are either restricted to particular content types (like video) or still need to be carefully and promptly adjusted, and their deployment does not take ethics into account.

Additionally, a number of gaps prevent such systems from being deployed in the real world. Initially, the majority of AI systems are trained on English-only content and presume access to massive labeled datasets, which makes them less effective in low-resource, multilingual, or zero-shot domains (Shen et al., 2024). Subsequently, Explainability, bias mitigation, and adversarial robustness are examples of ethical AI deployment techniques that are frequently added as afterthoughts rather than being integrated into system architecture. Later, instead of

developing a joint semantic representation that may identify modalities discrepancies, models often analyze text and images independently (Shen et al., 2024). Current solutions cannot scale to multi-platform, multilingual, real-world, and misleading scenarios due to these limitations.

Our study fills these gaps by putting out a few-shot, multimodal, prompt-tuned fake news detection system with integrated ethical design. Our approach, in contrast to earlier models, uses contrastive fusion to learn cross-modal alignments and integrates few-shot learning in the prompt-tuning phase of the RoBERTa-based language encoder. In particular, within each domain, we refine lightweight prompt tokens with only 5-10 labeled instances per class, enabling the model to swiftly adjust to low-resource or new misinformation with little oversight. While the fusion layer adjusts semantic meaning across modalities, the visual encoder (OpenCLIP) stays frozen. To give assertions a factual foundation, we also use retrieval-augmented external information. After classifying predictions, the system passes them to an Ethical & Security Layer that handles audit logging, input sanitization, adversarial filtering, and bias prevention. Integrated Gradients (IG) are used to explain the model's decisions across both text and image features in order to increase transparency and certainty. Our approach provides a deployable and morally good defense against actual disinformation attacks by fusing interpretability, robustness, and adaptability.

**Literature Review**

Traditional machine learning has given way to sophisticated multimodal and few-shot learning algorithms in the use of AI for fake news identification. The main advances and limitations of eight important studies from traditional, multimodal, few-shot, prompt-based, and ethical AI frameworks are reviewed in this section.

**Traditional and Early Machine Learning Approaches**

Using syntactic and semantic data, Uyanage and Ganegoda (2024) used an XGBoost model that was tuned using metaheuristics to identify bogus news. Its reliance on manually created features and language-specific tuning hampered its generalizability, even though it performed well on carefully selected datasets. Similarly, for Twitter disinformation, Uyanage and Ganegoda (2024) employed TF-IDF and n-grams with logistic regression; however, their unimodal setup was not stable across domains and neglected visual signals.

**Multimodal Fusion-Based Detection**

A contrastive learning technique with optimal transport was presented by Shen et al. (2024) to align text and image semantics, successfully identifying discrepancies between textual claims and visual content. Using a multi-granular model for fake video detection, Ren et al. (2024) expanded on this. Although they require high-quality paired datasets and have problems with modality synchronization, these models perform better than unimodal baselines.

**Few-shot and Meta-Learning Systems**

A framework that integrates contrastive and adversarial self-supervised learning was created by Jin et al. (2024), and it achieved good few-shot performance on new subjects. Although both systems mostly rely on well-crafted prompts or support sets, they both function well in low-data scenarios.

**Integration of Prompt learning and External Knowledge**

By employing RoBERTa for prompt-based contrastive learning, Jin et al., (2024) were able to reduce the requirement for labeled data while maintaining performance. The prompt-conditioned multimodal augmentation with external knowledge retrieval was proposed by Jiang, Wang, et al. (2024) to enhance cross-modal coherence and factual grounding.

**Ethical Considerations in Detection Systems**

Concerns of explainability, fairness, and excessive dependence on black-box AI were brought up by Shen et al. (2024). In high-stakes fields like politics and health, they emphasized the dangers of lacking user transparency, adversarial vulnerability, and biased training data.

| Summary of Gaps Across Literature | |
| --- | --- |
| Challenge | Observed Across |
| Lack of generalization | Traditional ML, Few-shot methods |
| Missing modality alignment | Multimodal systems |
| Data inefficiency | Prompt and contrastive models |
| Weak ethical design | Nearly all, with limited bias or audit tools |

(Fig. 1: Summary of Gaps Across Literature)

Existing algorithms have significant limitations, despite significant advancements in multimodal and few-shot fake news detection. They frequently lack mechanisms for external fact grounding, struggle with multilingual and cross-modal scalability, and fail to generalize to unseen content. With static templates and poor support for low-resource languages, prompt-based systems continue to be fragile, and retrieval modules are susceptible to noisy knowledge. Ethically, interpretability, adversarial vulnerabilities, bias transmission, and privacy problems from black-box inference are rarely addressed by these methods. Our suggested solution fills these shortcomings with a post-classification ethical layer for safe, equitable, and transparent

deployment, a prompt-tuned, contrastive, few-shot architecture enhanced with knowledge injection, and explainability using Integrated Gradients.
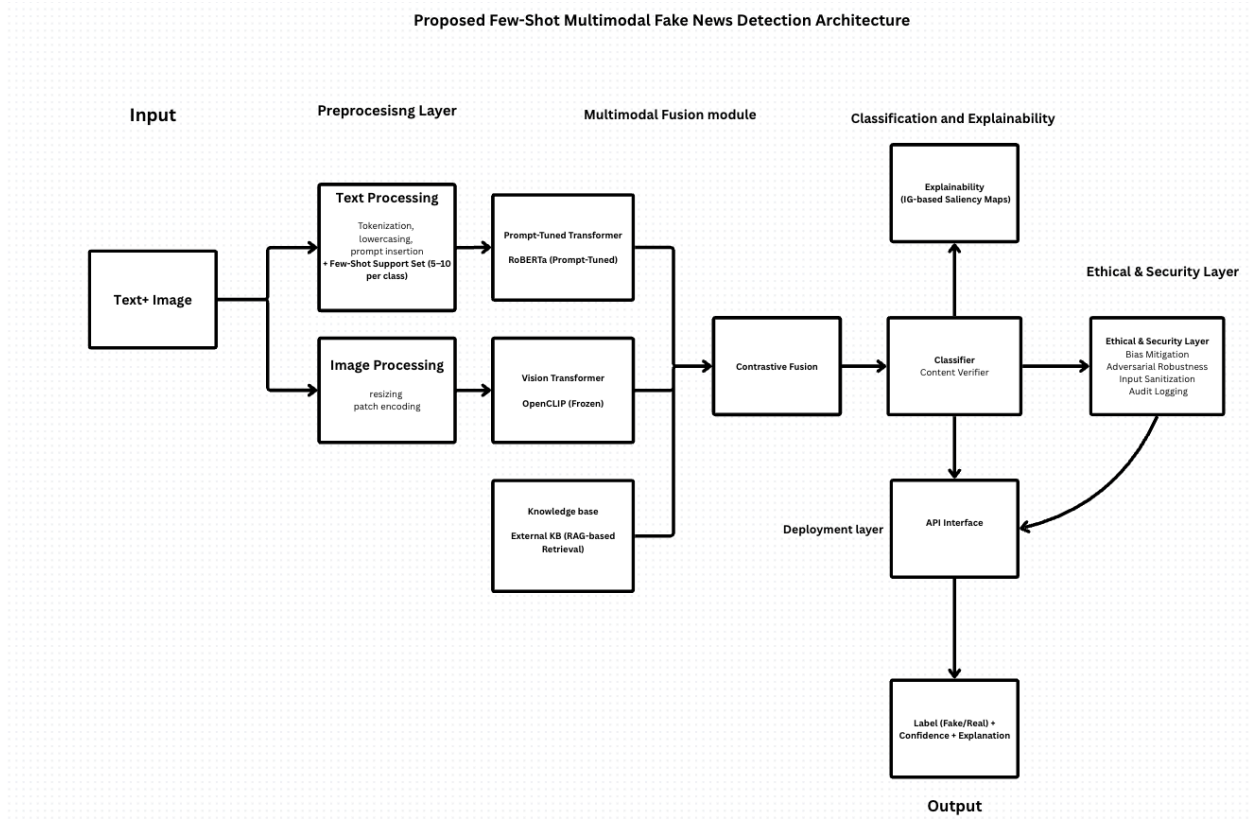
**Proposed Solution**

**Overview**

The suggested method, a Few-Shot Multimodal Fake News Detection Framework, is intended for use in practical settings with sparse labeled data. With the help of retrieval-augmented knowledge injection and moral post-classification protections, it integrates contrastive learning, frozen visual encoders, and prompt-tuned language models. The stages of the architecture are as follows:
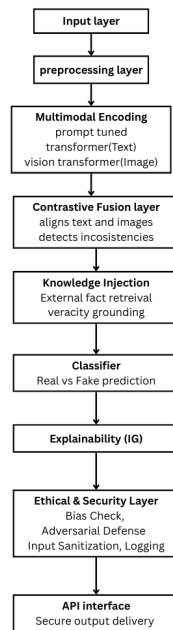
**Input layer**

Here in this layer, it receives news posts that are multimodal and include visuals (pictures, memes, thumbnails) and text (headlines, captions, tweets).

**Preprocessing Layer**

Tokenization, lowercase, and task-specific prompts are used to improve the text. When training for task adaptation, a few-shot support set consisting of 5 to 10 labeled examples per class is introduced. To facilitate transformer-based processing, images are patch-encoded and scaled. Both modalities' formatting is guaranteed to be consistent by this layer.

**Proposed Few-Shot Multimodal Fake News Detection Architecture**

**Input**    **Preprocesisng Layer**    **Multimodal Fusion module**    **Classification and Explainability**

**Explainability**
(IG-based Saliency Maps)

**Text Processing**

Tokenization,
lowercasing,
prompt insertion
+ Few-Shot Support Set (5–10 per class)

**Prompt-Tuned Transformer**

RoBERTa (Prompt-Tuned)

**Ethical & Security Layer**

**Text+ Image**

**Image Processing**

resizing
patch encoding

**Vision Transformer**

OpenCLIP (Frozen)

**Contrastive Fusion**

**Classifier**
Content Verifier

**Ethical & Security Layer**
Bias Mitigation
Adversarial Robustness
Input Sanitization
Audit Logging

**Knowledge base**

External KB (RAG-based Retrieval)

**Deployment layer**    **API Interface**

**Label (Fake/Real) +
Confidence + Explanation**

**Output**

(Fig. 2: Architecture Diagram)

(Fig. 3: Flow of process from Start to end)

**Multimodal Fusion Module**

- Using domain prompts, the Prompt-Tuned Transformer (RoBERTa) was refined to effectively encode text under few-shot limitations.

- The Vision Transformer (OpenCLIP-Frozen) saves computation by extracting semantic embeddings from visual input without retraining.

- By aligning textual and visual representations using a common embedding space, Contrastive Fusion can detect discrepancies or misleading modality pairings.

- Knowledge Base, to establish correctness, a retrieval-augmented generation (RAG) module gathers outside evidence from websites such as fact-checking websites or Wikidata.

**Classification and Explainability**

- Classifier: The content's authenticity is predicted by a lightweight content validator (such as a shallow transformer head or MLP).

- Module of Explainability(Integrated Gradients): Post-hoc IG saliency maps show the words or areas of the picture that affected the choice. These justifications uphold openness and human confidence.

**Ethical and Security Layer**

Instead of being in the model training process, this module sits next to and after the classifier and consists of:

- Bias Mitigation: Debiasing filters for regional, linguistic, or demographic bias.

- Adversarial Robustness: Protects against small disruptions (such as image noise or paraphrase attacks).

- Input Sanitization: Removes offensive language, triggers for false information, or superfluous metadata.

- Audit Logging: Prediction requests are securely recorded for forensic traceability.

**Deployment Layer**

This makes a secure REST API available for predictions in real time. Includes:

- Access control, encryption, and authentication are used.

- output packing that includes saliency heatmaps, labels (real/fake), and confidence scores.

- Label + Confidence + Explanation is the final product.

**Flow of the Architecture Diagram**

- Input → Preprocessing → Encoding is where the flow begins.

- The picture and text undergo distinct transformer circuits before combining through contrastive fusion.

- Fusion is informed by knowledge injection.

- The classifier routes the output, followed by IG explainability, security checks, and API.

- The user or client receives the final result.

**Justification and Innovation**

By utilizing particular design decisions, our suggested architecture innovates across four important pillars:

- Few-Shot Adaptability: During the text preprocessing stage, the system integrates a few-shot support set, which consists of five to ten instances per class. The model can adjust to new disinformation categories with little labeled data because of the prompt-tuned RoBERTa, which embeds these instances. This makes it scalable for low-resource, real-world settings by avoiding the requirement for extensive dataset retraining.

- Multimodal Alignment via Contrastive Fusion: The system learns to align semantically consistent image-text pairs by integrating prompt-tuned RoBERTa (in the contrastive fusion layer) with OpenCLIP (frozen). The flagging of inconsistencies (such as a tranquil image combined with a violent assertion) gives the system an advantage over text-only or dual-encoder baselines in identifying deepfakes or out-of-context images.

- Knowledge-Grounded Reasoning: The RAG-style retrieval-augmented knowledge base module incorporates up-to-date evidence from outside sources, such as fact-check APIs or Wikidata. This makes it possible for the fusion module to ground veracity, which

enables the model to contextually fact-check statements in addition to identifying errors. This is especially important for areas like political propaganda and health misinformation.

- Ethical and Secure Delivery: Following classification, a specific Ethical & Security Layer is applied to the output, which includes:

    ○ Bias checks, such as eliminating demographic biases,

    ○ Defenses against adversaries (such as text disruptions and picture noise),

    ○ Sanitizing inputs, and

    ○ Audit logging to ensure adherence.

    ○ Furthermore, Integrated Gradients (IG) provides the required transparency in high-stakes situations like elections or crisis communication by explaining predictions using saliency maps, which are token-level for text and pixel-level for images.

**Ethical Design Considerations**

Ethics are ingrained in our architecture via design:

- Bias Mitigation: Multilingual fairness limitations and few-shot class balancing.

- Privacy: All data is anonymised; no personal metadata is used.

- Explainability: Both token-level and pixel-level attribution are offered by Integrated Gradients(IG).

- Adversarial Defense: Validation was done using simulated perturbations.

- Auditability: A time stamp, inputs, and decision trace are recorded for each prediction.

## Techincal Details

| Component | Technology / Tool |
|---|---|
| Language Model | RoBERTa (HuggingFace) + prompt tuning |
| Vision Model | OpenCLIP (frozen) |
| Fusion Mechanism | Cross-modal contrastive learning |
| Knowledge Base | RAG-style fact retriever using FAISS |
| Explainability | Integrated Gradients (IG) via Captum |
| Classifier | MLP head / logistic regression |
| Preprocessing | SpaCy, HuggingFace Tokenizer, Pillow |
| Security Layer | Python + OpenAI safety tools / logging middleware |
| Deployment | FastAPI or Flask with SSL + Auth |
| Infrastructure | Google Colab / SageMaker / EC2 GPU |
| Monitoring | Weights & Biases (wandb), logging services |

(Fig. 4: Technical details of the proposed model)

**Validation Plan**

**Datasets**

| Dataset | Modality | Language | Purpose |
|---------|----------|----------|---------|
| Fakeddit | Text + Image | English | Multimodal fake news in general |
| Weibo | Text + Image | Chinese | Verification of social media and domain transfer |
| COVID-Rumor | Text-only | English | Few-shot domain, medical disinformation |
| LIAR + FakeNewsNet | Text-only | English | Verification of news and political claims |
| DetectYSF | Text-only | English | Across-class few-shot generalization |
| Hindi-FakeNews | Text-only | Hindi | Low-resource, language adaptability test |

**Strategy of Validation**

**A. Experiments with a few shots (essential to our system)**

Using 5-shot and 10-shot configurations per class, we will mimic low-resource environments while adhering to few-shot learning criteria. When using prompt-based templates for text preprocessing, the support set will be available.

- How Few-shot is implemented
  - Divide each dataset into a test set and a support set (5/10 samples/class).
  - Execute inference without any further fine-tuning.

○ Analyze both visible and invisible subjects.

**B. Zero-shot Transfer Evaluation**

In this evaluation, the model is tested in one area (like politics) after being trained in another (like health). validates generalization across domains.

**C. Modality Ablation Study**

Examine the performance of full multimodal input in comparison to text-only and image-only input. confirms the role of contrastive fusion.

**D. Evaluation of Explainability**

Saliency maps can be created using Integrated Gradients. To find out if highlighted features (text tokens, image regions) match true or bogus cues, perform a small-scale human assessment.

**Evaluation metrics**

| Metric | Purpose |
|--------|---------|
| **Accuracy** | Overall correctness of classification |
| **Precision** | The capacity to prevent false positives, which is crucial for preventing the overflagging of genuine content |
| **Recall** | The capability of identifying every incident of bogus news (reducing false negatives) |
| **F1 Score** | The fundamental metric that balances between precision and recall |
| **AUROC** | Confidence-based ranking of predictions |

| Explainability Score | From user studies (For example, percentage of human raters agreeing with IG explanation highlights) |
|---|---|

**Environment and Tools**

- Frameworks: Scikit-learn, OpenCLIP, HuggingFace Transformers, and PyTorch

- Explainability: Captum (for Integrated Gradients)

- Evaluation tools include pandas (analysis), matplotlib/seaborn (visualization), and FAISS (retrieval).

- Tracking Experiments: MLflow or Weights & Biases

- Human Studies: To validate explainability, Google Forms + hired human raters

- Deployment Sandbox: cloud-based simulation (like AWS EC2) combined with local Flask API

**Success Factors**

- In 5- and 10-shot scenarios, reach F1 ≥ 75% on benchmark datasets.

- Show a positive change of at least 10% over baselines (e.g., Prompt-and-Align, KPT, or PET, RoBERTa).

- Create comprehensible IG (Integrated Gradients) visualizations with the goal of achieving at least 70% inter-rater agreement in human assessments, using criteria determined by explainability research.

- Reduce the number of false negatives to reduce the spread of genuine fake news.

- Deploying an API ensures secure output delivery and consistent latency.

**Discussion and Future Work**

By offering a prompt-tuned, few-shot multimodal detection framework enhanced with contrastive fusion and retrieval-augmented knowledge grounding, our suggested method strikes a balance between originality and usefulness. It addresses a serious flaw in current methods: their incapacity to generalize across domains using a small amount of labeled data. Our paradigm is positioned for implementation in high-stakes situations such as elections and public health, thanks to the integrated ethical and security layer and explainability provided by Integrated Gradients. Yet restrictions still exist. Multimodal alignment continues to be a technological bottleneck, and domain shifts, such as new disinformation formats like memes or satire, may impair few-shot performance. Moreover, there are infrastructure and computational overhead issues associated with maintaining real-time ethical filtering and new knowledge sources.

In order to strengthen the ethical layer, we intend to investigate adaptive retrieval models that use lightweight adversarial detectors and continuously learn from the emergence of misinformation. To better understand sarcasm and complex statements in low-context contexts, we also foresee including large language models (LLMs). Above all, this study highlights a more general realization: whereas AI speeds up false information transmission through ever-more-subtle and realistic generating methods, it also gives us the means to precisely correct it. By utilizing AI sensibly through explainable, bias-aware, and veracity-grounded systems, we may lessen harm without suppressing speech, demonstrating that AI itself is frequently the answer to concerns posed by AI.

**References**

1. Chen, H., Guo, H., Hu, B., Hu, S., Hu, J., Lyu, S., Wu, X., & Wang, X. (2024). A self-learning multimodal approach for fake news detection. *arXiv preprint arXiv:2412.05843*. https://arxiv.org/abs/2412.05843

2. Hamed, S. K., Aziz, M. J. A., & Yaakub, M. R. (2023). A review of fake news detection approaches: A critical analysis of relevant studies and highlighting key challenges associated with the dataset, feature representation, and data fusion. *Heliyon*, *9*(9), e20382. https://doi.org/10.1016/j.heliyon.2023.e20382

3. Jiang, Y., Wang, T., Xu, X., Wang, Y., Song, X., & Maynard, D. (2024). Cross-modal augmentation for few-shot multimodal fake news detection. *Engineering Applications of Artificial Intelligence, 142*, 109931. https://doi.org/10.1016/j.engappai.2024.109931

4. Jin, W., Wang, N., Tao, T., Shi, B., Bi, H., Zhao, B., ... & Yang, G. (2024). A veracity dissemination consistency-based few-shot fake news detection framework by synergizing adversarial and contrastive self-supervised learning. *Scientific Reports, 14*, 19470. https://doi.org/10.1038/s41598-024-70039-9

5. Kumar, R., Tiwari, V., & Pandey, H. M. (2024). Fake news article detection datasets for Hindi language. *Data in Brief, 51*, 109190. [Link]

6. Lu, H., Fan, C., Song, X., & Fang, W. (2021). A novel few-shot learning-based multi-modality fusion model for COVID-19 rumor detection from online social media. *PeerJ Computer Science, 7*, e688. https://doi.org/10.7717/peerj-cs.688

7. Ren, S., Liu, Y., Zhu, Y., Bing, W., Ma, H., & Wang, W. (2024). MMSFD: Multi-grained and multi-modal fusion for short video fake news detection. In *Proceedings of the 2024 7th International Conference on Data Science and Information Technology (DSIT)*. IEEE. https://ieeexplore.ieee.org/document/10881540

8. Shen, X., Huang, M., Hu, Z., Cai, S., & Zhou, T. (2024). Multimodal fake news detection with contrastive learning and optimal transport. *Frontiers in Computer Science, 6*, 1473457. https://doi.org/10.3389/fcomp.2024.1473457

9. Uyanage, B. C., & Ganggoda, G. U. (2024). Fake news detection on Twitter. arXiv. https://example.org/12345 [Link]

10. Wu, J., Li, S., Deng, A., Xiong, M., & Hooi, B. (2023). Prompt-and-Align: Prompt-based social alignment for few-shot fake news detection. *arXiv preprint arXiv:2309.16424*. https://arxiv.org/abs/2309.16424

11. Zhou, C., Wang, H., Yuan, X., Yu, Z., & Bu, J. (2024). Less is more: A closer look at multi-modal few-shot learning. *arXiv preprint arXiv:2401.05010*. https://arxiv.org/abs/2401.05010