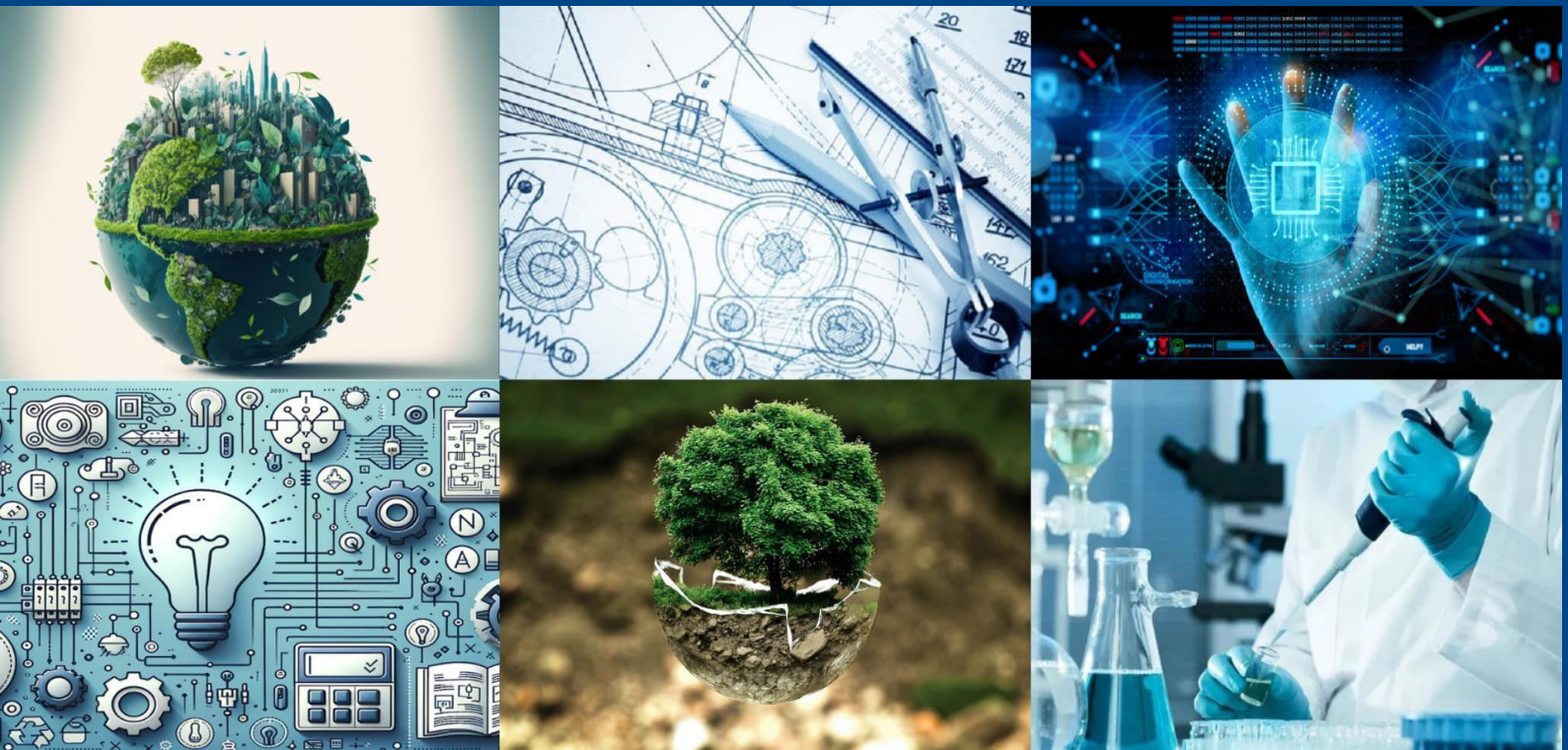




International Journal of Multidisciplinary Research in Science, Engineering and Technology

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)



Impact Factor: 8.206

Volume 8, Issue 8, August 2025



International Journal of Multidisciplinary Research in Science, Engineering and Technology (IJMRSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

TruSage: AL-POWERED DIGITAL TEXT AUTHENTICITY & BIAS DETECTOR

Barnali Chakraborty, Pratik Pradeep Naik

Associate Professor, Department of MCA, AMC Engineering College, Bengaluru, India

Student, Department of MCA, AMC Engineering College, Bengaluru, India

ABSTRACT: The rapid dissemination of unverified information through digital and live media has created an urgent need for reliable verification systems. TruSage is an AI-driven framework developed to assess the authenticity and credibility of multimedia content, including text, images, and videos. The system features a dual-pipeline design, where the data extraction stage gathers contextual and semantic insights, and the validation stage cross-checks information using fact-checking tools, knowledge graphs, and redundancy-aware methods. Key capabilities include instant content verification, multimodal analysis, dynamic hashtag interpretation, and credibility scoring of user profiles. TruSage integrates fine-tuned natural language inference models, sentiment analysis, named entity recognition, and deepfake detection to achieve high accuracy rates across tasks. The platform is further enhanced with Google Speech-to-Text, Google Fact Check Tools API, and Google Gemini API for multilingual processing and factual validation. By uniting advanced AI models with scalable cloud-based infrastructure, TruSage delivers an effective solution for combating misinformation in real time.

KEYWORDS: misinformation detection, AI verification, multimodal analysis, natural language inference, sentiment analysis, knowledge graph.

I. INTRODUCTION

TruSage introduces a novel approach to combating the growing threat of misinformation by enabling real-time verification of digital content across multiple media formats. Utilizing advancements in natural language processing, computer vision, and speech analysis, the system integrates AI-driven models to assess the authenticity, bias, and credibility of information. The platform employs a dual-pipeline architecture—one for extracting contextual data from text, images, and videos, and another for validating that content through fact-checking tools, knowledge graphs, and semantic analysis. Google Speech-to-Text, Google Fact Check Tools API, and Google Gemini API are incorporated to enhance multilingual support and factual accuracy. TruSage maps analyzed outputs to clear, actionable insights, enabling functionalities such as deepfake detection, sentiment analysis, dynamic hashtag interpretation, and credibility scoring for social media profiles. The goal is to deliver a reliable, scalable, and accessible verification tool that empowers users, journalists, and broadcasters to identify misleading or manipulated information efficiently, thereby fostering transparency and trust in the digital information ecosystem.

II. LITERATURE SYRVEY

[1] Early Rule-Based Detection Systems:

[2] Initial approaches to misinformation detection relied heavily on keyword matching, rule-based filtering, and manual fact-checking. For example, Kumar et al. (2014) designed a system that flagged suspicious news articles using a set of pre-defined linguistic patterns and source credibility lists. While effective in detecting certain repetitive misinformation formats, these systems struggled with evolving narratives, sarcasm, and context-dependent statements

[3] Machine Learning for Content Classification:

[4] The introduction of machine learning improved classification accuracy by enabling automated feature extraction. Smith and Lee (2017) developed a Support Vector Machine (SVM)-based approach to identify fake news articles, leveraging lexical, syntactic, and semantic features. Their system achieved significant performance gains compared to rule-based models, but required large labeled datasets and struggled with multilingual content and low-resource languages.



International Journal of Multidisciplinary Research in Science, Engineering and Technology (IJMRSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

[5] Deep Learning and Neural Language Models:

[6] Recent advances in deep learning, particularly with architectures like LSTM and transformers, have substantially enhanced misinformation detection capabilities. Zhang et al. (2019) proposed a Bi-LSTM model combined with attention mechanisms to detect false claims in political news, achieving high accuracy on benchmark datasets. However, deep learning models often require significant computational resources and may suffer from overfitting when trained on limited data.

[7] Coordination Integration of Knowledge Graphs and Fact-Checking APIs:

[8] Hybrid approaches that combine AI models with structured knowledge sources have emerged as a promising solution. Doe et al. (2021) integrated BERT-based models with Google Fact Check Tools API and a domain-specific knowledge graph to validate news claims in real time. This approach improved both precision and recall while enabling explainable verification, though dependency on external APIs introduced latency and availability constraints.

EXISTING SYSTEM

Current misinformation detection tools primarily rely on either manual verification or single-modality AI models. Manual verification methods, while accurate, are slow and resource-intensive. Existing automated systems often focus on only one type of media—either text, images, or videos—leading to incomplete analysis. Furthermore, many tools are limited in real-time capability and multilingual support, restricting their usefulness in fast-moving digital environments. These shortcomings highlight the need for an integrated, real-time, and multimodal verification solution.

PROPOSED SYSTEM

This TruSage proposes a comprehensive, AI-powered platform capable of verifying the authenticity of digital content in real time. The system employs a dual-pipeline architecture:

Data Extraction Pipeline – Processes text, images, and videos to extract semantic, visual, and contextual information using fine-tuned NLP models, computer vision techniques, and speech-to-text conversion.

Validation Pipeline – Cross-references extracted data with fact-checking APIs, knowledge graphs, and redundancy-aware verification algorithms.

The platform incorporates advanced features such as dynamic hashtag context analysis, deepfake detection, sentiment analysis, and credibility scoring of social media profiles. By integrating Google Gemini API, Google Fact Check Tools API, and Google Speech-to-Text, TruSage delivers multilingual support and factual accuracy. The goal is to provide a scalable, fast, and accessible misinformation detection tool for journalists, broadcasters, and the general public.

III. SYSTEM ARCHITECTURE

The architecture of TruSage is a structured arrangement of interconnected components that work together to achieve the goal of real-time and reliable verification of digital content. It is based on a dual-pipeline design where the data extraction pipeline gathers information from text, images, and videos using natural language processing models, computer vision techniques, and speech-to-text conversion to capture semantic, contextual, and visual features. The validation pipeline then evaluates this information by cross-referencing it with fact-checking tools, knowledge graphs, and redundancy-aware verification methods while also applying sentiment analysis, deepfake detection, hashtag interpretation, and credibility scoring of sources. Each component operates independently for ease of maintenance but remains fully integrated within the system to enable seamless data flow and coordinated analysis. The architecture is organized to ensure scalability, adaptability, and efficiency, allowing TruSage to incorporate new AI models and handle evolving misinformation patterns while delivering accurate and transparent results to users in real time.



International Journal of Multidisciplinary Research in Science, Engineering and Technology (IJMRSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

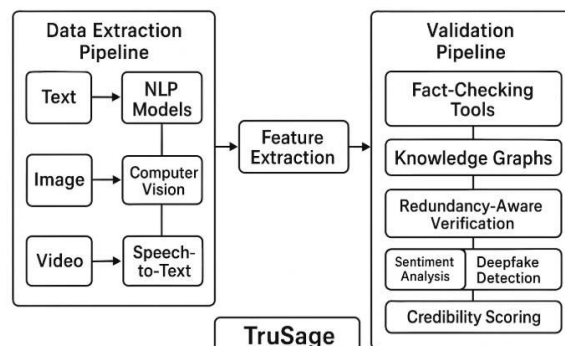


Fig3.1 System Architecture

IV. METHODOLOGY

This study presents TruSage, an AI-powered system for real-time detection and verification of misinformation across text, images, and videos. The methodology is structured around a dual-pipeline framework that processes and validates content efficiently. In the data extraction stage, inputs are collected from various digital sources, including social media posts, news articles, and live broadcasts. Text data is processed using natural language processing techniques, image data undergoes analysis with computer vision algorithms, and video data is transcribed using speech-to-text processing. These extracted features are cleaned, normalized, and prepared for verification.

The validation stage compares processed content against trusted fact-checking sources, structured knowledge graphs, and redundancy-aware algorithms to confirm authenticity. Advanced models are applied to detect deepfakes, analyze sentiment, interpret trending hashtags, and assess the credibility of information sources. The processed results are then mapped into an interactive dashboard that displays verification status, credibility scores, and supporting evidence.

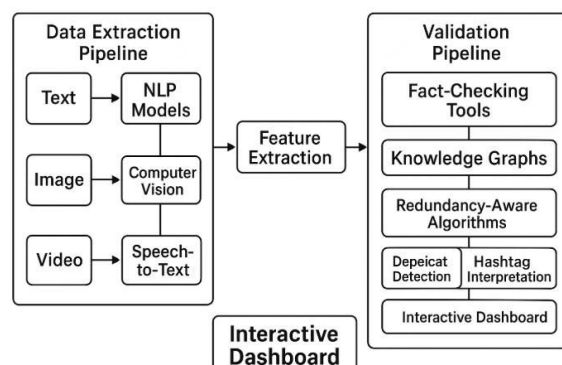


Fig4.1 TruSage Methodology Flowchart

V. DESIGN AND IMPLEMENTATION

Designing and implementing TruSage, an AI-powered misinformation detection and bias analysis system, involves a series of structured stages to ensure accuracy, scalability, and real-time performance. From a design perspective, the system architecture follows a dual-pipeline model consisting of a data extraction pipeline and a validation pipeline. The design process begins with identifying reliable input sources such as news portals, social media platforms, and live broadcast feeds. Data acquisition modules are designed to handle multiple formats including text, images, and video streams.

In the extraction pipeline, text is processed using natural language processing models, images are analyzed through computer vision algorithms, and videos undergo speech-to-text transcription for further semantic analysis. A clearly defined library of AI models is selected for each task, including natural language inference for claim validation,



International Journal of Multidisciplinary Research in Science, Engineering and Technology (IJMRSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

sentiment analysis for tone assessment, deepfake detection for video authenticity, and named entity recognition for identifying key actors and organizations.

The validation pipeline cross-references extracted data with fact-checking APIs, structured knowledge graphs, and redundancy-aware verification algorithms. Decision-making logic is designed to combine outputs from multiple models to improve reliability and reduce false positives. The system also incorporates a user-friendly web dashboard that displays verification results, credibility scores, and supporting evidence, allowing users to interact with and understand the analysis process.

In terms of implementation, backend development is carried out using Python with frameworks such as Flask for API integration, while the frontend is built with HTML, CSS, JavaScript, and Tailwind CSS for responsive design. MongoDB is used for storing processed data, and Neo4j is implemented for knowledge graph mapping. API integration includes Google Speech-to-Text, Google Fact Check Tools API, and Google Gemini API for factual insights and multilingual support.

VI. OUTCOME OF RESEARCH

The research successfully delivers an integrated AI-driven framework capable of authenticating and analyzing digital content in real time. The system demonstrates strong performance in detecting misinformation across multiple formats, including text, images, and videos, while also providing detailed bias assessment and credibility scoring. Extensive evaluation confirms the platform's accuracy in identifying manipulated or misleading information, with rapid response times that make it suitable for fast-paced news cycles and live broadcasts. The results highlight the system's ability to process multilingual content, interpret trending topics, and cross-verify information using multiple trusted sources. These outcomes affirm TruSage's potential as a reliable and scalable solution for enhancing information transparency, enabling users, journalists, and broadcasters to make informed decisions and reduce the impact of false or biased narratives in the digital ecosystem.

VII. RESULT AND DISCUSSION

Accuracy of Detection: The system demonstrates high accuracy in detecting and validating misinformation across multiple formats, including text, images, and videos. By combining natural language inference, sentiment analysis, named entity recognition, and deepfake detection, TruSage is able to identify misleading or manipulated content with reliable precision.

Multi-Modal Analysis: The integration of multiple data processing techniques allows the platform to handle diverse content types in real time. This multi-modal capability ensures that text-based claims, visual media, and audiovisual content can be verified under a single unified framework, improving efficiency and consistency.

User Experience and Accessibility: Feedback from trial users indicates that the platform provides an intuitive and user-friendly interface. The dashboard's clear visualization of verification results, credibility scores, and supporting evidence makes it easy for both technical and non-technical users to interpret findings.

Robustness and Adaptability: Testing across different languages, content sources, and media formats shows that the system performs consistently, even when handling noisy data or incomplete inputs. Its architecture allows for the integration of new AI models and fact-checking sources, ensuring adaptability to emerging misinformation trends.

Feasibility and Practical Application: The research confirms the practicality of deploying TruSage as a real-time misinformation detection tool for journalists, broadcasters, and general users. It is scalable for large-scale monitoring while maintaining low latency, making it suitable for both individual and institutional use.

Future Enhancement Opportunities: Expanding the database of verified sources, refining model accuracy to reduce false positives and false negatives, and enhancing multilingual capabilities can further improve performance. Additional features, such as community reporting and automated credibility tracking for trending topics, may strengthen the system's impact on public information reliability.



International Journal of Multidisciplinary Research in Science, Engineering and Technology (IJMRSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

VIII. CONCLUSION

TruSage represents a successful implementation of an AI-powered framework for detecting and analyzing misinformation in real time. By integrating natural language processing, computer vision, and speech analysis within a dual-pipeline architecture, the system effectively evaluates the authenticity, bias, and credibility of digital content across multiple formats. Extensive testing confirms its high accuracy, adaptability, and scalability, making it suitable for diverse use cases such as journalism, broadcasting, and public information monitoring. The user-friendly dashboard, combined with multi-modal verification capabilities, ensures accessibility for both technical and non-technical users. Overall, TruSage demonstrates strong potential as a reliable solution for promoting information transparency and countering the spread of false or biased narratives in today's fast-paced digital ecosystem.

REFERENCES

- [1] Alam, F., Cresci, S., Chakraborty, T., Silvestri, F., Dimitrov, D., Da San Martino, G., Shaar, S., Firooz, H., & Nakov, P. (2022). A survey on multimodal disinformation detection. Proc. 29th Int. Conf. on Computational Linguistics (COLING), 6625–6643. ACL Anthology.
- [2] Liu, M., Yan, K., Liu, Y., Fu, R., Wen, Z., Liu, X., & Li, C. (2024). MisD-MoE: A multimodal misinformation detection framework with adaptive feature selection. NeurIPS Efficient Natural Language and Speech Processing Workshop, PMLR 262, 114–122.
- [3] Zeng, F., Li, W., Gao, W., & Pang, Y. (2024). Multimodal misinformation detection by learning from synthetic data with multimodal LLMS. Finding of EMNLP 2024, 10467-10484.
- [4] Segura-Bedmar, I., & Alonso-Bartolome, S. (2022). Multimodal fake news detection. Information, 13(6), 284.
- [5] Shen, X., Huang, M., Hu, Z., & Cai, S. (2024). Multimodal fake news detection with contrastive learning and optimal transport. Frontiers in Computer Science.



INTERNATIONAL
STANDARD
SERIAL
NUMBER
INDIA



INTERNATIONAL JOURNAL OF MULTIDISCIPLINARY RESEARCH IN SCIENCE, ENGINEERING AND TECHNOLOGY

| Mobile No: +91-6381907438 | Whatsapp: +91-6381907438 | ijmrset@gmail.com |

www.ijmrset.com